

2. Algebraic review (focus on real, symmetric, positive)

* Transpose of a matrix

$$(A^T)^T = A$$

$$(A+B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

* A matrix A is symmetric $\Leftrightarrow A = A^T$

* Inverse of A : $AA^{-1} = A^{-1}A = I$

* A is invert \Leftrightarrow $\begin{cases} \det A \neq 0 \\ \text{rows/columns of } A \text{ are linearly indep.} \\ \text{the only sol of } Ax = 0 \text{ is } x = 0 \end{cases}$

* A and B are invertible $\Leftrightarrow (AB)^{-1} = B^{-1}A^{-1}$

$$\underline{(A^T)^{-1} = (A^{-1})^T}$$

* Def 2.5

App - real valued X_1, \dots, X_p are mutually independent $F_{\underline{X}}(\underline{x}) = \prod_{i=1}^p F_{X_i}(x_i), \forall \underline{x} \in \mathbb{R}^p$

* Def

The two random vectors \underline{X} and \underline{Y} are independent $\Leftrightarrow F_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) = F_{\underline{X}}(\underline{x}) F_{\underline{Y}}(\underline{y})$

* \underline{X} and \underline{Y} are independent

$$\underline{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \quad \underline{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_q \end{bmatrix}$$

(X_1, \dots, X_p) don't need to be independent

* Def (Population correlation matrix)

Let variance-covariance matrix $\Sigma_{\underline{X}} = V(\underline{X}) = C(\underline{X}, \underline{X}) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix}$ where $\underline{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix}$

Then consider population correlation matrix

$$\rho_{\underline{X}} = \begin{bmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \rho_{21} & 1 & \dots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \dots & 1 \end{bmatrix}_{p \times p}$$

a symmetric matrix

$$\rho_{ij} = \text{cor}(X_i, X_j) = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}} \sqrt{\sigma_{jj}}}$$

Let V : standard deviation matrix

$$V^{1/2} = \begin{bmatrix} \sqrt{\sigma_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{\sigma_{22}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\sigma_{pp}} \end{bmatrix}_{p \times p} = [\text{diag}(\sqrt{\sigma_{ii}})]$$

Then

$$\Sigma_{\underline{X}} = V^{1/2} \rho_{\underline{X}} V^{1/2}$$

$$\rho_{\underline{X}} = V^{-1/2} \Sigma_{\underline{X}} V^{-1/2}$$

* The mean vector + covariance matrix for linear combinations of random vectors

* Theorem 2.3

Let $\underline{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix}$ then $\underline{c}^T \underline{X} = [c_1 \ c_2 \ \dots \ c_p] \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = \sum_{i=1}^p c_i X_i$ is a linear combination of p random variables X_1, \dots, X_p

number vector \underline{c}
random vector \underline{X}

Then we have a) $E(\underline{c}^T \underline{X}) = \underline{c}^T E(\underline{X})$

b) $V(\underline{c}^T \underline{X}) = \underline{c}^T V(\underline{X}) \underline{c}$

* Theorem 2.3

Let $\underline{C} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix}_{p \times p}$

$\underline{d} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_p \end{bmatrix}$

Then $\underline{Z} = \underline{C}\underline{X} + \underline{d} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} + \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_p \end{bmatrix} = \begin{bmatrix} Z_{11} \\ Z_{21} \\ \vdots \\ Z_{p1} \end{bmatrix}$
 $Z_i = (c_{i1} X_1 + c_{i2} X_2 + \dots + c_{ip} X_p) + d_i = \sum_{j=1}^p c_{ij} X_j + d_i$

Then we have

$E(\underline{C}\underline{X} + \underline{d}) = \underline{C}E(\underline{X}) + \underline{d}$

$V(\underline{C}\underline{X}) = \underline{C}V(\underline{X})\underline{C}^T$

$V(\underline{C}\underline{X} + \underline{d}) = \underline{C}V(\underline{X})\underline{C}^T$

* Theorem:

$A_{p \times p}$: symmetric matrix
 $\underline{X}_{p \times 1} \propto (\underline{X}_1, \underline{Z}_n)$
 $E(\underline{X}^T A \underline{X}) = \underline{1}^T (A \Sigma_{\underline{X}}) \underline{1} + \underline{1}^T A \underline{1}$
 Quadratic form of \underline{X}

* Chapter 3: Sample geometry and Random sampling

* 3.1 Introduction

* 3.2 The geometry of the sample

- A single multivariate observation is a collection of measurements on p different variables taken on a trial.

- If we have n observations for those p variables, the entire data set is

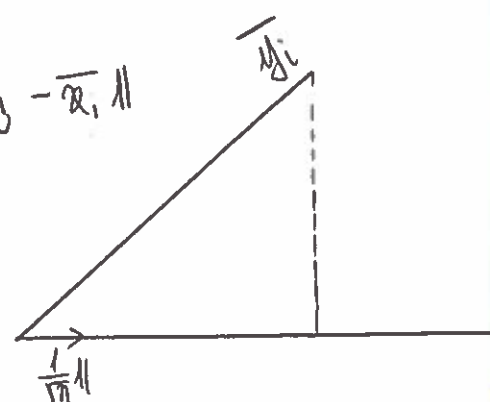
$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}_{n \times p} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = [Y_1 | Y_2 | \dots | Y_p] \quad \bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

where $x_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix}_{p \times 1}$, $i = 1, \dots, n$ (people)
 observed from the person i

where $Y_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix}_{n \times 1}$, $i = 1, \dots, p$
 observe for variable i by n people

* Def 3.1: For each $Y_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix}$, we define deviation vector $\underline{d}_i = Y_i - \bar{x}_i \mathbf{1}$

- Projection of Y_i to $\frac{1}{n} \mathbf{1}$ is $\bar{x}_i \mathbf{1}$, where $\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$.





4. Multivariate normal distribution

* Recall univariate normal dist $X \sim N(\mu, \sigma^2)$ $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$

$$\frac{(x-\mu)^2}{\sigma^2} = (x-\mu)(\sigma^2)^{-1}(x-\mu)^T$$

* Multivariate normal distribution

$\tilde{X}_{p \times 1} \sim N_p(\mu_{p \times 1}, \Sigma_{p \times p})$ where Σ p.d.
 Σ p.s.d. (in general case)

Pro pdf

$$f_{\tilde{X}}(\tilde{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\tilde{x}-\mu)^T \Sigma^{-1}(\tilde{x}-\mu)\right\}, \text{ where } \tilde{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$$

is a function in \mathbb{R}^p , depends on (x_1, \dots, x_p)

* Example 4.1 The bivariate normal ($p=2$)

$\tilde{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right)$ X_1 and X_2 are independent $\Leftrightarrow \rho = 0$

$$f_{\tilde{X}}(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right)\right]\right\}$$

* Theorem 4.1

Let $\tilde{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_p(\mu, \Sigma)$

then X_1 and X_2 are independent $\Leftrightarrow \Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix}$

* Theorem 4.2

Constant probability density contours = {all \tilde{x} such that $(\tilde{x}-\mu)^T \Sigma^{-1}(\tilde{x}-\mu) = c^2$ }

= surface of an ellipsoid centered at μ

(where (λ_i, e_i) are (eigenvalues, eigen vectors) of Σ having axes $\pm c\sqrt{\lambda_i} e_i$)

* Result 4.1/53

Σ positive definite $\Leftrightarrow \Sigma^{-1}$ positive definite

(λ_i, e_i) are (eigenvalue, eigen vector) of $\Sigma \Rightarrow \left(\frac{1}{\lambda_i}, e_i\right)$ are (eigenvalue, eigen vector) of Σ^{-1}

* Theorem 4.37

$\tilde{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \begin{matrix} P_1 \\ P_2 \end{matrix} \sim N_{P_1+P_2}(\mu, \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix})$

then $X_1 \sim N_{P_1}(\mu_{11}, \Sigma_{11})$

* Additional properties of multivariate normal distribution

Let $X \sim N_p(\mu, \Sigma)$, $\Sigma > 0$, then

$\Rightarrow X + c \sim N_p(\mu + c, \Sigma)$

$\Rightarrow AX \sim N_q(AM, A \Sigma A^T)$ $a^T X \sim N(a^T \mu, a^T \Sigma a)$

$\Rightarrow I_p \Sigma^{-1} > 0$ then $Y = \Sigma^{-1/2}(X - \mu) \sim N_p(0, I)$

$\Sigma = \sigma^2 I$ (diagonal matrix with elements σ^2) $\Rightarrow GX \sim N(G\mu, \sigma^2 I)$
 G is orthogonal ($G^T G = I$)
 G is $q \times p$ ($q < p$)

* Alternative multivariate normal distribution (MVN)

Let $Z \sim N(0, I)$ $Z = \begin{pmatrix} z_1 \\ \vdots \\ z_p \end{pmatrix}$ where $z_i \sim N(0, 1)$.

then $X = \mu + LZ \sim N(\mu, LL^T)$

* Def (for Chi Square)

Let $Z \sim N(0, I)$ $Z = \begin{bmatrix} z_1 \\ \vdots \\ z_p \end{bmatrix} \sim N(0, I)$ where $z_i \sim N(0, 1)$
 then $\chi_p^2 = Z^T Z = \sum_{i=1}^p z_i^2 \sim \chi_p^2$

* Theorem

$X \sim N_p(\mu, \Sigma_{p \times p})$ (Σ p.d.)
 then $(X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi_p^2$

When $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$
 $\Rightarrow n(\bar{X} - \mu)^T \Sigma^{-1} (\bar{X} - \mu) \sim \chi_p^2$

* Theorem
 $X_p \sim N_p(\mu, \Sigma_{p \times p}) \Sigma > 0 \Rightarrow P\left\{ (X - \mu)^T \Sigma^{-1} (X - \mu) \leq \chi_{p, 1-\alpha}^2 \right\} = 1 - \alpha$

* Theorem 4.67 (Conditional distribution)

$X \sim N_p(\mu, \Sigma)$ ($\Sigma > 0$)

$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ $P_1 + P_2 = p$

$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ $\Sigma_{21} = \Sigma_{12}^T$

Then:
 $X_1 | X_2 = x_2 \sim N_{P_1}\left(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}\right)$

4.3/ p168 * Maximum likelihood estimation of μ, Σ

* Def

Let $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$, $(\Sigma > 0)$

The likelihood is

$$L(\mu, \Sigma | X_1, \dots, X_n) = \prod_{j=1}^n \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)} \right)$$

$$= \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\sum_{j=1}^n \frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)}$$

* Rem 4.9/p 168

A $n \times n$ symmetric real sector
 then $x^T A x = \frac{1}{2} (x^T A x + x^T A x) = \frac{1}{2} (x^T A x + x^T A^T x) = \frac{1}{2} (x^T (A + A^T) x)$

* Theorem

$$\sum_{j=1}^n \frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) = \frac{1}{2} \text{tr} \left[\Sigma^{-1} \left(\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T \right) \right] + n \left[(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu) \right]$$

* Def 4.2 (Sufficiency)

$T(X)$ is sufficient for θ if $P(X=x | T(X)=t, \theta) = P(X=x | T(X)=t)$

* Theorem 4.8 (Factorization theorem for sufficiency) free of θ

T is sufficient for $\theta \iff f_{\theta}(x) = h(x) g_{\theta}(T(x))$

* Theorem 4.9

(S, \bar{X}) are jointly sufficient for (μ, Σ)

$$S = \frac{1}{(n-1)} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$$

* Theorem 4.10

Let $B > 0$

$b > 0$ be a scalar

$$\frac{1}{|B|} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} B)} \leq \frac{1}{|B|} (2b)^{bp} e^{-bp} \text{ for all } \Sigma > 0$$

equality $\iff \Sigma = \frac{1}{2b} B$

* Theorem 4.11

Let $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$

The MLE of (μ, Σ) is

$$\hat{\mu}_{MLE} = \bar{X}$$

$$\hat{\Sigma}_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$$

independent

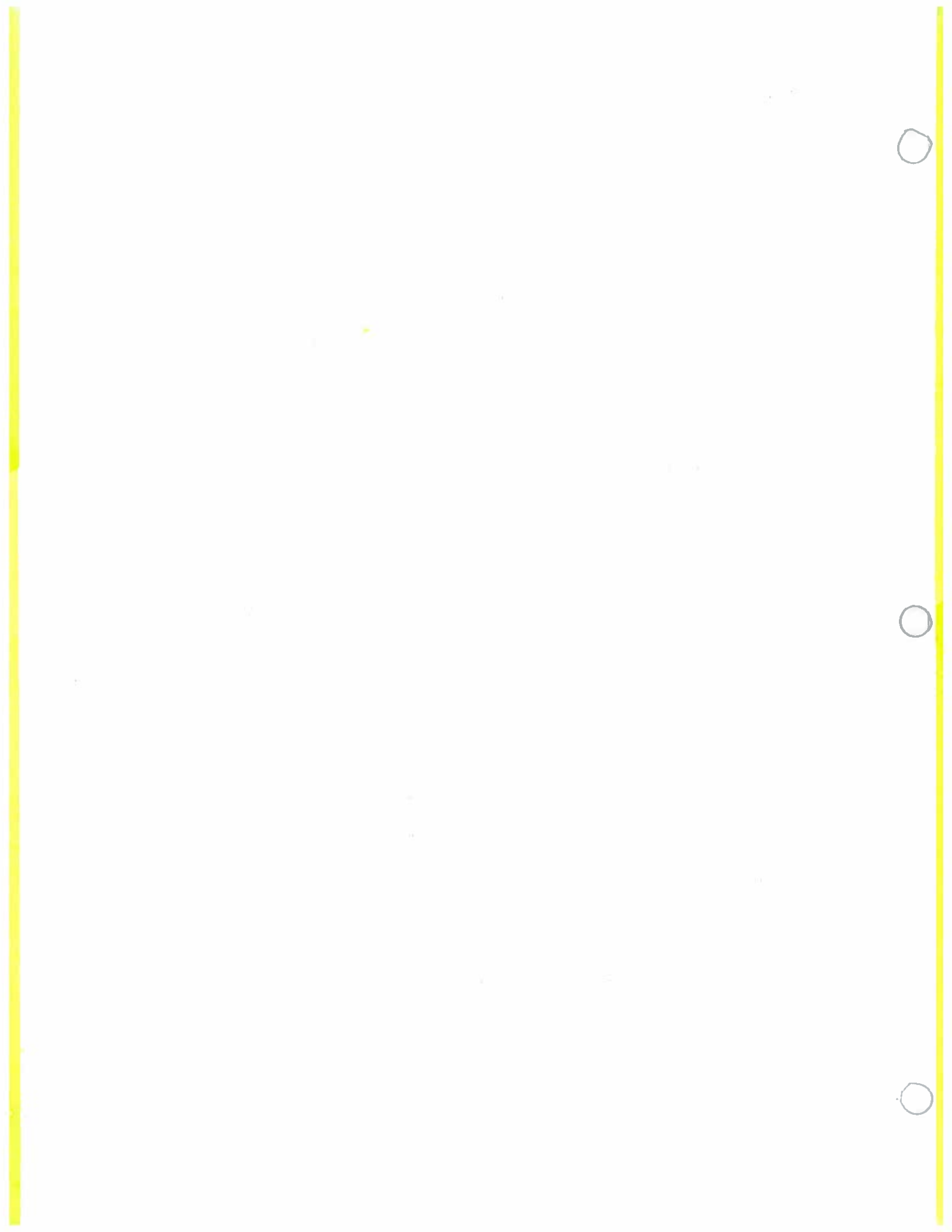
* Theorem 4.12

Let A symmetric $p \times p$

$B_{n \times p}$

$X \sim N_p(\mu, \Sigma)$

Then $B \hat{X}$ and $X^T A X$ are independent $\iff B \Sigma A = 0_{q \times p}$



$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 2 & 4 & 2 & 6 \\ 1 & 2 & 3 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 2 & 3 & 4 \\ 0 & 6 & 7 & 8 \\ 4 & 4 & 2 & 6 \\ 2 & 2 & 3 & 1 \end{bmatrix} \quad B_1 = B_0 C$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2 & 3 & 4 \\ 3 & 6 & 7 & 8 \\ 4 & 4 & 2 & 6 \\ 2 & 2 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 2 & 3 & 4 \\ 3 & 6 & 7 & 8 \\ 2 & 2 & 1 & 3 \\ 2 & 2 & 1 & 3 \end{bmatrix} \quad B_2 = D B_1 = D B_0 C$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2 & 3 & 4 \\ 3 & 6 & 7 & 8 \\ 2 & 2 & 1 & 3 \\ 2 & 2 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 4 & 4 & 4 & 7 \\ 3 & 6 & 7 & 8 \\ 2 & 2 & 1 & 3 \\ 2 & 2 & 1 & 3 \end{bmatrix} \quad B_3$$

$$\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & 4 & 4 & 7 \\ 3 & 6 & 7 & 8 \\ 2 & 2 & 1 & 3 \\ 2 & 2 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 8 & 6 & 7 & 3 \\ 4 & 4 & 4 & 7 \\ 3 & 6 & 7 & 8 \\ 3 & 6 & 7 & 8 \end{bmatrix} \quad B_4 = B_3 D$$

subtract row 2 from each of other row.

$$\begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 8 & 6 & 7 & 3 \\ 4 & 4 & 4 & 7 \\ 3 & 6 & 7 & 8 \\ 3 & 6 & 7 & 8 \end{bmatrix} = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 8 & 6 & 7 & 3 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad B_5$$

$$\begin{aligned} \text{row 1}(B_5) &= \text{row 1}(B_4) + (-1)\text{row 2}(B_4) \\ \text{row 2}(B_5) &= \text{row 2}(B_4) \\ \text{row 3}(B_5) &= \text{row 3}(B_4) - (-1)\text{row 2}(B_4) \end{aligned}$$

work with column: after A
works with rows: P A
then P a row / col remaining

Time	Monday Jan 15	Tuesday Jan 16	Wednesday Jan 17	Thursday Jan 18	Friday Jan 19	Saturday Jan 20	Sunday Jan 21
8:00AM		MAT 682 - M001 Section 8:00AM - 9:20AM Carnegie 120		MAT 682 - M001 Section 8:00AM - 9:20AM Carnegie 120	8:25 - 9:20 Can 124		
9:00AM		MAT 682 - M001 Section 8:00AM - 9:20AM Carnegie 120		MAT 682 - M001 Section 8:00AM - 9:20AM Carnegie 120	9:30 - 10:25 Can 124		
10:00AM		MAT 652 - M001 Section 9:30AM - 10:50AM Carnegie 122		MAT 652 - M001 Section 9:30AM - 10:50AM Carnegie 122	ENL 640 - M002 Section 10:35AM - 11:30AM MARSHALL SQUARE MALL 208B		
11:00AM			11am - 12pm Math Clinic <i>Office R</i>				
12:00PM		MAT 695 - M001 Section 12:30PM - 1:50PM Carnegie 100		MAT 695 - M001 Section 12:30PM - 1:50PM Carnegie 100	12:45 - 1:40 Can 124		
1:00PM					2:15 - 3:10 Can 124		
2:00PM					<i>Office R</i>		
3:00PM		MAT 526 - M001 Section 3:30PM - 4:50PM Carnegie 115		MAT 526 - M001 Section 3:30PM - 4:50PM Carnegie 115			
4:00PM							
5:00PM							
6:00PM							
7:00PM							
8:00PM							
9:00PM							

* Eigenvalues, eigenvectors and eigendecomposition.

Eigendecomposition
 why's useful?
 n inverse
 power
 Properties
 How to compute?
 Power iteration
 QR algorithm

* Eigenvalue, eigenvector

$v \neq 0, Av = \lambda v \Rightarrow \lambda$ eigenvalue
 v eigenvector

$Av = \lambda v \Leftrightarrow (A - \lambda I)v = 0 \Rightarrow v \neq 0$ is a solution of $[A - \lambda I]v = 0$

* Why is eigendecomposition is useful?

Assume A has k eigenvalues $\lambda_1, \dots, \lambda_k$ with k associated eigenvectors v_1, v_2, \dots, v_k

Then if $V = [v_1 \ v_2 \ \dots \ v_k]$ then $AV = V\Sigma$

$$\Sigma = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_k \end{bmatrix}$$

$A = V\Sigma V^{-1}$

$AV = V\Sigma$

Σ : triangular matrix

contains all eigenvalue of A .

V contains all eigenvectors of A

We say V and A are similar

\Rightarrow This can only be done if a matrix is diagonalizable.

\Rightarrow Def of a diagonalizable matrix

A is a diagonalizable matrix if we can eigendecompose it into n eigenvectors.
 $A \in \mathbb{C}^{n \times n}$

Matrix inverse with eigendecomposition

$A = V\Sigma V^{-1} \Rightarrow A^{-1} = V\Sigma^{-1}V^{-1}$

Power of a matrix with eigendecomposition

$A^n = V\Sigma^n V^{-1}$ for all integers n

* Remark: If a matrix A is a diagonal triangular matrix

\rightarrow all eigenvalue are on the diagonal.

* Properties of eigendecomposition

$\det A = \prod_{i=1}^n \lambda_i$

trace $A = \sum_{i=1}^n \lambda_i$

eigenvalues of A^{-1} are λ_i^{-1}

A^n are λ_i^n

The eigenvectors of $A^{-1} =$ The eigenvectors of A .

$\left. \begin{matrix} A \text{ is Hermitian} \\ A \text{ is full rank} \end{matrix} \right\} \Rightarrow \left\{ \begin{matrix} \text{eigenvectors are mutually orthogonal} \\ \text{eigenvalues are real} \end{matrix} \right.$
 (all row/columns are linearly independent)

A is invertible $\Leftrightarrow \lambda_i \neq 0, \forall i$

If the eigenvalues of A are distinct $\Rightarrow A$ can be eigendecomposed

Theorem 2.5

$A \in \mathbb{C}^{m \times m}$ Hermitian (symmetric if $A \in \mathbb{R}$)
 \rightarrow all eigenvalues are real
 \rightarrow all eigenvectors corresponding to "distinct" are orthogonal.



per
N...

J

B



§ 1.7 Aspects of multivariate analysis. 1.57 The organization of Data.

$p \gg 1$: # of variable/characters to record. (The values of those variables are all recorded for each distinct item/individual/unit.)

x_{jR} ← R^{th} variable

j^{th} item/trial

	Item	Variable 1	Variable 2	...	Variable p
$\left\{ \begin{array}{l} n \text{ items (people)} \\ p \text{ variables} \end{array} \right.$	Item 1				
	⋮				
	Item n				

$\Rightarrow [X]_{n \times p}$ n : # observations (people/items)
 p : # of variables

* Descriptive Statistics

+ The arithmetic average $\bar{x}_1 = \frac{1}{n} \sum_{j=1}^n x_{j1}$ ← sample mean for the first variable (sample mean)

$\bar{x}_R = \frac{1}{n} \sum_{j=1}^n x_{jR}$ (works with R^{th} column $\leftrightarrow R^{\text{th}}$ variable)

+ sample variance

$$s_{RR}^2 = \frac{1}{n} \sum_{j=1}^n (x_{jR} - \bar{x}_R)^2$$

each variable may have different units

+ sample standard deviation

$$\sqrt{s_{RR}^2}$$

+ sample covariance

$$s_{12} = \frac{1}{n} \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{j2} - \bar{x}_2)$$

$s_{12} > 0$ large \leftrightarrow large
 $s_{12} > 0$ large \leftrightarrow small
 $s_{12} \approx 0$ there is no particular association

$s_{12} = \frac{1}{n} \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{j2} - \bar{x}_2)$: A measure of linear association between the measurement of variable i^{th} and variable k .

+ Sample correlation coefficient

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11}} \sqrt{s_{22}}} = \frac{\sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{j2} - \bar{x}_2)}{\sqrt{\sum_{j=1}^n (x_{j1} - \bar{x}_1)^2} \sqrt{\sum_{j=1}^n (x_{j2} - \bar{x}_2)^2}}$$

• is a standardized version of sample covariance
 is the same whether n or $(n-1)$, is chosen as the common divisor for
 can be considered sample covariance, when we replace x_{jR} by $\frac{x_{jR} - \bar{x}_R}{\sqrt{s_{RR}}}$
 of standardized observations x_{jR} by $\frac{x_{jR} - \bar{x}_R}{\sqrt{s_{RR}}}$

① $-1 \leq r \leq 1$

② r measures the strength of the linear association.

$r < 0 \Rightarrow$ tendency for one value in the pair $<$ its average

Other value is smaller than its average.

$r > 0 \Rightarrow$ both value in one pair bigger or smaller than its average

③ $\bar{x}_{i\cdot}$ is unchanged when x_{ji} changed to $y_{ji} = a x_{ji} + b, j=1, n$
 1st variables \rightarrow x_{ji} changed to $y_{ji} = a x_{ji} + b, j=1, n$
 \rightarrow a, c have the same sign
 \rightarrow $y_{j\cdot} = c x_{j\cdot} + d, j=1, n$
 2nd variables

* The sum of squares of the deviations from the mean

$$W_{i\cdot} = \sum_{j=1}^n (x_{ji} - \bar{x}_{i\cdot})^2$$

* The sum of cross product deviation

$$W_{i\cdot} = \sum_{j=1}^n (x_{ji} - \bar{x}_{i\cdot}) (x_{j\cdot} - \bar{x}_{\cdot})$$

* Arrays of Basic Descriptive Statistics.

• Sample means

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{bmatrix}_{p \times 1}$$

• Sample variances and covariances

use to remind that n items are samples for elements.

$$S_{ij} = \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ s_{21} & s_{22} & \dots & s_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ s_{p1} & s_{p2} & \dots & s_{pp} \end{bmatrix}_{p \times p}$$

← symmetric matrix

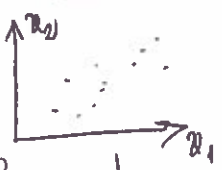
• Sample correlation

$$R_{ij} = \begin{bmatrix} 1 & r_{12} & \dots & r_{1p} \\ r_{21} & 1 & & \\ \vdots & & \ddots & \\ r_{p1} & r_{p2} & & 1 \end{bmatrix}_{p \times p}$$

← symmetric matrix

* Graphical Techniques.

• Plot pairs of variables is quite informative.



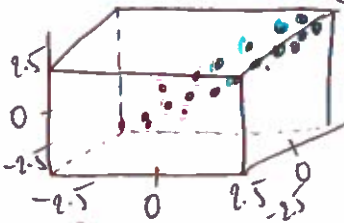
can help calculate $\bar{x}_1, \bar{x}_2, s_{11}, s_{12}, s_{22}$

or p points in n-dimensional space.

• The effect of unusual observation on sample correlations

• Graphic of n points in p dimensions. $(x_{j1}, x_{j2}, x_{j3}, \dots, x_{jp})$ for $j=1, n$

• Example 1.6: (looking for lower dimensional structure)



• some variables may have a much larger variance than the others

• We want to use standardized variables $Z_{j\cdot} = \frac{x_{j\cdot} - \bar{x}_{\cdot}}{\sqrt{s_{\cdot\cdot}}}$

• Example 1.7 (looking for group structure in three dimensions) (p18)

• female

• male

groups of female and male.

1.47 Data displays and pictorial representations

Powerful computers → help have goal of data analysis graphics

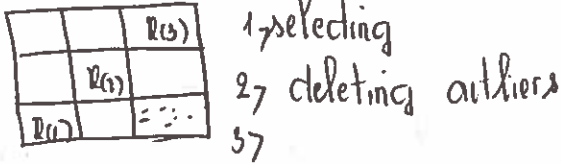
PC (principal component analysis)
 CA (canonical correlation analysis)

represent p -dimensional observations in few (e) dimensions s.t. the original distances (or similarities) between pair of observations are (nearly) preserved.

- 2-dimensional scatter plot
- graphs of growth curve
- stars
- chessboard game.

* Linking multiple 2-dimensional scatter plots.

* Example 1.8 (Linked scatter plots and brushing)



→ brushing: operation of highlighting points corresponding to a selected range of one of the variables

* Example 1.9 (Rotate plots in three dimensions)

Rotating and turning the 3-dimensional coordinate axes can help getting a better idea understanding of the three-dimensional aspects of the data. → to confirm what specimens are outliers



* Graphs of growth curves

Growth curve: the points are plotted and then connected by lines to produce a graph.

(height of a young child, for ex) is measure at each birth day)

* Example 1.10 (Arrays of growth curve)

for example a growth curve for weight of some female bears by year (each line for each bear)



we can also have → individual growth curves for weight for $n \neq$ bears



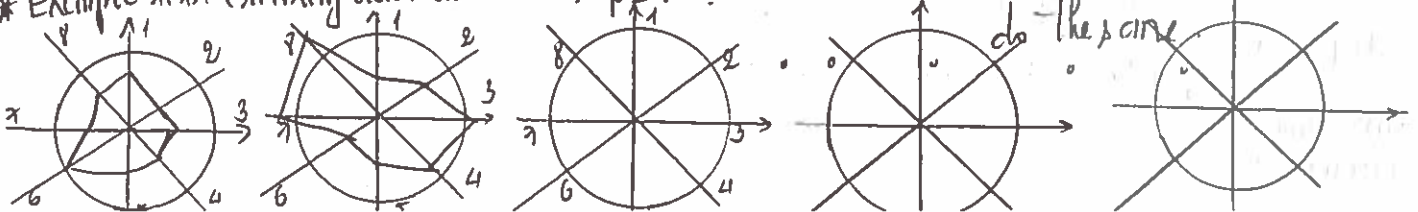
* Stars: each data unit consists of nonnegative observation on $p \geq 2$ variables

→ In two dimensions, → construct circles with a fixed radius.

• The observations on all variables were standardized.

has p equally spaced rays
 the length of the rays represent the values of the variable

* Example 1.11 (utility data as stars) $p=7$



* Chernoff Faces.

* 1.5 Distance.

* Euclidean distance.

$$P(x_1, \dots, x_p) \rightarrow Q(y_1, \dots, y_p) \rightarrow d(P, Q) := \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

* Statistical distance

* $P(x_1, x_2)$
 $Q(y_1, y_2)$

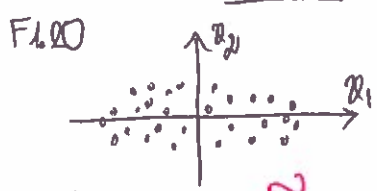
$$d(P, Q) = \sqrt{\frac{(x_1 - y_1)^2}{s_1^2} + \frac{(x_2 - y_2)^2}{s_2^2}} \quad (1.13)$$

* $P(x_1, \dots, x_p)$
 $Q(y_1, \dots, y_p)$

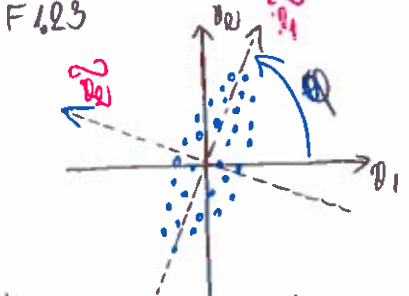
$$d(P, Q) = \sqrt{\sum_{i=1}^p \frac{(x_i - y_i)^2}{s_i^2}} = \sqrt{\frac{(x_1 - y_1)^2}{s_1^2} + \frac{(x_2 - y_2)^2}{s_2^2} + \dots + \frac{(x_p - y_p)^2}{s_p^2}}$$

If $s_1^2 = \dots = s_p^2$
 \Rightarrow the Euclidean distance formula is appropriate

The above distance does not include the most important cases we shall encounter,
 \rightarrow (because : assumption of independent coordinates.



when x_1 measurements do not vary independently of the x_2 measurements
 $(x_1, x_2) \leftarrow$ has tendency to be large or small together.



* So now we want to find a distance (statistical distance that is meaningful) when
 variability in the x_1 direction is different from the variability in the x_2 direction.
 the variables x_1 and x_2 are correlated
 (If we look at things in the right way)

If we axes $(\tilde{x}_1, \tilde{x}_2)$ then with $P(\tilde{x}_1, \tilde{x}_2)$
 coordinate system $O(0,0)$

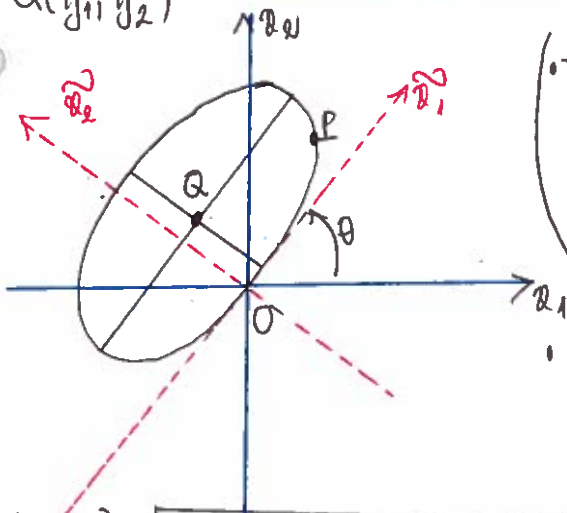
$$d(P, O) = \sqrt{\frac{\tilde{x}_1^2}{s_1^2} + \frac{\tilde{x}_2^2}{s_2^2}} \quad (1.17)$$

where $\begin{cases} \tilde{x}_1 = x_1 \cos \varphi + x_2 \sin \varphi \\ \tilde{x}_2 = -x_1 \sin \varphi + x_2 \cos \varphi \end{cases}$

Then $d(P, O) = \sqrt{a_{11} x_1^2 + 2a_{12} x_1 x_2 + a_{22} x_2^2}$
 in original coordinate

* In general

$$d(P, Q) = \sqrt{a_{11}(x_1 - y_1)^2 + 2a_{12}(x_1 - y_1)(x_2 - y_2) + a_{22}(x_2 - y_2)^2} \quad (1.10)$$



The coordinates of all points $P(x_1, x_2)$, that are a constant square c^2 from Q satisfy

$$a_{11}(x_1 - y_1)^2 + 2a_{12}(x_1 - y_1)(x_2 - y_2) + a_{22}(x_2 - y_2)^2 = c^2 \quad (1.11)$$

(ellipse center Q)

• Let $(O) = (0, 0, \dots, 0)$

$Q = (y_1, \dots, y_p) \leftarrow$ fixed point

$$d(O, P) = \sqrt{a_{11}x_1^2 + a_{22}x_2^2 + \dots + a_{pp}x_p^2 + 2a_{12}x_1x_2 + 2a_{13}x_1x_3 + \dots + 2a_{p-1,p}x_{p-1}x_p}$$

$$d(P, Q) = \sqrt{a_{11}(x_1 - y_1)^2 + a_{22}(x_2 - y_2)^2 + \dots + a_{pp}(x_p - y_p)^2 + 2a_{12}(x_1 - y_1)(x_2 - y_2) + 2a_{13}(x_1 - y_1)(x_3 - y_3) + \dots + 2a_{p-1,p}(x_{p-1} - y_{p-1})(x_p - y_p)}$$

where a_{ij} are the numbers, such that the distances are always nonnegative.

\Rightarrow The distances are completely determined by coefficients (weights) a_{ij} $i = 1, \dots, p$
 $j = 1, \dots, p$



* Review descriptive statistics.

before analyzing data

1st step: check if data is correct ex 0 1 2 gender is not correct

2nd step: makes sense? → incl out when correct

The sample mean of the k^{th} variable

$$\bar{x}_k = \frac{1}{n} \sum_{j=1}^n x_{jk}$$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

Sample variance - covariance.

$$s_{ii}^2 = s_{x_{ij}}^2 = \frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)^2$$

$$s_{ij} = \frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jp} - \bar{x}_j)$$

$$S_{p \times p} = \begin{bmatrix} s_{11} & s_{12} & s_{13} & \dots & s_{1p} \\ s_{21} & s_{22} & s_{23} & \dots & s_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ s_{p1} & s_{p2} & \dots & \dots & s_{pp} \end{bmatrix}$$

is a $p \times p$ symmetric matrix correlation between variables.

Sample correlation coefficient.

$$\lambda_{ik} = \frac{s_{ik}}{\sqrt{s_{ii} s_{kk}}} \quad -1 \leq \lambda_{ik} \leq 1$$

$$R_{p \times p} = \begin{bmatrix} 1 & \lambda_{12} & \dots & \lambda_{1p} \\ \lambda_{21} & 1 & \dots & \lambda_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{p1} & \lambda_{p2} & \dots & 1 \end{bmatrix}$$

* Distance and metric

Let \mathcal{S} be a simple space.

Distance function $d: \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ satisfies

$$\begin{cases} \forall x, y, z \in \mathcal{S} \\ d(x, y) \geq 0, \forall x, y & d(x, y) = 0 \Leftrightarrow x = y \\ d(x, y) = d(y, x) & \text{symmetry} \\ d(x, y) \leq d(x, z) + d(z, y) & \Delta \text{ inequality} \end{cases}$$

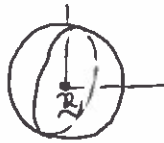
nonneg
coincide

* Common distance in \mathbb{R}^p .

Euclidean

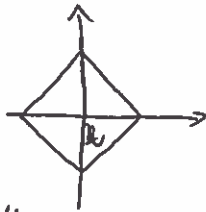
$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$$

all points on this curve equidistant from x



• City block

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

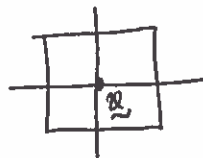


• Minkowski

$$d_\lambda(x, y) = \left[\sum_{i=1}^p |x_i - y_i|^\lambda \right]^{1/\lambda}, \lambda \geq 1$$

• Manhattan:

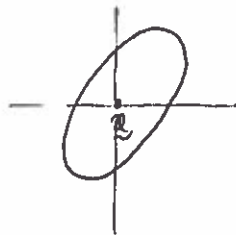
$$d(x, y) = \lim_{\lambda \rightarrow \infty} d_\lambda(x, y) = \max_i |x_i - y_i|$$



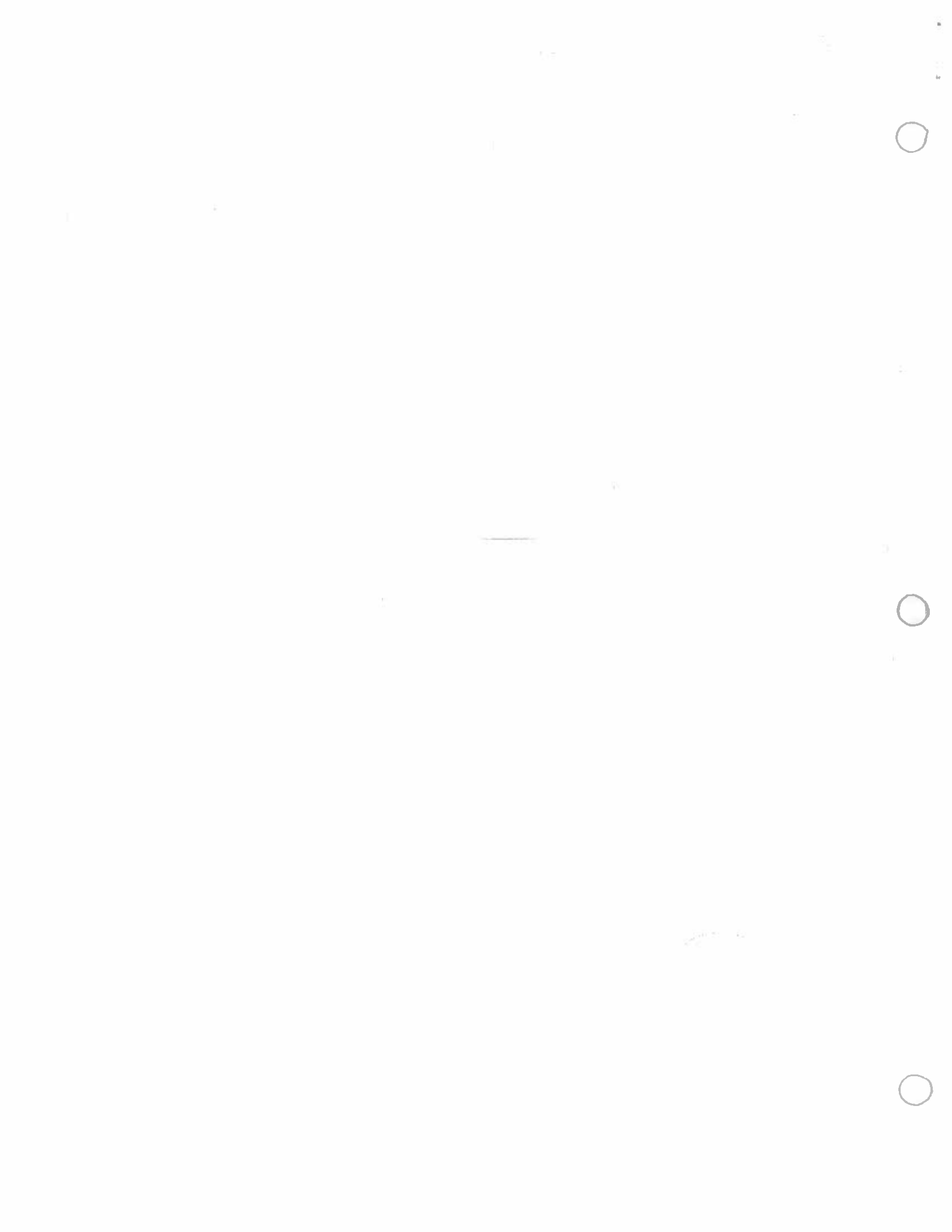
• Mahalanobis distance

$$d(x, y) = \sqrt{(x - y)_{1 \times p}^T S^{-1} (x - y)_{p \times 1}}$$

S: variance-covariance metric



← includes the structure of data



27 Algebraic review (focus on Real, symmetric, positive)

* Definition: Let $A = (A_{ij})_p$

Let $(A)_{ij} = (a_{ij})$ then transpose matrix of A : $A^T = (a_{ji})$

$$\begin{aligned}(A^T)^T &= A \\ (A+B)^T &= A^T + B^T \\ (AB)^T &= B^T A^T\end{aligned}$$

* Def 2.2

A matrix A is symmetric if $A^T = A$ all var, cov, cor matrices are symmetric

* Def 2.3

Let A be a matrix and suppose $\exists B$ s.t. $AB = BA = I$

then $A^{-1} = B$ is the inverse of A

A matrix is invertible (non singular)

When $A_{p \times p}$

- \Leftrightarrow
- (1) $\dim A = p =$ the size of matrix.
 - (2) rows / columns of A linearly independent
 - (3) the only sol for $Ax = 0$ is $x = 0$

Note: $x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ $x^T = (x_1, \dots, x_p)$

A and B are invertible $(AB)^{-1} = B^{-1} A^{-1}$

$$(A^T)^{-1} = (A^{-1})^T$$

* Def 2.4

The determinant of a matrix $A_{p \times p}$ satisfies

$$\begin{aligned}|AB| &= |A||B| \\ |A^T| &= |A| \\ |cA| &= c^p |A|\end{aligned}$$

$$|I| = 1$$

$$|A^{-1}| = |A|^{-1}$$

$$\det \begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ a_{p1} & a_{p2} & a_{p3} & \dots & a_{pp} \end{pmatrix} = - \det A$$

$$\det \begin{pmatrix} a_{11} & \dots & a_{1j} + c a_{1k} & \dots & a_{1p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ a_{p1} & \dots & a_{pj} + c a_{pk} & \dots & a_{pp} \end{pmatrix} = c \det A$$

column vector column

* Def 2.4

$A_{p \times p}$ is non singular iff $|A| \neq 0$

singular iff $|A| = 0$

vectors $x, y \Leftrightarrow x^T y = 0$
perpendicular

* Def 2.5

Let A be a non singular matrix satisfying $A^{-1} = A^T \Leftrightarrow A$ is orthogonal

$$AA^T = I$$

* Def 2.6

only applied for square matrices

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & \dots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{bmatrix}$$

$$\text{then } a_i^T a_j = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

The trace of a matrix $A_{p \times p}$ $\text{tr}(A) = \sum_{i=1}^p a_{ii}$

$$\text{We have } \text{tr}(A + cB) = \text{tr}(A) + c \text{tr}(B)$$

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{tr}(AA^T) = \text{tr}(A^T A) = \sum_{i,j=1}^p a_{ij}^2$$

$$\text{tr}(A \pm B) = \text{tr}(A) \pm \text{tr}(B)$$

$$\text{tr}(A_1 A_2 \dots A_n) = \text{tr}(A_{n+1} \dots A_n A_1 \dots A_{n-1})$$

* Def 2.7

Let A be a square matrix $| I_p$ we have
for some $\underline{x} \neq \underline{0}$ $A\underline{x} = \lambda\underline{x}, \lambda \in \mathbb{C}$

then \underline{x} is an eigenvector
 λ associated eigenvalue

* We have

$(A - \lambda I)\underline{x} = \underline{0}$

- i) A is singular $\Leftrightarrow \lambda = 0$ is an eigenvalue
- ii) $\text{rank}(A) = \#$ (nonzero) eigenvalues.
- iii) $\#$ of (independent) eigenvectors \leq size A
- iv) $\sum_{i=1}^p \lambda_i = \text{trace}(A)$ $\prod_{i=1}^p \lambda_i = \det A$

A matrix has size p has p eigenvalues.

* If A symmetric

- i) $\lambda \in \mathbb{R}, i = \overline{1, p}$
- ii) $\#$ linear independent eigenvectors = p
- iii) If $\underline{e}_i = \frac{\underline{x}_i}{\|\underline{x}_i\|}$ then $\langle \underline{e}_i, \underline{e}_j \rangle = \delta_{ij}$, \underline{e}_i and \underline{e}_j are orthogonal and \exists a set of p linear (independent) eigenvectors.

in some cases, we may have repeated eigenvalues.

* Theorem 2.17 (Spectral decomposition theorem for matrices) eigendecomposition

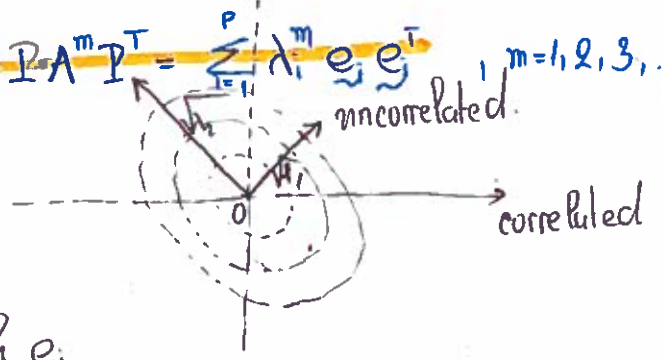
Let A be symmetric then \exists an orthogonal matrix $P = [\underline{e}_1, \dots, \underline{e}_p]$

(i.e. $P^T P = I$)

such that $A = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i^T = P \Lambda P^T = P \Lambda P^T \Rightarrow A^m = P \Lambda^m P^T = \sum_{i=1}^p \lambda_i^m \underline{e}_i \underline{e}_i^T, m=1, 2, 3, \dots$

where $\Lambda = \text{diag}(\lambda_i) = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}$

λ_i : eigenvalue of A associated with \underline{e}_i



positive definite \Rightarrow symmetric non-singular

* Def 2.8

Let A be symmetric, then $\underline{x}^T A \underline{x}$ is a quadratic form.

* Def 2.9

Suppose for all $\underline{x} \neq \underline{0}$ the of $\underline{x}^T A \underline{x} > 0$ then A is positive definite (pd)

If $\underline{x}^T A \underline{x} \geq 0$ then A positive semi-definite (psd)

* Theorem 2.2

positive semi definite $\Leftrightarrow \lambda_i \geq 0, \forall i$

Let A be a non-singular + symmetric then A is positive definite $\Leftrightarrow \lambda_i > 0, \forall i \in \{1, \dots, p\}$

$\{A \underline{x} = 0 \Rightarrow \underline{x}^T A \underline{x} = 0 \Rightarrow$ positive semi def

Actually when A is singular $\rightarrow A$ can not be positive definite

positive definite \Leftrightarrow non-singular symmetric $\lambda_i > 0, \forall i$

Proof

$\underline{x}^T A \underline{x} = \underline{x}^T P \Lambda P^T \underline{x} = (P^T \underline{x})^T \Lambda (P^T \underline{x}) = \underline{y}^T \Lambda \underline{y} = \sum_{i=1}^p \lambda_i y_i^2$ (can choose \underline{y} arbitrary $\underline{x} = P \underline{y}$)

$\rightarrow y_i^2 > 0$ choose $\underline{y} = \begin{pmatrix} 0 \\ \vdots \\ 1 \leftarrow i \\ \vdots \\ 0 \end{pmatrix}$ $\underline{x}^T A \underline{x} > 0$ for all $\underline{x} \neq 0 \Rightarrow \lambda_i > 0$

Conversely if $\lambda_i > 0, i=1, 2, \dots, p$ for $\underline{y} \neq 0$

$0 < \sum_{i=1}^p \lambda_i y_i^2 = \underline{x}^T A \underline{x} \Rightarrow A$ is positive definite

* Theorem 2.3

orthogonal matrix created by eigenvectors of A

diagonal matrix contains all eigenvalues of A

Let A be positive definite $A^{1/2} := P \Lambda^{1/2} P^T = \sum_{i=1}^p \sqrt{\lambda_i} \underline{e}_i \underline{e}_i^T$

well defined since positive def $\Leftrightarrow \lambda_i > 0, \forall i$

then $A^{1/2}$ is the unique p.d. matrix B satisfying $(A^{1/2})^2 = A$

* Theorem 2.4 (Sylvester's Criterion)

$A = \begin{bmatrix} A_m & | \\ \hline & P \end{bmatrix}^p$

Let A is symmetric

A is p.s. \Leftrightarrow all leading $m \times m$ ($m = \overline{1, p}$) matrices $A_{m \times m}$ have $\det A_{m \times m} > 0$

* Example 2.1

Let $R = \begin{bmatrix} 1 & p & p^2 \\ p & 1 & p \\ p^2 & p & 1 \end{bmatrix}$

$A_{(1)} = 1 \quad \det A_{(1)} = 1 > 0$

$A_{(2)} = \begin{pmatrix} 1 & p \\ p & 1 \end{pmatrix} \quad \det A_{(2)} = 1 - p^2 > 0$ when $|p| < 1$

$A_{(3)} = A \quad \det A_{(3)} = (1 - p^2) - p(p - p^3) + p^2(p^2 - p^2)$
 $= 1 - 2p^2 + p^4$
 $= (1 - p^2)^2 > 0, |p| \neq 1$

So by S.C R is p.d. $\Leftrightarrow |p| < 1$

$X_t = p X_{t-1} + Z_t$

* Example:

Consider any random vector $\underline{X} \in \mathbb{R}^n = \begin{bmatrix} X_{11} \\ X_{21} \\ \vdots \\ X_{n1} \end{bmatrix}$

$$\mu = E(\underline{X}) = \begin{bmatrix} E(X_{11}) \\ E(X_{21}) \\ \vdots \\ E(X_{n1}) \end{bmatrix}$$

$$S = E\left((\underline{X} - \mu)(\underline{X} - \mu)^T\right)$$

Then we have S is positive semi definite.

• Proof

$$\begin{aligned} \forall \underline{x} \in \mathbb{R}^n, \underline{x}^T S \underline{x} &= \underline{x}^T E\left((\underline{X} - \mu)(\underline{X} - \mu)^T\right) \underline{x} = E\left(\underline{x}^T (\underline{X} - \mu)(\underline{X} - \mu)^T \underline{x}\right) = \\ &= E\left(\left[\underline{x}^T (\underline{X} - \mu)\right] \left[\underline{x}^T (\underline{X} - \mu)\right]^T\right) = E\left(\left[\underline{x}^T (\underline{X} - \mu)\right]^2\right) \end{aligned}$$

* 25 Random vectors and matrices

* Def:

• A random vector $X_{1 \times p} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$ where $X_i, i=1, p$ are random variables

• A random matrix $X_{n \times p} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$ where $X_{ij}, i=1, n, j=1, p$ are random variables

* The expected of a random matrix

$$E(X) = \begin{bmatrix} E(X_{11}) & E(X_{12}) & \dots & E(X_{1p}) \\ E(X_{21}) & E(X_{22}) & \dots & E(X_{2p}) \\ \vdots & \vdots & \ddots & \vdots \\ E(X_{n1}) & E(X_{n2}) & \dots & E(X_{np}) \end{bmatrix}$$

• For f is a function $E(f(X)) = [E(f(X_{ij}))]_{ij}$

• Let X, Y be random matrices
 A, B be constant matrices
 Then $E(X+Y) = E(X) + E(Y)$
 $E(AX+B) = A E(X) + B$

26 Mean vectors and covariance matrices

* Def

• $X_{p \times 1}$ is a random vector $X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$ where $X_i, i=1, p$ are random variables

• $\mu_x = E(X) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}$
 $\mu_i = E(X_i) = \forall i=1, p$
 ↑ mean vector

• $\Sigma_x = V(X) = \text{cov}(X, X) = E((X - \mu_x)(X - \mu_x)^T) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} & \dots & \sigma_{2p} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} & \dots & \sigma_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \sigma_{p3} & \dots & \sigma_{pp} \end{bmatrix}_{p \times p} = V \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} V^T$
 ↑ covariance matrix
 covariance matrix
 $\sigma_{ii} = V(X_i) = E((X_i - \mu_i)^2)$
 $\sigma_{ij} = \text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$

* Def 2.4

Two real valued r.v.s X, Y are independent $\Leftrightarrow F_{X,Y}(x,y) = F_X(x) F_Y(y) \quad \forall x, y \in \mathbb{R}$

* Theorem 2.2

If X, Y are independent $\Leftrightarrow \text{cov}(X, Y) = 0$
 X and Y are not independent

• $\text{cov}(X, Y) = 0 \not\Rightarrow$ independent

EX	Y	0	1	2
0		0	1/4	0
1		1/4	0	1/4
2		0	1/4	0

* Def 2.5
 p -real valued r.v.s X_1, \dots, X_p are mutually independent $F_X(\underline{x}) = \prod F_{X_i}(x_i) \quad \forall \underline{x} \in \mathbb{R}^p$

* Def
 The two random vectors X and Y are independent $\iff F_{X,Y}(\underline{x}, \underline{y}) = F_X(\underline{x}) \cdot F_Y(\underline{y})$
 X_1, \dots, X_p don't need to be mutually independent

• If X and Y are independent then $\Rightarrow \Sigma_{X,Y} = \begin{bmatrix} \Sigma_X & 0 \\ 0 & \Sigma_Y \end{bmatrix}$

• Def of population correlation matrix
 When $\Sigma_X = V(X) = C(X, X) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix}$

Then we consider population correlation coefficient $\rho_{i,j} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii} \sigma_{jj}}}$

Then the population correlation matrix

$$\rho_X = \begin{bmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \rho_{21} & 1 & \dots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \dots & 1 \end{bmatrix} \quad \text{where } \rho_{i,j} = \text{correlation}(X_i, X_j) = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii} \sigma_{jj}}} = (\mathbf{V}^{1/2})^{-1} \Sigma (\mathbf{V}^{1/2})^{-1}$$

• Let $V_{p \times p}$ standard deviation matrix

$$V = \begin{bmatrix} \sqrt{\sigma_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{\sigma_{22}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\sigma_{pp}} \end{bmatrix} \quad V = [\text{diag}(\sqrt{\sigma_{ii}})]$$

$$\Sigma_X = V^{1/2} \rho_X V^{1/2}$$

$$\rho_X = V^{-1/2} \Sigma_X V^{-1/2}$$

$\frac{1}{2} \frac{d}{dt} \left(\frac{1}{2} m v^2 \right) = \frac{1}{2} m v \frac{dv}{dt}$
 $\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$

$\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$
 $\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$

$\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$
 $\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$

$\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$
 $\frac{1}{2} m v \frac{dv}{dt} = \frac{1}{2} m v \frac{dv}{dt}$

(P75) The mean vector and covariance matrix for linear combinations of r.v.s.

* Theorem 2.3 (Linear combination)

Let $\underline{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix}$ then $\underline{c}^T \underline{X} = [c_1 \ c_2 \ \dots \ c_p] \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = \sum_{i=1}^p c_i X_i$ is a linear combination of random variables X_1, \dots, X_p .

↑
number vector $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$ random vector
Then we have $\Rightarrow E(\underline{c}^T \underline{X}) = \underline{c}^T E(\underline{X})$
 $\Rightarrow V(\underline{c}^T \underline{X}) = \underline{c}^T \Sigma_X \underline{c}$

* Example 2.17

Let $X_1, \dots, X_p \stackrel{iid}{\sim} (\mu, \sigma^2)$

Then we have $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \mathbf{1}^T \underline{X}$ where $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$ $\mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$.

Note that $\underline{\mu} = \mu \mathbf{1} = \begin{bmatrix} \mu \\ \mu \\ \vdots \\ \mu \end{bmatrix}$ $\Sigma_X = \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}$ since $X_1, \dots, X_p \stackrel{iid}{\sim} (\mu, \sigma^2)$

Then we have

$$\bullet E(\bar{X}) = E\left(\frac{1}{n} \mathbf{1}^T \underline{X}\right) = \frac{1}{n} \mathbf{1}^T E(\underline{X}) = \frac{1}{n} \mathbf{1}^T \underline{\mu} = \frac{1}{n} \mathbf{1}^T \mu \mathbf{1} = \frac{1}{n} \mu \mathbf{1}^T \mathbf{1} = \frac{1}{n} \mu n = \mu.$$

$$\bullet V(\bar{X}) = V\left(\frac{1}{n} \mathbf{1}^T \underline{X}\right) = \left(\frac{1}{n} \mathbf{1}^T\right) V(\underline{X}) \left(\frac{1}{n} \mathbf{1}\right) = \frac{1}{n^2} \mathbf{1}^T V(\underline{X}) \mathbf{1} = \frac{1}{n^2} V(\mathbf{1}^T \underline{X})$$

$$= \frac{1}{n^2} \mathbf{1}^T (\sigma^2 \mathbf{I}) \mathbf{1} = \frac{\sigma^2}{n^2} \mathbf{1}^T \mathbf{1} = \frac{\sigma^2}{n^2} \mathbf{1}^T \mathbf{1} = \frac{\sigma^2}{n^2} n = \frac{\sigma^2}{n} \quad \square \text{ example.}$$

* Explain the idea of theorem 2.3 in 2+2 case.

$$\bullet E(a_1 X_1 + a_2 X_2) = a_1 E(X_1) + a_2 E(X_2) \Leftrightarrow E\left([a_1 \ a_2] \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}\right) = [a_1 \ a_2] \begin{bmatrix} E(X_1) \\ E(X_2) \end{bmatrix}$$

$$\bullet V(a_1 X_1 + a_2 X_2) = a_1^2 V(X_1) + a_2^2 V(X_2) + 2a_1 a_2 \text{cov}(X_1, X_2) = a_1^2 \sigma_{11} + a_2^2 \sigma_{22} + 2a_1 a_2 \sigma_{12}$$

$$\Leftrightarrow V\left([a_1 \ a_2] \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}\right) = [a_1 \ a_2] \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

* Theorem 2.3 (General case)

Now consider $C_{p \times p} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix}$ Then let $\underline{Z} = \underline{C} \underline{X} + \underline{d}$

$$\uparrow \text{constant}$$

$$\underline{d}_{k \times 1} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$

$$= \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} + \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$

$$\underline{Z}_1 = c_{11} X_1 + c_{12} X_2 + \dots + c_{1p} X_p = \sum_{i=1}^p c_{1i} X_i$$

$$\text{Then } \begin{cases} \underline{\mu}_Z = E(\underline{C} \underline{X} + \underline{d}) = \underline{C} E(\underline{X}) + \underline{d} \\ V_Z = V(\underline{C} \underline{X} + \underline{d}) = \underline{C} V(\underline{X}) \underline{C}^T \end{cases}$$

* Example 2.2

Let $X_{p \times 1} \sim (\mu_x, \Sigma_x)$

$$Z_{p \times 1} = \Sigma_x^{-1/2} (X - \mu_x) = \frac{X - \mu_x}{\sqrt{\Sigma_x}}$$

then $E(Z_p) = \Sigma_x^{-1/2} E(X - \mu_x) = \Sigma_x^{-1/2} [E(X) - \mu_x] = \Sigma_x^{-1/2} [0] = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

$$V(Z_p) = V(\Sigma_x^{-1/2} (X - \mu_x)) = \Sigma_x^{-1/2} V(X - \mu_x) (\Sigma_x^{-1/2})^T = \Sigma_x^{-1/2} \Sigma_x \Sigma_x^{-1/2} = I$$

↑
symmetric

* Theorem 2.4

Let $A_{p \times p}$ be symmetric constant matrix

$X_{p \times 1} \sim (\mu_x, \Sigma_x)$

Then $E(X^T A X) = \text{tr}(A \Sigma_x) + \mu_x^T A \mu_x$

$$E \left(\begin{bmatrix} x_1 & \dots & x_p \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \right) =$$

* Proof

• Assume $\mu_x = 0$

$$E(X^T A X) = E(\text{tr}(X^T A X)) = E(\text{tr}(A X X^T)) = \text{tr}(E(A X X^T)) = \text{tr}(A E(X X^T)) = E(\text{tr}(X X^T)) = \text{tr}(A \Sigma_x)$$

2.7 p78 Matrix inequalities and maximization.

* Cauchy Schwarz inequality

Let $b_{p \times 1}$ & $d_{p \times 1}$ vectors } Then $(b'd)^2 \leq (b'b)(d'd)$
 equality $\Leftrightarrow b = cd$ for some constant $c \neq 0$.

* Extended Cauchy Schwarz inequality.

Let $b_{p \times 1}$ & $d_{p \times 1}$ } Then $(b'd)^2 \leq (b'Bb)(d'B^{-1}d)$
 $B_{p \times p}$ positive definite matrix } equality $\Leftrightarrow b = cB^{-1}d$ for some constant $c \neq 0$
 or $d = cBb$

* Maximization Lemma

$B_{p \times p}$ positive definite
 $d_{p \times 1}$ given vector
 Then for any arbitrary nonzero vector $x_{p \times 1}$ } $\max_{x \neq 0} \frac{(x'd)^2}{x'Bx} = d'B^{-1}d$
 with maximum attained when $x = cB^{-1}d$ $c \neq 0$

* Theorem 2.5 (Extended Cauchy Schmatz) (p79).

Let \underline{x} and \underline{y} be $p \times 1$ vectors } Then $(\underline{x}^T \underline{y})^2 \leq (\underline{x}^T A \underline{x})(\underline{y}^T A^{-1} \underline{y})$
 $A_{p \times p}$ positive definite

* Theorem 2.6 Maximization Lemma: (p80)

Let A : positive definite } then $\max_{x \neq 0} \frac{(x^T d)^2}{x^T A x} = d^T A^{-1} d$
 d : constant vector
 $x \neq 0$ vector
 c is any constant } the maximum attained when $\underline{x} = c A^{-1} d$, for any constant $c \neq 0$.
 some kind of normalization

* Theorem 2.7 (Maximization of Quadratic Forms for Points on the unit sphere.

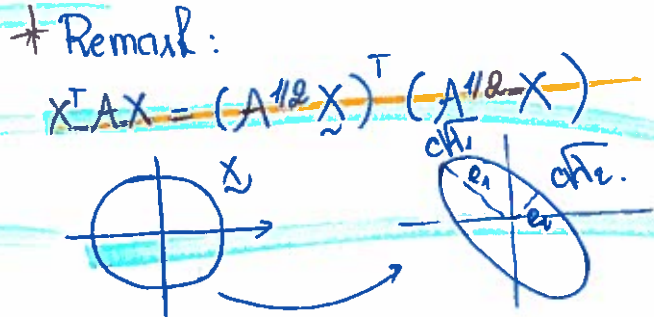
Let A be a positive definite, with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$
 associated eigenvectors $\underline{e}_1, \underline{e}_2, \dots, \underline{e}_p$

Then

$$\max_{x \neq 0} \frac{x^T B x}{x^T x} = \lambda_1$$

$$\min_{x \neq 0} \frac{x^T A x}{x^T x} = \lambda_p$$

$$\max_{x \perp e_1, \dots, e_{p-1}} \frac{x^T A x}{x^T x} = \lambda_{p+1}$$



C3. Sample geometry and random sampling

3.1 Introduction

There are n vector x_i , x_i is a $(p \times 1)$ vector

3.2 The geometry of the sample

* A single multivariate observation v (group of n people, p variables)

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

If obs, the entire data set is

$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

$$= \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$$

\rightarrow p variables observation of the 1st person.

$$= [y_1; y_2; \dots; y_p]$$

\uparrow observation of the i th variable (n people)

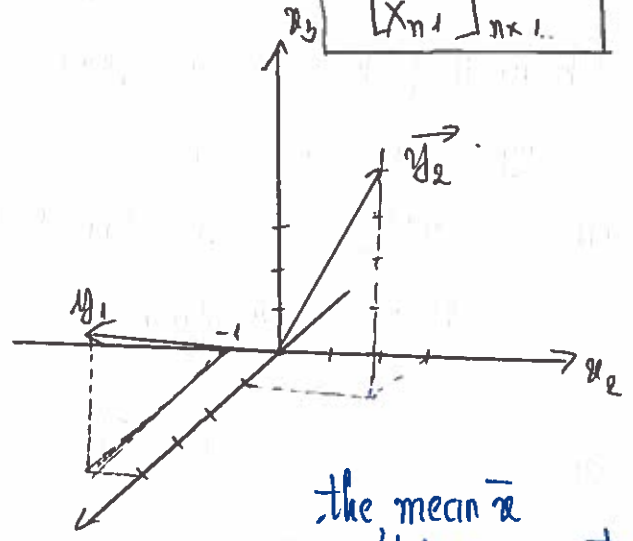
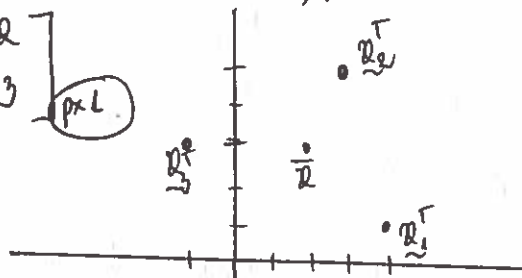
$$x_1 = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1p} \end{bmatrix}_{p \times 1}$$

$$y_1 = \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{bmatrix}_{n \times 1}$$

* Example 3.1 (Compute the mean vector).

$$X = \begin{bmatrix} 4 & 1 \\ -1 & 3 \\ 3 & 5 \end{bmatrix}_{3 \times 2}$$

$$\bar{x} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}_{2 \times 1}$$



if we consider $n=3$ observations (object) in $p=2$ dimensional space.

(then the mean is in the center)

* From the second kind, we can find a geometrical interpretation of the process of finding a sample mean

$$y_1 = \begin{bmatrix} 4 \\ -1 \\ 3 \end{bmatrix}$$

$$y_2 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix}$$

the mean \bar{x} can't be represented in this case

Then consider $u = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1}$ then $(\frac{1}{\sqrt{n}})u$ has length 1.

consider the projection of y_i on the unit vector $(\frac{1}{\sqrt{n}})u$ is z_i

Remind: the projection of x on y is $\frac{x^T y}{y^T y} y$

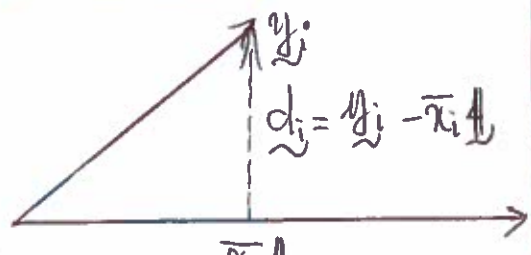
$$\frac{y_i^T (\frac{1}{\sqrt{n}})u}{\text{length}(\frac{1}{\sqrt{n}})u} \frac{1}{\sqrt{n}} u = \frac{y_i^T}{\sqrt{n}} \frac{1}{\sqrt{n}} u \frac{1}{\sqrt{n}} u = \frac{1}{n} (y_i^T u) u = \frac{1}{n} (x_{i1} + \dots + x_{in}) = \bar{x}_i$$

$$y_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix}$$

* Def 3.17

For each $y_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{pi} \end{bmatrix}$, we define deviation vector $d_i = y_i - \bar{x}_i \mathbf{1}$, $i = 1, \dots, p$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$



- Note: d_i and $\mathbf{1}$ are orthogonal
- Note: $\bar{x}_i = \frac{1}{n} \mathbf{1}^T y_i \Rightarrow n \bar{x}_i = \mathbf{1}^T y_i$
- Note: $\bar{x}_i \mathbf{1}$ is the projection of y_i onto $\frac{1}{n} \mathbf{1}$. $\text{Proj}_{\frac{1}{n} \mathbf{1}} y_i = \frac{\bar{x}_i \mathbf{1}}{\frac{1}{n}}$

* Prove that d_i and $\mathbf{1}$ are orthogonal

$$d_i^T \mathbf{1} = (y_i - \bar{x}_i \mathbf{1})^T \mathbf{1} = y_i^T \mathbf{1} - \bar{x}_i \mathbf{1}^T \mathbf{1} = \mathbf{1}^T y_i - \bar{x}_i n = 0$$

* The length of the deviation square $L_{d_i} = d_i^T d_i = \sum_{j=1}^n (x_{ji} - \bar{x}_i)^2$

For any i and k , $d_i^T d_k = \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = L_{d_i} L_{d_k} \cos \theta_{ik}$

$$\text{Then } \lambda_{ik} = \frac{d_i^T d_k}{\sqrt{d_i^T d_i} \sqrt{d_k^T d_k}} = \frac{\sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)}{\sqrt{\frac{1}{(n-1)} \sum_{j=1}^n (x_{ji} - \bar{x}_i)^2} \sqrt{\frac{1}{(n-1)} \sum_{j=1}^n (x_{jk} - \bar{x}_k)^2}} = \frac{L_{d_i} L_{d_k} \cos(\theta_{ik})}{\sqrt{L_{d_i}^2} \sqrt{L_{d_k}^2}} = \cos(\theta_{ik})$$

Then $\lambda_{ik} = \cos \theta_{ik} = \cos(\hat{d}_i, \hat{d}_k)$

$$* [S_n]_{ik} = \frac{1}{(n-1)} d_i^T d_k$$

$$[S_n] = \frac{1}{(n-1)} D^T D \text{ where } D = [d_1 | d_2 | \dots | d_p]$$

* 3.3 Random Samples and the expected values of the sample mean | covariance matrix

* Def 3.2

A sequence including n $p \times 1$ random vectors $\underline{X}_1, \dots, \underline{X}_n$ is a random sample if
- they are independent
- they have common distribution (i.e. dist of each $p \times 1$ vector is the same)

* Theorem 3.1

Let $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ be a random sample from $(\underline{\mu}, \underline{\Sigma})$. Then:

a) $E(\bar{X}) = \underline{\mu}$

b) $V(\bar{X}) = \frac{1}{n} \underline{\Sigma}$

c) $E(S_n) = \underline{\Sigma}$

Note that $s_{ij} = \frac{1}{(n-1)} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{j\ell} - \bar{x}_\ell)$ and $[S_n] = [s_{ij}]$

$$S_n = \frac{1}{(n-1)} \sum_{j=1}^n (\underline{X}_j - \bar{X})(\underline{X}_j - \bar{X})^T$$

3.4/p123. Generalized Variance.

Idea: $S = \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ s_{12} & s_{22} & \dots & s_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ s_{1p} & s_{2p} & \dots & s_{pp} \end{bmatrix} = [s_{ie}]_{p \times p}$

square symmetric

there are $p + \frac{1}{2} p(p-1)$ s_{ie}

$$s_{ii} = \frac{1}{(n-1)} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{ji} - \bar{x}_i)$$

where $\underline{x}_{ji} = \begin{bmatrix} x_{ji} \\ x_{2j} \\ \vdots \\ x_{pj} \end{bmatrix}$ $S \leftarrow \frac{1}{n-1}$
 $S_n \leftarrow \frac{1}{n}$ \uparrow variable

n persons \rightarrow

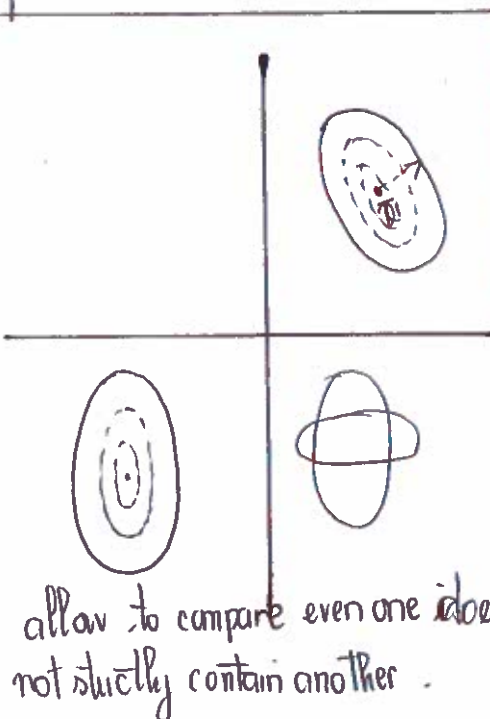
* Def 3.37

The generalized sample variance is $\text{Generalized var} = \det(S_n) = |S_n|$

* Note:

Let $V(d_1, \dots, d_p)$ be the volume spanned by d_1, \dots, d_p

Then $|S_n| = \frac{1}{(n-1)^p} |D|^2 = \frac{1}{(n-1)^p} V^2(d_1, \dots, d_p)$



* Recall: the quantity $(x - \bar{x})^T S^{-1} (x - \bar{x})$ has ellipse level curves

$\rightarrow \text{Volume}\{x \mid (x - \bar{x})^T S^{-1} (x - \bar{x}) \leq c\} = k_p |S|^{-1/2} c^{p/2} \propto |S|^{-1/2}$

$$k_p = \frac{2 \pi^{p/2}}{p} \Gamma(p/2)$$

* When is the generalized sample variance equal 0

(If some subset of x 's is a linear combination of the other)

EX: $x_1 = 0.5x_2 + 1.3x_3$
 or $x_1 = 0x_2 + 0x_3$
 or x_1 is a degenerated random variable

* Theorem 3.2

$$|S_n| = \left(\prod_{i=1}^p s_{ii} \right) |R| = \prod_{i=1}^p \lambda_i$$

where λ_i are eigenvalues of S_n (make sure units are the same)

* Def 3.2

The total (sample) variance is $\text{tr}(S_n) = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p s_{ii} = \frac{1}{n-1} \sum_{i=1}^p d_i^T d_i$

$$\text{tr}(R) = \sum_{i=1}^p \lambda_{ii} = \sum_{i=1}^p 1 = p \text{ (always)}$$

does not depend on correlation structure

4. Multivariate normal distribution

* Recall univariate normal pdf $X \sim N(\mu, \sigma^2)$ $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

• $\frac{(x-\mu)^2}{2\sigma^2} = (x-\mu)(\sigma^2)^{-1}(x-\mu)$ quadratic form

• Replace with $(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})$

* Then we have multivariate normal distribution

$X_{p \times 1} \sim N(\underline{\mu}_{p \times 1}, \underline{\Sigma}_{p \times p})$, $\underline{\Sigma}$ positive definite (has density)
positive semi-definite (generally)

with (*) has expression is the square of the generalized distance from $\underline{\mu}$

has pdf

$$f(\underline{x}) = \frac{1}{(2\pi)^{d/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})}$$

$\underline{\mu} \in \mathbb{R}^d$, $\underline{\Sigma} \in \mathbb{R}^{d \times d}$ (pd.)
 $\underline{x} \in \mathbb{R}^d$

* Example 4.1 The bivariate normal ($p=2$) $\underline{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$ X_1, X_2 ind $\Leftrightarrow \rho=0$

$$f(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right) \right]\right\}$$

• Note: X_1, X_2 independent $\Leftrightarrow \rho=0$

* Theorem 4.1

Let $\underline{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \sim N_p(\underline{\mu}, \underline{\Sigma})$

then X_1 and X_2 are independent $\Leftrightarrow \underline{\Sigma} = \begin{bmatrix} \underline{\Sigma}_{X_1} & | & 0 \\ \hline 0 & | & \underline{\Sigma}_{X_2} \end{bmatrix}$

* Theorem 4.2:

The contours of constant density for p-variate normal are

the ellipsoids such that $(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu}) = c^2$
center at $\underline{\mu}$

The axes of the ellipse are in the directions of the eigenvectors of $\underline{\Sigma}$
the length of the j^{th} longest axis is $\propto \sqrt{\lambda_j}$

major axis $\pm c\sqrt{\lambda_1} \underline{e}_1$, where $\underline{\Sigma} \underline{e}_j = \lambda_j \underline{e}_j$

* Result 4.4/153

If $\underline{\Sigma}$ positive definite $\Rightarrow \underline{\Sigma}^{-1}$ exists and is

($\underline{\Sigma}$ has eigenvalue $\lambda \Leftrightarrow \underline{\Sigma}^{-1}$ has eigenvalue $(\frac{1}{\lambda}, \underline{e}) = (\lambda^{-1}, \underline{e})$)
 (λ, \underline{e})

* Example 4.27

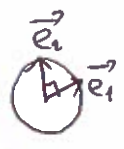
Let $\Sigma \in \mathbb{R}^2$ ($\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{11} \end{bmatrix}$) equal

• Compute the eigenvalues of Σ .

$$\det(\Sigma - \lambda I) = \det \begin{pmatrix} \sigma_{11} - \lambda & \sigma_{12} \\ \sigma_{12} & \sigma_{11} - \lambda \end{pmatrix} = (\sigma_{11} - \lambda)^2 - \sigma_{12}^2 \Rightarrow \lambda = \sigma_{11} \pm \sigma_{12}$$

• If $\sigma_{12} = 0 \Rightarrow \lambda_1 = \lambda_2$

e_1 and e_2 can be any orthogonal unit vectors

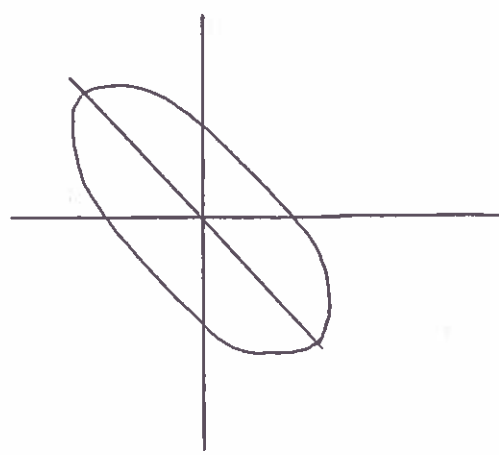
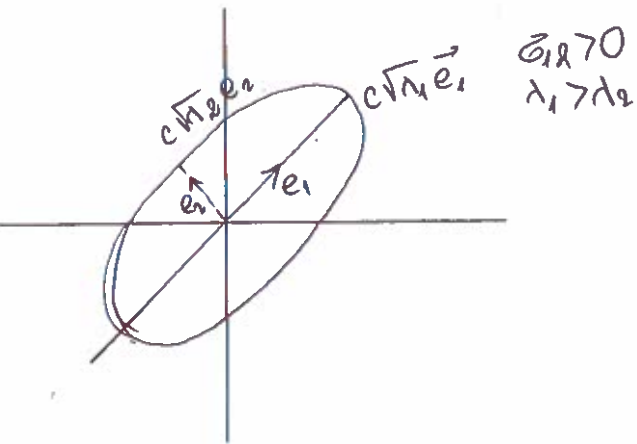


• If $\sigma_{12} > 0 \Rightarrow \lambda_1 > \lambda_2$ $\sigma_{12} < 0 \Rightarrow \lambda_1 < \lambda_2$

• Compute the eigenvector when $\sigma_{12} \neq 0$

$$(\lambda_1 I - \Sigma) \begin{pmatrix} e_1 \\ e_1 \end{pmatrix} = 0 \Leftrightarrow \begin{pmatrix} \sigma_{12} & -\sigma_{12} \\ -\sigma_{12} & \sigma_{12} \end{pmatrix} \Rightarrow e_1 \propto \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow e_1 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$$

$$(\lambda_2 I - \Sigma) \begin{pmatrix} e_2 \\ e_2 \end{pmatrix} = 0 \Leftrightarrow \begin{pmatrix} -\sigma_{12} & -\sigma_{12} \\ -\sigma_{12} & -\sigma_{12} \end{pmatrix} \begin{pmatrix} e_2 \\ e_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow e_2 \propto \begin{pmatrix} 1 \\ -1 \end{pmatrix} \Rightarrow e_2 = \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$$



* Properties of

Let $X \sim N_p(\mu, \Sigma)$ s.p.s.d ($\Sigma \succ 0$) then

a) $X + G \sim N_p(\mu + G, \Sigma)$

b) $AX \sim N_q(A\mu, A\Sigma A^T)$

c) If $\Sigma \succ 0$ (Σ is p.d), then $Y = \Sigma^{-1/2}(X - \mu) \sim N_p(0, I)$

d) If $\Sigma = \sigma^2 I$ (diagonal matrix with element σ^2)
 G is orthogonal ($G^T G = I$) $\Rightarrow GX \sim N(G\mu, \sigma^2 I)$
 in the meaning that $[I | 0]$

* Alternative definition of multivariate normal distribution (MVN)

Let $Z_1, Z_2, \dots, Z_p \stackrel{iid}{\sim} N(0, 1)$ Then let $Z = \begin{bmatrix} Z_1 \\ \vdots \\ Z_p \end{bmatrix} \sim N(0, I)$

Then $X = \mu + LZ \sim N_p(\mu, \Sigma = LL^T)$

* Definition 4.1 (for Chi Square)

Let $Z_1, \dots, Z_p \stackrel{iid}{\sim} N(0, 1)$ Let $Z = \begin{bmatrix} Z_1 \\ \vdots \\ Z_p \end{bmatrix} \sim N(0, I)$

Let $X = \sum_{i=1}^p Z_i^2 = Z^T Z \sim \chi_p^2$

* Theorem 4.4

Let $X \sim N_p(\mu, \Sigma)$ s.p.d.

Then $(X - \mu)^T (\Sigma)^{-1} (X - \mu) \sim \chi_p^2$

* Proof:

$(X - \mu)^T \Sigma^{-1} (X - \mu) = [\Sigma^{-1/2} (X - \mu)]^T [\Sigma^{1/2} (X - \mu)] = Y^T Y$ where $Y \sim N_p(0, I)$
 $\sim \chi_p^2$ theorem 4.3c

where $P(\chi_p^2 \leq \chi_p^2(c)) = c$

* Theorem 4.5

Let $Z \succ 0$ (Z p.d.) } Then $P\{(X-\mu)^T Z^{-1} (X-\mu) \leq \chi_{p, 1-\alpha}^2\} = 1-\alpha$

* Assessing Normality:

Recall QQ plots: plot pairs $(X_{(j)}, q_{(j)})$
 empirical theoretical
 \uparrow \uparrow
 j^{th} order statistics \leftarrow quantile

$q_{(j)} = F^{-1}(P_{(j)})$
 $0 < P_{(j)} < 1$
 $P_{(j)} < P_{(j+1)}$
 $P_{(j)} / \frac{1}{n} \rightarrow 1$

In R for normal

$q_{(j)} = \Phi^{-1}\left(\frac{j-0.375}{n+0.4}\right)$

If the iid data are normal, then the line should be straight

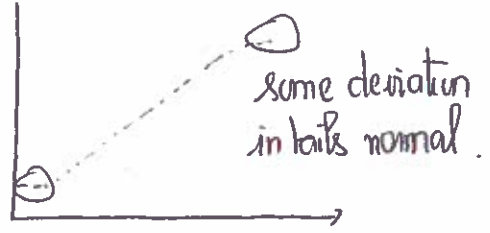
For univariate ps & data, apply theorem 4.4

plot $d_j^2 = j^{th}$ largest $(x_i - \bar{x}_i)^T S^{-1} (x_i - \bar{x}_i)$

$(d_j^2, \chi_p^2(P_{(j)}))$ or $(d_j^2, \sqrt{\chi_p^2(P_{(j)})})$

both should be line through origin with slope 1. (include (0,0) in plot)

Rayson test in MVN
 $P_{(j)} = \frac{j-0.5}{n}$



* Theorem 4.6

Let $X \sim N_p(\mu, \Sigma)$ (Σ p.d.)

Let $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ $\begin{matrix} p_1 \times 1 \\ p_2 \times 1 \end{matrix}$ $p_1 + p_2 = p$

$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$

$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ $\Sigma_{21} = \Sigma_{12}^T$

Then $X_1 | X_2 = x_2 \sim N_{p_1}(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})$

* Example

$$\text{Let } \mu = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & p & p^2 \\ p & 1 & p \\ p^2 & p & 1 \end{bmatrix}$$

Let $X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \sim N_3(\mu, \Sigma)$ Observe $X_1 = a$.

Find dist of $\begin{bmatrix} X_2 \\ X_3 \end{bmatrix} \Big|_{X_1 = a}$.

We have $\mu_{\begin{bmatrix} X_2 \\ X_3 \end{bmatrix} \Big|_{X_1 = a}} = \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (a - \mu_1) = 1 \begin{bmatrix} 1 & -p \\ 1 & -p^2 \end{bmatrix} + a \begin{bmatrix} p \\ p^2 \end{bmatrix}$

$$\begin{aligned} \Sigma_{\begin{bmatrix} X_2 \\ X_3 \end{bmatrix} \Big|_{X_1 = a}} &= \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} = \begin{bmatrix} 1 & p \\ p & 1 \end{bmatrix} - \begin{bmatrix} p \\ p^2 \end{bmatrix} [1]^{-1} [p \ p^2] = \\ &= \begin{bmatrix} 1 & p \\ p & 1 \end{bmatrix} - \begin{bmatrix} p^2 & p^3 \\ p^3 & p^4 \end{bmatrix} = \begin{bmatrix} 1-p^2 & p-p^3 \\ p-p^3 & 1-p \end{bmatrix} = (1-p^2) \begin{bmatrix} 1 & p \\ p & 1+p^2 \end{bmatrix} \end{aligned}$$

Then $\begin{bmatrix} X_2 \\ X_3 \end{bmatrix} \Big|_{X_1 = a} \sim N_2 \left(1 \begin{bmatrix} 1-p \\ 1-p^2 \end{bmatrix} + a \begin{bmatrix} p \\ p^2 \end{bmatrix}, (1-p^2) \begin{bmatrix} 1 & p \\ p & 1+p^2 \end{bmatrix} \right)$

* Find $X_3 \Big|_{\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix}}$ and $X_2 \Big|_{\begin{bmatrix} X_1 \\ X_3 \end{bmatrix} = \begin{bmatrix} a \\ c \end{bmatrix}}$

* Theorem 4.3.

$$\begin{matrix} 2) \\ \sim \end{matrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \begin{matrix} p_1 \\ p_2 \end{matrix} \sim N_{p_1+p_2} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Then $X_2 \sim N_{p_2}(\mu_2, \Sigma_{22})$

* Example 4.3.

$$X_2 | X_1 = a \sim N(\mu(1-p) + ap, (L-p^2))$$

$$X_3 | X_1 = a \sim N(\mu(L-p) + ap^2, (1-p^2)(L+p^2))$$

* 4.57 Maximum likelihood estimation of Σ and μ .

* Let $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$ ($\Sigma > 0$)

The likelihood is

$$L(\mu, \Sigma | X_1, \dots, X_n) = \prod_{j=1}^n \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \right) e^{-\frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)}$$

$$= \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)}$$

* Theorem 4.7

$$\sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) = \text{tr} \left[\Sigma^{-1} \left(\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T \right) \right] + n \left[(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu) \right]$$

$$= \text{tr} \left[\Sigma^{-1} (n-1)S \right] + n (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

* Result 4.9/168 Asymmetric Σ is a vector $\Rightarrow x^T A x = \text{tr}(x^T A x) = \text{tr}(A x x^T)$

* Proof

$$\sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) \stackrel{\text{insert } \bar{x}}{=} \sum_{j=1}^n (x_j - \bar{x} + \bar{x} - \mu)^T \Sigma^{-1} (x_j - \bar{x} + \bar{x} - \mu)$$

$$= \underbrace{\sum_{j=1}^n (x_j - \bar{x})^T \Sigma^{-1} (x_j - \bar{x}) + \sum_{j=1}^n (x_j - \bar{x})^T \Sigma^{-1} (\bar{x} - \mu) + \sum_{j=1}^n (\bar{x} - \mu)^T \Sigma^{-1} (x_j - \bar{x}) + \sum_{j=1}^n (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)}_{=0 \text{ because } \sum_{j=1}^n (x_j - \bar{x})^T = 0^T \text{ and } \sum_{j=1}^n (x_j - \bar{x}) = 0}$$

$$= \text{tr} \left(\sum_{j=1}^n (x_j - \bar{x})^T \Sigma^{-1} (x_j - \bar{x}) \right) + n (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

Σ^{-1} symmetric

$$= \text{tr} \left(\Sigma^{-1} \left(\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T \right) \right) + n (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

* Def 4.2 (Sufficiency)

$T(X)$ is sufficient for θ if $\mathbb{P}(X = x | T(x) = t, \theta) = \mathbb{P}(X = x | T(x) = t)$

* Theorem 4.8 (Factorization theorem)

T is sufficient for $\theta \iff f_{\theta}(x) = h(x) g_{\theta}(T(x))$

* Theorem 4.9

(S, \bar{X}) are jointly sufficient for μ, Σ

* Theorem 4.10 (this result will eventually allow us to obtain the maximum likelihood estimators of μ and Σ)

Let $B > 0$
 $b > 0$, a scalar

$$\left. \begin{array}{l} \frac{1}{|\Sigma|^b} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} B)} \leq \frac{1}{|B|^b} (2b)^{pb} e^{-bp} \text{ for all } \Sigma > 0 \\ \text{equality holds only for } \Sigma = \frac{1}{2b} B \end{array} \right\}$$

* Proof:

• Note that $\text{tr}(\Sigma^{-1} B) = \text{tr}(\Sigma^{-1} B^{1/2} B^{1/2}) = \text{tr}(B^{1/2} \Sigma^{-1} B^{1/2}) = \sum_{i=1}^p n_i$
 eigenvalues of $B^{1/2} \Sigma^{-1} B^{1/2}$

• $|B^{1/2} \Sigma^{-1} B^{1/2}| = |B| |\Sigma^{-1}|$
 $\Rightarrow |\Sigma^{-1}| = \frac{|B^{1/2} \Sigma^{-1} B^{1/2}|}{|B|} = \frac{\prod_{i=1}^p n_i}{|B|}$

$\Rightarrow \text{LHS} = \frac{1}{|\Sigma|^b} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} B)} = \frac{(\prod n_i)^b}{|B|^b} e^{-\frac{1}{2} \sum_{i=1}^p n_i} = \frac{1}{|B|^b} \prod n_i^b \prod e^{-\frac{1}{2} n_i} = \frac{1}{|B|^b} \prod (n_i^b e^{-\frac{1}{2} n_i})$

• Let $f(n) = n^b e^{-\frac{1}{2} n}$ (Gamma density) is maximized at $n = 2b$.

$\Rightarrow \text{LHS} \leq \frac{1}{|B|^b} \prod_{i=1}^p ((2b)^b e^{-b}) = \frac{1}{|B|^b} (2b)^{pb} e^{-bp}$

* Equality happens when $n_i = 2b, i = 1, \dots, p$
 i.e. the characteristic polynomial $B^{1/2} \Sigma^{-1} B^{1/2} = (\lambda - 2b)^p = 2b I$
 $B = 2b \Sigma \Leftrightarrow \Sigma = \frac{1}{2b} B$

* Theorem 4.11.

Let $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$

The MLE of μ is $\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

Σ is $\frac{n-1}{n} S = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$

* Proof: We have the likelihood is
 $L(\mu, \Sigma | X_1, \dots, X_n) = \frac{1}{(2\pi)^{np/2}} \left(\frac{1}{|\Sigma|^{n/2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} (\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T))} \right) e^{-\frac{n}{2} \frac{(\bar{X} - \mu)^T \Sigma^{-1} (\bar{X} - \mu)}{n}}$
 ≥ 0 unless $\mu = \bar{X}$
 (since $\Sigma > 0 \Rightarrow \Sigma^{-1} > 0$)

$\Rightarrow \hat{\mu} = \bar{X}$

• Now we want to maximize (*)
 Put $B = \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T = (n-1) S$

then (*) = $\frac{1}{|\Sigma|^{n/2}} e^{-\frac{1}{2} (\Sigma^{-1} \frac{(n-1) S}{b})} \leq \frac{1}{|B|^b} (2b)^{pb} e^{-bp}$ $b = \frac{n}{2}$

equality holds when $\Sigma = \frac{1}{2b} B = \frac{1}{2 \cdot \frac{n}{2}} (n-1) S = \frac{n-1}{n} S \quad \square$

max at $\hat{\mu} = \bar{X}$ since $\Sigma \succ 0 \Leftrightarrow \Sigma^{-1} \succ 0$

$-(\bar{X} - \mu)^T \Sigma^{-1} (\bar{X} - \mu) \leq 0$ then equal 0 only when $\bar{X} = \mu$ i.e. $\bar{X}_{ME} = \mu$

Follows from Theorem 4.10 with $b = \frac{n}{2}$ $B = \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^T = (n-1)S$

$$\hat{\Sigma} = \frac{1}{2b} B = \frac{1}{2(\frac{n}{2})} (n-1)S = \frac{n-1}{n} S$$

* Theorem 4.127

Let A be symmetric
 $p \times p$

B
 $q \times p$

Suppose $X \sim N_p(\mu, \Sigma)$

Then \bar{X} and $X^T A X$ are independent if

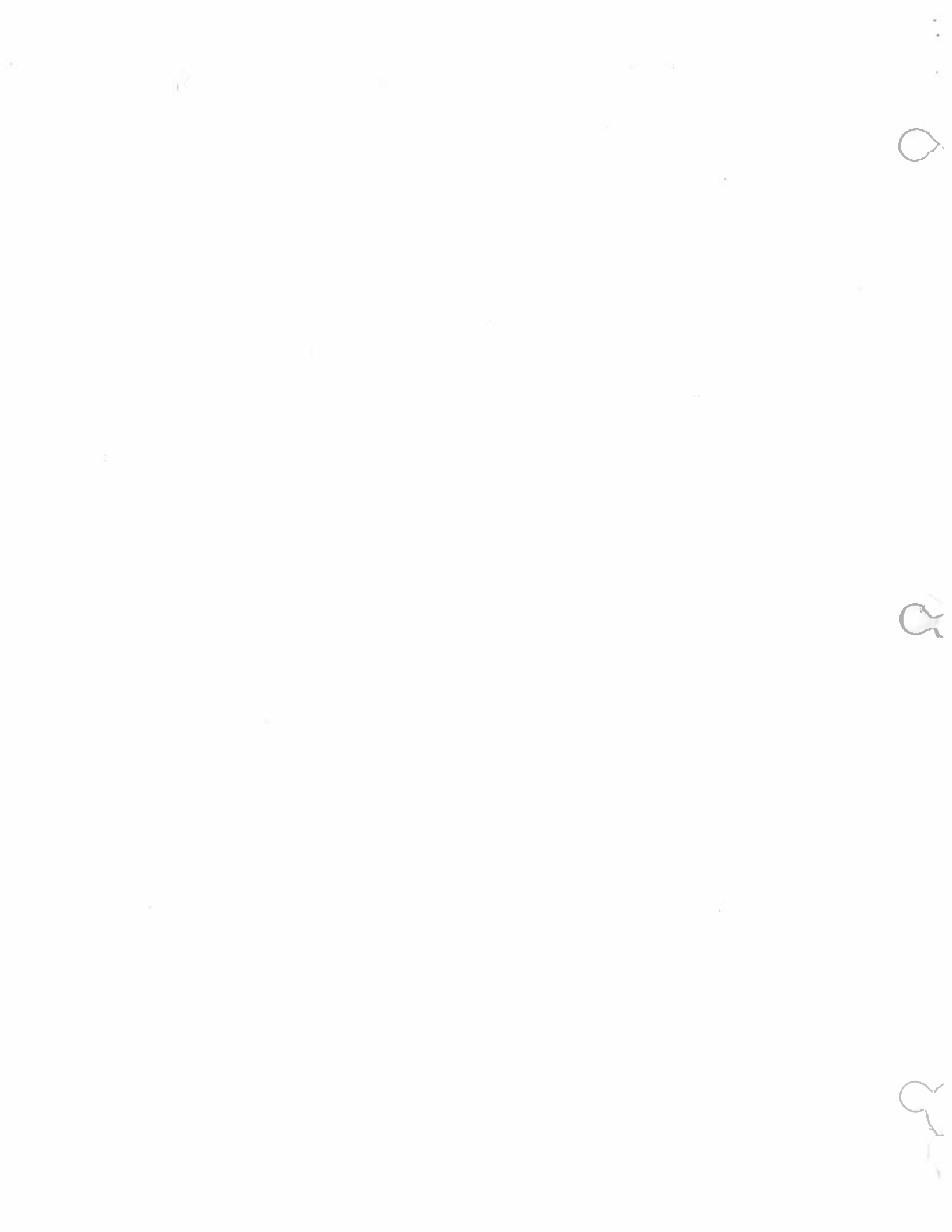
$$B \Sigma A = 0$$

* Theorem 4.137

Let $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$

then $\hat{\mu} = \bar{X}$ and $\hat{\Sigma} = \frac{n-1}{n} S$ are independent.

Σ positive $\Rightarrow \exists \Sigma^{-1/2} \Rightarrow \text{Var} = Q$



Q5: Inferences about a mean vector.

5.1 Introduction: p correlated variables must be analyzed jointly.

5.2 The plausibility of μ_0 as a value for a Normal population mean.

* Recall for univariate variables.

Want to test $H_0: \mu = \mu_0$
 $H_1: \mu \neq \mu_0$

\bullet If $X_1, \dots, X_n \sim \text{Normal}$, $\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i$ $S^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2$

Then the test statistic

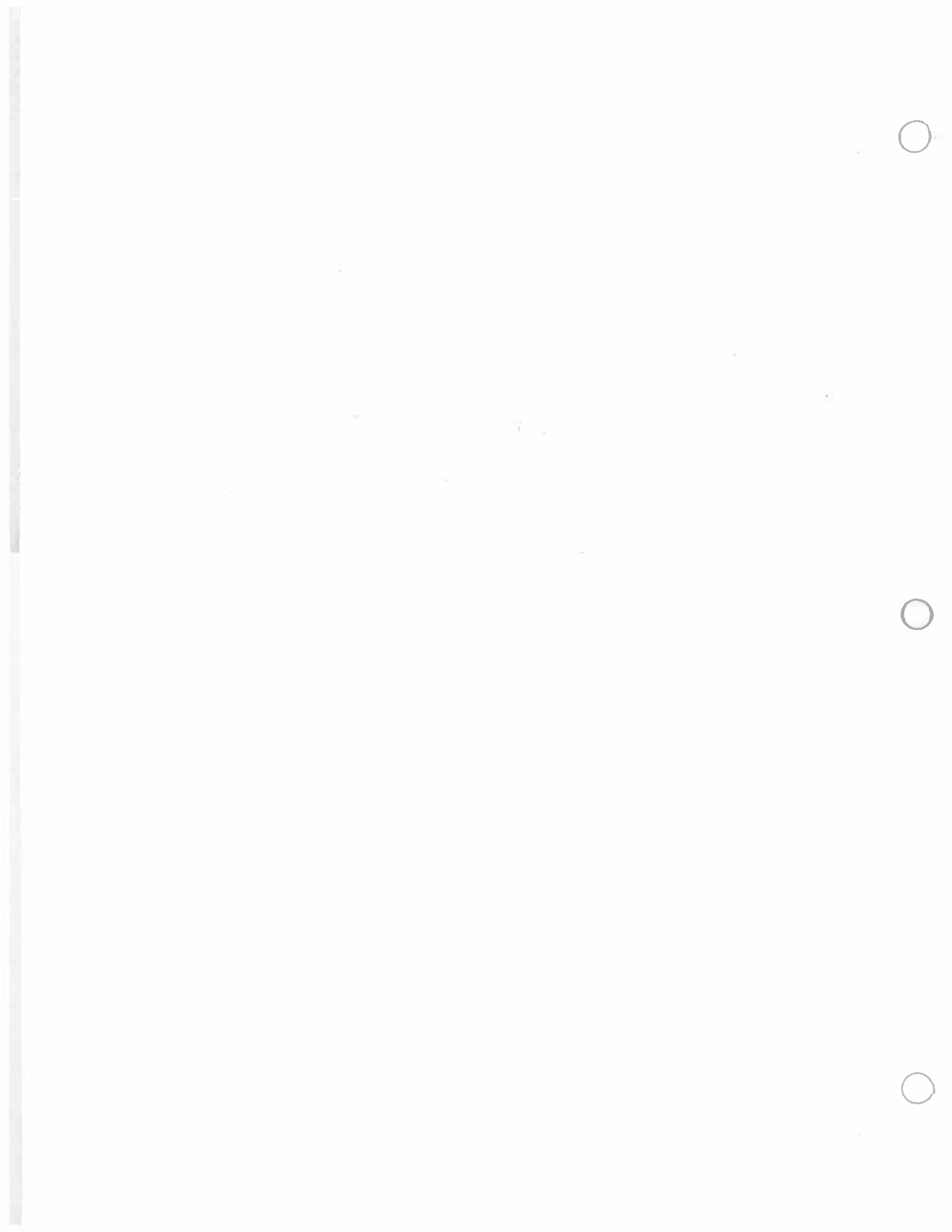
$t = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \sim t_{n-1}$ We reject H_0 when $|t_{\text{observed}}| > t_{\frac{\alpha}{2}, n-1}$.

$t^2 = n(\bar{X} - \mu_0)(S^2)^{-1}(\bar{X} - \mu_0) \Rightarrow$ Reject H_0 when $n(\bar{X} - \mu_0)(S^2)^{-1}(\bar{X} - \mu_0) > t_{\frac{\alpha}{2}, n-1}^2$.

* The confidence interval $CI = \left\{ \mu_0, \left| \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \right| < t_{\frac{\alpha}{2}, n-1} \right\} = \bar{X} \pm t_{\frac{\alpha}{2}, n-1} \frac{S}{\sqrt{n}}$

* Now, consider a problem with $(p \times 1)$ vector random vectors.

$$T^2 = n (\bar{\underline{X}} - \underline{\mu}_0)^T S^{-1} (\bar{\underline{X}} - \underline{\mu}_0)$$



* Chapter 2: Inference for the mean (c5 + c6 in the book).

* Wishart distribution.

Suppose that $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\underline{\mu}, \underline{\Sigma})$
 then $W_n = \sum_{i=1}^n X_i X_i^T \sim \text{Wishart}(\underline{\Sigma}, n)$

* Properties.

$W_1 \sim \text{Wishart}(\underline{\Sigma}, m_1)$

$W_2 \sim \text{Wishart}(\underline{\Sigma}, m_2)$

W_1 and W_2 are independent

Let $C_{q \times p}$: constant matrix

$a_{p \times 1}$: constant vector $a^T \underline{\Sigma} a \neq 0$

Then

a) $W_1 + W_2 \sim \text{Wishart}(\underline{\Sigma}, m_1 + m_2)$

b) $C W_1 C^T \sim W_p(C \underline{\Sigma} C^T, m_1)$

c) $\frac{a^T W_1 a}{a^T \underline{\Sigma} a} \sim \chi_{m_1}^2$

* In univariate case $X_1, X_2, \dots, X_n \sim N(0, \sigma^2)$

then the sample mean $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$ $\left(\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}\right)^2 \sim \chi_1^2$

* In multivariate case $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N_p(\underline{\mu}, \underline{\Sigma})$

Then $\bar{X} \sim N_p(\underline{\mu}, \frac{1}{n} \underline{\Sigma}) \Rightarrow \sqrt{n}(\bar{X} - \underline{\mu}) \sim N_p(0, \underline{\Sigma})$

$(n-1)S \sim W_p(\underline{\Sigma}, n-1)$

$\sqrt{n} \underline{\Sigma}^{-1/2} (\bar{X} - \underline{\mu}) \sim N_p(0, I)$

$n(\bar{X} - \underline{\mu})^T \underline{\Sigma}^{-1} (\bar{X} - \underline{\mu}) \sim \chi_p^2$

Replace $\underline{\Sigma}^{-1}$ by S^{-1}

$n(\bar{X} - \underline{\mu})^T S^{-1} (\bar{X} - \underline{\mu}) \stackrel{?}{\sim} \chi_p^2$
no.

> we can't use \bar{X} to infer about $\underline{\mu}$
 => introduce T^2 Hotelling's dist

* Page 174

When $X_1, X_2, \dots, X_n \sim N_p(\underline{\mu}, \underline{\Sigma})$. Then

a) $\bar{X} \sim N_p(\underline{\mu}, \frac{1}{n} \underline{\Sigma})$

$(n-1)S \sim W(\underline{\Sigma}, n-1)$

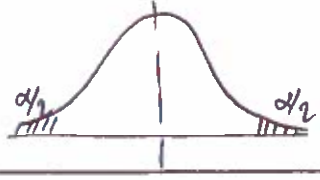
2) Hotelling's T^2

* Suppose that (univariate) $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

We use $t = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}}$ as a test statistic

(Remind test) $\begin{cases} \mu = \mu_0 \\ \mu \neq \mu_0 \end{cases}$

C.I. $\bar{X} \pm t_{\frac{1-\alpha}{2}, n-1} \frac{S}{\sqrt{n}}$ equivalent



* For multivariate distribution

Let $\underline{X} \sim N_p(\underline{0}, \Sigma)$
 $W \sim W_p(\Sigma, (n-1))$ independent.

Then $T^2 = \underline{X}^T \left(\frac{1}{n-1} W \right)^{-1} \underline{X} \sim T_{p, n-1}^2$: Hotelling T^2 distribution. ← not very user friendly

* Mahalanobis distance

$$d(\underline{x}, \underline{y}) = \sqrt{(\underline{x} - \underline{y})^T S^{-1} (\underline{x} - \underline{y})}$$

Reject H_0 when T^2 is large
 compute T^2 from data numbers.

then $T^2 = n (\bar{\underline{X}} - \underline{\mu}_0)^T S^{-1} (\bar{\underline{X}} - \underline{\mu}_0)$ n : sample size

$\bar{\underline{X}}, S$: sample mean, sample covariance matrix.

* Try to make connection to F

$$F = \frac{n}{n-1} \frac{n-p}{p} (\bar{\underline{X}} - \underline{\mu}_0)^T S^{-1} (\bar{\underline{X}} - \underline{\mu}_0) \sim F_{p, n-p}$$

$$T^2 \sim \frac{(n-1)p}{(n-p)} F_{p, n-p}$$

Note: $\frac{(n-p)}{n} T^2 \sim F_{p, n-p}$

2) S is non-singular

* Remind: Sample

$$\underline{X} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix} = \begin{bmatrix} X_{11} \\ X_{21} \\ \vdots \\ X_{n1} \end{bmatrix}^T = [Y_1 | Y_2 | \dots | Y_p]_{n \times p}$$

There are $n[X_i]_{p \times 1}$ represent for n people with p variables.

$$\bar{\underline{X}} = \begin{bmatrix} \bar{Y}_1 \\ \bar{Y}_2 \\ \vdots \\ \bar{Y}_p \end{bmatrix} = \frac{1}{n} \sum_{i=1}^n X_{ij}$$

$$S = \frac{1}{(n-1)} \sum_{i=1}^n (X_i - \bar{\underline{X}})(X_i - \bar{\underline{X}})^T$$

→ Two sample Hotelling's T^2 test.

Assume that the two independent groups of sample have the same sample variance.

Want to test

$$\begin{cases} H_0: \mu_1 = \mu_2 \\ H_1: \mu_1 \neq \mu_2 \end{cases}$$

$$\begin{cases} H_0: \mu_1 = \mu_2 \\ H_1: \mu_1 \neq \mu_2 \end{cases}$$

Test statistic $T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{X}_1 - \bar{X}_2)^T S^{-1} (\bar{X}_1 - \bar{X}_2)$

where S : pooled correlation matrix $S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{(n_1 + n_2 - 2)}$

* Note connection between Hotelling T^2 dist to F

$$F = \frac{(n_1 + n_2 - p - 1)}{(n_1 + n_2 - 2)p} T^2 \rightsquigarrow F_{p, n_1 + n_2 - p - 1}$$

We require $\begin{cases} n_1 + n_2 > p \\ S \text{ non singular} \end{cases}$

* One sample T^2 test:

We want to test $H_0: \mu = \mu_0$

$H_1: \mu \neq \mu_0$

At level α of significance, we reject H_0 when $T^2 = n (\bar{X} - \mu_0)^T S^{-1} (\bar{X} - \mu_0) > \frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha)$

4) Constant density Ellipses. (page 220-221) (to construct confidence interval)

Flury (1997) gives an interpretation of constant density ellipses in terms of Mahalanobis distance. Rejected if $T^2 > c^2 / p 285$

We wish to find region of squared Mahalanobis distance s.t.

$$\Pr\left(\frac{(\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu)}{c^2} \leq c^2\right) = (1-\alpha) \quad c^2 = \frac{n-1}{n} \frac{p}{n-p} F_{(1-\alpha), p, n-p}$$

c^2 objective

$F_{1-\alpha, p, n-p}$ is the $(1-\alpha)$ quantile of the F distribution with df $p, n-p$; sample size n

* we want to compare this constant density ellipse with normal CI

V * The Simultaneous confidence intervals

(care about within variable correlation) to do the test, we construct this interval and then compare.

Opt 2
12/06 12/10
Present Final

$$\left(\underline{a}^T \bar{X} \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{1-\alpha, p, n-p}} \underline{a}^T S \underline{a} \right)$$

where $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$ $\underline{a}^T S \underline{a}$ diag

Given that Σ positive definite

Let $\underline{a}_i = (0, \dots, 0, 1, \dots, 0)$

The the CI (Constant density ellipses).

$$\left(\bar{X}_i \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{1-\alpha, p, n-p}} S_{ii} \right)$$

* Want to compare the constant density ellipses with t interval.

• On at a time t intervals (regular t interval) ignore the correlation

Let $i = 1, \dots, p$

for μ_i , CI is $\left(\bar{X}_i \pm t_{n-1, \alpha/2} \sqrt{\frac{S_{ii}}{n}} \right)$ $S_{ii} = \text{diag} \Sigma$

• When the number m of specified component mean μ_i or $\underline{a}^T \mu$ is small

Bonferroni corrected C.I.s (χ^2 interval \rightarrow when we have very large sample)

$$\underline{a}^T \bar{X} \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{\underline{a}^T S \underline{a}}{n}}$$

5* Review of Anova

* Assume we have g random variables, $X_l \sim N(\mu_l, \sigma^2)$, $l = \overline{1, g}$

For each group $X_l \sim N(\mu_l, \sigma^2)$, we take a sample $X_{l,1}, \dots, X_{l,n_l}$, $l = \overline{1, g}$
 (The random samples are independent)

* For $l = \overline{1, g}$, $X_l \sim N(\mu_l, \sigma^2)$ then $\mu_l = \bar{X}_l = \mu + \tau_l$
 $X_l \sim N(\mu + \tau_l, \sigma^2)$
 τ_l group mean overall mean τ_l group effect (treatment effect)

* The response j in group l , $X_{l,j} \sim N(\mu + \tau_l, \sigma^2)$.

$X_{l,j} \sim N(\mu + \tau_l, \sigma^2)$

$X_{l,j} = \mu + \tau_l + \epsilon_{l,j} = \mu + \tau_l + N(0, \sigma^2)$, where $\epsilon_{l,j} \sim N(0, \sigma^2)$.

* $\sum_{l=1}^g n_l \tau_l = 0$

* We decompose $X_{l,j}$ as

$X_{l,j} = \bar{X} + (\bar{X}_l - \bar{X}) + (X_{l,j} - \bar{X}_l)$
 observation overall sample mean (an estimator of μ) estimated treatment effect (an estimator of τ_l) residual (an estimate of error $\epsilon_{l,j}$)

* $SST = SST_t + SSR$
 sum of square total SS^T treatment SS^R residual

$\sum_{l=1}^g \sum_{j=1}^{n_l} (X_{l,j} - \bar{X})^2 = \sum_{l=1}^g n_l (\bar{X}_l - \bar{X})^2 + \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{l,j} - \bar{X}_l)^2$

* One way ANOVA table (use to compare 2 means from two independent, unrelated groups using F distribution)

$H_0: \mu_1 = \mu_2$
 $H_1: \mu_1 \neq \mu_2$

Source	SS (sum of square)	df (degree of freedom)	MS (mean square)
treatment	$SST_t = \sum_{l=1}^g n_l (\bar{X}_l - \bar{X})^2$	$g - 1$	$\frac{SST_t}{(g - 1)}$
residual	$SSR = \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{l,j} - \bar{X}_l)^2$	$\sum_{l=1}^g n_l - g$	$\frac{SSR}{\sum_{l=1}^g n_l - g}$
total	$\sum_{l=1}^g \sum_{j=1}^{n_l} (X_{l,j} - \bar{X})^2$	$\sum_{l=1}^g n_l - 1$	

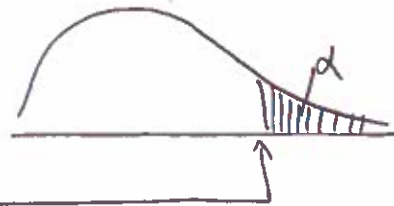
* When we want to test there is treatment effect or not treatment effect.

$$H_0: T_k = 0$$

$$H_1: T_k \neq 0$$

$$F_{\text{test statistic}} = \frac{\frac{SSt}{g-1}}{\frac{SSR}{\sum_{k=1}^g n_k - g}} = \frac{MSt}{MSR}$$

At level α , $F > F_{g-1, \sum_{k=1}^g n_k - g}(\alpha)$



67 MANOVA (Use for multiple comparison)

* The response $X_{ij} \sim N(\mu_j + \tau_j, \Sigma)$

$$X_{ij} = \mu_j + \tau_j + \epsilon_{ij} \quad \text{where } \epsilon_{ij} \sim N_p(0, \Sigma), \quad \sum_{k=1}^{n_j} n_k \tau_k = 0$$

* One way MANOVA (for comparing population mean vectors)

Source	SS	df
Treatment	$B = \sum_{k=1}^g n_k (\bar{X}_k - \bar{X})(\bar{X}_k - \bar{X})^T$	$g-1$
Residual	$W = \sum_{k=1}^g \sum_{j=1}^{n_k} (X_{kj} - \bar{X}_k)(X_{kj} - \bar{X}_k)^T$	$\sum_{k=1}^g n_k - g$
Total	$B + W = \sum_{k=1}^g \sum_{j=1}^{n_k} (X_{kj} - \bar{X})(X_{kj} - \bar{X})^T$	$\sum_{k=1}^g n_k - 1$

* We consider the test

$$H_0: \tau_j = 0, \quad p = 1, g$$

$$\Lambda^* = \frac{|W|}{|B+W|} \quad \leftarrow \text{Wilks' lambda.}$$

Reject H_0 when Λ^* is small.

* Let $\lambda_1, \dots, \lambda_p$ be the eigenvalues of (BW^{-1})

$$\text{then } \Lambda^* = \frac{|I|}{|BW^{-1} + I|} = \prod_{i=1}^p \frac{1}{\lambda_i + 1}$$

* Assumption and limitation.

1) The response variables are continuous.

2) The residuals follow the $N(0, \Sigma)$

3) The individuals are independent

4) The variance-covariance matrices of each group of residuals are equal

* Box's M test (to test that the variance-covariance

$$H_0: \Sigma_1 = \Sigma_2 = \dots = \Sigma_g$$

$$\Lambda_{\text{box}} = \prod_{i=1}^g \left(\frac{|S_i|}{|S_{\text{pooled}}|} \right)^{\frac{n_i-1}{2}}, \quad S_{\text{pooled}} = \frac{1}{\sum_{i=1}^g (n_i-1)} \left[(n_1-1)S_1 + (n_2-1)S_2 + \dots + (n_g-1)S_g \right]$$

$$\text{Let } M = -2 \ln \Lambda = \sum_{i=1}^g (n_i-1) \ln |S_{\text{pooled}}| - \sum_{i=1}^g (n_i-1) \ln |S_i|$$

* Extend to 2 way MANOVA

$[X_1, X_2, X_3]$

	<u>rate</u>	<u>additive</u>
	low high	low high
	rate x additive (LL, HH, LH, HL)	

Source

- ① Treatment 1
- ② Treatment 2
- ③ Treatment 1+2

Residual

Total

b+w

* Example 2-way ANOVA

variable	treatments			
	rate	additive	rate × additive	
$[X_1, X_2, X_3]$	low	low	l	h
	high	high	h	h
			l	l
			h	h

$H_0: \mu_{1,t_1} = \mu_{1,t_2}$

not significant.

by function.

* Multiple comparison in MANOVA

inference $\left\{ \begin{array}{l} \text{CI for } \mu \\ \text{hypotheses test} \end{array} \right.$ for $\mu_1, \mu_2, \dots, \mu_g$

$$\mu_1 = \mu_{11} - \mu_{12} \Rightarrow \mu_{11} = \mu_{12} \Rightarrow \mu_{11} - \mu_{12}$$

with residual diagonal of the residual matrix

$(t_{1-\alpha})$ (B adjusted interval) $\left(\begin{array}{c} \text{non} \\ \text{is} \end{array} \right)$

↑
approximate of normal Z

* Example MANOVA

4 variables Species
3 species

table (I wish # Species) X matrix

first observe that 4 variable are \neq
 \rightarrow then we want to compare them in pairs

Variable 1 \rightarrow Spec 1 & 2
 1 & 3
 2 & 3.

2 \rightarrow

3 \rightarrow

4 \rightarrow

} 3x4 CI

C3 →

17

β : unknown

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \quad & \quad & \quad \\ \quad & \quad & \quad \\ \quad & \quad & \quad \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_L \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$Y = X\beta + \epsilon$$

data

want to get $\hat{\beta}$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1L} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & \dots & x_{nL} \end{bmatrix}_{n \times (L+1)} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_L \end{bmatrix}_{(L+1) \times 1} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

$$Y = \beta^T X + \epsilon$$

$1 \times n$ $L \times (L+1)$ $(L+1) \times 1$ $1 \times n$

← different way ← keep covariance then OK

Theorem 1: (Least square)

* Least square
Maximum like

← try to minimize the square of the residual

$$\hat{\beta}_{LS} = \underset{\beta}{\text{argmin}} (Y - X\beta)^T (Y - X\beta)$$

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T Y$$

$$\begin{bmatrix} X \\ \vdots \\ X \end{bmatrix}$$

H: Hat matrix

Y: Y observed

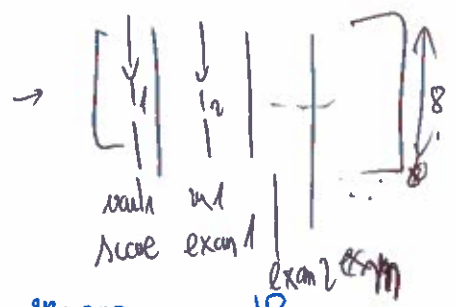
T2: L

T3: β : true mean

T4: symmetric → X^2
normality t

2.7 Multivariate Regression

Y : matrix
 $n \times m$



(n individuals
 n observations)

response variable

$$\begin{bmatrix} Y_{11} & \dots & Y_{1m} \\ \vdots & & \vdots \\ Y_{n1} & \dots & Y_{nm} \end{bmatrix} = \begin{bmatrix} X_{11} & X_{1r} \\ \vdots & \vdots \\ X_{n1} & X_{nr} \end{bmatrix} \begin{bmatrix} \beta_{01} \\ \vdots \\ \beta_{0r} \end{bmatrix} + \begin{bmatrix} \beta_{0m} \\ \vdots \\ \beta_{0n} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$\underbrace{\quad}_{\beta_0}$ $\underbrace{\quad}_{\beta_1}$

2 way $\left\{ \begin{array}{l} \text{least square } m \text{ respond variable} \\ \text{book Kim's course} \end{array} \right.$

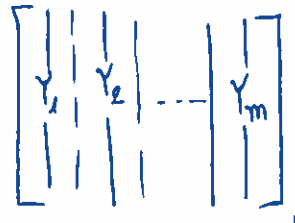
Read the prove

(us max likelihood esting)

want to get $\hat{\beta}$.

$$\begin{bmatrix} \epsilon_1^T \\ \epsilon_2^T \\ \vdots \\ \epsilon_n^T \end{bmatrix}$$

Y matrix



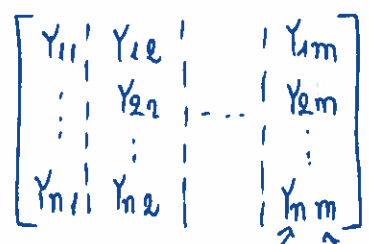
m response variable

$$n \times m \quad n \times (\lambda+1) \quad (\lambda+1) \times m \quad - \quad n \times m$$

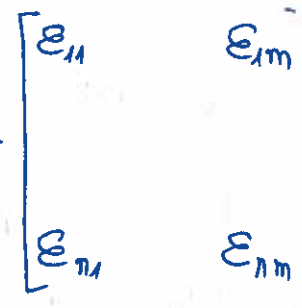
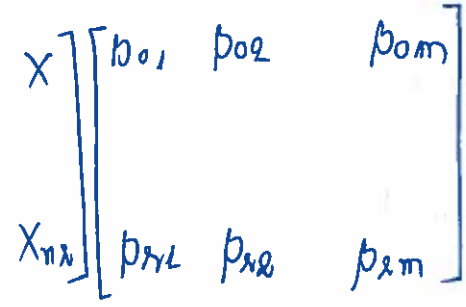
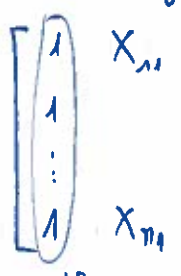
$$\underline{Y} = \underline{X} \underline{\beta} + \underline{\epsilon}$$

λ : # of predictor
 m : # of response

$$= \{ \epsilon_{11}, \dots, \epsilon_{nm} \}$$



person response variable



$$\epsilon_{(j)} \sim N(0_n, \sigma_{11}^2 I_n) \quad i = 1, \dots, m, \quad 1^{th} \text{ column}$$

$$\epsilon_{(i)} \sim N(0_m, \Sigma) \quad j^{th} \text{ row}$$

full model
 reduce model
 linear regression

$$\text{vec}(\epsilon) = \begin{bmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{n1} \\ \epsilon_{12} \\ \vdots \\ \epsilon_{n2} \\ \vdots \\ \epsilon_{1m} \\ \vdots \\ \epsilon_{nm} \end{bmatrix} \sim N(0_{nm}, \Sigma \otimes I_n)$$

* $\beta_2(0) \Rightarrow$ we want to test if only the upper part has affect on the data
 the lower part does not have affect on the data

use χ^2
 not friendly

lm command | resid $\leftrightarrow \epsilon$

$H_0: \beta_2 = 0$
 Reject H_0

$$l=2 \quad m=2$$

$$\hat{Y} = X \hat{\beta} \quad * \text{ want to test } \beta_2 = 0$$

$$\hat{\epsilon} = Y - \hat{Y} \quad \hat{\Sigma}_{NIE} = \frac{1}{n} (Y - \hat{Y})^T (Y - \hat{Y}) = \frac{1}{n} \hat{\epsilon}^T \hat{\epsilon}$$

all $\beta = 0$?

* Quicker way

directly compute residual without computing \hat{y} , $\hat{y} - y$

$$Y = M + \hat{Y}_e + \epsilon_e \quad (\text{relation between linear regression} \neq \text{MANOVA})$$

$$\text{MANOVA helps test } \begin{cases} \hat{Y}_{full} = X_1 \hat{\beta} \\ Y_{reduce} = \end{cases}$$

when MANOVA still be valid $\Rightarrow X_2$ is not important.

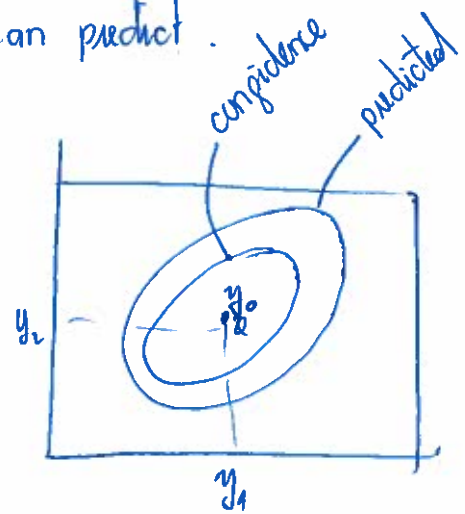
* Confidence interval
Predicted

\leftarrow new data comes in \Rightarrow we can predict.

$$\hat{Y}_0 = X_0 \hat{\beta}$$

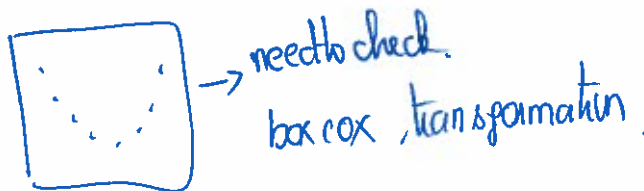
↑ new data ↑ the $\hat{\beta}$ that we got from previous data

predicted data \rightarrow predicted interval



• Pair data new $X_0 = \begin{bmatrix} X_{011} & X_{01v} \\ X_{021} & X_{02v} \end{bmatrix}$

• Residual plot



* $\left\{ \begin{array}{l} \text{partly reject the } H_0 \\ \text{type III F test} \end{array} \right.$

With MANOVA
we assume that they are
real factors
affect the variable.







$$\text{Inve}(S) = S^{-1}$$

$$t(S) = S^T$$

Chapter 2. Inference for the mean

Jianxuan Liu

Fall 2018

1. Wishart Distribution

Let $X_1, \dots, X_n \sim^{iid} N_p(\mathbf{0}, \Sigma)$, then

$$W_n = \sum_{i=1}^n X_i X_i^T \sim \text{Wishart}(\Sigma, n).$$

Let $X_1 \sim W_p(\Sigma, m_1)$ and $X_2 \sim W_p(\Sigma, m_2)$ be independent. Let C be a $q \times p$ constant matrix and a be a p constant vector s.t. $a^T \Sigma a \neq 0$. then

- $X_1 + X_2 \sim W_p(\Sigma, m_1 + m_2)$
- $CX_1C^T \sim W_p(C\Sigma C^T, m_1)$
- $\frac{a^T X a}{a^T \Sigma a} \sim \chi_{m_1}^2$

Let $X_1, \dots, X_n \sim^{iid} N_p(\mu, \Sigma)$. then

- $\bar{X} \sim N_p(\mu, \Sigma/n)$.
- $(n-1)S \sim W_p(\Sigma, n-1)$.

Since $\bar{X} \sim N_p(\mu, \Sigma/n)$, then $\sqrt{n}(\bar{X} - \mu) \sim N_p(\mathbf{0}, \Sigma)$ and $\sqrt{n}\Sigma^{-1/2}(\bar{X} - \mu) \sim N_p(\mathbf{0}, \mathbf{I})$. We also have $n(\bar{X} - \mu)^T \Sigma^{-1}(\bar{X} - \mu) \sim \chi_p^2$. $n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu)$ does not follow χ_p^2 .

2. Hotelling's T^2

Suppose that (in univariate case) $x_1, \dots, x_n \sim^{iid} N(\mu, \sigma^2)$. We use $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$ as a test statistic. The confidence interval for μ is

$$\bar{x} \pm \frac{s}{\sqrt{n}} t_{1-\alpha/2, n-1}$$

Let $\mathbf{x} \sim N_p(\mathbf{0}, \Sigma)$ and $W \sim W_p(\Sigma, n-1)$ be independent, then

$$T^2 = \mathbf{x}^T \left(\frac{1}{n-1} W \right)^{-1} \mathbf{x} \sim T_{p, n-1}^2$$

Consideration of the squared Mahalanobis distance leads us to consider the so called T^2 statistic (the nomenclature reflects that this relates to a t -statistic for one variable). This can be found as:

$$T^2 = n(\bar{\mathbf{x}} - \mu_0)^T \mathbf{S}^{-1}(\bar{\mathbf{x}} - \mu_0)$$

where n is the sample size, μ_0 is the hypothesized mean, $\bar{\mathbf{x}}$ and \mathbf{S} are the sample mean and covariance matrices respectively. It turns out that this statistic follows a T^2 distribution, however, given that there is a simple relationship between the T^2 and F distribution it is often easier to work with the latter. If $x_i, i = 1, \dots, n$ represent a sample from a p variate normal distribution with mean μ_0 and covariance Σ , provided Σ is positive definite and $n > p$, given sample estimators for mean and covariance $\bar{\mathbf{x}}$ and \mathbf{S} respectively, then:

$$F = \frac{n}{n-1} \frac{n-p}{p} (\bar{\mathbf{x}} - \mu_0)^T \mathbf{S}^{-1}(\bar{\mathbf{x}} - \mu_0)$$

follows an F -distribution with p and $(n - p)$ degrees of freedom. Note the requirement that $n > p$, i.e. that S is non-singular.

3. Two sample Hotelling's T^2 test

Analogous to the univariate context, we wish to determine whether the mean vectors are comparable, more formally:

$$H_0 : \mu_1 = \mu_2$$

The T^2 statistic proposed by Hotelling (1931), will be based this time on the distance between two mean vectors. It can be calculated as:

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{x}_1 - \bar{x}_2)^T S^{-1} (\bar{x}_1 - \bar{x}_2)$$

where S^{-1} is the inverse of the pooled covariance matrix given by:

$$S_{\text{pooled}} = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$$

given the sample estimates for covariance, S_1 and S_2 in the two samples. As before, there is a simple relationship between the test statistic, T^2 , and the F distribution.

If $x_{1i}, i = 1, \dots, n_1$ and $x_{2i}, i = 1, \dots, n_2$ represent independent samples from two p variate normal distribution with mean vectors μ_1 and μ_2 but with common covariance matrix Σ , provided Σ is positive definite and $n > p$, given sample estimators for mean and covariance \bar{x} and S respectively, then:

$$F_{\text{observed}} \leftarrow F = \frac{(n_1 + n_2 - p - 1)T^2}{(n_1 + n_2 - 2)p} \sim F_{p, n_1 + n_2 - p - 1}$$

has an F distribution on p and $(n_1 + n_2 - p - 1)$. Essentially, we compute the test statistic, and see whether it falls within the $(1 - \alpha)$ quantile of the F distribution on those degrees of freedom. Note again that to ensure non-singularity of S , we require that $n_1 + n_2 > p$.

We are going to consider an example using data from Flea Beetles reported by Lubischew (1962) and used in (Flury, 1997, page 307). It should be noted that in terms of practical computation, methods are based on the QR decomposition will be used, details are given in Seber (1984). However, for the purposes of understanding the principles behind the test, we follow the formula directly.

```
library(Flury)
?flea.beetles
data(flea.beetles)
```

Species	TG	Elytra	Second.Antenna	Third.Antenna
1 Oleracea				
2 Carduorum				

It can be seen that there is a factor "Species" denoting whether the beetles are from 'oleracea' or 'carduorum'. There are four numeric variables as follows:

- 'TG': Distance of the Transverse Groove to the posterior border of the prothorax (microns)
- 'Elytra': Length of the Elytra (in units of 0.01mm)
- 'Second.Antenna': Length of the second antennal joint (microns)
- 'Third.Antenna': Length of the third antennal joint (microns).

We need to estimate the mean for each sample, and calculate the difference between the two vectors:

```
mu <- by(flea.beetles[, -1], flea.beetles$Species, colMeans)
mudiff <- mu[[1]] - mu[[2]]
p <- dim(flea.beetles)[2] - 1 ## how many variables are we using
```

The next step is to extract the two covariance matrices:

Variables: TG, Elytra
Second.Antenna
Third.Antenna.

2 species.
Oleracea eardn
X1 X2
TG = X11
Elytra or X22
Second Ant f X23
Third Ant X24
mu_x1 = mu_x2
mu_x1 != mu_x2

CR2 - Hotelling T2 R
TG = X11
Elytra or X22
Second Ant f X23
Third Ant X24

$$\text{covmats} = [S_1, S_2] \quad \text{covmats}[[1]] = S_1 \quad \text{covmats}[[2]] = S_2.$$

```

covmats <- by(flea.beetles[,-1], flea.beetles$Species, cov)
covmats

## flea.beetles$Species: oleracea
##                TG      Elytra Second.Antenna Third.Antenna
## TG              187.59649 176.86257      48.37135      113.58187
## Elytra           176.86257 345.38596      75.97953      118.78070
## Second.Antenna  48.37135 75.97953       66.35673       16.24269
## Third.Antenna  113.58187 118.78070      16.24269      239.94152
## -----
## flea.beetles$Species: carduorum
##                TG      Elytra Second.Antenna Third.Antenna
## TG              101.83947 128.06316      36.98947      32.59211
## Elytra           128.06316 389.01053      165.35789      94.36842
## Second.Antenna  36.98947 165.35789      167.53684      66.52632
## Third.Antenna   32.59211 94.36842       66.52632      177.88158

```

and then to estimate the pooled covariance matrix S for the flea beetle data (where $N[1]$ gives n_1 , $N[2]$ gives n_2), can be calculated as:

```

N <- xtabs(-flea.beetles[,1]) # count of the level of categorical variables
pooledS <- ((N[1]-1) * covmats[[1]] + (N[2]-1) * covmats[[2]]) / (N[1] + N[2] - 2)
pooledS

```

	X_1	X_2	X_3	X_4
	TG	Elytra	Second.Antenna	Third.Antenna
## TG	143.55910	151.8034	42.52660	71.99253
## Elytra	151.80341	367.7878	121.87653	106.24467
## Second.Antenna	42.52660	121.8765	118.31408	42.06401
## Third.Antenna	71.99253	106.2447	42.06401	208.07290

```

Sinv <- solve(pooledS) # calculate the inverse of Spooled
Sinv

```

```

##                TG      Elytra Second.Antenna Third.Antenna
## TG              0.013257964 -0.0053492256  0.0015134494 -0.0021617878
## Elytra          -0.005349226  0.0066679441 -0.0047337699 -0.0005969439
## Second.Antenna  0.001513449 -0.0047337699  0.0130490933 -0.0007445297
## Third.Antenna  -0.002161788 -0.0005969439 -0.0007445297  0.0060093005

```

Having calculated the inverse of the pooled correlation matrix we also need the scaling factor $\frac{n_1 n_2}{n_1 + n_2}$.

Hotellings T^2 is then quite straightforward to calculate:

```

scaleFact <- (N[1]*N[2]) / (N[1]+N[2])
Hotellings <- t(mudiff) %*% Sinv %*% mudiff * scaleFact # T^2
Hotellings

```

```

##                [,1]
## [1,] 133.4873

```

which is the value of the T^2 statistic. We could work with this value directly, but it is more convenient to transform it into something we can compare with the F distribution.

```

test <- ((N[1] + N[2] - p - 1) * Hotellings) / ((N[1] + N[2] - 2) * p) # test = ((n1+n2-p-1)T^2) / ((n1+n-2)P)
test

```

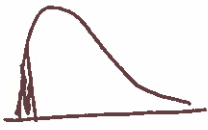
```

##                [,1]
## [1,] 30.666

```

Reject H_0 when $F_{\text{observed}} = \frac{(n_1+n_2-p-1) T^2}{(n_1+n_2-2) p} > F_{p, n_1+n_2-p-2}$

When we want to test $\begin{cases} H_0: \mu_1 = \mu_2 \\ H_1: \mu_1 \neq \mu_2 \end{cases}$



Stat 505 note
 we would reject H_0 at level α if $F_{observed} > F_{p(n_1+n_2-p), \alpha}$

* Equal covariance case:
 Reject when $T_{observed}^2 > T_{1-\alpha, p, n_1+n_2-2}^2$

and we compare the follows: *degree*

an check this as

```
pf(test, p, N[1]+N[2]-p-1, lower.tail = FALSE)
```

$P(Y_1, Y_2)$
 inverted = T $\Rightarrow P(X \leq 2)$
 lower tail = F $\Rightarrow P(X > 2)$

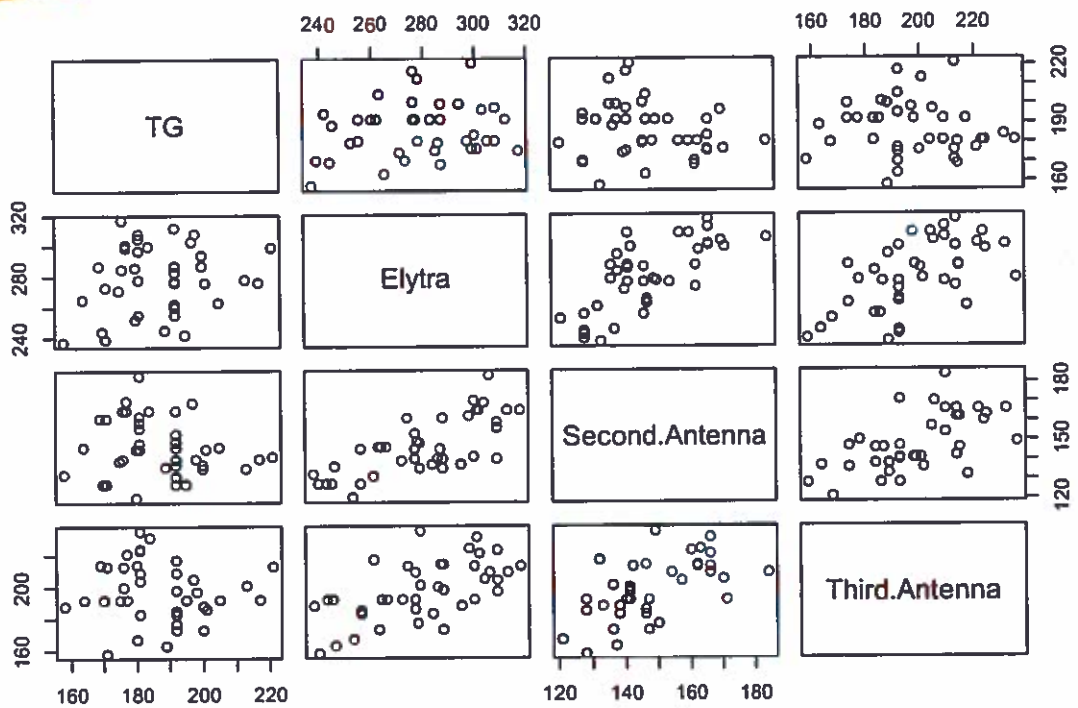
```
## [1,]
## [1,] 7.521799e-11
```

$\Rightarrow P(F > F_{crit}) \Rightarrow$ really small $\Rightarrow F_{obs} < F$

which gives us the area under the curve from our test statistic (30.666) to 0. Clearly in this case, we have reject H_0 , i.e. there is evidence that the mean vectors, $\bar{x}_{ultraxea} = (194.4737, 267.0526, 137.3684, 185.9474)$, $\bar{x}_{carduorum} = (179.55, 290.80, 157.20, 209.25)$, for the two species differ.

```
pairs(flea.beetles[, -1], col=flea.beetles[, 1])
```

$\Rightarrow \bar{x}_1 \neq \bar{x}_2$



4. Constant Density Ellipses

Flury (1997) gives an interpretation of constant density ellipses in terms of the Mahalanobis distance which is worth reading. Essentially, we wish to find a region of squared Mahalanobis distance such that:

$$Pr((\bar{x} - \mu)^T S^{-1} (\bar{x} - \mu) \leq c^2) = 1 - \alpha$$

and we can find c^2 as follows:

$$c^2 = \frac{n-1}{n} \frac{p}{n-p} F_{(1-\alpha), p, (n-p)}$$

where $F_{(1-\alpha), p, (n-p)}$ is the $(1-\alpha)$ quantile of the F distribution with p and $n-p$ represents the number of variables and n the sample size. Illustration: Firstly, we need a function to draw ellipses:

```
Ellipse <- function(covmat, centroid, csquare, resolution, plot = TRUE) {
  angles <- seq(0, by = (2 * pi)/resolution, length = resolution)
  sd <- covmat[1,2] / sqrt(covmat[1,1] * covmat[2,2])
  projmat <- matrix(0,2,2)
```

Ch 2. Constant Density Ellipse

length

```

projmat[1,1] <- sqrt(covmat[1,1] %*% (1+sd)/2)
projmat[1,2] <- -sqrt(covmat[1,1] %*% (1-sd)/2)
projmat[2,1] <- sqrt(covmat[2,2] %*% (1+sd)/2)
projmat[2,2] <- sqrt(covmat[2,2] %*% (1-sd)/2)
circle <- cbind(cos(angles), sin(angles))
Ellipse <- t(centroid + sqrt(csquare) * projmat %*% t(circle))
if (plot == TRUE) {lines(Ellipse)}
return(Ellipse)
}

```

compute c^2

It is possible to define a function which calculates c^2 and calls the Ellipse routine

```

cdellipse <- function (data, alpha=0.05, resolution=500)
{
  xbar <- colMeans(data)
  n <- dim(data)[1]
  p <- dim(data)[2]
  f <- qf(1-alpha, p, n-p)
  csquare <- ((n-1)/n) * (p / (n-p)) * f
  cat(csquare)
  Ellipse <- Ellipse(cov(data), xbar, csquare, resolution)
}

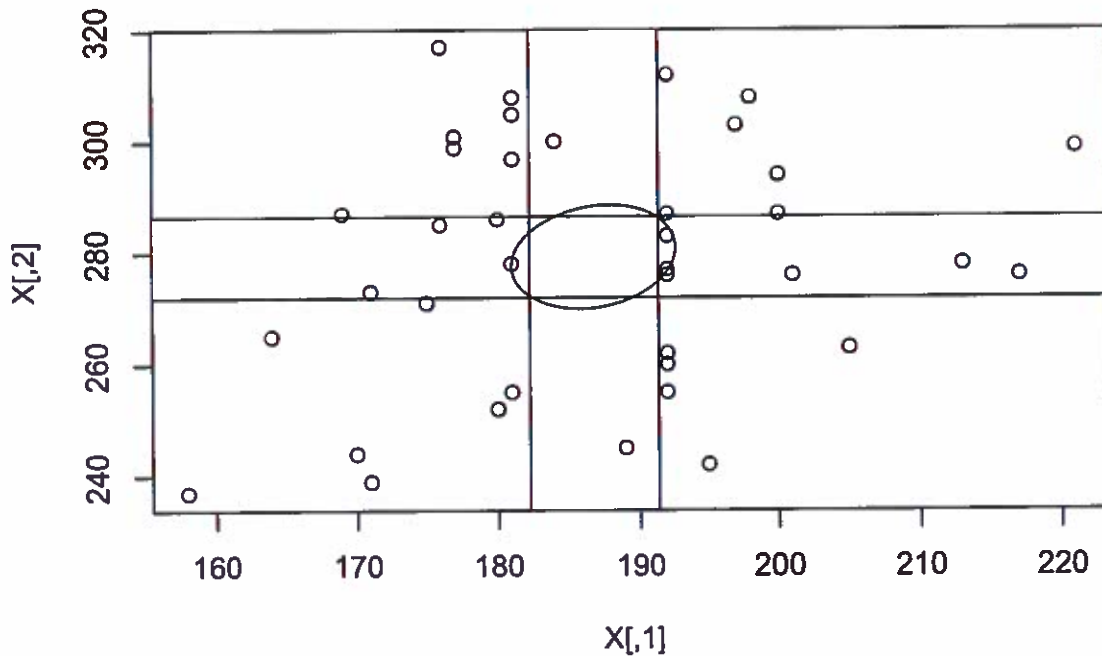
```

For illustrative purposes, we'll create a $n \times 2$ data object from our flea beetles, and plot the confidence ellipse for these. $\alpha = 0.05$

```

X <- cbind(flea.beetles[,2], flea.beetles[,3])
plot(X)
cdellipse(X, alpha = 0.05) # the cdelipse helps draw an ellipse covers 2 variables
## 0.1712725
#These can be contrasted with the univariate confidence intervals:
abline(v = confint(lm(X[,1]~1)))
abline(h = confint(lm(X[,2]~1))) } # univariate conf for the two above variables

```



It can be seen that the univariate confidence intervals and the constant density ellipse support different areas of parameter space. Ignoring the correlation structure in these data could lead to flaws in inference when assessing parameter uncertainty.

Results and comparison

Let $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ be a random sample from an $N_p(\mu, \Sigma)$ population with Σ positive definite.

- Simultaneously for all \mathbf{a} , the intervals

$$\left(\mathbf{a}^T \bar{\mathbf{x}} \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{1-\alpha, p, n-p} \mathbf{a}^T \mathbf{S} \mathbf{a}} \right)$$

are exact at $1 - \alpha$ for all $\mathbf{a}^T \mu$.

- Let $(\mathbf{a}_i = (0, \dots, 0, 1, 0, \dots, 0))$, then $\mathbf{a}_i^T \mu = \mu_i$ and

$$\left(\mathbf{a}_i^T \bar{\mathbf{x}} \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{1-\alpha, p, n-p} \mathbf{a}_i^T \mathbf{S} \mathbf{a}_i} \right)$$

that is

$$\left(\bar{x}_i \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{1-\alpha, p, n-p} S_{ii}} \right)$$

is simultaneously exact for all μ_i at $1 - \alpha$ level.

- One-at-a-Time intervals ignore the covariance structure of the p variables. This approach constructs confidence intervals for μ_j one at a time. For $j = 1, \dots, p$, the One-at-a-Time intervals are

$$\left(\bar{x}_j \pm t_{n-1, \alpha/2} \sqrt{s_{jj}/n} \right)$$

- When the number m of specified component mean μ_i , or linear combinations $\mathbf{a}^T \boldsymbol{\mu}$ is small, then a better approach is to calculate the Bonferroni corrected confidence intervals as given in the expression below:

$$\left(\bar{x}_j \pm t_{n-1, \alpha/2p} \sqrt{s_{jj}/n} \right)$$

- Asymptotic simultaneous intervals use the χ^2 approximation for large sample. If $n - p$ is large, then

$$\left(\mathbf{a}^T \bar{\mathbf{x}} \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{\mathbf{a}^T \mathbf{S} \mathbf{a}}{n}} \right)$$

contain $\mathbf{a}^T \boldsymbol{\mu}$ for every \mathbf{a} , with probability approximately $1 - \alpha$.

Lizard data example:

```
library(ICSNP) # For one and two sample Hotelling's T2 test (HotellingsT2)
```

```
## Warning: package 'ICSNP' was built under R version 3.4.4
```

```
## Loading required package: mvtnorm
```

```
## Warning: package 'mvtnorm' was built under R version 3.4.4
```

```
## Loading required package: ICS
```

```
## Warning: package 'ICS' was built under R version 3.4.3
```

```
library(ellipse) # For ellipse CI
```

```
## Warning: package 'ellipse' was built under R version 3.4.3
```

```
##
```

```
## Attaching package: 'ellipse'
```

```
## The following object is masked from 'package:graphics':
```

```
##
```

```
## pairs
```

```
#
```

```
# Import lizard data set
```

```
#
```

```
lizard <- read.table("Wichern_data/T1-3.dat")
```

```
names(lizard) <- c("Mass", "SVL", "HLS")
```

```
n <- dim(lizard)[1] # number of observations
```

```
p <- dim(lizard)[2] # number of variables
```

```
n
```

```
## [1] 25
```

```
p
```

```
## [1] 3
```

```
#
```

```
# Do the one sample hotelling T2-test manually
```

```
#
```

One sample Hotelling T2 test

$X = \begin{bmatrix} \text{mass} \\ \text{SVL} \\ \text{HLS} \end{bmatrix}$

we want to test if

$$\mu_x = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \end{bmatrix} \neq \mu_0 = \begin{bmatrix} 0 \\ 68 \\ 129 \end{bmatrix}$$

* One sample T² test manually

Do the one sample Hotelling T²-test manually

```
alpha <- 0.05

liz.xbar <- sapply(lizard, mean)
liz.cov <- cov(lizard)

liz.xbar

##      Mass      SVL      HLS
## 8.6866 68.4000 129.3200

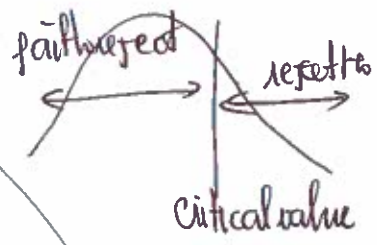
liz.cov

##      Mass      SVL      HLS
## Mass 7.186551 20.69952 33.53684
## SVL 20.699521 63.77083 102.08542
## HLS 33.536842 102.08542 185.83083

solve(liz.cov) # S^{-1}

##      Mass      SVL      HLS
## Mass 2.17163898 -0.64273969 -0.03882867
## SVL -0.64273969 0.32026006 -0.05993851
## HLS -0.03882867 -0.05993851 0.04531561

u0 <- c(9, 68, 129)
```



$$T2 <- n * mahalanobis(x = \bar{x}, center = \mu_0, cov = S) \quad T^2 = n (\bar{x} - \mu_0)^T S^{-1} (\bar{x} - \mu_0)$$

```
T2
## [1] 10.56927
crit.val <- p * (n-1) / (n-p) * qf(1-alpha, p, n-p)
```

$$\frac{(n-1)p}{(n-p)} F_{p, n-p, 1-\alpha} ?$$

crit.val # T2 > crit.val ⇒ reject H₀.

```
## [1] 9.978955
# -----
# Do the one sample hotelling T2-test using HotellingsT2
# -----
```

$$liz.T2 <- HotellingsT2(X = lizard, mu = u0) \quad liz.T2 = \frac{n-p}{p(n-1)} T2 = F \rightsquigarrow F_{p, n-p}$$

```
liz.T2

##
## Hotelling's one sample T2-test
##
## data: lizard
## T.2 = 3.2295, df1 = 3, df2 = 22, p-value = 0.04202
## alternative hypothesis: true location is not equal to c(9,68,129)
# Note that this is the rescale version of the above T2 statistic.
# So that liz.T2 = (n-p) / p / (n-1) * T2,
# where T2 is the above one covered in class
```

p-value < alpha ⇒ reject H₀

One sample # T² test using R command

$$a^T \bar{X} \pm \sqrt{\frac{p(n-1)}{n(n-p)} F_{L-\alpha, p, n-p} \underbrace{a^T S a}_{\text{diag of } S}} \quad \text{where } a_i = (0, \dots, \underset{\uparrow}{1}, \dots, 0)$$

```

# Simultaneous T2 intervals (Using the Hotelling T2 intervals)
# -----
ci.T2 <- rbind(liz.xbar - sqrt(p*(n-1)/(n-p)/n*qt(1-alpha, p, n-p)*diag(liz.cov)),
              liz.xbar + sqrt(p*(n-1)/(n-p)/n*qt(1-alpha, p, n-p)*diag(liz.cov)))
row.names(ci.T2)<-c("T2_L", "T2_U")

```

```

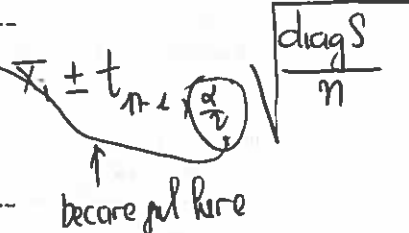
# -----
# One at a time t-intervals (the usual t-interval at level 1 - alpha
# with no correction for p = 3 variables). See page 229 - 231
# -----

```

```

ci.t <- rbind(liz.xbar - sqrt(diag(liz.cov)/n)*qt(1-alpha/2, n-1),
             liz.xbar + sqrt(diag(liz.cov)/n)*qt(1-alpha/2, n-1))
row.names(ci.t)<-c("t_L", "t_U")

```



```

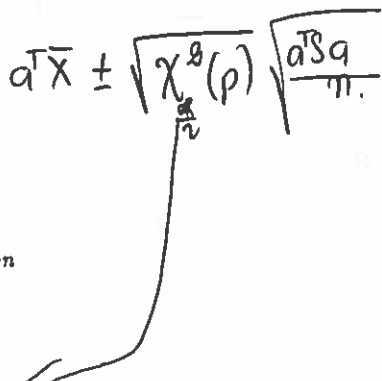
# -----
# Asmptotic Simultaneous intervals (Using the chisq approximation -
# i.e., the Hotelling T2 for large n) See page 224 - 235
# -----

```

```

ci.Asy <- rbind(liz.xbar - sqrt(qchisq(1-alpha/2,p)/n*diag(liz.cov)),
               liz.xbar + sqrt(qchisq(1-alpha/2,p)/n*diag(liz.cov)))
row.names(ci.Asy)<-c("Asy_L", "Asy_U")

```



```

# -----
# Bonferonni corrected (the usual t-interval with bonferroni correction
# i.e., level 1 - alpha / p) See page 232 - 234
# -----

```

```

ci.Bon <- rbind(liz.xbar - sqrt(diag(liz.cov)/n)*qt(1-alpha/2/p, n-1),
               liz.xbar + sqrt(diag(liz.cov)/n)*qt(1-alpha/2/p, n-1))
row.names(ci.Bon)<-c("Bon_L", "Bon_U")

```

```

# -----
# All together (_L indicates the lower endpoint and _U the upper endpoint)
# -----

```

```

all.ci <- round(rbind(ci.T2, ci.t, ci.Asy, ci.Bon), digits=2)
all.ci

```

##	Mass	SVL	HLS
## T2_L	6.99	63.35	120.71
## T2_U	10.38	73.45	137.93
## t_L	7.58	65.10	123.69

```

## t_U    9.79 71.70 134.95
## Asy_L  7.05 63.52 120.98
## Asy_U 10.33 73.28 137.66
## Bon_L  7.31 64.29 122.30
## Bon_U 10.07 72.51 136.34

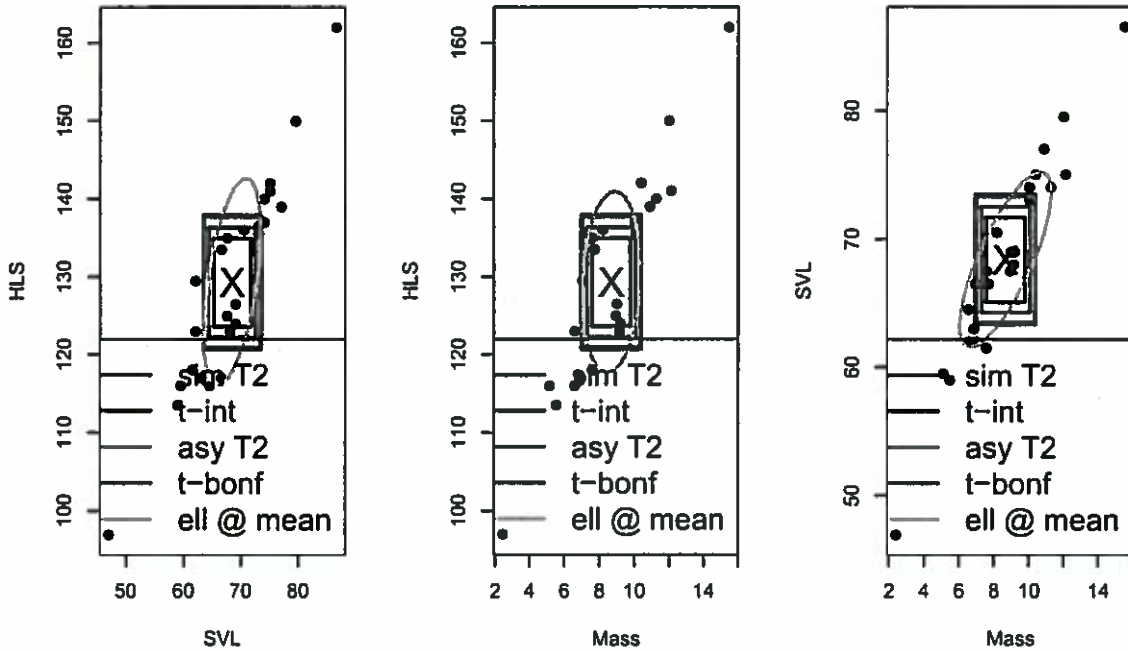
# -----
#
# Plot over scatter plots
#
# -----

old.par <- par(no.readonly = TRUE)

#dev.new(width = 15)
par(mfrow = c(1, p), oma = c(0, 0, 2, 0))
for(i in 1:p){
  plot(lizard[,-i], pch = 19)
  text(liz.xbar[-i][1], liz.xbar[-i][2], labels = "X", col = "black", cex = 2)
  lines(all.ci[,-i][c(1, 1, 2, 2, 1), 1], all.ci[,-i][c(1, 2, 2, 1, 1), 2],
        col = "red", lwd = 2)
  lines(all.ci[,-i][2+c(1, 1, 2, 2, 1), 1], all.ci[,-i][2+c(1, 2, 2, 1, 1), 2],
        col = "blue", lwd = 2)
  lines(all.ci[,-i][4+c(1, 1, 2, 2, 1), 1], all.ci[,-i][4+c(1, 2, 2, 1, 1), 2],
        col = "green", lwd = 2)
  lines(all.ci[,-i][6+c(1, 1, 2, 2, 1), 1], all.ci[,-i][6+c(1, 2, 2, 1, 1), 2],
        col = "purple", lwd = 2)
  lines(ellipse(x = solve(solve(liz.cov)[-i,-i]), centre = liz.xbar[-i],
                level = 1-alpha, npoints = 1000), lwd = 2, col = "orange")
  legend("bottomright", legend = c("sim T2", "t-int", "asy T2", "t-bonf", "ell @ mean"),
        col = c("red", "blue", "green", "purple", "orange"), lty = 1, lwd = 2, cex = 1.5)
}
title("2-dimensional view of CI's for means", outer = TRUE)

```

2-dimensional view of CI's for means



`#par(old.par)`

ell & mean is the ellipsoid evaluated at the mean of the third omitted variable. These ellipses may look small and misleading because they are not drawn in the direction of the principal components, but it gives an idea of the relative size and shape.

5. Review of ANOVA

One factor \Rightarrow one way ANOVA
 two factors \Rightarrow two-way ANOVA (not vector variables)

Analysis of variance (ANOVA) is a collection of statistical methodologies for comparing several means. When there is only one way (one factor) to classify the populations of interest, we use one-way ANOVA to analyze the data. When analyzing the effect of two factors, we use two-way ANOVA.

Suppose that $X_{11}, X_{12}, \dots, X_{1n_1}$ is a random sample from an $N(\mu_l, \sigma^2)$ population, $l = 1, 2, \dots, g$, and that the random samples are independent. We write

$$X_l = \mu + \tau_l$$

where

- X_l is the l^{th} group mean
- μ is the overall mean
- τ_l is the l^{th} group effect or treatment effect

The response $X_{lj} \sim N(\mu + \tau_l, \sigma^2)$ can be expressed as

$$X_{lj} = \mu + \tau_l + \epsilon_{lj}$$

where ϵ_{lj} are independent $N(0, \sigma^2)$ random variables and $\sum_{l=1}^g n_l \tau_l = 0$.

We decompose X_{lj} as

$$X_{lj} = \bar{X} + (\bar{X}_l - \bar{X}) + (X_{lj} - \bar{X}_l)$$

where

- X_{lj} : observation
- \bar{X} : overall sample mean, an estimate of μ
- $(\bar{X}_l - \bar{X})$: estimated treatment effect, an estimate of τ_l
- $(X_{lj} - \bar{X}_l)$: residual, an estimate of the error ϵ_{lj}

Recall decomposing $SST = SStr + SSR$

$$\underbrace{\sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X})^2}_{SST} = \underbrace{\sum_{l=1}^g n_l (\bar{X}_l - \bar{X})^2}_{SStr} + \underbrace{\sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X}_l)^2}_{SSR \text{ (within samples)}}$$

SStr within sample.

One-Way ANOVA table

Table 1: ANOVA Table for Comparing Univariate Population Means

Source	Sum of squares (SS)	df	MS	F
Treatments	$SStr = \sum_{l=1}^g n_l (X_l - \bar{X})^2$	$g - 1$	$SStr / (g - 1)$	$MStr / MSR$
Residuals	$SSR = \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X}_l)^2$	$\sum_{l=1}^g n_l - g$	$SSR / \sum_{l=1}^g n_l - 1$	
Total	$SST = \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X})^2$	$\sum_{l=1}^g n_l - 1$		

The usual F-test rejects $H_0 : \tau_l = 0, l = 1, \dots, g$ at level α if

$$F = \frac{SStr / (g - 1)}{SSR / (\sum_{l=1}^g n_l - g)} > F_{g-1, \sum_{l=1}^g n_l - g}(\alpha)$$

where $F_{g-1, \sum_{l=1}^g n_l - g}(\alpha)$ is the upper $\alpha \times 100\%$ of the F-distribution with df $g - 1$ and $\sum_{l=1}^g n_l - g$.

6. Multivariate Analysis of variance (MANOVA)

As with the univariate situation, t-tests are fine for comparing the means of two groups, but we would have "multiple comparison" problems if we tried to compare more than two. In an analogous way, in the multivariate context we have MANOVA.

The response $\mathbf{X}_{lj} \sim N_p(\mu + \tau_l, \Sigma)$ can be expressed as

$$\mathbf{X}_{lj} = \mu + \tau_l + \epsilon_{lj}$$

where ϵ_{lj} are independent $N_p(\mathbf{0}, \Sigma)$ random variables and $\sum_{l=1}^g n_l \tau_l = \mathbf{0}$.

MANOVA table

Table 2: MANOVA Table for Comparing Population Mean Vectors

Source	Sum of squares (SS)	df
Treatments	$\mathbf{B} = \sum_{l=1}^g n_l (\mathbf{X}_l - \bar{\mathbf{X}})(\mathbf{X}_l - \bar{\mathbf{X}})^T$	$g - 1$
Residuals	$\mathbf{W} = \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{X}_{lj} - \bar{\mathbf{X}}_l)(\mathbf{X}_{lj} - \bar{\mathbf{X}}_l)^T$	$\sum_{l=1}^g n_l - g$
Total	$\mathbf{B} + \mathbf{W} = \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{X}_{lj} - \bar{\mathbf{X}})(\mathbf{X}_{lj} - \bar{\mathbf{X}})^T$	$\sum_{l=1}^g n_l - 1$

matrix

$$\mathbf{W} = \sum^T \sum \quad \text{where } S \approx \mathbf{V}$$

test $> F \Rightarrow$

Test of $H_0 : \tau_l = 0, l = 1, \dots, g$ involves generalized variance

$$\Lambda^* = \frac{|W|}{|B+W|}$$

is too small. Λ^* is referred as Wilk's lambda. Let $\lambda_1, \dots, \lambda_r$

$$\Lambda^* = \frac{|I|}{|BW^{-1} + I|} =$$

Assumptions and Limitations

The following assumptions are made when using a MANOVA.

1. The response variables are continuous.
2. The residuals follow the multivariate-normal probability distribution with means equal to zero.
3. The individuals are independent.
4. The variance-covariance matrices of each group of residuals are equal. MANOVA makes the assumption that the within-cell (group) covariance matrices are equal. If the design is balanced so that there is an equal number of observations in each cell, the robustness of the MANOVA tests is guaranteed. If the design is unbalanced, you should test the equality of covariance matrices using Box's M test.

$$H_0 : \Sigma_1 = \Sigma_2 = \dots = \Sigma_g$$

Let

$$\Lambda_{Box} = \prod_{l=1}^g \left(\frac{|S_l|}{|S_{pooled}|} \right)^{(n_l-1)/2}$$

where

$$S_{pooled} = \frac{1}{\sum_l (n_l - 1)} \{ (n_1 - 1)S_1 + (n_2 - 1)S_2 + \dots + (n_g - 1)S_g \}$$

Let

$$M = -2 \ln \Lambda = \left[\sum_l (n_l - 1) \right] \ln |S_{pooled}| - \sum_l [(n_l - 1) \ln |S_l|]$$

Set

$$u = \left[\sum_l \frac{1}{n_l - 1} - \frac{1}{\sum_l (n_l - 1)} \right] \left[\frac{2p^2 + 3p - 1}{6(p+1)(g-1)} \right]$$

where p is the number of variables and g is the number of groups. Then

$$C = (1 - u)M = (1 - u) \left\{ \left[\sum_l (n_l - 1) \right] \ln |S_{pooled}| - \sum_l [(n_l - 1) \ln |S_l|] \right\}$$

has an approximate χ^2 distribution with degrees of freedom

$$\nu = \frac{1}{2}gp(p+1) - \frac{1}{2}p(p+1) = \frac{1}{2}p(p+1)(g-1).$$

We reject H_0 if $C > \chi_{p(p+1)(g-1)/2}^2(\alpha)$. Box's χ^2 approximation worked well if $n_l \leq 20$ and $p, g \leq 5$. Box's M test is sensitive to some forms of non-normality.

* Example Film quality

X_1	X_2	X_3	rate
tear	gloss	opacity	low
			high

Example

First of all we'll enter some data relating to the production of plastic film reported in Krzanowski (2000). Tear, gloss and opacity are measures of the manufactured films.

```
tear <- c(6.5, 6.2, 5.8, 6.5, 6.5, 6.9, 7.2, 6.9, 6.1, 6.3,
6.7, 6.6, 7.2, 7.1, 6.8, 7.1, 7.0, 7.2, 7.5, 7.6)
gloss <- c(9.5, 9.9, 9.6, 9.6, 9.2, 9.1, 10.0, 9.9, 9.5, 9.4,
9.1, 9.3, 8.3, 8.4, 8.5, 9.2, 8.8, 9.7, 10.1, 9.2)
opacity <- c(4.4, 6.4, 3.0, 4.1, 0.8, 5.7, 2.0, 3.9, 1.9, 5.7,
2.8, 4.1, 3.8, 1.6, 3.4, 8.4, 5.2, 6.9, 2.7, 1.9)
Y <- cbind(tear, gloss, opacity)
```

We now need to put in information on the rate of extrusion, and the amount of additive used (gl() is a command which specifically creates these kind of experimental factors).

```
rate <- factor(gl(2,10), labels=c("Low", "High"))
additive <- factor(gl(2, 5, len=20), labels=c("Low", "High"))
```

There are three conventional ANOVA that could be considered here, but to consider the three responses together we may wish to conduct a MANOVA. However, we can use manova() to fit the multivariate ANOVA, and use summary.aov() to extract the results of the univariate analyses. There are three matrices of interest in MANOVA:

- Between-group B
- Within-group W
- Total B + W

Two way

$$X_{ij} = \mu + \tau_i + \beta_j + \text{interaction} + \epsilon_{ij}$$

Wilk's Lambda is the ratio $\Lambda^* = \frac{|W|}{|B+W|}$

```
fit <- manova(Y ~ rate * additive)
summary.aov(fit) # univariate ANOVA tables
```

```
## Response tear :
##           Df Sum Sq Mean Sq F value Pr(>F)
## rate      1  1.7405  1.74050  15.7868 0.001092 **
## additive  1  0.7605  0.76050   6.8980 0.018330 *
## rate:additive 1  0.0005  0.00050   0.0045 0.947143
## Residuals 16  1.7640  0.11025
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Response gloss :
##           Df Sum Sq Mean Sq F value Pr(>F)
## rate      1  1.3005  1.30050   7.9178 0.01248 *
## additive  1  0.6125  0.61250   3.7291 0.07139 .
## rate:additive 1  0.5445  0.54450   3.3151 0.08740 .
## Residuals 16  2.6280  0.16425
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Response opacity :
##           Df Sum Sq Mean Sq F value Pr(>F)
## rate      1  0.421  0.4205  0.1036 0.7517
## additive  1  4.901  4.9005  1.2077 0.2881
## rate:additive 1  3.960  3.9605  0.9760 0.3379
```

In this example, we have two cases for factors

- rate & additive have interaction
- rate & additive don't have interaction

* With problem about Egyptian skulls we have 4 variable

- X_1 breadth of skull
- X_2 basibremath of sk.
- X_3 brainpanerlar
- X_4 nasal

with factors are time periods

- 4000 BC
- 3300 BC
- 1850 BC

⇒ they don't have interaction

* Care about interaction

$p\text{value} < \alpha$ the action is significant \Rightarrow affect

$p\text{value} > \alpha$ the interaction is not significant \Rightarrow we

```
## Residuals      16 64.924  4.0578
```

A call to `summary()` will give the MANOVA table.

```
summary(fit, test="Wilks") # ANOVA table of Wilks' lambda
```

```
##           Df  Wilks approx F num Df den Df  Pr(>F)
## rate      1 0.38186  7.5543      3   14 0.003034 **
## additive  1 0.52303  4.2556      3   14 0.024745 *
## rate:additive 1 0.77711  1.3385      3   14 0.301782
## Residuals  16
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

As with Hotellings T^2 , Wilk's Lambda has to be converted into an F statistic, a calculation that has been done by the software.

In any case, the interaction term is not significant. We can fit the model without interactions, which as anticipated suggests that both additive and extrusion rate have an effect on the outcome measures. We will use `by()` to examine the various group means.

```
fit <- manova(Y ~ rate + additive)
summary(fit, test = "Wilks")
```

(care about rate and additive without caring about the interaction)

```
##           Df  Wilks approx F num Df den Df  Pr(>F)
## rate      1 0.38684  7.9253      3   15 0.00212 **
## additive  1 0.55384  4.0279      3   15 0.02753 *
## Residuals 17
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
by(Y, rate, summary) ## group means according to extrusion rate
```

```
## INDICES: Low
```

```
##      tear      gloss      opacity
## Min. :5.800  Min. : 9.100  Min. :0.800
## 1st Qu.:6.225  1st Qu.: 9.425  1st Qu.:2.250
## Median :6.500  Median : 9.550  Median :4.000
## Mean   :6.490  Mean   : 9.570  Mean   :3.790
## 3rd Qu.:6.800  3rd Qu.: 9.825  3rd Qu.:5.375
## Max.   :7.200  Max.   :10.000  Max.   :6.400
```

```
## INDICES: High
```

```
##      tear      gloss      opacity
## Min. :6.60  Min. : 8.300  Min. :1.600
## 1st Qu.:6.85  1st Qu.: 8.575  1st Qu.:2.725
## Median :7.10  Median : 9.150  Median :3.600
## Mean   :7.08  Mean   : 9.060  Mean   :4.080
## 3rd Qu.:7.20  3rd Qu.: 9.275  3rd Qu.:4.925
## Max.   :7.60  Max.   :10.100  Max.   :8.400
```

```
by(Y, additive, summary) ## group means according to additive
```

```
## INDICES: Low
```

```
##      tear      gloss      opacity
## Min. :5.800  Min. :8.300  Min. :0.80
## 1st Qu.:6.500  1st Qu.:8.650  1st Qu.:2.85
## Median :6.550  Median :9.250  Median :3.60
## Mean   :6.590  Mean   :9.140  Mean   :3.44
```

```
## 3rd Qu.:6.775 3rd Qu.:9.575 3rd Qu.:4.10
## Max. :7.200 Max. :9.900 Max. :6.40
```

```
## -----
```

```
## INDICES: High
```

```
##      tear      gloss      opacity
## Min. :6.10 Min. : 8.80 Min. :1.900
## 1st Qu.:6.90 1st Qu.: 9.20 1st Qu.:2.175
## Median :7.05 Median : 9.45 Median :4.550
## Mean :6.98 Mean : 9.49 Mean :4.430
## 3rd Qu.:7.20 3rd Qu.: 9.85 3rd Qu.:5.700
## Max. :7.60 Max. :10.10 Max. :8.400
```

```
by(Y, list(rate,additive), summary) ## group means by both.
```

```
## : Low
```

```
## : Low
```

```
##      tear      gloss      opacity
## Min. :5.8 Min. :9.20 Min. :0.80
## 1st Qu.:6.2 1st Qu.:9.50 1st Qu.:3.00
## Median :6.5 Median :9.60 Median :4.10
## Mean :6.3 Mean :9.56 Mean :3.74
## 3rd Qu.:6.5 3rd Qu.:9.60 3rd Qu.:4.40
## Max. :6.5 Max. :9.90 Max. :6.40
```

```
## -----
```

```
## : High
```

```
## : Low
```

```
##      tear      gloss      opacity
## Min. :6.60 Min. :8.30 Min. :1.60
## 1st Qu.:6.70 1st Qu.:8.40 1st Qu.:2.80
## Median :6.80 Median :8.50 Median :3.40
## Mean :6.88 Mean :8.72 Mean :3.14
## 3rd Qu.:7.10 3rd Qu.:9.10 3rd Qu.:3.80
## Max. :7.20 Max. :9.30 Max. :4.10
```

```
## -----
```

```
## : Low
```

```
## : High
```

```
##      tear      gloss      opacity
## Min. :6.10 Min. : 9.10 Min. :1.90
## 1st Qu.:6.30 1st Qu.: 9.40 1st Qu.:2.00
## Median :6.90 Median : 9.50 Median :3.90
## Mean :6.68 Mean : 9.58 Mean :3.84
## 3rd Qu.:6.90 3rd Qu.: 9.90 3rd Qu.:5.70
## Max. :7.20 Max. :10.00 Max. :5.70
```

```
## -----
```

```
## : High
```

```
## : High
```

```
##      tear      gloss      opacity
## Min. :7.00 Min. : 8.8 Min. :1.90
## 1st Qu.:7.10 1st Qu.: 9.2 1st Qu.:2.70
## Median :7.20 Median : 9.2 Median :5.20
## Mean :7.28 Mean : 9.4 Mean :5.02
## 3rd Qu.:7.50 3rd Qu.: 9.7 3rd Qu.:6.90
## Max. :7.60 Max. :10.1 Max. :8.40
```

High levels of extrusion rate lead to higher levels of tear and opacity but lower levels of gloss. High levels of

additive lead to higher levels of tear, gloss and opacity.

Multiple Comparisons in MANOVA

If the F-test rejects the null that all g groups have equal mean vectors, we can use multiple comparisons to determine which specific pairs of groups have differing mean vectors (and which components of those mean vectors differ). This is most simply done via Bonferroni confidence intervals for all the differences $\tau_{ki} - \tau_{li}$ for all components $i = 1, \dots, p$ and all difference $l < k = 1, \dots, g$. With family confidence level $1 - \alpha$, the set of confidence intervals defined by

$$\bar{x}_{ki} - \bar{x}_{li} \pm t_{\frac{\alpha}{pg(g-1)}} \sqrt{\frac{w_{ii}}{n-g} \left(\frac{1}{n_k} + \frac{1}{n_l} \right)}$$

probably W in table 2.
we need
qt(probability degree, lower.tail = FALSE, log.p = FALSE)

Here w_{ii} is the i^{th} diagonal element of w . Any such interval that does not contain zero would indicate a significant difference in that component of the mean vector between that pair of groups.

4 variables | 1 Factor

Example: Manova on three groups, using Fisher's classic Iris data

(cond, independent, normality)

This dataset consists of 50 cases of each of 3 species, namely Iris setosa, Iris virginica, and Iris versicolor. Each case has 4 measurements on the length and width of its petals and sepals.

```
attach(iris) variables: length, width, petal, sepal ~ over 3 species (3 groups)
iris.manova <- manova(cbind(Sepal.Length, Sepal.Width, Petal.Length, Petal.Width) ~ Species)
```

Doing the MANOVA test using Wilks' Lambda:

```
summary(iris.manova, test="Wilks") by default
##           Df      Wilks approx F num Df den Df      Pr(>F)
## Species    2 0.023439   199.15     8    288 < 2.2e-16 ***
## Residuals 147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the other test statistics:

```
summary(iris.manova, test="Roy")
##           Df      Roy approx F num Df den Df      Pr(>F)
## Species    2 32.192    1167     4    145 < 2.2e-16 ***
## Residuals 147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

⇒ Reject H₀

```
summary(iris.manova, test="Hotelling-Lawley")
##           Df Hotelling-Lawley approx F num Df den Df      Pr(>F)
## Species    2   32.477    580.53     8    286 < 2.2e-16 ***
## Residuals 147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(iris.manova, test="Pillai")
##           Df Pillai approx F num Df den Df      Pr(>F)
## Species    2  1.1919    53.466     8    290 < 2.2e-16 ***
```

```
## Residuals 147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

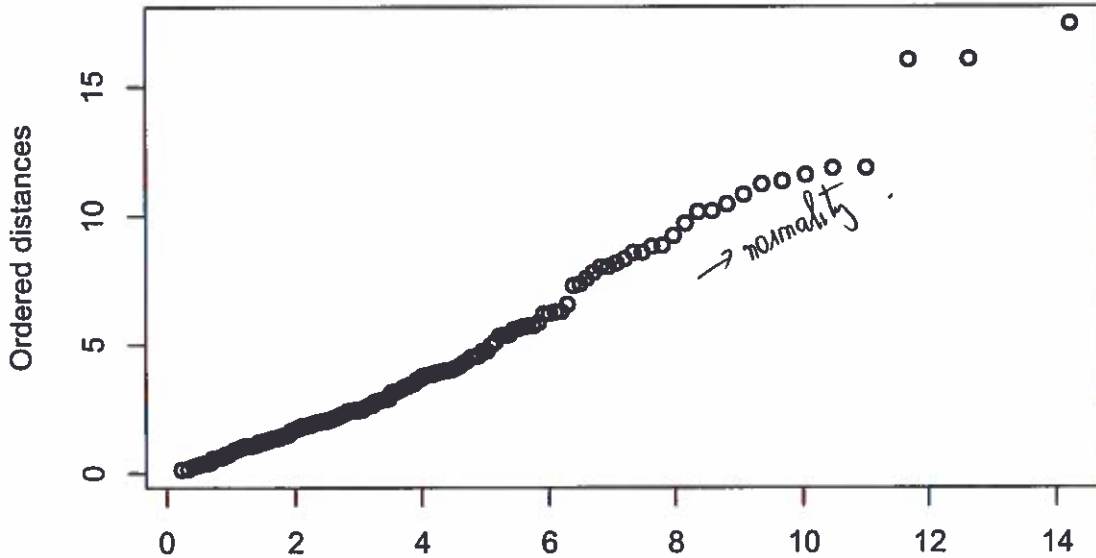
We have strong evidence that the mean vectors differ across the 3 species.

Checking model assumption of normality:

```
chisplot <- function(x) {
  if (!is.matrix(x)) stop("x is not a matrix")
  ### determine dimensions
  n <- nrow(x)
  p <- ncol(x)
  xbar <- apply(x, 2, mean)
  S <- var(x)
  S <- solve(S)
  index <- (1:n)/(n+1)
  xcent <- t(t(x) - xbar)
  di <- apply(xcent, 1, function(x,S) x %*% S %*% x,S)
  quant <- qchisq(index,p)
  plot(quant, sort(di), ylab = "Ordered distances",
       xlab = "Chi-square quantile", lwd=2,pch=1)
}
```

to satisfy independent, normality, variance unequal or equal, continuous. There is QQ plot function in R that we can use

```
chisplot(residuals(iris.manova))
```



Check SE - V equal

Chi-square quantile

There is no strong evidence against normality - we are safe. Now we examine the sample covariance matrices for each group:

```
by(iris[, -5], Species, var)
```

```
## Species: setosa
##           Sepal.Length Sepal.Width Petal.Length Petal.Width
## Sepal.Length  0.12424898 0.099216327  0.016355102  0.010330612
```



```

## Sepal.Width      0.09921633 0.143689796 0.011697959 0.009297959
## Petal.Length     0.01635510 0.011697959 0.030159184 0.006069388
## Petal.Width      0.01033061 0.009297959 0.006069388 0.011106122
## -----
## Species: versicolor
##           Sepal.Length Sepal.Width Petal.Length Petal.Width
## Sepal.Length  0.26643265 0.08518367 0.18289796 0.05577959
## Sepal.Width   0.08518367 0.09846939 0.08265306 0.04120408
## Petal.Length  0.18289796 0.08265306 0.22081633 0.07310204
## Petal.Width   0.05577959 0.04120408 0.07310204 0.03910612
## -----
## Species: virginica
##           Sepal.Length Sepal.Width Petal.Length Petal.Width
## Sepal.Length  0.40434286 0.09376327 0.30328980 0.04909388
## Sepal.Width   0.09376327 0.10400408 0.07137959 0.04762857
## Petal.Length  0.30328980 0.07137959 0.30458776 0.04882449
## Petal.Width   0.04909388 0.04762857 0.04882449 0.07543265

```

An R function to test the equality of several covariance matrices:

Need to install the 'asbio' package from the Internet:

```
library(asbio) # load the 'asbio' package once it is installed
```

```
## Warning: package 'asbio' was built under R version 3.4.4
```

```
## Loading required package: tcltk
```

```
Kullback(Y=iris[,-5],X=Species)
```

```
##
```

```
## Kullback test for equal covariance matrices
```

```
##           Chi* df  P(Chi>Chi*)
```

```
## setosa 73.33162 20 5.158075e-08
```

Another R function to test the equality of several covariance matrices: Need to install the 'biotools' package from the Internet:

```
library(biotools) # load the 'biotools' package once it is installed
```

```
## Loading required package: rpanel
```

```
## Package `rpanel', version 1.1-3: type help(rpanel) for summary information
```

```
## Loading required package: tkrplot
```

```
## Warning: package 'tkrplot' was built under R version 3.4.4
```

```
## Loading required package: MASS
```

```
## Loading required package: lattice
```

```
## Loading required package: SpatialEpi
```

```
## Warning: package 'SpatialEpi' was built under R version 3.4.4
```

```
## Loading required package: sp
```

```
## ---
```

```
## biotools version 3.1
```

```
##
```

X quant 2

Use the box M test to check the equality of cov matrices

```
boxM(iris[, -5], Species)
```

```
##
## Box's M-test for Homogeneity of Covariance Matrices
##
## data: iris[, -5]
## Chi-Sq (approx.) = 140.94, df = 20, p-value < 2.2e-16
The covariance matrices may not be equal across groups.
```

Differences in Means for doing Bonferroni multiple comparisons:

```
means.by.grps <- cbind(tapply(Sepal.Length, Species, mean),
                        tapply(Sepal.Width, Species, mean),
                        tapply(Petal.Length, Species, mean),
                        tapply(Petal.Width, Species, mean) )
```

$$\sum_{150 \times 4}$$

The matrix W:

```
my.n <- nrow(iris[, -5])
W <- (my.n - 1) * var(residuals(iris.manova))
```

need W

$$\text{Var}(\text{residuals}) = \frac{\sum^T \sum}{n-1}$$

We want $W = \sum^T \sum$

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
## Sepal.Length	38.9562	13.6300	24.6246	5.6450
## Sepal.Width	13.6300	16.9620	8.1208	4.8084
## Petal.Length	24.6246	8.1208	27.2226	6.2718
## Petal.Width	5.6450	4.8084	6.2718	6.1566

Calculation of the Bonferroni CIs:

```
my.alpha <- 0.10 # family confidence level for CIs is 90% here
```

```
my.m <- 3 # number of groups
my.q <- 4 # number of variables
group.sample.sizes <- c(50,50,50)
```

```
for (k in 1:my.q) {
```

```
pair.mean.diffs <- cbind( t(combn(my.m,2)),
                        combn(tapply(iris[,k], Species, mean), 2, FUN=diff) )
```

```
t.val <- qt(my.alpha/(my.q*my.m*(my.m-1)), df=my.n-my.m, lower=F)
```

```
#print(paste("For i=", i, "i.prime=", j, "and k=", k, "t est is"
```

```
CI.L <- pair.mean.diffs[,3] -
t.val*sqrt((diag(W)[k]/(my.n-my.m))*
(1/group.sample.sizes[pair.mean.diffs[,1]] +
1/group.sample.sizes[pair.mean.diffs[,2]])) )
```

```
CI.U <- pair.mean.diffs[,3] +
t.val*sqrt((diag(W)[k]/(my.n-my.m))*
(1/group.sample.sizes[pair.mean.diffs[,1]] +
1/group.sample.sizes[pair.mean.diffs[,2]])) )
```

```
my.table.mat <- cbind(pair.mean.diffs, round(CI.L, 3), round(CI.U, 3),
rep(k, times=nrow(pair.mean.diffs))) )
```

```
my.table<-as.data.frame(my.table.mat)
names(my.table)=c('grp1','grp2','diff.samp.means',
                  'lower.CI','upper.CI','variable');
print(my.table)
}
```

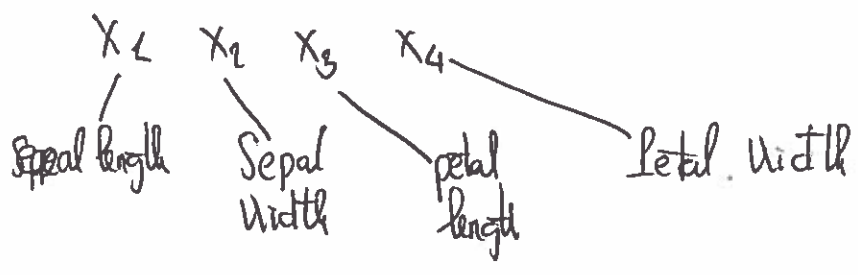
Compare groups

	grp1	grp2	diff.samp.means	lower.CI	upper.CI	variable
X ₁ sepal length	1	2	0.930	0.655	1.205	1
	1	3	1.582	1.307	1.857	1
	2	3	0.652	0.377	0.927	1
X ₂	1	2	-0.658	-0.840	-0.476	2
	1	3	-0.454	-0.636	-0.272	2
	2	3	0.204	0.022	0.386	2
X ₃	1	2	2.798	2.568	3.028	3
	1	3	4.090	3.860	4.320	3
	2	3	1.292	1.062	1.522	3
X ₄	1	2	1.08	0.971	1.189	4
	1	3	1.78	1.671	1.889	4
	2	3	0.70	0.591	0.809	4

note not contain 0 → reject, there is no different
 → different between μ_1 & μ_2

μ_{X_1, F_1} μ_{X_1, F_2} $\neq 0$
 μ_{X_1, F_1} μ_{X_1, F_3}
 μ_{X_1, F_2} μ_{X_1, F_3}

lower CI and upper CI of the $\bar{\mu}_i - \hat{\mu}_i$



1 Factor	
Iris setosa	1
Iris virginica	2
Iris versicolor	3



18



19



we want to find β

Chapter 3. Multivariate Regression

Jianxuan Liu
Fall 2018

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & \dots & X_{1r} \\ 1 & X_{21} & & X_{2r} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & & X_{nr} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_r \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$n \times 1$ $n \times (r+1)$ $(r+1) \times 1$ $n \times 1$

1. Classical Linear Regression Model

$$\tilde{Y} = X\beta + \epsilon$$

where X is a $n \times (r + 1)$ data matrix, $\beta = [\beta_0, \beta_1, \dots, \beta_r]$ unknown parameters, Y is a $n \times 1$ vector and the error terms have the following properties:

- $E(\epsilon) = 0$
- $Cov(\epsilon) = E(\epsilon\epsilon^T) = \sigma^2 I$

The least square estimate

Theorem 1. Let X have full rank $r + 1 \leq n$. The least squares estimate of β is

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{\beta}_{LS} = \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta)$$

Let $\hat{Y} = X\hat{\beta} = HY$ denote the fitted values of Y , where $H = X(X^T X)^{-1} X^T$. Then the residuals

$$\hat{\epsilon} = Y - \hat{Y} = [I - X(X^T X)^{-1} X^T] Y = [I - H] Y$$

satisfy $X^T \hat{\epsilon} = 0$ and $\hat{Y}^T \hat{\epsilon} = 0$. $\hat{\epsilon}$ is real and \hat{Y} is approximated.

Theorem 2. The least squares estimator $\hat{\beta} = (X^T X)^{-1} X^T Y$ is unbiased, i.e.

$$E(\hat{\beta}) = \beta$$

and

$$Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

The residuals $\hat{\epsilon}$ have the properties

- $E(\hat{\epsilon}) = 0$
- $Cov(\hat{\epsilon}) = \sigma^2 [I - H]$
- $E(\hat{\epsilon}^T \hat{\epsilon}) = (n - r - 1) \sigma^2$
- $E(s^2) = \sigma^2$ where $s^2 = \frac{\hat{\epsilon}^T \hat{\epsilon}}{n - r - 1} = \frac{Y^T [I - H] Y}{n - r - 1}$
- $\hat{\beta}$ and $\hat{\epsilon}$ are uncorrelated.

Theorem 3. Given that X has full rank $r + 1$ and $\epsilon \sim N_n(0, \sigma^2 I)$, the maximum likelihood estimator of β is the same as the least squares estimator, $\hat{\beta} = (X^T X)^{-1} X^T Y$. Moreover, $\hat{\beta} \sim N_{r+1}(\beta, \sigma^2 (X^T X)^{-1})$ and $\hat{\beta}$ is independent of the residuals $\hat{\epsilon} = Y - \hat{Y}$. Further,

$$n\hat{\sigma}^2 = \hat{\epsilon}^T \hat{\epsilon} \sim \sigma^2 \chi_{n-r-1}^2$$

where $\hat{\sigma}^2$ is the maximum likelihood estimator of σ^2 .

Theorem 4. The simultaneous $100(1 - \alpha)\%$ confidence region for β_i are

$$\hat{\beta}_i \pm \sqrt{(r + 1) F_{r+1, n-r-1}(\alpha) Var(\hat{\beta}_i)}, i = 0, 1, \dots, r$$

For a new observation $\mathbf{x}_0^T = [1, x_{01}, \dots, x_{0r}]$, the estimate of the expected values of Y_0 is $\hat{Y}_0 = \mathbf{x}_0^T \hat{\beta}$. If the error terms are normally distributed, then a $100(1 - \alpha)\%$ confidence interval of $E(Y_0 | \mathbf{x}_0)$ is

$$\mathbf{x}_0^T \hat{\beta} \pm t_{n-r-1, \alpha/2} \sqrt{(\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0) s^2}$$

A $100(1 - \alpha)\%$ prediction interval of Y_0 is

$$\mathbf{x}_0^T \hat{\beta} \pm t_{n-r-1, \alpha/2} \sqrt{(1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0) s^2}$$

2. Multivariate Regression

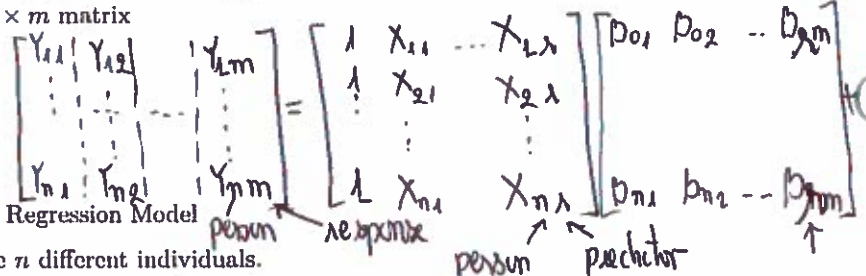
In multivariate regression we wish to predict or explain a set of m response (or dependent) variables Y_1, \dots, Y_m via a set of r predictor (or independent) variables x_1, \dots, x_r . The classical linear regression model can be written in matrix-vector form as

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

$n \times m$ $n \times (\lambda + 1)$ $(\lambda + 1) \times m$ $n \times m$

where

- \mathbf{x} is a $n \times (r + 1)$ data matrix
- $\beta = [\beta_{(1)}; \beta_{(2)}; \dots; \beta_{(m)}]$
- $\mathbf{Y} = [\mathbf{Y}_{(1)}; \mathbf{Y}_{(2)}; \dots; \mathbf{Y}_{(m)}]$ is a $n \times m$ matrix
- $\epsilon = [\epsilon_{(1)}; \epsilon_{(2)}; \dots; \epsilon_{(m)}]$
- $E(\epsilon_{(i)}) = \mathbf{0}$
- $Cov(\epsilon_{(i)}, \epsilon_{(j)}) = E(\epsilon_{(i)} \epsilon_{(j)}^T) = \sigma_{ij}^2 \mathbf{I}$



Further Explanation of the Multivariate Regression Model

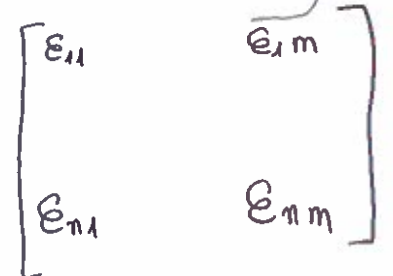
- The n rows of \mathbf{Y} correspond to the n different individuals.
- The m columns of \mathbf{Y} correspond to the m different response variables.
- Note that the first row of β is a row of intercept terms corresponding to the m response variables.
- Then the $(i + 1, j)$ entry of β measures the marginal effect of the i -th predictor variable on the j -th response variable.

Theorem 5.: The least squares estimate of β is

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

In particular, for the i^{th} response, $\hat{\beta}_{(i)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_{(i)}$

- $E(\hat{\beta}_{(i)}) = \beta_{(i)}$, or $E(\hat{\beta}) = \beta$
- $Cov(\hat{\beta}_{(i)}, \hat{\beta}_{(k)}) = \sigma_{ik} (\mathbf{X}^T \mathbf{X})^{-1}$, $i, k = 1, 2, \dots, m$
- $E(\hat{\epsilon}) = \mathbf{0}$ and $E(\hat{\epsilon}^T \hat{\epsilon}) = (n - r - 1) \Sigma$



Theorem 6.: Let $\mathbf{P} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{I}_n - \mathbf{H}$, then the maximum likelihood estimators are $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ and $\hat{\Sigma} = \frac{1}{n} \mathbf{Y}^T \mathbf{P} \mathbf{Y}$.

Proof:

Assume $\text{rank}(\mathbf{X}) = r$ and $n > m + r$.

$$L(\mathbf{Y}; \beta, \Sigma) = \prod_{i=1}^n |2\pi\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(y_i - \mathbf{x}_i\beta)\Sigma^{-1}(y_i - \mathbf{x}_i\beta)^T\right\}$$

implies

$$l(\mathbf{Y}; \beta, \Sigma) = -\frac{n}{2} \log |2\pi\Sigma| - \frac{1}{2} \text{tr}\{(\mathbf{Y} - \mathbf{X}\beta)\Sigma^{-1}(\mathbf{Y} - \mathbf{X}\beta)^T\}.$$

Write

$$l(\beta, \Sigma) = -\frac{n}{2} \log |2\pi\Sigma| - \frac{1}{2} \text{tr}\{\Sigma^{-1}\hat{\Sigma}\} - \frac{1}{2} \text{tr}\{\Sigma^{-1}(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta)\}.$$

Here, $(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta)$ is independent of Σ . $l(\beta, \Sigma)$ is maximized when $(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta)$ is minimized which is at $\beta = \hat{\beta}$. Plugging this into the log-likelihood we get

$$l(\hat{\beta}, \Sigma) = -\frac{np}{2} \log |2\pi| - \frac{n}{2} (\log |\Sigma| + \text{tr}\{\Sigma^{-1}\hat{\Sigma}\})$$

which is maximized at $\Sigma = \hat{\Sigma}$. Note that $(\hat{\beta}, \hat{\Sigma})$ is sufficient for (β, Σ) .

Theorem 7.: $\hat{\beta}$ and $\hat{\Sigma}$ are independent.

Proof*:

$$\hat{\Sigma} = \frac{1}{n} \mathbf{Y}^T \mathbf{P} \mathbf{Y} = \frac{1}{n} (\mathbf{Y} - \beta \mathbf{X})^T \mathbf{P} (\mathbf{Y} - \beta \mathbf{X})$$

$$\hat{\beta} - \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} (\mathbf{Y} - \beta \mathbf{X})$$

where $\mathbf{Y} - \beta \mathbf{X} \sim N_m(0, \Sigma)$. Note that $\mathbf{P} \mathbf{Y}$ and $\hat{\beta} - \beta$ are independent because $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{P} = \mathbf{0}$.

The Multivariate Regression Model Assumptions

- We assume that all of the nm elements of ϵ have mean 0.
- Any single row of ϵ has covariance matrix Σ (generally non-diagonal). This implies that the response variables within an individual multivariate observation may be correlated.
- However, we also assume that response values from different individuals are uncorrelated.
- For doing inference about the multivariate regression model, we further assume that each column of ϵ has a multivariate normal distribution.

Likelihood Ratio Tests for Regression Parameters

We may we may wish to test whether a set of several predictor variables is not related to the set of response variables.

$$H_0 : \beta_{(2)} = \mathbf{0}$$

where $\beta = [\beta_{(1)}; \beta_{(2)}]^T$. $\beta_{(1)}$ is the first $q + 1$ of the predictors are related to the set of response variables, and the last $r - q$ are useless in predicting the set of response variables. The likelihood ratio Λ can be expressed in terms of generalized variances

$$\Lambda = \frac{\max_{\beta_{(1)}, \Sigma} L(\beta_{(1)}, \Sigma)}{\max_{\beta, \Sigma} L(\beta, \Sigma)} = \frac{L(\hat{\beta}_{(1)}, \hat{\Sigma})}{L(\hat{\beta}, \hat{\Sigma})} = \left(\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|} \right)^{n/2}$$

Equivalently, Wilk's lambda statistic is

$$\Lambda^{2/n} = \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|}$$

Under H_0 , $Y = X_1\beta_{(1)} + \epsilon$ and the likelihood ratio test of H_0 is equivalent to rejecting H_0 for large values of

$$-2 \ln \Lambda = -n \ln \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|} = -n \ln \frac{|n\hat{\Sigma}|}{|n\hat{\Sigma} + n(\hat{\Sigma} - \hat{\Sigma}_1)|}$$

For n larger, the modified statistic

$$- \left[n - r - 1 - \frac{1}{2}(m - r + q + 1) \right] \ln \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|}$$

can be approximated as a $\chi^2_{m(r-q)}$.

Example 1: The Computer Data

```
comput <- read.table("computerdata.txt", header=T)
#comput <- read.table(file.choose(), header=T)
```

```
attach(comput)
```

```
# y1 = CPU time (in hours)
# y2 = disk input/output capacity
# x1 = customer orders (in thousands)
# x2 = add-delete items (in thousands)
```

```
# Fitting the multivariate linear regression model:
```

```
comp.mod.y1 <- lm(y1 ~ x1 + x2)
comp.mod.y2 <- lm(y2 ~ x1 + x2)
```

```
Beta.hat <- cbind(coef(comp.mod.y1), coef(comp.mod.y2))
Beta.hat
```

```
##           [,1]      [,2]
## (Intercept) 8.4236890 14.141491
## x1          1.0789825  2.253854
## x2          0.4198885  5.665367
```

```
# A quicker way to get a summary:
Y <- as.matrix(comput[,c("y1", "y2")])
comp.mod <- lm(Y ~ x1 + x2)
summary(comp.mod)
```

```
## Response y1 :
```

```
##
```

```
## Call:
```

```
## lm(formula = y1 ~ x1 + x2)
```

```
##
```

```
## Residuals:
```

```
##      1      2      3      4      5      6      7
## -1.0632 -1.0315  1.1007 -0.9151  0.4650  1.1066  0.3375
```

```
##
```

```

## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  8.42369    3.44328   2.446  0.0707 .
## x1           1.07898    0.02749  39.249 2.52e-06 ***
## x2           0.41989    0.14447   2.906  0.0438 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.204 on 4 degrees of freedom
## Multiple R-squared:  0.9979, Adjusted R-squared:  0.9969
## F-statistic: 966.5 on 2 and 4 DF,  p-value: 4.265e-06
##
##
## Response y2 :
##
## Call:
## lm(formula = y2 ~ x1 + x2)
##
## Residuals:
##      1      2      3      4      5      6      7
## -2.6350  0.4754  1.4750 -0.9790  0.7887 -0.3555  1.2304
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 14.14149    5.06912   2.79  0.0493 *
## x1           2.25385    0.04047  55.69 6.22e-07 ***
## x2           5.66537    0.21269  26.64 1.18e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.772 on 4 degrees of freedom
## Multiple R-squared:  0.9993, Adjusted R-squared:  0.999
## F-statistic: 2932 on 2 and 4 DF,  p-value: 4.646e-07
# Getting the matrix of fitted values:

X.mat <- cbind(rep(1,times=nrow(comput)),x1,x2) X =  $\begin{bmatrix} 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \\ 1 & x_1 & x_2 \end{bmatrix}$ 

Y.hat <- X.mat %*% Beta.hat
Y.hat

##           [,1]      [,2]
## [1,] 142.5632 304.4350
## [2,] 169.9315 395.6246
## [3,] 153.6993 326.7250
## [4,] 147.4151 308.3790
## [5,] 172.3350 361.6113
## [6,] 158.9934 369.8555
## [7,] 108.1625 227.8696
# Getting the matrix of residuals:

resid.mat <- cbind (y1,y2) - Y.hat
resid.mat

##           y1      y2

```

```
## [1,] -1.0631528 -2.6350279
## [2,] -1.0314649  0.4754450
## [3,]  1.1006666  1.4749621
## [4,] -0.9151495 -0.9789776
## [5,]  0.4649604  0.7887049
## [6,]  1.1065938 -0.3555287
## [7,]  0.3375464  1.2304222
```

Note we could get these matrices from the individual regression models as well:

```
cbind(fitted(comp.mod.y1), fitted(comp.mod.y2))
```

```
##      [,1]      [,2]
## 1 142.5632 304.4350
## 2 169.9315 395.6246
## 3 153.6993 326.7250
## 4 147.4151 308.3790
## 5 172.3350 361.6113
## 6 158.9934 369.8555
## 7 108.1625 227.8696
```

```
cbind(resid(comp.mod.y1), resid(comp.mod.y2))
```

```
##      [,1]      [,2]
## 1 -1.0631528 -2.6350279
## 2 -1.0314649  0.4754450
## 3  1.1006666  1.4749621
## 4 -0.9151495 -0.9789776
## 5  0.4649604  0.7887049
## 6  1.1065938 -0.3555287
## 7  0.3375464  1.2304222
```

Testing about x2 in the model:

Full model:

```
my.n <- length(y1) # number of individuals
my.r <- ncol(X.mat) - 1 # total number of predictors
my.m <- ncol(resid.mat) # total number of responses
```

```
Sigma..full <- (my.n-1)*var(resid.mat)/my.n
# or
Sigma..full <- t(resid.mat)%*%resid.mat/my.n
# or
Sigma..full <- crossprod(resid.mat)/my.n
```

Reduced model (without x2)

```
red.mod <- lm(Y ~ x1)
Beta.hat.redu <- cbind(coef(red.mod))
```

Note: If we were testing the null hypothesis that the *entire set* of predictors was useless in the model, our reduced model would contain ONLY an intercept term. We could fit such a reduced model using:

```
Beta.hat.redu <- cbind(coef(lm(y1 ~ 1)), coef(lm(y2 ~ 1)))
```

```

X.mat.redu <- cbind(rep(1,times=nrow(comput)),x1)
my.q <- ncol(X.mat.redu) - 1 # number of predictors in the reduced model

Y.hat.redu <- X.mat.redu %*% Beta.hat.redu

resid.mat.redu <- cbind (y1,y2) - Y.hat.redu

Sigma..redu <- (my.n-1)*var(resid.mat.redu)/my.n

Wilks <- det(Sigma..full)/det(Sigma..redu)
Wilks

## [1] 0.004002217
my.test.stat <- -(my.n - my.r - 1 - 0.5*(my.m - my.r + my.q + 1)) * log(Wilks )
my.test.stat

## [1] 16.56272
p.value <- pchisq(my.test.stat, df = my.r*(my.r - my.q), lower.tail=F )
p.value

## [1] 0.0002531925
## A quicker way to test H_0
# Estimated Error Covariance Matrix
hat_epsilon2 <- crossprod(Y - comp.mod$fitted.values)
#SigmaTilde <- hat_epsilon2 / (my.n - my.r - 1) # Sums-of-Squares and Crossproducts
SigmaHat <- hat_epsilon2 / my.n #MLE of \Sigma
SigmaHat

##          y1          y2
## y1 0.8282650 0.7455714
## y2 0.7455714 1.7951036

SigmaHat1 <- crossprod(Y - red.mod$fitted.values) / my.n #MLE of \Sigma for redu model
Wilks <- det(SigmaHat)/det(SigmaHat1)
Wilks

## [1] 0.004002217
# modified statistic for large n
-(my.n - my.r - 1 - 0.5*(my.m - my.r + my.q + 1)) *log(Wilks)

## [1] 16.56272
## Even quicker.
anova(comp.mod, red.mod, test="Wilks")

## Analysis of Variance Table
##
## Model 1: Y ~ x1 + x2
## Model 2: Y ~ x1
##   Res.Df Df Gen.var.      Wilks approx F num Df den Df    Pr(>F)
## 1      4      1.6885
## 2      5  1 21.3521 0.0040022  373.29      2      3 0.0002532 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Confidence Ellipsoids and Prediction Ellipsoids in Multivariate Regression

Suppose we have a new individual whose values of the predictor variables are known but whose values for the response variables are not available. A point prediction of $[Y_1, Y_2, \dots, Y_m]$ for this individual is simply $\mathbf{x}_0^T \hat{\boldsymbol{\beta}}$ where $\mathbf{x}_0 = [1, \mathbf{x}_{10}, \dots, \mathbf{x}_{r0}]$ contains the known values of the predictor variables for that individual.

An m -dimensional $100(1 - \alpha)\%$ **Confidence ellipsoid** for the individual responses \mathbf{Y}_{0i} are

$$\mathbf{X}_0^T \hat{\boldsymbol{\beta}}_{(i)} \pm \sqrt{\left(\frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha)\right)} \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 \left(\frac{n}{n-r-1} \hat{\sigma}_{ii}\right)}, i = 1, 2, \dots, m$$

An m -dimensional $100(1 - \alpha)\%$ **prediction ellipsoid** for the individual responses \mathbf{Y}_{0i} are

$$\mathbf{X}_0^T \hat{\boldsymbol{\beta}}_{(i)} \pm \sqrt{\left(\frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha)\right)} \sqrt{(1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0) \left(\frac{n}{n-r-1} \hat{\sigma}_{ii}\right)}, i = 1, 2, \dots, m$$

Suppose we have a new site with 130 thousand orders and 7.5 thousand add-delete items. Let's get a 95% prediction ellipse for that site's (CPU time, disk input/output capacity). The predictor values of interest (including a "1" for the intercept term):

```

pred.mlm <- function(object, newdata, level=0.95,
                    interval = c("confidence", "prediction")){
  form <- as.formula(paste("-", as.character(formula(object))[3]))
  xnew <- model.matrix(form, newdata)
  fit <- predict(object, newdata)
  Y <- model.frame(object)[,1]
  X <- model.matrix(object)
  n <- nrow(Y)
  m <- ncol(Y)
  p <- ncol(X) - 1
  sigmas <- colSums((Y - object$fitted.values)^2) / (n - p - 1)
  fit.var <- diag(xnew %*% tcrossprod(solve(crossprod(X)), xnew))
  if(interval[1]=="prediction") fit.var <- fit.var + 1
  const <- qf(level, df1=m, df2=n-p-m) * m * (n - p - 1) / (n - p - m)
  vmat <- (n/(n-p-1)) * outer(fit.var, sigmas)
  lwr <- fit - sqrt(const) * sqrt(vmat)
  upr <- fit + sqrt(const) * sqrt(vmat)
  if(nrow(xnew)==1L){
    ci <- rbind(fit, lwr, upr)
    rownames(ci) <- c("fit", "lwr", "upr")
  } else {
    ci <- array(0, dim=c(nrow(xnew), m, 3))
    dimnames(ci) <- list(1:nrow(xnew), colnames(Y), c("fit", "lwr", "upr") )
    ci[, ,1] <- fit
    ci[, ,2] <- lwr
    ci[, ,3] <- upr
  }
  ci
}

newdata <- data.frame(x1=130, x2=7.5)
# confidence interval
pred.mlm(comp.mod, newdata)

```



```

##          y1      y2
## fit 151.8406 349.6327
## lwr 146.9515 342.4351
## upr 156.7297 356.8303
# prediction interval
pred.mlm(comp.mod, newdata, interval="prediction")

##          y1      y2
## fit 151.8406 349.6327
## lwr 142.4323 335.7821
## upr 161.2488 363.4833
# confidence interval (multiple new observations)
newdata <- data.frame(x1=c(130,140),x2=c(7.5,8.5))
pred.mlm(comp.mod, newdata)

## , , fit
##
##          y1      y2
## 1 151.8406 349.6327
## 2 163.0503 377.8366
##
## , , lwr
##
##          y1      y2
## 1 146.9515 342.4351
## 2 157.7044 369.9666
##
## , , upr
##
##          y1      y2
## 1 156.7297 356.8303
## 2 168.3962 385.7067
# prediction interval (multiple new observations)
pred.mlm(comp.mod, newdata, interval="prediction")

## , , fit
##
##          y1      y2
## 1 151.8406 349.6327
## 2 163.0503 377.8366
##
## , , lwr
##
##          y1      y2
## 1 142.4323 335.7821
## 2 153.3968 363.6250
##
## , , upr
##
##          y1      y2
## 1 161.2488 363.4833
## 2 172.7038 392.0483

```

Checking Model Assumptions in Multivariate Regression

- The model assumptions should be checked in multivariate regression using techniques similar to those used in simple linear regression or multiple linear regression.
- To check the normality of the error terms, a normal QQ plot of the residual vectors $\epsilon_{(1)}, \epsilon_{(2)}, \dots, \epsilon_{(m)}$ for each response variable can be examined.
- For each response variable, the residual vector can be plotted against the vector of fitted values to look for outliers or unusual patterns.
- Transformations of one or more response variables may be tried if violations of the model assumptions are apparent.

Checking model assumptions:

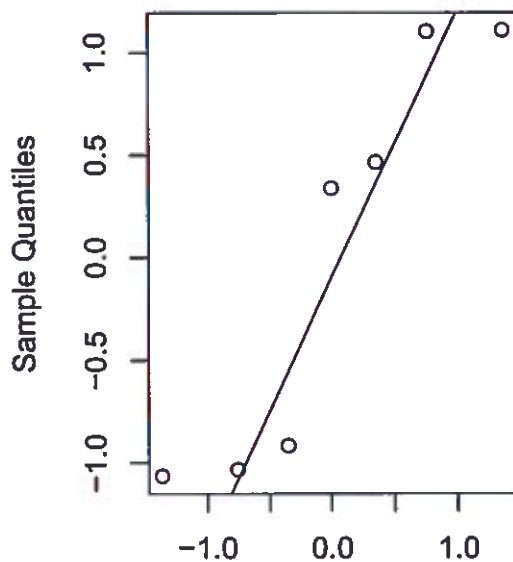
```
par(mfrow=c(1,2))
```

```
qqnorm(resid.mat[,1], main = "Normal Q-Q plot, y1")
```

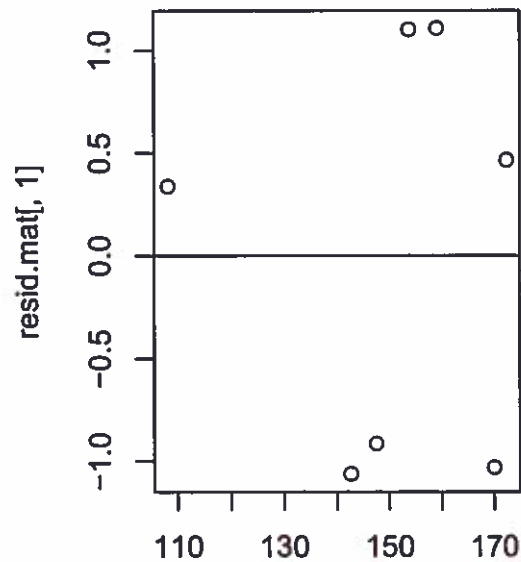
```
qqline(resid.mat[,1])
```

```
plot(Y.hat[,1], resid.mat[,1], main = "Residual plot vs. fitted values, y1"); abline(h=0)
```

Normal Q-Q plot, y1



Residual plot vs. fitted values, y



Theoretical Quantiles

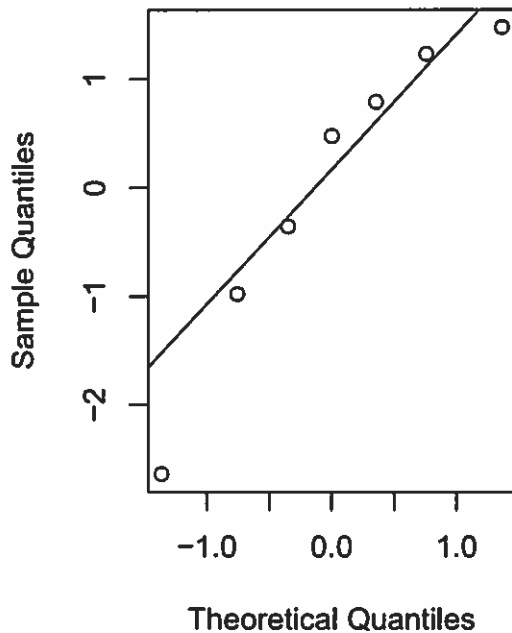
Y.hat[, 1]

```
qqnorm(resid.mat[,2], main = "Normal Q-Q plot, y2")
```

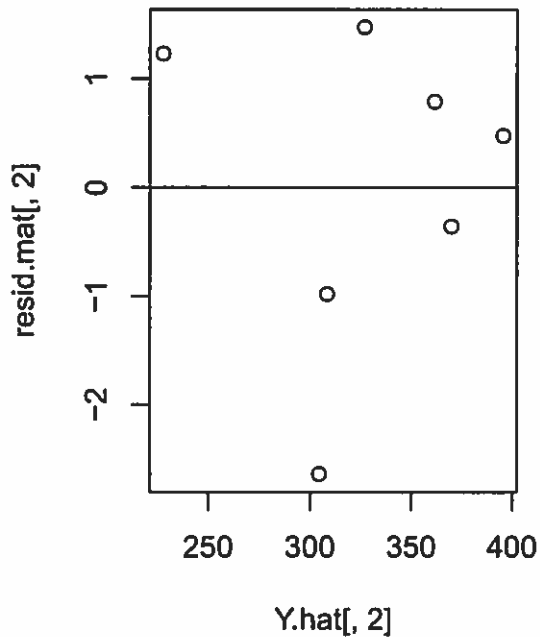
```
qqline(resid.mat[,2])
```

```
plot(Y.hat[,2], resid.mat[,2], main = "Residual plot vs. fitted values, y2"); abline(h=0)
```

Normal Q-Q plot, y2



Residual plot vs. fitted values, y



Example 2: The sales performance data:

```
salesdata <- read.table("salesmat.txt", header=T)
```

```
attach(salesdata)
```

```
## The following objects are masked from comput:
```

```
##
```

```
## x1, x2
```

```
# Recall the variables:
```

```
# X1 = Sales growth
```

```
# X2 = Sales profitability
```

```
# X3 = New account sales
```

```
# X4 = Creativity Test
```

```
# X5 = Mechanical Reasoning Test
```

```
# X6 = Abstract Reasoning Test
```

```
# X7 = Mathematics Test
```

```
# X8 = Historical Facts Test
```

```
# X9 = Sports Trivia Test
```

```
# X10 = Music Trivia Test
```

```
# We will try to predict the "performance" variables (x1, x2, x3) using the "test score" variables (x4,
```

```
# Fitting the multivariate linear regression model:
```

```
sales.mod.1 <- lm(x1 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
```

```
sales.mod.2 <- lm(x2 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
```

```
sales.mod.3 <- lm(x3 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
```

```
# A quick summary:
sales.mod <- lm(cbind(x1,x2,x3) ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
summary(sales.mod)
```

```
## Response x1 :
##
## Call:
## lm(formula = x1 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.9553 -0.9130  0.4151  1.2096  2.0329
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 69.27323    1.67778  41.289 < 2e-16 ***
## x4           0.38292    0.07790   4.915 1.40e-05 ***
## x5           0.29869    0.10143   2.945 0.00525 **
## x6           0.69855    0.13863   5.039 9.38e-06 ***
## x7           0.45469    0.03288  13.830 < 2e-16 ***
## x8          -0.01418    0.02128  -0.666 0.50877
## x9           0.02917    0.02515   1.160 0.25273
## x10          -0.01121    0.02521  -0.445 0.65886
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.71 on 42 degrees of freedom
## Multiple R-squared:  0.9536, Adjusted R-squared:  0.9459
## F-statistic: 123.3 on 7 and 42 DF, p-value: < 2.2e-16
##
##
## Response x2 :
##
## Call:
## lm(formula = x2 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7704 -1.5003 -0.1821  1.1856  3.6953
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 76.48366    1.87834  40.719 < 2e-16 ***
## x4           0.17506    0.08722   2.007 0.0512 .
## x5           0.83095    0.11356   7.317 5.13e-09 ***
## x6          -0.48542    0.15520  -3.128 0.0032 **
## x7           0.78419    0.03681  21.306 < 2e-16 ***
## x8          -0.02167    0.02382  -0.910 0.3681
## x9          -0.03491    0.02816  -1.240 0.2220
## x10          -0.02057    0.02822  -0.729 0.4701
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```

## Residual standard error: 1.914 on 42 degrees of freedom
## Multiple R-squared: 0.9693, Adjusted R-squared: 0.9642
## F-statistic: 189.4 on 7 and 42 DF, p-value: < 2.2e-16
##
##
## Response x3 :
##
## Call:
## lm(formula = x3 ~ x4 + x5 + x6 + x7 + x8 + x9 + x10)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.4800 -0.6304 -0.0103  0.6852  2.7140
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 84.057314   1.327509   63.320 < 2e-16 ***
## x4           0.531052   0.061640    8.615 7.88e-11 ***
## x5          -0.059393   0.080257   -0.740  0.463
## x6           0.672190   0.109689    6.128 2.60e-07 ***
## x7           0.229394   0.026012   8.819 4.16e-11 ***
## x8          -0.010272   0.016834   -0.610  0.545
## x9          -0.009681   0.019902   -0.486  0.629
## x10          0.009646   0.019947    0.484  0.631
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.353 on 42 degrees of freedom
## Multiple R-squared: 0.9298, Adjusted R-squared: 0.9181
## F-statistic: 79.46 on 7 and 42 DF, p-value: < 2.2e-16
Beta.hat <- cbind(coef(sales.mod.1), coef(sales.mod.2), coef(sales.mod.3) )
Beta.hat

##              [,1]      [,2]      [,3]
## (Intercept) 69.27322626 76.48366386 84.057313608
## x4           0.38292240  0.17505706  0.531051725
## x5           0.29869074  0.83095235 -0.059393242
## x6           0.69854730 -0.48542385  0.672190486
## x7           0.45468653  0.78419335  0.229393734
## x8          -0.01417885 -0.02166938 -0.010271937
## x9           0.02917020 -0.03491088 -0.009680862
## x10          -0.01120951 -0.02057208  0.009646095
X.mat <- cbind(rep(1,times=nrow(salesdata)), x4, x5, x6, x7, x8, x9, x10)
Y.hat <- cbind(fitted(sales.mod.1), fitted(sales.mod.2), fitted(sales.mod.3))

resid.mat <- cbind(resid(sales.mod.1), resid(sales.mod.2), resid(sales.mod.3))

#### Testing whether the set (x8, x9, x10) is useless,
#### given the presence of x4, x5, x6, x7 as predictors in the model:

## Full model:

my.n <- length(x1) # number of individuals

```

```

my.p <- ncol(X.mat) - 1 # total number of predictors
my.r <- ncol(resid.mat) # total number of responses

Sigma.full <- (my.n-1)*var(resid.mat)/my.n

## Reduced model (without x8, x9, x10)

Beta.hat.redu <- cbind(coef(lm(x1 - x4 + x5 + x6 + x7) ),
                      coef(lm(x2 - x4 + x5 + x6 + x7) ), coef(lm(x3 - x4 + x5 + x6 + x7)) )

X.mat.redu <- cbind(rep(1,times=nrow(salesdata)), x4, x5, x6, x7)
my.p.redu <- ncol(X.mat.redu) - 1 # number of predictors in the reduced model

Y.hat.redu <- X.mat.redu %*% Beta.hat.redu

resid.mat.redu <- cbind (x1,x2,x3) - Y.hat.redu

Sigma.redu <- (my.n-i)*var(resid.mat.redu)/my.n

my.test.stat <- -(my.n - my.p - 1 - 0.5*(my.r - my.p + my.p.redu + 1))*
log( det(Sigma.full)/det(Sigma.redu) )
my.test.stat

## [1] 6.440941
p.value <- pchisq(my.test.stat, df = my.r*(my.p - my.p.redu), lower.tail=F )
p.value

## [1] 0.6951036
## new data
newdata <- data.frame(x4=10, x5=12,x6=9, x7=25, x8=15,x9=12,x10=40 )
# confidence interval
pred.mlm(sales.mod, newdata)

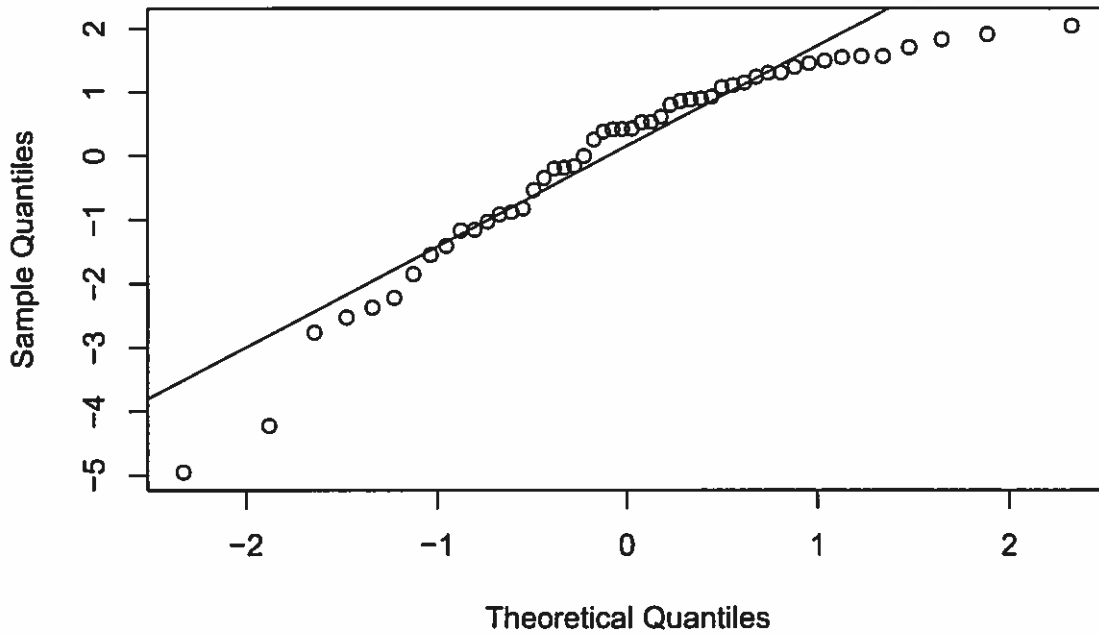
##          x1          x2          x3
## fit 94.02981 101.87483 100.55526
## lwr 91.75158  99.32427  98.75266
## upr 96.30804 104.42539 102.35786
# prediction interval
pred.mlm(sales.mod, newdata, interval="prediction")

##          x1          x2          x3
## fit 94.02981 101.87483 100.55526
## lwr 88.00355  95.12822  95.78712
## upr 100.05606 108.62143 105.32341
### Checking model assumptions:

qqnorm(resid.mat[,1], main = "Normal Q-Q plot, Sales growth")
qqline(resid.mat[,1])

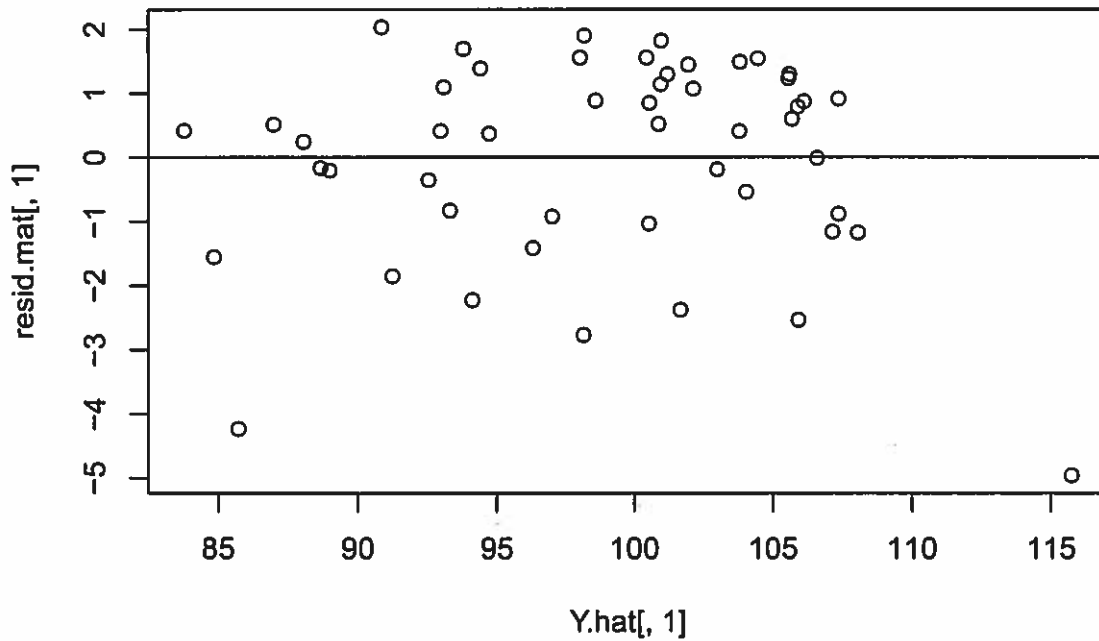
```

Normal Q-Q plot, Sales growth



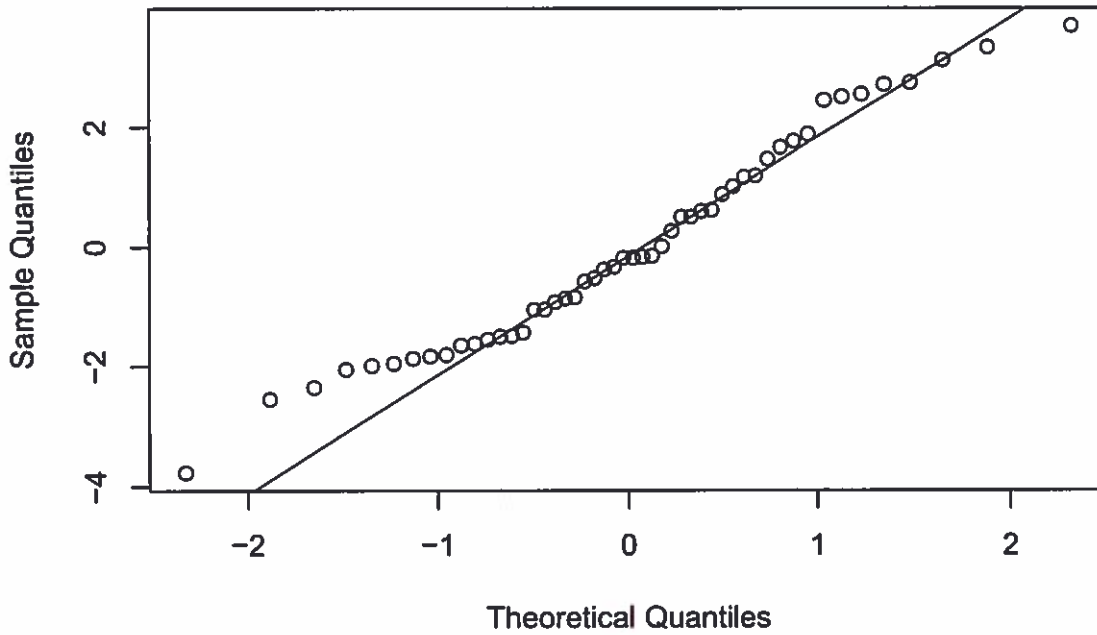
```
plot(Y.hat[,1],resid.mat[,1], main = "Residual plot vs. fitted values, Sales growth");  
abline(h=0)
```

Residual plot vs. fitted values, Sales growth



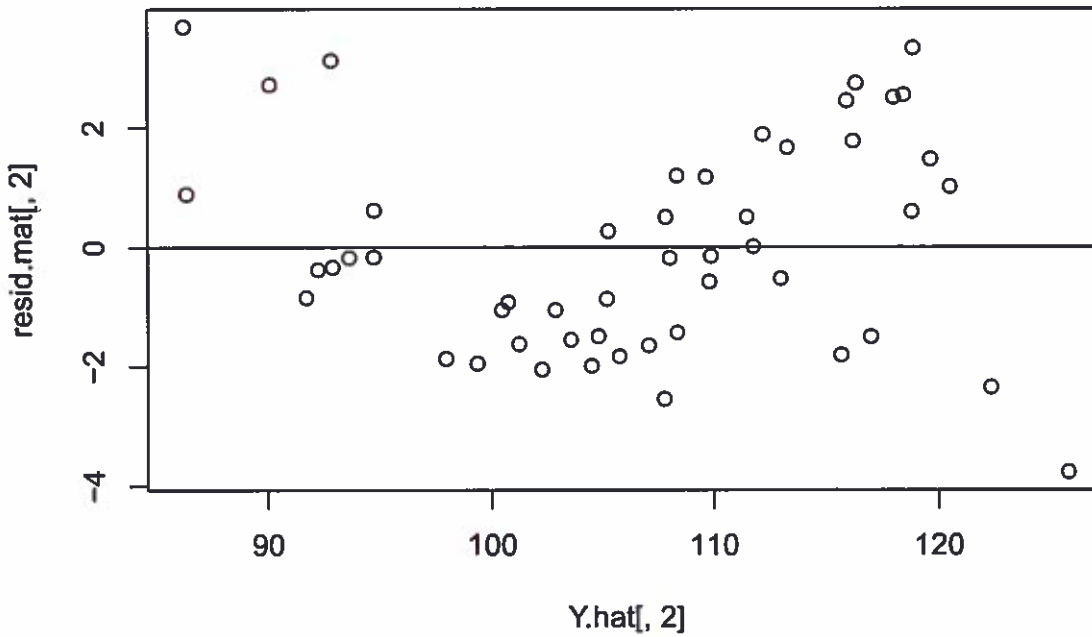
```
qqnorm(resid.mat[,2], main = "Normal Q-Q plot, Sales profitability")  
qqline(resid.mat[,2])
```


Normal Q-Q plot, Sales profitability



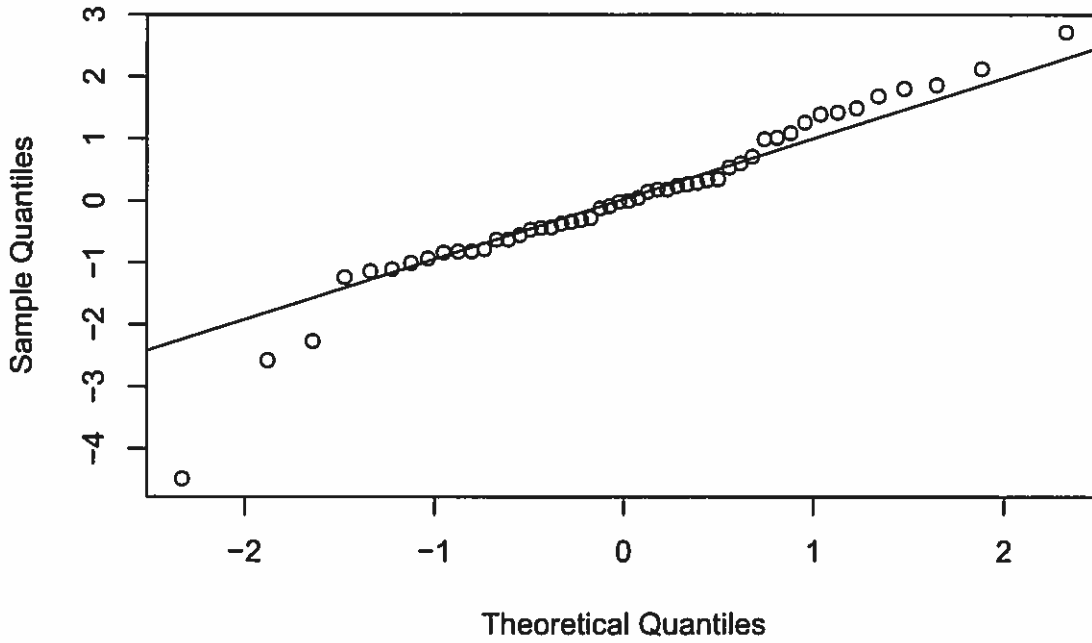
```
plot(Y.hat[,2],resid.mat[,2], main = "Residual plot vs. fitted values, Sales profitability");  
abline(h=0)
```

Residual plot vs. fitted values, Sales profitability



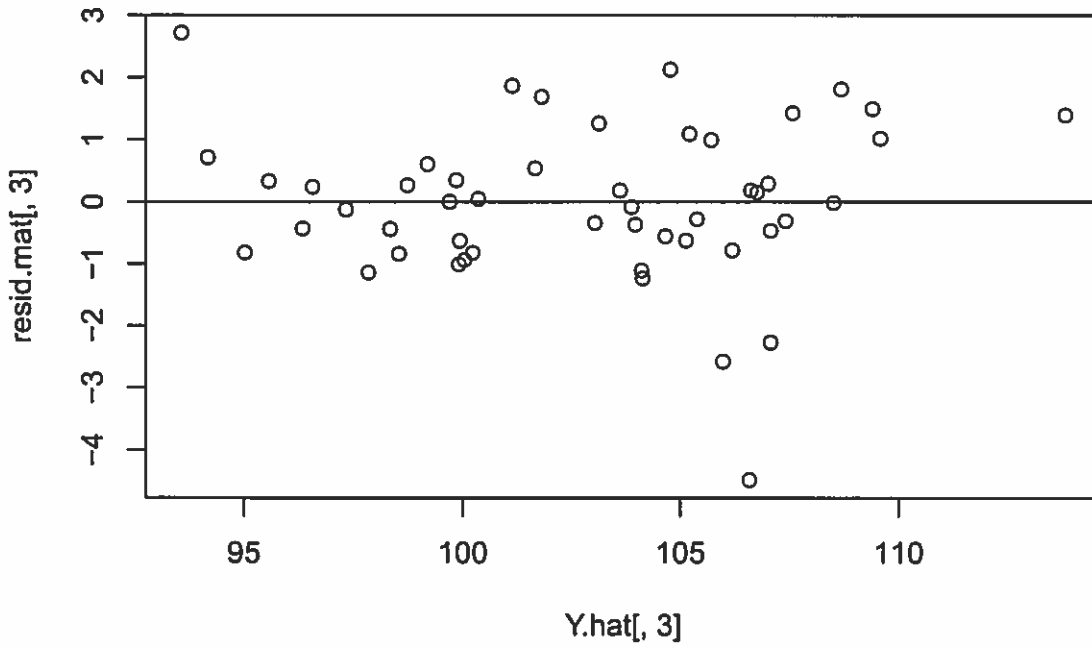
```
qqnorm(resid.mat[,3], main = "Normal Q-Q plot, New account sales")  
qqline(resid.mat[,3])
```

Normal Q-Q plot, New account sales



```
plot(Y.hat[,3],resid.mat[,3], main = "Residual plot vs. fitted values, New account sales");  
abline(h=0)
```

Residual plot vs. fitted values, New account sales





Note: At first we have larger amount of random vectors $X_1, X_2, \dots, X_q, \dots, X_n$ $n > q$.

Chapter 4. Principal Component Analysis

Jianxuan Liu

Fall 2018

One of the aims in multivariate data analysis is to summarize the data in fewer than the original number of dimensions without losing essential information. More than a century ago, Pearson (1901) considered this problem, and Hotelling (1933) proposed a solution to it: instead of treating each variable separately, he considered combinations of the variables. Clearly, the average of all variables is such a combination, but many others exist. Two fundamental questions arise: 1. How should one choose these combinations? 2. How many such combinations should one choose?

There is no single strategy that always gives the right answer. This chapter will describe many ways of tackling at least the first problem.

$$\tilde{X} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{q1} & X_{q2} & \dots & X_{qn} \end{bmatrix} \Rightarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \Rightarrow Y = a^T X$$

1. Principal Components Analysis

random variables

The goal of principal components analysis (PCA) is to describe most of the variation in the multivariate data set X_1, X_2, \dots, X_q using a smaller or more concise set of variables. PCA can be thought of as a form of data reduction.

Data reduction is to use a lower-dimensional summarization of the data that still contains the relevant information that is in the full data. The hope is to be able to summarize what makes observations similar and what makes them different using relatively few indices. Also, we'd like each of these indices to say something distinct about how the observations vary.

In PCA, we create a new set of "variables" Y_1, Y_2, \dots, Y_q , each of which is a linear combination of X_1, X_2, \dots, X_q ,

choose the first m variables $Y_i, i = 1, m$ so that we can capture

$$\begin{cases} Y_1 = a_{11}X_1 + a_{12}X_2 + \dots + a_{1q}X_q \\ Y_2 = a_{21}X_1 + a_{22}X_2 + \dots + a_{2q}X_q \\ \vdots \\ Y_q = a_{q1}X_1 + a_{q2}X_2 + \dots + a_{qq}X_q \end{cases} = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_q^T \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_q \end{bmatrix}$$

The principal components are those uncorrelated linear combinations Y_1, Y_2, \dots, Y_q whose variance are as large as possible.

We further require these new "variables" (or indices) Y 's to be uncorrelated. This assures us that the information in Y_2 , say, does not overlap with the information in Y_1 . However, having q of these indices does not give us any data reduction. We would like to choose only the first m of these (where $m < q$) to focus on. Thus we choose the first linear combination, Y_1 , so that it is the linear combination of X_1, X_2, \dots, X_q that accounts for as much of the variation in the original data as possible. Then the second linear combination, Y_2 , is the linear combination that accounts for as much of the variation in the original data as possible (provided that it is uncorrelated with Y_1). And the third linear combination, Y_3 , is the linear combination that accounts for as much of the variation in the original data as possible (provided that it is uncorrelated with Y_1 and Y_2), and so on. Eventually, it becomes a practical issue of how many of these indices (which are called principal components) we really need.

Idea: We first have big amount of variables. Then we choose the first q important vectors. Through those important q important vectors, X_1, \dots, X_q , we construct Y_1, \dots, Y_q and then.

2. Mathematics Behind PCA

Singular value decomposition (SVD)

Suppose X is a $p \times n$ matrix with rank r . Then there exists an $p \times n$ matrix

$$\Gamma = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}$$

$\text{rank}(A)$
||

for which the diagonal entries in D are the first r singular values of X , $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, and there exist an $p \times p$ orthogonal matrix U and an $n \times n$ orthogonal matrix V such that

$$X = U \Gamma V^T$$

orthogonal orthogonal

$$U^T U = I \\ V^T V = I$$

The singular values of X are the square roots of the eigenvalues of $X^T X$, denoted by $\sigma_i, i = 1, \dots, p$ and they are arranged in decreasing order. That is, $\sigma_i = \sqrt{\lambda_i}$ for $1 \leq i \leq n$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

- the columns of U are called **left singular vectors** of X
- the columns of V are called **right singular vectors** of X

Example: Find a singular value decomposition of $X = \begin{bmatrix} 1 & -1 \\ -2 & 2 \\ 2 & -2 \end{bmatrix}$.

Solution:

$$* A = X^T X = \begin{bmatrix} 1 & -2 \\ -1 & 2 \\ 2 & -2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -2 & 2 \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} 9 & -9 \\ -9 & 9 \end{bmatrix}$$

* Find eigenvalue of A :

$$\det(\lambda I - A) = \det \begin{pmatrix} \lambda - 9 & -9 \\ -9 & \lambda - 9 \end{pmatrix} = (\lambda - 9)^2 - 81 = 0$$

$$\Rightarrow \lambda_1 = 18 \quad \lambda_2 = 0$$

\Rightarrow corresponding unit eigenvectors

$$\vec{v}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \quad \vec{v}_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$\Rightarrow V = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

* Singular values are $\sigma_1 = \sqrt{18} \quad \sigma_2 = 0$

$$\Rightarrow \Gamma = \begin{bmatrix} \sqrt{18} & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} D$$

$$* \text{To construct } U: \quad X V_1 = \begin{bmatrix} \frac{2\sqrt{2}}{\sqrt{2}} \\ -\frac{4\sqrt{2}}{\sqrt{2}} \\ \frac{4\sqrt{2}}{\sqrt{2}} \end{bmatrix} \quad X V_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

orthogonal unit vector u_1, u_2

$$u_1^T X = 0 \Rightarrow x_1 - 2x_2 + 2x_3 = 0 \quad (1)$$

A basis for the solution of (1)

$$\text{is } w_1 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad w_2 = \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix}$$

$$\rightarrow u_1 = \begin{bmatrix} \frac{2\sqrt{15}}{15} \\ \frac{1\sqrt{15}}{15} \\ 0 \end{bmatrix} \quad u_2 = \begin{bmatrix} \frac{-2\sqrt{15}}{15} \\ \frac{4\sqrt{15}}{15} \\ \frac{5\sqrt{15}}{15} \end{bmatrix}$$

Then

$$X = \begin{bmatrix} 1 & -2 \\ -2 & 2 \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{2}{\sqrt{15}} & \frac{-2}{\sqrt{15}} \\ -\frac{2}{\sqrt{3}} & \frac{4}{\sqrt{15}} & \frac{5}{\sqrt{15}} \\ \frac{2}{\sqrt{3}} & 0 & \frac{5}{\sqrt{15}} \end{bmatrix} \begin{bmatrix} \sqrt{18} & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$Y = a^T X$$

$$L(e_i, \lambda_i) = e_i^T \Sigma e_i + \lambda (1 - e_i^T e_i) \quad \left| \begin{array}{l} \text{Choose } e_i \text{ to maximize} \\ e_i^T \Sigma e_i \text{ s.t. } e_i^T e_i = 1 \end{array} \right.$$

Principal component analysis can capture the maximum sample variance

- First step: Choose coefficients $a_{11}, a_{12}, \dots, a_{1q}$ such that the sample variance of

$$Y_1 = a_{11}X_1 + a_{12}X_2 + \dots + a_{1q}X_q$$

$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1q} \end{bmatrix}$$

is as large as possible.

Of course, we could arbitrarily inflate this sample variance by making $a_{11}, a_{12}, \dots, a_{1q}$ arbitrarily large. However, to make this maximization problem meaningful, we include the constraint that

$$a_1^T a_1 = \sum_{i=1}^q a_{1i}^2 = 1 \quad \vec{a}_1: \text{new dimension}$$

- Second step: Choose coefficients $a_{21}, a_{22}, \dots, a_{2q}$ such that the sample variance of

$$Y_2 = a_{21}X_1 + a_{22}X_2 + \dots + a_{2q}X_q$$

is as large as possible, subject to $a_2^T a_2 = 1$. We also want Y_2 to be uncorrelated with Y_1 , so we add the further constraint that

$$a_1^T a_2 = \sum_{i=1}^q a_{1i} a_{2i} = 0$$

- j^{th} step: Choose coefficients $a_{j1}, a_{j2}, \dots, a_{jq}$ such that the sample variance of

$$Y_j = a_{j1}X_1 + a_{j2}X_2 + \dots + a_{jq}X_q$$

is as large as possible, subject to $a_j^T a_j = 1$ and $a_j^T a_{j'} = 0$ for all $j' < j$

Mathematically, we can obtain q linear combinations Y_1, Y_2, \dots, Y_q in this way. This optimization problem is solved using the method of Lagrange multipliers. In practice, we use a software package like R or SAS to find the coefficients for us.

Formally, consider $X \sim (\mu, \Sigma)$ (could be any continuous distribution), with $\Sigma = PAP^T$. Let r be the rank of Λ and $k = 1, \dots, r$

- The k^{th} principal component score is the scalar

$$W_k = a_k^T (X - \mu)$$

\vec{a}_k : new k^{th} dimension

- The k -dimensional principal component vector is the scalar

$$W^{(k)} = [W_1, \dots, W_k]^T = P_k^T (X - \mu) = [a_1^T (X - \mu)]^T [a_k (X - \mu)]$$

- The k^{th} principal component projection (vector) is the scalar

$$P_k = a_k a_k^T (X - \mu)$$

We use the letter W for the principal component score, indicating that they are weighted random variables. The k^{th} principal component score, W_k represents the contribution of X in the direction a_k . The vector a_k are sometimes called the loadings as they represent the "load" or "weight" each variable is accorded in the projection.

Mathematically, W_k is obtained by projecting the centered X onto a_k . Collecting the contributions of the first k scores into one object leads to the principal component vector $W^{(k)}$, a k -dimensional principal component vector which summarizes the contributions of X along the first k eigen directions, the eigenvectors of Σ .

The principal component projection (vector) P_k is a q -dimensional random vector which points in the same or opposite direction as the k^{th} eigenvector a_k , and the length or Euclidean norm of P_k is the absolute value of W_k .

Example

Let $\mathbf{X} = [X_1, X_2]^T$ be a two-dimensional random vector with mean μ and covariance matrix Σ given by

$$\mu = [0, -1]^T, \text{ and } \Sigma = \begin{pmatrix} 2.4 & -0.5 \\ -0.5 & 1 \end{pmatrix}.$$

The eigenvalues and eigenvectors of Σ are

$$(\lambda_1, \mathbf{a}_1) = (2.5602, \begin{bmatrix} -0.9523 \\ 0.3052 \end{bmatrix}) \text{ and } (\lambda_2, \mathbf{a}_2) = (0.8398, \begin{bmatrix} -0.3052 \\ -0.9523 \end{bmatrix})$$

Sigma=matrix(c(2.4, -0.5, -0.5, 1), nrow=2)

```
Sigma
##      [,1] [,2]
## [1,] 2.4 -0.5
## [2,] -0.5 1.0
eigen(Sigma, symmetric=T)

## eigen() decomposition
## $values
## [1] 2.5602325 0.8397675
##
## $vectors
##      [,1] [,2]
## [1,] -0.9522955 -0.3051774
## [2,] 0.3051774 -0.9522955
```

The proportion of total variance accounted by the first PC, W_1 is

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{2.560}{2.560 + 0.839} =$$

$$\left. \begin{aligned} \lambda_1 + \lambda_2 &= 3.4 \\ \text{Var}(X_1) + \text{Var}(X_2) &= 3.4 \end{aligned} \right\} \rightarrow \text{equal}$$

$$W_1 = \mathbf{a}_1^T (\mathbf{X} - \mu) = \begin{bmatrix} -0.952 & 0.3052 \end{bmatrix} \begin{bmatrix} X_1 - 0 \\ X_2 + 1 \end{bmatrix}$$

The eigenvectors show the axes along which the data vary most, with the first vector $\sqrt{\lambda_1} \mathbf{a}_1$ pointing along the direction in which the data have the largest variance. The first eigenvalue is considerably bigger than the second, so much more variability exists along this direction. The two principal component scores are

$$W_1 = -0.9523X_1 + 0.3052(X_2 + 1) \text{ and } W_2 = -0.3052X_1 - 0.9523(X_2 + 1).$$

$$\text{Var}(X_1) > \text{Var}(W_2) \Rightarrow$$

The first PC score is heavily weighted in the direction of the first variable, implying that the first variable contributes more strongly than the second to the variance of \mathbf{x} .

$$\text{Var}(W_1) = 0.9523^2 \underbrace{\text{Var}(X_1)}_{2.4} + 0.3052^2 \underbrace{\text{Var}(X_2)}_1 \quad \Bigg| \quad \text{Var}(W_2) =$$

$$= 1.13 \quad \quad \quad = 1.00$$

PCA and SVD

PCA is an application of the SVD. The principal components are equal to the right singular vectors if we first scale the data by subtracting the column mean and dividing each column by its standard deviation (scale() in R).

Below we plot the first left and right singular vectors along with the original data.

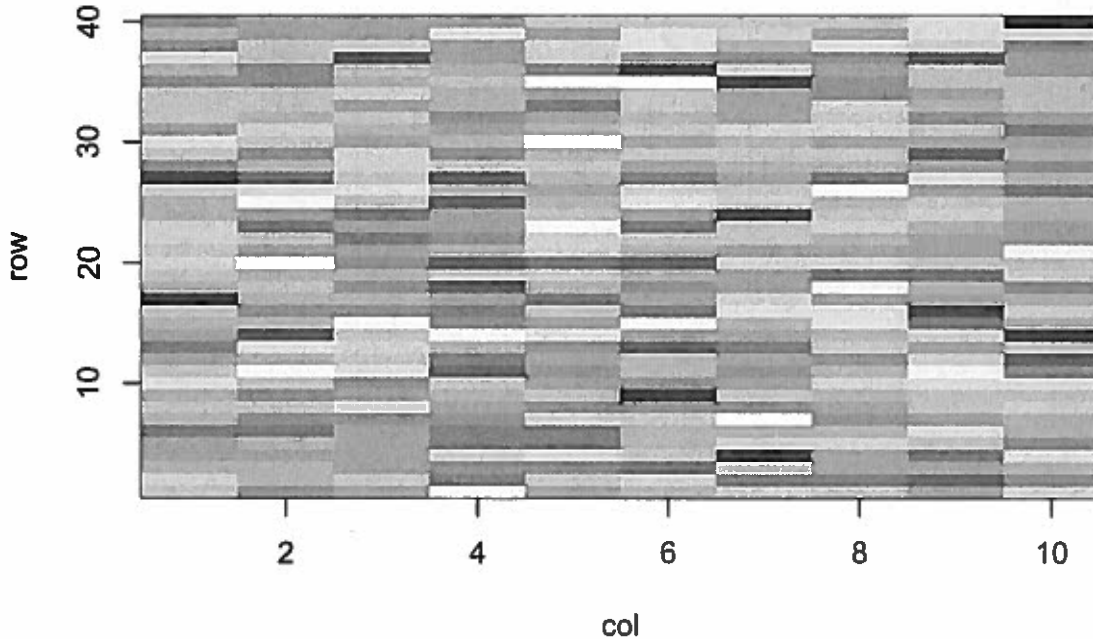
```
library(dplyr) ← working on data name
```

```
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
## filter, lag
## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union
```


40x10

```
set.seed(1)
DX=matrix(rnorm(400), nrow = 40)
image(1:10, 1:40, t(DX)[, nrow(DX):1], ylab="row", xlab="col", main="original data")
```

original data



```
hh=dist(DX) %>% hclust
DXOrdered=DX[hh$order, ]
par(mfrow = c(1, 3))
## Complete data
image(t(DXOrdered)[, nrow(DXOrdered):1], main="Original Data")
## Show the row means
plot(rowMeans(DXOrdered), 40:1, , xlab = "Row Mean", ylab = "Row", pch = 19)
## Show the column means
plot(colMeans(DXOrdered), xlab = "Column", ylab = "Column Mean", pch = 19)
```

* Connection between SVD and PCA

• In SVD $X = U \Gamma V^T$ Γ : contains singular value of X

• In PCA $X \approx (\hat{\mu}, \hat{\Sigma})$ $\hat{\Sigma} = P \Lambda P^T$

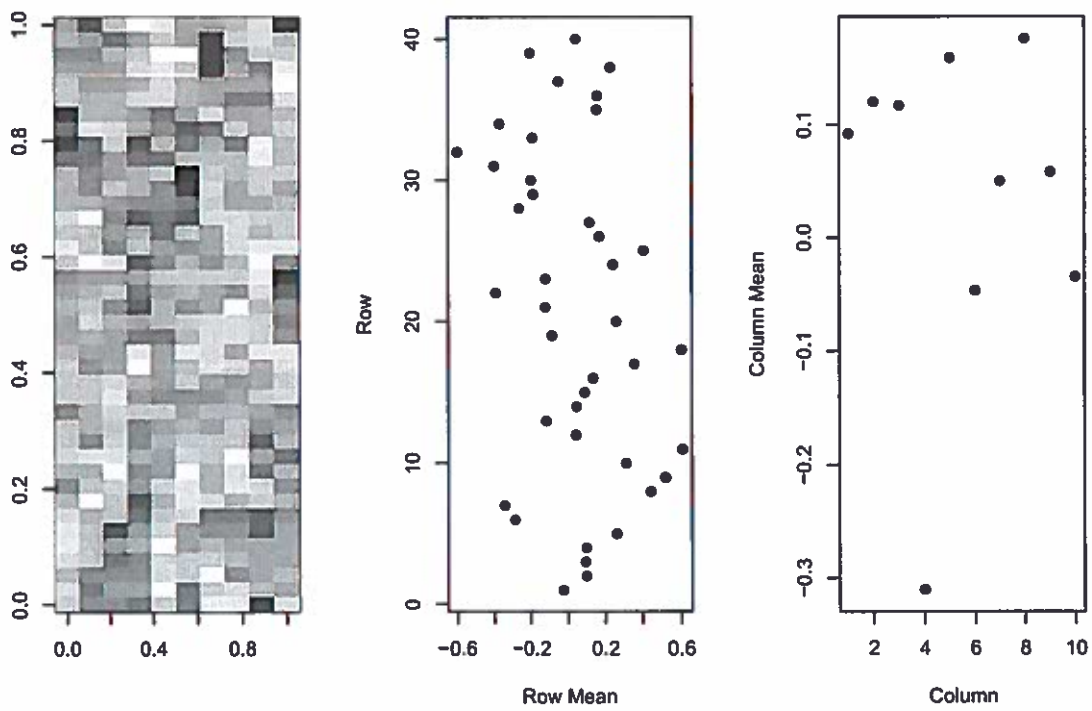
$$\text{cov}(X) = X^T X = (V \Gamma U^T)(U^T \Gamma V^T) = V \Gamma^2 V^T \propto \hat{\Sigma} \propto P \Lambda P^T$$

Γ^2 : eigenvalue of $X^T X$

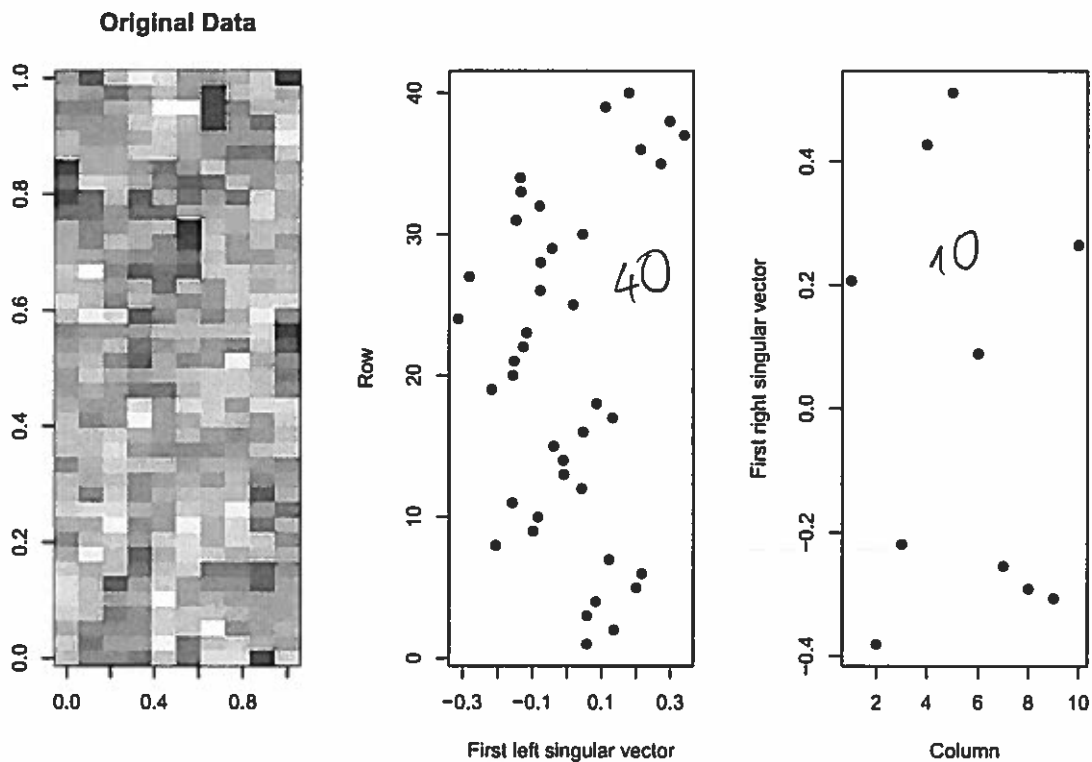
$$\Rightarrow \Gamma = \Lambda^{1/2}$$

$$\Rightarrow P = \frac{1}{\sigma} \alpha^T (X - \hat{\mu}) = V$$

Original Data



```
svd1 = svd(scale(DXOrdered))
par(mfrow = c(1, 3))
image(t(DXOrdered)[, nrow(DXOrdered):1], main = "Original Data")
plot(svd1$u[, 1], 40:1, , ylab = "Row", xlab = "First left singular vector",
     pch = 19)
plot(svd1$v[, 1], xlab = "Column", ylab = "First right singular vector", pch = 19)
```



We can see how the first left and right singular vectors pick up the mean shift in both the rows and columns of the matrix.

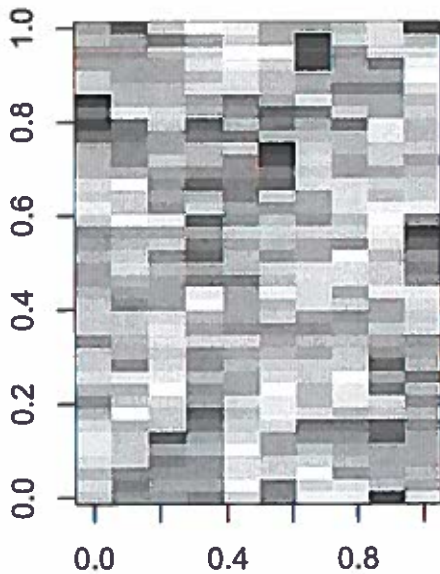
If we believed that the first left and right singular vectors, call them u_1 and v_1 , captured all of the variation in the data, then we could approximate the original data matrix with

$$X \approx u_1 v_1^T.$$

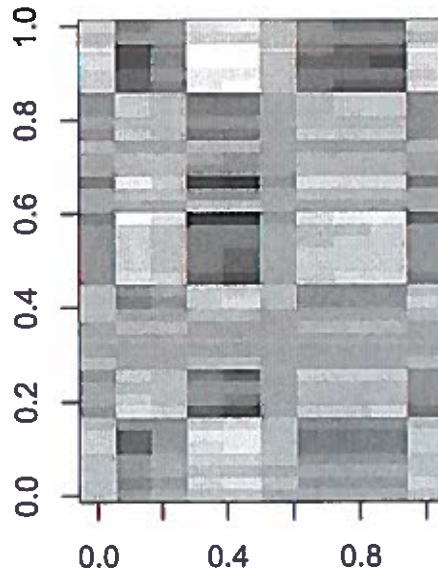
Thus, we would reduce 400 numbers in the original matrix to $40 + 10 = 50$ numbers in the compressed matrix, a nearly 90% reduction in information. Here's what the original data and the approximation would look like.

```
## Approximate original data with outer product of first singular vectors
approx =with(svd1, outer(u[, 1], v[, 1]))
## Plot original data and approximated data
par(mfrow = c(1, 2))
image(t(DXOrdered)[, nrow(DXOrdered):1], main = "Original Matrix")
image(t(approx)[, nrow(approx):1], main = "Approximated Matrix")
```

Original Matrix



Approximated Matrix



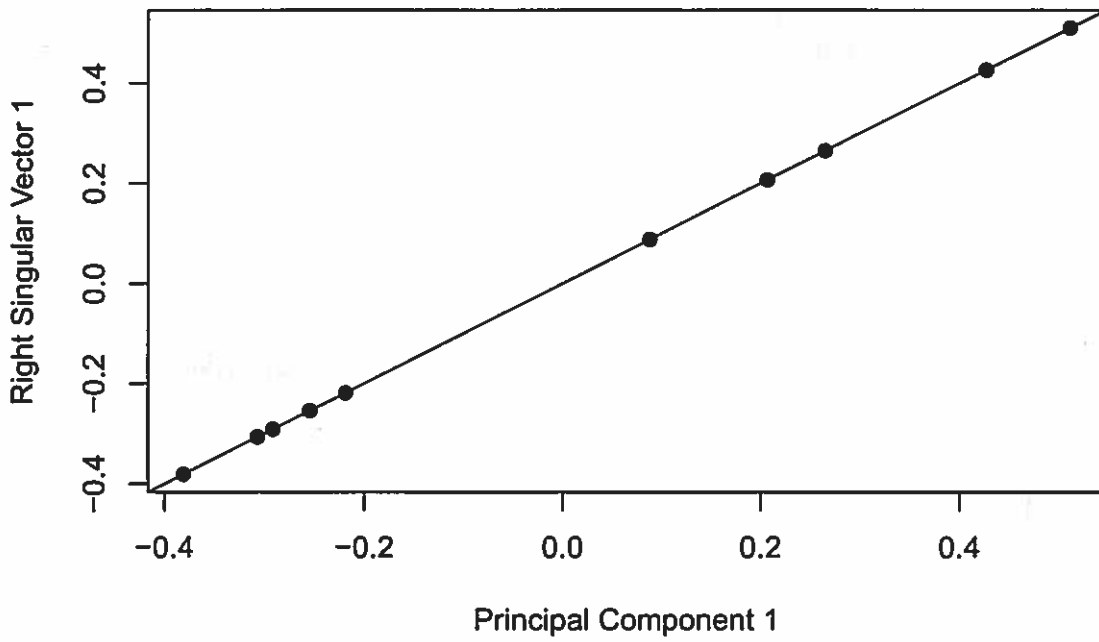
Obviously, the two matrices are not identical, but the approximation seems reasonable in this case. This is not surprising given that there was only one real feature in the original data.

As we mentioned above, the SVD has a close connection to principal components analysis (PCA). Here, we show that the first right singular vector from the SVD is equal to the first principal component vector returned by PCA.

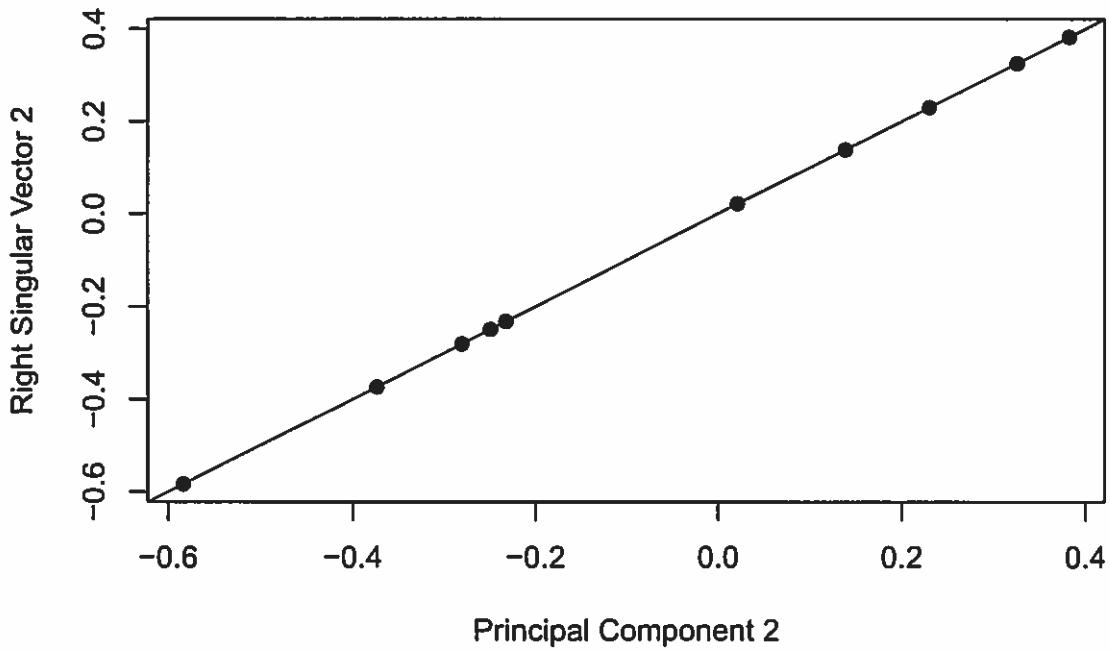
```
svd1=svd(scale(DXOrdered))  
pcal=prcomp(DXOrdered, scale = TRUE)  
plot(pcal$rotation[, 1], svd1$v[, 1], pch = 19,  
      xlab = "Principal Component 1", ylab = "Right Singular Vector 1")  
abline(c(0, 1))
```

have to scale

* `summary(pcal)` → help us know the proportion of variance capture by PCA 1.



```
plot(pca1$rotation[, 2], svd1$w[, 2], pch = 19,
     xlab = "Principal Component 2", ylab = "Right Singular Vector 2")
abline(c(0, 1))
```



$$a_i = \begin{bmatrix} a_{i1} \\ a_{i2} \\ a_{i3} \\ \vdots \\ a_{iq} \end{bmatrix}$$

Y_1, Y_2, \dots, Y_p are uncorrelated

Results

1. Let Σ be the covariance matrix associated with the random vector $\mathbf{X}^T = [X_1, X_2, \dots, X_q]$. Let Σ have the eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{a}_1), (\lambda_2, \mathbf{a}_2), \dots, (\lambda_q, \mathbf{a}_q)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_q \geq 0$. Then the i^{th} principal component is given by

$$Y_i = a_{i1}X_1 + a_{i2}X_2 + \dots + a_{iq}X_q, i = 1, 2, \dots, q.$$

With these choices,

$$E(Y_i) = \mathbf{a}_i^T \boldsymbol{\mu} = \sum_{j=1}^q a_{ij} \mu_j, i = 1, 2, \dots, q$$

$$\text{Var}(Y_i) = \mathbf{a}_i^T \Sigma \mathbf{a}_i = \lambda_i, i = 1, 2, \dots, q \quad \text{same indices}$$

$$\text{Cov}(Y_i, Y_k) = \mathbf{a}_i^T \Sigma \mathbf{a}_k = 0, i \neq k. \quad \neq \text{indices}$$

If some λ_i are equal, then choices of the corresponding coefficient vectors, \mathbf{a}_i , and hence Y_i are not unique.

2. Let Σ be the covariance matrix associated with the random vector $\mathbf{X}^T = [X_1, X_2, \dots, X_q]$. Let Σ have the eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{a}_1), (\lambda_2, \mathbf{a}_2), \dots, (\lambda_q, \mathbf{a}_q)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_q \geq 0$. Let $Y_1 = \mathbf{a}_1^T \mathbf{X}, Y_2 = \mathbf{a}_2^T \mathbf{X}, \dots, Y_q = \mathbf{a}_q^T \mathbf{X}$ be the principal components. Then

$$\sigma_{11} + \sigma_{22} + \dots + \sigma_{qq} = \sum_{i=1}^q \text{Var}(X_i) = \lambda_1 + \lambda_2 + \dots + \lambda_q = \sum_{i=1}^q \text{Var}(Y_i)$$

The trace of a covariance matrix is equal to the sum of its eigenvalues

$$\text{trace}(\Sigma) = \sum_{i=1}^q \lambda_i$$

3. If $Y_1 = \mathbf{a}_1^T \mathbf{X}, Y_2 = \mathbf{a}_2^T \mathbf{X}, \dots, Y_q = \mathbf{a}_q^T \mathbf{X}$ are the principal components obtained from the covariance matrix Σ , then

$$\rho_{Y_i, X_k} = \frac{a_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, i, k = 1, 2, \dots, q$$

are the correlation coefficients between the components Y_i and the variables X_k . Here $(\lambda_1, \mathbf{a}_1), (\lambda_2, \mathbf{a}_2), \dots, (\lambda_q, \mathbf{a}_q)$ are the eigenvalue-eigenvector pairs for Σ .

Choosing the Number of Components

To explain all the variation in the original data, we would need (in general) all q principal components. But it is practically sufficient to explain most of the variation in the original data. This can usually be done using merely the "first few" principal components. If we will retain the first $m < q$ components, how can we choose m ?

Possible Rules

1. Retain the first m components sufficient to explain a specified percentage (70%? 80%? 90%?) of the total variance of the original variables.
2. Keep only the components whose eigenvalues are at least $\sum_{i=1}^q \frac{\lambda_i}{q}$, which is the average eigenvalue and also the average sample variance of the original variables.
3. When PCA is done on the correlation matrix, this average is 1, so Kaiser (1958) suggested keeping components with eigenvalues at least 1. Jolliffe (1972) preferred using 0.7 as the threshold.
4. Cattell (1965) introduced the scree diagram, which plots λ_i against i , for $i = 1, \dots, q$. Look for the "elbow" in the curve and choose the corresponding number of components.
5. Jolliffe (1986) suggested a modified scree plot that plotted $\log(\lambda_i)$ against i , for $i = 1, \dots, q$.

3 An illustration with the Sydney heptathlon data

Before doing anything else with these data, it needs to be noted that in the three running events, better performance is indicated by a lower measure (time), whereas in the jumping and throwing events good performance is indicated by a higher measure (distance). It seems sensible to readjust data so that large values are good and small values bad (Timed racing events are obviously the opposite).

```
library(HSAUR)
```

```
## Loading required package: tools
```



```
data("heptathlon", package = "HSAUR")
head(heptathlon)
```

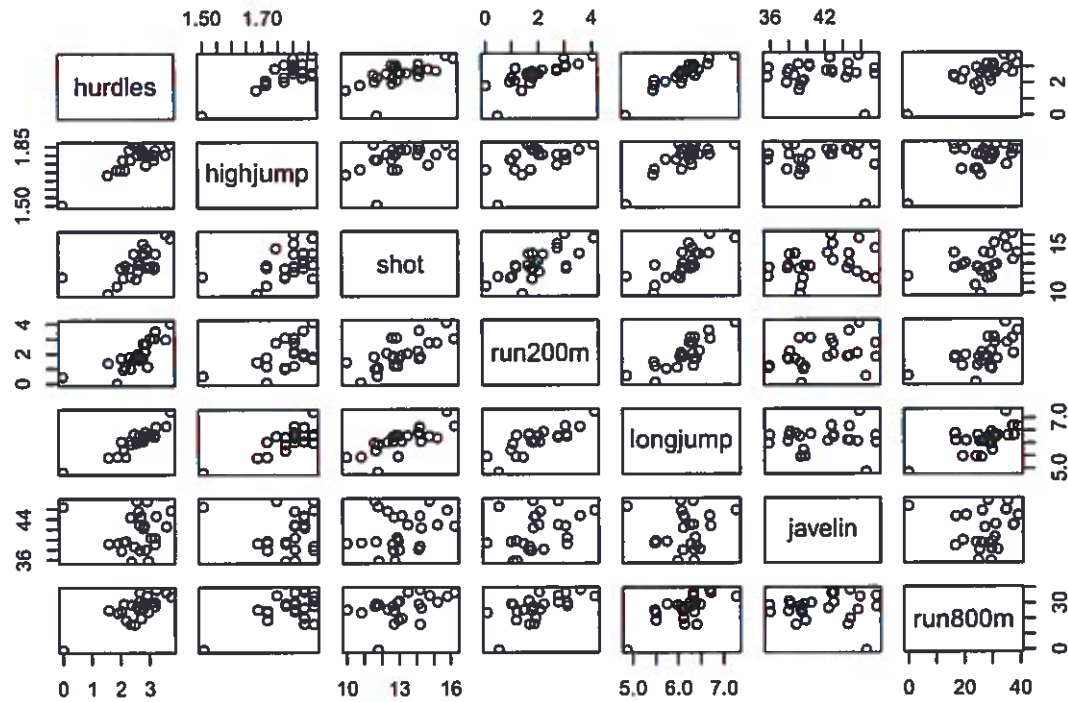
```
##           hurdles highjump shot run200m longjump javelin
## Joyner-Kersey (USA) 12.69  1.86 15.80  22.56   7.27 45.66
## John (GDR)         12.85  1.80 16.23  23.65   6.71 42.56
## Behmer (GDR)       13.20  1.83 14.20  23.10   6.68 44.54
## Sablovskaitė (URS) 13.61  1.80 15.23  23.92   6.25 42.78
## Choubenkova (URS)  13.51  1.74 14.76  23.93   6.32 47.46
## Schulz (GDR)       13.75  1.83 13.50  24.65   6.33 42.82
```

```
##           run800m score
## Joyner-Kersey (USA) 128.51 7291
## John (GDR)         126.12 6897
## Behmer (GDR)       124.20 6858
## Sablovskaitė (URS) 132.24 6540
## Choubenkova (URS)  127.90 6540
## Schulz (GDR)       125.79 6411
```

This example has a score column

```
# Readjust data so that large values are good and small values bad
# (Timed racing events are obviously the opposite)
heptathlon$hurdles <- max(heptathlon$hurdles) - heptathlon$hurdles
heptathlon$run200m <- max(heptathlon$run200m) - heptathlon$run200m
heptathlon$run800m <- max(heptathlon$run800m) - heptathlon$run800m
```

```
plot(heptathlon[, -which(names(heptathlon) %in% "score")])
```



```
# Check out the correlation
round(cor(heptathlon[, -which(names(heptathlon) %in% "score")]), digits = 3)
```

```
##           hurdles highjump shot run200m longjump javelin run800m
## hurdles   1.000   0.811 0.651  0.774   0.912  0.008  0.779
```

```
## highjump 0.811 1.000 0.441 0.488 0.782 0.002 0.591
## shot 0.651 0.441 1.000 0.683 0.743 0.269 0.420
## run200m 0.774 0.488 0.683 1.000 0.817 0.333 0.617
## longjump 0.912 0.782 0.743 0.817 1.000 0.067 0.700
## javelin 0.008 0.002 0.269 0.333 0.067 1.000 -0.020
## run800m 0.779 0.591 0.420 0.617 0.700 -0.020 1.000
```

```
summary(heptathlon[, -which(names(heptathlon) %in% "score")])
```

```
##      hurdles      highjump      shot      run200m
## Min.   :0.00   Min.   :1.500   Min.   :10.00   Min.   :0.000
## 1st Qu.:2.35   1st Qu.:1.770   1st Qu.:12.32   1st Qu.:1.380
## Median :2.67   Median :1.800   Median :12.88   Median :1.780
## Mean   :2.58   Mean   :1.782   Mean   :13.12   Mean   :1.961
## 3rd Qu.:2.95   3rd Qu.:1.830   3rd Qu.:14.20   3rd Qu.:2.690
## Max.   :3.73   Max.   :1.860   Max.   :16.23   Max.   :4.050
##      longjump      javelin      run800m
## Min.   :4.880   Min.   :35.68   Min.   : 0.00
## 1st Qu.:6.050   1st Qu.:39.06   1st Qu.:24.95
## Median :6.250   Median :40.28   Median :28.69
## Mean   :6.152   Mean   :41.48   Mean   :27.38
## 3rd Qu.:6.370   3rd Qu.:44.54   3rd Qu.:31.19
## Max.   :7.270   Max.   :47.50   Max.   :39.23
```

Scaling (using cor) is done with scale.

```
hep.pr <- prcomp(heptathlon[, -which(names(heptathlon) %in% "score")], scale. = TRUE)
```

```
attributes(hep.pr)
```

```
## $names
## [1] "sdev" "rotation" "center" "scale" "x"
##
## $class
## [1] "prcomp"
summary(hep.pr)
```

```
## Importance of components:
```

```
##          PC1    PC2    PC3    PC4    PC5    PC6
## Standard deviation 2.1119 1.0928 0.72181 0.67614 0.49524 0.27010
## Proportion of Variance 0.6372 0.1706 0.07443 0.06531 0.03504 0.01042
## Cumulative Proportion 0.6372 0.8078 0.88223 0.94754 0.98258 0.99300
##          PC7
## Standard deviation 0.2214
## Proportion of Variance 0.0070
## Cumulative Proportion 1.0000
```

Note that 2 PC's have $\lambda > 1$ (or 0.7) first 4 PC's have total variance greater than 90%. Let us look at the scree plot.

```
par(mfrow = c(1, 2))
plot(1:(length(hep.pr$sdev)), (hep.pr$sdev)^2, type='b',
main="Scree Plot", xlab="Number of Components", ylab="Eigenvalue Size")
# OR use
screepLOT(hep.pr)
```

Command of principal component
prcomp

prcomp \$ sdev :

$\sqrt{\lambda_1}$ $\sqrt{\lambda_2}$...

standard deviations
the principal components

prcomp \$ rotation : Factor loadings
the matrix of variable loadings

prcomp \$ x

center, scale

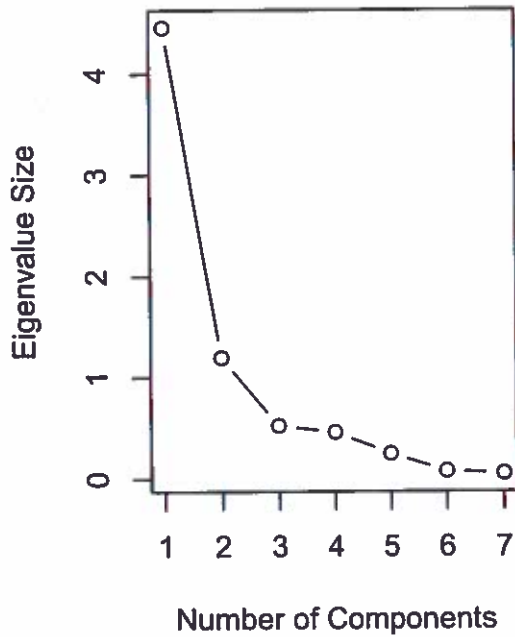
principal component

no cor argument

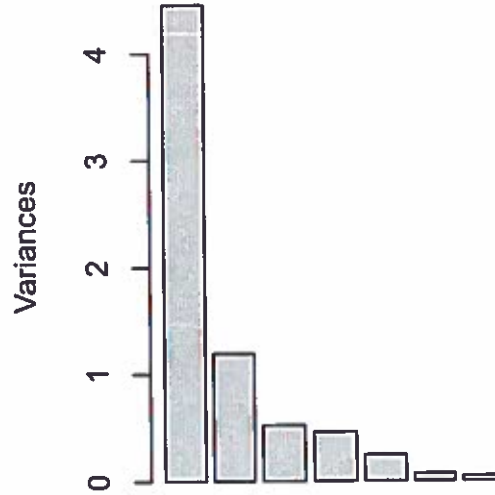
prcomp (data, cor = TRUE)
prcomp (data, scale = TRUE)

* make scree plot

Scree Plot



hep.pr



Mean and sd dev
hep.pr\$center

```
## hurdles highjump shot run200m longjump javelin run800m
## 2.5800 1.7820 13.1176 1.9608 6.1524 41.4824 27.3760
```

hep.pr\$scale

```
## hurdles highjump shot run200m longjump javelin
## 0.73664781 0.07794229 1.49188438 0.96955712 0.47421233 3.54565612
## run800m
## 8.29108809
```

a

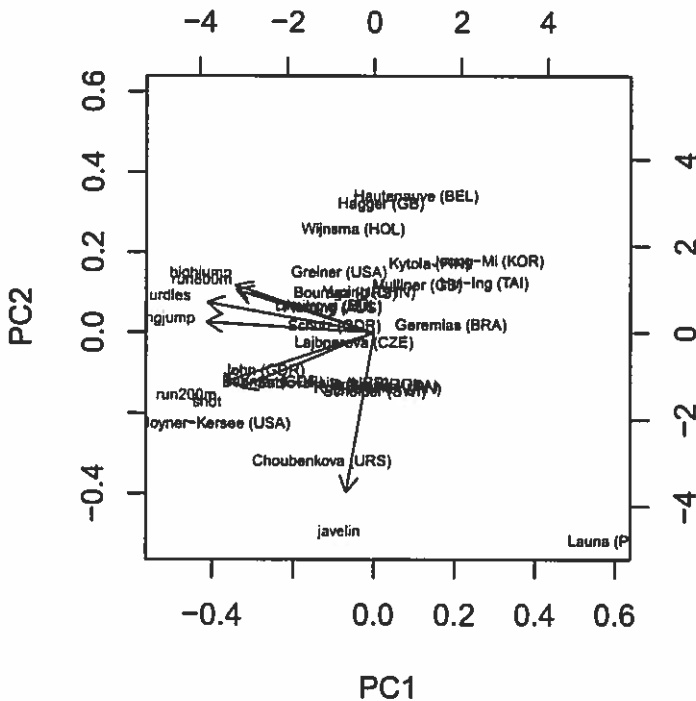
Factor loadings
hep.pr\$rotation

	PC1	PC2	PC3	PC4	PC5
## hurdles	-0.4528710	0.15792058	-0.04514996	0.02653873	-0.09494792
## highjump	-0.3771992	0.24807386	-0.36777902	0.67999172	0.01879888
## shot	-0.3630725	-0.28940743	0.67618919	0.12431725	0.51165201
## run200m	-0.4078950	-0.26038545	0.08359211	-0.36106580	-0.64983404
## longjump	-0.4562318	0.05587394	0.13931653	0.11129249	-0.18429810
## javelin	-0.0754090	-0.84169212	-0.47156016	0.12079924	0.13510669
## run800m	-0.3749594	0.22448984	-0.39585671	-0.60341130	0.50432116
##	PC6	PC7			
## hurdles	-0.78334101	0.38024707			
## highjump	0.09939981	-0.43393114			
## shot	-0.05085983	-0.21762491			
## run200m	0.02495639	-0.45338483			
## longjump	0.59020972	0.61206388			
## javelin	-0.02724076	0.17294667			

```
## run800m 0.1555520 -0.09830963
# Note that the scores with prcomp are 'x', so hep.pr$x
# Correlation of pc's to data
round(cor(heptathlon[, -which(names(heptathlon) %in% "score")], hep.pr$x), digits = 4)

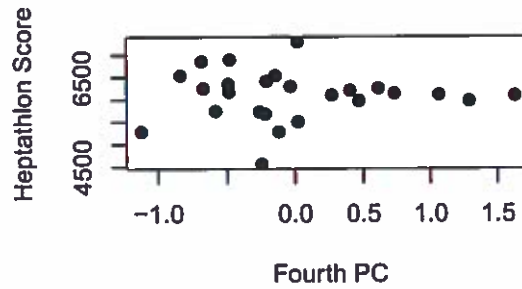
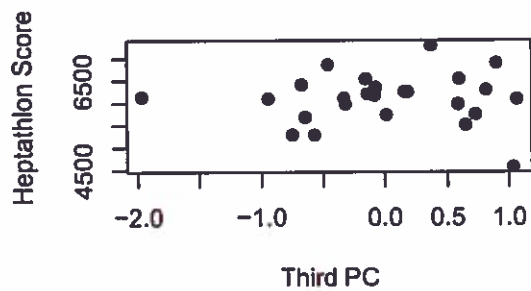
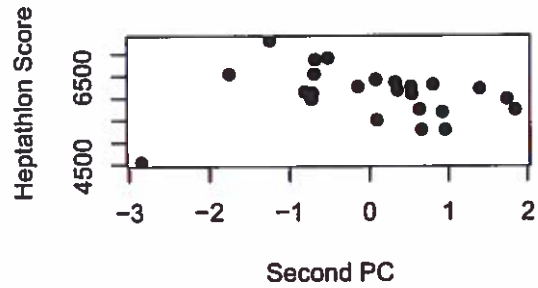
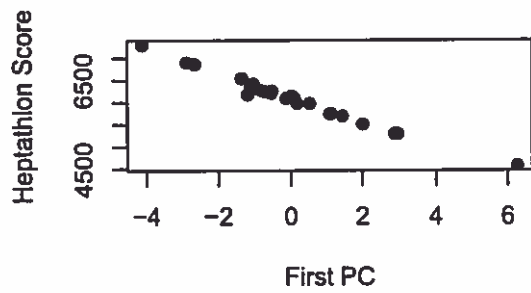
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## hurdles -0.9564  0.1726 -0.0326  0.0179 -0.0470 -0.2116  0.0842
## highjump -0.7966  0.2711 -0.2655  0.4598  0.0093  0.0268 -0.0961
## shot      -0.7668 -0.3163  0.4881  0.0841  0.2534 -0.0137 -0.0482
## run200m  -0.8614 -0.2846  0.0603 -0.2441 -0.3218  0.0067 -0.1004
## longjump  -0.9635  0.0611  0.1006  0.0752 -0.0913  0.1594  0.1355
## javelin   -0.1593 -0.9198 -0.3404  0.0817  0.0669 -0.0074  0.0383
## run800m   -0.7919  0.2453 -0.2857 -0.4080  0.2498  0.0420 -0.0218

par(mfrow = c(1, 1))
biplot(hep.pr, cex = 0.5)
```



Compare the first PC to the overall score. Note that the first is highly correlated with the score with rapidly declining association

```
par(mfrow = c(2, 2))
plot(hep.pr$x[,1], heptathlon[, "score"], pch = 19, xlab = "First PC", ylab = "Heptathlon Score")
plot(hep.pr$x[,2], heptathlon[, "score"], pch = 19, xlab = "Second PC", ylab = "Heptathlon Score")
plot(hep.pr$x[,3], heptathlon[, "score"], pch = 19, xlab = "Third PC", ylab = "Heptathlon Score")
plot(hep.pr$x[,4], heptathlon[, "score"], pch = 19, xlab = "Fourth PC", ylab = "Heptathlon Score")
```



```
cor(hep.pr$x, heptathlon$score)
```

```
##           [,1]
## PC1 -0.991097775
## PC2 -0.097885776
## PC3 -0.005162860
## PC4  0.005161233
## PC5  0.045780044
## PC6  0.030648671
## PC7  0.006510970
```

Example: Bumpus bird data

The data contains the body measurements of female sparrows:

- X1=total length
- X2=alar length
- X3=length of beak and head
- X4=length of humerus
- X5=length of keel and sternum; all in mm)

Birds 1 to 21 survived a severe storm near Brown University in Rhode Island while the remainder died. (Original source Bumpus 1898.)

```
bumpbird=read.table("bumpusbird.txt", header=T) #"bumpusbird.txt"
names(bumpbird)
```

```
## [1] "ID" "X1" "X2" "X3" "X4" "X5"
```

```

names(bumpbird) = c("ID", "tot.length", "alar.length",
                    "beak.head.length", "humerus.length",
                    "keel.stern.length")

attach(bumpbird)
bird.pc=princomp(bumpbird[, -1], cor=T)
summary(bird.pc, loadings=T)

## Importance of components:
##              Comp.1   Comp.2   Comp.3   Comp.4   Comp.5
## Standard deviation   1.9015726 0.7290433 0.62163056 0.5491498 0.4056199
## Proportion of Variance 0.7231957 0.1063008 0.07728491 0.0603131 0.0329055
## Cumulative Proportion 0.7231957 0.8294965 0.90678139 0.9670945 1.0000000
##
## Loadings:
##              Comp.1 Comp.2 Comp.3 Comp.4 Comp.5
## tot.length      -0.452      0.690  0.420 -0.374
## alar.length     -0.462 -0.300  0.341 -0.548  0.530
## beak.head.length -0.451 -0.325 -0.454  0.606  0.343
## humerus.length  -0.471 -0.185 -0.411 -0.388 -0.652
## keel.stern.length -0.398  0.876 -0.178      0.192

```

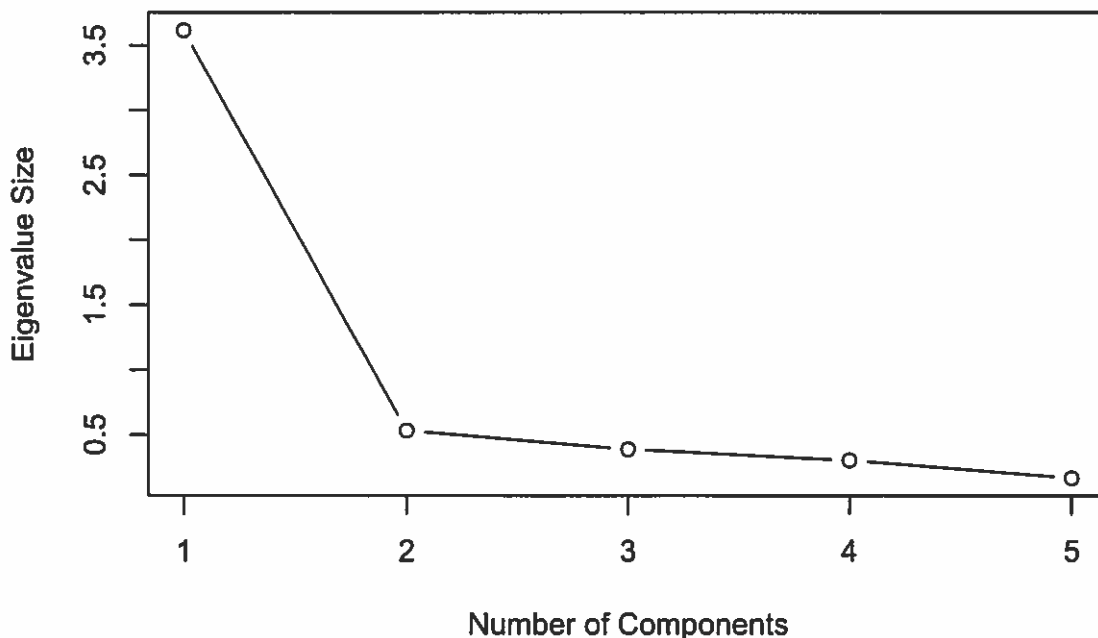
The first component seems to be a very general size component, but the second component really picks up the variability in keel and sternum length. Let's make a scree plot

```

plot(1:(length(bird.pc$sdev)), (bird.pc$sdev)^2, type='b',
     main="Scree Plot", xlab="Number of Components", ylab="Eigenvalue Size")

```

Scree Plot

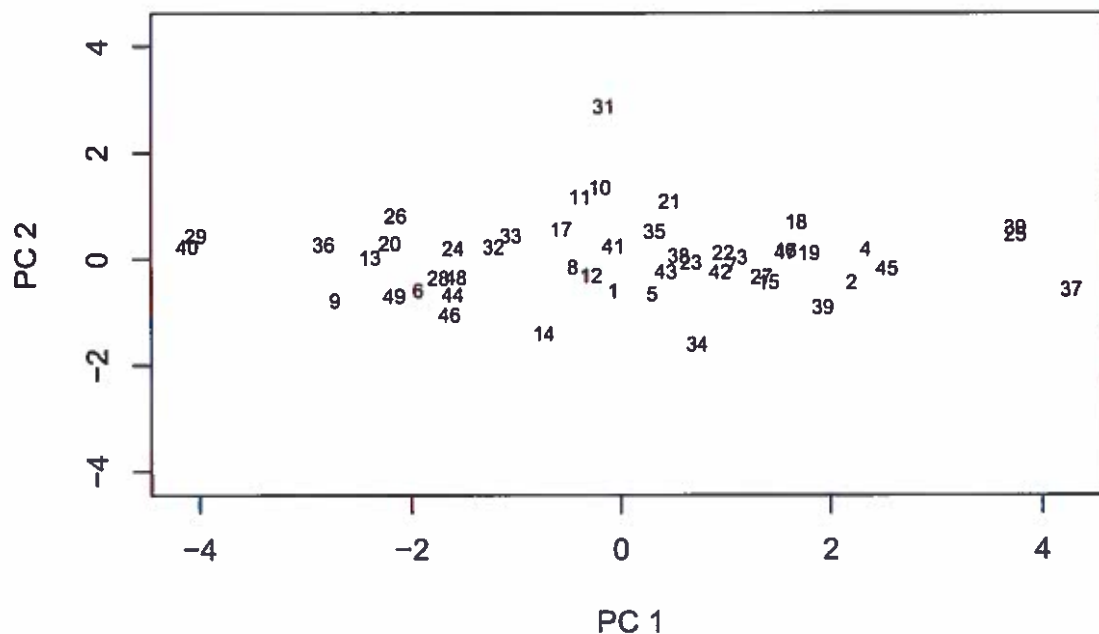


Where does the "elbow" occur? What seems to be a reasonable number of PCs to use? We plot the PC scores for the sample data in the space of the first two principal components:

```

plot(bird.pc$scores[,1], bird.pc$scores[,2], ylim=range(bird.pc$scores[,1]),
xlab="PC 1", ylab="PC 2", type='n', lwd=2)
# labeling points with IDs for birds:
text(bird.pc$scores[,1], bird.pc$scores[,2], labels=ID, cex=0.7, lwd=2)

```

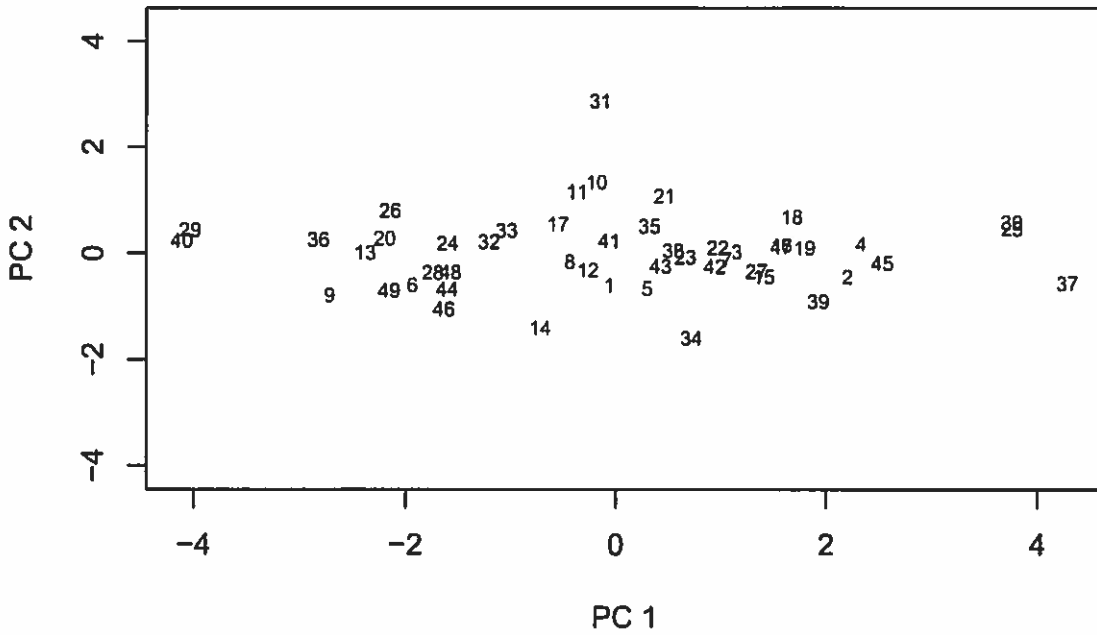


Interestingly, birds 1 to 21 survived a storm in RI while birds 22 to 49 died. Let's see the survivors (red) and the deceased (blue) separated by color.

```

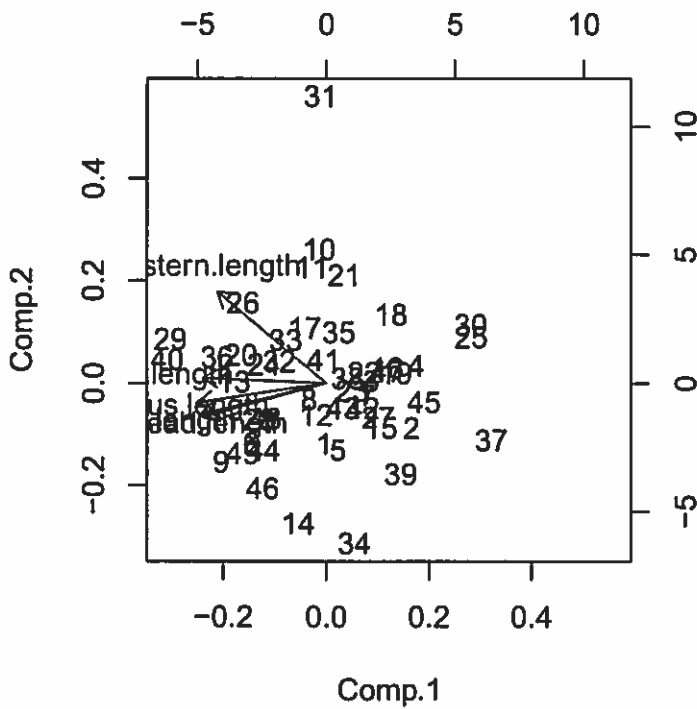
plot(bird.pc$scores[,1], bird.pc$scores[,2], ylim=range(bird.pc$scores[,1]),
xlab="PC 1", ylab="PC 2", type='n', lwd=2)
# labeling points with IDs for birds:
text(bird.pc$scores[,1], bird.pc$scores[,2], labels=ID, cex=0.7, lwd=2,
col=c(rep("red", times = 21), rep("blue", times=28)) )

```

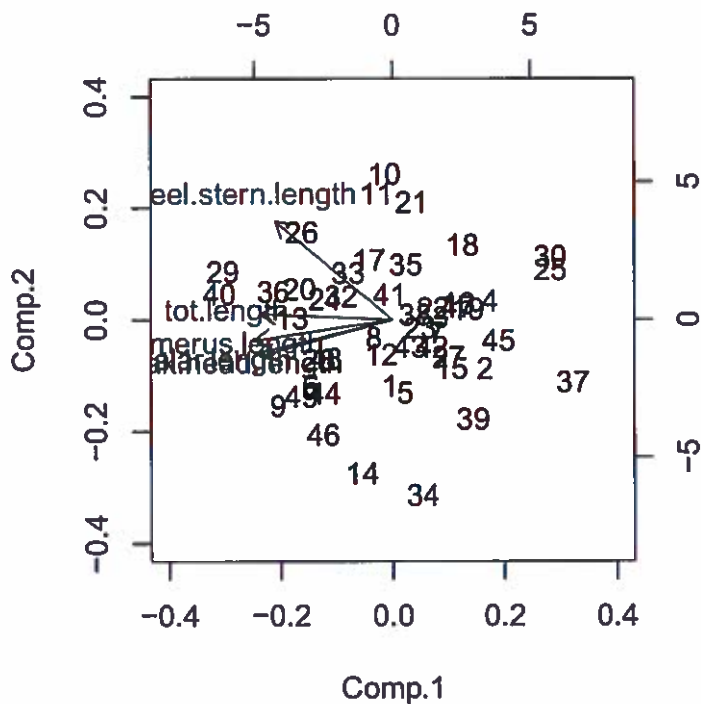



The medium-sized birds seem to have had the best survival chances. The biplot can add information about the variables to the plot of the first two PC scores:

```
biplot(bird.pc,xlabs=ID)
```



```
# Expanding the plotting window a bit:
biplot(bird.pc,xlabs=ID, xlim=c(-0.4,0.4), ylim=c(-0.4,0.4))
```



4. Statistical Inference with Principal Components

The principal components solutions are based on the eigenvalues and eigenvectors of the sample covariance matrix S . The eigenvalues and eigenvectors of the true covariance matrix Σ are unknown. We can (for large samples) perform inference about these unknown eigenvalues and eigenvectors. These results assume the data vectors X_1, X_2, \dots, X_n are a random sample from a multivariate normal distribution. This assumption could be checked with plots such as the chi-squared plot. If this multivariate normal assumption is close to correct, the inferences will be approximately valid.

Let the variance of the i^{th} population principal component (which is the i^{th} eigenvalue of the true covariance matrix Σ) be denoted by ξ_i . Under multivariate normality, the large-sample distribution of the variance of the i^{th} (sample) principal component λ_i is $N(\xi_i, 2\xi_i^2/n)$. Furthermore, the λ_i 's are independent. Thus a large-sample $100(1 - \alpha)\%$ confidence interval for ξ_i is

$$\left(\frac{\lambda_i}{1 + z_{\alpha/2}\sqrt{2/n}}, \frac{\lambda_i}{1 - z_{\alpha/2}\sqrt{2/n}} \right).$$

If we want simultaneous intervals for m such eigenvalues, replace $z_{\alpha/2}$ with $z_{\alpha/2m}$ in the formula to get Bonferroni confidence intervals.

We can write a R function to calculate the confidence intervals for eigenvalue(s) of Sigma:

```

evalue.CI=function(datamatrix, labelvec, m=length(labelvec), conf.level=0.95)
{
  alpha=1-conf.level
  n=nrow(datamatrix)
  z=qnorm(alpha/(2*m),lower=F)
  lambdas<- princomp(datamatrix)$sdev^2
  print(lambdas)
  LCL=lambdas[labelvec]/(1+z*sqrt(2/n))

```

```

UCL=lambdas[labelvec]/(1-z*sqrt(2/n))
CIs=cbind(LCL,UCL)
return(CIs)
}

```

In the bird example, the 95% CI for variance of the first population principle component is

```

evaluate.CI(bumpbird[,-1],1)

```

```

##      Comp.1      Comp.2      Comp.3      Comp.4      Comp.5
## 34.60482313  4.52812344  0.61804191  0.30640029  0.07593901

##           LCL           UCL
## Comp.1 24.78904 57.29015

```

The joint 95% CIs for variances of the first two population principle components are

```

evaluate.CI(bumpbird[,-1],c(1,2))

```

```

##      Comp.1      Comp.2      Comp.3      Comp.4      Comp.5
## 34.60482313  4.52812344  0.61804191  0.30640029  0.07593901

##           LCL           UCL
## Comp.1 23.818879 63.243476
## Comp.2  3.116757  8.275559

```

Note that the first interval is wider because the FAMILYWISE confidence coefficient is 95% here.

We can also calculate the 99% CI for variance of the first population PC:

```

evaluate.CI(bumpbird[,-1],1, conf.level=0.99)

```

```

##      Comp.1      Comp.2      Comp.3      Comp.4      Comp.5
## 34.60482313  4.52812344  0.61804191  0.30640029  0.07593901

##           LCL           UCL
## Comp.1 22.7604 72.15292

```

Check multivariate normality

#Function will produce a chi-square plot to test multivariate normality

```

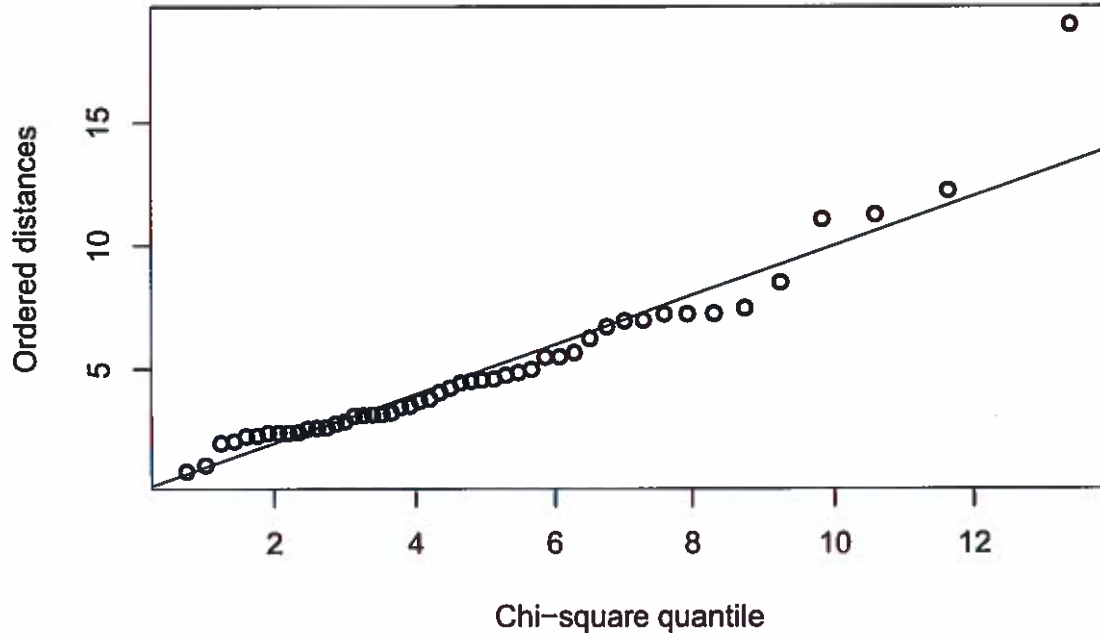
chisplot <- function(x) {
  if (!is.matrix(x)) stop("x is not a matrix")

  ### determine dimensions
  n <- nrow(x)
  p <- ncol(x)
  #
  xbar <- apply(x, 2, mean)
  S <- var(x)
  S <- solve(S)
  index <- (1:n)/(n+1)
  #
  xcent <- t(t(x) - xbar)
  di <- apply(xcent, 1, function(x,S) x %*% S %*% x,S)
  #
  quant <- qchisq(index,p)
  plot(quant, sort(di), ylab = "Ordered distances",
       xlab = "Chi-square quantile", lwd=2,pch=1)
  abline(a=0,b=1)
}

```

}

```
chisplot(as.matrix(bumpbird[,-1]))
```



Impossible to measure the quantities

C5 Factor Analysis

Jianzuan Liu
Fall 2018

measured quantities ← hidden quantities
observable quantities exhibit relationship ← underlying latent

It is not always possible to measure the quantities of interest directly (e.g. intelligence quotient). Behavioral scientist Charles Spearman is credited with being the originator and pioneer of the classical theory of mental tests, the theory of intelligence and what is now called Factor Analysis. The appeal of Factor Analysis lies in the ease of use and the recognition that there is an association between the hidden quantities and the measured quantities. The aim of Factor Analysis is:

1. to exhibit the relationship between the measured and the underlying variables;
2. to estimate the underlying variables, called the hidden or latent variables.

FA and PCA have similar themes, i.e., to explain covariation between variables via linear combinations of other variables. However, there are distinctions between the two approaches:

- FA assumes a statistical model that describes covariation in observed variables via linear combinations of latent variables
- PCA finds uncorrelated linear combinations of observed variables that explain maximal variance (no latent variables here)
- FA refers to a statistical model, whereas PCA refers to the eigenvalue decomposition of a covariance (or correlation) matrix.

1. The Orthogonal Factor Model

Let the observable random vector X , with p components, has mean μ and covariance matrix Σ . The factor model postulates that X is linearly dependent upon a few unobservable random variables F_1, F_2, \dots, F_m , called common factors (or factor scores or vector of latent or hidden variables). In particular, the factor analysis model is unobservable $X = LF + \mu + \epsilon$.

- L is a $p \times m$ linear transformation matrix, called the factor loadings. λ_{jk} is the loading of the j^{th} variable on the k^{th} common factor

- $F = (F_1, \dots, F_m)^T$ where F_k is the score on the k^{th} common factor (Assumption 1)

$$F = \begin{pmatrix} F_1 \\ F_2 \\ \vdots \\ F_m \end{pmatrix}$$

$$E(F) = 0$$

$$Cov(F) = E(FF^T) = I$$

- $\epsilon = (\epsilon_1, \dots, \epsilon_p)^T$ is a p -dimensional random vector of latent error terms such that (Assumption 2)

$$E(\epsilon) = 0$$

$$Cov(\epsilon) = E(\epsilon\epsilon^T) = \Psi = \begin{bmatrix} \psi_1 & 0 & \dots & 0 \\ 0 & \psi_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \psi_p \end{bmatrix}$$

- If F and ϵ are uncorrelated, (Assumption 3)

$$Cov(\epsilon, F) = E(\epsilon F^T) = 0$$

$$X = LF + \mu + \varepsilon$$

The covariance structure

$$\Sigma = \text{Cov}(X) = LL^T + \Psi$$

$$\begin{aligned} \text{Cov}(X) &= E((X - \mu)(X - \mu)^T) \\ &= E((LF + \varepsilon)(LF + \varepsilon)^T) \\ &= E(LFF^T L^T + LF\varepsilon^T + \varepsilon LF + \varepsilon\varepsilon^T) \\ &= L E(FF^T) L^T + \underbrace{L E(F\varepsilon^T)}_{=0} + \underbrace{E(\varepsilon F^T) L^T}_{=0} + E(\varepsilon\varepsilon^T) \\ &= L \underbrace{E(FF^T)}_{\substack{\text{non orthogonal} \\ \text{if } \text{Cov}(F) = \Phi, \text{ then } \Sigma = L\Phi L^T + \Psi}} L^T + \Psi \end{aligned}$$

This implies that the covariance between X and F has the form

$$\text{Cov}(X, F) = E((X - \mu)F^T) = E((LF + \varepsilon)F^T) = L$$

The portion of variance of the i^{th} variable that is explained by the m common factors is called the **communality** of the i^{th} variable:

$$\underbrace{\sigma_{ii}}_{\text{Var}(X_i)} = \underbrace{h_i^2}_{\text{communality}} + \underbrace{\psi_i}_{\text{uniqueness specific variance}}$$

is the proportion of variance of X_i explained by the m common factor \leftarrow If X_i is informative \Rightarrow communality is high.

where

- σ_{ii} is the variance of X_i (i.e., the i^{th} diagonal of Σ)
- $h_i^2 = (LL^T)_{ii} = (l_{i1}^2 + l_{i2}^2 + \dots + l_{im}^2)$ is the **communality** of X_i . Note that the communality h_i^2 is the sum of squared loadings for X_i
- ψ_i is the **specific variance** (or **uniqueness**) of X_i

Example *Idea of this example Given $\text{Cov}(X, \Sigma)$, we want to find L*

We consider a one-factor model. The factor might be physical fitness, which we could obtain from the performance of two different sports, or the talent for language, which might arise from oral and written communication skills. We assume that the two-dimensional random vector X has mean zero and covariance matrix

$$\Sigma = \begin{pmatrix} 1.25 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$$

Calculate the communality and uniqueness for a one-factor model.

Minimum variance is the preferred choice

$$\begin{bmatrix} l_1 \\ l_2 \end{bmatrix} =$$

saved

$$\tilde{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} l_1 \\ l_2 \end{bmatrix} F + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad F \text{ is a scalar}$$

$$\begin{aligned} \Sigma &= LL^T + \Phi \\ \sigma_{11} &= l_1^2 + \phi_{11} \\ \sigma_{22} &= l_2^2 + \phi_{22} \\ \sigma_{12} &= l_1 l_2 \end{aligned} \quad \left. \begin{aligned} \text{Cov}(X_1, X_2) &= l_{11} l_{21} + \dots + l_{1m} l_{2m} \\ 1.25 = \sigma_{11} &> l_1^2 \\ 0.5 = \sigma_{12} &> l_1 l_2 \\ \sigma_{12} = l_1 l_2 &= 0.5 \end{aligned} \right\} \Rightarrow \begin{cases} l_1 = 1 \\ l_2 = 0.5 \\ \text{or } l_1 = \frac{3}{4} \\ l_2 = \frac{2}{3} \end{cases}$$

Consider $\begin{cases} l_1 = 1 \\ l_2 = 0.5 \end{cases}$

Sol. $\begin{bmatrix} l_1 \\ l_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \Rightarrow LL^T = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 0.25 \end{bmatrix}$

$\Sigma = LL^T + \Phi \quad \Phi = E(\epsilon)$

$\text{tr}(\Phi) = 0.25 + 0.25 = 0.5$

$\begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 0.25 \end{bmatrix} + \begin{bmatrix} 0.25 & 0 \\ 0 & 0.25 \end{bmatrix}$

$\begin{bmatrix} l_1^* \\ l_2^* \end{bmatrix} = \begin{bmatrix} 3/4 \\ 2/3 \end{bmatrix} \quad LL^T = \begin{bmatrix} 9/16 & 1/2 \\ 1/2 & 4/9 \end{bmatrix} \Rightarrow \begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} = \begin{bmatrix} 9/16 & 1/2 \\ 1/2 & 4/9 \end{bmatrix} + \begin{bmatrix} 1/16 & 0 \\ 0 & 1/18 \end{bmatrix} \Rightarrow \text{tr}(\Phi) = \frac{11}{16} + \frac{1}{18} = 0.74$

Example: Chicken-bone data (from Example 9.14 in Johnson & Wichern)

For $l_1^* = \frac{3}{4} \quad l_2^* = \frac{2}{3}$
 $1.25 = \left(\frac{3}{4}\right)^2 + \frac{1}{16}$
 \uparrow
 $\text{var}(X_1)$ communality u

The correlation matrix for chicken-bone measurements.

```
var.names <- c("SkullL", "SkullB", "FemurL", "TibiaL", "HumerusL", "UlnaL")
CorMat <- read.table("Wichern_data/E9-14.dat", fill = TRUE, col.names = var.names)
row.names(CorMat) <- var.names
n <- dim(CorMat)
# fill the upper triangle
# for(j in 1:(n-1)){CorMat[j, (j+1):n] <- CorMat[(j+1):n, j]}
CorMat <- as.matrix(CorMat)
```

$0.5 = \left(\frac{2}{3}\right)^2 + \frac{1}{18}$
 $\text{var}(X_2)$ communality u unique

The estimated factor loadings were extracted by the maximum likelihood procedure

```
# Enter the estimated factor loadings 2 factors model
L <- matrix(c(0.602, 0.200, 0.467, 0.154, 0.926, 0.143,
             1.000, 0.000, 0.874, 0.476, 0.894, 0.327), ncol=2, byrow=T)
as.data.frame(cbind(var.names, L))
```

```
## var.names V2 V3
## 1 SkullL 0.602 0.2
## 2 SkullB 0.467 0.154
## 3 FemurL 0.926 0.143
## 4 TibiaL 1 0
## 5 HumerusL 0.874 0.476
## 6 UlnaL 0.894 0.327
```

```
# the specific variances
psi <- diag(diag(CorMat - L %*% t(L)))
round(diag(psi), digits = 3)
```

```
## [1] 0.598 0.758 0.122 0.000 0.010 0.094 ← uniqueness!
```

```
# the communalities
hi2 <- apply(L, 1, function(x) sum(x^2))
round(hi2, digits = 3)
```

```
## [1] 0.402 0.242 0.878 1.000 0.990 0.906
```

\uparrow
~~uniqueness~~

\uparrow
~~for the load~~

2. Factor Loadings

The main two approaches to estimating the factor loadings divide into nonparametric methods, and methods which rely on the normality of the data. For normal data, we expect the latter methods to be better; in practice, methods based on assumptions of normality still work well if the distribution of the data does not deviate too much from the normal distribution.

A tool for finding the factors is the covariance matrix. In a parametric framework, there are two main methods for determining the factors:

- ~~Principal Component Factor Analysis~~
- ~~Maximum likelihood method~~

① Principal Component Factor Analysis < to estimate the factor

Note that the parameters of interest are the factor loadings L and specific variances on the diagonal of Ψ . Let $X \sim N_p(\mu, \Sigma)$. For $m < p$ common factors, the PCA solution estimates L and Ψ as

$$\hat{L} = [\lambda_1^{1/2} \mathbf{v}_1, \lambda_2^{1/2} \mathbf{v}_2, \dots, \lambda_m^{1/2} \mathbf{v}_m]$$

$$\hat{\Psi}_i = \sigma_{ii} - \hat{h}_i^2$$

where $\Sigma = V\Lambda V^T$ is the spectral decomposition of Σ . The communality of the i^{th} variable is estimated as

$$\hat{h}_i^2 = \sum_{k=1}^m \hat{l}_{ik}^2$$

Proportion of total sample variance explained by the k-th factor is

$$R_k^2 = \frac{\sum_{i=1}^p \hat{l}_{ik}^2}{\sum_{i=1}^p \sigma_{ii}} = \frac{(\lambda_k^{1/2} \mathbf{v}_k)^T (\lambda_k^{1/2} \mathbf{v}_k)}{\sum_{i=1}^p \sigma_{ii}} = \frac{\lambda_k}{\sum_{i=1}^p \sigma_{ii}}$$

$k = 1, m$ ← to determine how many factors that we want to keep.

the proportion of variance explained by each factor

```
VxF <- apply(L, 2, function(x) sum(x^2))
```

```
totVar <- sum(diag(CorMat))
```

```
propVxF <- VxF / totVar
```

```
round(propVxF, digits = 3)
```

```
## [1] 0.667 0.070
```

← 1st factor explains 71% of the sample variance.

```
# residual matrix R = LL^T - Psi
```

```
Rmat <- CorMat - L %*% t(L) - psi
```

```
round(Rmat, digits = 4)
```

```
##          SkullL SkullB FemurL TibiaL HumerusL UlnaL
## SkullL  0.0000      NA      NA      NA      NA      NA
## SkullB  0.1931  0.0000      NA      NA      NA      NA
## FemurL -0.0171 -0.0325  0.0000      NA      NA      NA
## TibiaL  0.0000  0.0000  0.0000      0      NA      NA
## HumerusL -0.0003  0.0005 -0.0004      0      0      NA
## UlnaL   -0.0006 -0.0179  0.0034      0      0      0
```

iterative method

② Principal Factor Analysis or Principal Axis Factoring (modified method 2(1))

Principal Factor Analysis or Principal Axis Factoring is a modified approach. In Principal Component Factor Analysis, the underlying covariance matrix is Σ . In contrast, Principal Factor Analysis is based on the scaled covariance matrix of the common factors. Assume we are applying Factor Analysis to a sample correlation matrix R

$$R - \Psi = LL^T$$

and we have some initial estimate of the specific variance $\hat{\psi}_i$. We can use $\hat{\psi}_i = 1/r^{ii}$ where r^{ii} is the i^{th} diagonal of R^{-1} . The iterated principal axis factoring algorithm:

1. Form $\tilde{R} = R - \hat{\Psi}$ given current $\hat{\psi}_i$ estimates
2. Update $\tilde{L} = [\tilde{\lambda}_1^{1/2} \tilde{v}_1, \tilde{\lambda}_2^{1/2} \tilde{v}_2, \dots, \tilde{\lambda}_m^{1/2} \tilde{v}_m]$ where $\tilde{R} = \tilde{V} \tilde{\Lambda} \tilde{V}^T$ is the eigenvalue decomposition of \tilde{R}
3. Update $\hat{\psi}_i = 1 - \sum_{k=1}^m \tilde{l}_{ik}^2$

③ Maximum Likelihood Method for Factor Analysis

(MLE of multivariate normal distribution) Let X_1, X_2, \dots, X_n be a random sample from a normal population with mean μ and covariance Σ . Then

$$\hat{\mu} = \bar{x}, \text{ and } \hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^T = \frac{n-1}{n} S$$

are the maximum likelihood estimators of μ and Σ , respectively. Proof.

$X \sim N(\mu, \Sigma)$

* The likelihood

$$L(\mu, \Sigma | X) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} \ln \left[\Sigma^{-1} \left(\sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^T \right) + n(\bar{X} - \mu)(\bar{X} - \mu)^T \right] \right\}$$

Σ^{-1} is positive definite, so $(X_j - \bar{X}) \Sigma^{-1} (X_j - \bar{X})^T > 0$ unless $\mu = \bar{X}$

Thus the likelihood $L(\mu, \Sigma)$ is maximized w.r.t μ at $\hat{\mu} = \bar{x}$

Let X_1, X_2, \dots, X_n be a random sample from $N_p(\mu, LL^T + \Psi)$, where $\Sigma = LL^T + \Psi$. The maximum likelihood estimators \hat{L} , $\hat{\Psi}$ and $\hat{\mu}$ maximize the log-likelihood function for a sample of n observations

$$\mathcal{L}(\mu, L, \Psi) = \frac{np \log(2\pi)}{2} + \frac{n \log(|\Sigma^{-1}|)}{2} - \frac{\sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}{2}$$

subject to $\hat{L}^T \hat{\Psi}^{-1} \hat{L}$ being diagonal. Use an iterative algorithm to maximize \mathcal{L} . The maximum likelihood estimates of the communalities are

$$\hat{h}_i^2 = \hat{l}_{i1}^2 + \hat{l}_{i2}^2 + \dots + \hat{l}_{im}^2, i = 1, 2, \dots, p.$$

The proportion of total sample variance due to j^{th} factor is

$$\frac{\hat{l}_{1j}^2 + \hat{l}_{2j}^2 + \dots + \hat{l}_{pj}^2}{s_{11} + s_{22} + \dots + s_{pp}}$$

Example: Stock Price Data $p=5$ $m=2$

Data are weekly gains in stock prices for 100 consecutive weeks for five companies: Allied Chemical, Du Pont, Union Carbide, Exxon and Texaco. Note that the first three are chemical companies and the last two are oil companies. Fit an $m = 2$ orthogonal factor model

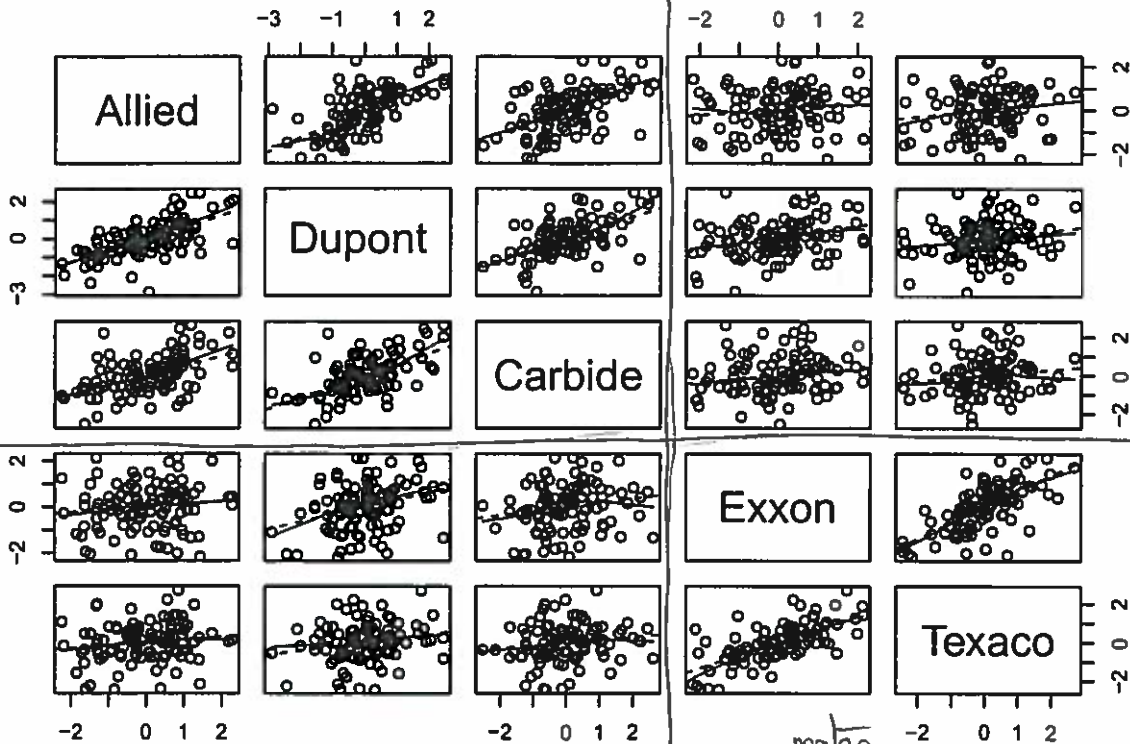
```
stocks <- read.table("stocks.txt", col.names = paste("x", 1:5, sep = ""))
names(stocks)=c("Allied", "Dupont", "Carbide", "Exxon", "Texaco") (5 companies)
head(stocks);dim(stocks)
```

	chemical company			oil companies	
##	Allied	Dupont	Carbide	Exxon	Texaco
## 1	0.0130338	-0.0078431	-0.0031889	-0.0447693	0.0052151
## 2	0.0084862	0.0166886	-0.0062100	0.0119560	0.0134890
## 3	-0.0179153	-0.0086393	0.0100360	0.0000000	-0.0061428
## 4	0.0215589	-0.0034858	0.0174353	-0.0285917	-0.0069534
## 5	0.0108225	0.0037167	-0.0101345	0.0291900	0.0409751
## 6	0.0101713	-0.0121978	-0.0083768	0.0137083	0.0029895

```
## [1] 103 5
# plot standardized data on pairwise scatter plots
stockss <- scale(stocks, center = T, scale = T)
pairs(stockss, labels = c("Allied", "Dupont", "Carbide", "Exxon", "Texaco"),
      panel = function(x, y) { panel.smooth(x, y)
        abline(lsfitted(x, y), lty = 2)})
```

positive

hard to say the relationship



```
# perform factor analysis with two factor
stocks.fac = factanal(x = stocks, factors = 2,
  scores = "regression",
  rotation = "varimax",
  method = "mls")
```

Another function may be fa

```
stocks.fac
##
## Call:
## factanal(x = stocks, factors = 2, scores = "regression", rotation = "varimax", method = "mls")
##
## Uniquenesses:
## Allied Dupont Carbide Exxon Texaco
## 0.417 0.275 0.542 0.005 0.530
##
## Loadings: Factor 1 Factor 2
## Allied 0.763 0.232
## Dupont 0.819 0.108
## Carbide 0.668 0.991
## Exxon 0.113 0.677
## Texaco 0.108
##
## SS loadings 1.725 1.507
## Proportion Var 0.345 0.301
## Cumulative Var 0.345 0.646
##
## Test of the hypothesis that 2 factors are sufficient.
## The chi square statistic is 1.97 on 1 degree of freedom.
## The p-value is 0.16
```

the first factor is affected by the first 3 companies

← sum of square of the loading (don't care!)
 ← R_i in page 4
 just cumulative

} This only have us
 decide # of factors

p-value = 0.167 → big ⇒ → we see that since
 6 factors is eqn enough.

The proportions of total variance accounted by the first and second factors are 0.345 and 0.301.

The first factor is dominated by chemical companies while the second factor by oil companies

Now we calculate the residuals:

residual
LL^T - Φ
`stocks.pred <- stocks.fac$loadings%*% t(stocks.fac$loadings) + diag(stocks.fac$uniquenesses)`
`residuals <- var(stocks) - stocks.pred`

Now list factor scores

`stocks.fac$scores`



Example: Air pollution data

Use air pollution variables X_1, X_2, X_5 and X_6 . Compare the factorization obtained by the principal component and the maximum likelihood methods for a factor model with $m = 1$ and $m = 2$.

`library(psych)`

`## Warning: package 'psych' was built under R version 3.4.4`

```
air <- read.table("Wichern_data/T1-5.dat")
names(air) <- c("Wind", "Solar", "CO", "NO", "NO2", "O3", "HC")
air.4 <- air[,c(1,2,5,6)]
```

```
# principal component solution to a factor model with m = 1, no rotation
air.pc1 <- principal(air.4, nfactors = 1, covar = TRUE, rotate = "none")
# MLE of L and \Psi for m = 1 no rotation
air.fa1 <- fa(air.4, nfactors = 1, covar = TRUE, fm = "ml", rotate = "none")
```

```
# principal component solution to a factor model with m = 1, varimax
air.pc1.v <- principal(air.4, nfactors = 1, covar = TRUE, rotate = "varimax")
# MLE of L and \Psi for m = 1, varimax
air.fa1.v <- fa(air.4, nfactors = 1, covar = TRUE, fm = "ml", rotate = "varimax")
```

```
# principal component solution to a factor model with m = 2, no rotation
air.pc2 <- principal(air.4, nfactors = 2, covar = TRUE, rotate = "none")
# MLE of L and \Psi for m = 2 no rotation
air.fa2 <- fa(air.4, nfactors = 2, covar = TRUE, fm = "ml", rotate = "varimax")
```

```
# principal component solution to a factor model with m = 1, varimax
air.pc2.v <- principal(air.4, nfactors = 2, covar = TRUE, rotate = "varimax")
# MLE of L and \Psi for m = 1, varimax
air.fa2.v <- fa(air.4, nfactors = 2, covar = TRUE, fm = "ml", rotate = "varimax")
```

	PC			ML		
	PC1	ψ_i	h_i^2	ML1	ψ_i	h_i^2
Wind	-0.1750	2.469	0.031	-0.1747	2.469	0.031
Solar	17.3247	0.371	300.145	17.2962	1.358	299.158
NO2	0.4214	11.186	0.178	0.4207	11.187	0.177
O3	1.9587	27.142	3.837	1.9555	27.154	3.824

(a) Principal component and ML solutions for 1 factor with no rotation

	PC				ML			
	PC1	PC2	ψ_i	h_i^2	ML1	ML2	ψ_i	h_i^2
Wind	-0.1750	-0.4048	2.306	0.194	-0.1747	-0.3975	2.311	0.189
Solar	17.3247	-0.6086	0.001	300.515	17.2962	-0.5976	1.001	299.515
NO2	0.4214	0.7422	10.635	0.728	0.4207	0.7288	10.655	0.708
O3	1.9587	5.1867	0.239	30.739	1.9555	5.0932	1.213	29.765

(b) Principal component and ML solutions for 2 factors with no rotation

	PC				ML		
	PC1	ψ_i	h_i^2		ML1	ψ_i	h_i^2
Wind	-0.1750	2.469	0.031		-0.1747	2.469	0.031
Solar	17.3247	0.371	300.145		17.2962	1.358	299.158
NO2	0.4214	11.186	0.178		0.4207	11.187	0.177
O3	1.9587	27.142	3.837		1.9555	27.154	3.824

(c) Principal component and ML solutions for 1 factor with varimax rotation

	PC				ML			
	RC1	RC2	ψ_i	h_i^2	ML1	ML2	ψ_i	h_i^2
Wind	-0.0758	-0.4344	2.306	0.194	-0.0829	-0.4262	2.311	0.189
Solar	16.9895	3.4457	0.001	300.515	17.0038	3.2225	1.001	299.515
NO2	0.2368	0.8200	10.635	0.728	0.2500	0.8035	10.655	0.708
O3	0.6960	5.5004	0.239	30.739	0.7870	5.3987	1.213	29.765

(d) Principal component and ML solutions for 2 factors with varimax rotation

Table 1: Principal component and maximum likelihood solutions with and without varimax with 1 and 2 factors (factor loadings, uniquenesses, and communalities)

	PC	ML
	Factor 1	Factor 1
SS loadings	304.1895	303.1895
Proportion Var	0.8808	0.8779
Cumulative Var	0.8808	0.8779

(a) Principal component and ML solutions for 1 factor with no rotation

	PC		ML	
	Factor 1	Factor 2	Factor 1	Factor 2
SS loadings	304.1895	27.9874	303.1895	26.9874
Proportion Var	0.8808	0.0810	0.8779	0.0781
Cumulative Var	0.8808	0.9618	0.8779	0.9560

(b) Principal component and ML solutions for 2 factors with no rotation

	Factor 1	Factor 1
SS loadings	304.1895	303.1895
Proportion Var	0.8808	0.8779
Cumulative Var	0.8808	0.8779

(c) Principal component and ML solutions for 1 factor with varimax rotation

	PC		ML	
	Factor 1	Factor 2	Factor 1	Factor 2
SS loadings	289.1885	42.9884	289.8194	40.3575
Proportion Var	0.8374	0.1245	0.8392	0.1169
Cumulative Var	0.8374	0.9618	0.8392	0.9560

(d) Principal component and ML solutions for 2 factors with varimax rotation

use / no use
rotation
=> diagonal

Table 2: Principal component and maximum likelihood solutions with and without varimax with 1 and 2 factors (variance explained)

In both the 1 and 2 factor models, the principal component and maximum likelihood solutions are extremely similar. The only apparent differences come in the size of the uniquenesses (and thus to some degree communalities), having largest differences between the principal component and maximum likelihood solutions for Solar and O3 which have the largest variances

```
# variances
round(diag(cov(air.4)), digits = 2)
```

```
## Wind Solar NO2 O3
## 2.50 300.52 11.36 30.98
```

To decide the adequate number of factors that we should use. p: data size, m: # of factors.

4. Asymptotic Results and the Number of Factors

We consider normal models and formulate hypothesis tests for the number of factors. In PCA, very small eigenvalues are evidence that the corresponding PC score is negligible. In FA, the likelihood of the data drives the hypothesis testing for the adequacy of a m-factor model.

$$H_0 : \Sigma_{p \times p} = L_{p \times m} L_{m \times p}^T + \Psi_{p \times p}$$

against $H_1 : \Sigma$ any other positive definite matrix. When Σ does not have any special form, the maximum of the likelihood function with $\hat{\Sigma} = (n-1)S/n$ is proportional to $|S_n|^{-n/2} e^{-np/2}$, where $S_n = (n-1)S/n$.

Under H_0 , the MLE for $\hat{\mu} = \bar{x}$ and $\hat{\Sigma} = \hat{L}\hat{L}^T + \hat{\Psi}$ where \hat{L} and $\hat{\Psi}$ are the MLE of L and Ψ .

The likelihood ratio statistic for testing H_0 is

$$-2 \ln \Lambda = -2 \ln \left(\frac{\text{maximized likelihood under } H_0}{\text{maximized likelihood}} \right) = -2 \ln \left(\frac{|\hat{\Sigma}|}{|S_n|} \right)^{-n/2} + n [tr(\hat{\Sigma}^{-1} S_n) - p]$$

with degrees of freedom $\nu = \frac{1}{2}p(p-1) - [p(m+1) - \frac{1}{2}m(m-1)] = \frac{1}{2}[(p-m)^2 - p - m]$. $tr(\hat{\Sigma}^{-1} S_n) = p$ provided that $\hat{\Sigma} = \hat{L}\hat{L}^T + \hat{\Psi}$ (See Johnson Wichern Page 527). Thus we have the likelihood ratio statistic

$$-2 \ln \Lambda = n \ln \left(\frac{|\hat{\Sigma}|}{|S_n|} \right), \xrightarrow{n \rightarrow \infty} \chi_{\nu}^2$$

which is approximately χ_{ν}^2 as $n \rightarrow \infty$.

In the stocks data, we test the adequacy of the m-factor model with the tests

$$H_{0,m} : \Sigma = LL + \Psi, H_{1,m} : \Sigma \neq LL + \Psi$$

Starting with $m = 1$,

```
stocks.fac=factanal(x = stocks, factors = 1,
  scores = "regression",
  rotation = "varimax",
  method = "mle")
```

when have cov matrix
need to add one option n.obs =
(when we don't have the raw data)

```
stocks.fac
```

```
##
## Call:
## factanal(x = stocks, factors = 1, scores = "regression", rotation = "varimax", method = "mle")
##
## Uniquenesses:
## Allied Dupont Carbide Exxon Texaco
```



```
## 0.488 0.234 0.560 0.880 0.922
##
## Loadings:
##      Factor1
## Allied 0.715
## Dupont 0.875
## Carbide 0.663
## Exxon 0.346
## Texaco 0.279
##
##      Factor1
## SS loadings 1.915
## Proportion Var 0.383
##
## Test of the hypothesis that 1 factor is sufficient.
## The chi square statistic is 62.22 on 5 degrees of freedom.
## The p-value is 4.22e-12
```

Since $p\text{-value} < 0.0001$, we reject the null hypothesis and conclude that a one factor model is not sufficient to describe the correlations among the weekly returns for the five stocks.

Two factors, $m = 2$

```
stocks.fac=factanal(x = stocks, factors = 2,
                    scores = "regression",
                    rotation = "varimax",
                    method = "mle")
```

stocks.fac

```
##
## Call:
## factanal(x = stocks, factors = 2, scores = "regression", rotation = "varimax", method = "mle")
##
## Uniquenesses:
## Allied Dupont Carbide Exxon Texaco
## 0.417 0.275 0.542 0.005 0.530
##
## Loadings:
##      Factor1 Factor2
## Allied 0.763
## Dupont 0.819 0.232
## Carbide 0.668 0.108
## Exxon 0.113 0.991
## Texaco 0.108 0.677
##
##      Factor1 Factor2
## SS loadings 1.725 1.507
## Proportion Var 0.345 0.301
## Cumulative Var 0.345 0.646
##
## Test of the hypothesis that 2 factors are sufficient.
## The chi square statistic is 1.97 on 1 degree of freedom.
## The p-value is 0.16
```

Since $p\text{-value} = 0.16$, we do not reject the null hypothesis that a two-factor model is sufficient to describe the correlations among the weekly returns for the five stocks.

Can write a function with # numbers of factors and give output with \neq p values to find # of factors

$$P = \begin{bmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{bmatrix}$$

orthogonal / 30° rotation

4. Factor Rotation

In factor analysis, the concept of **factor rotation** is commonly used. Unlike a rotation matrix, which is an orthogonal transformation, a factor rotation is a linear transformation which need not be orthogonal but which satisfies certain criteria, such as goodness-of-fit.

Let \hat{L} is the $p \times m$ matrix of estimated factor loadings obtained by any method (principal component, maximum likelihood, etc), then $\hat{L}^* = \hat{L}\mathbf{R}$ where $(\mathbf{R}\mathbf{R}^T = \mathbf{R}^T\mathbf{R} = \mathbf{I})$ is a $p \times m$ matrix of "rotated" loadings. That is, multiplying a matrix of factor loadings by any orthogonal matrix leads to the same approximation to the covariance (or correlation) matrix.

- The estimated residual matrix also remains unchanged

$$\hat{L}^* \hat{L}^{*T} = \hat{L} (\mathbf{R}\mathbf{R}^T) \hat{L}^T = \hat{L}\hat{L}^T = \mathbf{S}_n - \hat{L}\hat{L}^T - \hat{\Psi} = \mathbf{S}_n - \hat{L}\mathbf{R}\mathbf{R}^T\hat{L}^T - \hat{\Psi} = \mathbf{S}_n - \hat{L}^* \hat{L}^{*T} - \hat{\Psi}$$

review about rotation matrix (MHAO/SC)

- The specific variances $\hat{\Psi}_i$ and therefore the communalities also remain unchanged

- Since only the loadings change by rotation, we rotate factors to see if we can better interpret results. There are many choices for a rotation matrix \mathbf{R} , and to choose, we first establish a mathematical criterion and then see which \mathbf{R} can best satisfy the criterion. One possible objective is to have each one of the p variables load highly on only one factor and have moderate to negligible loads on all other factors.
- Many popular orthogonal factor rotation methods try to maximize

$$V(\mathbf{L}, \mathbf{R} | \gamma) = \frac{1}{p} \sum_{k=1}^m \left[\sum_{j=1}^p (\hat{l}_{jk}^* / \hat{h}_j)^4 - \frac{\gamma}{p} \left(\sum_{j=1}^p (\hat{l}_{jk}^* / \hat{h}_j)^2 \right) \right]$$

where \hat{l}_{jk}^* is the rotated loading of the j^{th} variable on the k^{th} factor. \hat{h}_j is the square-root of the communality for X_j

Changing the γ parameter corresponds to different criteria

- $\gamma = 1$ corresponds to varimax criterion overall factor
- $\gamma = 0$ corresponds to quartimax criterion preserve overall factor almost equal weights
- $\gamma = m/2$ corresponds to equamax criterion
- $\gamma = p(m-1)/(p+m-2)$ corresponds to parsimax criterion

The procedure finds the orthogonal transformation of the loading matrix that maximizes the sum of those variances, summing across all m rotated factors.

- The varimax rotation will destroy an "overall" factor. After rotation each of the p variables should load highly on at most one of the rotated factors.
- The quartimax rotation try to preserve an overall factor such that each of the p variables has a high loading on that factor and create other factors such that each of the p variables has a high loading on at most one factor.

The above four rotations produce uncorrelated factors.

Example: Visualization of 2D Clockwise Rotation

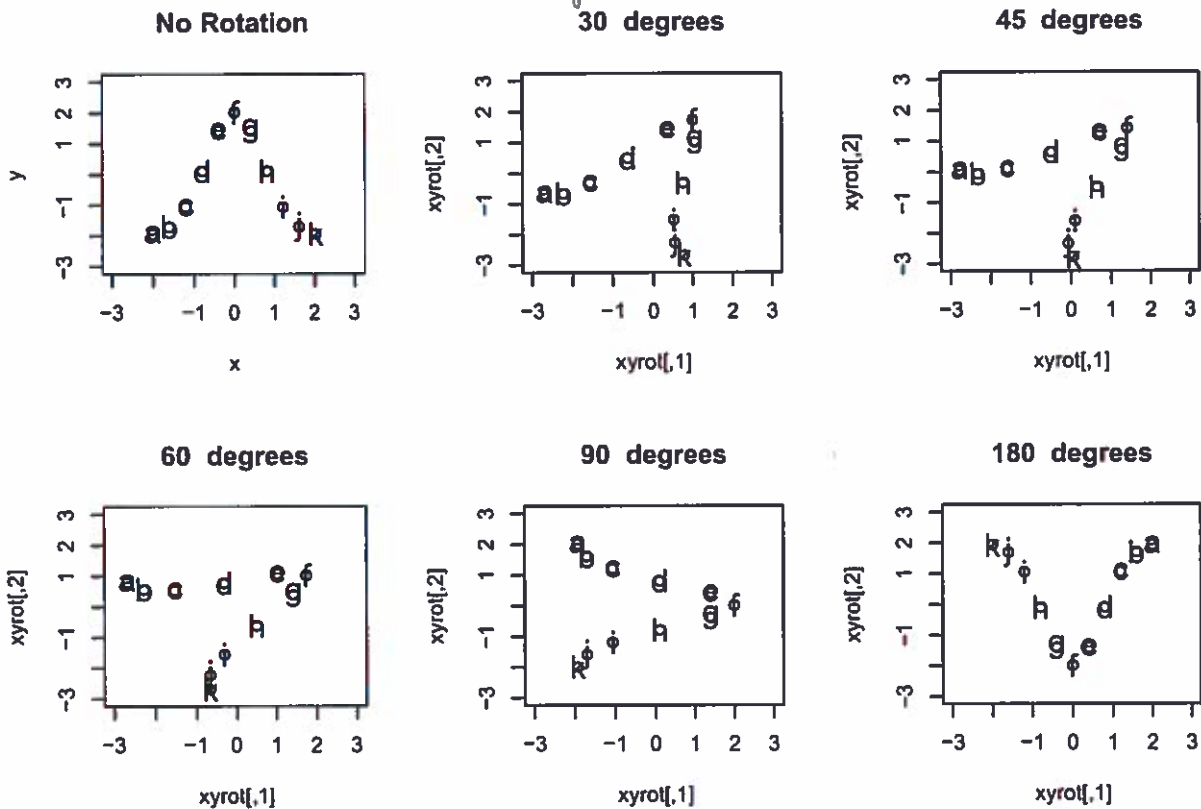
```
rotmat2d <- function(theta){
  matrix(c(cos(theta), sin(theta), -sin(theta), cos(theta)), 2, 2)
}
x <- seq(-2, 2, length=11)
```

```

y <- 4*exp(-x^2) - 2
xy <- cbind(x,y)
rang <- c(30,45,60,90,180)
par(mfrow=c(2,3))
plot(x,y,xlim=c(-3,3),ylim=c(-3,3),main="No Rotation")
text(x,y,labels=letters[1:11],cex=1.5)
for(j in 1:5){
  rmat <- rotmat2d(rang[j]*2*pi/360)
  xyrot <- xy%*%rmat
  plot(xyrot,xlim=c(-3,3),ylim=c(-3,3))
  text(xyrot,labels=letters[1:11],cex=1.5)
  title(paste(rang[j]," degrees"))
}

```

orthogonal rotation



Oblique transformation

Orthogonal rotations produce uncorrelated factors. However, PROMAX is a non-orthogonal (oblique) transformation that is not a rotation and can produce correlated factors

In PROMAX transformation,

1. we first perform a varimax rotate to obtain loadings L^* .
2. Construct another $p \times m$ matrix Q such that

$$q_{ij} = \begin{cases} |l_{ij}^*|^{m-1} l_{ij}^* & \text{for } l_{ij}^* \neq 0 \\ 0 & \text{for } l_{ij}^* = 0 \end{cases}$$

where $m > 1$ is selected by trial and error, usually $m < 4$.

3. Find a matrix U such that each column of L^*U is close to the corresponding column of Q . Choose the j -th column of U to minimize

$$(q_j - L^*u_j)(q_j - L^*u_j)^T$$

which yields

$$U = (L^{*T}L^*)^{-1}L^{*T}Q$$

4. Rescale U so that the transformed factors have unit variance. Compute $D^2 = \text{diag}[(U^TU)^{-1}]$ and $M = UD$.
5. The PROMAX factors are obtained from $L^*F = L^*MM^{-1}F = L_{promax}F_{promax}$. The PROMAX transformation yields factors with loadings

$$L_{promax} = L^*M$$

The correlation matrix for the transformed factors is $(M^TM)^{-1}$.

$$\text{Var}(F_{promax}) = \text{Var}(M^{-1}F) = M^{-1}\text{Var}(F)(M^{-1})^T = M^{-1}I(M^{-1})^T = (M^TM)^{-1}$$

In the stocks example, now we change the rotation to promax from varimax

```
promax(stocks.fac$loadings, m = 2)
```

```
## $loadings
##
## Loadings:
##      Factor1 Factor2
## Allied    0.778
## Dupont    0.814    0.123
## Carbide   0.672
## Exxon          0.995
## Texaco       0.676
##
##      Factor1 Factor2
## SS loadings    1.722    1.468
## Proportion Var  0.344    0.294
## Cumulative Var  0.344    0.638
##
## $rotmat
##      [,1] [,2]
## [1,] 1.0234906 -0.1377833
## [2,] -0.1066074  1.0197617
```

Of course, we could also have done it directly, if we had started using "promax"

```
stocks.facp <- factanal(x = stockss, factors = 2,
                        scores = "regression",
                        rotation = "promax",
                        method = "mle")

stocks.facp
mod1=factanal(x = stockss, factors = 2, scores = "regression",
              rotation = "varimax", method = "mle")
mod1$loadings

##
## Loadings:
```

```
##          Factor1 Factor2
## Allied  0.763
## Dupont  0.819  0.232
## Carbide 0.668  0.108
## Exxon   0.113  0.991
## Texaco  0.108  0.677
##
##          Factor1 Factor2
## SS loadings  1.725  1.507
## Proportion Var  0.345  0.301
## Cumulative Var  0.345  0.646
```

```
tcrossprod(solve(mod1$rotmat)) # correlation between rotated factor scores
```

```
##          [,1]      [,2]
## [1,] 1.000000e+00 4.163336e-17
## [2,] 4.163336e-17 1.000000e+00
```

when we use varimax
 $L L^T \leftarrow F_1 \text{ and } F_2 \text{ are uncorrelated}$

```
mod2=factanal(x = stockss, factors = 2, scores = "regression",
              rotation = "promax", method = "mle")
```

```
mod2$loadings
```

```
##
## Loadings:
##          Factor1 Factor2
## Allied   0.786  -0.104
## Dupont   0.821
## Carbide  0.679
## Exxon    0.997
## Texaco   0.676
##
##          Factor1 Factor2
## SS loadings  1.753  1.471
## Proportion Var  0.351  0.294
## Cumulative Var  0.351  0.645
```

when we use promax
 then F_1 and F_2 are correlated.

(can use (*) in page 19)

```
tcrossprod(solve(mod2$rotmat)) # correlation between rotated factor scores
```

```
##          [,1]      [,2]
## [1,] 1.0000000 0.2787533
## [2,] 0.2787533 1.0000000
```

F_1 and F_2 are correlated
 off diagonal elements are not zero.

5. Factor Scores

In Factor Analysis, interest is usually centered on the parameters in the factor model. However, the estimated values of the common factors, called **factor scores**, may also be required. Factor scores are estimates of values for the unobserved random factor vectors $F_j, j = 1, \dots, n$.

Estimating Factor Scores via the Weighted Least Squares Method

Suppose that the mean vector μ , the factor loading L and the specific variance Ψ are known for the factor model

$$X = LF + \mu + \epsilon$$

where $\epsilon^T = [\epsilon_1, \dots, \epsilon_p]$ are the errors, $Var(\epsilon_i) = \psi_i, i = 1, \dots, p$ need not be equal. The sum of the squares of the errors, weighted by the reciprocal of their variances, is

sum of squares $\rightarrow \sum_{i=1}^p \left(\frac{\epsilon_i^2}{\psi_i} \right) = \epsilon^T \Psi^{-1} \epsilon = (x - \mu - Lf)^T \Psi^{-1} (x - \mu - Lf)$

The estimate \hat{f} that minimize the weighted sum of the squares of the errors is

weighted least square $\hat{f}_i^{(WLS)} = (\hat{L}^T \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}^T \hat{\Psi}^{-1} (x_i - \bar{x})$

where x_i is the i^{th} subject's vector of data, $\bar{x} = \sum_{i=1}^n x_i / n$ is the sample mean.

Note that if the factor loadings are estimated by the principal component method, then it is typical to use

$$\hat{f}_i = (\hat{L}^T \hat{L})^{-1} \hat{L}^T (x_i - \bar{x})$$

which is the unweighted least squares estimate,

Estimating Factor Scores via the Regression Method

center data

Using the maximum likelihood method, the joint distribution of $(X - \mu, F)$ is multivariate normal with mean vector 0_{p+m} and covariance matrix

$$(X - \mu, F) \propto N_{p+m} \left(0, \Sigma_{(m+p) \times (m+p)} = \begin{pmatrix} \text{var}(X - \mu) & \\ & \text{var}(F) \end{pmatrix} = \begin{pmatrix} (LL^T + \Psi)_{p \times p} & L_{p \times m} \\ L_{m \times p}^T & I_{m \times m} \end{pmatrix} \right)$$

orthogonal factor model

which implies that the conditional distribution of F given X has

- $E(F | X) = L^T (LL^T + \Psi)^{-1} (X - \mu)$
 - $Var(F | X) = I_{m \times m} - L^T (LL^T + \Psi)^{-1} L$
- $X = LF + \mu + \epsilon$

Proof $E(F|X) = L^T \Sigma^{-1} (X - \mu) = L^T (LL^T + \Psi)^{-1} (X - \mu)$
 $Var(F|X) = I - L^T \Sigma^{-1} L = I - L^T [(LL^T) + \Psi]^{-1} L$

Assume that Σ, Ψ and $(I_{m \times m} + L^T \Psi^{-1} L)$ are invertible matrices. The underlying Factor model is $\Sigma = L^T L + \Psi$, then the regression estimate of the factor scores have the form

$$\hat{f}_i^{(REG)} = \hat{L}^T (\hat{L} \hat{L}^T + \hat{\Psi})^{-1} (x_i - \mu) = (I_{m \times m} + \hat{L}^T \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}^T \hat{\Psi}^{-1} (x_i - \mu)$$

We want to show

$$L^T (LL^T + \Psi)^{-1} = (I_{m \times m} + L^T \Psi^{-1} L)^{-1} L^T \Psi^{-1} \quad (3)$$

By assumption, $\Sigma^{-1} = (LL^T + \Psi)^{-1}$
 and $(I_{m \times m} + L^T \Psi^{-1} L)^{-1}$ exist

Because $(I_{m \times m} + B)^{-1} (I_{m \times m} + B) (I_{m \times m} + B)^{-1} = (I_{m \times m} + B)^{-1} (I_{m \times m} + B) [I_{m \times m} - (I_{m \times m} + B)^{-1} B]$
 it follows that
 $(L^T + \Psi)^{-1} (LL^T + \Psi) = [\Psi^{-1} - \Psi^{-1} L (I_{m \times m} + B)^{-1} L^T \Psi^{-1}] (L^T + \Psi)$

Let $B = L^T \Psi^{-1} L$
 $(I_{m \times m} + B) (I_{m \times m} + B)^{-1} B = (I_{m \times m} + B) [I_{m \times m} - (I_{m \times m} + B)^{-1} B] = B [(I_{m \times m} + B) - I_{m \times m}] = 0$
 which implies that $(I_{m \times m} + B)^{-1} B = I_{m \times m} - (I_{m \times m} + B)^{-1} \quad (2)$

cause $(\mathbf{L}\mathbf{L}^T + \Psi)$ is invertible we conclude that $(\mathbf{L}\mathbf{L}^T + \Psi)^{-1} = \Psi^{-1} - \Psi^{-1}\mathbf{L}(\mathbf{I}_{m \times m} + \mathbf{L}^T\Psi^{-1}\mathbf{L})^{-1}\mathbf{L}^T\Psi^{-1}$
 $\Rightarrow (\mathbf{L}\mathbf{L}^T + \Psi)^{-1}\mathbf{L} = \Psi^{-1}\mathbf{L}(\mathbf{I}_{m \times m} + \mathbf{L}^T\Psi^{-1}\mathbf{L})^{-1}$
 the desired expression for $\hat{\beta}_i^{REG}$ follows from (6) \Rightarrow the equality.

Note that there is a simple relationship between the weighted least squares estimate and the regression estimate

$$\hat{\beta}_i^{(WLS)} = (\hat{\mathbf{L}}^T \hat{\Psi}^{-1} \hat{\mathbf{L}})^{-1} (\mathbf{I}_{m \times m} + \hat{\mathbf{L}}^T \hat{\Psi}^{-1} \hat{\mathbf{L}}) \hat{\beta}_i^{(REG)} = [\mathbf{I}_{m \times m} + (\hat{\mathbf{L}}^T \hat{\Psi}^{-1} \hat{\mathbf{L}})^{-1}] \hat{\beta}_i^{(REG)}$$

The relationship between the weighted least squares estimate and the regression estimate implies that

$$\|\hat{\beta}_i^{(WLS)}\|^2 > \|\hat{\beta}_i^{(REG)}\|^2$$

where $\|\cdot\|$ denotes the Euclidean norm.

Example: Decathlon data

```
library(ade4)
```

```
## Warning: package 'ade4' was built under R version 3.4.4
```

```
data(olympic)
```

```
decathlon <- cbind(olympic$tab, olympic$score)
```

```
colnames(decathlon) <- c("run100", "long.jump", "shot", "high.jump", "run400", "hurdle", "discus", "pole", "javelin")
```

```
head(decathlon)
```

```
## run100 long.jump shot high.jump run400 hurdle discus pole.vaule javelin
## 1 11.25 7.43 15.48 2.27 48.90 15.13 49.28 4.7 61.32
## 2 10.87 7.45 14.97 1.97 47.71 14.46 44.36 5.1 61.76
## 3 11.18 7.44 14.20 1.97 48.29 14.81 43.66 5.2 64.16
## 4 10.62 7.38 15.02 2.03 49.06 14.72 44.80 4.9 64.04
## 5 11.02 7.43 12.92 1.97 47.44 14.40 41.20 5.2 57.46
## 6 10.83 7.72 13.58 2.12 48.34 14.18 43.06 4.9 52.18
## run1500 score
## 1 268.95 8488
## 2 273.02 8399
## 3 263.20 8328
## 4 285.11 8306
## 5 256.64 8286
## 6 274.07 8272
```

```
# resign running events (so higher score means better performance)
```

```
decathlon[,c(1,5,6,10)] <- (-1)*decathlon[,c(1,5,6,10)]
```

```
head(decathlon)
```

```
## run100 long.jump shot high.jump run400 hurdle discus pole.vaule javelin
## 1 -11.25 7.43 15.48 2.27 -48.90 -15.13 49.28 4.7 61.32
## 2 -10.87 7.45 14.97 1.97 -47.71 -14.46 44.36 5.1 61.76
## 3 -11.18 7.44 14.20 1.97 -48.29 -14.81 43.66 5.2 64.16
## 4 -10.62 7.38 15.02 2.03 -49.06 -14.72 44.80 4.9 64.04
## 5 -11.02 7.43 12.92 1.97 -47.44 -14.40 41.20 5.2 57.46
## 6 -10.83 7.72 13.58 2.12 -48.34 -14.18 43.06 4.9 52.18
## run1500 score
## 1 -268.95 8488
## 2 -273.02 8399
## 3 -263.20 8328
## 4 -285.11 8306
## 5 -256.64 8286
## 6 -274.07 8272
```

Another way to solve
is taking the maximum - the minimum.

```
# check variable scales
apply(decathlon,2,mean)
```

```
##      run100  long.jump      shot  high.jump  run400  hurdle
## -11.196364  7.133333  13.976364  1.982727 -49.276667 -15.048788
##      discus  pole.vaule  javelin  run1500  score
##  42.353939  4.739394  59.438788 -276.038485 7856.909091
```

```
apply(decathlon,2,sd)
```

```
##      run100  long.jump      shot  high.jump  run400  hurdle
##  0.2433210  0.3043401  1.3319906  0.0939838  1.0696602  0.5067652
##      discus  pole.vaule  javelin  run1500  score
##  3.7191312  0.3344206  5.4959984  13.6570975 415.0694493
```

*Min long distance
has bigger s.d.*

```
# scree plot for factor analysis results
```

```
famods <- vector("list", 6)
```

```
for(k in 1:6){
```

```
  famods[[k]] <- factanal(x=decathlon[,1:10], factors=k)
```

```
}
```

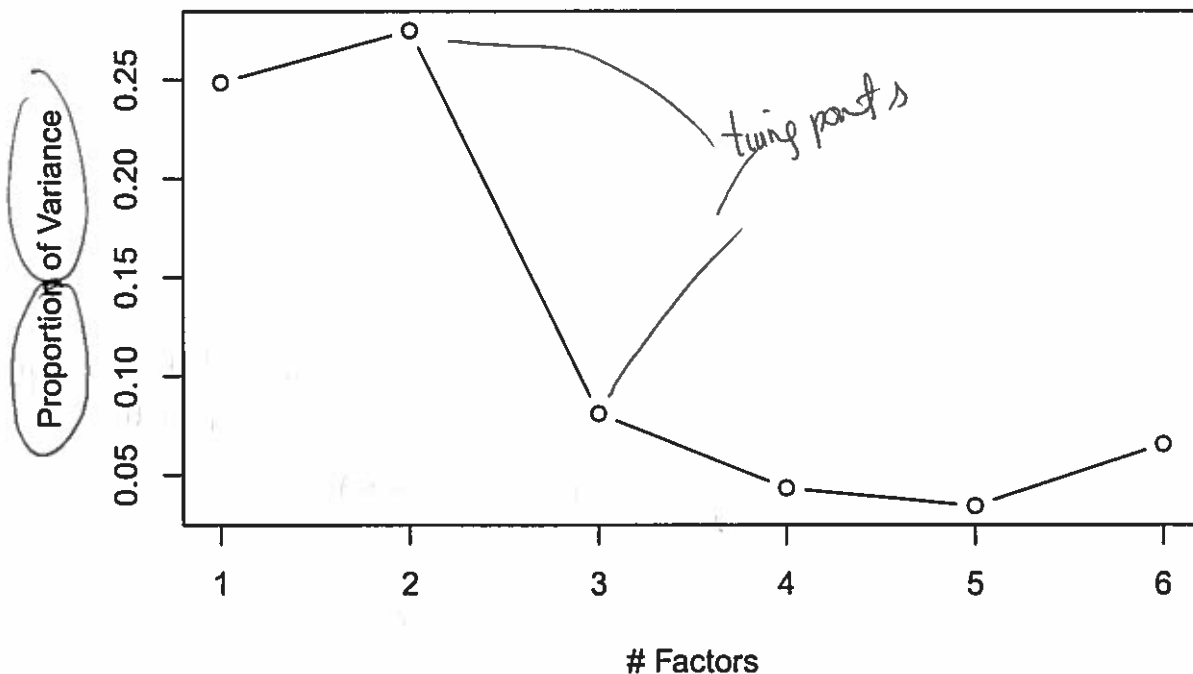
```
vafs <- sapply(famods, function(x) sum(x$loadings^2)) / nrow(famods[[1]]$loadings)
```

```
vaf.scree <- vafs - c(0, vafs[1:5])
```

```
#dev.new(width=10, height=5, noRStudioGD=TRUE)
```

```
plot(1:6, vaf.scree, type="b", xlab="# Factors", ylab="Proportion of Variance",
     main="FA Scree Plot")
```

FA Scree Plot



```
# FA on correlation matrix (w/ varimax rotation)
```

```
famod <- factanal(x=decathlon[,1:10], factors=2)
```

```
names(famod)
```

```
## [1] "converged" "loadings" "uniquenesses" "correlation"
```

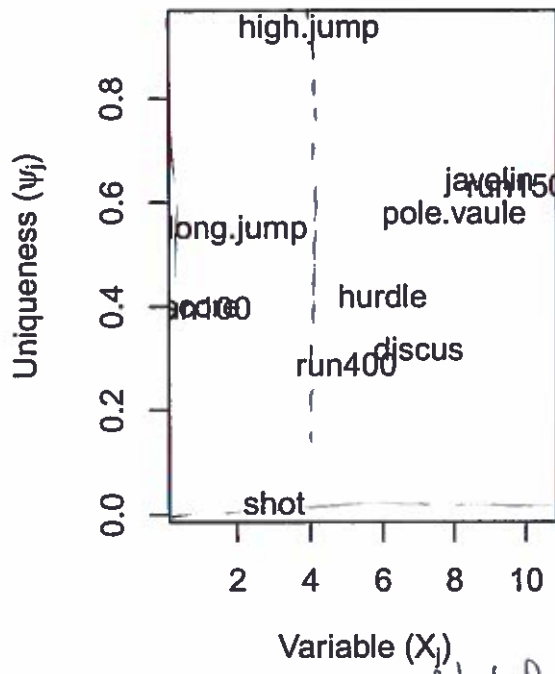
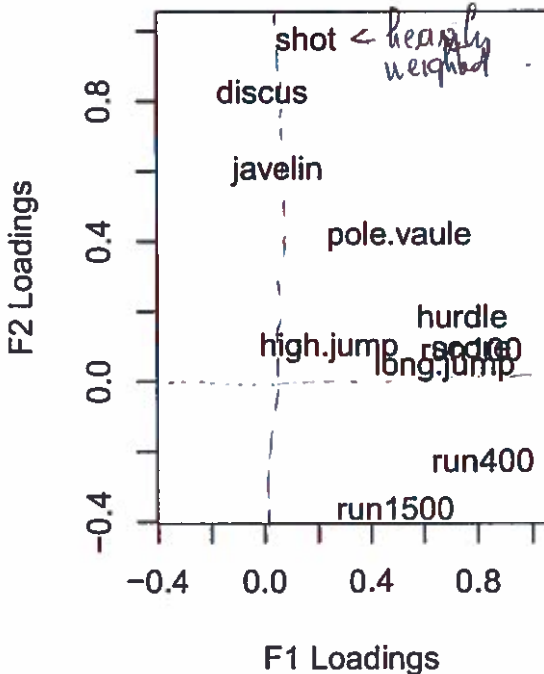


```
## [5] "criteria"      "factors"      "dof"          "method"
## [9] "rotmat"        "STATISTIC"   "PVAL"        "n.obs"
## [13] "call"

# plot factor loadings and uniquenesses
#dev.new(width=10, height=5, noRStudioGD=TRUE)
par(mfrow=c(1,2))
plot(famod$loadings, xlab="F1 Loadings", ylab="F2 Loadings",
     type="n", main="Factor Loadings", xlim=c(-0.35, 1), ylim=c(-0.35, 1))
text(famod$loadings, labels=colnames(decathlon))
plot(famod$uniquenesses, xlab=expression("Variable ("*X[j]*")"),
     ylab=expression("Uniqueness ("*psi[j]*")"),
     type="n", main="Factor Uniquenesses", xlim=c(0.5,10.5))
text(famod$uniquenesses, labels=colnames(decathlon))
```

Factor Loadings

Factor Uniquenesses



```
# refit model and get FA scores (NOT GOOD IDEA!!)
famodW <- factanal(x=decathlon[,1:10], factors=2, scores="Bartlett")
famodR <- factanal(x=decathlon[,1:10], factors=2, scores="regression")
round(cor(famodW$scores, famodR$scores), 4)
```

weighted least square approach

$cor(\hat{\beta}_{WLS}, \hat{\beta}_{REG})$

```
##          Factor1 Factor2
## Factor1  0.9998  0.0000
## Factor2  0.0000  0.9998
```

```
# check correlation of FA scores w/ overall decathlon score
round(cor(decathlon$score, famodR$scores), 4)
```

```
##          Factor1 Factor2
## [1,] 0.7952  0.5022
```

overall score (X_1, \dots, X_{10}, Y)
 correlation of Y and F_1

```
round(cor(decathlon$score, famodW$scores), 4)
```

① $Y \approx X_1 + \dots + X_{10}$

② $X = LF + \epsilon$

→ 2 factor model

$\begin{matrix} | F_1 \\ | F_2 \end{matrix}$

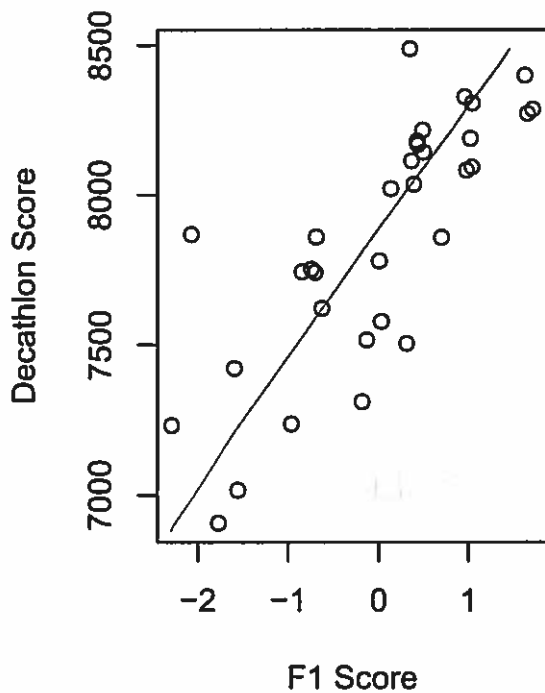
$Y \approx F_1 + F_2$

reduce from 10 to 2

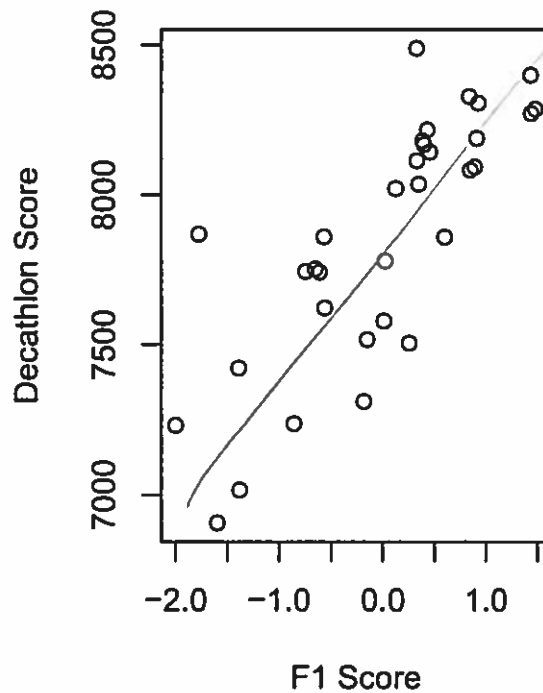
$$R^2 = \frac{SSR}{SST}$$

```
##      Factor1 Factor2
## [1,] 0.7854 0.4865
#dev.new(width=10, height=5, noRStudioGD=TRUE)
par(mfrow=c(1,2))
plot(famodW$scores[,1], decathlon$score, xlab="F1 Score",
      ylab="Decathlon Score", main="Weighted Least Squares Method")
plot(famodR$scores[,1], decathlon$score, xlab="F1 Score",
      ylab="Decathlon Score", main="Regression Method")
```

Weighted Least Squares Method



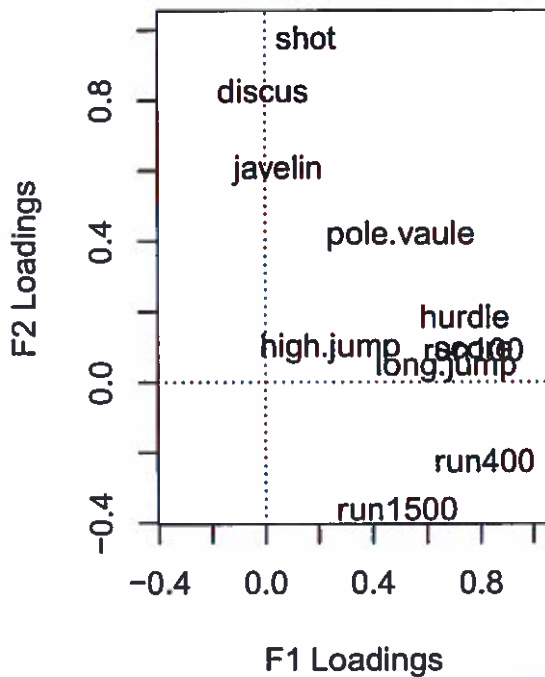
Regression Method



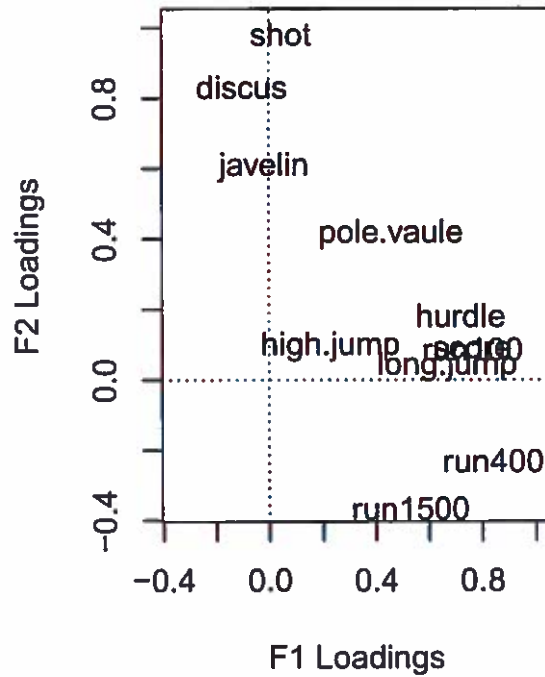
```
# try promax rotation
famod.promax <- promax(famod$loadings)
#dev.new(width=10, height=5, noRStudioGD=TRUE)
par(mfrow=c(1,2))
plot(famod$loadings, xlab="F1 Loadings", ylab="F2 Loadings",
      type="n", main="Varimax Factor Loadings", xlim=c(-0.35, 1), ylim=c(-0.35, 1))
text(famod$loadings, labels=colnames(decathlon))
abline(0,0,lty=3)
abline(v=0,lty=3)
plot(famod.promax$loadings, xlab="F1 Loadings", ylab="F2 Loadings",
      type="n", main="Promax Factor Loadings", xlim=c(-0.35, 1), ylim=c(-0.35, 1))
text(famod.promax$loadings, labels=colnames(decathlon))
abline(0,0,lty=3)
abline(v=0,lty=3)
```

prefer uncorrelated version

Varimax Factor Loadings



Promax Factor Loadings



```
# compare loadings
oldFALoadings <- famod$loadings
newFALoadings <- famod$loadings %*% famod.promax$rotmat
sum((newFALoadings - famod.promax$loadings)^2)

## [1] 1.55904e-31

# compare reproduced data before and after rotation
oldFAScores <- famodR$scores
newFAScores <- oldFAScores %*% t(solve(famod.promax$rotmat))
Xold <- tcrossprod(oldFAScores, oldFALoadings)
Xnew <- tcrossprod(newFAScores, newFALoadings)
sum((Xold - Xnew)^2)

## [1] 4.431144e-30

# population and sample factor score covariance matrix (after rotation)
tcrossprod(solve(famod.promax$rotmat)) # population

##          [,1]      [,2]
## [1,] 1.0000000 0.1182015
## [2,] 0.1182015 1.0000000
cor(newFAScores) # sample

##          [,1]      [,2]
## [1,] 1.0000000 0.1320303
## [2,] 0.1320303 1.0000000
```

← promax

$L^* = LR.$

correlated

6. Comparison of PCA and FA (double check)

There are many parallels between Principal Component Analysis and Factor Analysis, but there are also important differences. The common goal is to represent the data in a smaller number of simpler components, but the two methods are based on different strategies and models to achieve these goals. We list some of the differences between the two approaches.

Table 1: Comparison of PCA and FA

	PCA	FA
Aim	Fewer and simpler components	Fewer factor loadings
Model	Not explicit	Strick k-factor model $\mathbf{X} = \boldsymbol{\mu} + \mathbf{LF} + \boldsymbol{\epsilon}$
Covariance matrix	$\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Lambda}\mathbf{P}^T$	factor $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$
Solution method	Singular value decomposition $\mathbf{X} = \mathbf{U}\boldsymbol{\Gamma}\mathbf{V}^T$; Spectral decomposition $\boldsymbol{\Sigma} = \mathbf{P}\boldsymbol{\Lambda}\mathbf{P}^T$?	PC-based loadings or ML-based loadings with orthogonal or oblique rotations
Scores	Projection onto eigenvectors ranked by variance; uniquely defined uncorrelated component scores	PC-based scores or ML regression scores; not unique, include rotations; rotated scores correlated
Data description	Approximate \mathbf{X} $Y_1 = \mathbf{a}_1^T \mathbf{X}, Y_2 = \mathbf{a}_2^T \mathbf{X}, \dots, Y_q = \mathbf{a}_q^T \mathbf{X}$	Complete, includes specific factor $\mathbf{X} - \boldsymbol{\mu} = \mathbf{LF} + \boldsymbol{\epsilon}$

Two type of factor analysis

Exploratory

← It is exploratory when you do not have a pre-decided idea of the structure or how many dimensions are in the set of variable

Confirmatory

← helps confirm when we want to test about the structure of the # of dimension underlying a set of variables.

Chapter 6. Canonical Correlation

Jianxuan Liu

Fall 2018

Canonical Correlation Analysis (CCA) seeks to identify and quantify the associations between two sets of variables. CCA connects two sets of variables by finding linear combinations of variables that maximally correlate. The idea is first to determine the pair of linear combinations having the largest correlation. Next, we determine the pair of linear combinations having the largest correlation among all pairs uncorrelated with the initially selected pair, and so on. The pairs of linear combinations are called **canonical variables** and their correlations are called **canonical correlations**

There are two typical purposes of CCA:

ANOVA

- Data reduction: explain covariation between two sets of variables using small number of linear combinations
- Data interpretation: find features (i.e., canonical variates) that are important for explaining covariation between sets of variables

1. Canonical variables and canonical correlations

The first set, of p variables, is represented by the $p \times 1$ random vectors $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$. The second group, of q variables, is represented by the $q \times 1$ random vector $\mathbf{Y} = (Y_1, Y_2, \dots, Y_q)^T$. For the random vectors \mathbf{X} and \mathbf{Y} , let

$$\begin{aligned} E(\mathbf{X}) &= \boldsymbol{\mu}_X, E(\mathbf{Y}) = \boldsymbol{\mu}_Y \\ \text{Var}(\mathbf{X}) &= \boldsymbol{\Sigma}_X, \text{Var}(\mathbf{Y}) = \boldsymbol{\Sigma}_Y \\ \text{Cov}(\mathbf{X}, \mathbf{Y}) &= E[(\mathbf{X} - \boldsymbol{\mu}_X)(\mathbf{Y} - \boldsymbol{\mu}_Y)^T] = \boldsymbol{\Sigma}_{XY} \end{aligned}$$

- Define new variables U and V via linear combinations of \mathbf{X} and \mathbf{Y}

canonical variable $\rightarrow U = \mathbf{a}^T \mathbf{X} = a_1 X_1 + a_2 X_2 \dots \quad V = \mathbf{b}^T \mathbf{Y} =$

for some pair of coefficient vectors \mathbf{a} and \mathbf{b} . Then U and V have the following properties

$$\text{Var}(U) = \mathbf{a}^T \boldsymbol{\Sigma}_X \mathbf{a}, \text{Var}(V) = \mathbf{b}^T \boldsymbol{\Sigma}_Y \mathbf{b}, \text{Cov}(U, V) = \mathbf{a}^T \boldsymbol{\Sigma}_{XY} \mathbf{b}$$

CCA seek coefficient vectors \mathbf{a} and \mathbf{b} such that

$$\text{Cor}(U, V) = \frac{\text{Cov}(U, V)}{\sqrt{\text{Var}(U)\text{Var}(V)}}$$

is as large as possible, subject to U and V having unit variance,

- The first pair of canonical variables is the pair of linear combinations of U_1 and V_1 having unit variance, which maximize the correlation $\text{Cor}(U, V)$.
- The second pair of canonical variables is the pair of linear combinations of U_2 and V_2 having unit variance, which maximize the correlation $\text{Cor}(U, V)$ among all choices that are uncorrelated with the first pair of canonical variables (U_1, V_1)
- The k^{th} pair of canonical variables is the pair of linear combinations of U_k and V_k having unit variance, which maximize the correlation $\text{Cor}(U, V)$ among all choices that are uncorrelated with the previous $k - 1$ pairs of canonical variable pairs.

The correlation between the k^{th} pair of canonical variables is called the k^{th} canonical correlation.

The k^{th} pair of canonical variates is given by

$$U_k = \underbrace{u_k^T \Sigma_X^{-1/2}}_{a_k^T} X, \quad V_k = \underbrace{v_k^T \Sigma_Y^{-1/2}}_{b_k^T} Y$$

where $a_k^T = u_k^T \Sigma_X^{-1/2}$, $b_k^T = v_k^T \Sigma_Y^{-1/2}$.

- u_k is the k^{th} eigenvector of $\Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1/2}$
- v_k is the k^{th} eigenvector of $\Sigma_Y^{-1/2} \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} \Sigma_Y^{-1/2}$

The k^{th} canonical correlation is given by

square here

$$\text{Cor}(U_k, V_k) = \rho_k$$

where ρ_k^2 is the k^{th} eigenvalue of $\Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1/2}$ or $\Sigma_Y^{-1/2} \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} \Sigma_Y^{-1/2}$.

assume $X \sim \mathbb{R}^p$, $Y \sim \mathbb{R}^q$
 $\circ (N_{X, \Sigma_X}) \quad Y \sim (N_{Y, \Sigma_Y})$

$$\text{cov}(X, Y) = \Sigma_{XY}$$

$$\text{let } u = a^T X, \quad v = b^T Y$$

$$r(u, v) = \rho = \frac{a^T \Sigma_{XY} b}{\sqrt{a^T \Sigma_X a} \sqrt{b^T \Sigma_Y b}}$$

$$\text{define } \underline{c} = \Sigma_X^{-1/2} a \Rightarrow a = \Sigma_X^{1/2} \underline{c} \quad a^T = \underline{c}^T \Sigma_X^{-1/2}$$

$$\underline{d} = \Sigma_Y^{-1/2} b$$

$$= \frac{\underline{c}^T \Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1/2} \underline{d}}{\|\underline{c}\| \|\underline{d}\|} \quad \|\cdot\|: \text{Euclidean norm}$$

$$|\underline{c}^T \Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1/2} \underline{d}| \leq \|\underline{c}^T \Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1/2}\| \|\underline{d}\|$$

$$\rho = \frac{\underline{c}^T \Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1/2} \underline{d}}{\|\underline{c}\| \|\underline{d}\|}$$

ie maximum occur when \underline{c} is the eigen vector of corresponding to the largest eigenvalue of

$$\Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1/2}$$

ie maximized correlation is the largest eigenvalue of $\Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1/2}$

(ρ, a) are ~~eigenvalue~~ eigenvalue, eigenvector of $\Sigma_X^{-1} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1}$

* Use this

Alternatively, the canonical variables coefficient vectors a and b and the corresponding correlations can be found by solving the eigenvalue equations

$$\Sigma_X^{-1} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} a = \rho^2 a$$

$$e_v \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} b = \rho^2 b$$

This is often more practical to use when computing the coefficients and the correlations.

The canonical variables have the following properties:

- $Var(U_k) = Var(V_k) = 1$
- $Cov(U_k, U_l) = Cor(U_k, U_l) = 0, k \neq l$
- $Cov(V_k, V_l) = Cor(V_k, V_l) = 0, k \neq l$
- $Cov(U_k, V_l) = Cor(U_k, V_l) = 0, k \neq l$

Let $U = A^T X$. Show that $cov(U) = A^T \Sigma_X A = I$
 $A = e^T \Sigma^{-1/2} = e^T P_r \Lambda_r^{-1/2} P_r^T$ where e is an orthogonal mat
 $\Sigma_X = P_r \Lambda_r P_r^T$
 $P_r^T X$ is the set of all p_r derived from X alone
 matrix $\Lambda_r^{-1/2} P_r^T X$ has their row $\frac{1}{\sqrt{\lambda_i}} P_r^T X$ which is the i^{th} pc scaled to have unit variance

Example: Open/closed book exams

```
library(bootstrap)
data(scor) # first two are closed book, last three are open book
S=cov(scor)
S11=S[1:2,1:2]; S22=S[3:5,3:5]; S12=S[1:2,3:5]; S21=S[3:5,1:2]
e1=eigen(solve(S11)%*%S12)%*%solve(S22)%*%S21
e2=eigen(solve(S22)%*%S21)%*%solve(S11)%*%S12
# canonical vectors for 2 closed book exams
a=-e1$vector[,1]
# canonical vectors for 3 open book exams
b=e2$vector[,1]
# (a,b) is the first pair
```

$$cov(\Lambda_r^{-1/2} P_r^T X) = \Lambda_r^{-1/2} P_r^T \Sigma_X P_r \Lambda_r^{-1/2}$$

$$= \Lambda_r^{-1/2} P_r^T (P_r \Lambda_r P_r^T) P_r \Lambda_r^{-1/2}$$

$$= \Lambda_r^{-1/2} \Lambda_r \Lambda_r^{-1/2} = I$$

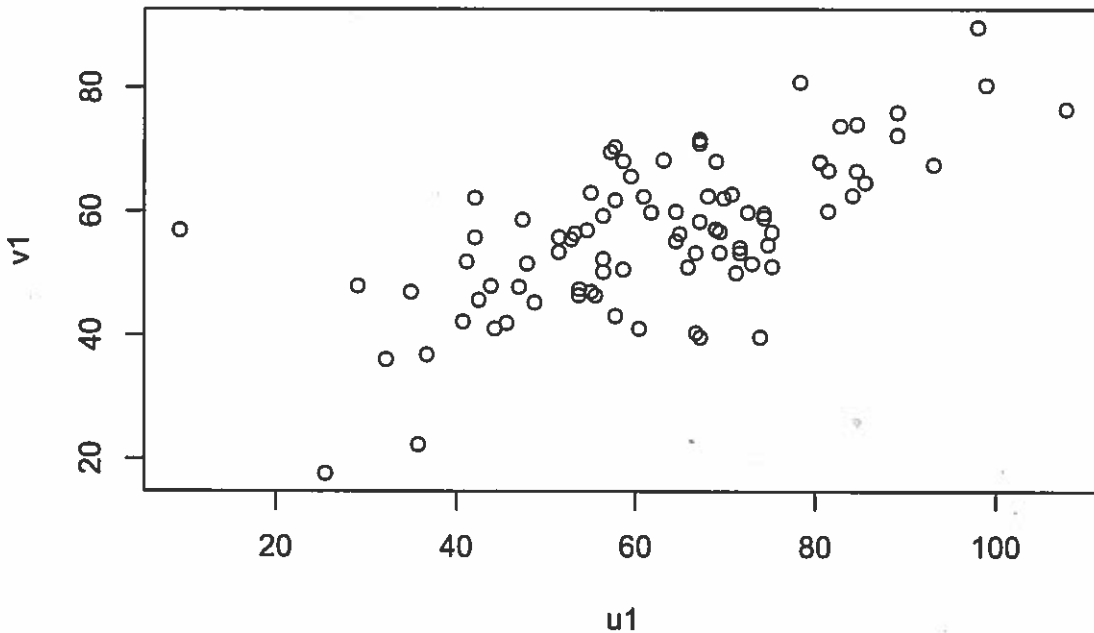
* Show $cov(V) = E(VV^T) = E(Y^T \Sigma_Y^{-1} (Y - \mu_Y)(Y - \mu_Y)^T \Sigma_Y^{-1} Y)$

$$= Y^T \Sigma_Y^{-1} \Sigma_Y \Sigma_Y^{-1} Y$$

$$= I \text{ since } cov(Y) = E(Y Y^T) = \Lambda_r \Lambda_r^T = \Sigma_Y$$

X [closed b score]
 [closed b score]

$\begin{bmatrix} Y \\ Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \rightarrow \begin{cases} \text{op b score} \\ \text{op b score} \\ \text{op b score} \end{cases}$




```

cor(t(rbind(u1,v1)))

##          [,1]      [,2]
## [1,] 1.0000000 0.6630521
## [2,] 0.6630521 1.0000000
sqrt(e1$values[1]) #cor(u1, v1)

## [1] 0.6630521
sqrt(e2$values[1])

```

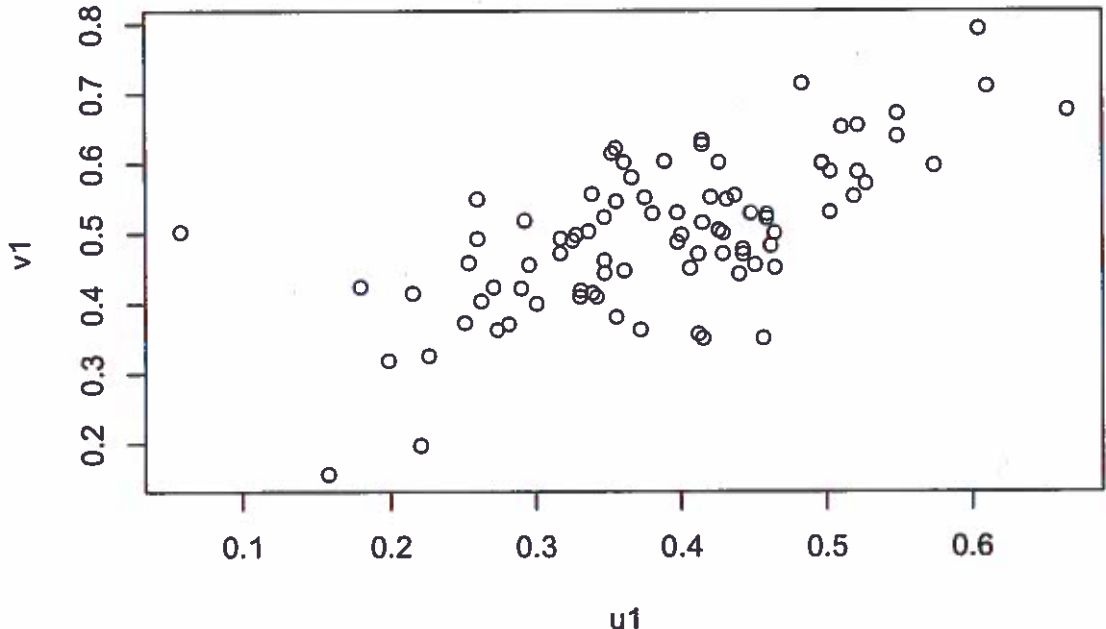
be build in
mchm

```

## [1] 0.6630521
## use build-in function cancel
u1=fancor(scov[,1:2],scov[,3:5])$xcoef[,1]*%t(scov[,1:2])
v1=fancor(scov[,1:2],scov[,3:5])$ycoef[,1]*%t(scov[,3:5])
plot(u1,v1)

```

first elat
X
Y



```

\cor(t(rbind(u1,v1)))

##          [,1]      [,2]
## [1,] 1.0000000 0.6630521
## [2,] 0.6630521 1.0000000
sqrt(e2$values[1])

## [1] 0.6630521
## cor(u1,u2) , cor(v1,v2)
u2=-e1$vector[,2]*%t(scov[,1:2]) aX
cor(t(rbind(u1,u2))) #cor(u1, u2)

##          [,1]      [,2]
## [1,] 1.000000e+00 -3.700511e-18
## [2,] -3.700511e-18 1.000000e+00

```

```

v2=e2$variables[,2]%*%t(scor[,3:5])
cor(t(rbind(v1,v2))) #cor(v1, v2)

##           [,1]      [,2]
## [1,] 1.000000e+00 3.490312e-16
## [2,] 3.490312e-16 1.000000e+00
cor(t(rbind(u1,v2))) #cor(u1, v2)

##           [,1]      [,2]
## [1,] 1.000000e+00 2.512154e-16
## [2,] 2.512154e-16 1.000000e+00

```

Covariance of Original and Canonical Variables

Let $U = A^T X$ and $V = B^T Y$ where $A = [a_1, \dots, a_p]$ and $B = [b_1, \dots, b_q]$. Then

- $U = (U_1, \dots, U_p)^T$ contains the p canonical variates from X .
- $V = (V_1, \dots, V_q)^T$ contains the q canonical variates from Y .

The canonical variates and original variables have covariance matrices

$$Cov(U, X) = Cov(A^T X, X) = A^T \Sigma_X$$

$$Cov(U, Y) = Cov(A^T X, Y) = A^T \Sigma_{XY}$$

$$Cov(V, X) = Cov(B^T Y, X) = B^T \Sigma_{YX}$$

$$Cov(V, Y) = Cov(B^T Y, Y) = B^T \Sigma_Y$$

The canonical variates and original variables have correlation matrices

$$Cor(U, X) = Cov(A^T X, \tilde{\Sigma}_X^{-1/2} X) = A^T \Sigma_X \tilde{\Sigma}_X^{-1/2}$$

$$Cor(U, Y) = Cov(A^T X, \tilde{\Sigma}_Y^{-1/2} Y) = A^T \Sigma_{XY} \tilde{\Sigma}_Y^{-1/2}$$

$$Cor(V, X) = Cov(B^T Y, \tilde{\Sigma}_X^{-1/2} X) = B^T \Sigma_{YX} \tilde{\Sigma}_X^{-1/2}$$

$$Cor(V, Y) = Cov(B^T Y, \tilde{\Sigma}_Y^{-1/2} Y) = B^T \Sigma_Y \tilde{\Sigma}_Y^{-1/2}$$

given that $Var(U_k) = Var(V_l) = 1$ for all k, l . Here $\tilde{\Sigma}_X^{-1/2}$ is a diagonal matrix containing X variances and $\tilde{\Sigma}_Y^{-1/2}$ is a diagonal matrix containing Y variances.

3. Canonical Variables Interpretation

Canonical variables are, in general, artificial. They have no physical meaning. the canonical coefficients a and b have units proportional to the original variables X and Y . It is common to standardize all the variables before performing the canonical correlation analysis. We may interpret the coefficients of the canonical variables similarly to how we interpret the coefficients of principal components.

For checking the condition $\Sigma_{XY} \neq 0$ since the canonical correlation can only be applied when $\Sigma_{XY} \neq 0$.

4. Large Sample Inferences

When $\Sigma_{XY} = 0$, $a^T X$ and $b^T Y$ have covariance $a^T \Sigma_{XY} b = 0$ for all vectors a and b . Consequently, all the canonical correlations must be zero and there is no point in pursuing a canonical correlation analysis. To test

$$H_0 : \Sigma_{XY} = 0_{p \times q} \text{ (All canonical correlations are zero)}$$

$$H_1 : \Sigma_{XY} \neq 0_{p \times q} \text{ (At least one canonical correlation significantly differs from zero)}$$

Let

$$Z_j = [X_j, Y_j]^T, j = 1, 2, \dots, n$$

be a random sample from an $N_{p+q}(\mu, \Sigma)$ population with

$$\Sigma = \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_Y \end{bmatrix}$$

Then the likelihood ratio test of $H_0 : \Sigma_{XY} = 0_{p \times q}$ versus $H_1 : \Sigma_{XY} \neq 0_{p \times q}$ rejects H_0 for large values of

$$-2 \ln \Lambda = n \ln \left(\frac{|\hat{S}_X| |\hat{S}_Y|}{|\hat{S}|} \right) = -n \ln \prod_{i=1}^p (1 - \hat{\rho}_i^2)$$

where \hat{S} is the unbiased estimator of Σ . For $n \rightarrow \infty$, the test statistic is approximately distributed as a χ_{pq}^2 . Bartlett suggests replacing the multiplicative factor n in the likelihood ratio statistic with the factor $-[n - 1 - \frac{1}{2}(p+q+1) \ln \prod_{i=1}^p (1 - \hat{\rho}_i^2)]$ to improve the χ^2 approximation to the sampling distribution of $-2 \ln \Lambda$. Thus, reject H_0 at significance level α if

$$-[n - 1 - \frac{1}{2}(p+q+1) \ln \prod_{i=1}^p (1 - \hat{\rho}_i^2)] > \chi_{pq}^2(\alpha)$$

Not necessary to have the raw data

Example: The Los Angeles Depression Study Data

We don't have the raw data

Canonical correlation analysis using a sample correlation matrix rather than the raw data matrix.

Data include $n=294$ individuals with variables

X are "health variables":

- CESD: A numerical measure of depression
- Health: A measure of general perceived health status

Y are "personal (~demographic) variables":

- Gender: Low=Male, High=Female
- Age
- Income
- Education Level

Suppose the sample correlation matrix R is given

```
R22 <- matrix( c(
1, .044, -.106, -.18,
.044, 1, -.208, -.192,
-.106, -.208, 1, .492,
-.18, -.192, .492, 1 ),
ncol=4, byrow=T)
```

$$\begin{bmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{bmatrix}$$

use formulas in page 3 to comp

+ From covariance matrix

\Rightarrow compute a, b, ρ

$$\begin{aligned} \mu &= a^T X \\ \eta &= b^T Y \end{aligned}$$

```
R11 <- matrix( c(
1,.212,.212,1),
ncol=2,byrow=T)
```

```
R12 <- matrix( c(
.124,-.164,-.101,-.158,
.098,.308,-.27,-.183),
ncol=4, byrow=T)
```

```
R21 <- t(R12)
```

because from page 3

Finding the E1 and E2 matrices:

```
E1 <- solve(R11) %*% R12 %*% solve(R22) %*% R21
E2 <- solve(R22) %*% R21 %*% solve(R11) %*% R12
eigen(E1)
```

contain
(g, a)

contain
(g, b)

```
## eigen() decomposition
## $values
## [1] 0.16763669 0.06806171
##
## $vectors
##      [,1]      [,2]
## [1,] 0.4560253 0.9476307
## [2,] -0.8899668 -0.3193681
```

$$a_1 = \begin{bmatrix} 0.46 & -0.89 \end{bmatrix}$$

```
## eigen() decomposition
## $values
## [1] 1.676367e-01+0.000000e+00i 6.806171e-02+0.000000e+00i
## [3] 3.483919e-19+6.105072e-18i 3.483919e-19-6.105072e-18i
##
## $vectors
##      [,1]      [,2]      [,3]
## [1,] 0.02473006+0i 0.4522578+0i 0.6926400+0.0000000i
## [2,] 0.90283535+0i -0.4598545+0i 0.2071287-0.1099990i
## [3,] -0.40965230+0i -0.4745706+0i 0.4675316-0.3579553i
## [4,] 0.12830336+0i -0.5989820+0i 0.0297314+0.3429957i
##
##      [,4]
## [1,] 0.6926400+0.0000000i
## [2,] 0.2071287+0.1099990i
## [3,] 0.4675316+0.3579553i
## [4,] 0.0297314-0.3429957i
```

← complex

The canonical correlations are:

```
canon.corr <- sqrt(eigen(E1)$values)
canon.corr
## [1] 0.4094346 0.2608864

canon.corr <- sqrt(eigen(E2)$values)
canon.corr
## [1] 4.094346e-01+0.00000e+00i 2.608864e-01+0.00000e+00i
## [3] 1.797693e-09+1.69803e-09i 1.797693e-09-1.69803e-09i
```

ρ_1

$\text{cor}(u_i, v_i) = \rho_i$ which are the eigen value of E_1
the same.

First canonical variate:

$$u_1 = 0.46CESD - 0.89Health = a_1 X$$

$$v_1 = 0.02Gender + 0.90Age - 0.41Education + 0.13Income = b_1 Y$$

Second canonical variate:

$$u_2 = -0.95CESD - 0.32Health = a_2 X$$

$$v_2 = -0.45Gender + 0.46Age + 0.47Education + 0.60Income = b_2 Y \quad (* \text{ flip the sides of } b_2)$$

Test for the significance of the first canonical correlation: The null hypothesis is that the first (and smaller) canonical correlations are zero.

part 4 page 6

```
n <- 294;
p <- 2;
q <- 4
test.stat <- -( (n-1) - 0.5*(p+q+1) ) * sum(log(1-eigen(E1)$values))
test.stat
```

```
## [1] 73.52575
P.value <- pchisq(test.stat, df = p*q, lower.tail=F)
P.value
```

```
## [1] 9.728813e-13
```

Reject H_0 The first pair has non zero correlation.

Since the P-value is small, we conclude that there is at least one nonzero canonical correlation.

* Test for the significance of the second canonical correlation: The null hypothesis is that the second (and smaller) canonical correlations are zero in general, but there's only two here.

```
test.stat <- -( (n-1) - 0.5*(p+q+1) ) * sum(log(1-eigen(E1)$values[-1]))
test.stat
```

```
## [1] 20.40647
P.value <- pchisq(test.stat, df = (p-1)*(q-1), lower.tail=F)
P.value
```

↑ exclude the first eigenvalue

```
## [1] 0.0001398029
```

because exclude 1st pair

The P-value is again very small, so we conclude there are at least two nonzero canonical correlations. In this case, that means exactly two nonzero canonical correlations.

Example: Decathlon Example

This example: Input raw data (value of $X \in Y$)

```
library(ade4)

## Warning: package 'ade4' was built under R version 3.4.4
data(olympic)
decathlon <- cbind(olympic$tab,olympic$score)
colnames(decathlon) <- c("run100", "long.jump", "shot", "high.jump", "run400", "hurdle",
  "discus", "pole.vault", "javelin", "run1500", "score")
# resign running events (so higher score means better performance)
#head(decathlon)
decathlon[,c(1,5,6,10)] <- (-1)*decathlon[,c(1,5,6,10)]
#head(decathlon)
```

to standardize data so that larger \rightarrow better performance

run100 longjump short highjump min400 hurdle discus pole.vault
javelin min

min500 score.

```
## 1 268.95 8488
## 2 273.02 8399
## 3 263.20 8328
## 4 285.11 8306
## 5 256.64 8286
## 6 274.07 8272
```

```
decathlon[,c(1,5,6,10)] <- (-1)*decathlon[,c(1,5,6,10)]
head(decathlon)
```

```
## run100 long.jump shot high.jump run400 hurdle discus pole.vault javelin
## 1 -11.25 7.43 15.48 2.27 -48.90 -15.13 49.28 4.7 61.32
## 2 -10.87 7.45 14.97 1.97 -47.71 -14.46 44.36 5.1 61.76
## 3 -11.18 7.44 14.20 1.97 -48.29 -14.81 43.66 5.2 64.16
## 4 -10.62 7.38 15.02 2.03 -49.06 -14.72 44.80 4.9 64.04
## 5 -11.02 7.43 12.92 1.97 -47.44 -14.40 41.20 5.2 57.46
## 6 -10.83 7.72 13.58 2.12 -48.34 -14.18 43.06 4.9 52.18
```

```
## run1500 score
## 1 -268.95 8488
## 2 -273.02 8399
## 3 -263.20 8328
## 4 -285.11 8306
## 5 -256.64 8286
## 6 -274.07 8272
```

```
# check variable scales
apply(decathlon,2,mean)
```

```
## run100 long.jump shot high.jump run400 hurdle
## -11.223529 7.095000 13.850882 1.974412 -49.366176 -15.107647
## discus pole.vault javelin run1500 score
## 41.905294 4.676471 58.840588 -276.191471 7782.852941
```

```
apply(decathlon,2,sd)
```

```
## run100 long.jump shot high.jump run400 hurdle
## 0.2872322 0.3738680 1.5019268 0.1044811 1.1755463 0.6056555
## discus pole.vault javelin run1500 score
## 4.5007105 0.4930172 6.4387360 13.4781327 594.5827227
```

```
# separate into running/jumping vs throwing/arm events
X <- as.matrix(decathlon[,c("shot", "discus", "javelin", "pole.vault")])
Y <- as.matrix(decathlon[,c("run100", "run400", "run1500",
" hurdle", "long.jump", "high.jump")])
```

```
n <- nrow(X)
p <- ncol(X)
q <- ncol(Y)
```

CCA via Covariance

```
# canonical correlations of covariance (unstandardized data)
cca <- cancor(X, Y)
```

```
# cca (the normal way)
Sx <- cov(X)
```

check the center tendency

check the spread out tendency

Build in function

Manually

```

Sy <- cov(Y)
Sxy <- cov(X,Y)
Sxeig <- eigen(Sx, symmetric=TRUE)
Sxisqrt <- Sxeig$vector %>% diag(1/sqrt(Sxeig$values)) %>% t(Sxeig$vector)
Syeig <- eigen(Sy, symmetric=TRUE)
Sysisqrt <- Syeig$vector %>% diag(1/sqrt(Syeig$values)) %>% t(Syeig$vector)
Xmat <- Sxisqrt %>% Sxy %>% solve(Sy) %>% t(Sxy) %>% Sxisqrt
Ymat <- Sysisqrt %>% t(Sxy) %>% solve(Sx) %>% Sxy %>% Sysisqrt
Keig <- eigen(Xmat, symmetric=TRUE)
Yeig <- eigen(Ymat, symmetric=TRUE)

```

$$\lambda_{mat} = \Sigma_X^{-1/2} \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX} \Sigma_X^{-1/2}$$

```

# compare correlations (same)
cca$cor

```

```
## [1] 0.7702006 0.5033532 0.4184145 0.3052556
```

```
rho <- sqrt(Keig$values)
rho
```

$$\rho = \text{Eigenvalue of } X = \sqrt{\rho^2}$$

```
## [1] 0.7702006 0.5033532 0.4184145 0.3052556
```

```
sqrt(Yeig$values[1:p])
```

```
## [1] 0.7702006 0.5033532 0.4184145 0.3052556
```

```
# compare linear combinations (different!)
```

```

Ahat <- Sxisqrt %>% Keig$vector
Bhat <- Sysisqrt %>% Yeig$vector
sum((cca$xcoef - Ahat)^2)

```

```
## [1] 6.710414
```

```
sum((cca$ycoef[,1:p] - Bhat[,1:p])^2)
```

```
## [1] 42.98483
```

NOTE*: we need to multiply R's xcoef and ycoef by $\sqrt{n-1}$ to obtain the results we are expecting.

```
# compare linear combinations (same!)
```

```

Ahat <- Sxisqrt %>% Keig$vector
Bhat <- Sysisqrt %>% Yeig$vector
sum((cca$xcoef * sqrt(n-1) - Ahat)^2)

```

```
## [1] 3.031301e-28
```

```
sum((cca$ycoef[,1:p] * sqrt(n-1) - Bhat[,1:p])^2)
```

```
## [1] 2.414499e-25
```

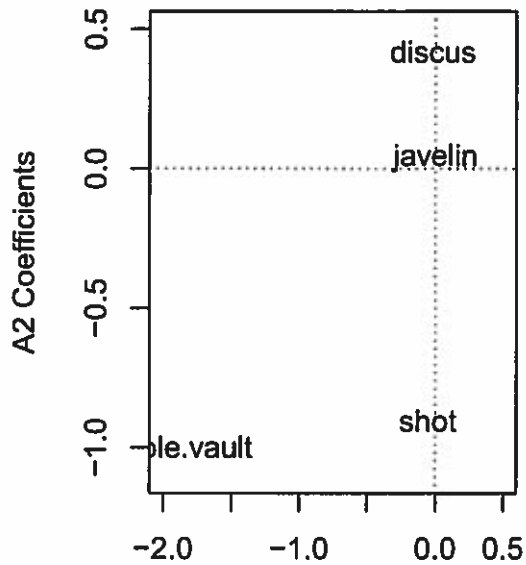
```

# plot coefficients
par(mfrow=c(1,2))
plot(Ahat[,1:2], xlab="A1 Coefficients", ylab="A2 Coefficients",
     type="n", main="X Coefficients", xlim=c(-2, 0.5), ylim=c(-1.1, 0.5))
text(Ahat[,1:2], labels=colnames(X))
abline(0,0,lty=3)
abline(v=0,lty=3)
plot(Bhat[,1:2], xlab="B1 Coefficients", ylab="B2 Coefficients",
     type="n", main="Y Coefficients", xlim=c(-2, 0.2), ylim=c(-2, 6))
text(Bhat[,1:2], labels=colnames(Y))

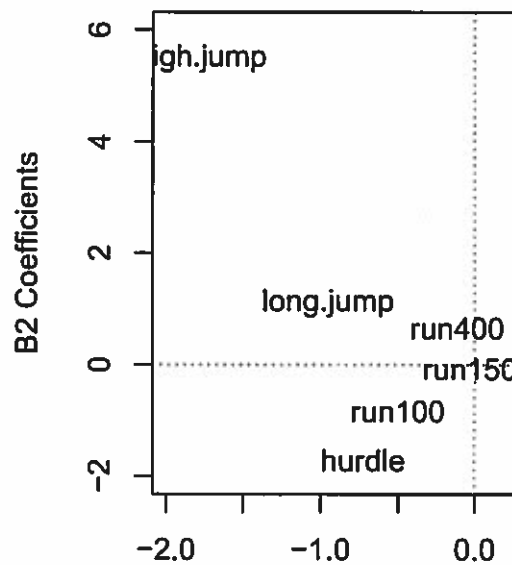
```

```
abline(0,0,lty=3)
abline(v=0,lty=3)
```

X Coefficients



Y Coefficients



A1 Coefficients

B1 Coefficients

If $X = \{x_{ij}\}_{n \times p}$ and $Y = \{y_{ij}\}_{n \times q}$, then

- $\hat{U} = \hat{A}^T X = \{\hat{u}_{ij}\}_{n \times p}$ where columns of \hat{U} contain the canonical variables for the X set.
- $\hat{V} = \hat{B}^T Y = \{\hat{v}_{ij}\}_{n \times q}$ where columns of \hat{V} contain the canonical variables for the Y set.

define canonical variates

```
U <- X %*% Ahat
```

```
V <- Y %*% Bhat
```

canonical variable covariances

```
round(cov(U),4)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]  1   0   0   0
## [2,]  0   1   0   0
## [3,]  0   0   1   0
## [4,]  0   0   0   1
```

```
round(cov(V),4)
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]  1   0   0   0   0   0
## [2,]  0   1   0   0   0   0
## [3,]  0   0   1   0   0   0
## [4,]  0   0   0   1   0   0
## [5,]  0   0   0   0   1   0
## [6,]  0   0   0   0   0   1
```



```

round(cov(U,V),4)

##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,] 0.7702 0.0000 0.0000 0.0000 0 0
## [2,] 0.0000 0.5034 0.0000 0.0000 0 0
## [3,] 0.0000 0.0000 0.4184 0.0000 0 0
## [4,] 0.0000 0.0000 0.0000 0.3053 0 0

# covariance of original and canonical variables (U and X)
Ainv <- solve(Ahat)
sum( ( cov(U, X) - crossprod(Ahat, Sx) )^2 )

## [1] 3.396329e-30
sum( ( Sx - crossprod(Ainv) )^2 )

## [1] 4.364327e-27
Sxhat <- matrix(0, p, p)
for(j in 1:p) Sxhat <- Sxhat + outer(Ainv[j,], Ainv[j,])
sum( (Sx - Sxhat)^2 )

## [1] 4.364327e-27
# covariance of original and canonical variables (V and Y)
Binv <- solve(Bhat)
sum( ( cov(V, Y) - crossprod(Bhat, Sy) )^2 )

## [1] 1.696269e-28
sum( ( Sy - crossprod(Binv) )^2 )

## [1] 3.027024e-26
Syhat <- matrix(0, q, q)
for(j in 1:q) Syhat <- Syhat + outer(Binv[j,], Binv[j,])
sum( (Sy - Syhat)^2 )

## [1] 3.027024e-26
# covariance of original and canonical variables (U and Y)
sum( (cov(U, Y) - crossprod(Ahat, Sxy))^2 )

## [1] 2.071712e-29
# covariance of original and canonical variables (V and X)
sum( (cov(V, X) - crossprod(Bhat, t(Sxy)))^2 )

## [1] 2.943246e-28
# covariance of canonical variables (U and V)
rhomat <- cbind(diag(rho), matrix(0, p, q-p))
sum( (cov(U, V) - rhomat)^2 )

## [1] 1.241068e-27
sum( (Sxy - crossprod(Ainv, rhomat) %*% Binv)^2 )

## [1] 1.355523e-25

```

```
Sxyhat <- matrix(0, p, q)
for(j in 1:p) Sxyhat <- Sxyhat + rho[j] * outer(Ainv[j,], Binv[j,])
sum( (Sxy - Sxyhat)^2 )
```

```
## [1] 1.37319e-25
```

Error of approximation matrices provide a descriptive measure of how well the first r pairs of canonical variables explain the covariation in the data. The error of approximation matrices are

$$S_X = \sum_{j=1}^r [a^{(j)}][a^{(j)}]^T$$

$$S_Y = \sum_{j=1}^r [b^{(j)}][b^{(j)}]^T$$

$$S_{XY} = \sum_{j=1}^r \rho_j [a^{(j)}][b^{(j)}]^T$$

```
# error of approximation matrices (with r=2)
Ainv <- solve(Ahat)
Binv <- solve(Bhat)
r <- 2
Ex <- Sx - crossprod(Ainv[1:r,])
Ey <- Sy - crossprod(Binv[1:r,])
Exy <- Sxy - crossprod(diag(rho[1:r]) %*% Ainv[1:r,], Binv[1:r,])
```

```
# get norms of error matrices
sqrt(mean(Ex^2))
```

```
## [1] 6.561393
```

```
sqrt(mean(Ey^2))
```

```
## [1] 18.37339
```

```
sqrt(mean(Exy^2))
```

```
## [1] 1.725392
```

CCA via Correlation (standardized CCA)

```
# standardize data
```

```
Xs <- scale(X)
```

```
Ys <- scale(Y)
```

```
# canonical correlations of correlations (standardized data)
```

```
ccas <- cancel(Xs, Ys)
```

```
# cca (the normal way)
```

```
Sx <- cov(Xs)
```

```
Sy <- cov(Ys)
```

```
Sxy <- cov(Xs, Ys)
```

```
Sxeig <- eigen(Sx, symmetric=TRUE)
```

```
Sxisqrt <- Sxeig$eigenvectors %*% diag(1/sqrt(Sxeig$values)) %*% t(Sxeig$eigenvectors)
```

```

Syeig <- eigen(Sy, symmetric=TRUE)
Syisqrt <- Syeig$vectors %*% diag(1/sqrt(Syeig$values)) %*% t(Syeig$vectors)
Xmat <- Sxisqrt %*% Sxy %*% solve(Sy) %*% t(Sxy) %*% Sxisqrt
Ymat <- Syisqrt %*% t(Sxy) %*% solve(Sx) %*% Sxy %*% Syisqrt
Xeig <- eigen(Xmat, symmetric=TRUE)
Yeig <- eigen(Ymat, symmetric=TRUE)

# compare correlations (same)
cca$cor

## [1] 0.7702006 0.5033532 0.4184145 0.3052556
sqrt(Xeig$values)

## [1] 0.7702006 0.5033532 0.4184145 0.3052556
sqrt(Yeig$values[1:p])

## [1] 0.7702006 0.5033532 0.4184145 0.3052556
# compare linear combinations (different?)
Ahat <- Sxisqrt %*% Xeig$vectors
Bhat <- Syisqrt %*% Yeig$vectors
sum((ccas$xcoef * sqrt(n-1) - Ahat)^2)

## [1] 3.332536e-29
sum((ccas$ycoef[,1:p] * sqrt(n-1) - Bhat[,1:p])^2)

## [1] 11.59453
# note that the signing is arbitrary!!
ccas$ycoef[,1:p] * sqrt(n-1)

##           [,1]      [,2]      [,3]      [,4]
## run100   -0.1439138 -0.2404940  0.5274876 -0.13754449
## run400   -0.1373435  0.7655659 -1.2826821  0.96359176
## run1500  0.3023537 -1.0519285 -0.1514027 -0.52923644
## hurdle   -0.4396044 -1.0374417  0.6303782  0.49905604
## long.jump -0.3564702  0.4110878 -0.0253127 -1.09325282
## high.jump -0.1855627  0.5731149 -0.2615838 -0.09007821
Bhat[,1:p]

##           [,1]      [,2]      [,3]      [,4]
## [1,]  0.1439138 -0.2404940 -0.5274876 -0.13754449
## [2,]  0.1373435  0.7655659  1.2826821  0.96359176
## [3,] -0.3023537 -1.0519285  0.1514027 -0.52923644
## [4,]  0.4396044 -1.0374417 -0.6303782  0.49905604
## [5,]  0.3564702  0.4110878  0.0253127 -1.09325282
## [6,]  0.1855627  0.5731149  0.2615838 -0.09007821
Bhat[,1:p] <- Bhat[,1:p] %*% diag(c(-1,1,-1,1))
sum((ccas$ycoef[,1:p] * sqrt(n-1) - Bhat[,1:p])^2)

## [1] 1.132493e-28
# plot coefficients
par(mfrow=c(1,2))

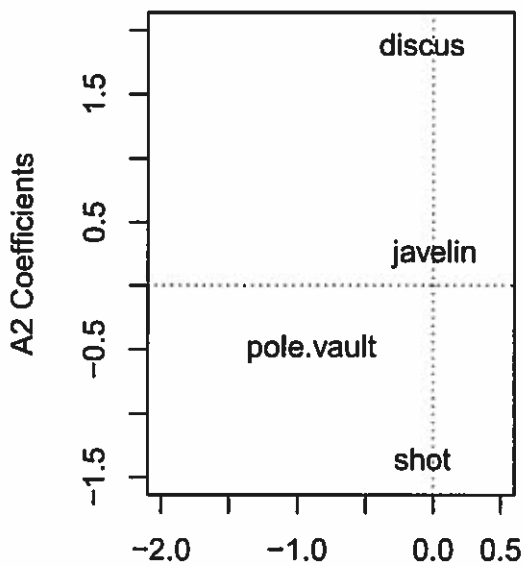
```

```

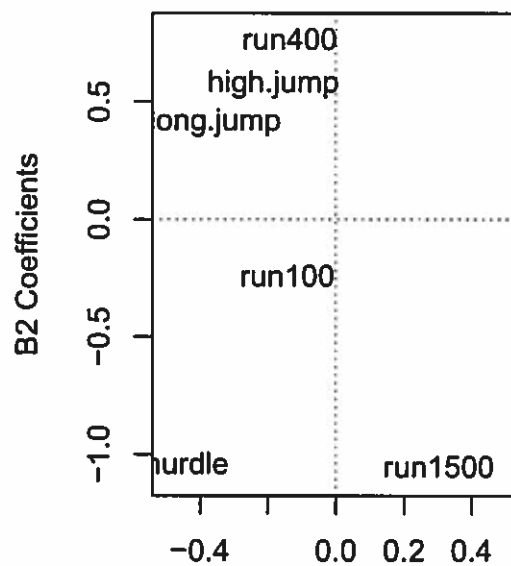
plot(Ahat[,1:2], xlab="A1 Coefficients", ylab="A2 Coefficients",
     type="n", main="X Coefficients", xlim=c(-2, 0.5), ylim=c(-1.5, 2))
text(Ahat[,1:2], labels=colnames(X))
abline(0,0,lty=3)
abline(v=0,lty=3)
plot(Bhat[,1:2], xlab="B1 Coefficients", ylab="B2 Coefficients",
     type="n", main="Y Coefficients", xlim=c(-0.5, 0.5), ylim=c(-1.1, 0.8))
text(Bhat[,1:2], labels=colnames(Y))
abline(0,0,lty=3)
abline(v=0,lty=3)

```

X Coefficients



Y Coefficients



A1 Coefficients

B1 Coefficients

```

# define canonical variates
U <- Xs %*% Ahat
V <- Ys %*% Bhat

# canonical variable covariances
round(cov(U),4)

##      [,1] [,2] [,3] [,4]
## [1,]  1  0  0  0
## [2,]  0  1  0  0
## [3,]  0  0  1  0
## [4,]  0  0  0  1

round(cov(V),4)

##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]  1  0  0  0  0  0
## [2,]  0  1  0  0  0  0
## [3,]  0  0  1  0  0  0

```

```

## [4,] 0 0 0 1 0 0
## [5,] 0 0 0 0 1 0
## [6,] 0 0 0 0 0 1
round(cov(U,V),4)

##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,] 0.7702 0.0000 0.0000 0.0000 0 0
## [2,] 0.0000 0.5034 0.0000 0.0000 0 0
## [3,] 0.0000 0.0000 0.4184 0.0000 0 0
## [4,] 0.0000 0.0000 0.0000 0.3053 0 0

# covariance of original and canonical variables (U and Xs)
Ainv <- solve(Ahat)
sum( ( cov(U, Xs) - crossprod(Ahat, Sx) )^2 )

## [1] 2.759323e-31
sum( ( Sx - crossprod(Ainv) )^2 )

## [1] 6.569732e-30
Sxhat <- matrix(0, p, p)
for(j in 1:p) Sxhat <- Sxhat + outer(Ainv[j,], Ainv[j,])
sum( ( Sx - Sxhat )^2 )

## [1] 6.569732e-30

# covariance of original and canonical variables (V and Ys)
Binv <- solve(Bhat)
sum( ( cov(V, Ys) - crossprod(Bhat, Sy) )^2 )

## [1] 2.406961e-31
sum( ( Sy - crossprod(Binv) )^2 )

## [1] 3.136492e-29
Syhat <- matrix(0, q, q)
for(j in 1:q) Syhat <- Syhat + outer(Binv[j,], Binv[j,])
sum( ( Sy - Syhat )^2 )

## [1] 3.136492e-29

# covariance of original and canonical variables (U and Ys)
sum( (cov(U, Ys) - crossprod(Ahat, Sxy))^2 )

## [1] 5.477785e-32

# covariance of original and canonical variables (V and Xs)
sum( (cov(V, Xs) - crossprod(Bhat, t(Sxy)))^2 )

## [1] 1.336149e-31

# covariance of canonical variables (U and V)
rhomat <- cbind(diag(rho), matrix(0, p, q-p))
sum( (cov(U, V) - rhomat)^2 )

## [1] 1.272906e-29
sum( (Sxy - crossprod(Ainv, rhomat) %*% Binv)^2 )

```

```

## [1] 7.505349e-30
Sxyhat <- matrix(0, p, q)
for(j in 1:p) Sxyhat <- Sxyhat + rho[j] * outer(Ainv[j,], Binv[j,])
sum( (Sxy - Sxyhat)^2 )

## [1] 7.283289e-30
# error of approximation matrices (with r=2)
Ainv <- solve(Ahat)
Binv <- solve(Bhat)
r <- 2
Ex <- Sx - crossprod(Ainv[1:r,])
Ey <- Sy - crossprod(Binv[1:r,])
Exy <- Sxy - crossprod(diag(rho[1:r]) %*% Ainv[1:r,], Binv[1:r,])

# get norms of error matrices
sqrt(mean(Ex^2))

## [1] 0.2432351
sqrt(mean(Ey^2))

## [1] 0.2296716
sqrt(mean(Exy^2))

## [1] 0.07458264

```

Canonical correlation analysis.

→ explore the relationship between two ^(vector) multivariate sets of variables

- x [exercise : climbing rate, fast can run, amount weight lifted,
- health : blood pressure, cholesterol measure, glucose level, body max.. ○
- sale performance
- optical variable

Chapter 7. Discrimination and Classification

Jianxuan Liu

Fall 2018

Main objectives of discrimination and classification are:

1. Separate distinct sets of objects: discrimination.
2. Classify new objects into well defined populations: classification.

Examples of discrimination:

- A tumor is benign or malignant, and the correct diagnosis needs to be obtained.
- In the finance and credit risk area, one wants to assess whether a company is likely to go bankrupt in the next few years or whether a client will default on mortgage repayments.

Example of classification:

Using information on prisoners eligible for parole (good behavior, history of drug use, job skills, etc) can we successfully allocate a prisoners eligible for parole into two groups: those who will commit another crime or those who will not commit another crime?

Discriminant analysis and classification are termed "supervised learning" while clustering is called "unsupervised learning" in the machine learning literature. Discriminant analysis and classification are slightly different actions, but they are used interchangeably.

In practice, the two objectives often overlap: A function of the p variables that serves as a discriminant is also used for classifying a new object into one of the populations. An allocation or classification rule can often serve as a discriminant.

The setup is the usual one: p variables are measured on n sample units or subjects. We wish to find a function of the variables that will optimize the discrimination between units belonging to different populations (minimize classification errors).

1. Classifying Two Populations

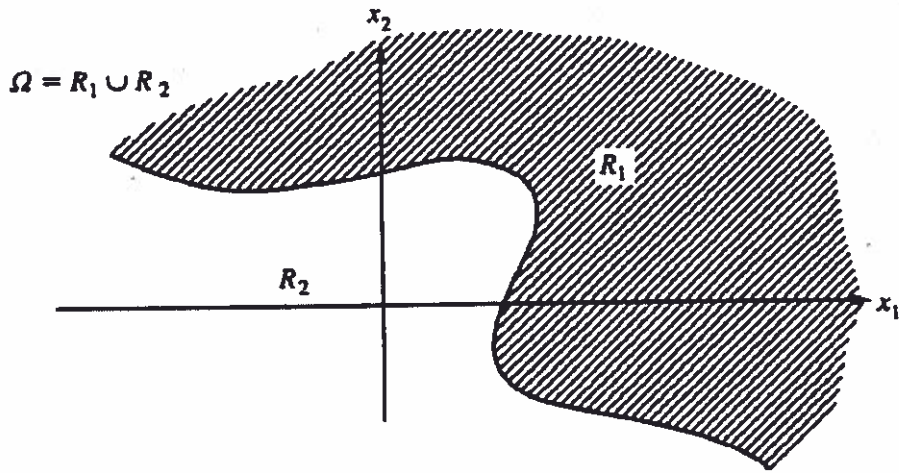
We wish to separate two populations and also allocate new units to one of the two populations. Let $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ denote a p -dimensional random vector. Let $f_1(\mathbf{x})$ denote the probability density function (pdf) for population π_1 and $f_2(\mathbf{x})$ denote the probability density function (pdf) for population π_2 . Given a realization $\mathbf{X} = \mathbf{x}$, we want to assign \mathbf{x} to π_1 or π_2 . We want to find some classification rule to determine whether a realization $\mathbf{X} = \mathbf{x}$ should be assigned to population π_1 or π_2 .

Let Ω denote the sample space (collection of all possible values for \mathbf{x}).

- R_1 is the set of values of \mathbf{x} for which we classify objects into π_1 and
- $R_2 = \Omega - R_1$ is the set of values of \mathbf{x} for which we classify objects into π_2 .

Since every object belongs into one of the two populations, we have that

$$\Omega = R_1 \cup R_2 \quad R_1 \cap R_2 = \phi$$

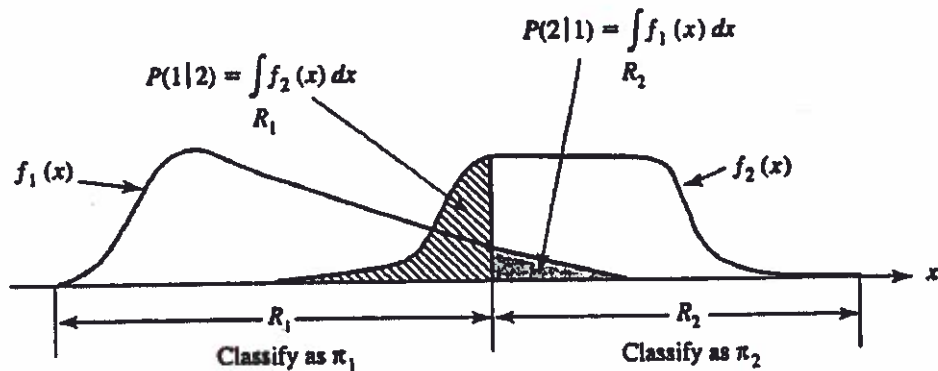


Ignore for now the prior probabilities of each population and the potentially different misclassification costs. The probability of misclassifying an object into π_2 when it belongs in π_1 is

$$P(2|1) = P(X \in R_2 | \pi_1) = \int_{R_2} f_1(x) dx$$

and probability of misclassifying an object into π_1 when it belongs in π_2 is

$$P(1|2) = P(X \in R_1 | \pi_2) = \int_{R_1} f_2(x) dx$$



In the above graph ($p = 1$ variable), the solid shaded region is $P(2|1)$ and the striped region is $P(1|2)$.

Probability of Misclassification

Let p_1, p_2 denote the prior probabilities of π_1, π_2 , respectively, with the constraint that $p_1 + p_2 = 1$.

The overall probabilities of the four outcomes have the form

$$\begin{aligned}
 P(\pi_1 \text{ and } R_1) &= P(\text{correctly classify as } \pi_1) = P(X \in R_1 | \pi_1)P(\pi_1) = P(1|1)p_1 \\
 P(\pi_2 \text{ and } R_2) &= P(\text{correctly classify as } \pi_2) = P(X \in R_2 | \pi_2)P(\pi_2) = P(2|2)p_2
 \end{aligned}$$

$$P(\pi_1 \text{ and } R_2) = \frac{P(R_2 | \pi_1) P(\pi_1)}{P(\text{misclassify } \pi_1 \text{ as } \pi_2)} = P(X \in R_2 | \pi_1) P(\pi_1) = P(2 | 1) p_1$$

$$P(\text{misclassify } \pi_2 \text{ as } \pi_1) = P(X \in R_1 | \pi_2) P(\pi_2) = P(1 | 2) p_2$$

Further, let $c(1|2)$ and $c(2|1)$ be the costs of misclassifying an object into π_2 and π_1 , respectively.

In many real world cases, costs of misclassification are not equal, for example:

- (a). π_1 and π_2 are diseased and healthy
- (b). π_1 and π_2 are guilty and not guilty
- (c). π_1 and π_2 are buy and not buy stock

We can make a cost matrix to tabulate our misclassification costs:

		Classify as:	
		π_1	π_2
Truth	π_1	0	$c(2 1)$
	π_2	$c(1 2)$	0

(ECM)

Classification rules are often evaluated in terms of the expected cost of misclassification or ECM:

$$ECM = c(2|1)P(2|1)p_1 + c(1|2)P(1|2)p_2 \quad \text{cost}(\pi_1) P(R_2 \text{ and } \pi_1) + P(R_1)$$

There is no cost when units are correctly classified. We seek rules that minimize the ECM. This leads to an **optimal classification rule**: classify an object into π_1 if

$$R_1 \text{ iff } \frac{f_1(x)c(2|1)p_1}{f_2(x)c(1|2)p_2} > 1. \quad \text{when the "cost" that we assign } R_2 \text{ while } \pi_1 \text{ is bigger than another one.}$$

Equivalently, the regions R_1, R_2 that minimize the ECM are defined by the values of x for which

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \underbrace{\left(\frac{c(1|2)}{c(2|1)} \right)}_{\text{ratio of misclassification cost}} \underbrace{\left(\frac{p_2}{p_1} \right)}_{\text{ratio of prior prob}}$$

$$R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

just need to know the ratio of the cost

$$ECM = c(2|1) \int_{R_2} f_1(x) dx p_1 + c(1|2) \int_{R_1} f_2(x) dx p_2$$

$$\text{We have } \Omega = R_1 + R_2 \Rightarrow 1 = \int_{R_1} f_1(x) dx + \int_{R_2} f_2(x) dx$$

Then

$$ECM = c(2|1) p_1 \left(1 - \int_{R_1} f_1(x) dx \right) + c(1|2) p_2 \left(\int_{R_1} f_2(x) dx \right)$$

$$= \int_{R_1} \left[c(1|2) p_2 f_2(x) - c(2|1) p_1 f_1(x) \right] dx + c(2|1) p_1$$

ECM is minimized when R_1 is chosen to be the region where (*) ≤ 0

$$\Rightarrow c(1|2) p_2 f_2(x) \leq c(2|1) p_1 f_1(x)$$

$$\Rightarrow R_1: \frac{f_1}{f_2} > \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}$$

Remark: Implementation of the minimum ECM rule require for a new unit requires evaluation of f_1 and f_2 at the new vector of observations \mathbf{x}_0 , but it does not require knowing the two costs or the two prior probabilities, just their ratio.

Special Cases of Minimum ECM

When the prior probabilities or the misclassification costs are equal, the classification rule above simplifies correspondingly.

- If $p_1 = p_2$, then $R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > \frac{c(1|2)}{c(2|1)}$, and $R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{c(1|2)}{c(2|1)}$
- If $c(1|2) = c(2|1)$, then $R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > \frac{p_2}{p_1}$, and $R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1}$
- If $p_1 = p_2$ and $c(1|2) = c(2|1)$, then $R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > 1$, and $R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < 1$

If \mathbf{x} is on the boundary between R_1 and R_2 then toss a coin, or randomly classify in some way.

Other Criteria for Choosing a Classification Rule

1. Highest posterior probability criteria

We might consider classifying a new unit with observation \mathbf{x}_0 into the population with the highest posterior probability $P(\pi_i | \mathbf{x}_0)$. By Bayes rule

$$P(\pi_1 | \mathbf{x}_0) = \frac{P(\mathbf{x}_0 | \pi_1)p_1}{P(\mathbf{x}_0 | \pi_1)p_1 + P(\mathbf{x}_0 | \pi_2)p_2} = \frac{p_1 f_1(\mathbf{x}_0)}{p_1 f_1(\mathbf{x}_0) + p_2 f_2(\mathbf{x}_0)}$$

$P(\pi_1 | \mathbf{x}_0) > 0.5$
 \Rightarrow assign to π_1
 $P(\pi_1 | \mathbf{x}_0) < 0.5 \Rightarrow$ assign to π_2

Clearly, $P(\pi_2 | \mathbf{x}_0) = 1 - P(\pi_1 | \mathbf{x}_0)$. Using the posterior probability criterion, we classify a unit with measurements \mathbf{x}_0 into π_1 when $P(\pi_1 | \mathbf{x}_0) > P(\pi_2 | \mathbf{x}_0)$

2. Minimum total probability of misclassification (TPM) criteria

The **total probability of misclassification (TPM)** ignores the cost of misclassification and is defined as the probability of either misclassifying a π_1 observation or misclassifying a π_2 observation:

$$TPM = p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}.$$

TPM is equivalent to ECM when costs are equal. An optimal rule in this sense would minimize TPM. The optimal TPM regions are the same as the special case of $c(1|2) = c(2|1)$. That is

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > \frac{p_2}{p_1} \quad \text{and} \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1}$$

2. Two Multivariate Normal Populations with $\Sigma_1 = \Sigma_2 = \Sigma$

($n \times p$) matrix n observations p variables

Let $\mathbf{X} = (X_1, \dots, X_p)^T$ denote a random vector. We now assume that $f_1(\mathbf{x}) \sim N(\mu_1, \Sigma)$ and $f_2(\mathbf{x}) \sim N(\mu_2, \Sigma)$ denote the pdf for population π_1 and π_2 , respectively. Given a realization $\mathbf{X} = \mathbf{x}$, we want to find some classification rule to determine whether a realization $\mathbf{X} = \mathbf{x}$ should be assigned to population π_1 or π_2 . The multivariate normal densities have the form

$$f_k(\mathbf{x}) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k)\right\}, \quad k = 1, 2$$

which implies that

$$\begin{aligned}
 f^* &= \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \\
 &= \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2) \right\} \\
 &= \exp \left\{ (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1}(\mu_1 + \mu_2) \right\}.
 \end{aligned}$$

Then, R_1 is given by the set of \mathbf{x} values for which:

$$\begin{cases}
 R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \\
 R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right).
 \end{cases}$$

Given the definitions of R_1, R_2 , an allocation that minimizes the ECM is the following: allocate \mathbf{x}_0 to π_1 if

$$\ln \left(\frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} \right) = (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x}_0 - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1}(\mu_1 + \mu_2) > \ln \left[\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \right] \Leftrightarrow \ln(f^*) > \ln \left(\frac{c(1|2)}{c(2|1)} \right)$$

allocate \mathbf{x}_0 to π_2 otherwise.

Since μ_1, μ_2, Σ are typically unknown, in practice we use $\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2$ as estimators of the population means, and \mathbf{S}_{pooled} as estimator of the common covariance matrix Σ . Given n_1 independent observations from π_1 and n_2 independent observations from π_2 , we can estimate the needed parameters $\hat{\mu}_1 = \bar{\mathbf{x}}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbf{x}_{1j}$, $\hat{\mu}_2 = \bar{\mathbf{x}}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} \mathbf{x}_{2j}$ and

$$\hat{\Sigma} = \mathbf{S}_{pooled} = \frac{1}{n_1 + n_2 - 2} \left[\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T \right].$$

$S_1 = \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (X_{1j} - \bar{x}_1)(X_{1j} - \bar{x}_1)$

The estimated classification rule replaces f^* with its sample estimate:

$$\hat{f}^* = \exp \left\{ (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 - \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \right\}.$$

• If $\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) = 1$, then $\ln(1) = 0$ and the rule becomes

$$R_1 : \hat{y} > \hat{m} \text{ and } R_2 : \hat{y} < \hat{m},$$

where the scalar variables

$$\begin{aligned}
 \hat{y} &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 = \mathbf{a}^T \mathbf{x} & \mathbf{a} &= \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \\
 \hat{m} &= \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) = \frac{1}{2}(\hat{y}_1 + \hat{y}_2)
 \end{aligned}$$

with

- $\hat{y}_1 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} \mathbf{x}_1 = \hat{\mathbf{a}}^T \mathbf{x}_1$ and
- $\hat{y}_2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} \mathbf{x}_2 = \hat{\mathbf{a}}^T \mathbf{x}_2$.

The coefficient vector $\hat{\mathbf{a}} = \mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$ is unique only up to a multiplicative constant. Therefore, $\hat{\mathbf{a}}$ is frequently "scaled" or "normalized" to ease the interpretation of its elements. Two most commonly methods for normalization are

1. $\hat{a}^* = \frac{\hat{a}}{\sqrt{\hat{a}^T \hat{a}}}$, \hat{a}^* has unit length

2. $\hat{a}^* = \frac{\hat{a}}{\hat{a}_1}$, the first element of the new coefficient vector \hat{a}^* is 1.

Fisher's Linear Discriminant Function

R. A. Fisher considered finding the linear combination $Y = \mathbf{a}^T \mathbf{x}$ that best separates the groups:

$$\text{separation} = \frac{|\bar{y}_1 - \bar{y}_2|}{s_y}$$

where $s_y^2 = \frac{\sum_{j=1}^{n_1} (\bar{y}_{1j} - \bar{y}_1)^2 + \sum_{j=1}^{n_2} (\bar{y}_{2j} - \bar{y}_2)^2}{n_1 + n_2 - 2}$ is the pooled estimate of the variance / \bar{y}_1 is the mean of the Y scores for the observations from π_1 / \bar{y}_2 is the mean of the Y scores for the observations from π_2 / The objective is to select the linear combination of the \mathbf{x} to achieve maximum separation of the sample means \bar{y}_1 and \bar{y}_2 .

Result: The linear combination $\hat{y} = \hat{\mathbf{a}}^T \mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} \mathbf{x}_0$ maximize the ratio

$$\text{separation}^2 = \frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(\hat{\mathbf{a}}^T \bar{\mathbf{x}}_1 - \hat{\mathbf{a}}^T \bar{\mathbf{x}}_2)^2}{\hat{\mathbf{a}}^T \mathbf{S}_{pooled} \hat{\mathbf{a}}} = \frac{(\hat{\mathbf{a}}^T \mathbf{d})^2}{\hat{\mathbf{a}}^T \mathbf{S}_{pooled} \hat{\mathbf{a}}}$$

over all possible coefficient vectors $\hat{\mathbf{a}}$ where $\mathbf{d} = \bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$. The maximum of the ratio is

$$D^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2).$$

which is equal to the squared distance between the multivariate means.

$$\max_{\mathbf{x}^T \mathbf{D} \mathbf{x}} \frac{(\mathbf{x}^T \mathbf{d})^2}{\mathbf{x}^T \mathbf{D} \mathbf{x}} = \mathbf{d}^T \mathbf{D}^{-1} \mathbf{d} \quad \text{when } \mathbf{x} = \mathbf{C} \mathbf{D}^{-1} \mathbf{d} \quad \mathbf{C} \neq \mathbf{0}$$

A test of $H_0: \mu_1 = \mu_2$ is also a test for the hypothesis of no separation. If Hotelling's T^2 fails to reject H_0 , then the data will not provide a useful classification rule.

Example: Classifying gender using dental data

```
library(reshape) # to create data frame from dental data
## Warning: package 'reshape' was built under R version 3.4.4
library(MASS) # has lda and gda functions
library(heavy) # has dental data
## Warning: package 'heavy' was built under R version 3.4.3
library(klaR) # visualize the LDA partitions
## Warning: package 'klaR' was built under R version 3.4.4
data(dental)
d2 = cast(melt(dental, id=c("Subject", "age", "Sex")), Subject+Sex+age)
names(d2)[3:6] = c("d8", "d10", "d12", "d14")
d2 = d2[, -1]
g = 2
p = ncol(d2) - 1
```

Time series
depends on time.
[=] 1st subject.

reshaped

Recall: $\text{Dim} \text{Sex} \text{ d8 d10 d12 d14}$

use their
measure to
detect their
gender

= 4 times.

```

# pooled covariances matrix
Sp=matrix(0, p, p)
nx=rep(0, g)
lev=levels(d2$Sex)
for(k in 1:g){
  x=d2[d2$Sex==lev[k],1:p+1]
  nx[k]=nrow(x)
  Sp=Sp + cov(x) * (nx[k] - 1)
}
Sp=Sp / (sum(nx) - g)
round(Sp,3)

##      [,1] [,2] [,3] [,4]
## [1,] 5.415 2.717 3.910 2.710
## [2,] 2.717 4.185 2.927 3.317
## [3,] 3.910 2.927 6.456 4.131
## [4,] 2.710 3.317 4.131 4.986

# fit lda model
# assume equal prior
ldamod=lda(d2$Sex ~ ., data=d2, prior=rep(1/2, 2))
names(ldamod)

## [1] "prior" "counts" "means" "scaling" "lev" "svd" "N"
## [8] "call" "terms" "xlevels"

# check the LDA coefficients/scalings
# a_1, ... a_p are given as "Coefficients of linear discriminants".
ldamod$scaling

##          LD1
## d8 -0.05161940
## d10 0.22969447
## d12 0.05015623
## d14 -0.59205628

crossprod(ldamod$scaling, Sp) %*% ldamod$scaling

## LD1
## LD1 1

# create the (centered) discriminant scores
mu.k=ldamod$means
mu=colMeans(mu.k)
dscores=scale(d2[,-1], center=mu, scale=F) %*% ldamod$scaling
sum((dscores - predict(ldamod)$x)^2)

## [1] 0

# plot the scores and coefficients
spid=as.integer(d2$Sex)
par(mfrow=c(1,2))
plot(dscores, xlab="LD1", ylab="LD2", pch=spid, col=spid,
      main="Discriminant Scores", xlim=c(-2, 30), ylim=c(-3, 3))
abline(h=0, lty=3)
abline(v=0, lty=3)
legend("bottomright", lev, pch=1:3, col=1:3, bty="n")

```

fit gender by teeth measure
we have to know prior probab
prior probabilities

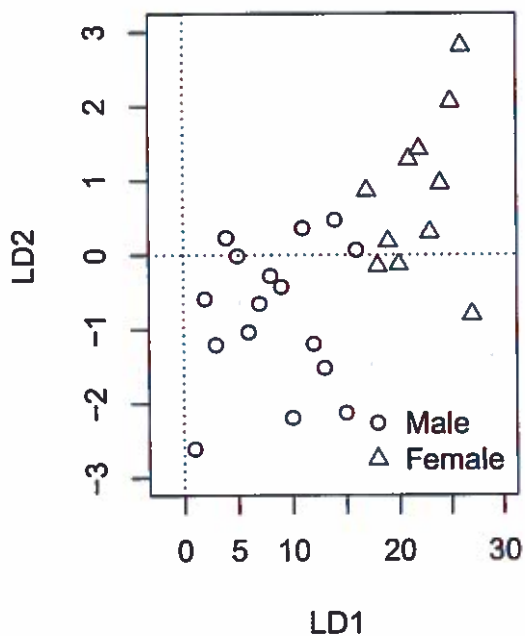
$Y = \sum a_i X_i$
↑
X mean e.

```

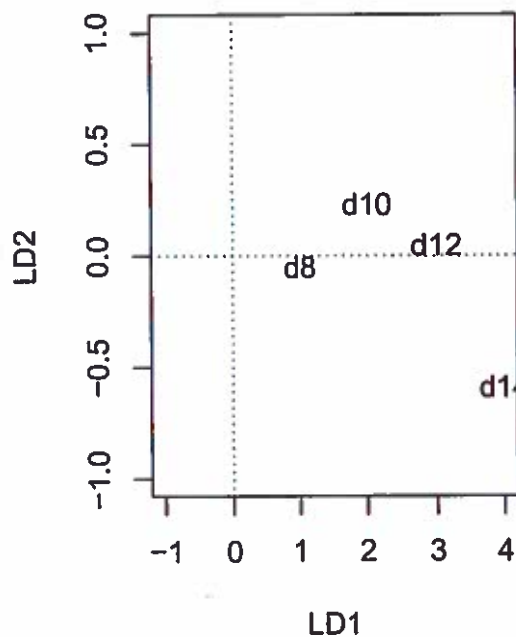
plot(ldamod$scaling, xlab="LD1", ylab="LD2", type="n",
     main="Discriminant Coefficients", xlim=c(-1, 4), ylim=c(-1, 1))
text(ldamod$scaling, labels=rownames(ldamod$scaling))
abline(h=0, lty=3)
abline(v=0, lty=3)

```

Discriminant Scores



Discriminant Coefficients

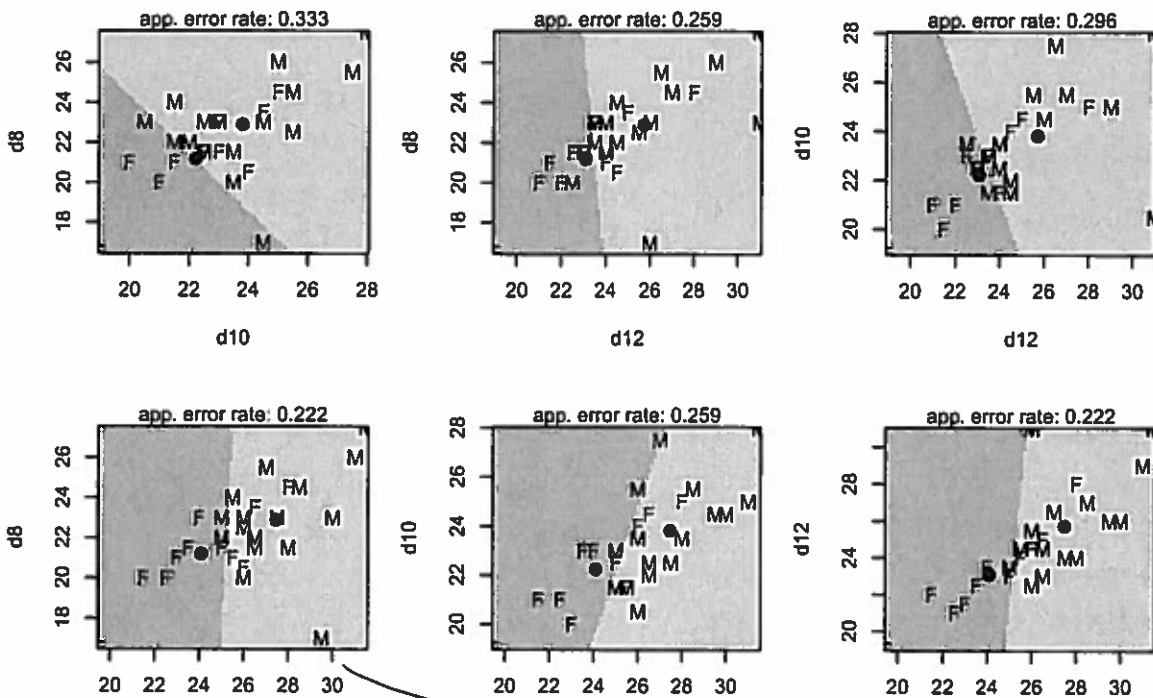


```

# visualize the LDA partitions
partimat(Sex ~ ., data=d2, method="lda")

```

Partition Plot $d(8, 10, 12, 14)$



First, have to check the assumption of $\Sigma_1 \neq \Sigma_2$ or $\Sigma_1 = \Sigma_2$ depends on miss copying then use only $d(8) \neq d(14)$.

3. Two Multivariate Normal Populations with $\Sigma_1 \neq \Sigma_2$

When $\Sigma_1 \neq \Sigma_2$, we can no longer use a simple linear classification rule. Under normality, terms involving $|\Sigma_i|^{1/2}$ do not cancel and the exponential terms do not combine easily. Using the definition of R_1 given earlier and expressing the likelihood ratio in the log scale, we now have:

$$R_1: \ln\left(\frac{p_1}{p_2}\right) > \ln\left(\frac{c(1|2)p_1}{c(2|1)p_2}\right) R_1: (\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1})x - k - \frac{1}{2}x^T(\Sigma_1^{-1} - \Sigma_2^{-1})x \geq \ln\left[\frac{c(1|2)}{c(2|1)}\right] \left(\frac{p_2}{p_1}\right)$$

$$R_2: (\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1})x - k - \frac{1}{2}x^T(\Sigma_1^{-1} - \Sigma_2^{-1})x \geq \ln\left[\frac{c(1|2)}{c(2|1)}\right] \left(\frac{p_2}{p_1}\right)$$

where $k = \frac{1}{2} \ln\left(\frac{|\Sigma_1|}{|\Sigma_2|}\right) + \frac{1}{2}(\mu_1^T \Sigma_1^{-1} \mu_1 - \mu_2^T \Sigma_2^{-1} \mu_2)$.

The classification regions are now quadratic functions of x . The classification rule for a new observation x_0 is now the following: allocate x_0 to π_1 if

$$R_1: (\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1})x_0 - k - \frac{1}{2}x_0^T(\Sigma_1^{-1} - \Sigma_2^{-1})x_0 \geq \ln\left[\frac{c(1|2)}{c(2|1)}\right] \left(\frac{p_2}{p_1}\right) \quad \checkmark$$

and allocate x_0 to π_2 otherwise. In practice, we estimate the classification rule by substituting the unobservable population parameters $(\mu_1, \mu_2, \Sigma_1, \Sigma_2)$ by sample estimates $(\bar{x}_1, \bar{x}_2, S_1, S_2)$.

Quadratic Classification Rule

Allocate x_0 to π_1 if

$$(\bar{x}_1^T S_1^{-1} - \bar{x}_2^T S_2^{-1})x_0 - k - \frac{1}{2}x_0^T(S_1^{-1} - S_2^{-1})x_0 \geq \ln\left[\frac{c(1|2)}{c(2|1)}\right] \left(\frac{p_2}{p_1}\right)$$

where $k = \frac{1}{2} \ln \left(\frac{|S_1|}{|S_2|} \right) + \frac{1}{2} (\bar{x}_1^T S_1^{-1} \bar{x}_1 - \bar{x}_2^T S_2^{-1} \bar{x}_2)$, and allocate x_0 to π_2 otherwise.

Revisit the wine data.

Check equal variance assumption

```
library(biotoools)
```

```
## Loading required package: rpanel
```

```
## Loading required package: tcltk
```

```
## Warning: running command '/usr/bin/otool' -L '/Library/Frameworks/
## R.framework/Resources/library/tcltk/libs//tcltk.so' had status 1
```

```
## Package 'rpanel', version 1.1-3: type help(rpanel) for summary information
```

```
## Loading required package: tkrplot
```

```
## Warning: package 'tkrplot' was built under R version 3.4.4
```

```
## Loading required package: lattice
```

```
## Loading required package: SpatialEpi
```

```
## Warning: package 'SpatialEpi' was built under R version 3.4.4
```

```
## Loading required package: sp
```

```
## ---
```

```
## biotoools version 3.1
```

```
##
```

```
boxM(d2[, -1], d2[, 1])
```

```
##
```

```
## Box's M-test for Homogeneity of Covariance Matrices
```

```
##
```

```
## data: d2[, -1]
```

```
## Chi-Sq (approx.) = 17.335, df = 10, p-value = 0.06727
```

At $\alpha = 10\%$, we reject $H_0: \Sigma_1 = \Sigma_2$. Now we fit a qda model

```
# fit qda model
```

```
qdamod=qda(Sex ~ ., data=d2, prior=rep(1/2, 2))
```

```
names(qdamod)
```

```
## [1] "prior" "counts" "means" "scaling" "ldet" "lev" "N"
```

```
## [8] "call" "terms" "xlevels"
```

```
# check the QDA coefficients/scalings
```

```
dim(qdamod$scaling)
```

```
## [1] 4 4 2
```

```
dnames <- dimnames(qdamod$scaling)
```

```
dnames
```

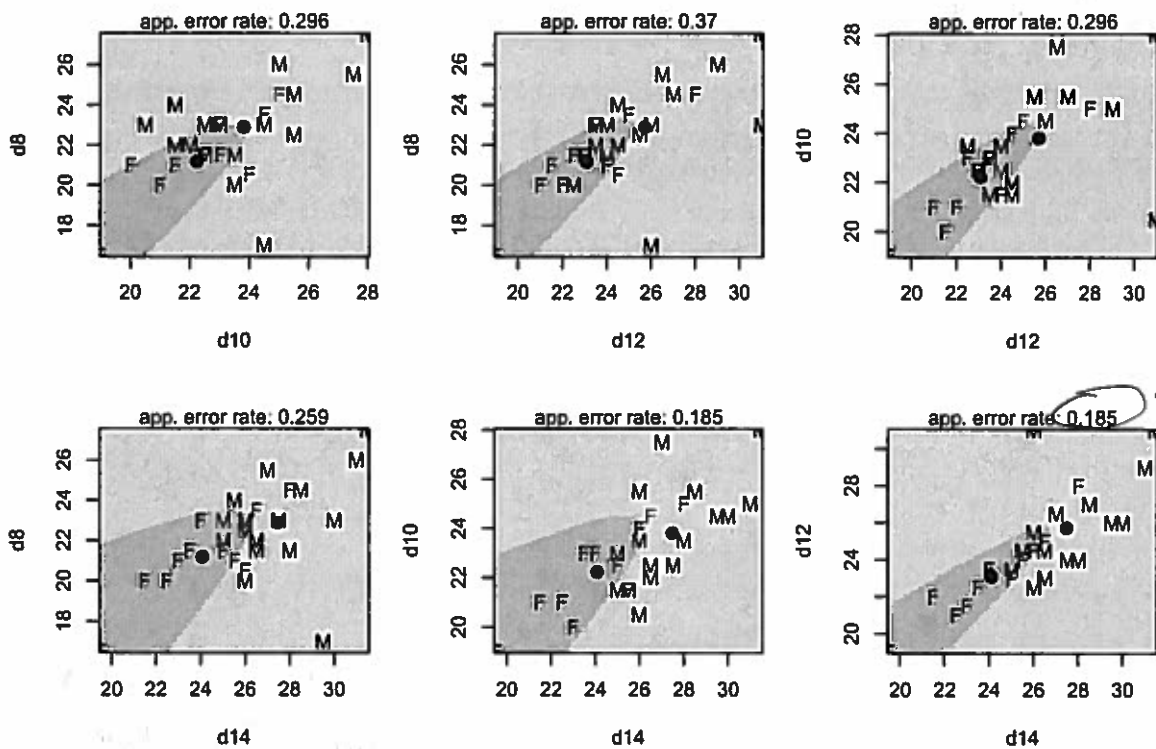
↑
coefficient.

? qda vs ldda ?

← give equation . (3 dimensional coefficients)

```
## [[1]]
## [1] "d8" "d10" "d12" "d14"
##
## [[2]]
## [1] "1" "2" "3" "4"
##
## [[3]]
## [1] "Male" "Female"
# visualize the QDA partitions
partimat(Sex ~ ., data=d2, method="qda")
```

Partition Plot



4. Evaluating Classification Functions

One important way of judging the performance of any classification procedure is to calculate its "error rate" or misclassification probabilities. In general, the population parameters are unknown in practice, so we focus on approaches that can estimate the error rates from the observed data.

The Total Probability of Misclassification (TPM) is defined as

$$TPM = p_1 \int_{R_2} f_1(x) dx + p_2 \int_{R_1} f_2(x) dx$$

want to minimize TPM

for any classification rule (region) that partitions $\Omega = R_1 \cup R_2$. The Optimum Error Rate (OER) is the minimum possible value of TPM

$$OER = \min_{R_1, R_2} TPM(R_1, R_2) \text{ subject to } \Omega = R_1 \cup R_2$$

which is obtained when

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1}$$

and

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1}$$

Note that if $c(2 | 1) = c(1 | 2)$, then minimizing TPM is same as minimizing ECM.

The error rates require knowledge of the (typically unknown) parameters that define the densities $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$. For example, in LDA, calculating OER requires μ_1, μ_2 and Σ . We replace these quantities in the allocation rules and evaluate the ~~actual~~ error rate (AER),

practically impossible to know f_1, f_2
theoretical possible

$$AER(\hat{R}_1, \hat{R}_2) = p_1 \int_{\hat{R}_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{\hat{R}_1} f_2(\mathbf{x}) d\mathbf{x}$$

where \hat{R}_1 and \hat{R}_2 are the classification regions determined by samples of size n_1 and n_2 respectively.

The AER indicates how the sample classification function will perform in future samples. Like the OER, it cannot be calculate because it depends on unknown densities f_1 and f_2 .

There is a measure of performance that does not depend on the form of populations and that can be calculated for any classification procedue. This measure is called the **apparent error rate (APER)** which is estimated using the observed (training) sample of data. The APER can be calculated from the **confussion matrix**.

Table 1: Confussion Matrix

		Classified as		
		π_1	π_2	
Truth	π_1	n_{C2} <i>correct</i>	n_{M1} <i>mis</i>	n_1 <i>Total</i>
	π_2	n_{M2} <i>mis</i>	n_{C2} <i>correct</i>	n_2

where - n_{1C} is the number of π_1 items correctly classified as π_1 items

- n_{1M} is the number of π_1 items misclassified as π_2 items
- n_{2C} is the number of π_2 items correctly classified
- n_{2M} is the number of π_2 items misclassified

Then

$$APER = \frac{n_{1M} + n_{2M}}{n_1 + n_2}$$

Expected actual error rate

$$E(AER) = \frac{n_{1M}^{(H)} + n_{2M}^{(H)}}{n_1 + n_2}$$

which is the **total proportion of misclassified sample observations**.

Leave-one-out Cross-Validation

Lachenbruch and Mickey proposed a better approach to estimate the AER:

1. Population 1 (for $i = 1, \dots, n_1$)
 - (a). Hold out the i^{th} observation from π_1 and build classification rule
 - (b). Use classification rule from Step 1(a) to classify the i^{th} observation
2. Population 2 (for $i = 1, \dots, n_2$)
 - (a). Hold out the i^{th} observation from π_2 and build classification rule
 - (b). Use classification rule from Step 2(a) to classify the i^{th} observation

Training data — fit the model — error of training data
 Testing data — — — testing data .

An (almost) unbiased estimate of the expected AER is given by

$$\hat{E}(AER) = \frac{n_{1M}^* + n_{2M}^*}{n_1 + n_2}$$

where n_{1M}^* and n_{2M}^* are the number of misclassified observations using the above "leave-one-out" procedure.

Revisit the wine data example: LDA

```
# make confusion matrix (and APER)
confusion=table(d2$Sex, predict(ldamod)$class)
confusion
```

```
##
##      Male Female
## Male      12      4
## Female     3      8

n=sum(confusion)
aper=(n - sum(diag(confusion))) / n
aper
```

APER

Use confusion to get APER
 confusionCV to get E(AER)

leave one out method

```
## [1] 0.2592593
```

```
# use CV to get expected AER
ldamodCV=lda(Sex ~ ., data=d2, prior=rep(1/2, 2), CV=TRUE)
confusionCV=table(d2$Sex, ldamodCV$class)
confusionCV
```

```
##
##      Male Female
## Male      10      6
## Female     4      7

eaer=(n - sum(diag(confusionCV))) / n
eaer
```

$$\frac{n_{1M} + n_{2M}}{n_1 + n_2} \quad E(AER)$$

when $c(2/1) = c(1/2)$

```
## [1] 0.3703704
```

```
# split into separate matrices for each flower
X1=subset(d2, Sex=="Male")
X2=subset(d2, Sex=="Female")
```

```
# split into training and testing
set.seed(1)
id1=sample.int(n=16, size=12)
id2=sample.int(n=11, size=8)
Xtrain=rbind(X1[id1,], X2[id2,])
Xtest=rbind(X1[-id1,], X2[-id2,])
```

training data .

```
# fit lda to training and evaluate on testing
ldatrain=lda(Sex ~ ., data=Xtrain, prior=rep(1/2, 2))
confusionTest=table(Xtest$Sex, predict(ldatrain, newdata=Xtest)$class)
confusionTest
```

```
##
##      Male Female
## Male      1      3
## Female     1      2
```

```

n=sum(confusionTest)
aer=(n - sum(diag(confusionTest))) / n
aer

## [1] 0.5714286 ← high error
# split into training and testing (100 splits)
nrep=100
aer=rep(0, nrep)
set.seed(1)
for(k in 1:nrep){
  #cat("rep:", k, "\n")
  id1=sample.int(n=16, size=12)
  id2=sample.int(n=11, size=8)
  Xtrain=rbind(X1[id1,], X2[id2,])
  Xtest=rbind(X1[-id1,], X2[-id2,])
  ldatrain=lda(Sex ~ ., data=Xtrain, prior=rep(1/2, 2))
  confusionTest=table(Xtest$Sex, predict(ldatrain, newdata=Xtest)$class)
  confusionTest
  n=sum(confusionTest)
  aer[k]=(n - sum(diag(confusionTest))) / n
}
mean(aer)

```

[1] 0.3614286

Using QDA

```

# make confusion matrix (and APER)
confusion=table(d2$Sex, predict(qdamod)$class)
confusion

```

```

##
##      Male Female
## Male      12      4
## Female     1     10

```

```

n=sum(confusion)
aper=(n - sum(diag(confusion))) / n
aper

```

[1] 0.1851852

```

# use CV to get expected AER
qdamodCV=qda(Sex ~ ., data=d2, prior=rep(1/2, 2), CV=TRUE)
confusionCV=table(d2$Sex, qdamodCV$class)
confusionCV

```

```

##
##      Male Female
## Male      9      7
## Female     3      8

```

```

eaer = (n - sum(diag(confusionCV))) / n
eaer

```

[1] 0.3703704

Using QDA
**

E(AER)

```

# split into training and testing
set.seed(1)
id1=sample.int(n=16, size=12)
id2=sample.int(n=11, size=8)
Xtrain=rbind(X1[id1,], X2[id2,])
Xtest=rbind(X1[-id1,], X2[-id2,])

# fit qda to training and evaluate on testing
qdatrain=qda(Sex ~ ., data=Xtrain, prior=rep(1/2, 2))
confusionTest= table(Xtest$Sex, predict(qdatrain, newdata=Xtest)$class)
confusionTest

##
##           Male Female
## Male         4      0
## Female        2      1
n=sum(confusionTest)
aer=(n - sum(diag(confusionTest))) / n
aer

## [1] 0.2857143
# split into training and testing (100 splits)
nrep=100
aer=rep(0, nrep)
set.seed(1)
for(k in 1:nrep){
  #cat("rep:", k, "\n")
  id1=sample.int(n=16, size=12)
  id2=sample.int(n=11, size=8)
  Xtrain=rbind(X1[id1,], X2[id2,])
  Xtest=rbind(X1[-id1,], X2[-id2,])
  qdatrain=qda(Sex ~ ., data=Xtrain, prior=rep(1/2, 2))
  confusionTest=table(Xtest$Sex, predict(qdatrain, newdata=Xtest)$class)
  confusionTest
  n=sum(confusionTest)
  aer[k]=(n - sum(diag(confusionTest))) / n
}
mean(aer)

## [1] 0.3971429

```

Example: Classifying gender using dental data

```

library(reshape) # to create data frame from dental data
library(MASS) # has lda and qda functions
library(heavy) # has dental data
data(dental)
d2=cast(melt(dental, id=c("Subject", "age", "Sex")), Subject+Sex+age)
names(d2)[3:6]=c("d8", "d10", "d12", "d14")
# these functions estimate pi_i using sample proportions
# you can provide other pi_i if needed
f1=lda(x=as.matrix(d2[,3:6]), grouping=d2[,2], CV=T)

```

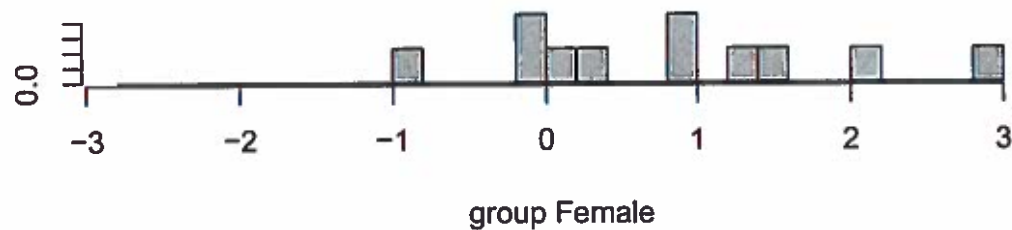
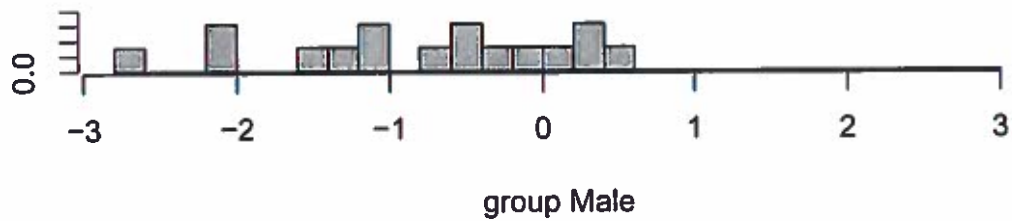
```

f1=lda(Sex~d8+d10+d12+d14,data=d2,CV=T)
f2=qda(Sex~d8+d10+d12+d14,data=d2,CV=T)
sum(f1$class!=d2$Sex)/length(d2$Sex) # CV error linear

## [1] 0.2962963
sum(f2$class!=d2$Sex)/length(d2$Sex) # CV error quadratic

## [1] 0.3333333
f1=lda(Sex~d8+d10+d12+d14,data=d2) # refit without CV
f2=qda(Sex~d8+d10+d12+d14,data=d2) # refit without CV
plot(f1) # uses discriminant functions

```

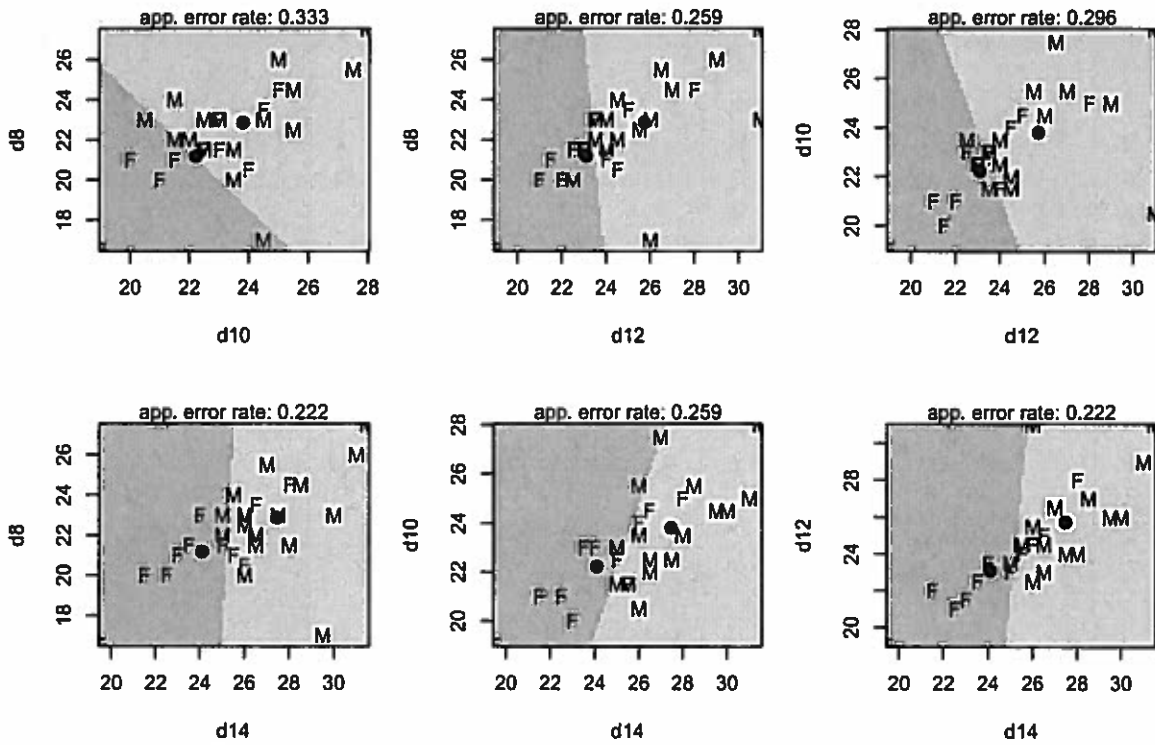


```

library(klaR) # provides panel of bivariate plots w/ regions
partimat(Sex~d8+d10+d12+d14,data=d2,method="lda")

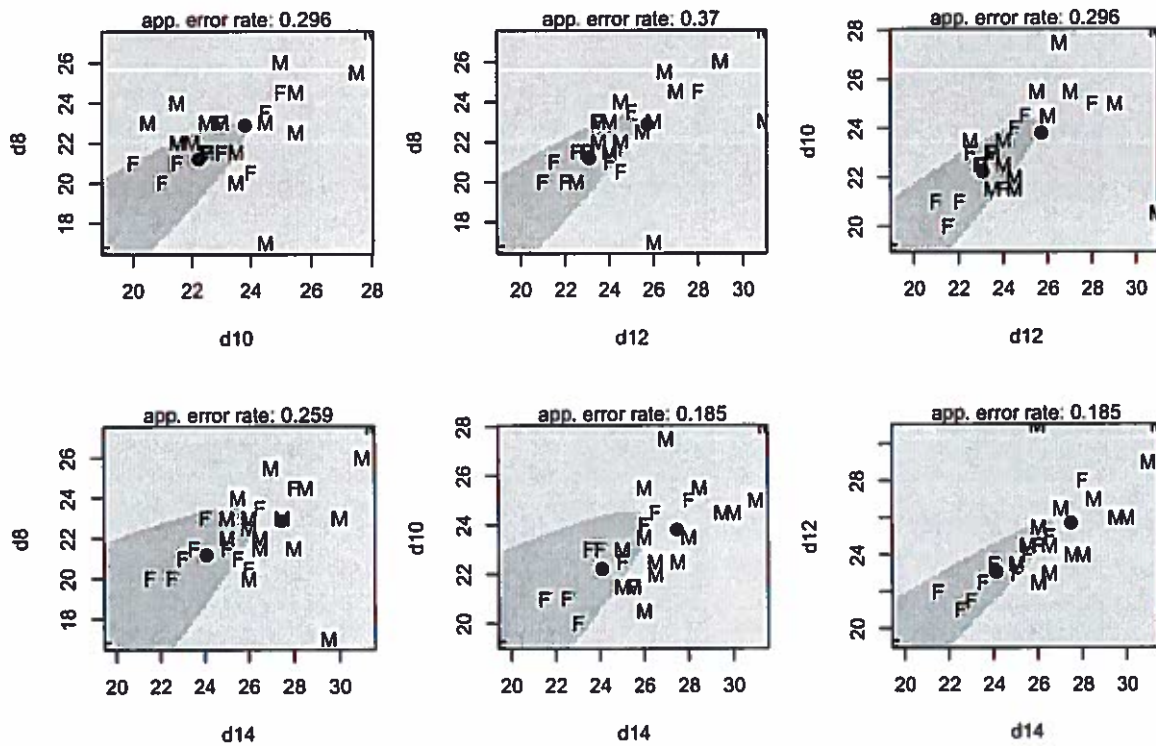
```

Partition Plot



`partimat(Sex-d8+d10+d12+d14,data=d2,method="qda")`

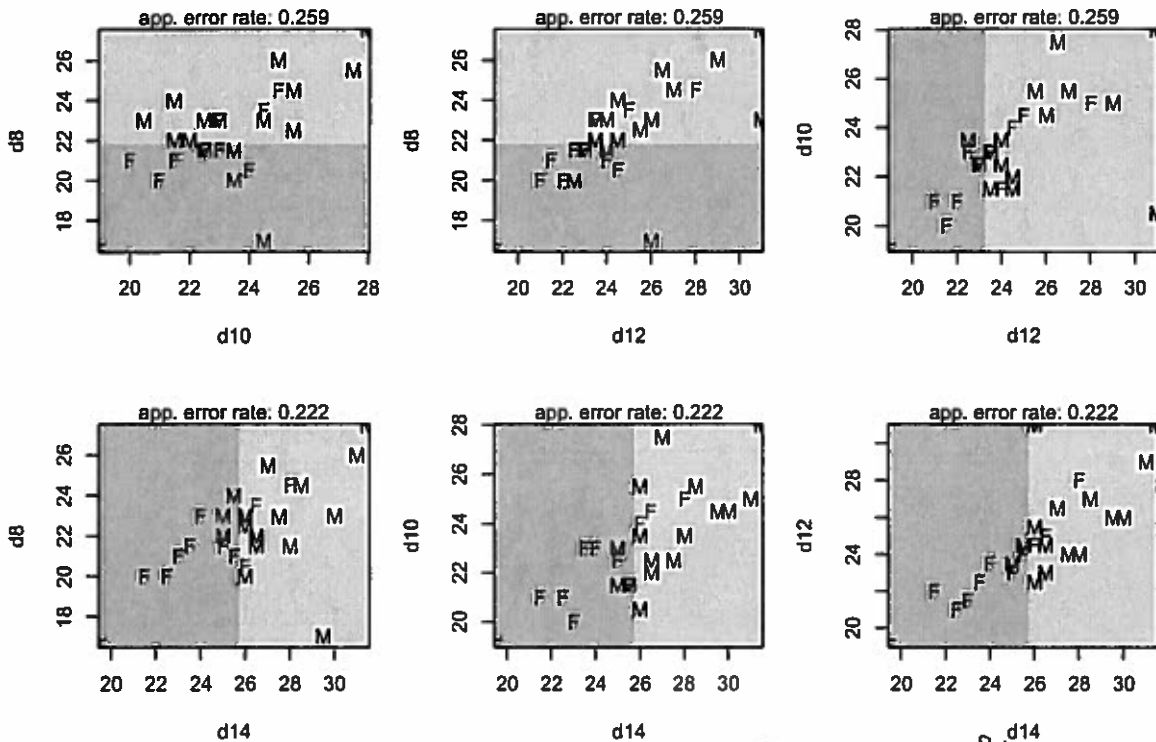
Partition Plot



```
partimat(Sex~d8+d10+d12+d14,data=d2,method="rpart") # classification tree
```

```
## Loading required namespace: rpart
```

Partition Plot



5. Classification with Several Populations

Let $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ denote a p -dimensional random vector. Let $f_k(\mathbf{x})$ denote the probability density function (pdf) for population π_k for $k = 1, 2, \dots, g$. Given a realization $\mathbf{X} = \mathbf{x}$, we want to assign \mathbf{x} to π_1, π_2 or π_g . We want to find some classification rule to determine whether a realization $\mathbf{X} = \mathbf{x}$ should be assigned to population π_1 or π_2 , or π_g .

Let Ω denote the sample space (collection of all possible values for \mathbf{x}). $R_1 \subset \Omega$ is the set of values of \mathbf{x} for which we classify objects into π_1 , $R_2 \subset \Omega$ is the set of values of \mathbf{x} for which we classify objects into π_2 , $\dots, R_g \subset \Omega$ is the set of values of \mathbf{x} for which we classify objects into π_g . The classification rule partitions the sample space and the classification regions are mutually exclusive, i.e., $\Omega = R_1 \cup R_2 \cup \dots \cup R_g$ and $R_k \cap R_l = \emptyset$ for all $k \neq l$.

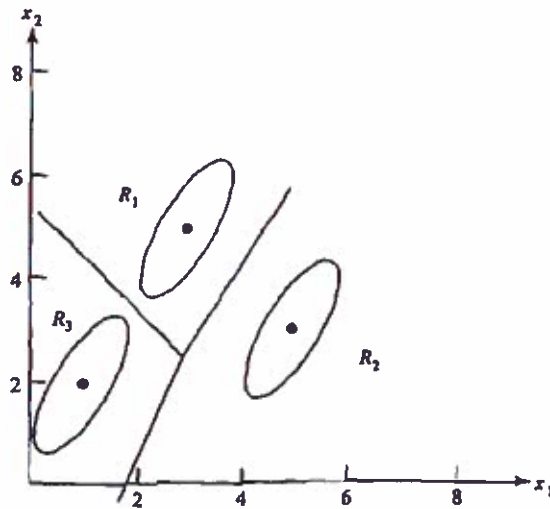


Figure 11.10 The classification regions R_1 , R_2 , and R_3 for the linear minimum TPM rule ($p_1 = \frac{1}{4}$, $p_2 = \frac{1}{2}$, $p_3 = \frac{1}{4}$).

The conditional probability $P(l | k)$ of classifying an object as π_l when the object really belongs to π_k is given by

$$P(l | k) = P(\mathbf{X} \in R_l | \pi_k) = \int_{R_l} f_k(\mathbf{x}) d\mathbf{x}$$

for all $k \neq l$ with $k, l = 1, 2, \dots, g$. Note that $P(k | k) = 1 - \sum_{l \neq k} P(l | k)$ by definition. Let $c(l | k)$ denote the cost of allocating an object to π_l when the object really belongs to π_k , and let p_k denote the prior probability of π_k . Then the conditional expected cost of misclassifying an object from π_k is

$$ECM = \sum_{l \neq k} P(l | k) c(l | k)$$

Incorporating the prior probabilities, the overall ECM is given by

$$ECM = \sum_{k=1}^g p_k ECM(k) = \sum_{k=1}^g p_k \left[\sum_{l \neq k} P(l | k) c(l | k) \right]$$

The classification regions $\{R_1, R_2, \dots, R_g\}$ that minimize the ECM are defined by allocating $\mathbf{X} = \mathbf{x}$ to the population π_k that minimizes

$$ECM = \sum_{l \neq k} p_l f_l(\mathbf{x}) c(k | l)$$

Fisher's method for Discriminating among Several Populations

Fisher also developed his discriminant analysis for $g > 2$ populations. The motivation behind the Fisher discriminant analysis is the need to obtain a reasonable representation of the populations that involves only a few linear combination of the observations, such as $\mathbf{a}_1^T \mathbf{x}$, $\mathbf{a}_2^T \mathbf{x}$, $\mathbf{a}_3^T \mathbf{x}$. His approach offers a simple and useful procedure for classification, which also provides nice visualizations. We will be able to plot the linear combinations to visualize the discriminants.

Let $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ denote a p -dimensional random vector and let $f_k(\mathbf{x}) \sim (\mu_k, \Sigma)$ denote the pdf for population π_k . Here we do not need the multivariate normality assumption but assume homogeneity of covariance matrix assumption. Let $\bar{\mu} = \frac{1}{g} \sum_{k=1}^g \mu_k$ denote the mean of the combined populations, and

$$B_{\mu} = \sum_{k=1}^g (\mu_k - \bar{\mu})(\mu_k - \bar{\mu})^T$$

denote "Between" group sums of cross products matrix. We consider the linear combination $Y = \mathbf{a}^T \mathbf{X}$ which has expected value

$$E(Y | \pi_k) = \mathbf{a}^T E(\mathbf{X} | \pi_k) = \mathbf{a}^T \boldsymbol{\mu}_k$$

and variance

$$\text{Var}(Y | \pi_k) = \mathbf{a}^T \text{Cov}(\mathbf{X} | \pi_k) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$$

We define the overall mean

$$\bar{\mu}_Y = \frac{1}{g} \sum_{k=1}^g \mu_{Y_k} = \frac{1}{g} \sum_{k=1}^g \mathbf{a}^T \boldsymbol{\mu}_k = \mathbf{a}^T \bar{\boldsymbol{\mu}}$$

and form the ratio of the between group separation over the variance of Y

$$\begin{aligned} F^* &= \frac{\sum_{k=1}^g (\mu_{Y_k} - \bar{\mu}_Y)^2}{\sigma_Y^2} \\ &= \frac{\sum_{k=1}^g (\mathbf{a}^T \boldsymbol{\mu}_k - \mathbf{a}^T \bar{\boldsymbol{\mu}})^2}{\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}} \\ &= \frac{\mathbf{a}^T \sum_{k=1}^g (\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}})^T \mathbf{a}}{\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}} \\ &= \frac{\mathbf{a}^T \mathbf{B}_\mu \mathbf{a}}{\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}} \end{aligned}$$

Note that higher F^* values relate to more separation between groups.

In general, $\boldsymbol{\Sigma}$ and the $\boldsymbol{\mu}_i$ are unavailable. We use the sample data to estimate the "Between" and "Within" group sums of cross products matrices.

$$\mathbf{B} = \sum_{k=1}^g (\bar{\mathbf{x}}_k - \bar{\mathbf{x}})(\bar{\mathbf{x}}_k - \bar{\mathbf{x}})^T$$

$$\mathbf{W} = \sum_{k=1}^g (n_k - 1) \mathbf{S}_k = \sum_{k=1}^g \sum_{i=1}^{n_k} (\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)(\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)^T$$

where the sample mean vectors $\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{x}_{ki}$ and the "overall mean" vector $\bar{\mathbf{x}} = \frac{1}{g} \sum_{k=1}^g \bar{\mathbf{x}}_k$.

Consequently, $\frac{\mathbf{W}}{n_1 + \dots + n_g - g} = \mathbf{S}_{pooled}$ is the estimate of $\boldsymbol{\Sigma}$.

Then the sample k^{th} discriminant is the linear combination

$$\hat{Y}_k = \hat{\mathbf{a}}_k^T \mathbf{X}$$

where $\hat{\mathbf{a}}_k$ are scaled to make the \hat{Y}_k have unit variance, i.e., $\hat{\mathbf{a}}_k^T \hat{\boldsymbol{\Sigma}} \hat{\mathbf{a}}_k = 1$ where $\hat{\boldsymbol{\Sigma}} = \mathbf{S}_p = \frac{1}{n-g} \mathbf{W}$ with $n = \sum_{k=1}^g n_k$.

Properties of Discriminants

Let $\mathbf{Y} = [Y_1, Y_2, \dots, Y_g]^T = \mathbf{A}^T \mathbf{X}$, where $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_g]$, then \mathbf{Y} contains the g discriminants, that is, $\mathbf{Y} = [Y_1, Y_2, \dots, Y_g]^T$. Columns of \mathbf{A} contain the linear combination weights. \mathbf{Y} has the following properties:

- mean of \mathbf{Y} : $E(\mathbf{Y} | \pi_k) = \mathbf{A}^T E(\mathbf{X} | \pi_k) = \mathbf{A}^T \boldsymbol{\mu}_k = \boldsymbol{\mu}_{kY}$
- covariance matrix of \mathbf{Y} : $\text{Cov}(\mathbf{Y}) = \mathbf{A}^T \text{Cov}(\mathbf{X} | \pi_k) \mathbf{A} = \mathbf{A}^T \boldsymbol{\Sigma} \mathbf{A} = \mathbf{I}_g$, the discriminants have unit variance and are uncorrelated.

6. Using Fisher's Discriminants to Classify Objects

Fisher's Discriminants were derived for the purpose of obtaining a low dimensional representation of the data that separates the populations as much as possible. The discriminants provide the basis for a classification rule. Let $Y_k = \mathbf{a}_k^T \mathbf{X}$ be the k^{th} discriminant, $k \leq s$, then the mean of \mathbf{Y} is

$$\boldsymbol{\mu}_{iY} = [\mu_{iY_1}, \mu_{iY_2}, \dots, \mu_{iY_s}]^T = [\mathbf{a}_1^T \boldsymbol{\mu}_i, \mathbf{a}_2^T \boldsymbol{\mu}_i, \dots, \mathbf{a}_s^T \boldsymbol{\mu}_i]^T$$

under population π_i and the covariance matrix is \mathbf{I} for all populations. Because the components of \mathbf{Y} have unit variances and zero covariances, it is appropriate to measure the distance using

$$D_i = (\mathbf{y} - \boldsymbol{\mu}_{iY})^T (\mathbf{y} - \boldsymbol{\mu}_{iY}) = \sum_{j=1}^s (y_j - \mu_{iY_j})^2 = \sum_{j=1}^s [\mathbf{a}_j^T (\mathbf{x} - \boldsymbol{\mu}_i)]^2$$

where $\boldsymbol{\mu}_{iY} = \mathbf{A}^T \boldsymbol{\mu}_i$, $y_j = \mathbf{a}_j^T \mathbf{x}$ and $\mu_{iY_j} = \mathbf{a}_j^T \boldsymbol{\mu}_i$.

We allocate \mathbf{x} to population π_k if

$$\sum_{j=1}^r (y_j - \mu_{kY_j})^2 = \sum_{j=1}^r [\mathbf{a}_j^T (\mathbf{x} - \boldsymbol{\mu}_k)]^2 \leq \sum_{j=1}^r [\mathbf{a}_j^T (\mathbf{x} - \boldsymbol{\mu}_i)]^2$$

for all $i \neq k$, $r \leq s$.

Example: Wine Recognition data

```
## Wine Recognition data in 'rattle'
library(rattle)

## Warning: package 'rattle' was built under R version 3.4.4

## Rattle: A free graphical interface for data science with R.
## Version 5.2.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.

library(MASS)
data(wine)
attach(wine)
head(wine)
```

##	Type	Alcohol	Malic	Ash	Alcalinity	Magnesium	Phenols	Flavanoids
## 1	1	14.23	1.71	2.43	15.6	127	2.80	3.06
## 2	1	13.20	1.78	2.14	11.2	100	2.65	2.76
## 3	1	13.16	2.36	2.67	18.6	101	2.80	3.24
## 4	1	14.37	1.95	2.50	16.8	113	3.85	3.49
## 5	1	13.24	2.59	2.87	21.0	118	2.80	2.69
## 6	1	14.20	1.76	2.45	15.2	112	3.27	3.39

##	Nonflavanoids	Proanthocyanins	Color	Hue	Dilution	Proline
## 1	0.28	2.29	5.64	1.04	3.92	1065
## 2	0.26	1.28	4.38	1.05	3.40	1050
## 3	0.30	2.81	5.68	1.03	3.17	1185
## 4	0.24	2.18	7.80	0.86	3.45	1480
## 5	0.39	1.82	4.32	1.04	2.93	735
## 6	0.34	1.97	6.75	1.05	2.85	1450

table(Type)

function table in R

Table(one column) → give # counts of elements.

```
## Type
## 1 2 3
## 59 71 48
g <- 3 # of groups
p <- ncol(wine) - 1 # dimension
# pooled covariances matrix
Sp <- matrix(0, p, p)
nx <- rep(0, g)
lev <- levels(wine$Type)
for(k in 1:g){
  x <- wine[wine$Type==lev[k], 1:p+1]
  nx[k] <- nrow(x)
  Sp <- Sp + cov(x) * (nx[k] - 1)
}
Sp <- Sp / (sum(nx) - g)
round(Sp, 3)
```

##	Alcohol	Malic	Ash	Alcalinity	Magnesium	Phenols
## Alcohol	0.262	0.008	-0.013	-0.097	0.018	0.027
## Malic	0.008	0.888	0.021	0.415	-0.889	-0.014
## Ash	-0.013	0.021	0.066	0.484	0.690	0.016
## Alcalinity	-0.097	0.415	0.484	8.007	3.238	0.109
## Magnesium	0.018	-0.889	0.690	3.238	180.658	0.577
## Phenols	0.027	-0.014	0.016	0.109	0.577	0.191
## Flavanoids	0.023	-0.009	0.030	0.246	0.667	0.161
## Nonflavanoids	-0.001	0.010	0.007	0.047	-0.287	-0.008
## Proanthocyanins	0.017	0.017	0.003	0.093	1.300	0.092
## Color	0.246	-0.259	0.010	-0.102	1.772	0.199
## Hue	0.001	-0.042	0.002	-0.010	0.126	-0.001
## Dilution	-0.005	0.047	0.011	0.226	-0.281	0.060
## Proline	12.237	-33.058	-0.501	-32.833	476.249	8.395
##	Flavanoids	Nonflavanoids	Proanthocyanins	Color	Hue	
## Alcohol	0.023	-0.001	0.017	0.246	0.001	
## Malic	-0.009	0.010	0.017	-0.259	-0.042	
## Ash	0.030	0.007	0.003	0.010	0.002	
## Alcalinity	0.246	0.047	0.093	-0.102	-0.010	
## Magnesium	0.667	-0.287	1.300	1.772	0.126	
## Phenols	0.161	-0.008	0.092	0.199	-0.001	
## Flavanoids	0.275	-0.015	0.128	0.287	-0.004	
## Nonflavanoids	-0.015	0.012	-0.009	-0.002	0.001	
## Proanthocyanins	0.128	-0.009	0.246	0.229	-0.006	
## Color	0.287	-0.002	0.229	2.285	-0.041	
## Hue	-0.004	0.001	-0.006	-0.041	0.024	
## Dilution	0.068	-0.010	0.042	-0.066	-0.003	
## Proline	3.426	-0.493	11.482	68.094	4.491	
##	Dilution	Proline				
## Alcohol	-0.005	12.237				
## Malic	0.047	-33.058				
## Ash	0.011	-0.501				
## Alcalinity	0.226	-32.833				
## Magnesium	-0.281	476.249				

```

## Phenols          0.060      8.395
## Flavonoids       0.068      3.426
## Nonflavanoids    -0.010     -0.493
## Proanthocyanins  0.042     11.482
## Color            -0.066     68.094
## Hue              -0.003      4.491
## Dilution        0.161    -10.943
## Proline          -10.943  29707.682

# fit lda model
# assume equal prior
ldamod <- lda(Type ~ ., data=wine, prior=rep(1/3, 3))
names(ldamod)

## [1] "prior" "counts" "means" "scaling" "lev" "svd" "N"
## [8] "call" "terms" "xlevels"

# check the LDA coefficients/scalings

# a_1, ... a_p are given as "Coefficients of linear discriminants".
ldamod$scaling

##                LD1          LD2
## Alcohol        -0.356369042  0.892051370
## Malic           0.181293364  0.296139458
## Ash             -0.243541326  2.362184415
## Alkalinity      0.146777356 -0.154421898
## Magnesium       -0.002185082 -0.000346807
## Phenols         0.615457266 -0.065102824
## Flavonoids      -1.685049247 -0.402774674
## Nonflavanoids   -1.580605983 -1.548924728
## Proanthocyanins 0.117537507 -0.313796985
## Color           0.368045785  0.233950076
## Hue             -0.897641355 -1.469888074
## Dilution       -1.153187021  0.112797172
## Proline         -0.002535348  0.002992344

crossprod(ldamod$scaling, Sp) %*% ldamod$scaling

##                LD1          LD2
## LD1  1.000000e+00 -9.436896e-16
## LD2 -1.082467e-15  1.000000e+00
# unique variance property

# create the (centered) discriminant scores
mu.k <- ldamod$means
mu <- colMeans(mu.k)
dscores <- scale(wine[, -1], center=mu, scale=F) %*% ldamod$scaling
sum((dscores - predict(ldamod)$x)^2)

## [1] 2.746106e-28

# plot the scores and coefficients
spid <- as.integer(wine$Type)
par(mfrow=c(1,2))
plot(dscores, xlab="LD1", ylab="LD2", pch=spid, col=spid,
      main="Discriminant Scores", xlim=c(-10, 10), ylim=c(-3, 3))
abline(h=0, lty=3)
abline(v=0, lty=3)

```

discriminant
score

unique variance property

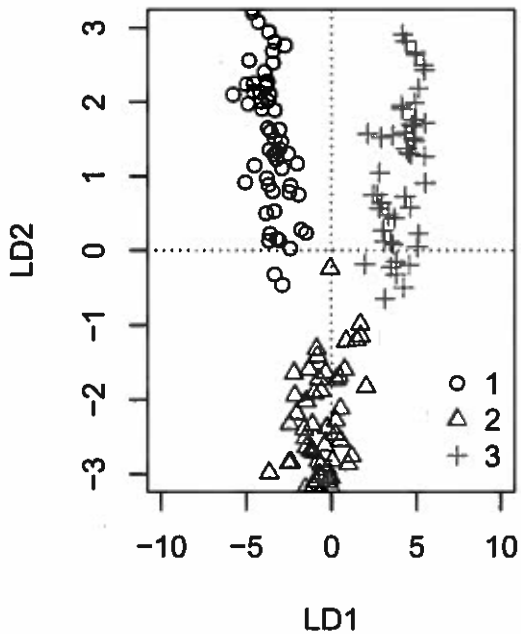
? DE
↑ not Fish.

```

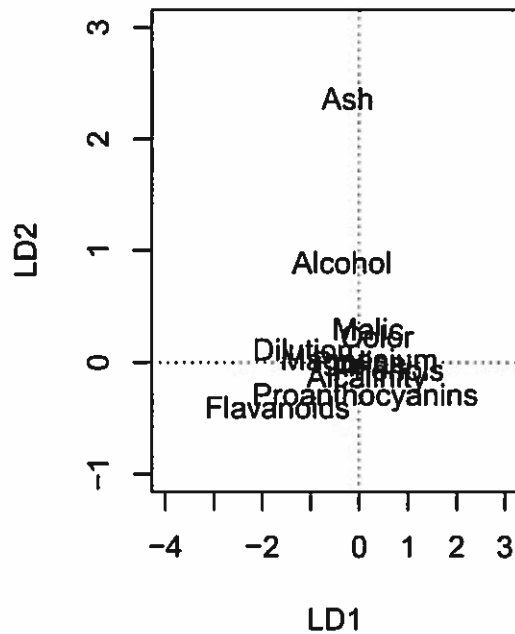
legend("bottomright",lev,pch=1:3,col=1:3,bty="n")
plot(ldamod$scaling, xlab="LD1", ylab="LD2", type="n",
      main="Discriminant Coefficients", xlim=c(-4, 3), ylim=c(-1, 3))
text(ldamod$scaling, labels=rownames(ldamod$scaling))
abline(h=0, lty=3)
abline(v=0, lty=3)

```

Discriminant Scores



Discriminant Coefficients

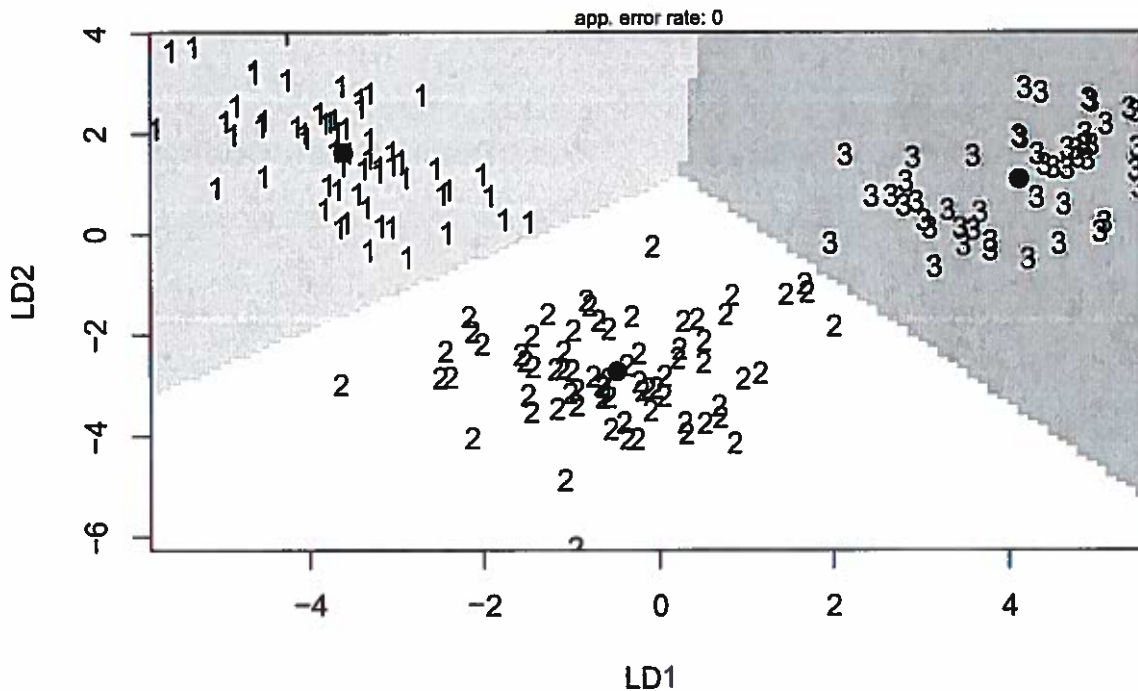


```

# visualize the LDA partitions
type <- factor(Type)
library(klaR)
partimat(x=dscores[,2:1], grouping=type, method="lda")

```


Partition Plot



```
# visualize the LDA partitions (for all pairs)
#partimat(x=wine[,-1], grouping=type, method="lda")
```

```
##### QDA #####
# Check equal variance assumption
library(biotools)
boxM(wine[,-1],wine[,1])
```

```
##
## Box's M-test for Homogeneity of Covariance Matrices
##
## data: wine[, -1]
## Chi-Sq (approx.) = 684.2, df = 182, p-value < 2.2e-16
# fit qda model
qdamod <- qda(Type ~ ., data=wine, prior=rep(1/3, 3))
names(qdamod)
```

```
## [1] "prior" "counts" "means" "scaling" "ldet" "lev" "N"
## [8] "call" "terms" "xlevels"
```

```
# check the QDA coefficients/scalings
dim(qdamod$scaling)
```

```
## [1] 13 13 3
```

```
dnames <- dimnames(qdamod$scaling)
dnames
```

```
## [[1]]
## [1] "Alcohol" "Malic" "Ash"
```

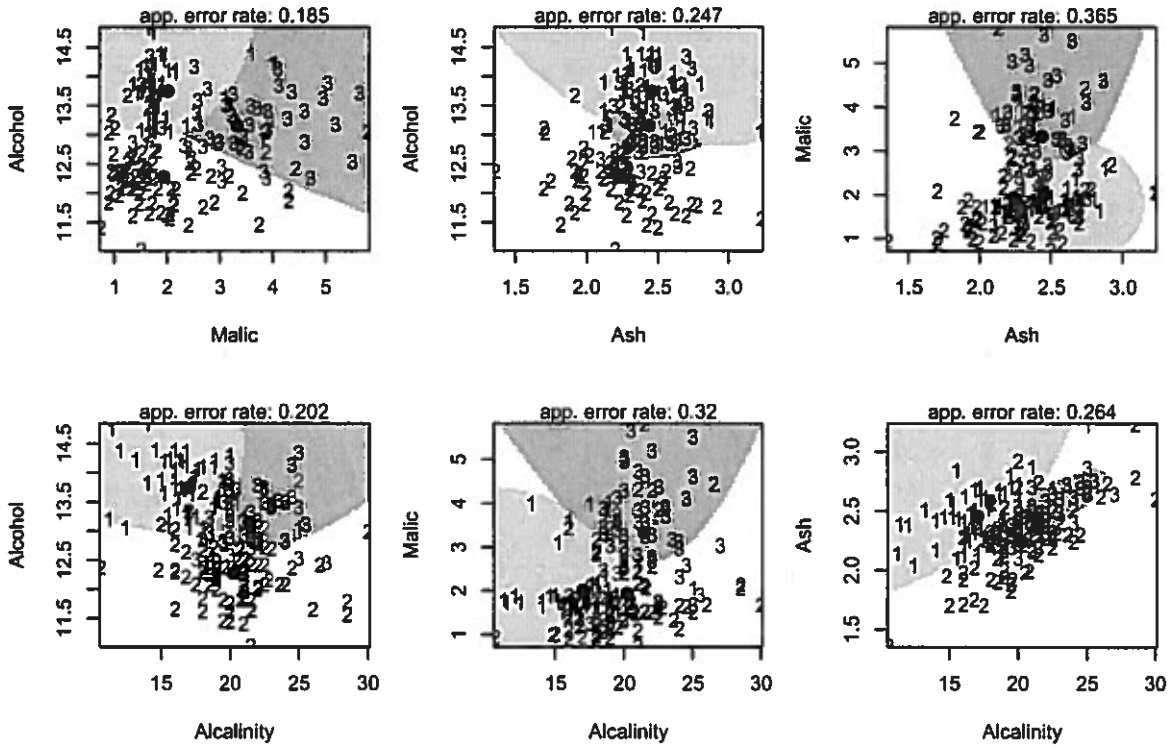
ed equal
variance
assumption

→ the 3 type of wines ≠ not equal
at least 2 of them have ≠ variance
prior equally likely belong to 1/2/3.

```
## [4] "Alcalinity"      "Magnesium"      "Phenols"
## [7] "Flavanoids"      "Nonflavanoids" "Proanthocyanins"
## [10] "Color"           "Hue"            "Dilution"
## [13] "Proline"
##
## [[2]]
## [1] "1" "2" "3" "4" "5" "6" "7" "8" "9" "10" "11" "12" "13"
##
## [[3]]
## [1] "1" "2" "3"
# visualize the QDA partitions
#partimat(Type ~ ., data=wine, method="qda") # all pairs
partimat(Type ~ ., data=wine[,2:5], method="qda")
```

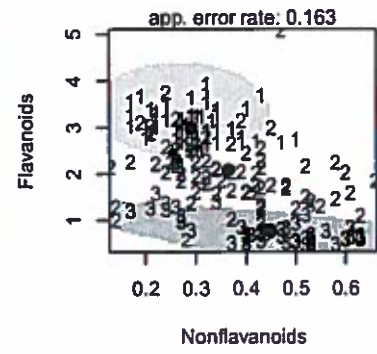
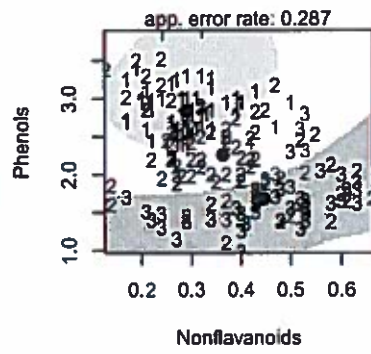
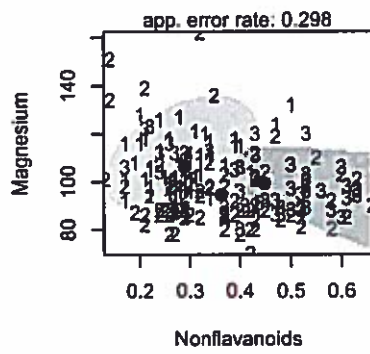
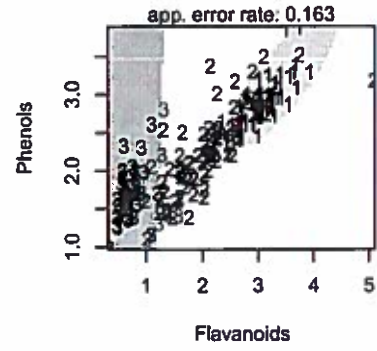
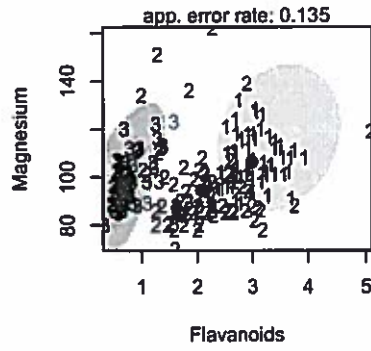
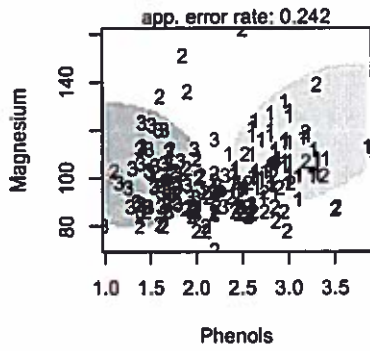
Partition Plot

if we only use 2 variables to classify



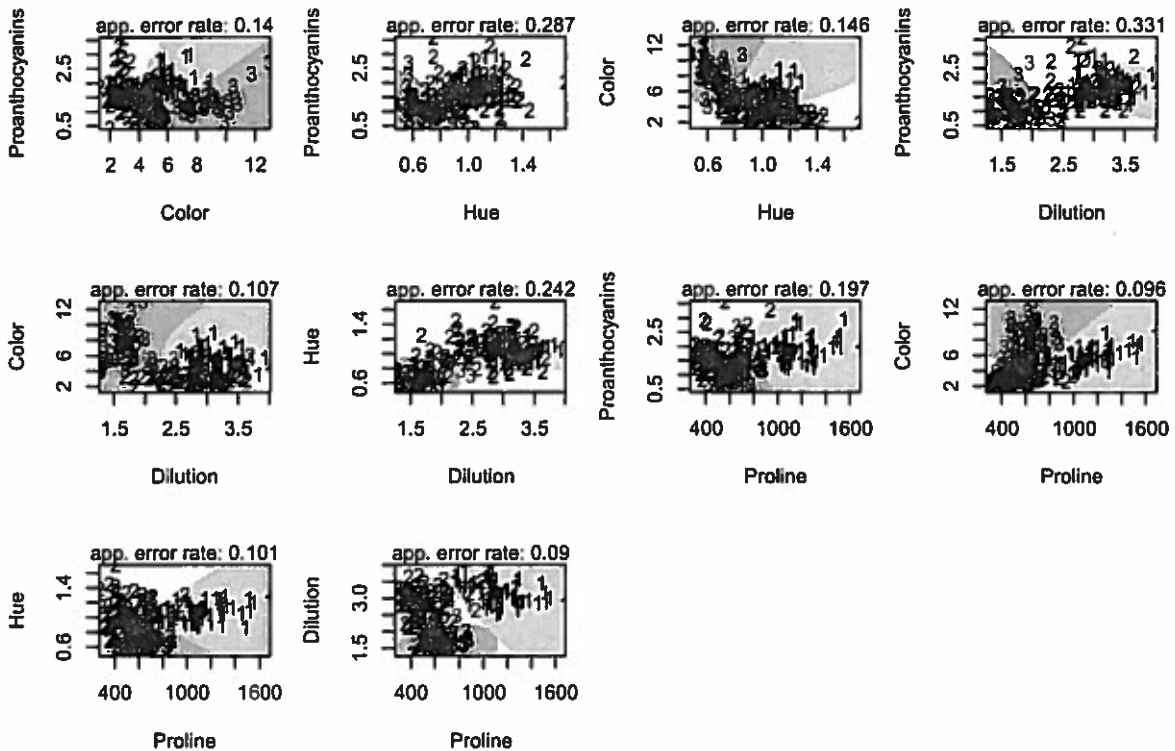
```
partimat(Type ~ ., data=wine[,6:9], method="qda")
```

Partition Plot



```
partimat(Type ~ ., data=wine[,10:14], method="qda")
```

Partition Plot



```

##### CV LDA #####
# make confusion matrix (and APER)
confusion <- table(wine$Type, predict(ldamod)$class)
confusion

##
##      1  2  3
## 1 59  0  0
## 2  0 71  0
## 3  0  0 48

n <- sum(confusion)
aper <- (n - sum(diag(confusion))) / n
aper

## [1] 0
# use CV to get expected AER
ldamodCV <- lda(Type ~ ., data=wine, prior=rep(1/3, 3), CV=TRUE)
confusionCV <- table(wine$Type, ldamodCV$class)
confusionCV

##
##      1  2  3
## 1 59  0  0
## 2  1 69  1
## 3  0  0 48

eaer <- (n - sum(diag(confusionCV))) / n
eaer
    
```

direct

leave one out method

```

## [1] 0.01123596
# OR
f1 <- lda(Type~., data=wine, CV=T)
sum(f1$class!=wine$Type)/length(wine$Type) # CV error linear

## [1] 0.01123596
# split into separate matrices for each flower
X1 <- subset(wine, Type=="1")
X2 <- subset(wine, Type=="2")
X3 <- subset(wine, Type=="3")

# split into training and testing
set.seed(1)
id1 <- sample.int(n=59, size=41)
id2 <- sample.int(n=71, size=50)
id3 <- sample.int(n=48, size=34)
Xtrain <- rbind(X1[id1,], X2[id2,], X3[id3,])
Xtest <- rbind(X1[-id1,], X2[-id2,], X3[-id3,])

# fit lda to training and evaluate on testing
ldatrain <- lda(Type ~ ., data=Xtrain, prior=rep(1/3, 3))
confusionTest <- table(Xtest$Type, predict(ldatrain, newdata=Xtest)$class)
confusionTest

##
##      1  2  3
##  1 18  0  0
##  2  0 21  0
##  3  0  0 14

n <- sum(confusionTest)
aer <- (n - sum(diag(confusionTest))) / n
aer

## [1] 0 ← testing error rate ← good job! lda

# split into training and testing (100 splits)
nrep <- 100
aer <- rep(0, nrep)
set.seed(1)
for(k in 1:nrep){
  #cat("rep:", k, "\n")
  id1 <- sample.int(n=59, size=41)
  id2 <- sample.int(n=71, size=50)
  id3 <- sample.int(n=48, size=34)
  Xtrain <- rbind(X1[id1,], X2[id2,], X3[id3,])
  Xtest <- rbind(X1[-id1,], X2[-id2,], X3[-id3,])
  ldatrain <- lda(Type ~ ., data=Xtrain, prior=rep(1/3, 3))
  confusionTest <- table(Xtest$Type, predict(ldatrain, newdata=Xtest)$class)
  confusionTest
  n <- sum(confusionTest)
  aer[k] <- (n - sum(diag(confusionTest))) / n
}
mean(aer)

```

Split
cing
+
sting

```

## [1] 0.01622642 ← mis class
##### CV QDA #####
# make confusion matrix (and APER)
confusion <- table(wine$Type, predict(qdamod)$class)
confusion

##
##      1  2  3
## 1 59  0  0
## 2  1 70  0
## 3  0  0 48

n <- sum(confusion)
aper <- (n - sum(diag(confusion))) / n
aper

## [1] 0.005617978

# use CV to get expected AER
qdamodCV <- qda(Type ~ ., data=wine, prior=rep(1/3, 3), CV=TRUE)
confusionCV <- table(wine$Type, qdamodCV$class)
confusionCV

##
##      1  2  3
## 1 59  0  0
## 2  1 70  0
## 3  0  0 48

eaer <- (n - sum(diag(confusionCV))) / n
eaer

## [1] 0.005617978

# OR
f2 <- qda(Type ~ ., data=wine, CV=T)
sum(f2$class != wine$Type) / length(wine$Type) # CV error quadratic

## [1] 0.005617978

# split into training and testing
set.seed(1)
id1 <- sample.int(n=59, size=41)
id2 <- sample.int(n=71, size=50)
id3 <- sample.int(n=48, size=34)
Xtrain <- rbind(X1[id1,], X2[id2,], X3[id3,])
Xtest <- rbind(X1[-id1,], X2[-id2,], X3[-id3,])

# fit qda to training and evaluate on testing
qdatrain <- qda(Type ~ ., data=Xtrain, prior=rep(1/3, 3))
confusionTest <- table(Xtest$Type, predict(qdatrain, newdata=Xtest)$class)
confusionTest

##
##      1  2  3
## 1 17  1  0
## 2  0 21  0
## 3  0  0 14

```

leave on
out

qda → qda.
↑
smaller
error
logans
of the data

```

n <- sum(confusionTest)
aer <- (n - sum(diag(confusionTest))) / n
aer

## [1] 0.01886792
# split into training and testing (100 splits)
nrep <- 100
aer <- rep(0, nrep)
set.seed(1)
for(k in 1:nrep){
  #cat("rep:", k, "\n")
  id1 <- sample.int(n=59, size=41)
  id2 <- sample.int(n=71, size=50)
  id3 <- sample.int(n=48, size=34)
  Xtrain <- rbind(X1[id1,], X2[id2,], X3[id3,])
  Xtest <- rbind(X1[-id1,], X2[-id2,], X3[-id3,])
  qdatrain <- qda(Type = ., data=Xtrain, prior=rep(1/3, 3))
  confusionTest <- table(Xtest$Type, predict(qdatrain, newdata=Xtest)$class)
  confusionTest
  n <- sum(confusionTest)
  aer[k] <- (n - sum(diag(confusionTest))) / n
}
mean(aer)

```

```
## [1] 0.01528302
```

```
##### visualize LDA and QDA results via PCA #####
```

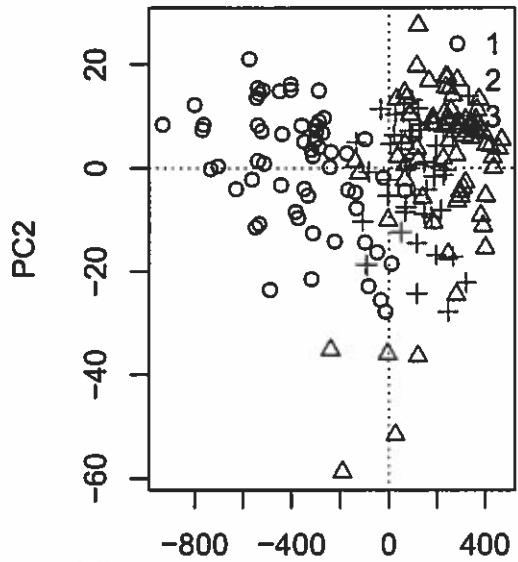
```

ldaaid <- as.integer(predict(ldamod)$class)
qdaaid <- as.integer(predict(qdamod)$class)
pcamod <- princomp(wine[, -1])
par(mfrow=c(1,2))
← plot(pcamod$scores[,1:2], xlab="PC1", ylab="PC2", pch=ldaaid, col=ldaaid, main="LDA Results")
legend("topright", lev, pch=1:3, col=1:3, bty="n")
abline(h=0, lty=3)
abline(v=0, lty=3)
plot(pcamod$scores[,1:2], xlab="PC1", ylab="PC2", pch=qdaaid, col=qdaaid, main="QDA Results")
legend("topright", lev, pch=1:3, col=1:3, bty="n")
abline(h=0, lty=3)
abline(v=0, lty=3)

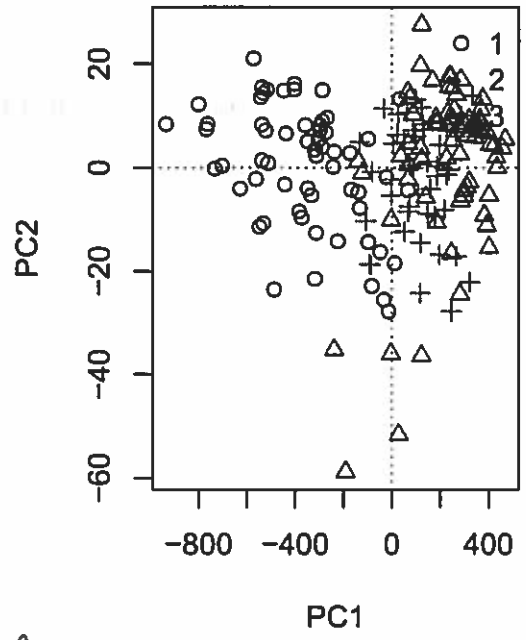
```

principle comp
pc scores apply linear regression on the
new scores

LDA Results



QDA Results



PC1 first PC1
separate wine 1 to wine 3.

7. Logistic Regression and Classification

* When the number of groups $g = 2$, an alternative to model-based classification is logistic regression. Assume that $Y_i \in \{0, 1\}$, There are two classification probabilities $p(Y = 1)$ and $1 - p(Y = 0)$. The model is

$$\log \frac{p}{1-p} = \beta_0 + \beta^T x$$

? model based vs model free
in machine learning.

where $\beta^T = (\beta_1, \dots, \beta_p)^T$. The odds ratio

$$\text{odds} = \frac{p}{1-p} \leftarrow \text{success prop} \leftarrow \text{failure prop}$$

which is the ratio of the probability of 1 to be the probability of 0. The logistic model can be written as

$$p(x) = \frac{\exp(\beta_0 + \beta^T x)}{1 + \exp(\beta_0 + \beta^T x)}$$

* Note that x is allocated to the $Y = 1$ group when the odds ratio is greater than 1. Equivalently, we have simple linear discriminant rule

$$\hat{\beta}_0 + \hat{\beta}^T x > 0$$

which is similar to Fisher's linear rule.

A more powerful version is additive logistic regression

$$\log \frac{p}{1-p} = \beta_0 + g_1(x_1) + \dots + g_p(x_p)$$

This can provide nonlinear boundaries in \mathcal{R}_p . Additive logistic regression is available in gam. There is also locally-weighted logistic regression and kernel logistic regression for nonlinear classification. Logistic regression can be thought of as providing a "nonparametric" linear discriminant rule.

< Can be used in final project .

math read writing

Example:

In the satgradu data, the first three columns are the SAT scores on the 3 tests, the last column is an indicator of whether the student successfully graduated (1 = graduated, 0 = did not graduate)

```
satgradu <- read.table("satgradu.txt", header=T)
attach(satgradu)
glm.fit.sat <- glm(gradu ~ math + reading + writing, data=satgradu, family=binomial)
summary(glm.fit.sat)
```

```
##
## Call:
## glm(formula = gradu ~ math + reading + writing, family = binomial,
## data = satgradu)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.95784  0.04024  0.65174  0.78397  1.09748
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.4940606  2.5348993  -0.589   0.556
## math         0.0012308  0.0047303   0.260   0.795
## reading      0.0037698  0.0034071   1.106   0.269
## writing       0.0002277  0.0052656   0.043   0.966
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 44.987  on 39  degrees of freedom
## Residual deviance: 43.116  on 36  degrees of freedom
## AIC: 51.116
##
## Number of Fisher Scoring iterations: 4
```

not any are significant

biggest far away from 0 => more significant in predicted gradient

=> more important factor to predict graduation rate

The z-tests tell us which explanatory variables are most important in predicting whether a student graduates. We see the reading score is the most important explanatory variable in predicting graduation.

Now let us predict whether a new applicant with SAT scores of: math = 550, reading = 610, writing = 480 will graduate:

```
newobs <- rbind( c(550,610,480) )
dimnames(newobs) <- list(NULL,c('math','reading','writing'))
newobs <- data.frame(newobs)
predict(glm.fit.sat,newdata=newobs, type='response')
```

← predict for the new data

```
##      1
## 0.830873 ← high probability of graduate
```

Making predictions for several new individuals at once:

```
newobs <- rbind( c(300,420,280), c(510,480,470), c(780,760,710) )
dimnames(newobs) <- list(NULL,c('math','reading','writing'))
newobs <- data.frame(newobs)
predict(glm.fit.sat,newdata=newobs, type='response')
```

```
##      1      2      3
## 0.6276845 0.7408203 0.9236341
```

Check the misclassification rate of logistic regression rule:

1. Simple plug-in misclassification rate:

```
group.probs<-predict(glm.fit.sat, satgradu, type='response')
pred.group <- rep(0,times=nrow(satgradu))
pred.group[group.probs > 0.5] <- 1
table(pred.group,gradu)
```

```
##          gradu
## pred.group 0 1
##           1 10 30
```

The plug-in misclassification rate for logistic regression here is $10/40 = 0.25$, but all the students have been predicted to graduate! If we change the cutoff to 0.6, then two people are predicted to not graduate, and the misclassification rate is still $(1 + 9)/40 = 0.25$

```
group.probs<-predict(glm.fit.sat, satgradu, type='response')
pred.group <- rep(0,times=nrow(satgradu))
pred.group[group.probs > 0.6] <- 1
table(pred.group,gradu)
```

```
##          gradu
## pred.group 0 1
##           0 1 1
##           1 9 29
```

2. cross-validation rate of logistic regression rule:

```
correct<-rep(0,times=nrow(satgradu) )
for (j in 1:nrow(satgradu) ) {
  glm.fit.no.j<-glm(gradu ~ math + reading + writing, data=satgradu, family=binomial, subset=-j)
  mypred<-(predict(glm.fit.no.j,newdata=satgradu[j,1:3],typ='response') > 0.5)
  correct[j] <- (mypred==gradu[j])
}
cv.misclass <- 1-mean(correct)
cv.misclass
## [1] 0.275
```

The cross-validation misclassification rate for logistic regression here is 0.275.

Logistic regression with high dimensional covariates

reduce # of factors that affect the outcome

to shrink some parameters.

The **LASSO** (least absolute shrinkage and selection operator) and **elastic net** are two regularized regression approaches. "Regularization" broadly means to impose constraints to solve overparameterized problems. LASSO and elastic net both shrink logistic regression coefficients toward zero, and both are formulated as penalized regression. The LASSO maximizes

$$L_{\lambda}(\beta) = \frac{1}{2} \log \left[\prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \right] - \lambda \sum_{j=1}^p |\beta_j|$$

where $p_i = \frac{\exp(\beta_0 + \beta^T x_i)}{1 + \exp(\beta_0 + \beta^T x_i)}$ and $\beta = (\beta_1, \dots, \beta_p)^T$. Often λ is chosen through k-fold cross-validation, i.e. chosen to minimize prediction error. Typically covariates are standardized to have unit variance before using LASSO.

```
library(glmnet)
```

```
## Warning: package 'glmnet' was built under R version 3.4.4
```

```
## Loading required package: Matrix
```

```
##
```

```
## Attaching package: 'Matrix'
```

```
## The following object is masked from 'package:reshape':
```

```
##
```

```
## expand
```

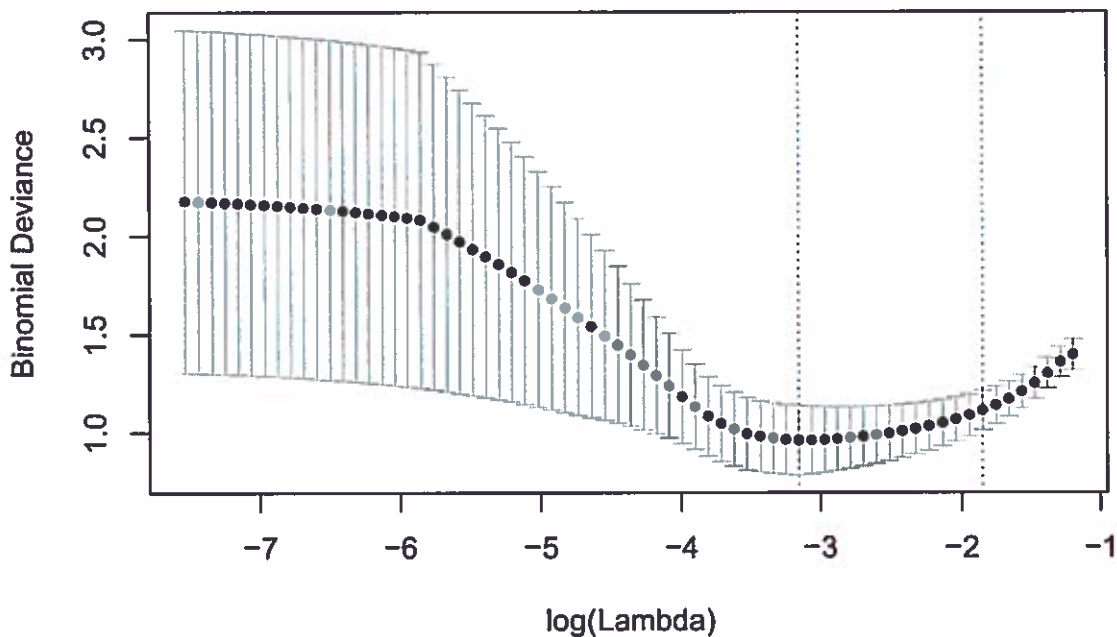
```
## Loading required package: foreach
```

```
## Warning: package 'foreach' was built under R version 3.4.3
```

```
## Loaded glmnet 2.0-16
```

```
f <- cv.glmnet(as.matrix(d2[,3:6]), d2[,2], family="binomial", grouped=FALSE)
plot(f)
```

4 4 4 4 4 4 4 4 3 3 2 1 1 1 1 1 1 0



The plot shows that the log of the optimal value of lambda (i.e. the one that minimizes the root mean square error) is approximately -3. The exact value can be viewed by examining the variable 'lambda_min' in the code below.

```
#min value of lambda
lambda_min <- f$lambda.min
lambda_min
```

```
## [1] 0.04257696
```

The cv.glmnet function finds the value of lambda that gives the simplest model but also lies within one standard error of the optimal value of lambda. This value of lambda (lambda.1se) is what we'll use in the rest of the computation.

*ex dental data set
predict gender based on teeth measurement.*

```

#best value of lambda
lambda_1se <- f$lambda.1se
lambda_1se

```

the best lambda.

```

## [1] 0.1566142 ←
#regression coefficients
coef(f,s=lambda_1se)

```

*direct
plugin.*

```

## 5 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept) -5.8427230
## V1           .
## V2           .
## V3           .
## V4           -0.2395536 ← only the last variable being used.
f2 <- glmnet(as.matrix(d2[,3:6]),d2[,2],family="binomial")
pred.sex <- predict(f2,newx=as.matrix(d2[,3:6]),s=lambda_1se,type="class")
table(pred.sex,d2$Sex)

##
## pred.sex Male Female
## Female    0      6
## Male     16      5

```

*know which factor is
the most important, how
to decide a person
has or does not
have diabetes.*

The misclassification error is $5/27 = 0.185$

```

cerror=0
for(i in 1:dim(d2)[1]){
learn  
not lambda=cv.glmnet(as.matrix(d2[-i,3:6]),d2[-i,2],family="binomial",grouped=FALSE)$lambda_1se
f=glmnet(as.matrix(d2[-i,3:6]),d2[-i,2],family="binomial")
pred.sex=predict(f,newx=as.matrix(d2[i,3:6]),s=lambda_1se,type="class")
if(d2[i,2]!=pred.sex){cerror=cerror+1}
}

```

```

## Warning in lognet(x, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has fewer than 8 observations; dangerous
## ground

```

```

## Warning in lognet(x, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has fewer than 8 observations; dangerous
## ground

```

```
cerror/dim(d2)[1]
```

```
## [1] 0.2222222
```

We use cross-validation to find $\hat{\lambda}$. Note that LASSO can shrink coefficients to zero. That is why it is also useful for variable selection. It provides an "automatic" stepwise procedure. As $\lambda \rightarrow \infty$, the usual MLE estimates from logistic regression are obtained.

Example Pima Indian Diabetes dataset from the 'mlbench' package

```

library(mlbench)
data("PimaIndiansDiabetes")
attach(PimaIndiansDiabetes)

```

```
## The following object is masked from package:datasets:
```

```
##
```

```
## pressure
```

```
head(PimaIndiansDiabetes)
```

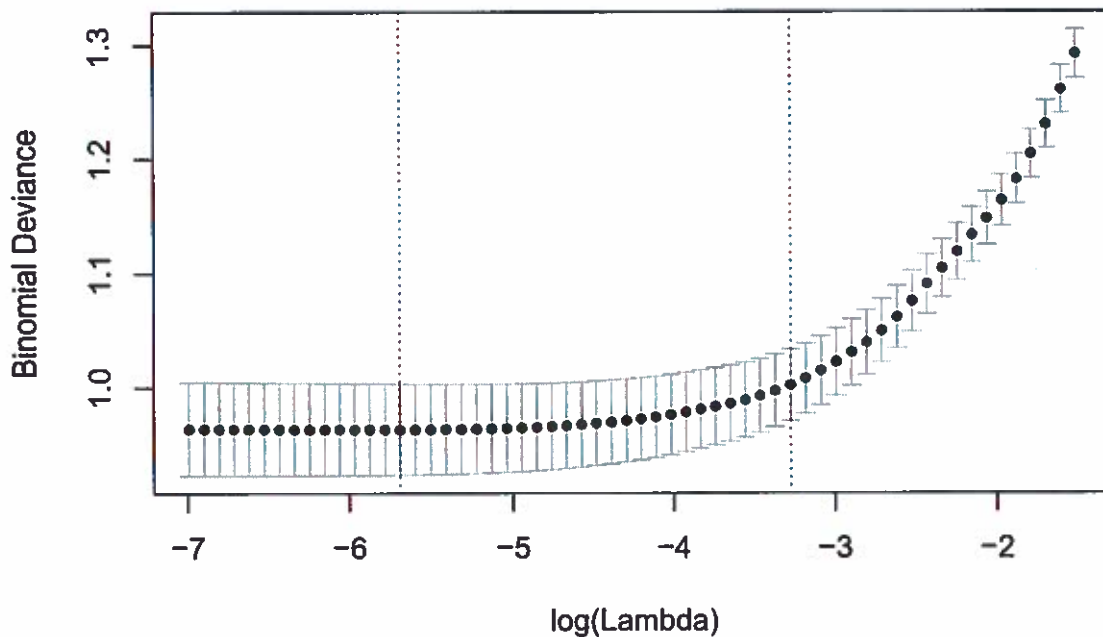
```
##      weekpregnant glucose pressure triceps insulin mass pedigree age diabetes
## 1      6      148      72      35      0 33.6   0.627 50      pos
## 2      1      85      66      29      0 26.6   0.351 31      neg
## 3      8      183      64      0      0 23.3   0.672 32      pos
## 4      1      89      66      23     94 28.1   0.167 21      neg
## 5      0     137      40      35    168 43.1   2.288 33      pos
## 6      5     116      74      0      0 25.6   0.201 30      neg
```

```
library(glmnet)
```

```
f <- cv.glmnet(as.matrix(PimaIndiansDiabetes[,-9]),
               PimaIndiansDiabetes[,9],
               family="binomial",grouped=FALSE)
```

```
plot(f)
```

7 7 7 7 7 7 7 7 7 6 6 5 5 5 3 3 2 1 1



```
#min value of lambda
```

```
lambda_min <- f$lambda.min
```

```
lambda_min
```

```
## [1] 0.003380156
```

```
#best value of lambda
```

```
lambda_1se <- f$lambda.1se
```

```
lambda_1se
```

```
## [1] 0.03797011
```

```
#regression coefficients
```

```
coef(f,s=lambda_min)
```

```

## 9 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept) -8.0790626793
## pregnant    0.1162781653
## glucose      0.0336913048
## pressure     -0.0111001234
## triceps . ← be removed from logistic model
## insulin     -0.0008656466
## mass         0.0838838054
## pedigree     0.8561335099
## age          0.0136150356
f2 <- glmnet(as.matrix(PimaIndiansDiabetes[,-9]),
             PimaIndiansDiabetes[,9],family="binomial")
pred.sex <- predict(f2,newx=as.matrix(PimaIndiansDiabetes[,-9]),
                   s=lambda_min,type="class")
table(pred.sex,PimaIndiansDiabetes$diabetes)

##
## pred.sex neg pos
##      neg 445 115
##      pos   55 153
#The misclassification error is (55+115)/768

cverror=0
for(i in 1:dim(PimaIndiansDiabetes)[1]){
  lambda=cv.glmnet(as.matrix(PimaIndiansDiabetes[-i,-9]),
                  PimaIndiansDiabetes[-i,9],
                  family="binomial",
                  grouped=FALSE)$lambda_min
  f=glmnet(as.matrix(PimaIndiansDiabetes[-i,-9]),
           PimaIndiansDiabetes[-i,9],
           family="binomial")
  pred.diab=predict(f,newx=as.matrix(PimaIndiansDiabetes[i,-9]),
                   s=lambda_min,type="class")
  if(PimaIndiansDiabetes[i,9]!=pred.diab){cverror=cverror+1}
}
#The misclassification error
cverror/dim(PimaIndiansDiabetes)[1]

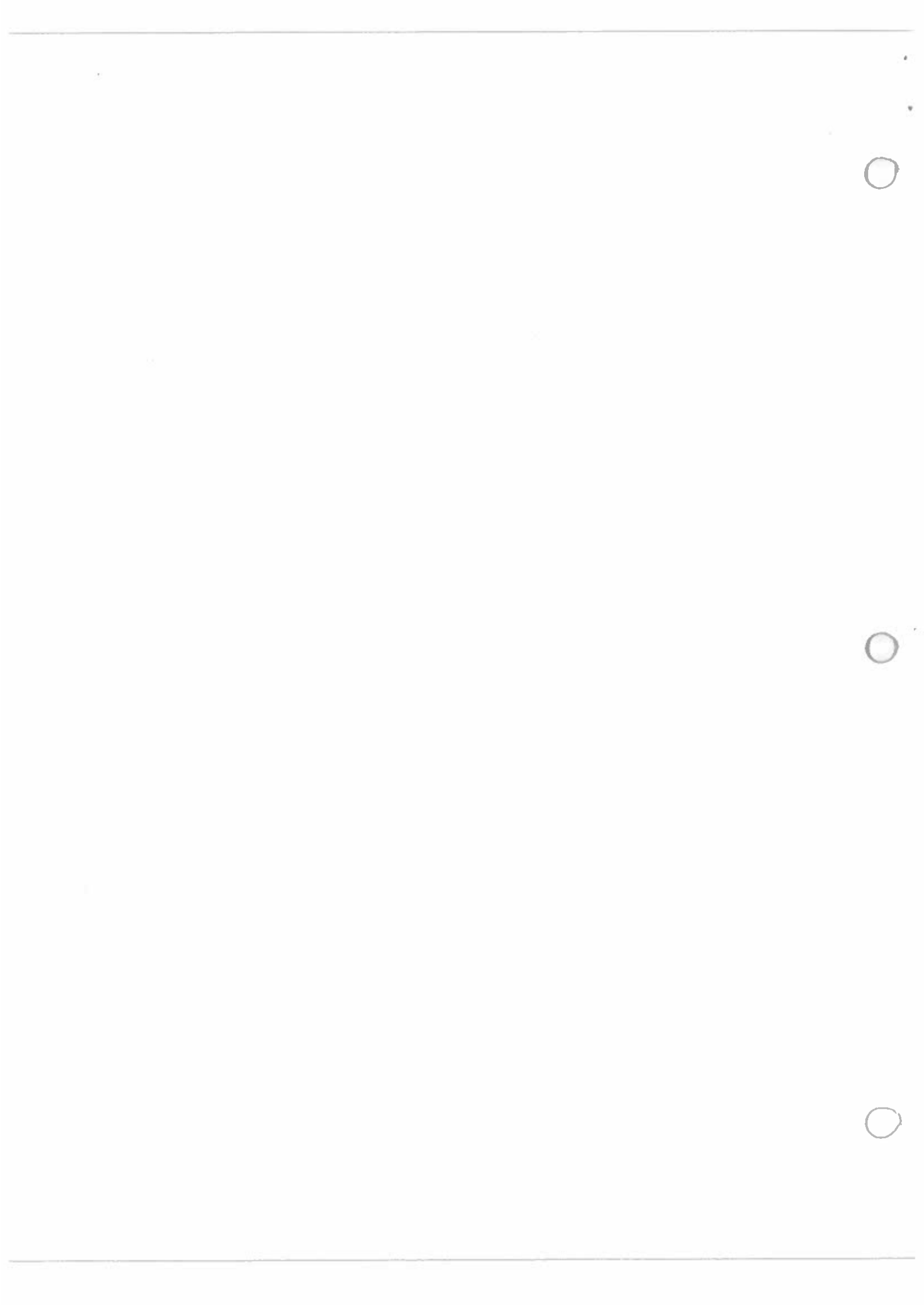
## [1] 0.2226562

```

direct phy. method

leave one out method

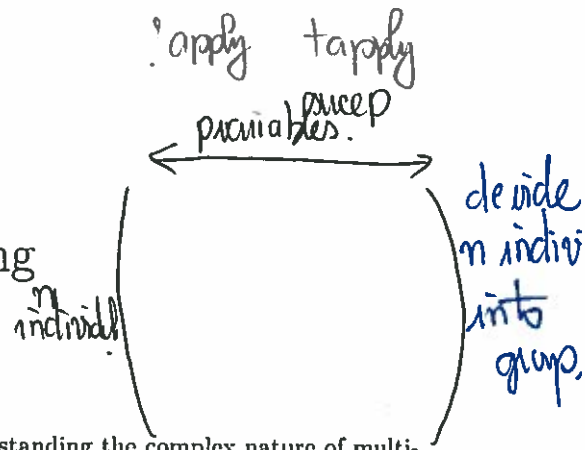
? glmnet vs cv.glmnet



Chapter 8. Clustering

Jianxuan Liu

Fall 2018



Rudimentary, exploratory procedures are often quite helpful in understanding the complex nature of multivariate relationships. In this chapter, we will discuss some displays based on certain measures of distance and suggested step-by-step rules for grouping objects. Searching the data for a structure of "natural" groupings is an important exploratory technique. Groupings can provide an informal means for assessing dimensionality, identifying outliers and suggesting interesting hypotheses concerning relationships.

The major goal of cluster analysis is to separate individual observations, or items, into groups, or clusters, on the basis of the values for the p variables measured on each individual. Often in clustering the items are called objects. We wish to create clusters such that the objects within each cluster are similar and objects in different clusters are dissimilar. The dissimilarity between any two objects is typically quantified using a distance measure (like Euclidean distance).

Applications of Cluster Analysis:

- In marketing, researchers attempt to find distinct clusters of the consumer population so that several distinct marketing strategies can be used for the clusters.
- In ecology, scientists classify plants or animals into various groups based on some measurable characteristics.
- Researchers in genetics may separate genes into several classes based on their expression ratios measured at different time points.

There are three major classes of clustering methods – from oldest to newest, they are:

1. **Hierarchical methods:** cluster the data in a series of n steps, typically joining observations together step by step to form clusters.
2. **Partitioning methods:** first determine k , and then typically attempt to find the partition into k clusters that optimizes some objective function of the data.
3. **Model-based clustering:** takes a statistical approach, formulating a model that categorizes the data into subpopulations and using maximum likelihood to estimate the model parameters.

To summarize, the basic objective in cluster analysis is to discover natural groupings of the items or variables. In turn, we must develop a quantitative scale on which to measure the association (similarity) between objects.

1. Similarity Measures

Most efforts to produce a rather simple group structure from a complex data set require a measure of "closeness", or "similarity". There is often a great deal of subjectivity involved in the choice of a similarity measure. Important considerations include the nature of the variables, scales of measurement and subject matter knowledge.

When items are clustered, proximity is usually indicated by some sort of distance. By contrast, variables are usually grouped on the basis of correlation coefficients or like measure of association.

Consider two p -dimensional observations (items) $x = [x_1, x_2, \dots, x_p]$ and $y = [y_1, y_2, \dots, y_p]$, then the distance of x, y is

1. Euclidean Distance

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2} = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

2. Manhattan Distance:

$$d(\mathbf{x}, \mathbf{y}) = |x_1 - y_1| + |x_2 - y_2| + \dots + |x_p - y_p| = \sum_{i=1}^p |x_i - y_i|$$

3. Statistical Distance

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{A} (\mathbf{x} - \mathbf{y})}$$

Ordinarily, $\mathbf{A} = \mathbf{S}^{-1}$ where \mathbf{S} contain the sample variances and covariances. However, without prior knowledge of the distinct groups, these sample quantities cannot be computed. For this reason, Euclidean distance is often preferred for clustering.

- Minkowski Metric

$$d(\mathbf{x}, \mathbf{y}) = \left[\sum_{i=1}^p |x_i - y_i|^m \right]^{1/m}$$

For $m = 1$, $d(\mathbf{x}, \mathbf{y})$ measure the "city-block" distance between two points in p dimensions. For $m = 2$, $d(\mathbf{x}, \mathbf{y})$ becomes the Euclidean distance. In general, varying m changes the weight given to larger and smaller differences.

- Canberra Metric

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{(x_i + y_i)}$$

- Czekanowski Coefficient

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{2 \sum_{i=1}^p \min(x_i, y_i)}{\sum_{i=1}^p (x_i + y_i)}$$

Note that Canberra Metric and Czekanowski coefficient are defined for nonnegative variables only.

Example: Characteristics matching

Characteristics Matching. If the characteristic is present when value is 1. If the characteristic is absent when value is 0.

Variables	1	2	3	4	5
Item 1	1	0	0	1	1
Item 2	1	1	0	1	0

Let x_{ij} be the score (1 or 0) of the j^{th} variable on the i^{th} item and x_{kj} be the score of the k^{th} item, $j = 1, 2, \dots, p$. Consequently, the squared Euclidean distance $\sum_{i=1}^p (x_{ij} - x_{kj})^2$ provides a count of the number of mismatches. A large distance corresponds to many mismatches, that is, dissimilar items. From the preceding display, the square of the distance between item i and k would be

$$\sum_{i=1}^p (x_{ij} - x_{kj})^2 = (1 - 1)^2 + (0 - 1)^2 + (0 - 0)^2 + (1 - 1)^2 + (1 - 0)^2 = 2$$

Although a distance based on $\sum_{i=1}^p (x_{ij} - x_{kj})^2$ might be used to measure similarity, it suffers from weighting the 1 - 1 and 0 - 0 matches equally. In some cases, a 1 - 1 match is a stronger indication of similarity than a

0 – 0 match. For instance, in grouping people, the evidence that two persons both read ancient Greek is stronger evidence of similarity than the absence of this ability. Thus, it might be reasonable to discount the 0 – 0 matches or even disregard them completely. To allow for differential treatment of the 1 – 1 matches and the 0 – 0 matches, several schemes for defining similarity coefficients have been suggested. To introduce these schemes, let us arrange the frequencies of matches and mismatches for items i and k in the form of a contingency table:

		Item k		Totals
		1	0	
Item i	1	a	b	$a + b$
	0	c	d	$c + d$
Totals		$a + c$	$b + d$	$p = a + b + c + d$

In this table, a represents the frequency of 1 – 1 matches, b is the frequency of 1 – 0 matches, and so forth. Given the foregoing five pairs of binary outcomes, $a = 2$ and $b = c = d = 1$. Here is a list of common similarity coefficients defined in terms of the frequencies

- Euclidean distance 1 – 1 and 0 – 0 are equal: $\frac{a + d}{p}$
- Double weight to 1 – 1 and 0 – 0: $\frac{2(a + d)}{2(a + d) + b + c}$
- Double weight for unmatched pairs: $\frac{a + d}{a + d + 2(b + c)}$
- Only 1 – 1: $\frac{a}{p}$
- Exclude 0 – 0 entirely: $\frac{a}{a + b + c}$
- Ratio of matches to mismatches with 0-0 matches excluded: $\frac{a}{b + c}$

2. Hierarchical Clustering Methods

Hierarchical clustering, as is denoted by the name, involves organizing your data into a kind of hierarchy. The common approach is called an **agglomerative approach**. Agglomerative hierarchical methods start with the individual objects. Thus, there are initially as many clusters as objects. The most similar objects are first grouped, and these initial groups are merged according to their similarities. Eventually, as the similarity decreases, all subgroups are fused in to a single cluster.

Imagine there is all these little particles floating around (your data points), and you start kind of grouping them together into little balls. And then the balls get grouped up into bigger balls, and the bigger balls get grouped together into one big massive cluster. That's the agglomerative approach to clustering. The algorithm is recursive and goes as follows:

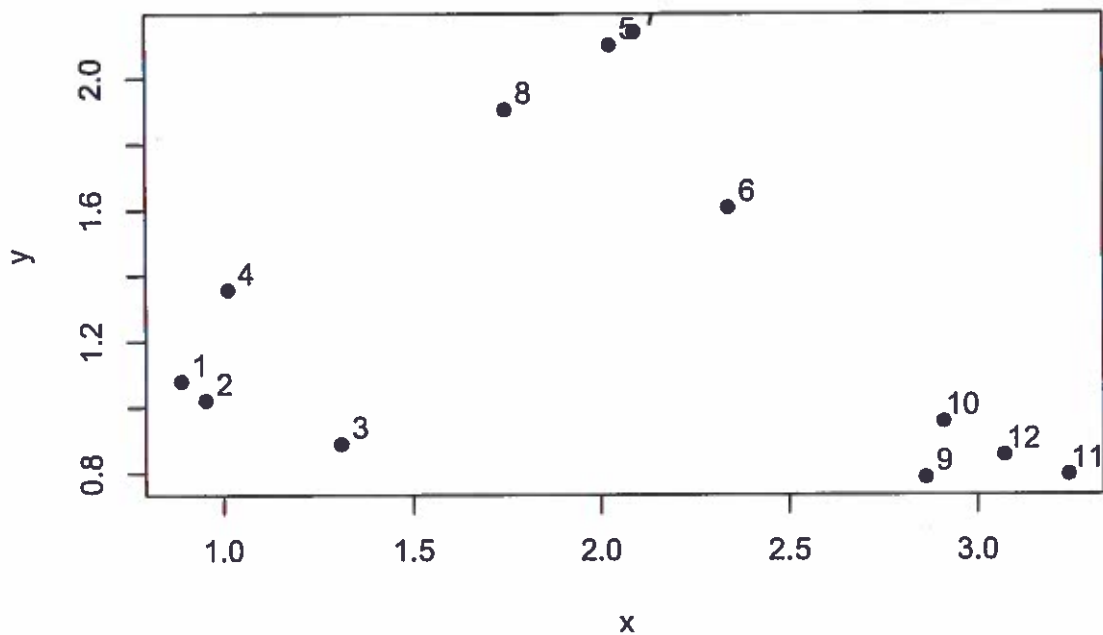
1. Start with n clusters and on $N \times N$ distance matrix $D = \{d_{ik}\}$
2. Search for nearest/most similar pairs d_{uv} where U and V are most similar.
3. Merge clusters U and V to (UV) . Update the distance matrix
4. Repeat steps 1-3 $N - 1$ times

Example

Here is a simple example demonstrating how hierarchical clustering works. First we'll simulate some data in three separate clusters.

```
set.seed(123)
x = rnorm(12, rep(1:3, each = 4), 0.2)
y = rnorm(12, rep(c(1, 2, 1), each = 4), 0.2)
plot(x, y, col = 2, pch = 19, cex = 1, main='Simulated Cluster Data')
text(x + 0.05, y + 0.05, labels = as.character(1:12))
```

Simulated Cluster Data



The first step in the basic clustering approach is to calculate the distance between every point with every other point. The result is a distance matrix, which can be computed with the `dist()` function in R. Here is just a piece of the distance matrix associated with the figure above.

```
df=data.frame(x=x, y=y)
round(dist(df),4)
```

```
##      1      2      3      4      5      6      7      8      9     10
## 2  0.0879
## 3  0.4650 0.3818
## 4  0.3046 0.3406 0.5551
## 5  1.5278 1.5198 1.4056 1.2548
## 6  1.5474 1.5070 1.2565 1.3521 0.5861
## 7  1.6044 1.5955 1.4749 1.3324 0.0778 0.5896
## 8  1.1913 1.1871 1.1059 0.9151 0.3398 0.6667 0.4175
## 9  1.9964 1.9232 1.5543 1.9347 1.5571 0.9710 1.5577 1.5801
## 10 2.0267 1.9580 1.6006 1.9387 1.4457 0.8633 1.4394 1.5018 0.1767
## 11 2.3741 2.3021 1.9354 2.3006 1.7856 1.2134 1.7717 1.8647 0.3823 0.3710
## 12 2.1957 2.1246 1.7606 2.1185 1.6264 1.0476 1.6168 1.6913 0.2200 0.1908
##      11
```

```
## 2
## 3
## 4
## 5
## 6
## 7
## 8
## 9
## 10
## 11
## 12 0.1828
```

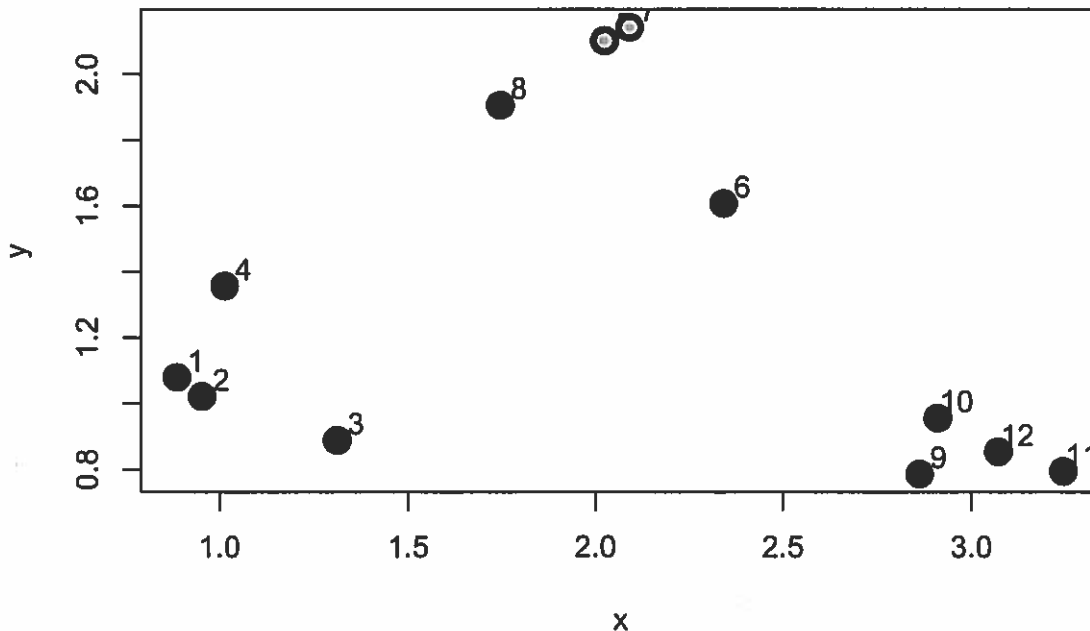
The default distance metric used by the 'dist()' function is Euclidean distance. First an agglomerative clustering approach attempts to find the two points that are closest together. In other words, we want to find the smallest non-zero entry in the distance matrix.

```
rdistxy = as.matrix(dist(df))
## Remove the diagonal from consideration
diag(rdistxy)=diag(rdistxy) + 100000
# Find the index of the points with minimum distance
ind=which(rdistxy == min(rdistxy), arr.ind = TRUE)
ind

##   row col
## 7   7   5
## 5   5   7
```

Now we can plot the points and show which two points are closest together according to our distance metric.

```
plot(x, y, col = "blue", pch = 19, cex = 2)
text(x + 0.05, y + 0.05, labels = as.character(1:12))
points(x[ind[1, ]], y[ind[1, ]], col = "orange", pch = 19, cex = 1)
```



The next step for the algorithm is to start drawing the tree, the first step of which would be to "merge" these two points together.

```

par(mfrow = c(1, 2))
plot(x, y, col = "blue", pch = 19, cex = 1, main = "Data")
text(x + 0.05, y + 0.05, labels = as.character(1:12))
points(x[ind[1, ]], y[ind[1, ]], col = "orange", pch = 19, cex = 1)
# Make a cluster and cut it at the right height
library(dplyr)

##
## Attaching package: 'dplyr'

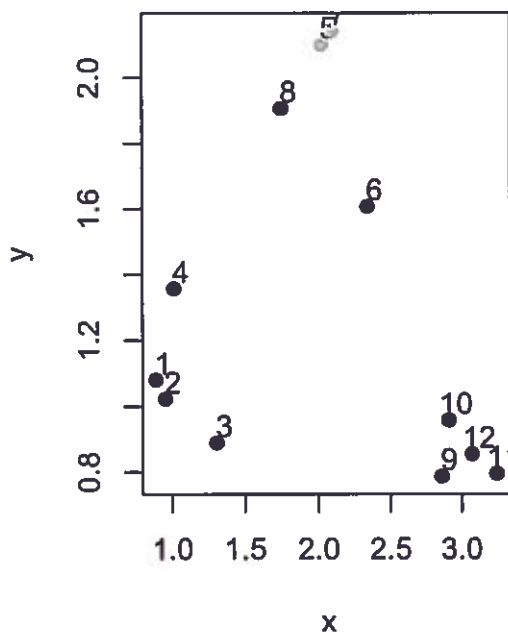
## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

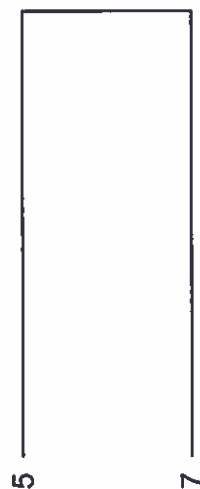
hcluster=dist(df) %>% hclust
dendro=as.dendrogram(hcluster)
cutDendro=cut(dendro, h = (hcluster$height[1] + 0.00001))
plot(cutDendro$lower[[1]], yaxt = "n", main = "Begin building tree")

```

Data



Begin building tree



Now that we've merged the first two "leaves" of this tree, we can turn the algorithm crank and continue to build the tree. Now, the two points we identified in the previous iteration will get "merged" into a single point. We need to search the distance matrix for the next two closest points, ignoring the first two that we already merged.

```

nextmin=rdistxy[order(rdistxy)][3]
ind= which(rdistxy == nextmin,arr.ind=TRUE)
ind

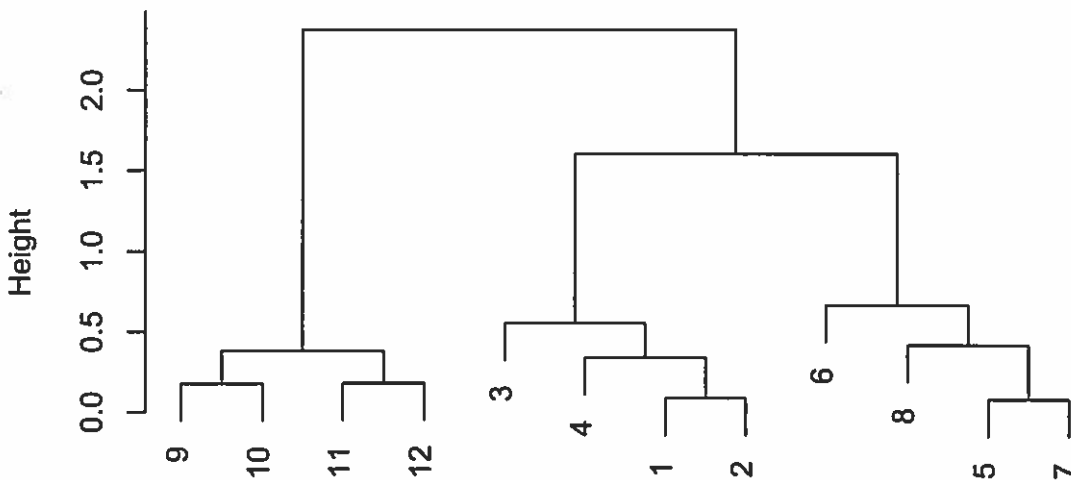
```

```
## row col
## 2 2 1
## 1 1 2
```

We continue with the next pair of points and the merged tree leaves. And on and on in this manner. If we were to continue in this fashion - identifying the two closest points and merging them, we'd end up with a dendrogram that looks like this one. Here, we call the `hclust()` to run the clustering algorithm.

```
hClustering = data.frame(x=x,y=y) %>% dist %>% hclust
plot(hClustering)
```

Cluster Dendrogram



```
hclust (*, "complete")
```

From the tree/dendrogram it's clear that there are three clusters each with four points.

Refined Dendrograms

It is desirable to make prettier dendrograms with color coded each of the cluster members by their cluster membership.

```
myplclust <- function(hclust, lab = hclust$labels,
                      lab.col = rep(1,length(hclust$labels)),
                      hang = 0.1, ...)
{
  y=rep(hclust$height, 2)
  x=as.numeric(hclust$merge)
  y=y[which(x < 0)]
  x=x[which(x < 0)]
  x=abs(x)
  y=y[order(x)]
  x=x[order(x)]
}
```

```

plot(hclust, labels = FALSE, hang = hang, ...)
text(x = x, y = y[hclust$order] - (max(hclust$height) * hang),
     labels = lab[hclust$order], col = lab.col[hclust$order],
     srt = 90, adj = c(1, 0.5), xpd = NA, ...) }

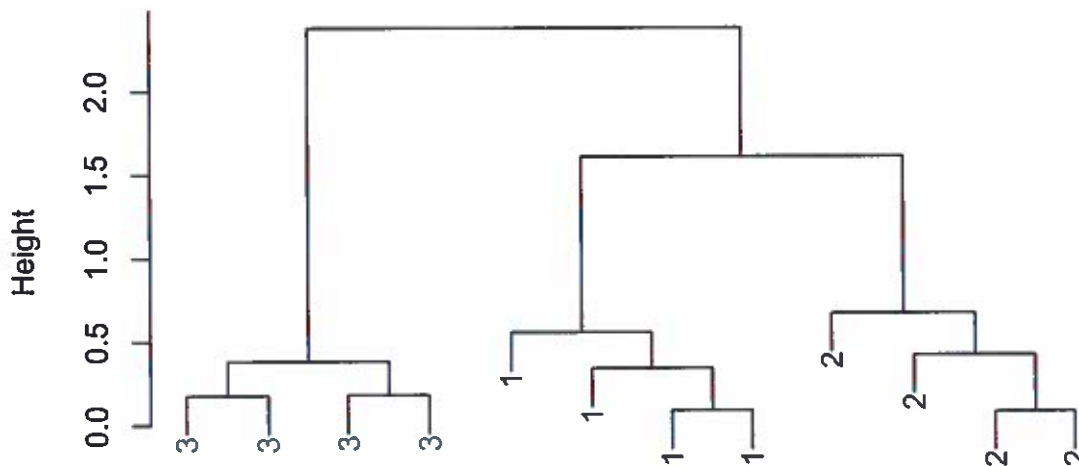
```

```

hClustering=data.frame(x = x, y = y) %>% dist %>% hclust
myplclust(hClustering, lab = rep(1:3, each = 4), lab.col = rep(1:3, each = 4))

```

Cluster Dendrogram



`hclust (*, "complete")`

3. Linkage Methods in Hierarchical Clustering

1. Single linkage clustering:

Single linkage clustering is sometimes called "nearest neighbor" clustering. At each step, joins the clusters whose minimum distance between objects is smallest, i.e., joins the clusters A and B with the smallest

$$d_{AB} = \min_{i \in A, j \in B} (d_{ij})$$

2. Complete linkage algorithm:

Complete linkage clustering is sometimes called "farthest neighbor" clustering. At each step, joins the clusters whose maximum distance between objects is smallest, i.e., joins the clusters A and B with the smallest

$$d_{AB} = \max_{i \in A, j \in B} (d_{ij})$$

3. Average linkage algorithm:

At each step, joins the clusters whose average distance between objects is smallest, i.e., joins the clusters A and B with the smallest

$$d_{AB} = \frac{1}{n_A n_B} \sum_{i \in A} \sum_{j \in B} d_{ij}$$

where n_A and n_B are the number of objects in clusters A and B , respectively

Example

To illustrate the linkage algorithms, we consider the hypothetical distances between pairs of five objects as follows:

$$D = \{d_{ij}\} = \begin{matrix} & A & B & C & D & E \\ \begin{matrix} A \\ B \\ C \\ D \\ E \end{matrix} & \begin{pmatrix} 0 & & & & \\ 9 & 0 & & & \\ 3 & 7 & 0 & & \\ 6 & 5 & 9 & 0 & \\ 11 & 10 & 2 & 8 & 0 \end{pmatrix} \end{matrix}$$

Single linkage clustering:

Treating each object as a cluster, we commence clustering by merging the two closest items. Since $\min_{i,k} = d_{EC} = 2$, object E and C are merged to form the cluster (CE). To implement the next level of clustering, we need the distance between the cluster (CE) and the remaining objects A, B and D. The nearest neighbor distances are

$$\begin{aligned} d_{(CE)A} &= \min\{d_{CA}, d_{EA}\} = \min\{3, 11\} = 3 \\ d_{(CE)B} &= \min\{d_{CB}, d_{EB}\} = \min\{7, 10\} = 7 \\ d_{(CE)D} &= \min\{d_{CD}, d_{ED}\} = \min\{9, 8\} = 8 \end{aligned}$$

Deleting the rows and columns of D corresponding to objects C and E, and adding a row and column for the cluster (CE), we obtain the new distance matrix

$$\begin{matrix} & CE & A & B & D \\ \begin{matrix} CE \\ A \\ B \\ D \end{matrix} & \begin{pmatrix} 0 & & & \\ 3 & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

The smallest distance between pairs of clusters in row $d_{(CE)A} = 3$, and we merge cluster A with cluster CE to get the next cluster, (ACE). Calculating

$$\begin{aligned} d_{(ACE)B} &= \min\{d_{(CE)B}, d_{AB}\} = \min\{7, 9\} = 7 \\ d_{(ACE)D} &= \min\{d_{(CE)D}, d_{AD}\} = \min\{8, 6\} = 6 \end{aligned}$$

We find that the distance matrix for the next level of clustering is

$$\begin{matrix} & ACE & B & D \\ \begin{matrix} A \\ B \\ D \end{matrix} & \begin{pmatrix} 0 & & \\ 7 & 0 & \\ 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

The minimum nearest neighbor distance between pairs of clusters is $d_{DB} = 5$, and we merge objects D and B to get the cluster (BD). At this point we have two distinct clusters (ACE) and (BD). Their nearest neighbor distance is

$$d_{(ACE)(BD)} = \min\{d_{(ACE)B}, d_{(ACE)D}\} = \min\{7, 6\} = 6$$

The final distance matrix becomes

$$\begin{array}{cc} & \begin{array}{cc} ACE & BD \end{array} \\ \begin{array}{c} ACE \\ BD \end{array} & \left(\begin{array}{cc} 0 & \\ 6 & 0 \end{array} \right) \end{array}$$

Complete linkage clustering:

At the first stage, object C and E are merged, since they are most similar. This gives the cluster (CE). At stage 2, we compute

$$\begin{aligned} d_{(CE)A} &= \max\{d_{CA}, d_{EA}\} = \max\{3, 11\} = 11 \\ d_{(CE)B} &= \max\{d_{CB}, d_{EB}\} = \max\{7, 10\} = 10 \\ d_{(CE)D} &= \max\{d_{CD}, d_{ED}\} = \max\{9, 8\} = 9 \end{aligned}$$

and the modified distance matrix becomes

$$\begin{array}{cccc} & \begin{array}{cccc} CE & A & B & D \end{array} \\ \begin{array}{c} CE \\ A \\ B \\ D \end{array} & \left(\begin{array}{cccc} 0 & & & \\ 11 & 0 & & \\ 10 & 9 & 0 & \\ 9 & 6 & 5 & 0 \end{array} \right) \end{array}$$

The next merger occurs between the most similar groups, B and D, to give the cluster (BD). At stage 3, we have

$$\begin{aligned} d_{(BD)(CE)} &= \max\{d_{B(CE)}, d_{D(CE)}\} = \max\{10, 9\} = 10 \\ d_{(BD)A} &= \max\{d_{BA}, d_{DA}\} = \max\{9, 6\} = 6 \end{aligned}$$

and the distance matrix

$$\begin{array}{ccc} & \begin{array}{ccc} CE & BD & A \end{array} \\ \begin{array}{c} CE \\ BD \\ A \end{array} & \left(\begin{array}{ccc} 0 & & \\ 10 & 0 & \\ 11 & 9 & 0 \end{array} \right) \end{array}$$

The next merger produces the cluster (ABD). At the final stage, the groups (CE) and (ABD) are merged as the single cluster (ABCDE) at level

$$d_{(ABD)(CE)} = \max\{d_{A(CE)}, d_{(BD)(CE)}\} = \max\{11, 10\} = 11$$

The final distance matrix becomes

$$\begin{array}{cc} & \begin{array}{cc} CE & ABD \end{array} \\ \begin{array}{c} CE \\ ABD \end{array} & \left(\begin{array}{cc} 0 & \\ 11 & 0 \end{array} \right) \end{array}$$

It's important to consider standardizing observations.

← explore this

Standardization of Observations

If the variables in the data set are of different types or are measured on very different scales, then some variables may play an inappropriately dominant role in the clustering process. In this case, it is recommended to standardize the variables in some way before clustering the objects. Possible standardization approaches:

- Divide each column by its sample standard deviation, so that all variables have standard deviation 1.
- Divide each variable by its sample range (max - min); Milligan and Cooper (1988) found that this approach best preserved the clustering structure.
- Convert data to z-scores by (for each variable) subtracting the sample mean and then dividing by the sample standard deviation - a common option in clustering software packages.

Example: Food data:

```
food=read.table("food.txt",head=T) #food.txt
dim(food);head(food)
```

```
## [1] 27 6
```

```
## Food Energy Protein Fat Calcium Iron
## 1 BB 340 20 28 9 2.6
## 2 HR 245 21 17 9 2.7
## 3 BR 420 15 39 7 2.0
## 4 BS 375 19 32 9 2.5
## 5 BC 180 22 10 17 3.7
## 6 CB 115 20 3 8 1.4
```

```
attach(food)
```

Let's first scale the data by dividing each variable by its standard deviation:

```
std=apply(food[,-1], 2, sd) # finding standard deviations of variables
food.std=sweep(food[,-1], 2, std, FUN="/")
```

? see the head(food)
after standardizing.

Calculating pairwise Euclidean distances between the (standardized) objects:

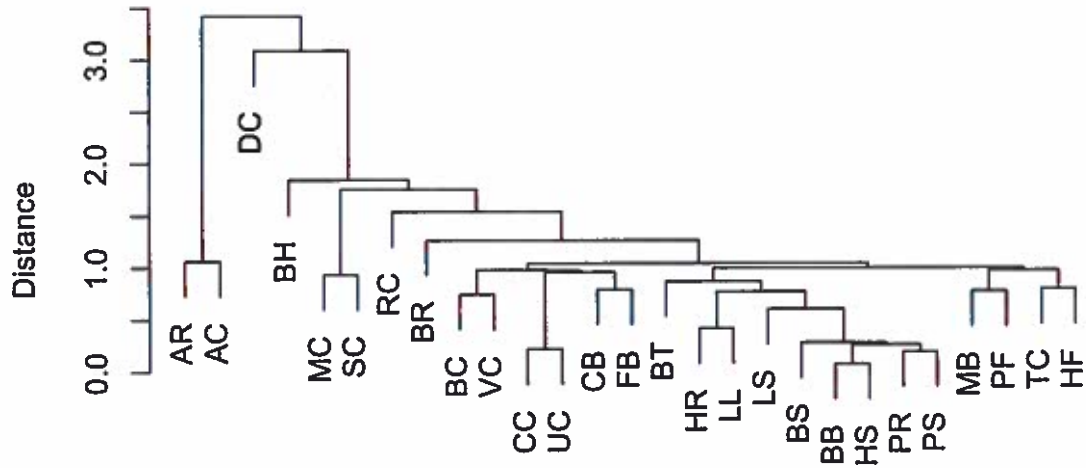
```
dist.food=dist(food.std)
```

! distinguish apply & sweep

Single linkage:

```
food.single.link=hclust(dist.food, method="single") ←
# Plotting the single linkage dendrogram:
plot(food.single.link, labels=Food, ylab="Distance")
```

Cluster Dendrogram

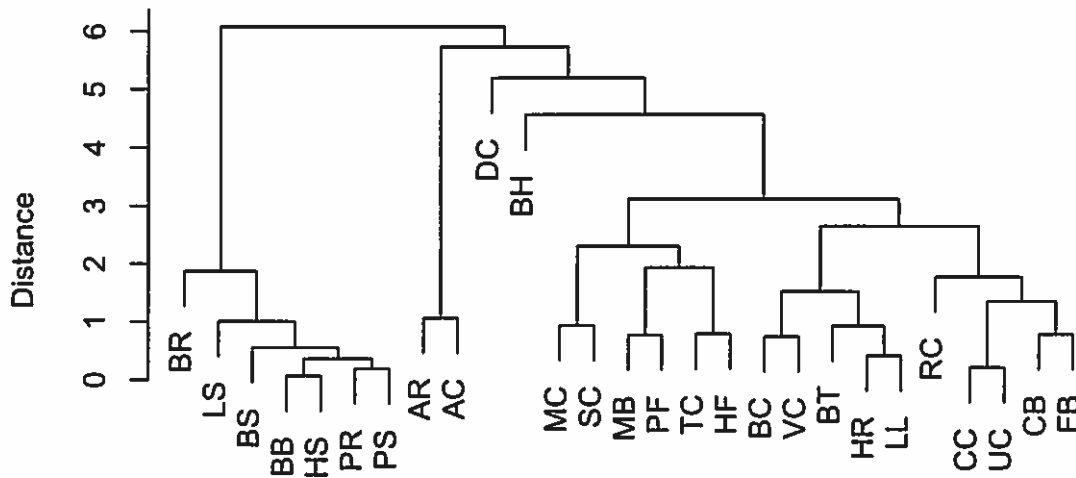


```
dist.food  
hclust (*, "single")
```

Complete linkage:

```
food.complete.link = hclust(dist.food, method='complete')  
# Plotting the complete linkage dendrogram:  
plot(food.complete.link, labels=Food, ylab="Distance")
```

Cluster Dendrogram



```
dist.food
hclust (*, "complete")
```

Average linkage:

```
food.avg.link = hclust(dist.food, method='average')
# Plotting the average linkage dendrogram:
plot(food.avg.link, labels=Food, ylab="Distance")
```

? with diabetes data, which algorithm is the best one ?

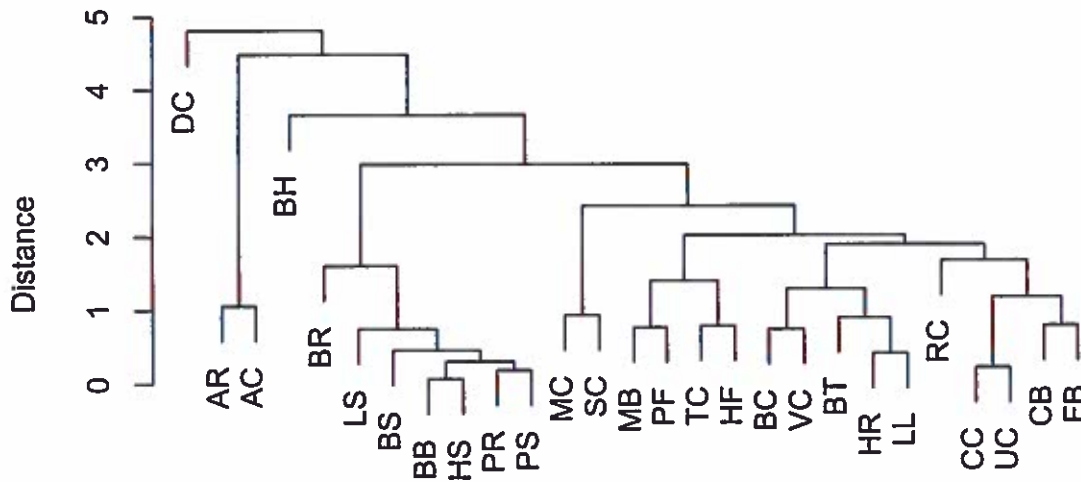
? cut the cluster dendrogram.

then the junction

cut ← `cutree` (food.avg.link, k = ...) will return the result.

the # of cluster that we want to have)

Cluster Dendrogram



```
dist.food  
hclust (*, "average")
```

Pros and Cons of Hierarchical Clustering

- An advantage of hierarchical clustering methods is their computational speed for small data sets.
- Another advantage is that the dendrogram gives a picture of the clustering solution for a variety of choices of k (the number of clusters) at once.
- On the other hand, a major disadvantage is that once two clusters have been joined, they can never be split apart later in the algorithm, even if such a move would improve the clustering.
- The so-called partitioning methods of cluster analysis do not have this restriction.
- In addition, hierarchical methods can be less efficient than partitioning methods for large data sets, when n is much greater than k .

27 Partitioning methods

Partitioning methods fix the number of clusters k and seek the best possible partition for that k . The goal is to choose the partition which gives the optimal value for some clustering criterion, or objective function. In reality, we cannot search all possible partitions to try to optimize the clustering criterion, but the algorithms are designed to search intelligently among the partitions. For a fixed k , partitioning methods are able to investigate far more possible partitions than a hierarchical method is. In practice, it is recommended to run a partitioning method for several choices of k and examine the resulting clusterings.

4. K-means Cluster (we know the # of clusters that we want to use)

The goal of K-means, the most well-known partitioning method, is to find the partition of n objects into k clusters that minimizes a within-cluster sum of squares criterion. In the traditional K-means approach, "closeness" to the cluster centers is defined in terms of squared Euclidean distance, defined by:

$$d_E^2(\mathbf{x}, \bar{\mathbf{x}}) = (\mathbf{x} - \bar{\mathbf{x}})^T(\mathbf{x} - \bar{\mathbf{x}}) = \sum_m (x_{im} - \bar{x}_{cm})^2$$

where $\mathbf{x} = (x_1, \dots, x_q)^T$ is any particular observation and $\bar{\mathbf{x}}$ is the centroid (multivariate mean vector) for cluster c .

The goal is to minimize the sum (over all objects within all clusters) of these squared Euclidean distances:

$$WSS = \sum_{c=1}^k \sum_{i \in c} d_E^2(\mathbf{x}_i, \bar{\mathbf{x}}_c)$$

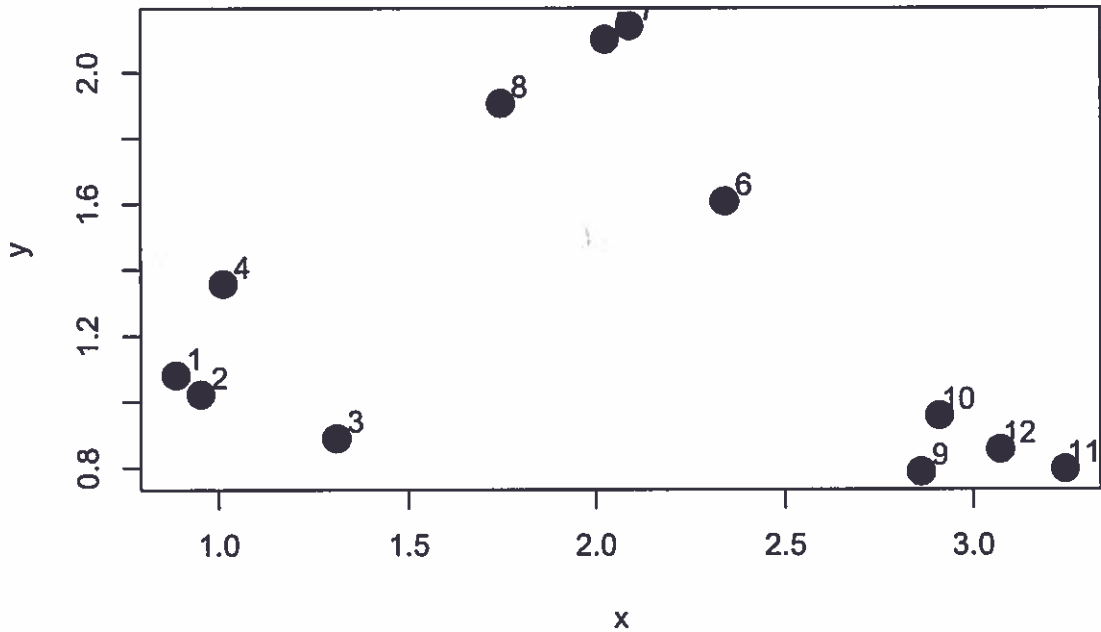
In practice, K-means will not generally achieve the global minimum of this criterion over the whole space of partitions. In fact, only under certain conditions will it achieve the local minimum (Selim and Ismail, 1984)

The K-means algorithm (MacQueen, 1967) begins by randomly allocating the n objects into k clusters (or randomly specifying k centroids). One at a time, the algorithm moves each object to the cluster whose centroid is closest to it, using the measure of closeness $d_E^2(\mathbf{x}, \bar{\mathbf{x}})$. When an object is moved, the centroids are immediately recalculated for the cluster gaining the object and the cluster losing it. The method repeatedly cycles through the objects until no reassignments of objects take place. The final clustering result will somewhat depend on the initial configuration of the objects. In practice, it is good to rerun the algorithm a few times (with different starting points) to make sure the result is stable. The R function `kmeans()` performs K-means clustering.

Example

We simulate some data from three clusters and plot the dataset below.

```
set.seed(123)
x=rnorm(12, mean = rep(1:3, each = 4), sd = 0.2)
y=rnorm(12, mean = rep(c(1, 2, 1), each = 4), sd = 0.2)
plot(x, y, col = "blue", pch = 19, cex = 2)
text(x + 0.05, y + 0.05, labels = as.character(1:12))
```



Note: The 'kmeans()' function in R implements the K-means algorithm and can be found in the 'stats' package, which comes with R and is usually already loaded when you start R. Two key parameters that you have to specify are x, which is a matrix or data frame of data, and centers which is either an integer indicating the number of clusters or a matrix indicating the locations of the initial cluster centroids. The data should be organized so that each row is an observation and each column is a variable or feature of that observation.

```
df= data.frame(x, y)
km =kmeans(df, centers = 3)
names(km)

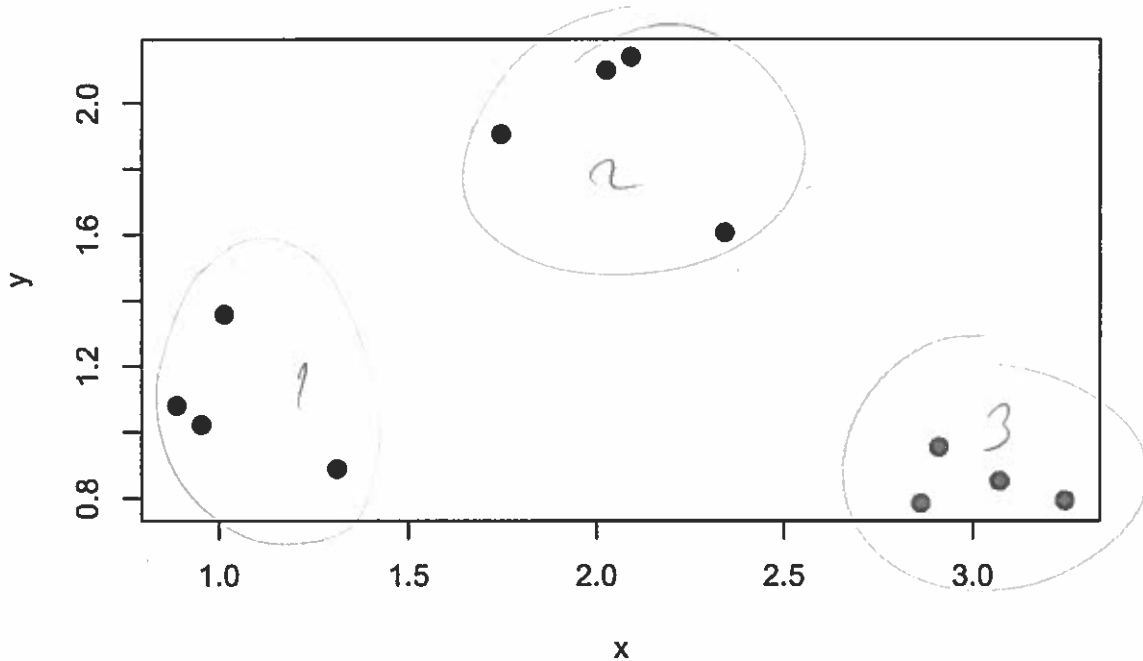
## [1] "cluster"      "centers"      "totss"       "withinss"
## [5] "tot.withinss" "betweenss"   "size"        "iter"
## [9] "ifault"

km$cluster

## [1] 3 3 3 3 1 1 1 1 2 2 2 2

plot(df, col =(km$cluster +1) , main="K-Means result with 3 clusters", pch=20, cex=2)
```

K-Means result with 3 clusters



Example: revisit food data

Consider the `food.std` data frame given above. We first consider a K-means cluster with $k = 5$

```
food.k5 <- kmeans(food.std, centers=5, iter.max=100, nstart=25)
food.k5
```

kmeans with food.std

```
## K-means clustering with 5 clusters of sizes 8, 8, 2, 8, 1
##
## Cluster means:
##   Energy Protein      Fat Calcium   Iron
## 1 3.377951 4.410004 2.56506304 0.1121302 1.6599848
## 2 1.414170 4.116004 0.57741679 0.6743833 0.6930864
## 3 0.568138 2.116802 0.08883335 0.9995610 3.9018198
## 4 1.759993 5.380205 0.77729183 0.2771219 1.9509099
## 5 1.778519 5.174405 0.79950017 4.7030629 1.7113245
##
## Clustering vector:
## [1] 1 4 1 1 4 2 4 4 1 1 1 1 1 4 4 2 3 3 2 2 2 2 2 5 4 4
##
## Within cluster sum of squares by cluster:
## [1] 4.3254549 10.2035285 0.5626614 13.0477424 0.0000000
## (between_SS / total_SS = 78.4 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"
## [5] "tot.withinss" "betweenss"    "size"         "iter"
## [9] "ifault"
```

*≠ last page
kmeans with food data*


```

# lry k=4:
food.k4 <- kmeans(food.std, centers=4, iter.max=100, nstart=25)
food.k4

## K-means clustering with 4 clusters of sizes 14, 2, 3, 8
##
## Cluster means:
##      Energy Protein      Fat  Calcium   Iron
## 1 1.619723 4.872005 0.68528586 0.2544670 1.388618
## 2 0.568138 2.116802 0.08883335 0.9995610 3.901820
## 3 1.498567 4.312004 0.68105570 2.9175222 1.140883
## 4 3.377951 4.410004 2.56506304 0.1121302 1.659985
##
## Clustering vector:
## [1] 4 1 4 4 1 1 1 1 4 4 4 4 1 1 1 2 2 1 1 1 3 1 3 1 1
##
## Within cluster sum of squares by cluster:
## [1] 28.9804747 0.5626614 6.9589520 4.3254549
## (between_SS / total_SS = 68.6 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"
## [5] "tot.withinss" "betweenss"    "size"         "iter"
## [9] "ifault"

```

5. Determining the Optimal Clusters

At some point we need to choose a single value of k to get a clustering solution. A variety of criteria have been proposed to pick the best value of k .

1. Average silhouette width:

- The average silhouette width is based on the difference between the average dissimilarity of objects to other objects in their own cluster and the average dissimilarity of objects to the objects in a “neighbor cluster.”
- The larger the average silhouette width, the better the clustering of the objects.
- We could calculate the average silhouette width for clusterings based on several values of k and choose the k with the largest average silhouette width.
- The silhouette function in the cluster package of R gives the average silhouette width for any clustering result and distance matrix

2. Dunn index and the Davies-Bouldin Index: These are implemented with the `clv.Dunn` and `clv.Davies.Bouldin` functions in the `clv` package.

3. Within cluster sum-of-squares (WSS):

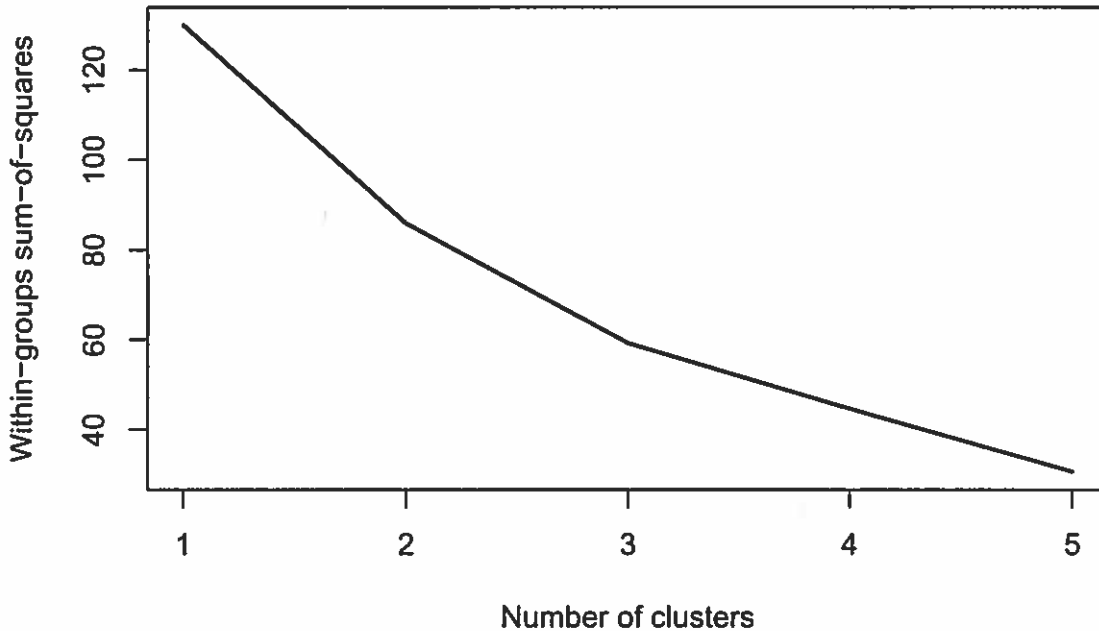
- A common way to choose k in K-means clustering
- As k increases, the corresponding WSS will decrease, and at some point will level off.
- The “best” choice of k usually occurs near the “elbow” in this plot.

In the `food.std` data, we use the WSS method to determine the optimal number of k clusters.

```

my.data.matrix = food.std
my.k.choices = 2:5
n = length(my.data.matrix[,1])
wss1 = (n-1)*sum(apply(my.data.matrix,2,var))
wss = numeric(0)
for(i in my.k.choices) {
W = sum(kmeans(my.data.matrix,i)$withinss)
wss = c(wss,W)
}
wss = c(wss1,wss)
plot(c(1,my.k.choices),wss,type='l',xlab='Number of clusters',
ylab='Within-groups sum-of-squares', lwd=2)

```



From the plot, we see that the optimal k is 3. We then recluster the data with $k = 3$.

```

food.k3= kmeans(food.std, centers=3, iter.max=100, nstart=25)
food.k3

## K-means clustering with 3 clusters of sizes 16, 2, 9
##
## Cluster means:
##   Energy Protein      Fat Calcium  Iron
## 1 1.546941 4.762805 0.63293763 0.7624857 1.313442
## 2 0.568138 2.116802 0.08883335 0.9995610 3.901820
## 3 3.271597 4.468804 2.44785237 0.1124862 1.680901
##
## Clustering vector:
## [1] 3 3 3 3 1 1 1 1 3 3 3 3 3 1 1 1 2 2 1 1 1 1 1 1 1 1
##
## Within cluster sum of squares by cluster:
## [1] 52.2876687 0.5626614 6.4094695
## (between_SS / total_SS = 54.4 %)
##

```

```
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"
## [5] "tot.withinss" "betweenss"    "size"         "iter"
## [9] "ifault"
```

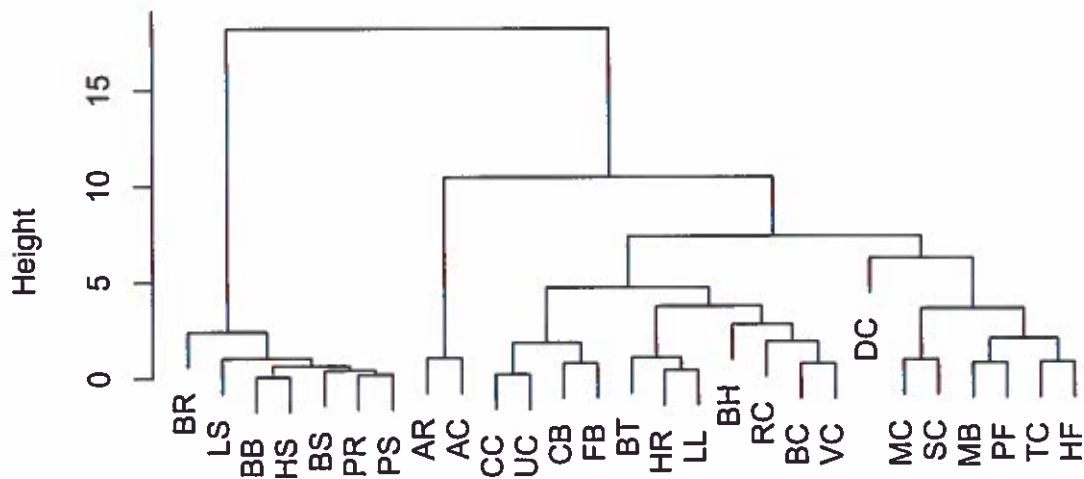
Ward's Method (Hierarchical clustering + K means)

The method of Ward (1963) is a hybrid of hierarchical clustering and K-means.

- It begins with n clusters and joins clusters together, one step at a time.
- At each step, the method searches over all possible ways to join a pair of clusters so that the K-means criterion WSS is minimized for that step.
- It begins with each object as its own cluster (so that $WSS = 0$) and concludes with all objects in one cluster.
- The R function `hclust` performs Ward's method if the option `method = 'ward'` is specified

```
food.ward = hclust(dist.food, method='ward.D')
plot(food.ward, labels=Food)
```

Cluster Dendrogram



```
dist.food
hclust (*, "ward.D")
```

Example: States data

```
##### NON-HIERARCHICAL (K MEANS) CLUSTERING #####
library(maps)
library(RColorBrewer) #Creates nice looking color palettes especially for thematic maps
```

```
head(state.x77)
```

```
##           Population Income Illiteracy Life Exp Murder HS Grad Frost
## Alabama      3615   3624         2.1   69.05   15.1   41.3   20
## Alaska        365   6315         1.5   69.31   11.3   66.7  152
## Arizona      2212   4530         1.8   70.55    7.8   58.1   15
## Arkansas     2110   3378         1.9   70.66   10.1   39.9   65
## California   21198  5114         1.1   71.71   10.3   62.6   20
## Colorado     2541   4884         0.7   72.06    6.8   63.9  166
##           Area
## Alabama      50708
## Alaska     566432
## Arizona     113417
## Arkansas     51945
## California  156361
## Colorado    103766
```

```
my.k.choices = 2:6
```

```
n = length(state.x77[,1])
```

```
wss1 = (n-1)*sum(apply(state.x77,2,var))
```

```
wss = numeric(0)
```

```
for(i in my.k.choices) {
```

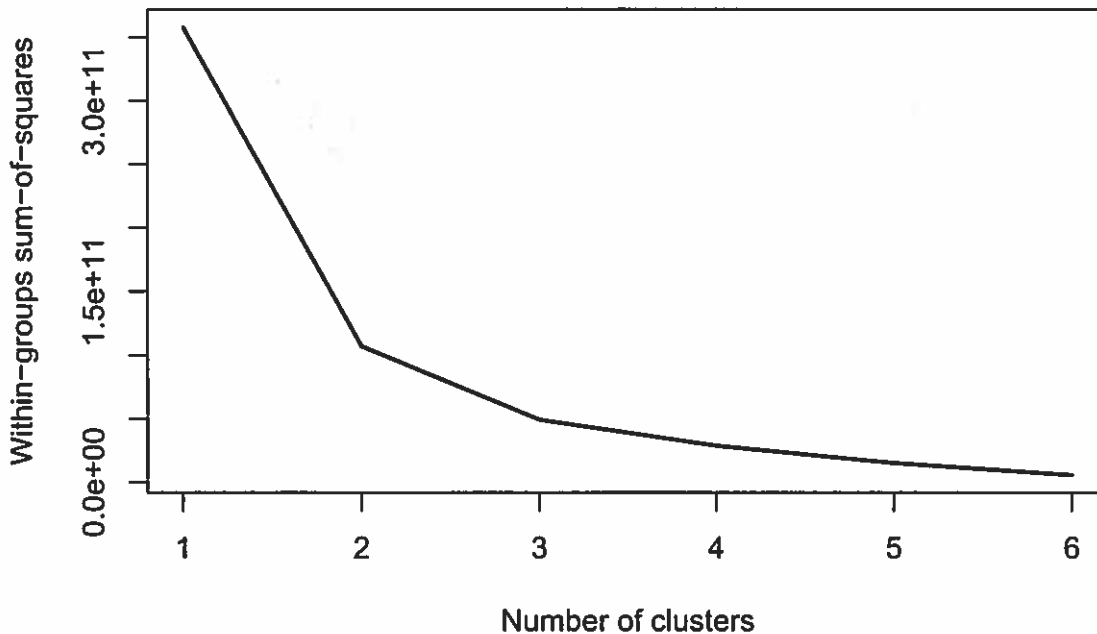
```
  W = sum(kmeans(state.x77,i)$withinss)
```

```
  wss = c(wss,W)
```

```
}
```

```
wss = c(wss1,wss)
```

```
plot(c(1,my.k.choices),wss,type='l',xlab='Number of clusters',
      ylab='Within-groups sum-of-squares', lwd=2)
```



```
km1list <- vector("list", 5)
```

```
for(k in 2:6){
```

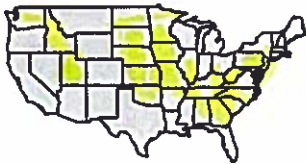
```
  set.seed(1)
```

similarity

```
kmlist[[k-1]] <- kmeans(state.x77, k, nstart=5000)
}

par(mfrow=c(2,2))
for(k in 3:6){
  map(database = "state")
  title(paste0("K=",k," Clusters: State Data"))
  cols=brewer.pal(k, "Paired")
  for(j in 1:k){
    ix=names(which(kmlist[[k-1]]$cluster==j))
    if(length(ix) > 1) map(database = "state", regions = ix, col = cols[j], fill=T, add=TRUE)
  }
}
} good by WSSS.
```

K=3 Clusters: State Data



K=4 Clusters: State Data



K=5 Clusters: State Data



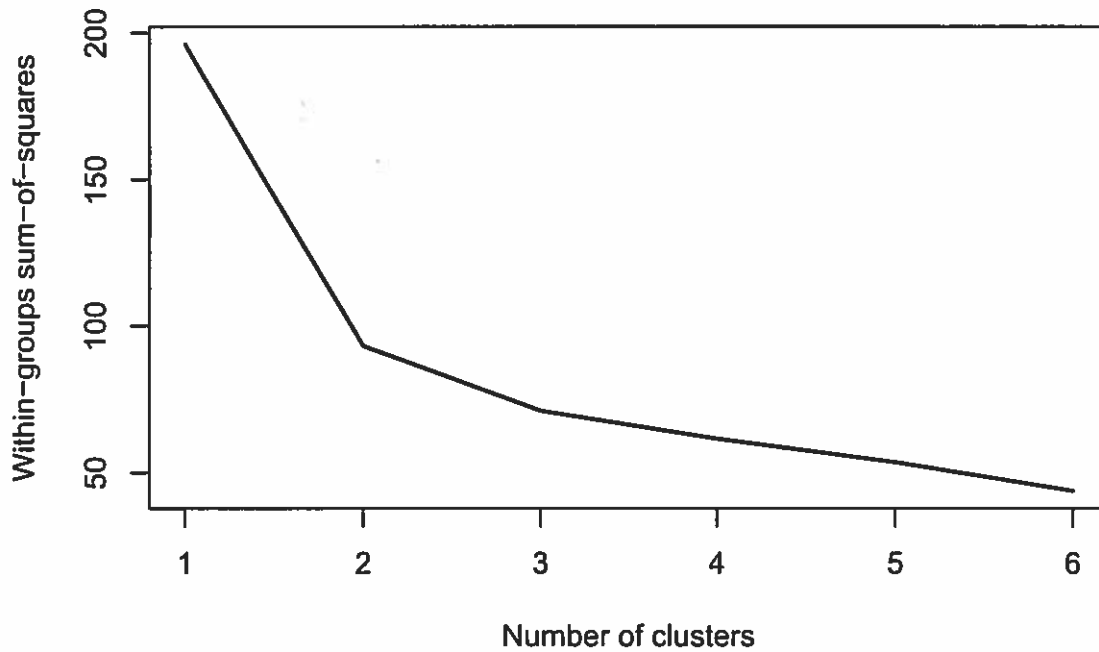
K=6 Clusters: State Data



Look at positive factors.

Look at variables "Income", "Illiteracy", "Life Exp", "HS Grad" only. We also standardized the data.

```
vars <- c("Income", "Illiteracy", "Life Exp", "HS Grad")
state.x77sub=scale(state.x77[,vars])
my.k.choices = 2:6
n = length(state.x77sub[,1])
wss1 = (n-1)*sum(apply(state.x77sub,2,var))
wss = numeric(0)
for(i in my.k.choices) {
  W = sum(kmeans(state.x77sub,i)$withinss)
  wss = c(wss,W)
}
wss = c(wss1,wss)
plot(c(1,my.k.choices),wss,type='l',xlab='Number of clusters',
      ylab='Within-groups sum-of-squares', lwd=2)
```



```

kmlist <- vector("list", 5)
for(k in 2:6){
  set.seed(1)
  kmlist[[k-1]] <- kmeans(state.x77sub, k, nstart=5000)
}

par(mfrow=c(2,2))
for(k in 3:6){
  map(database = "state")
  title(paste0("K= ",k," Clusters: State standardized Data"))
  cols=brewer.pal(k, "Paired")
  for(j in 1:k){
    ix=names(which(kmlist[[k-1]]$cluster==j))
    if(length(ix) > 1) map(database = "state", regions = ix, col = cols[j], fill=T, add=TRUE)
  }
}

```

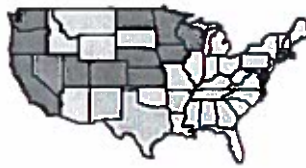
K= 3 Clusters: State standardized Data

K= 4 Clusters: State standardized Data



K= 5 Clusters: State standardized Data

K= 6 Clusters: State standardized Data



6. K-medoids Clustering

The K-medoids algorithm (Kaufman and Rousseeuw, 1987) is a robust alternative K-means.

- It attempts to minimize the criterion

$$\text{Crit}_{\text{Med}} = \sum_{c=1}^k \sum_{i \in c} d(x_i, m_c)$$

where m_c is a medoid, or "most representative object" for cluster c .

*not the centroid
← mode*

- The algorithm begins (in the "build step") by selecting k such representative objects.
- It proceeds by assigning each object to the cluster with the closest medoid.
- Then (in the "swap step"), if swapping any non-medoid object with a medoid results in a decrease in the criterion Crit_{Med} , the swap is made.
- The algorithm stops when no swap can decrease Crit_{Med} .

Like K-means, the K-medoids algorithm does not globally minimize its criterion in general. The R function `pam()` in the cluster package performs K-medoids clustering. An advantage of K-medoids is that (unlike `kmeans`) the function can accept a dissimilarity matrix, as well as a raw data matrix. This is because the criterion to be minimized is a direct sum of pairwise dissimilarities between objects. The `pam()` function also produces tools called the silhouette plot and average silhouette width to guide the choice of k .

? how to compute dissimilarity matrix in R.

Example: cars data.

Let's cluster the cars into k groups using the K-medoids approach. Here we will focus on the variables that are continuous in nature rather than discrete:

```
# Loading the "cluster" package:
library(cluster)
```

```

## Warning: package 'cluster' was built under R version 3.4.4
##
## Attaching package: 'cluster'
## The following object is masked from 'package:maps':
##
##   votes.repub
# Continuous variables
cars= mtcars[,c(1,3,4,5,6,7)]
# Standardizing by dividing through by the sample range of each variable
samp.range=function(x){
  myrange = diff(range(x))
  return(myrange)
}
my.ranges = apply(cars,2,samp.range)
cars.std = sweep(cars,2,my.ranges,FUN="/")
# K-medoids directly on the (standardized) data matrix:
cars.kmed.3=pam(cars.std, k=3, diss=F)
# Or you can do K-medoids by inputting the distance matrix:
#dist.cars=dist(cars,method = "euclidean", diag = FALSE)
#cars.kmed.3=pam(dist.cars, k=3, diss=T)
cars.kmed.3$clustering # printing the "clustering vector"

##           Mazda RX4      Mazda RX4 Wag      Datsun 710
##           1              1                  1
##   Hornet 4 Drive  Hornet Sportabout      Valiant
##           2              2                  2
##   Duster 360      Merc 240D      Merc 230
##           2              1                  1
##   Merc 280      Merc 280C      Merc 450SE
##           1              1                  2
##   Merc 450SL      Merc 450SLC  Cadillac Fleetwood
##           2              2                  3
## Lincoln Continental  Chrysler Imperial      Fiat 128
##           3              3                  1
##   Honda Civic      Toyota Corolla      Toyota Corona
##           1              1                  1
##   Dodge Challenger      AMC Javelin      Camaro Z28
##           2              2                  2
##   Pontiac Firebird      Fiat X1-9      Porsche 914-2
##           2              1                  1
##   Lotus Europa      Ford Pantera L      Ferrari Dino
##           1              2                  1
##   Maserati Bora      Volvo 142E
##           2              1

cars.kmed.3$silinfo$avg.width #printing the average silhouette width

## [1] 0.3980396
cars.3.clust=lapply(1:3, function(nc) row.names(cars)[cars.kmed.3$clustering==nc])
cars.3.clust # printing the clusters in terms of the car names
## [[1]]

```



```

iX4"      "Mazda RX4 Wag"  "Datsun 710"  "Merc 240D"
i0"       "Merc 280"      "Merc 280C"  "Fiat 128"
iivic"    "Toyota Corolla" "Toyota Corona" "Fiat X1-9"
: 914-2"  "Lotus Europa"  "Ferrari Dino" "Volvo 142E"

```

sportcars

```

## [1] "Hornet 4 Drive"      "Hornet Sportabout" "Valiant"
## [4] "Duster 360"          "Merc 450SE"        "Merc 450SL"
## [7] "Merc 450SLC"        "Dodge Challenger"  "AMC Javelin"
## [10] "Camaro Z28"         "Pontiac Firebird"  "Ford Pantera L"
## [13] "Maserati Bora"
##
## [[3]]
## [1] "Cadillac Fleetwood" "Lincoln Continental" "Chrysler Imperial"

```

highend cars.

From choosing $k = 3$, we see that Cluster 1 seems to be mostly compact cars, Cluster 2 is sports cars, Cluster 3 is large luxury sedans.

We write a function to calculate the average silhouette width and here is a variety of choices of k .

```

my.k.choices = 2:8
avg.sil.width= rep(0, times=length(my.k.choices))
for (ii in (1:length(my.k.choices))) {
  avg.sil.width[ii] <- pam(cars.std, k=my.k.choices[ii])$silinfo$avg.width
}
print( cbind(my.k.choices, avg.sil.width) )

```

```

##      my.k.choices avg.sil.width
## [1,]           2    0.4569631
## [2,]           3    0.3980396
## [3,]           4    0.3389684
## [4,]           5    0.4023427
## [5,]           6    0.3745073
## [6,]           7    0.3779733
## [7,]           8    0.3600124

```

A LARGE average silhouette width indicates that the observations are properly clustered. So maybe $k = 2$ is the best choice of k here?

How do we visualize the Clusters? There is built-in plots available with the pam() function

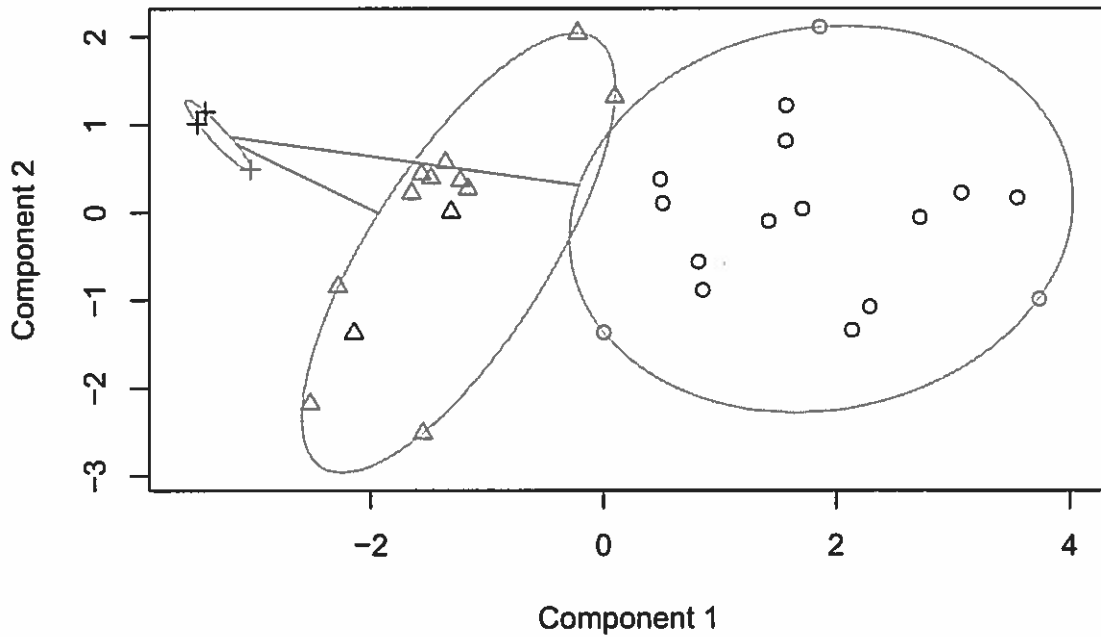
```

# The "clusplot":
plot(cars.kmed.3, which.plots=1)

```

↑
 put these people into clusters
 (# of clusters depends on # that we want)

```
clusplot(pam(x = cars.std, k = 3, diss = F))
```

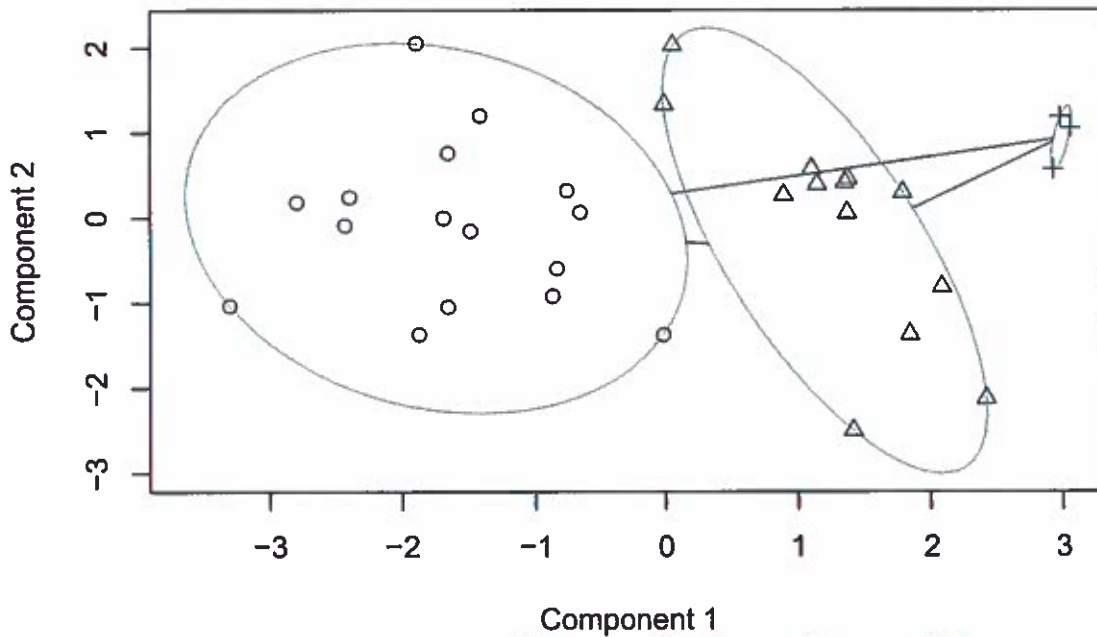


These two components explain 88.93 % of the point variability.

The `clusplot` (in the “cluster” library) can actually be used with any clustering partition by entering the data set and the clustering vector, e.g.:

```
clusplot(cars[, -1], cars.kmed.3$cluster)
```

CLUSPLOT(cars[, -1])



These two components explain 89.64 % of the point variability.

The "silhouette plot":

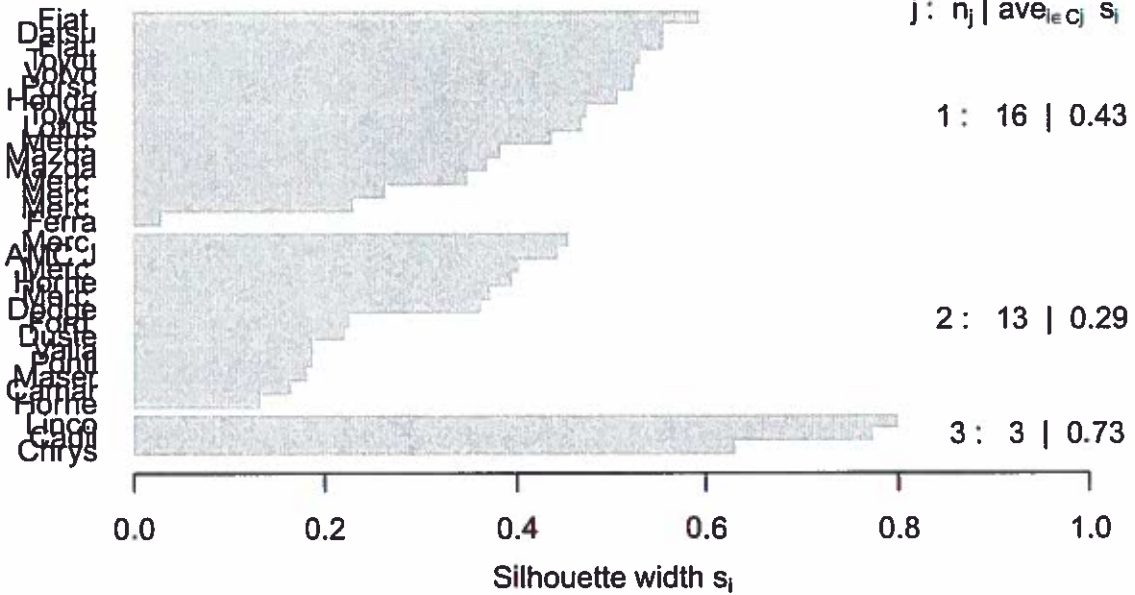
```
plot(cars.kmed.3, which.plots=2)
```

Silhouette plot of pam(x = cars.std, k = 3, diss = F)

n = 32

3 clusters C_j

$j: n_j | \text{ave}_{i \in C_j} s_i$



Average silhouette width : 0.4

3 Model-based clustering

This shows which observations are "best clustered".

7. Model-based Clustering

EM algorithm
Bayesian approach (it's useful to read this algorithm)

Neither hierarchical nor partitioning methods assume a specific statistical model for the data. They are strictly exploratory tools, and no formal inference about a wider population is possible. Model-based clustering assumes that the population generating the data consists of k subpopulations, which correspond to the k clusters. Therefore, the distribution for the data is assumed to be composed of k densities. This idea was originally proposed by Scott and Symons (1971) but fully developed in recent years by Banfield and Raftery (1993) and Fraley and Raftery (2002).

Clustering Model Setup: let $\gamma = [\gamma_1, \gamma_2, \dots, \gamma_n]^T$ be a vector of cluster labels, such that $\gamma_i = j$ if observation \mathbf{x}_i is from the j -th subpopulation. Suppose the subpopulation densities are denoted by $f_j(\mathbf{x}; \theta_j)$, where θ_j contains the set of unknown parameters for the j -th density. Then the likelihood, given the observed data, is:

$$L(\theta_1, \theta_2, \dots, \theta_k, \gamma \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \prod_{i=1}^n f_{\gamma_i}(\mathbf{x}_i; \theta_{\gamma_i})$$

Fitting the model amounts to choosing $\theta_1, \theta_2, \dots, \theta_k, \gamma$ to maximize this likelihood. The estimated γ is the "clustering vector" that defines which cluster each object is assigned to.

The Multivariate Normality Assumption:

Assume that each subpopulation ($j = 1, \dots, k$) follows a multivariate normal density having mean vectors μ_j and covariance matrices Σ_j , for $j = 1, \dots, k$, as its parameters. Then the likelihood becomes

$$L(\theta_1, \theta_2, \dots, \theta_k, \gamma) \propto \prod_{j=1}^k \prod_{i \in I_j} |\Sigma_j|^{1/2} \exp\left[-\frac{1}{2}(\mathbf{x}_i - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mu_j)\right].$$

The MLE of μ_j is $\bar{\mathbf{x}}_j$, the sample mean vector for the observations in subpopulation j . Replacing μ_j with $\bar{\mathbf{x}}_j$, the log-likelihood function is a constant plus

$$-\frac{1}{2} \sum_{j=1}^k \text{trace}(\mathbf{W}_j \Sigma_j^{-1}) + n \ln |\Sigma_j|,$$

where \mathbf{W}_j is a matrix containing the sums of squares and cross products of variables for observations in subpopulation j . We can assume a certain structure for the covariance matrices Σ_j for $j = 1, \dots, k$ and then determine computationally the value of γ that maximizes this (log) likelihood. \square

We consider a few possible covariance structures:

1. A simple (maybe unrealistic!) assumption is that each subpopulation has the same covariance structure and that all the $\Sigma_j = \sigma^2 \mathbf{I}$. In this case, γ is chosen so that the total within-group sum-of-squares $\text{trace}(\sum_{j=1}^k \mathbf{W}_j)$ is minimized. This tends to produce clusters that are spherical and roughly of equal size.
2. A slightly more complicated assumption is that each subpopulation has the same covariance structure, i.e., $\Sigma_j = \Sigma$ for all $j = 1, \dots, k$. This tends to produce clusters that are elliptical with roughly the same directional slope.
3. An extremely unrestrictive assumption is that each subpopulation may have a completely different covariance structure, $\Sigma_j, j = 1, \dots, k$. This may produce clusters that are different in size, shape, and orientation. We might consider assumptions that are less restrictive than the equal-covariances assumption yet more parsimonious than the unstructured-covariances assumption.

The covariance structure we assume leads to a clustering solution in which the sizes, shapes, and orientations of the clusters might be the same or different. In practice, the R function 'Mclust()' in the mclust package considers many such models, letting the covariance assumptions and the number of clusters k vary. Usually the Bayesian information Criterion (BIC) is used to choose the best of all these competing models and thus determine the model-based clustering result.

Example: USArrests data

In this example, we will perform a model-based clustering of the USArrests data. The data contains 50 states based on these 4 variables. Use

```
help(USArrests)
```

to see detail.

The R function 'Mclust()' performs model-based clustering for a range of models and a variety of values of k :

```
library(mclust)
```

```
## Warning: package 'mclust' was built under R version 3.4.3
```

```
## Package 'mclust' version 5.4
```

```
## Type 'citation("mclust")' for citing this R package in publications.
```

```
##
```

```
## Attaching package: 'mclust'
```

```
## The following object is masked from 'package:maps':
```

```
##
```

```
## map
```

```
arrest.clus <- Mclust(USArrests)
```

$G = 2$ ← or what ever number that we believe that it is the true # of clusters.

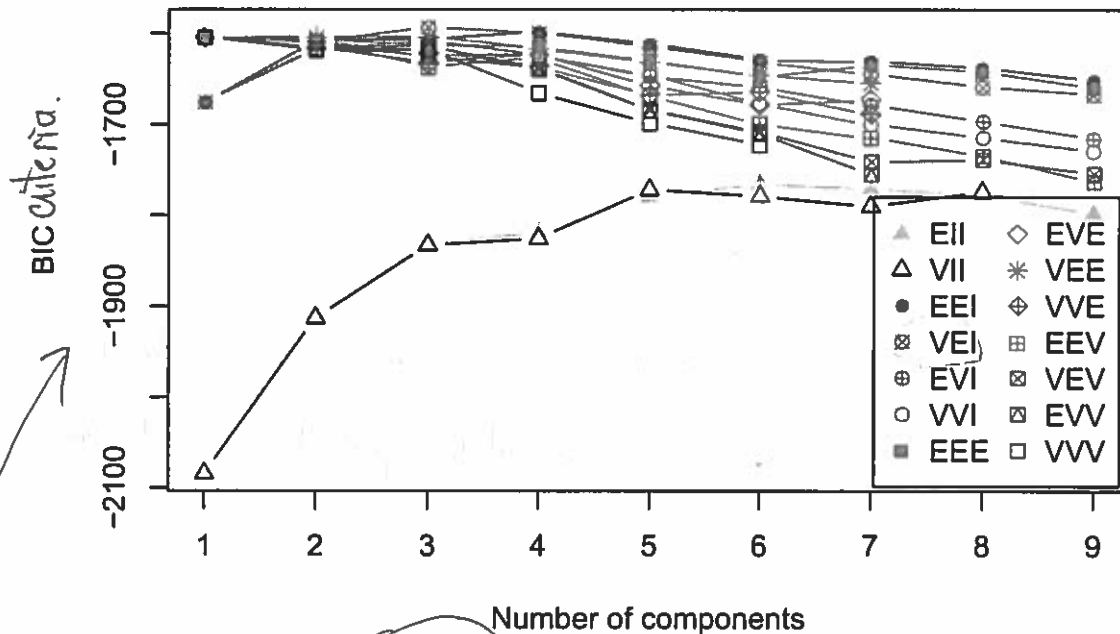
By default, the models considered are: - "EII": spherical, equal volume

- "VII": spherical, unequal volume
- "EEI": diagonal, equal volume and shape
- "VEI": diagonal, varying volume, equal shape
- "EVI": diagonal, equal volume, varying shape
- "VVI": diagonal, varying volume and shape
- "EEE": ellipsoidal, equal volume, shape, and orientation
- "EEV": ellipsoidal, equal volume and equal shape
- "VEV": ellipsoidal, equal shape
- "VVV": ellipsoidal, varying volume, shape, and orientation

We plot the BIC values: (Bayesian information criterion)

```
plot(arrest.clus, data=USArrests, what="BIC")
```

↑
BIC: Bayesian information criterion.



The best solution is VEI with 3 clusters.
 summary(arrest.clus)

can't determine # of clusters.

```
## -----
## Gaussian finite mixture model fitted by EM algorithm
## -----
##
## Mclust VEI (diagonal, equal shape) model with 3 components:
##
## log.likelihood n df      BIC      ICL
##      -757.5594 50 20 -1593.359 -1597.95
##
## Clustering table
## 1 2 3
## 20 10 20
```

To see the clustering vector:

```
clus.vec.3 <- arrest.clus$classification
arrest.3.clus <- lapply(1:3, function(nc) row.names(USArrests)[clus.vec.3==nc])
arrest.3.clus
```

```
## [[1]]
## [1] "Alabama"      "Alaska"      "Arizona"     "California"
## [5] "Colorado"     "Florida"     "Georgia"     "Illinois"
## [9] "Louisiana"    "Maryland"    "Michigan"    "Mississippi"
## [13] "Missouri"     "Nevada"      "New Mexico"  "New York"
## [17] "North Carolina" "South Carolina" "Tennessee"  "Texas"
##
## [[2]]
## [1] "Idaho"      "Iowa"      "Maine"      "Minnesota"
## [5] "New Hampshire" "North Dakota" "South Dakota" "Vermont"
## [9] "West Virginia" "Wisconsin"
```

```
##
## [[3]]
## [1] "Arkansas"      "Connecticut"  "Delaware"     "Hawaii"
## [5] "Indiana"       "Kansas"       "Kentucky"     "Massachusetts"
## [9] "Montana"       "Nebraska"     "New Jersey"   "Ohio"
## [13] "Oklahoma"     "Oregon"       "Pennsylvania" "Rhode Island"
## [17] "Utah"         "Virginia"     "Washington"   "Wyoming"
```

This gives the probabilities of belonging to each cluster for every object:

```
round(arrest.clus$z,2)
```

```
##           [,1] [,2] [,3]
## Alabama   1.00 0.00 0.00
## Alaska    1.00 0.00 0.00
## Arizona   1.00 0.00 0.00
## Arkansas  0.41 0.00 0.59
## California 1.00 0.00 0.00
## Colorado  1.00 0.00 0.00
## Connecticut 0.00 0.02 0.98
## Delaware  0.37 0.00 0.63
## Florida   1.00 0.00 0.00
## Georgia   1.00 0.00 0.00
## Hawaii    0.00 0.00 1.00
## Idaho     0.00 0.64 0.36
## Illinois  1.00 0.00 0.00
## Indiana   0.00 0.00 1.00
## Iowa      0.00 1.00 0.00
## Kansas    0.00 0.00 1.00
## Kentucky  0.01 0.00 0.99
## Louisiana 1.00 0.00 0.00
## Maine     0.00 1.00 0.00
## Maryland  1.00 0.00 0.00
## Massachusetts 0.00 0.00 1.00
## Michigan  1.00 0.00 0.00
## Minnesota 0.00 0.85 0.15
## Mississippi 1.00 0.00 0.00
## Missouri  0.66 0.00 0.34
## Montana   0.00 0.02 0.98
## Nebraska  0.00 0.10 0.90
## Nevada    1.00 0.00 0.00
## New Hampshire 0.00 1.00 0.00
## New Jersey 0.01 0.00 0.99
## New Mexico 1.00 0.00 0.00
## New York  1.00 0.00 0.00
## North Carolina 1.00 0.00 0.00
## North Dakota 0.00 1.00 0.00
## Ohio      0.00 0.00 1.00
## Oklahoma  0.00 0.00 1.00
## Oregon    0.01 0.00 0.99
## Pennsylvania 0.00 0.00 1.00
## Rhode Island 0.00 0.00 1.00
## South Carolina 1.00 0.00 0.00
## South Dakota 0.00 0.99 0.01
## Tennessee 1.00 0.00 0.00
```

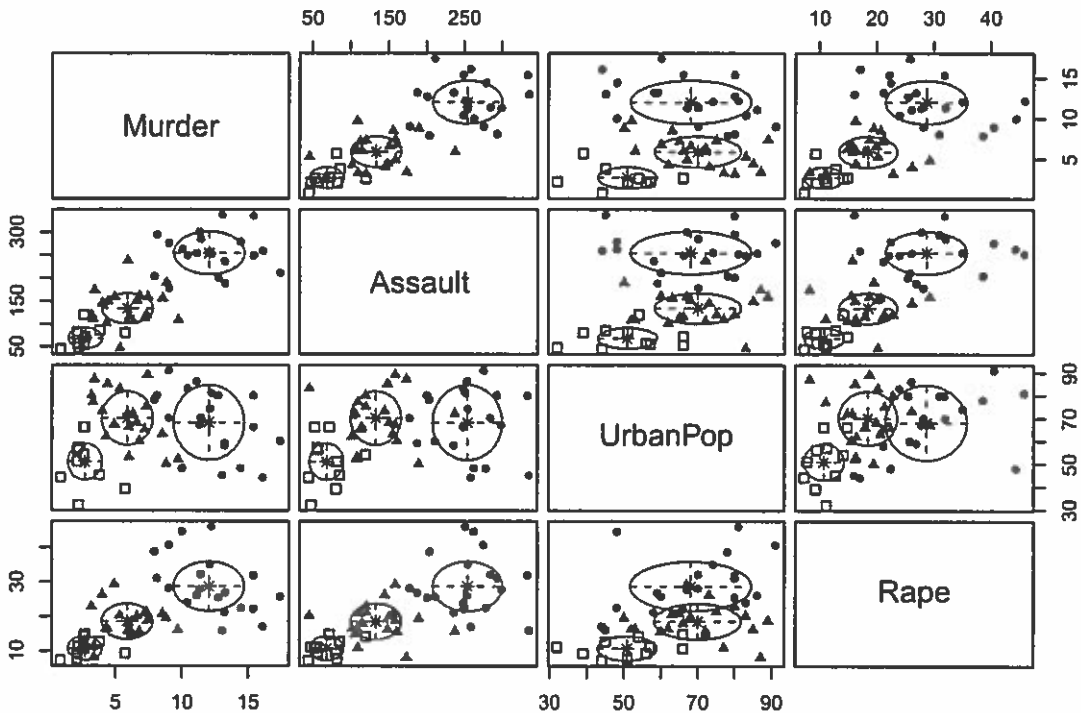
conclude that Alabama is in the 1st cluster.

It is hard to say which cluster Arkansas is in.

```
## Texas      1.00 0.00 0.00
## Utah       0.00 0.00 1.00
## Vermont    0.00 1.00 0.00
## Virginia   0.02 0.00 0.98
## Washington 0.00 0.00 1.00
## West Virginia 0.00 0.96 0.04
## Wisconsin  0.00 0.99 0.01
## Wyoming    0.00 0.00 1.00
```

To visualize the clusters via a scatterplot matrix:

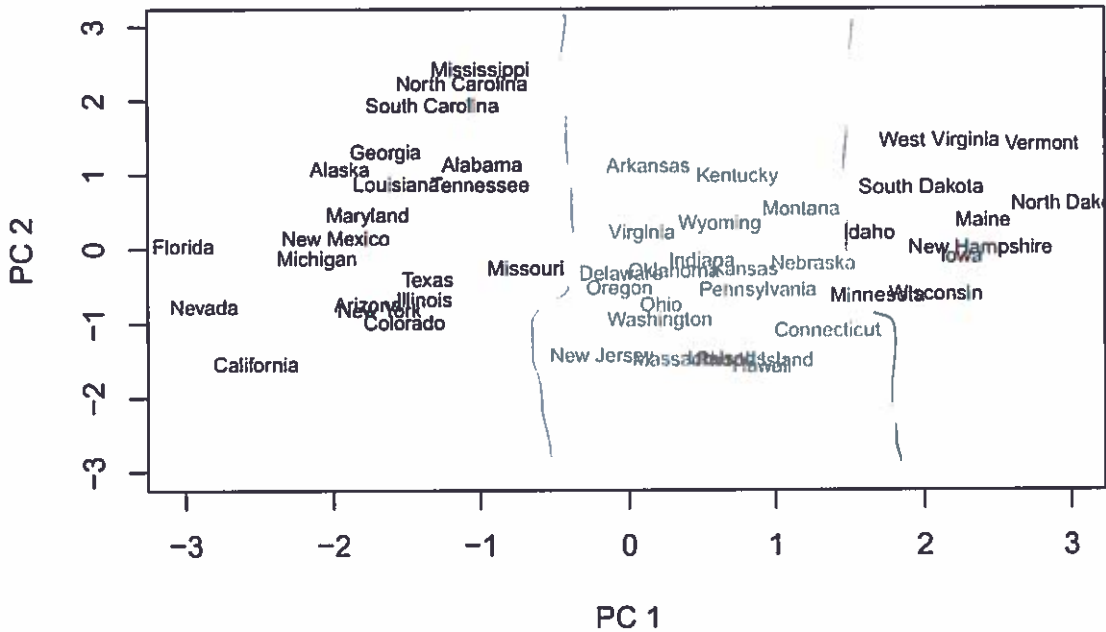
```
plot(arrest.clus, what="classification")
```



To visualize the clusters via a plot of the scores on the first 2 principal components, with the clusters separated by color:

```
arrests.pc <- princomp(USArrests, cor=T)
# Setting up the colors for the 5 clusters on the plot:
my.color.vector <- rep("blue", times=nrow(USArrests))
my.color.vector[arrest.clus$classification==2] <- "red"
my.color.vector[arrest.clus$classification==3] <- "green"

# Plotting the PC scores:
plot(arrests.pc$scores[,1], arrests.pc$scores[,2], ylim=range(arrests.pc$scores[,1]),
      xlab="PC 1", ylab="PC 2", type='n', lwd=2)
text(arrests.pc$scores[,1], arrests.pc$scores[,2], labels=row.names(USArrests),
      cex=0.7, lwd=2, col=my.color.vector)
```

To review the PCA:

```
summary(arrests.pc,loadings=T)
```

```
## Importance of components:
##                Comp.1  Comp.2  Comp.3  Comp.4
## Standard deviation  1.5748783 0.9948694 0.5971291 0.41644938
## Proportion of Variance 0.6200604 0.2474413 0.0891408 0.04335752
## Cumulative Proportion 0.6200604 0.8675017 0.9566425 1.00000000
##
## Loadings:
##      Comp.1  Comp.2  Comp.3  Comp.4
## Murder  -0.536  0.418 -0.341  0.649
## Assault  -0.583  0.188 -0.268 -0.743
## UrbanPop -0.278 -0.873 -0.378  0.134
## Rape     -0.543 -0.167  0.818
```

Note PC1 is an overall "lack-of-crime" index and PC2 is a "rural" index.

*lack of
crime because
we have negative
rjn here.*

Chapter 9. Multidimensional Scaling

Jianxuan Liu

Fall 2018

This chapter discuss methods for displaying multivariate data in low-dimensional space. The objective is to “fit” the original data into a low-dimensional coordinate system such that any distortion caused by a reduction in dimensionality is minimized

Distortion generally refers to the similarities or dissimilarities (distances) among the original data points.

Specifically, multidimensional scaling techniques deal with the following problem: For a set of observed similarities or distances between every pair of n items, find a representation of the items in few dimensions such that the interitem proximities “nearly match” the original similarities or distances. Multidimensional Scaling can be viewed as a way of generating a geometric representation of some observed proximity matrix.

The proximity matrix could contain similarity values for pairs of observations or dissimilarity values, but we will typically work with dissimilarities (i.e., distances).

With multidimensional scaling, we begin with a distance matrix and produce a “possible data set” that could have yielded such a distance matrix. Note that MDS can work on either objects or variables.

Example: Sports rankings data

```
source("Sports.txt")
head(sportsranks)

## Baseball Football Basketball Tennis Cycling Swimming Jogging
##      1      3      7      2      4      5      6
##      1      3      2      5      4      7      6
##      1      3      2      5      4      7      6
##      4      7      3      1      5      6      2
##      2      3      1      7      6      5      4
##      6      3      4      7      2      1      5

dim(sportsranks)

## [1] 130  7

cor(sportsranks)

##           Baseball  Football  Basketball  Tennis  Cycling
## Baseball  1.0000000  0.2816173  0.06225502 -0.25624133 -0.39268842
## Football  0.28161733  1.0000000  0.22412213 -0.27650958 -0.48593004
## Basketball 0.06225502  0.2241221  1.00000000 -0.18105420 -0.35726588
## Tennis    -0.25624133 -0.2765096 -0.18105420  1.00000000 -0.09455839
## Cycling   -0.39268842 -0.4859300 -0.35726588 -0.09455839  1.00000000
## Swimming  -0.49738160 -0.4740188 -0.35368023 -0.04047178  0.35106904
## Jogging   -0.28615216 -0.3892162 -0.33632916 -0.11179076  0.07558857
##           Swimming  Jogging
## Baseball -0.49738160 -0.28615216
## Football -0.47401883 -0.38921615
## Basketball -0.35368023 -0.33632916
## Tennis    -0.04047178 -0.11179076
```

```
## Cycling      0.35106904  0.07558857
## Swimming    1.00000000  0.05342708
## Jogging     0.05342708  1.00000000
```

D=dist(t(sportsranks))

Here we have

$$d_{ij} = \| \mathbf{z}_{(i)} - \mathbf{z}_{(j)} \|$$

for $1 \leq i, j \leq 7$ where $\mathbf{z}_{(1)}, \dots, \mathbf{z}_{(7)} \in \mathcal{R}^{130}$. We want to find $\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(7)} \in \mathcal{R}^2$, that is, $k = 2$ such that

$$d_{ij} = \| \mathbf{z}_{(i)} - \mathbf{z}_{(j)} \| \approx \| \mathbf{x}_i - \mathbf{x}_j \|$$

Here the traditional idea of objects and variables have been reversed. The $n = 7$ objects are different sport activities, and the $p = 130$ measurements taken on each sport are simply 130 individual rankings.

MDS is a family of different algorithms, each designed to arrive at some optimal low-dimensional configuration (usually $k = 2$ or $k = 3$)

With multidimensional scaling methods include

1. Classical MDS
2. Metric MDS
3. Non-metric MDS

1. Classical MDS

Suppose we have Euclidean distance matrix $\mathbf{D} = (d_{ij})$. The objective of classical multidimensional scaling is to find $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ such that

$$\| \mathbf{x}_i - \mathbf{x}_j \| = d_{ij}$$

. Such a solution is not unique, because if \mathbf{X} is the solution, then $\mathbf{x}_i^* = \mathbf{x}_i + \mathbf{c}$, $\mathbf{c} \in \mathcal{R}^k$ also satisfies

$$\| \mathbf{x}_i^* - \mathbf{x}_j^* \| = \| (\mathbf{x}_i + \mathbf{c}) - (\mathbf{x}_j + \mathbf{c}) \| = d_{ij}$$

Any location \mathbf{c} can be used, but the assumption of centered configuration, i.e.,

$$\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$$

serves well for the purpose of dimension reduction. The classical MDS finds the centered configuration $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{R}^k$ for some $k \geq n - 1$ so that their pairwise distances are the same as those corresponding distances in \mathbf{D} . We may find the $n \times n$ Gram matrix $\mathbf{B} = \mathbf{X}^T \mathbf{X}$ rather than \mathbf{X} . The Gram matrix is the matrix of inner products. Denote the ij^{th} element of \mathbf{B} as b_{ij} .

Theorem: Let $[\mathbf{x}_1, \dots, \mathbf{x}_n]$ be centered data, $d_{ij}^2 = (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j)$, and $\mathbf{B} = \mathbf{X}^T \mathbf{X}$.

1. Then the $n \times n$ matrix \mathbf{B} has entries

$$\begin{aligned} b_{ij} &= -\frac{1}{2} \left(d_{ij}^2 - \frac{2}{n} \sum_{i=1}^n d_{ij}^2 + \frac{1}{n^2} \sum_{i,j=1}^n d_{ij}^2 \right) \\ &= -\frac{1}{2} (d_{ij}^2 - d_j^2 - d_i^2 + d_{..}^2) \end{aligned}$$

where

- d_j^2 is the average of $\{d_{ij} : i = 1, \dots, n\}$ for each j
- d_i^2 is the average of $\{d_{ij} : j = 1, \dots, n\}$ for each i
- d^2 is the average of $\{d_{ij} : i, j = 1, \dots, n\}$

2. If $\mathbf{B} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ is the spectral decomposition of \mathbf{B} , then $\mathbf{X} = \mathbf{\Lambda}^{1/2}\mathbf{V}^T$ defines a configuration in k variables which minimize the STRESS.

$\mathbf{X} = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n]$ be centered data $d_{ij}^2 = (\underline{x}_i - \underline{x}_j)^T(\underline{x}_i - \underline{x}_j)$ $\mathbf{B} = \mathbf{X}^T\mathbf{X}$

* Note that $d_{ij}^2 = \underline{x}_i^T \underline{x}_i + \underline{x}_j^T \underline{x}_j - 2 \underline{x}_i^T \underline{x}_j = d_{ji}^2$ (0)

So we have $\frac{1}{n} \sum_{j=1}^n d_{ij}^2 = \underline{x}_i^T \underline{x}_i + \frac{1}{n} \sum_{j=1}^n \underline{x}_j^T \underline{x}_j \Rightarrow \underline{x}_i^T \underline{x}_i = \frac{1}{n} \sum_{j=1}^n d_{ij}^2 - \frac{1}{n} \sum_{j=1}^n \underline{x}_j^T \underline{x}_j$

$$\Rightarrow \frac{1}{n^2} \sum_{i,j=1}^n d_{ij}^2 = \frac{2}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i \quad (2)$$

* Consider $\mathbf{B} = \mathbf{X}^T\mathbf{X}$

So we have $b_{ij} = \underline{x}_i^T \underline{x}_j$ (1)

$$\begin{aligned} \text{So we have } b_{ij} &= \frac{1}{2} \left(d_{ij}^2 - \underline{x}_i^T \underline{x}_i - \underline{x}_j^T \underline{x}_j \right) - \frac{1}{n} \sum_{j=1}^n d_{ij}^2 \\ &= \frac{1}{2} \left(d_{ij}^2 - \underbrace{\frac{1}{n} \sum_{j=1}^n d_{ij}^2} + \underbrace{\frac{1}{n} \sum_{j=1}^n \underline{x}_j^T \underline{x}_j} + \frac{1}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i \right) \\ &= \frac{1}{2} \left(d_{ij}^2 - \frac{2}{n} \sum_{j=1}^n d_{ij}^2 + \frac{1}{n} \sum_{i,j=1}^n d_{ij}^2 \right) \end{aligned}$$

The space which \mathbf{X} lies is the eigenspace where the first coordinate contains the largest variation, and is identified with \mathcal{R}^k . If we wish to reduce the dimension to $k \leq p$, then the first k rows of \mathbf{X}_k best preserves the distances d_{ij} among all other linear dimension reduction of \mathbf{X} (to k). Then

$$\mathbf{X}_{(k)} = \mathbf{\Lambda}_k^{1/2} \mathbf{V}_k^T$$

where $\mathbf{\Lambda}_k$ is the first $k \times k$ sub matrix of $\mathbf{\Lambda}$, \mathbf{V}_k is collected through the first k columns of \mathbf{V} .

Classical MDS gives configurations $\mathbf{X}_{(k)}$ in \mathcal{R}^k for any dimension $k = 1, \dots, p$. The coordinates are given by the principal order of largest-to-smallest variances. Dimension reduction from $\mathbf{X} = \mathbf{X}_{(p)}$ to $\mathbf{X}_{(k)}$, $k < p$ is the same as PCA.

We can reduce the dimension of the solution by restricting attention to the k largest eigenvalues. If the distances are not Euclidean, \mathbf{B} is not positive definite and some eigenvalues of \mathbf{B} will be negative. In this case, we can still choose the dimension corresponding to the k largest positive eigenvalues.

Determining the Amount of Data Reduction

When using MDS to "reduce the dimensionality" from p to k , what is a proper choice of k ?

If there are k "relatively large" eigenvalues of \mathbf{B} , this is evidence that a k -dimensional solution is appropriate. We could base the choice of k on the sizes of the first few eigenvalues $\lambda_1, \lambda_2, \dots$ (listed in a decreasing order). We calculate (for each possible k)

$$P_k = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n |\lambda_i|}$$

Values of k that yield a P_k near 1 would give a good representation. The 'cmdscale' function in R prints this criterion when the "eig=T" option is specified.

There is another option in choosing k . For each possible value of k , we minimize

$$\phi = \sum_{i,j} (d_{ij}^2 - \hat{d}_{ij}^2)$$

where d_{ij}^2 is the Euclidean distance between the i -th and j -th observations in the full p -dimensional space, and \hat{d}_{ij}^2 is the Euclidean distance between the i -th and j -th observations in the reduced k -dimensional space. As k increases, the minimum value of ϕ will decrease monotonically. We can plot this minimum against the various values of k and pick the k value at the "elbow" of the plot. Takane et al. (1977) suggested a scaled version of ϕ called SStress that always lies between 0 and 1:

$$SStress = \left[\frac{\sum_{i<j} (d_{ij}^2 - \hat{d}_{ij}^2)^2}{\sum_{i<j} d_{ij}^4} \right]^{1/2}$$

Values of SStress below 0.1 represent a good fit.

Examples: European distances

```
# Classical MDS using the cmdscale function:
# The default number of dimensions is k = 2:
class(eurodist)
```

```
## [1] "dist"
```

```
euro.mds.2 <- cmdscale(eurodist, eig=T)
euro.mds.2
```

```
## $points
##           [,1]      [,2]
## Athens    2290.274680 1798.80293
## Barcelona -825.382790  546.81148
## Brussels   59.183341 -367.08135
## Calais     -82.845973 -429.91466
## Cherbourg -352.499435 -290.90843
## Cologne   293.689633 -405.31194
## Copenhagen 681.931545 -1108.64478
## Geneva     -9.423364  240.40600
## Gibraltar -2048.449113  642.45854
## Hamburg    561.108970 -773.36929
## Hook of Holland 164.921799 -549.36704
## Lisbon    -1935.040811  49.12514
## Lyons     -226.423236  187.08779
## Madrid    -1423.353697  305.87513
## Marseilles -299.498710  388.80726
## Milan      260.878046  416.67381
## Munich     587.675679  81.18224
```

```

## Paris          -156.836257 -211.13911
## Rome           709.413282  1109.36665
## Stockholm      839.445911 -1836.79055
## Vienna         911.230500  205.93020
##
## $eig
## [1] 1.953838e+07 1.185656e+07 1.528844e+06 1.118742e+06 7.893472e+05
## [6] 5.816552e+05 2.623192e+05 1.925976e+05 1.450845e+05 1.079673e+05
## [11] 5.139484e+04 -3.259629e-09 -9.496124e+03 -5.305820e+04 -1.322166e+05
## [16] -2.573360e+05 -3.326719e+05 -5.162523e+05 -9.191491e+05 -1.006504e+06
## [21] -2.251844e+06
##
## $x
## NULL
##
## $ac
## [1] 0
##
## $GOF
## [1] 0.7537543 0.8679134

```

The first number in the GOF section is what we called P_k . Now change the number of dimensions k :

```

euro.mds.4 <- cmdscale(eurodist, k=4, eig=T)
euro.mds.4

```

```

## $points
##           [,1]      [,2]      [,3]      [,4]
## Athens    2290.274680 1798.80293  53.79314 -103.826958
## Barcelona -825.382790  546.81148 -113.85842  84.585831
## Brussels  59.183341  -367.08135 177.55291  38.797514
## Calais    -82.845973 -429.91466 300.19274 106.353695
## Cherbourg -352.499435 -290.90843 457.35294 111.449150
## Cologne   293.689633 -405.31194 360.09323 -636.202379
## Copenhagen 681.931545 -1108.64478  26.09257 151.693056
## Geneva    -9.423364  240.40600 -344.20659 656.121110
## Gibraltar -2048.449113 642.45854 167.86631  78.621423
## Hamburg   561.108970  -773.36929  80.91722  48.548472
## Hook of Holland 164.921799 -549.36704 270.82327 116.886334
## Lisbon    -1935.040811  49.12514 -483.02056 -315.241752
## Lyons     -226.423236 187.08779 -358.43234 -257.737009
## Madrid    -1423.353697 305.87513 253.26763  2.478812
## Marseilles -299.498710 388.80726 -109.17417 12.651217
## Milan     260.878046 416.67381 -171.52428 20.926369
## Munich    587.675679  81.18224 -75.88485 13.080496
## Paris     -156.836257 -211.13911 131.30852 27.089432
## Rome       709.413282 1109.36665 -179.83052 -109.895049
## Stockholm 839.445911 -1836.79055 -541.35188 -108.755016
## Vienna    911.230500  205.93020  98.02313 62.375253
##
## $eig
## [1] 1.953838e+07 1.185656e+07 1.528844e+06 1.118742e+06 7.893472e+05
## [6] 5.816552e+05 2.623192e+05 1.925976e+05 1.450845e+05 1.079673e+05
## [11] 5.139484e+04 -3.259629e-09 -9.496124e+03 -5.305820e+04 -1.322166e+05
## [16] -2.573360e+05 -3.326719e+05 -5.162523e+05 -9.191491e+05 -1.006504e+06

```

```
## [21] -2.251844e+06
##
## $x
## NULL
##
## $ac
## [1] 0
##
## $GOF
## [1] 0.8173197 0.9411060
```

Using the P_k criterion for a variety of values of k to choose the appropriate amount of dimension reduction:

```
# consider values of k from 1 to 12
max.k <- 12
Pk <- rep(0,max.k)
SStress <- rep(0,max.k)
for (kk in 1:max.k){
my.mds.kk <- cmdscale(eurodist,k=kk,eig=T)
Pk[kk] <- my.mds.kk$GOF[1]
#SStress[kk] <- ( sum( eurodist^2 - (dist(my.mds.kk$points))^2 )/sum(eurodist^4) )^0.5
}
```

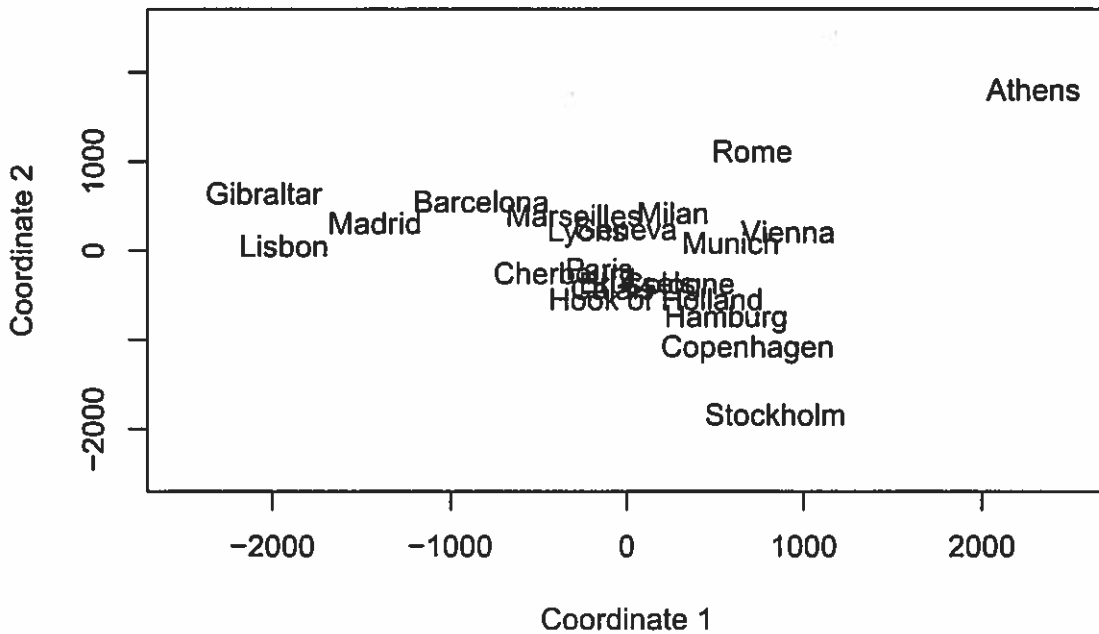
```
## Warning in cmdscale(eurodist, k = kk, eig = T): only 11 of the first 12
## eigenvalues are > 0
```

```
cbind(1:max.k,Pk)
```

```
##           Pk
## [1,]  1 0.4690928
## [2,]  2 0.7537543
## [3,]  3 0.7904600
## [4,]  4 0.8173197
## [5,]  5 0.8362709
## [6,]  6 0.8502358
## [7,]  7 0.8565337
## [8,]  8 0.8611578
## [9,]  9 0.8646411
## [10,] 10 0.8672332
## [11,] 11 0.8684672
## [12,] 12 0.8684672
```

It looks like 2 or 3 dimensions would be reasonable, and 4 gives quite a good fit. A 2-D representation of the solution for $k=2$:

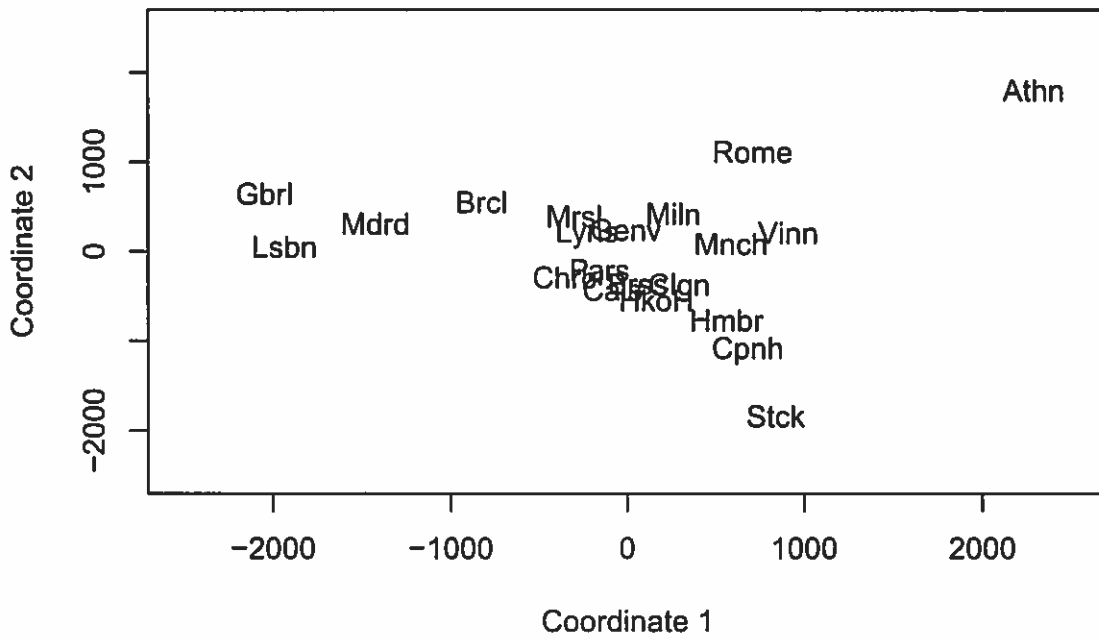
```
plot(euro.mds.2$points[,1], euro.mds.2$points[,2], type='n',
      xlab="Coordinate 1", ylab="Coordinate 2", xlim=c(-2500,2500), ylim=c(-2500,2500) )
text(euro.mds.2$points[,1], euro.mds.2$points[,2], labels=labels(eurodist) )
```



use abbreviations:

```
euro.abb <- abbreviate(labels(eurodist))
```

```
plot(euro.mds.2$points[,1], euro.mds.2$points[,2], type='n',
      xlab="Coordinate 1", ylab="Coordinate 2", xlim=c(-2500,2500), ylim=c(-2500,2500) )
text(euro.mds.2$points[,1], euro.mds.2$points[,2], labels=euro.abb )
```

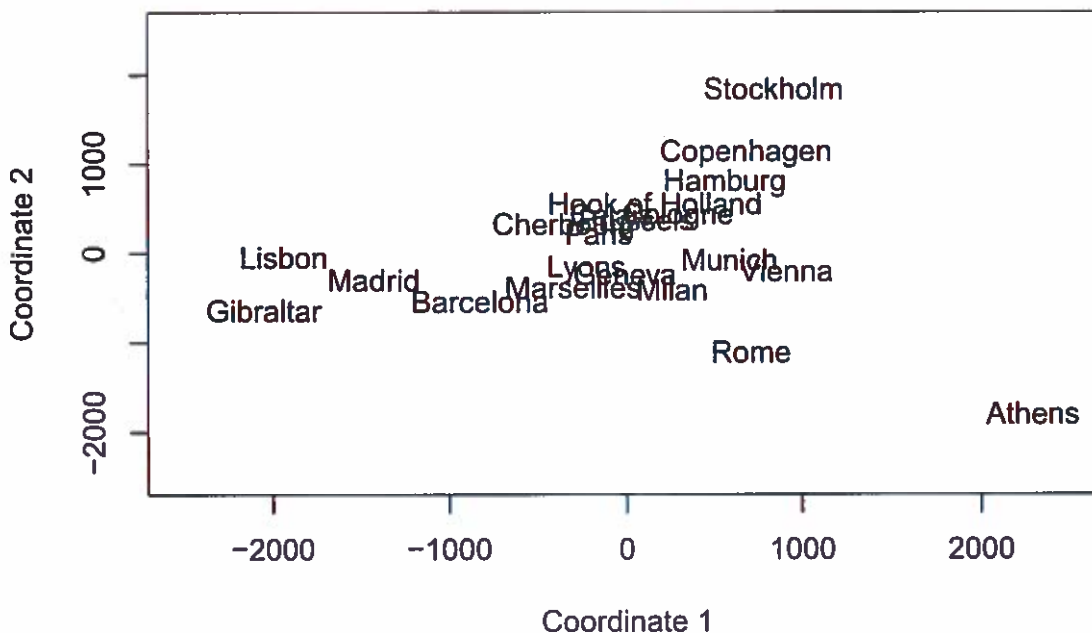


Maybe rotating across the x-axis would produce a better reflection of reality

```
plot(euro.mds.2$points[,1], -euro.mds.2$points[,2], type='n',
```



```
xlab="Coordinate 1", ylab="Coordinate 2", xlim=c(-2500,2500), ylim=c(-2500,2500) )
text(euro.mds.2$points[,1], -euro.mds.2$points[,2], labels=labels(eurodist) )
```



Summary:

- Classical MDS gives configurations $\mathbf{X}_{(k)}$ in \mathcal{R}^k for any dimension $1 \leq k \leq p$
- Configuration is centered
- Coordinates are given by the principal scores, ordered from largest-to-smallest variation.
- Dimension reduction from $\mathbf{X} = \mathbf{X}_{(k)}$ to $\mathbf{X}_{(k)}$ ($k < p$) is same as PCA (cutting some PC scores out).
- Leads to exact solution if the dissimilarity is based on Euclidean distances
- Can also be used for non-Euclidean distances, in fact, for any dissimilarities.

2. Metric MDS

In metric MDS, the dissimilarities $\{d_{ij}\}$ are quantitative measurements, usually Euclidean, but other distance metrics are possible. The function f is usually taken to be a parametric monotonic function, such as $f(d_{ij}) = \alpha + \beta d_{ij}$, where α and β are unknown positive coefficients. Since f is a parametric function, the distances $\{d_{ij}\}$ can be fitted to $\{f(d_{ij})\}$ by least squares method which results in metric LS scaling. If the dissimilarities are Euclidean distances and f is taken to be the identity function, then classical MDS can be viewed as an example of metric LS scaling.

Given a (low) dimension k and a monotone function f , metric MDS seeks to find an optimal configuration $\mathbf{X} \in \mathcal{R}^k$ that gives

$$f(d_{ij}) \approx \hat{d}_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|^2$$

as close as possible.

- The function f can be taken to be a parametric monotonic function, such as $f(d_{ij}) = \alpha + \beta d_{ij}$

- "As close as possible" is now explicitly stated by square loss

$$stress = \mathcal{L}(\hat{d}_{ij}) = \left(\frac{\sum_{i < j} [\hat{d}_{ij} - f(\hat{d}_{ij})]^2}{\sum_{i < j} d_{ij}^2} \right)^{1/2}$$

and the metric MDS minimizes $\mathcal{L}(\hat{d}_{ij})$ over all \hat{d}_{ij} and α and β .

Sammon Mapping

The Sammon nonlinear mapping, which has become a popular tool for pattern recognition, is generalization of the metric LS MDS. The Sammon mapping minimize the stress

$$\text{Sammon's stress}(\hat{d}_{ij}) = \frac{1}{\sum_{u < v} d_{uv}} \sum_{i < j} \frac{(\hat{d}_{ij} - d_{ij})^2}{d_{ij}}$$

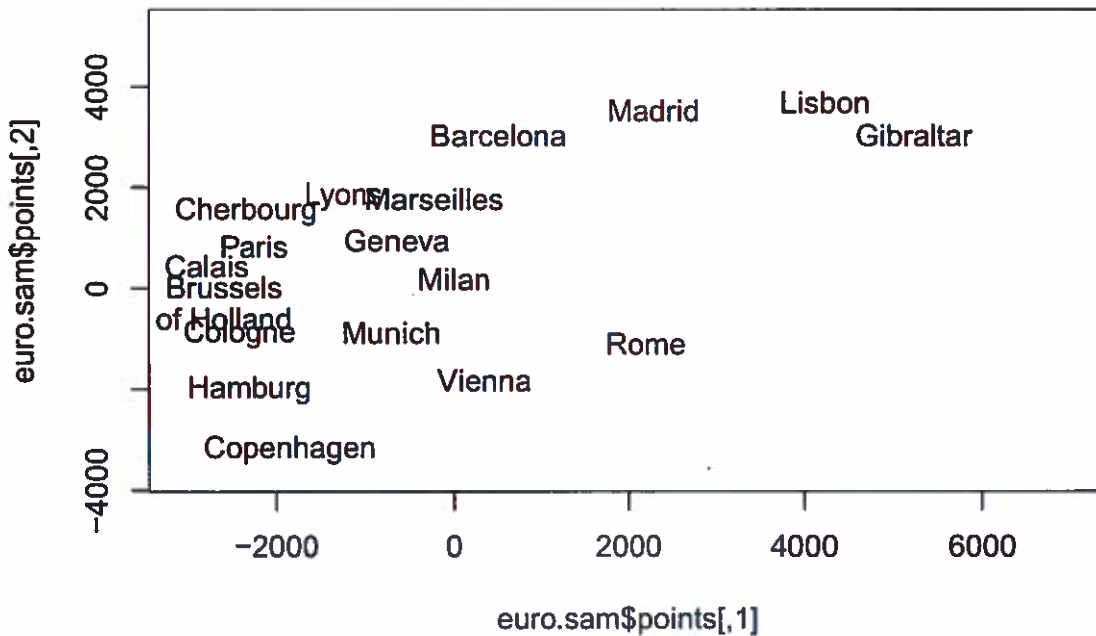
This weighting system normalizes the squared-errors in pairwise distances by using the distance in the original space. As a result, Sammon mapping preserves the small d_{ij} , giving them a greater degree of importance in the fitting procedure than for larger values of d_{ij} . Optimal solution is found by numerical computation.

library(MASS)

```
##
## Attaching package: 'MASS'
## The following objects are masked _by_ '.GlobalEnv':
##
##   lda, qda
euro.sam <- sammon(dist(eurodist))

## Initial stress      : 0.06361
## stress after 10 iters: 0.03351, magic = 0.092
## stress after 20 iters: 0.01925, magic = 0.500
## stress after 30 iters: 0.01864, magic = 0.500
## stress after 40 iters: 0.01843, magic = 0.500
## stress after 50 iters: 0.01835, magic = 0.500

plot(euro.sam$points, type = "n")
text(euro.sam$points[,1], -euro.sam$points[,2], labels=labels(eurodist) )
```



3. Non-metric MDS

In many applications of MDS, dissimilarities are known only by their rank order, and the spacing between successively ranked dissimilarities if of no interest or is unavailable. This may happen because the data collected involve only ordinal information. In non-metric MDS, a.k.a. ordinal MDS, we assume that f is an arbitrary function that satisfies the monotonicity constraint $f(d_{ij}) \leq f(d_{kl})$ whenever $d_{ij} \leq d_{kl}$ for all $i, j, u, v = 1, \dots, n$.

Kruskal (1964) considers construction of $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{R}^k$ from \mathbf{D} based on ranks; any monotone transformation of the distances in \mathbf{D} gives the same answer. A data matrix $[\mathbf{x}_1, \dots, \mathbf{x}_n]$ with corresponding interpoint distances d_{ij} is constructed such that

$$d_{ij} \approx \|\mathbf{x}_i - \mathbf{x}_j\|.$$

Kruskal's non-metric MDS minimizes the

$$stress(\bar{d}_{ij}, d_{ij}) = \left(\frac{\sum_{i < j} (\bar{d}_{ij} - d_{ij})^2}{\sum \bar{d}_{ij}^2} \right)^{1/2}$$

```
f=isoMDS(eurodist) # Kruskal's non-metric multidimensional scaling
```

```
## initial value 7.505733
```

```
## final value 7.505688
```

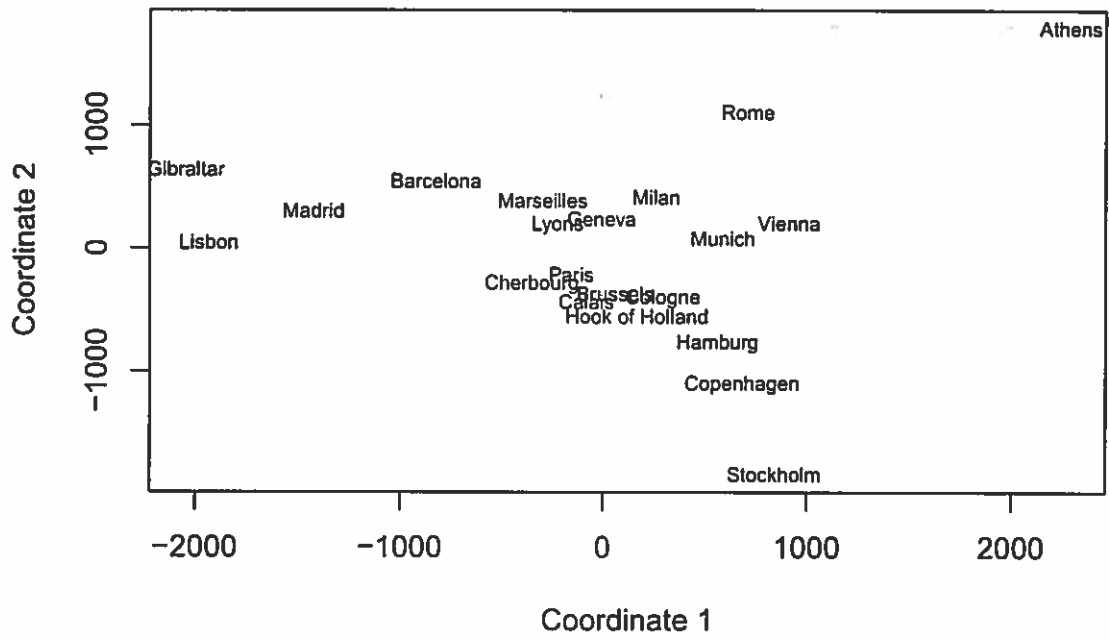
```
## converged
```

```
x=f$points[,1]; y=f$points[,2]
```

```
plot(x,y,xlab="Coordinate 1", ylab="Coordinate 2",main="Nonmetric MDS", type="n")
```

```
text(x,y,labels=labels(eurodist),cex=0.7)
```

Nonmetric MDS





25/28

Tran Le

MATH 755 MULTIVARIATE STATISTICAL ANALYSIS

Assignment 1

Due on in class Thursday September 6th

- 2 ~~Q1.~~ Prove the cyclic property of the trace for arbitrary square matrices of common finite size ($n \geq 2$). I.e., prove

$$\text{tr}(\mathbf{A}_1 \mathbf{A}_2 \cdots \mathbf{A}_n) = \text{tr}(\mathbf{A}_{k+1} \cdots \mathbf{A}_n \mathbf{A}_1 \mathbf{A}_2 \cdots \mathbf{A}_k)$$

- W ~~Q2.~~ Johnson and Wichern 2.4

- 3 ~~Q3.~~ We have for a symmetric matrix A and positive integer m , $A^m = P \Lambda^m P^T$. Show that if A is non-singular, this holds for all integers m .

- 4 ~~Q4.~~ Johnson and Wichern 2.8, 2.9 (Do not use a computer)

- 4 ~~Q5.~~ Johnson and Wichern 2.12, 2.13

- 4 ~~Q6.~~ Johnson and Wichern 2.15

- 4 ~~Q7.~~ Johnson and Wichern 2.16 (Note that non-negative definite means positive semi definite)

2



Q17

Prove the cyclic property of the trace for arbitrary square matrices of common finite size ($n \geq 2$). Prove $\text{tr}(A_1 A_2 \dots A_k A_{k+1} \dots A_n) = \text{tr}(A_{k+1} \dots A_n A_1 A_2 \dots A_k)$

* Note that for all B and C are square matrices, we have $\text{tr}(BC) = \text{tr}(CB)$

We have

$$\begin{aligned} \text{tr}(A_1 A_2 \dots A_k A_{k+1} \dots A_n) &= \text{tr}(C_1 B_1) = \text{tr}(A_2 A_3 \dots A_k A_{k+1} \dots A_n A_1) = \\ &= \text{tr}(C_2 B_2) = \text{tr}(A_3 A_4 \dots A_n A_1 A_2) = \text{tr}(C_3 B_3) = \text{tr}(A_4 \dots A_n A_1 A_2 A_3) = \dots \end{aligned}$$

Keep putting $\begin{cases} B_i = A_i \\ C_i = A_{i+1} \dots A_n A_1 \dots A_{i-1} \end{cases} \quad i = 1, 2, \dots, k$

and apply $\text{tr}(B_i C_i) = \text{tr}(C_i B_i)$

then we have what we need to prove \square

prove otherwise incomplete

Q2 J&W 2.4

When A^{-1} and B^{-1} exist, prove that
 a) $(A^{-1})^{-1} = (A^{-1})'$
 b) $(AB)^{-1} = B^{-1}A^{-1}$

a) We need to prove $(A^{-1})^{-1} = (A^{-1})' \Rightarrow$ need to prove $\begin{cases} (A^{-1})^{-1} (A^{-1})' = I \\ (A^{-1})' A^{-1} = I \end{cases}$

We have

$$\begin{cases} A^{-1} A = I = I' = (A^{-1} A)' = A' (A^{-1})' \\ A A^{-1} = I = I' = (A A^{-1})' = (A^{-1})' A' \end{cases} \Rightarrow \text{what we need to prove } \square 2a.$$

b) We need to prove $(AB)^{-1} = B^{-1}A^{-1} \Rightarrow$ need to prove $\begin{cases} (AB)^{-1} (B^{-1}A^{-1}) = I \\ (B^{-1}A^{-1})(AB) = I \end{cases}$

We have

$$\begin{cases} (AB)^{-1} (B^{-1}A^{-1}) = A (B B^{-1}) A^{-1} = A I A^{-1} = A A^{-1} = I \\ (B^{-1}A^{-1})(AB) = B^{-1} (A^{-1} A) B = B^{-1} I B = B^{-1} B = I \end{cases} \Rightarrow \text{what we need to prove } \square 2b.$$

Q37 We have for a symmetric matrix A a positive integer m .
 Spectral decomposition Λ : eigen values of A
 $A^m = P \Lambda^m P^T$ where P is an orthogonal matrix.
 Show that if A is non-singular then this holds for all integers m .

* We have for $i = -1$, then

$$A^{-1} = (P \Lambda P^T)^{-1} \overset{\text{Orthogonal}}{P^T = P^{-1}} (P \Lambda P^T)^{-1} = P \Lambda^{-1} P^{-1} \overset{P^{-1} = P^T}{=} P \Lambda^{-1} P^T$$

* Assume the hypothesis is true for $i = -k$, which means \uparrow k is a positive integer, we have $A^{-k} = P \Lambda^{-k} P^T$

* Then when $i = -(k+1)$, we have

$$A^{-(k+1)} = [A^{(k+1)}]^{-1} = (A^k A)^{-1} = A^{-1} (A^k)^{-1} = A^{-1} A^{-k} = (P \Lambda^{-1} P^T) (P \Lambda^{-k} P^T) = P \Lambda^{-1} (\underbrace{P^T P}_{=I}) \Lambda^{-k} P^T = P \Lambda^{-1} \Lambda^{-k} P^T = P \Lambda^{-(k+1)} P^T$$

This So by induction \Rightarrow this hypothesis is true for all negative integers.
 \rightarrow all integers. \square

what about $k=0$?

AS 3/4

Q47 J&W 28, 2.9 (Don't use computer)

2.87 Given $A = \begin{bmatrix} 1 & 2 \\ 2 & -2 \end{bmatrix}$ a) Find eigenvalues λ_1 and λ_2
 eigenvectors \vec{e}_1 and \vec{e}_2
 b) Determine the spectral decomposition of A

a) $\begin{cases} \lambda_1 + \lambda_2 = \text{trace } A = -1 \\ \lambda_1 \lambda_2 = \det A = -6 \end{cases} \Rightarrow \begin{cases} \lambda_1 = 2 \\ \lambda_2 = -3 \end{cases}$

• Solve $(A - \lambda_1 I)\vec{x}_1 = \vec{0} \Leftrightarrow \begin{pmatrix} -1 & 2 \\ 2 & -4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{cases} -x_1 + 2x_2 = 0 \\ 2x_1 - 4x_2 = 0 \end{cases}$

$\Rightarrow \begin{pmatrix} 2 \\ 1 \end{pmatrix}$ is an associative eigenvector of $\lambda_1 = 2 \Rightarrow \vec{e}_1 = \begin{pmatrix} 2 \\ 1 \\ \sqrt{5} \end{pmatrix}$

• Solve $(A - \lambda_2 I)\vec{x}_2 = \vec{0} \Leftrightarrow \underbrace{\begin{pmatrix} 4 & 2 \\ 2 & 1 \end{pmatrix}}_{\text{non singular}} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{cases} 4x_1 + 2x_2 = 0 \\ 2x_1 + 1x_2 = 0 \end{cases}$

$\Rightarrow \begin{pmatrix} 1 \\ -2 \end{pmatrix}$ is an associative eigenvector of $\lambda_2 = -3 \Rightarrow \vec{e}_2 = \begin{pmatrix} 1 \\ -2 \\ \sqrt{5} \end{pmatrix}$

b) Determine the spectral decomposition of A

$A = P \Lambda P^T = \begin{pmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 1/\sqrt{5} & -2/\sqrt{5} \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & -3 \end{pmatrix} \begin{pmatrix} 2/\sqrt{5} & 1/\sqrt{5} \\ 1/\sqrt{5} & -2/\sqrt{5} \end{pmatrix}$

$= 2 \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix} \begin{bmatrix} 2/\sqrt{5} & 1/\sqrt{5} \end{bmatrix} - 3 \begin{bmatrix} 1/\sqrt{5} \\ -2/\sqrt{5} \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & -2/\sqrt{5} \end{bmatrix}$
 $= 2 \begin{bmatrix} 4/5 & 2/5 \\ 2/5 & 1/5 \end{bmatrix} - 3 \begin{bmatrix} 1/5 & -2/5 \\ -2/5 & 4/5 \end{bmatrix} \quad \square 2.8.$

4/4

2.9, Let $A = \begin{bmatrix} 1 & 2 \\ 2 & -2 \end{bmatrix}$ as in exercise 2.8

- a) Find A^{-1}
- b) Compute the eigenvalues and eigenvectors of A^{-1}
- c) Write the spectral decomposition of A^{-1} . Compare it with that of A from ex 2.8.

a) $A^{-1} = \frac{1}{\det A} \begin{pmatrix} -2 & -2 \\ -2 & 1 \end{pmatrix} = -\frac{1}{6} \begin{pmatrix} -2 & -2 \\ -2 & 1 \end{pmatrix} = \begin{pmatrix} 1/3 & 1/3 \\ 1/3 & -1/6 \end{pmatrix}$ symmetric \rightarrow eigenvalues real
 non-singular \rightarrow eigenvalues $\in \mathbb{R}$.

b) eigenvalues $\begin{cases} \lambda_1 \lambda_2 = -1/6 \\ \lambda_1 + \lambda_2 = 1/3 - 1/6 = 1/6 \end{cases} \rightarrow \begin{cases} \lambda_1 = 1/2 \\ \lambda_2 = -1/3 \end{cases}$

• Eigenvectors associated with $\lambda_1 = 1/2$

$$(A - \lambda_1 I) \mathbf{x} = \mathbf{0} \Leftrightarrow \begin{pmatrix} 1/3 - 1/2 & 1/3 \\ 1/3 & -1/6 - 1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} -1/6 & 1/3 \\ 1/3 & -2/3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \text{row 2} = 2 \text{ row 1}$$

$\Rightarrow \begin{pmatrix} 2 \\ 1 \end{pmatrix}$ is an associated eigenvector of $\lambda_1 = 1/2$ choose $e_1 = \begin{pmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{pmatrix}$

• eigenvector associated with $\lambda_2 = -1/3$

$$(A - \lambda_2 I) \mathbf{x} = \mathbf{0} \Rightarrow \begin{pmatrix} 2/3 & 1/3 \\ 1/3 & 1/6 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \text{row 1} = 2 \text{ row 2}$$

$\begin{pmatrix} 1 \\ -2 \end{pmatrix}$ is an associated eigenvector of $\lambda_2 = -1/3$ choose $e_2 = \begin{pmatrix} 2/\sqrt{5} \\ -2/\sqrt{5} \end{pmatrix}$

4/4

* Observe:

eigenvalues of $A^{-1} = \frac{1}{\text{eigenvalues of } A}$

eigenvectors of $A^{-1} =$ eigenvector of A

c) $A^{-1} = \sum_{i=1}^2 \lambda_i e_i e_i^T = \frac{1}{2} \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix} \begin{bmatrix} 2\sqrt{5} & \sqrt{5} \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 2/\sqrt{5} \\ -2/\sqrt{5} \end{bmatrix} \begin{bmatrix} \sqrt{5} & -2\sqrt{5} \end{bmatrix}$

$= \frac{1}{2} \begin{bmatrix} 4/5 & 2/5 \\ 2/5 & 1/5 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 4/5 & -2/5 \\ -2/5 & 4/5 \end{bmatrix}$ where Λ is a diagonal matrix, we have when

\Rightarrow If we have $A = P \Lambda P^T$ then $A^{-1} = P \Lambda^{-1} P^T$ $\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \dots \end{bmatrix}$ then $\Lambda^{-1} = \begin{bmatrix} 1/\lambda_1 & & \\ & 1/\lambda_2 & \\ & & \dots \end{bmatrix}$

Q5 J & W 2.12, 2.13

symmetric

2.12 Show that the determinant of a square matrix $p \times p$ A can be expressed as the product of its eigenvalue $\lambda_1, \lambda_2, \dots, \lambda_p$

that is $\det(A) = \prod_{i=1}^p \lambda_i$ (the result is also true when A is not symmetric) (A diagonalizable)
 ← prove by using eigenvalue decomposition

2.13 Show that $\det Q = \pm 1$ if Q is a $p \times p$ orthogonal matrix

2.127 Since A nonsingular + symmetric, by Spectral decomposition theorem, $A = P \Lambda P^T$, where P orthogonal matrix

Λ : diagonal matrix contains p eigenvalues of A

$$\begin{aligned} \Rightarrow \det(A) &= \det(P \Lambda P^T) = \det(P) \det(\Lambda P^T) = \det(P) \det(\Lambda) \det(P^T) = \\ &= \det(P) \det(P^T) \det(\Lambda) = \det(P P^T) \det(\Lambda) = \det(I) \det(\Lambda) \\ &= \det(\Lambda) = \prod_{i=1}^p \lambda_i \end{aligned}$$

since Λ is a diagonal matrix $\Rightarrow \det \Lambda$ is the product of all entries in the diagonal which are all eigenvalues of A \square

2.137. Show that if Q is an orthogonal matrix $\Leftrightarrow \det Q = 1$ or $\det Q = -1$

• Q is orthogonal $\Leftrightarrow Q Q^T = I$

$$\Rightarrow 1 = \det(Q Q^T) = |Q Q^T| = \underset{|Q| = |A|}{|Q|} \underset{|A^T| = |A|}{|Q^T|} = |Q| |Q| = |Q|^2 \Rightarrow |Q| = \pm 1$$

• Observe: the converse is not true.

$$\text{ex } Q = \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

true.

4/4

Q6 2.15

A quadratic form $x'Ax$ is said to be positive definite $\stackrel{\text{def}}{\iff} A$ is positive definite.
 Is the quadratic form $3x_1^2 + 3x_2^2 - 2x_1x_2$ positive definite?

We have

$$3x_1^2 + 3x_2^2 - 2x_1x_2 = [x_1 \ x_2] \begin{bmatrix} 3 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \tilde{x}' A \tilde{x} \text{ where } A = \begin{bmatrix} 3 & -1 \\ -1 & 3 \end{bmatrix}$$

Consider A , we have $A_{(1)} = 3 \quad \det(A_{(1)}) > 0$

$$A_{(2)} = A \quad \det(A_{(2)}) = 9 - 1 = 8 > 0$$

$\Rightarrow A$ is positive definite \iff the quadratic form is positive definite \square

\Rightarrow Yes!

Q7 2.16

Consider an arbitrary $n \times p$ matrix A . Show that $A'A$ is necessarily
 Then $A'A$ is a symmetric $p \times p$ matrix non-negative definite.

Since $A'A$ is symmetric, $x'A'A x$ is well defined and we want to prove that

$$x'A'A x \geq 0, \forall x \neq 0$$

Put $y := Ax$ then $y'y = (Ax)'Ax = x'A'A x \Rightarrow x'A'A x \geq 0$.

Since $y'y = \sum_{i=1}^n y_i^2 \geq 0$

equality happens when $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$

when x is a solution of $Ax = 0$

good detail.

$\Rightarrow A'A$ is positive semi-definite.

4/4

MATH 755 MULTIVARIATE STATISTICAL ANALYSIS

Assignment 1

Due on in class Thursday September 20th

Let $\mathbf{j} = \mathbf{1}$ denote the column vector of ones, and $\mathbf{J} = \mathbf{jj}^T = \mathbf{11}^T$ be the square matrix of ones. All constants and random variables are taken to be real valued. The population variance σ^2 is taken to be positive ($\sigma^2 > 0$).

Q 1. Consider the sample variance S^2 for a random sample $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2)$.

- (a) Write the sample variance as a quadratic form
- (b) Using the appropriate theorems show that it is unbiased for σ^2 , namely show that $E(S^2) = \sigma^2$.

Q 2. Let \mathbf{X} be a p -by-1 random vector with variance-covariance matrix Σ .

- (a) Explain why Σ must be positive-semi-definite. (What would it imply if Σ wasn't positive?)
- (b) If Σ is positive-semi-definite, but not positive-definite, what does this say about \mathbf{X} ? (Consider $\mathbf{c}^T \mathbf{X}$ for some constant column vector \mathbf{c} .)

Q 3. Let \mathbf{A} be a p -by- p positive definite matrix.

- (a) Show that there exists a p -by-1 random vector \mathbf{X} such that $\Sigma_{\mathbf{X}} = \mathbf{A}$. (You may assume that there exists a random vector \mathbf{Z} where $\mathbf{Z} = (Z_1, \dots, Z_p)$ and $Z_1, \dots, Z_p \stackrel{iid}{\sim} (0, 1)$.)
- (b) Find a transformation of \mathbf{X} such that the variance-covariance matrix of this transformed variable $\mathbf{Y} = \mathbf{C}\mathbf{X}$ is the same as the correlation matrix of \mathbf{X} ($\Sigma_{\mathbf{Y}} = \rho_{\mathbf{X}}$). Comment on the relationship of \mathbf{Y} to \mathbf{X} .

Q 4. Suppose that $\mathbf{Y} \sim (\mu, \Sigma)$ is a 3-by-1 random vector where $\mu = \mu \mathbf{j}$ and for some ρ ,

$$\Sigma_{\mathbf{Y}} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix}$$

- (a) Write Σ as a linear combination of \mathbf{I} and \mathbf{J} .
- (b) For what values of ρ is Σ a valid variance-covariance matrix..
- (c) Compute the expectation $E \left[\sum_{i=1}^3 (Y_i - \bar{Y})^2 \right]$. Comment on this as a function of ρ .

Q 5. Let \mathbf{X} be a 3-by-1 random vector with variance covariance matrix

$$\Sigma_{\mathbf{X}} = \begin{bmatrix} 5 & -4 & 1 \\ -4 & 13 & 9 \\ 1 & 9 & 10 \end{bmatrix}$$

- (a) Find a unit vector $\mathbf{c} = (c_1, c_2, c_3)^T$ such that the variance of $\mathbf{c}^T \mathbf{X}$ is maximized/minimized. (You may use a computer for convenience)
- (b) What is the largest/smallest possible variance for such linear combinations?

Q 6. Johnson and Wichern 2.41

Q 7. Let $\mathbf{X} \sim (\mu, \Sigma)$. Without using a computer compute

- (a) $|\Sigma|$

- (b) $\text{tr}\Sigma$
- (c) $V(\mathbf{e}_i^T \mathbf{X}), i = 1, 2, 3$
- (d) Σ
- (e) $\Sigma^{1/2}$

where $\lambda_1 = 3$, $\lambda_2 = 7$, and $\lambda_3 = 1$ are the eigenvalues of Σ , and $\mathbf{e}_1^T = [1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3}]$, $\mathbf{e}_2^T = [2/\sqrt{6}, -1/\sqrt{6}, -1/\sqrt{6}]$, and $\mathbf{e}_3^T = [0, 1/\sqrt{2}, -1/\sqrt{2}]$ are the associated eigenvectors.

Q 8. Let $\mathbf{X} \sim N_3(\boldsymbol{\mu}, \Sigma)$ where $\boldsymbol{\mu}^T = [1, 0, -2]$ and

$$\Sigma = \begin{bmatrix} 2 & -0.5 & 1 \\ -0.5 & 1 & 0 \\ 1 & 0 & 2 \end{bmatrix}$$

- (a) Which components of \mathbf{X} are independent?
 - (b) Find the conditional distribution of $X_2 | (X_1, X_3)^T$
 - (c) Find the conditional distribution of $(X_1, X_3)^T | X_2$
 - (d) Comment on how your answer to part (a) affects parts (b) and (c).
- Q 9. Using the data in Table 4.6 from the book (T4-6), check for univariate and multivariate normality for the females. Comment.
- Q 10. Let $\mathbf{X}_1 \sim N_p(\boldsymbol{\mu}_1, \Sigma)$ and $\mathbf{X}_2 \sim N_p(\boldsymbol{\mu}_2, \Sigma)$. Given independent random samples of sizes n_1 and n_2 respectively, obtain the MLE's of $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, and Σ .
- Q 11. Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} N_p(\boldsymbol{\mu}, \Sigma)$
- (a) Find the mle of k - a scalar- where $\Sigma = k\Sigma_0$ and both $\boldsymbol{\mu}$ and Σ_0 are known.
 - (b) Find the mle of k - a scalar- where $\boldsymbol{\mu} = k\boldsymbol{\mu}_0$ and both $\boldsymbol{\mu}_0$ and Σ are known.

Q17

Consider the sample variance S^2 for a random sample $X_1, \dots, X_n \stackrel{iid}{\sim} (N, \sigma^2)$

a) Write the sample variance as a quadratic form

b) Use the appropriate theorem to show that it is unbiased for σ^2 ; $E(S^2) = \sigma^2$.

a) We have

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} [X_1 - \bar{X} \quad X_2 - \bar{X} \quad \dots \quad X_n - \bar{X}] \begin{bmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{bmatrix}$$

Consider

$$\begin{bmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} - \begin{bmatrix} \bar{X} \\ \bar{X} \\ \vdots \\ \bar{X} \end{bmatrix} = \underline{X} - \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} = \underline{X} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \underline{X} = \mathbf{I} \underline{X} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \underline{X} = \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \underline{X}, \text{ where } \underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

$$\begin{aligned} \text{Then } S^2 &= \frac{1}{n-1} \left[\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \underline{X} \right]^T \left[\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \underline{X} \right] \\ &= \frac{1}{n-1} \left[\underline{X}^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \underline{X} \right] = \end{aligned}$$

note that

$$\begin{aligned} \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) &= \mathbf{I}^2 + \frac{1}{n^2} (\mathbf{1} \mathbf{1}^T)^T (\mathbf{1} \mathbf{1}^T) - \frac{2}{n} \mathbf{I} \mathbf{1} \mathbf{1}^T = \mathbf{I}^2 + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{1} \mathbf{1}^T - \frac{2}{n} \mathbf{1} \mathbf{1}^T \\ &= \mathbf{I}^2 + \frac{1}{n} \mathbf{1} \mathbf{1}^T - \frac{2}{n} \mathbf{1} \mathbf{1}^T = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \stackrel{\text{put}}{=} \mathbf{K} \end{aligned}$$

$$\text{Then } S^2 = \frac{1}{n-1} \left(\underline{X}^T \mathbf{K} \underline{X} \right)$$

$$\text{where } \mathbf{K}^T = \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)^T = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T = \mathbf{K} \Rightarrow \mathbf{K} \text{ symmetric}$$

The quadratic form for S^2 is $S^2 = \frac{1}{n-1} \left(\underline{X}^T \mathbf{K} \underline{X} \right)$

b) Use the appropriate theorem to show that $E(S^2) = \sigma^2$.

$$\begin{aligned} E(S^2) &= \frac{1}{n-1} E \left(\underline{X}^T \mathbf{K} \underline{X} \right) = \frac{1}{n-1} \left\{ \text{tr}(\mathbf{K} \Sigma_X) + \mathbf{1}^T \mathbf{K} \mathbf{1} \right\} = \frac{1}{n-1} \left\{ \text{tr} \left[\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \sigma^2 \mathbf{I} + \mathbf{1}^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \mathbf{1} \right] \right\} \\ &= \frac{1}{n-1} \left\{ \text{tr} \left[\mathbf{I} \sigma^2 \mathbf{I} \right] - \text{tr} \left[\frac{1}{n} \mathbf{1} \mathbf{1}^T \sigma^2 \mathbf{I} \right] \right\} + \mathbf{1}^T \mathbf{I} \mathbf{1} - \frac{1}{n} \mathbf{1}^T \mathbf{1} \mathbf{1}^T \mathbf{1} \\ &= \frac{1}{n-1} \left\{ \left[n \sigma^2 - \frac{1}{n} \text{tr} \begin{bmatrix} \sigma^2 & \sigma^2 & \dots & \sigma^2 \\ \sigma^2 & \sigma^2 & \dots & \sigma^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma^2 & \sigma^2 & \dots & \sigma^2 \end{bmatrix} \right] \right\} + \underbrace{\mathbf{1}^T \mathbf{1} - \mathbf{1}^T \mathbf{1}}_{=0} \\ &= \frac{1}{n-1} \left\{ n \sigma^2 - \frac{1}{n} n \sigma^2 \right\} = \frac{1}{n-1} (n-1) \sigma^2 = \sigma^2 \quad \square \end{aligned}$$



Q2 > Let $X_{p \times 1}$: random vector with covariance-covariance matrix Σ .

a) Explain why Σ must be positive semidefinite (what would imply if Σ wasn't positive?)

b) If Σ is positive semidefinite, but not positive definite, what does this say about X ?

(Hint: consider $c^T X$ for some constant column vector c)

a) First, we want to prove that Σ is positive semi-definite.

We have $\Sigma = E((X - \mu_x)(X - \mu_x)^T)$. Then for any $u_{p \times 1} \neq 0$, we want to prove $u^T \Sigma u \geq 0$.

$$u^T \Sigma u = u^T E((X - \mu_x)(X - \mu_x)^T) u = E\left(\underbrace{u^T (X - \mu_x)}_{y^T} \underbrace{[(X - \mu_x)^T u]}_{\text{put } = y} \right) = E(y^T y) = \sum_{i=1}^p y_i^2$$

Then Σ_x is positive-semi-definite

put = $y \leftarrow$ real value random variable

* Explain why Σ must be positive-semi-definite

If Σ wasn't positive $\iff \exists c_{p \times 1}$, a constant vector $c^T \Sigma c < 0$.

However, $c^T \Sigma c = c^T V(X) c = V(c^T X) = V(Y)$ where $Y = c^T X$

can't be negative \Rightarrow (contradiction).

b) If Σ is positive-semi-definite, but not positive definite, then

$c^T \Sigma c = c^T V(X) c = V(c^T X) = V(Y)$, for $c_{p \times 1} \neq 0$.

This means $V(c^T X) = V(\sum_{i=1}^n c_i X_i) = 0$

This happens when $\exists i$ so that $c_i \neq 0$ so that $\sum_{i=1}^n c_i X_i = 0 \Rightarrow$ there are some X_i are not linearly independent.

If all X_i are linearly independent, $V(\sum_{i=1}^n c_i X_i) = 0 \iff \sum_{i=1}^n c_i V(X_i) = 0$, this can be happened when $\exists i$ so that $V(X_i) = 0$

To sum up, $X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$ can $\left[\begin{array}{l} \text{have some } \{X_i\} \text{ are not linearly independent} \\ \text{have one/some } X_i \text{ so that } V(X_i) = 0 \end{array} \right. \square$

Q3

Let $A_{p \times p}$ positive definite matrix.

a) Show that there exist $X_{p \times 1}$ random vector so that $\Sigma_X = A$

(You may assume that there exists a random vector $Z = \begin{bmatrix} z_1 \\ \vdots \\ z_p \end{bmatrix} \sim N(0, I)$)

b) Find a transformation of X such that the variance-covariance matrix of this transformed variable $Y = CX$ is the same as the correlation matrix of X ($\Sigma_Y = \rho_X$).

Comment on the relationship of Y to X

a) When A is a positive definite matrix, then we have there exist a unique positive def matrix L such that $L^2 = A$, $L = A^{1/2} = P \cdot \Lambda^{-1/2} P^T$

(Remind: P : orthogonal matrix contains eigenvectors of A $A = P \Lambda P$: spectral decomposition of A
 Λ : diagonal matrix, contains eigenvalues of A)

Then for $Z \sim N_p(0, I)$

Put $X = \underbrace{L}_{p \times 1} + \underbrace{L}_{p \times p} Z_{p \times 1} \stackrel{\text{a theorem}}{\sim} N_p(0, \underbrace{L L^T}_{\substack{\uparrow \\ L \text{ is positive def} \Rightarrow \text{symmetric}}}) = N_p(0, L^2) = N_p(0, A) \quad \square a.$

* Find $Y = CX$ so that $\Sigma_Y = \rho_X$

• Note that when $Y = CX$,

then $\Sigma_Y = V(Y) = V(CX) = C V(X) C^T = C \Sigma_X C^T$

• Note that when $V^{1/2} = \begin{pmatrix} \sqrt{\lambda_{11}} & & \\ & \ddots & \\ & & \sqrt{\lambda_{pp}} \end{pmatrix}$, then $\Sigma_X = V^{1/2} \rho_X V^{1/2} \Leftrightarrow \rho_X = V^{-1/2} \Sigma_X V^{-1/2}$
 ↑
 standard deviation matrix of X

Then let $C := V^{-1/2}$

then $\Sigma_Y = C \Sigma_X C^T = V^{-1/2} \Sigma_X (V^{-1/2})^T = V^{-1/2} \Sigma_X V^{-1/2} = \rho_X$

Then $Y = V^{-1/2} X$

V is diagonal \Rightarrow symmetric.

* Comment on the relationship of Y to X .

$Y = V^{-1/2} X$ then $\Sigma_Y = \rho_X$

$\begin{bmatrix} Y_1 \\ \vdots \\ Y_p \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{\lambda_{11}}} & & \\ & \ddots & \\ & & \frac{1}{\sqrt{\lambda_{pp}}} \end{bmatrix} \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \Leftrightarrow Y_i = \frac{1}{\sqrt{\lambda_{ii}}} X_i, \forall i = 1, p$

Q4) Suppose $Y \sim (\mu, \Sigma)$ $\mu = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \mathbf{1}\mathbf{1}^T$ $\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix}$

a) Write Σ_Y as a linear combination of \mathbf{I} and \mathbf{J}

b) For what value of ρ is Σ a valid variance-covariance matrix

c) Compute the expectation $E\left(\sum_{i=1}^3 (Y_i - \bar{Y})^2\right)$. Comment on this as a function of ρ .

$$\begin{aligned} \text{a) } \Sigma_Y &= \sigma^2 \begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix} = \sigma^2 \left\{ (1-\rho) \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} + \rho \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \right\} = \\ &= \sigma^2 \left\{ (1-\rho) \mathbf{I} + \rho \mathbf{1}\mathbf{1}^T \right\} = \sigma^2 \left\{ (1-\rho) \mathbf{I} + \rho \mathbf{J} \right\} \end{aligned}$$

b) By Sylvester's criterion

Σ_Y is a variance-covariance matrix
 $\Rightarrow \Sigma_Y$ is positive-semi-positive $\Leftrightarrow \begin{cases} \Sigma_Y \text{ is symmetric} \\ \text{all leading } m \times m \ (m=1, 2, 3) \ \Sigma_{n \times m} \text{ have } \det(\Sigma_{m \times m}) \geq 0 \end{cases}$

• $\det(\Sigma_{11}) = \det[\sigma^2] = \sigma^2 > 0$

• $\det(\Sigma_{22}) = \det[\sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}] = (\sigma^2)^2 \det \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} = \sigma^4 (1-\rho^2)$

• $\det(\Sigma_{33}) = (\sigma^2)^3 \det \begin{pmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{pmatrix} = \sigma^6 \det \begin{pmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ 1+2\rho & 1+2\rho & 1+2\rho \end{pmatrix} = \sigma^6 \det \begin{pmatrix} 1-\rho & 0 & \rho \\ 0 & 1-\rho & \rho \\ 0 & 0 & 1+2\rho \end{pmatrix} = \sigma^6 (1-\rho)^2 (1+2\rho)$

So we need $\begin{cases} 1-\rho^2 \geq 0 \\ (1-\rho)(1+2\rho) \geq 0 \end{cases} \Leftrightarrow \begin{cases} -1 \leq \rho \leq 1 \\ -\frac{1}{2} \leq \rho \leq 1 \end{cases} \Rightarrow \frac{1}{2} \leq \rho \leq 1$

c7 Compute $E\left(\sum_{i=1}^3 (Y_i - \bar{Y})^2\right)$. Comment on this as a function of ρ

From problem 1) $\sum_{i=1}^3 (Y_i - \bar{Y})^2 = \underline{Y}^T K \underline{Y}$, where $K = I - \frac{1}{n} \mathbb{1} \mathbb{1}^T$

Note that K is symmetric & $K^2 = K$, $\underline{1}^T K \underline{1} = 0$

when $\underline{1} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$

$$\text{Then } E\left(\sum_{i=1}^3 (Y_i - \bar{Y})^2\right) = E(\underline{Y}^T K \underline{Y}) = \text{tr}(K \Sigma_Y) + \underline{1}^T K \underline{1}$$

$$= \text{tr}\left\{\left(I - \frac{1}{n} \mathbb{1} \mathbb{1}^T\right) \left(\sigma^2[(1-\rho)I + \rho \mathbb{1} \mathbb{1}^T]\right)\right\} + \underbrace{\underline{1}^T \left(I - \frac{1}{n} \mathbb{1} \mathbb{1}^T\right) \underline{1}}_{=0 \text{ (similar to problem 1)}}$$

$$= \text{tr}\left(I \sigma^2(1-\rho)I - \frac{1}{n} \mathbb{1} \mathbb{1}^T \sigma^2(1-\rho)I + \rho I \mathbb{1} \mathbb{1}^T - \frac{\rho}{n} \mathbb{1} \mathbb{1}^T \mathbb{1} \mathbb{1}^T\right)$$

$$= n \sigma^2(1-\rho) - \frac{1}{n} \sigma^2(1-\rho) \text{tr}\left[\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & & \\ & \times & \\ & & \times \end{pmatrix}\right] + \rho \text{tr}\left[\begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix}\right] - \frac{\rho}{n} \text{tr}\left[\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}\right]$$

$$= n \sigma^2(1-\rho) - \frac{1}{n} \sigma^2(1-\rho) n + \rho n - \frac{\rho}{n} n n$$

$$= n \sigma^2(1-\rho) - \sigma^2(1-\rho) + n\rho - n\rho$$

$$= (n-1) \sigma^2(1-\rho)$$

$$\stackrel{n=3}{=} 2 \sigma^2(1-\rho)$$

* Comment:

$$E\left(\sum_{i=1}^3 (Y_i - \bar{Y})^2\right) = 2 \sigma^2(1-\rho) \geq 0 \text{ (since } -\frac{1}{2}\rho \leq 1)$$

Q57

Let X be a $3 \times L$ random vector,

$$\Sigma_X = \begin{bmatrix} 5 & -4 & 1 \\ -4 & 13 & 9 \\ 1 & 9 & 10 \end{bmatrix}$$

a7 Find a unit vector $\underline{c} = \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix}$ s.t $V(\underline{c}^T X)$ is maximized / minimized.

(May use a computer for convenience)

b7 What is the largest / smallest possible variance for such linear combinations?

a7 We have $V(\underline{c}^T X) = \underline{c}^T V(X) \underline{c} = \underline{c}^T \Sigma_X \underline{c}$

we want $V(\underline{c}^T X)$ to be maximized / minimized \Rightarrow

\rightarrow we want to find \underline{c} associative with maximum eigenvalue and minimum eigenvalue

* The two commands

$A = \text{matrix}(c(5, -4, 1, -4, 13, 9, 1, 9, 10), ncol=3, nrow=3, byrow=TRUE))$

$\text{eigen}(A)$

gives us three eigenvalues: $2.1e+01$, $7e+00$, $2.13e-14$.

$$\underline{c}_I = \begin{bmatrix} 0.1543 \\ -0.7715 \\ -0.6172 \end{bmatrix}$$

are the two unit vectors when

$V(\underline{c}^T X)$ is maximized

$$\underline{c}_{II} = \begin{bmatrix} 0.5773 \\ 0.5773 \\ -0.5773 \end{bmatrix}$$

$V(\underline{c}^T X)$ is minimized.

b7 The largest possible variance $\underline{c}_I^T \Sigma_X \underline{c}_I = 21$

The smallest possible variance $\underline{c}_{II}^T \Sigma_X \underline{c}_{II} = 0$

since $\underline{c}^T \Sigma_X \underline{c} = \begin{bmatrix} \lambda_1^{\max} \\ \lambda_2 \\ \lambda_3^{\min} \end{bmatrix}$.

Q7 Johnson and Wichern 2.4 L.

$$\text{Let } \underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} \quad \mu_X = \begin{bmatrix} 3 \\ 2 \\ -2 \\ 0 \end{bmatrix} \quad \Sigma_X = \begin{bmatrix} 3 & & & \\ & 3 & & \\ & & 3 & \\ & & & 3 \end{bmatrix} \quad \text{Let } A = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix}$$

a) Find $E(AX)$ b) Find $\text{Cov}(AX)$

c) Which pairs of linear combinations have zero covariances.

$$\text{a) } E(AX) = A E(X) = A \mu_X = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ -2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 9 \\ 3 \end{bmatrix}$$

$$\text{b) } \text{Cov}(AX) = A \text{Cov}(X) A^T = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & 1 \\ 0 & -2 & 1 \\ 0 & 0 & -3 \end{bmatrix} = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 18 & 0 \\ 0 & 0 & 36 \end{bmatrix}$$

3x3

c) Which pairs of linear combinations have zero covariances.

We have from the $\text{Cov}(AX)$ that all pairs (there of them) are have zero covariances.

Q77

Let $X \sim (\lambda_1, \lambda_2, \lambda_3)$. Without using computer, compute

(where $\lambda_1=3, \lambda_2=7, \lambda_3=1$ are eigenvalues of X)

$e_1 = \begin{pmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix}$ $e_2 = \begin{pmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{pmatrix}$ $e_3 = \begin{pmatrix} 0 \\ 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$ are associated eigenvectors.

compute, $a) |Z|$ $b) \text{tr}(Z)$ $c) V(e_i^T X), i=1,2,3$ $d) Z$ $e) Z^{112}$.

$a) |Z| = \prod_{i=1}^3 \lambda_i = 3 \cdot 7 \cdot 1 = 21$

$b) \text{tr}(Z) = 3 + 7 + 1 = 11$

$c) \forall i=1,2,3$

$$V(e_i^T X) = e_i^T V(X) e_i = e_i^T Z e_i = e_i^T P \Lambda P^T e_i = e_i^T [e_1^T e_2^T e_3^T] \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix} \begin{bmatrix} e_i^T \\ e_i^T \\ e_i^T \end{bmatrix} e_i$$

Note that $e_i^T e_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$

$\Rightarrow V(e_i^T X) = \lambda_i, \forall i=1,2,3$.

For example:

$$V(e_1^T X) = [e_1^T [e_1^T e_2^T e_3^T] \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix} \begin{bmatrix} e_1^T \\ e_2^T \\ e_3^T \end{bmatrix} e_1 = [1 \ 0 \ 0] \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} =$$

Similarly for case $i=2,3$. \Rightarrow done.

$d)$ Compute Z

$$Z = P \Lambda P^T = \sum_{i=1}^3 \lambda_i e_i e_i^T = 3 \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} [1/\sqrt{3} \ 1/\sqrt{3} \ 1/\sqrt{3}] + 7 \begin{bmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{bmatrix} [2/\sqrt{6} \ -1/\sqrt{6} \ -1/\sqrt{6}] + 1 \begin{bmatrix} 0 \\ 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} [0 \ 1/\sqrt{2} \ -1/\sqrt{2}]$$

$$= 3 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix} + 7 \begin{bmatrix} 4/6 & -2/6 & -2/6 \\ -2/6 & 1/6 & 1/6 \\ -2/6 & 1/6 & 1/6 \end{bmatrix} + 1 \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1/2 & -1/2 \\ 0 & -1/2 & 1/2 \end{bmatrix} = \begin{bmatrix} 17/3 & -4/3 & -4/3 \\ -4/3 & 8/3 & 5/3 \\ -4/3 & 5/3 & 8/3 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 17 & -4 & -4 \\ -4 & 8 & 5 \\ -4 & 5 & 8 \end{bmatrix}$$

$$e_7 \leq \lambda^2 = P \cdot \Lambda^2 \cdot P^T = \sum_{i=1}^3 \lambda_i e_i e_i^T$$

$$= \sqrt{3} \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} [1/\sqrt{3} \ 1/\sqrt{3} \ 1/\sqrt{3}] + \sqrt{7} \begin{bmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{bmatrix} [2/\sqrt{6} \ -1/\sqrt{6} \ -1/\sqrt{6}] + 1 \begin{bmatrix} 0 \\ 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} [0 \ 1/\sqrt{2} \ -1/\sqrt{2}]$$

$$= \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix} + \begin{bmatrix} \frac{4\sqrt{7}}{6} & \frac{-2\sqrt{7}}{6} & \frac{-2\sqrt{7}}{6} \\ \frac{-2\sqrt{7}}{6} & \frac{\sqrt{7}}{6} & \frac{\sqrt{7}}{6} \\ \frac{-2\sqrt{7}}{6} & \frac{\sqrt{7}}{6} & \frac{\sqrt{7}}{6} \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{1}{2} & -\frac{1}{2} \\ 0 & -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$= \frac{1}{6} \begin{bmatrix} 2\sqrt{3} + 4\sqrt{7} & 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} - 2\sqrt{7} \\ 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} + \sqrt{7} + 3 & 2\sqrt{3} + \sqrt{7} - 3 \\ 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} + \sqrt{7} - 3 & 2\sqrt{3} + \sqrt{7} + 3 \end{bmatrix} \quad \square$$

Q87 Let $X \sim N_3(\mu, \Sigma)$ $\mu = \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix}$ $\Sigma = \begin{bmatrix} 2 & -0.5 & 1 \\ -0.5 & 1 & 0 \\ 1 & 0 & 2 \end{bmatrix}$

a) Which components of X are independent?

b) Find the conditional distribution of $X_2 | (X_1, X_3)^T$

c) Find the conditional distribution of $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2$

d) Comment on how your answer to part a) affects parts b) and c)

a) Since $\Sigma = [\sigma_{ij}]_{i,j}$ contains σ_{ij} , the population covariance of X_i and X_j is σ_{ij} . $\sigma_{23} = \sigma_{32} = 0 \Rightarrow X_2$ and X_3 are independent.

b) Find the conditional distribution $X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}^T$

We consider $X = \begin{bmatrix} X_2 \\ X_1 \\ X_3 \end{bmatrix} = \begin{bmatrix} Y \\ Z \end{bmatrix}$ then we have $\mu_X = \begin{bmatrix} \mu_2 \\ \mu_1 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ -2 \end{bmatrix}$

$\Sigma_X = \begin{bmatrix} \sigma_{22} & \sigma_{21} & \sigma_{23} \\ \sigma_{12} & \sigma_{11} & \sigma_{13} \\ \sigma_{32} & \sigma_{31} & \sigma_{33} \end{bmatrix} = \begin{bmatrix} 2 & -0.5 & 1 \\ -0.5 & 1 & 0 \\ 1 & 0 & 2 \end{bmatrix}$

Then $f_{X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}}$

and we also have

$\mu_{\begin{pmatrix} X_1 \\ X_3 \end{pmatrix}} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$

$\Sigma_{\begin{pmatrix} X_1 \\ X_3 \end{pmatrix}} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$

Then we have the conditional distribution of $X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}$ has normal distribution with

mean $\mu_{X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}} = \mu_2 + \Sigma_{YZ} \Sigma_{ZZ}^{-1} \left[\begin{pmatrix} x_1 \\ x_3 \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_3 \end{pmatrix} \right] =$

$= 0 + [-0.5 \ 0] * \frac{1}{3} \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 - 1 \\ x_3 + 2 \end{bmatrix} = \frac{1}{3} [-1 \ 0.5] \begin{bmatrix} x_1 - 1 \\ x_3 + 2 \end{bmatrix} =$
 $= -\frac{1}{3}x_1 + \frac{1}{6}x_3 + \frac{2}{3}$

covariance $\Sigma_{X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}} = 1 - [-0.5 \ 0] \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} = 1 - [-0.5 \ 0] * \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} -0.5 \\ 0 \end{bmatrix}$
 $= 1 - \frac{1}{3} [-1 \ 0.5] \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} = 1 - \frac{1}{3} (0.5) = 1 - \frac{1}{6} = \frac{5}{6}$

So $X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix} \sim N\left(-\frac{1}{3}x_1 + \frac{1}{6}x_3 + \frac{2}{3}, \frac{5}{6}\right)$

e) Find the conditional distribution of $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2$.

Consider $X = \begin{pmatrix} X_1 \\ X_3 \\ X_2 \end{pmatrix}$, has $\mu_X = \begin{pmatrix} 1 \\ -2 \\ 0 \end{pmatrix}$ $\Sigma_X = \begin{pmatrix} \sigma_{11} & \sigma_{13} & \sigma_{12} \\ \sigma_{31} & \sigma_{33} & \sigma_{32} \\ \sigma_{21} & \sigma_{23} & \sigma_{22} \end{pmatrix} = \begin{pmatrix} 2 & 1 & -0.5 \\ 1 & 2 & 0 \\ -0.5 & 0 & 1 \end{pmatrix}$

Then $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2$ has normal distribution with

$$\text{mean} = \begin{pmatrix} \mu_1 \\ \mu_3 \end{pmatrix} + \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} (1)^{-1} (x_2 - \mu_2) = \begin{pmatrix} 1 \\ -2 \end{pmatrix} + \begin{pmatrix} -0.5 \\ 0 \end{pmatrix} (x_2 - 0) = \begin{pmatrix} 1 - 0.5x_2 \\ -2 \end{pmatrix}$$

$$\text{Covariance} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} (1)^{-1} [-0.5 \ 0] = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} 0.25 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1.75 & 1 \\ 1 & 2 \end{bmatrix}$$

Then $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2 = x_2 \sim \text{Normal} \left(\begin{pmatrix} 1 - 0.5x_2 \\ -2 \end{pmatrix}, \begin{bmatrix} 1.75 & 1 \\ 1 & 2 \end{bmatrix} \right)$

d) Comment on how the answer to part a ~~one~~ affects part b and c.

We have the result of part a is X_2 and X_3 are independent.

then from part b) the conditional distribution $X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}$ is still depended on $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix}$.

from part c) $\begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2 = x_2$ is still depended on X_2

⇒ This means that even X_2 and X_3 are independent $X_2 | X_3$ is not depended on X_3

$X_3 | X_2$ is not depended on X_2 .

← still depended on X_3 .

but when we consider $X_2 | \begin{pmatrix} X_1 \\ X_3 \end{pmatrix}$ ← vector created by X_1 and X_3

vector created by X_1 & $X_3 \iff \begin{pmatrix} X_1 \\ X_3 \end{pmatrix} | X_2 \rightarrow$ is still depended on X_2 .

```

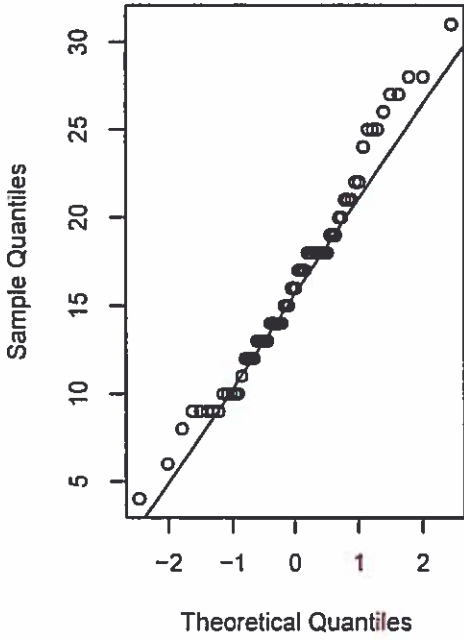
library(mvnTest)
setwd("C:/Users/tranl/Desktop/MAT755/Homework")
data<-read.table(file="C:/Users/tranl/Desktop/MAT755/Wichern_data_textbook/T4-6.DAT")
head(data)
#For each of teenagers, the gender(male=1, female=2)
#and socioeconomic status (low= 1, medium=2)
# The scores labeled independence (indep), support(supp), benevolence(benev)
#conformity(comform) and leadership (leader)

#Name the variables
names(data)<-c("Indep", "Supp", "Benev", "Conform", "Leader", "Gender", "Socio")
head(data)
#We want to check univariate and multivariate normality for female
is.female<-which(data$Gender==2)
data.female<-data[is.female, ]

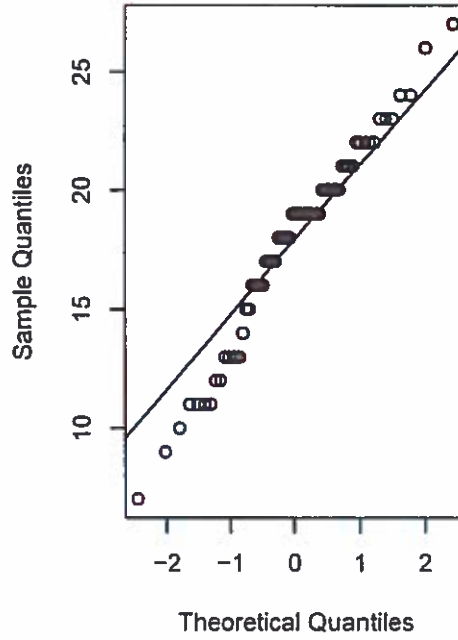
dev.new(width=15, height=15)
par(mfrow=c(2,3))
for (i in c("Indep", "Supp", "Benev", "Conform", "Leader"))
{
  qqnorm(data.female[,i], main="QQ plot for female")
  qqline(data.female[,i])
}

```

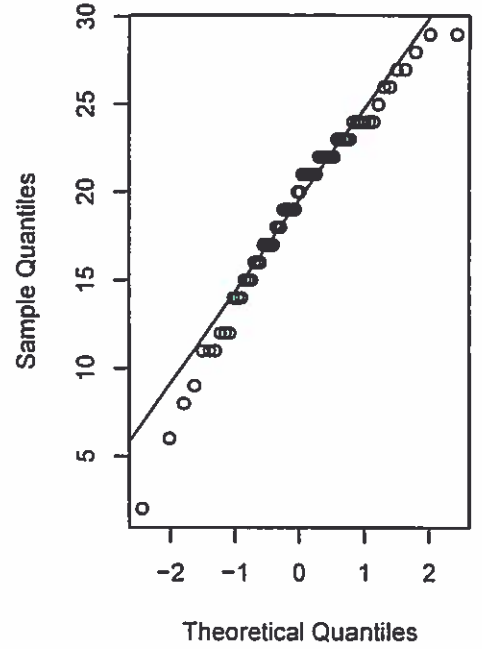
QQ plot for female



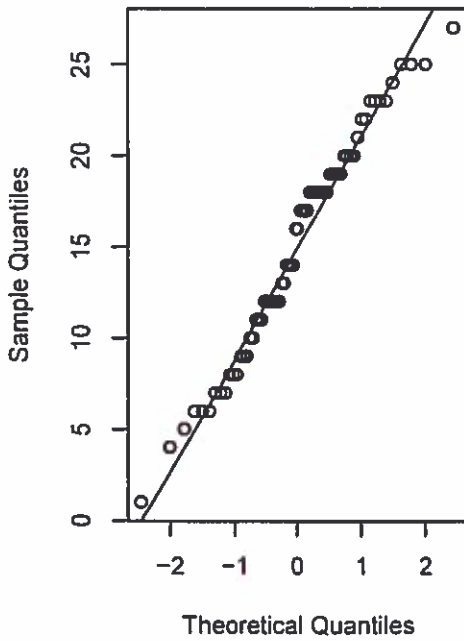
QQ plot for female



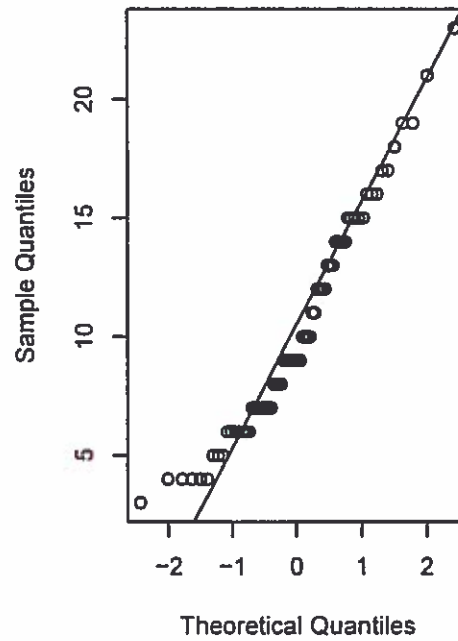
QQ plot for female



QQ plot for female



QQ plot for female



Q107

Let $X_1 \sim N_p(\mu_1, \Sigma)$ Given independent random samples of size n_1 and n_2 respectively,
 $X_2 \sim N_p(\mu_2, \Sigma)$ obtain the MLE of μ_1, μ_2, Σ .

* We have the likelihood is

$$L(\mu_1, \mu_2, \Sigma | X_{1,1}, \dots, X_{1,n_1}, X_{2,1}, \dots, X_{2,n_2}) = L(\mu_1, \Sigma | X_{1,1}, \dots, X_{1,n_1}) L(\mu_2, \Sigma | X_{2,1}, \dots, X_{2,n_2})$$

$$= \frac{1}{(2\pi)^{\frac{n_1 p}{2}} |\Sigma|^{\frac{n_1}{2}}} e^{-\frac{1}{2} \text{tr}[\Sigma^{-1}(n_1-1)S_1]} - \frac{n_1}{2} (\bar{X}_1 - \mu_1)^T \Sigma^{-1} (\bar{X}_1 - \mu_1) \cdot \frac{1}{(2\pi)^{\frac{n_2 p}{2}} |\Sigma|^{\frac{n_2}{2}}} e^{-\frac{1}{2} \text{tr}[\Sigma^{-1}(n_2-1)S_2]} - \frac{n_2}{2} (\bar{X}_2 - \mu_2)^T \Sigma^{-1} (\bar{X}_2 - \mu_2)$$

constant

= (*)

Since Σ is positive definite $\Rightarrow \Sigma^{-1}$ positive definite
 $(\bar{X}_i - \mu_i)^T \Sigma^{-1} (\bar{X}_i - \mu_i) \geq 0$ unless when $\mu_i = \bar{X}_i$

\Rightarrow these terms are maximized when $\mu_1 = \bar{X}_1$ and $\mu_2 = \bar{X}_2$

So we have $\hat{\mu}_{1MLE} = \bar{X}_1$ $\hat{\mu}_{2MLE} = \bar{X}_2$

* Now we want to maximize (*)

$$(*) = \frac{1}{|\Sigma|^{\frac{n_1+n_2}{2}}} e^{-\frac{1}{2} \text{tr}[\Sigma^{-1} \{ \underbrace{(n_1-1)S_1 + (n_2-1)S_2}_{\text{put this equal B}} \}]}$$

$$= \frac{1}{|\Sigma|^{\frac{n_1+n_2}{2}}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} B)} \leq \frac{1}{|B|^b} (2b)^{pb} e^{-bp} \quad \text{for all } \Sigma > 0$$

$b = \frac{n_1+n_2}{2}$

the equality happens when $\Sigma = \frac{1}{2b} B =$

$$= \frac{1}{2 \frac{n_1+n_2}{2}} [(n_1-1)S_1 + (n_2-1)S_2]$$

So we have $\hat{\Sigma}_{MLE} = \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2}$ \square

Q117 Suppose $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ known

a) Find the MLE of k , a scalar, where $\Sigma = k \Sigma_0$, μ, Σ_0

We want to find the likelihood of k .

$$L(k | X_1, \dots, X_n, \mu, \Sigma_0) = \frac{1}{(2\pi)^{\frac{np}{2}} |k\Sigma_0|^{\frac{n}{2}}} e^{-\frac{1}{2} \sum_{j=1}^n (X_j - \mu)^T \Sigma^{-1} (X_j - \mu)}$$

$$= \frac{1}{(2\pi)^{\frac{np}{2}}} \frac{1}{|k\Sigma_0|^{\frac{n}{2}}} e^{-\frac{1}{2} \sum_{j=1}^n (X_j - \mu)^T \Sigma_0^{-1} (X_j - \mu)} \quad (*)$$

* The log likelihood

$$l(k | X_1, \dots, X_n, \mu, \Sigma_0) \propto -\frac{np}{2} \log k - \frac{1}{2k} \sum_{j=1}^n (X_j - \mu)^T \Sigma_0^{-1} (X_j - \mu) \quad (\text{cancel})$$

$$\frac{\partial l}{\partial k} = -\frac{n}{2} + \frac{1}{2k^2} \sum_{j=1}^n (X_j - \mu)^T \Sigma_0^{-1} (X_j - \mu)$$

$$\frac{\partial^2 l}{\partial k^2} = \frac{-n}{k^3} - \frac{1}{k^3} \sum_{j=1}^n (X_j - \mu)^T \Sigma_0^{-1} (X_j - \mu)$$

$$\frac{\partial l}{\partial k} = 0 \Rightarrow \hat{k} = \frac{1}{n} \sum_{j=1}^n (X_j - \mu)^T \Sigma_0^{-1} (X_j - \mu) \quad \text{is a maximum since } \left. \frac{\partial^2 l}{\partial k^2} \right|_{k=\hat{k}} < 0$$

b) Find the MLE of k , a scalar, where $\mu = k\mu_0$, μ_0, Σ_0 are known.

$$L(k | X_1, \dots, X_n, \mu_0, \Sigma_0) \propto \frac{1}{|k\Sigma_0|^{\frac{n}{2}}} e^{-\frac{1}{2} \sum_{j=1}^n (X_j - k\mu_0)^T \Sigma_0^{-1} (X_j - k\mu_0)}$$

$$l(k | X_1, \dots, X_n, \mu_0, \Sigma_0) \propto \frac{n}{2} \log k + \frac{1}{|k\Sigma_0|^{\frac{n}{2}}}$$

MATH 755 MULTIVARIATE STATISTICAL ANALYSIS

Assignment 1

Due on in class Thursday September 6th

Let $\mathbf{j} = \mathbf{1}$ denote the column vector of ones, and $\mathbf{J} = \mathbf{j}\mathbf{j}^T = \mathbf{1}\mathbf{1}^T$ be the square matrix of ones. All constants and random variables are taken to be real valued. The population variance σ^2 is taken to be positive ($\sigma^2 > 0$).

Q 1. Consider the sample variance S^2 for a random sample $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2)$.

- (a) Write the sample variance as a quadratic form
- (b) Using the appropriate theorems show that it is unbiased for σ^2 , namely show that $E(S^2) = \sigma^2$.

(a) Note that $\mathbf{I}\mathbf{X} = \mathbf{X}$ and $1/n\mathbf{J}\mathbf{X} = \bar{X}\mathbf{1}$. Thus the sum of squared deviations must be

$$\begin{aligned} \sum (X_i - \bar{X})^2 &= ((\mathbf{I} - 1/n\mathbf{J})\mathbf{X})^T ((\mathbf{I} - 1/n\mathbf{J})\mathbf{X}) \\ &= ((\mathbf{I} - 1/n\mathbf{J})\mathbf{X})^T ((\mathbf{I} - 1/n\mathbf{J})\mathbf{X}) \end{aligned}$$

Note that $\mathbf{I} - 1/n\mathbf{J}$ is symmetric (since \mathbf{I} and \mathbf{J} are symmetric) and idempotent since

$$\begin{aligned} (\mathbf{I} - 1/n\mathbf{J})^2 &= \mathbf{I} - 2/n\mathbf{J} + 1/n^2\mathbf{J}^2 \\ &= \mathbf{I} - 2/n\mathbf{J} + 1/n^2(n\mathbf{J}) \\ &= \mathbf{I} - 2/n\mathbf{J} + 1/n\mathbf{J} \\ &= \mathbf{I} - 1/n\mathbf{J} \end{aligned}$$

so continuing along we have

$$\begin{aligned} \sum (X_i - \bar{X})^2 &= \mathbf{X}^T (\mathbf{I} - 1/n\mathbf{J})^T (\mathbf{I} - 1/n\mathbf{J}) \mathbf{X} \\ &= \mathbf{X}^T (\mathbf{I} - 1/n\mathbf{J})^2 \mathbf{X} \\ &= \mathbf{X}^T (\mathbf{I} - 1/n\mathbf{J}) \mathbf{X} \end{aligned}$$

and thus we have that

$$S^2 = \frac{1}{n-1} \mathbf{X}^t (\mathbf{I} - 1/n\mathbf{J}) \mathbf{X}$$

(b) Note that $\boldsymbol{\mu} = \mu\mathbf{1}$ and $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}$, then by Theorem 2.4 we have

$$\begin{aligned} E[S^2] &= \frac{1}{n-1} E[\mathbf{X}^t (\mathbf{I} - 1/n\mathbf{J}) \mathbf{X}] \\ &= \frac{1}{n-1} \left[\text{tr}((\mathbf{I} - 1/n\mathbf{J})\sigma^2\mathbf{I}) + \boldsymbol{\mu}^T \mathbf{I} \boldsymbol{\mu} - \frac{1}{n} \boldsymbol{\mu}^T \mathbf{J} \boldsymbol{\mu} \right] \\ &= \frac{1}{n-1} \left[\sigma^2 \text{tr}(\mathbf{I} - 1/n\mathbf{J}) + \boldsymbol{\mu}^T \mathbf{I} \boldsymbol{\mu} - \frac{1}{n} \boldsymbol{\mu}^T (n\boldsymbol{\mu}) \right] \\ &= \frac{1}{n-1} [\sigma^2(\text{tr}\mathbf{I} - 1/n\text{tr}\mathbf{J}) + \boldsymbol{\mu}^T \mathbf{I} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{I} \boldsymbol{\mu}] \\ &= \frac{1}{n-1} [\sigma^2(n - 1/n(n))] \\ &= \frac{1}{n-1} [\sigma^2(n - 1)] \\ &= \sigma^2 \end{aligned}$$

Q 2. Let \mathbf{X} be a p -by-1 random vector with variance-covariance matrix Σ .

- (a) Explain why Σ must be positive-semi-definite. (What would it imply if Σ wasn't positive?)
- (b) If Σ is positive-semi-definite, but not positive-definite, what does this say about \mathbf{X} ? (Consider $\mathbf{c}^T \mathbf{X}$ for some constant column vector \mathbf{c} .)

(a) Let \mathbf{c} be a p -by-1 constant vector. Then the variance of the linear combination $Z = \mathbf{c}^T \mathbf{X}$ is

$$V(Z) = V(\mathbf{c}^T \mathbf{X}) = \mathbf{c}^T \Sigma_{\mathbf{X}} \mathbf{c}$$

If $\Sigma_{\mathbf{X}}$ is not positive-semidefinite, then there exists a \mathbf{c} such that

$$0 > \mathbf{c}^T \Sigma_{\mathbf{X}} \mathbf{c} = V(Z)$$

which is impossible since for all RV's Z

$$V(Z) = E[(Z - \mu_Z)^2] \geq 0$$

(b) If Σ is positive-semi-definite, but not positive-definite, then there exists a non-zero constant vector \mathbf{c} such that the random variable Z has

$$V(Z) = V(\mathbf{c}^T \mathbf{X}) = \mathbf{c}^T \Sigma_{\mathbf{X}} \mathbf{c} = 0$$

This can happen in two ways. Either

- (i) One (or more) random variables are degenerate (i.e., $V(X_i) = 0$)
 - (ii) One (or more) subset of random variables are linearly dependent (almost surely)
- (i) is technically a subcase of (ii) but is distinct in interpretation and ease of resolution. In either case we have redundant variables.

Q 3. Let \mathbf{A} be a p -by- p positive definite matrix.

- (a) Show that there exists a p -by-1 random vector \mathbf{X} such that $\Sigma_{\mathbf{X}} = \mathbf{A}$. (You may assume that there exists a random vector \mathbf{Z} where $\mathbf{Z} = (Z_1, \dots, Z_p)$ and $Z_1, \dots, Z_p \stackrel{\text{iid}}{\sim} (0, 1)$.)
- (b) Find a transformation of \mathbf{X} such that the variance-covariance matrix of this transformed variable $\mathbf{Y} = \mathbf{C}\mathbf{X}$ is the same as the correlation matrix of \mathbf{X} ($\Sigma_{\mathbf{Y}} = \rho_{\mathbf{X}}$). Comment on the relationship of \mathbf{Y} to \mathbf{X} .

(a) Since \mathbf{A} is positive definite it can be factored via its spectral value decomposition as $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$, and it has a symmetric square root, namely $\mathbf{A}^{1/2} = \mathbf{P}\mathbf{\Lambda}^{1/2}\mathbf{P}^T$.

Let us consider first a random vector \mathbf{Y} with mean vector $\boldsymbol{\mu} = \mathbf{0}$ and variance-covariance matrix $\Sigma_{\mathbf{Y}} = \mathbf{I}$. We consider the random variable \mathbf{X} via the following transformation $\mathbf{X} = \mathbf{A}^{1/2}\mathbf{Y}$. Then

$$\begin{aligned}\Sigma_{\mathbf{X}} &= V(\mathbf{X}) = V(\mathbf{A}^{1/2}\mathbf{Y}\mathbf{A}^{1/2}) \\ &= \mathbf{A}^{1/2}\Sigma_{\mathbf{Y}}[\mathbf{A}^{1/2}]^T \\ &= \mathbf{A}^{1/2}\mathbf{I}\mathbf{A}^{1/2} \\ &= \mathbf{A}\end{aligned}$$

as desired

(b) Let us consider transformations $\mathbf{Y} = \mathbf{C}\mathbf{X}$ where \mathbf{C} is a p -by- p constant matrix. Then

$$V(\mathbf{Y}) = V(\mathbf{C}\mathbf{X}) = \mathbf{C}\Sigma_{\mathbf{X}}\mathbf{C}^T$$

but

$$\rho_{\mathbf{X}} = \mathbf{\Lambda}_{\mathbf{X}}^{-1/2}\Sigma_{\mathbf{X}}\mathbf{\Lambda}_{\mathbf{X}}^{-1/2}$$

So that if $\mathbf{C} = \mathbf{\Lambda}^{-1/2}$, so that $\mathbf{Y} = \mathbf{\Lambda}_{\mathbf{X}}^{-1/2}\mathbf{X}$ has the desired property.

Note that $\mathbf{\Lambda}_{\mathbf{X}}^{-1/2}$ is the diagonal matrix of the reciprocal of the standard deviations. So $Y_i = \frac{1}{\sigma_i}X_i, i = 1, 2, \dots, p$ are the un-centered normalized observations.

(Student answers could include other shifted versions)

Q 4. Suppose that $\mathbf{Y} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a 3-by-1 random vector where $\boldsymbol{\mu} = \mu \mathbf{j}$ and for some ρ ,

$$\boldsymbol{\Sigma}_{\mathbf{Y}} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix}$$

- (a) Write $\boldsymbol{\Sigma}$ as a linear combination of \mathbf{I} and \mathbf{J} .
 (b) For what values of ρ is $\boldsymbol{\Sigma}$ a valid variance-covariance matrix..
 (c) Compute the expectation $E \left[\sum_{i=1}^3 (Y_i - \bar{Y})^2 \right]$. Comment on this as a function of ρ .

(a) $\boldsymbol{\Sigma} = \sigma^2 [(1 - \rho)\mathbf{I} + \rho\mathbf{J}]$

(b) Checking Sylvester's criterion we have

$$\begin{aligned} \det(R_1) &= \sigma^2 > 0 && \text{for all } \rho \\ \det(R_2) &= \sigma^4(1 - \rho^2) > 0 && \text{if } |\rho| < 1 \\ \det(R_3) &= \sigma^6(1 - \rho^2)(1 + 2\rho) > 0 && \text{if } -1/2 < \rho < 1 \text{ or } \rho < -1 \end{aligned}$$

So $\boldsymbol{\Sigma}$ is positive definite only if $-1/2 < \rho < 1$. If $\rho = -1/2, 1$ then $\boldsymbol{\Sigma}$ is positive semi-definite. Thus $\boldsymbol{\Sigma}$ is a valid variance covariance matrix if and only if $-1/2 \leq \rho \leq 1$

(c) We should first write out the sum as a quadratic form. Namely $\mathbf{Y}^T(\mathbf{I} - 1/n\mathbf{J})\mathbf{Y} = \sum_{i=1}^3 (Y_i - \bar{Y})^2$
 Then by Theorem 2.4 we have ($n=3$)

$$\begin{aligned} E \left[\sum_{i=1}^3 (Y_i - \bar{Y})^2 \right] &= E [\mathbf{Y}^T(\mathbf{I} - 1/n\mathbf{J})\mathbf{Y}] \\ &= \text{tr}((\mathbf{I} - 1/n\mathbf{J})\sigma^2((1 - \rho)\mathbf{I} + \rho\mathbf{J})) + \boldsymbol{\mu}^T(\mathbf{I} - 1/n\mathbf{J})\boldsymbol{\mu} \\ &= \sigma^2 \text{tr} \left[(1 - \rho)\mathbf{I} - \frac{1}{n}(1 - \rho)\mathbf{J} + \rho\mathbf{J} - \frac{1}{n}\rho\mathbf{J}^2 \right] + \boldsymbol{\mu}^T\mathbf{I}\boldsymbol{\mu} - \frac{1}{n}\boldsymbol{\mu}^T\mathbf{J}\boldsymbol{\mu} \\ &= \sigma^2 \text{tr} \left[(1 - \rho)\mathbf{I} - \frac{1}{n}(1 - \rho)\mathbf{J} + \rho\mathbf{J} - \frac{1}{n}\rho n\mathbf{J} \right] + \boldsymbol{\mu}^T\boldsymbol{\mu} - \frac{1}{n}\boldsymbol{\mu}^T(n\boldsymbol{\mu}) \\ &= \sigma^2 \text{tr} \left[(1 - \rho)\mathbf{I} - \frac{1}{n}(1 - \rho)\mathbf{J} + \rho\mathbf{J} - \frac{1}{n}\rho n\mathbf{J} \right] + \boldsymbol{\mu}^T\boldsymbol{\mu} - \boldsymbol{\mu}^T\boldsymbol{\mu} \\ &= \sigma^2 \text{tr} \left[(1 - \rho)\mathbf{I} - \frac{1}{n}(1 - \rho)\mathbf{J} + \rho\mathbf{J} - \rho\mathbf{J} \right] \\ &= \sigma^2 \text{tr} \left[(1 - \rho)\mathbf{I} - \frac{1}{n}(1 - \rho)\mathbf{J} \right] \\ &= \sigma^2 \left[(1 - \rho)\text{tr}\mathbf{I} - \frac{1}{n}(1 - \rho)\text{tr}\mathbf{J} \right] \\ &= \sigma^2 \left[(1 - \rho)n - \frac{1}{n}(n)(1 - \rho) \right] \\ &= \sigma^2(1 - \rho)(n - 1) \\ &= 2\sigma^2(1 - \rho) \end{aligned}$$

This expectation and thus the mean of the sample variance is strictly decreasing to 0 as ρ increases to 1.

Q 5. Let \mathbf{X} be a 3-by-1 random vector with variance covariance matrix

$$\Sigma_{\mathbf{X}} = \begin{bmatrix} 5 & -4 & 1 \\ -4 & 13 & 9 \\ 1 & 9 & 10 \end{bmatrix}$$

- (a) Find a unit vector $\mathbf{c} = (c_1, c_2, c_3)^T$ such that the variance of $\mathbf{c}^T \mathbf{X}$ is maximized/minimized. (You may use a computer for convenience)
- (b) What is the largest/smallest possible variance for such linear combinations?

We want the eigenvectors associated with the largest/smallest eigenvalue. From R, they eigen vectors/values are

e_1	e_2	e_3
0.1543033	0.8017837	0.5773503
-0.7715167	-0.2672612	0.5773503
-0.6172134	0.5345225	-0.5773503

or

e_1	e_2	e_3
$1/\sqrt{42}$	$3/\sqrt{14}$	$1/\sqrt{3}$
$-5/\sqrt{42}$	$-1/\sqrt{14}$	$1/\sqrt{3}$
$-4/\sqrt{42}$	$2/\sqrt{14}$	$-1/\sqrt{3}$

with associated eigenvalues $\lambda_1 = 21, \lambda_2 = 7, \lambda_3 = 0$ (R will return $2.131628e - 14$ for the smallest, but this must be 0, since the determinant is 0)

Therefore, the constant vectors producing the maximum/minimum variances (21/0) of the linear combinations are $(1, -5, 4)/\sqrt{42}$ and $(1, 1, -1)/\sqrt{3}$ respectively.

(One can note then, that $X_3 = X_1 + X_2$ almost surely.)

Q 6. Johnson and Wichern 2.41

We have

$$(a) E(\mathbf{AX}) = \mathbf{AEX}$$

$$= \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ -2 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 7 & 3 \end{bmatrix}$$

$$(b) \text{Cov}(\mathbf{AX}) = \mathbf{ACovXA}^T$$

$$= \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & 1 \\ 0 & -2 & 1 \\ 0 & 0 & -3 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 1 & -2 & 0 \\ 1 & 1 & 1 & -3 \end{bmatrix} \begin{bmatrix} 3 & 3 & 3 \\ -3 & 3 & 3 \\ 0 & -6 & 3 \\ 0 & 0 & -9 \end{bmatrix}$$

$$= \begin{bmatrix} 6 & 0 & 0 \\ 0 & 18 & 0 \\ 0 & 0 & 36 \end{bmatrix}$$

(c) It is clear that the three pairs are all uncorrelated.

Q 7. Let $\mathbf{X} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Without using a computer compute

- (a) $|\boldsymbol{\Sigma}|$
- (b) $\text{tr}\boldsymbol{\Sigma}$
- (c) $V(\mathbf{e}_i^T \mathbf{X}), i = 1, 2, 3$
- (d) $\boldsymbol{\Sigma}$
- (e) $\boldsymbol{\Sigma}^{1/2}$

where $\lambda_1 = 3$, $\lambda_2 = 7$, and $\lambda_3 = 1$ are the eigenvalues of $\boldsymbol{\Sigma}$, and $\mathbf{e}_1^T = [1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3}]$, $\mathbf{e}_2^T = [2/\sqrt{6}, -1/\sqrt{6}, -1/\sqrt{6}]$, and $\mathbf{e}_3^T = [0, 1/\sqrt{2}, -1/\sqrt{2}]$ are the associated eigenvectors.

(a) $|\boldsymbol{\Sigma}| = \prod_{i=1}^3 \lambda_i = (7)(3)(1) = 21$

(b) $\text{tr}\boldsymbol{\Sigma} = \sum_{i=1}^3 \lambda_i = 7 + 3 + 1 = 11$

(c) We have

$$\begin{aligned} V(\mathbf{e}_i^T \mathbf{X}) &= \mathbf{e}_i^T \boldsymbol{\Sigma} \mathbf{e}_i \\ &= \mathbf{e}_i^T \left[\sum_{j=1}^3 \lambda_j \mathbf{e}_j \mathbf{e}_j^T \right] \mathbf{e}_i \\ &= \sum_{j=1}^3 \lambda_j \mathbf{e}_i^T \mathbf{e}_j \mathbf{e}_j^T \mathbf{e}_i \\ &= \sum_{j=1}^3 \lambda_j \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} \\ &= \lambda_i \end{aligned}$$

(d) We have

$$\begin{aligned} \boldsymbol{\Sigma} &= \sum_{j=1}^3 \lambda_j \mathbf{e}_j \mathbf{e}_j^T = 3 \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix} \\ &\quad + 7 \begin{bmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{bmatrix} \begin{bmatrix} 2/\sqrt{6} & -1/\sqrt{6} & -1/\sqrt{6} \end{bmatrix} + \begin{bmatrix} 0 \\ 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \\ &= 3/3 \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + 7/6 \begin{bmatrix} 4 & -2 & -2 \\ -2 & 1 & 1 \\ -2 & 1 & 1 \end{bmatrix} + 1/2 \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} \\ &= \frac{1}{3} \begin{bmatrix} 17 & -4 & -4 \\ -4 & 8 & 5 \\ -4 & 5 & 8 \end{bmatrix} \end{aligned}$$

(e) We have

$$\begin{aligned}
 \Sigma^{1/2} &= \sum_{i=j}^3 \sqrt{\lambda_j} \mathbf{e}_j \mathbf{e}_j^T = \sqrt{3} \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix} \\
 &+ \sqrt{7} \begin{bmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{bmatrix} \begin{bmatrix} 2/\sqrt{6} & -1/\sqrt{6} & -1/\sqrt{6} \end{bmatrix} + \begin{bmatrix} 0 \\ 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \\
 &= \sqrt{3}/3 \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \sqrt{7}/6 \begin{bmatrix} 4 & -2 & -2 \\ -2 & 1 & 1 \\ -2 & 1 & 1 \end{bmatrix} + 1/2 \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} \\
 &= \frac{1}{6} \left\{ \begin{bmatrix} 2\sqrt{3} & 2\sqrt{3} & 2\sqrt{3} \\ 2\sqrt{3} & 2\sqrt{3} & 2\sqrt{3} \\ 2\sqrt{3} & 2\sqrt{3} & 2\sqrt{3} \end{bmatrix} + \begin{bmatrix} 4\sqrt{7} & -2\sqrt{7} & -2\sqrt{7} \\ -2\sqrt{7} & \sqrt{7} & \sqrt{7} \\ -2\sqrt{7} & \sqrt{7} & \sqrt{7} \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 3 & -3 \\ 0 & -3 & 3 \end{bmatrix} \right\} \\
 &= \frac{1}{6} \begin{bmatrix} 2\sqrt{3} + 4\sqrt{7} & 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} - 2\sqrt{7} \\ 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} + \sqrt{7} + 3 & 2\sqrt{3} + \sqrt{7} - 3 \\ 2\sqrt{3} - 2\sqrt{7} & 2\sqrt{3} + \sqrt{7} - 3 & 2\sqrt{3} + \sqrt{7} + 3 \end{bmatrix}
 \end{aligned}$$

Q 8. Let $\mathbf{X} \sim N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\mu}^T = [1, 0, -2]$ and

$$\boldsymbol{\Sigma} = \begin{bmatrix} 2 & -0.5 & 1 \\ -0.5 & 1 & 0 \\ 1 & 0 & 2 \end{bmatrix}$$

- Which components of \mathbf{X} are independent?
- Find the conditional distribution of $X_2|(X_1, X_3)^T$
- Find the conditional distribution of $(X_1, X_3)^T|X_2$
- Comment on how your answer to part (a) affects parts (b) and (c).

(a) Checking the covariances clearly only X_2 and X_3 are independent.

Using the appropriate theorem for (b) and (c)

(b) We want $X_2|(X_1, X_3)^T = (a, c)^T$. Then

$$\begin{aligned} E\left(X_2 \mid \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} = \begin{bmatrix} a \\ c \end{bmatrix}\right) &= 0 + [-0.5 \ 0] \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \left(\begin{bmatrix} a \\ c \end{bmatrix} - \begin{bmatrix} 1 \\ -2 \end{bmatrix} \right) \\ &= \frac{1}{3} [-0.5 \ 0] \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} a-1 \\ c+2 \end{bmatrix} \\ &= \frac{1}{3} [-1 \ 0.5] \begin{bmatrix} a-1 \\ c+2 \end{bmatrix} \\ &= \frac{1}{3} [-1 \ 0.5] \begin{bmatrix} a-1 \\ c+2 \end{bmatrix} \\ &= \frac{1}{3} [(1-a) + 1/2(c+2)] = \frac{1}{3} [2 + 1/2c - a] \\ &= 2/3 + 1/6c - 1/3a \end{aligned}$$

$$\begin{aligned} \text{Cov}\left(X_2 \mid \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} = \begin{bmatrix} a \\ c \end{bmatrix}\right) &= 1 - [-0.5 \ 0] \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} \\ &= 1 - \frac{1}{3} [-0.5 \ 0] \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} \\ &= 1 - \frac{1}{3} [-0.5 \ 0] \begin{bmatrix} -1 \\ 0.5 \end{bmatrix} \\ &= 1 - \frac{1}{3} 0.5 = 5/6 \end{aligned}$$

So that $X_2|(X_1, X_3)^T = (a, c)^T \sim N(2/3 + 1/6c - 1/3a, 5/6)$

(c) We want $(X_1, X_3)^T | X_2 = b$. Then

$$\begin{aligned} E\left(\begin{bmatrix} X_1 \\ X_3 \end{bmatrix} | X_2 = b\right) &= \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} [1]^{-1} [b - 0] \\ &= \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} [b] \\ &= \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} -0.5b \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 - 1/2b \\ -2 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \text{Cov}\left(\begin{bmatrix} X_1 \\ X_3 \end{bmatrix} | X_2 = b\right) &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} [1]^{-1} [-0.5 \quad 0] \\ &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} [-0.5 \quad 0] \\ &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} 1/4 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 7/4 & 1 \\ 1 & 2 \end{bmatrix} \end{aligned}$$

So that $(X_1, X_3)^T | X_2 = b \sim N_2\left(\begin{bmatrix} 1 - 1/2b \\ -2 \end{bmatrix}, \begin{bmatrix} 7/4 & 1 \\ 1 & 2 \end{bmatrix}\right)$

(d) In part (b) the conditional mean is affected by both values of X_1 and X_3 even though X_3 is independent of X_1 because X_1 and X_3 are correlated. In part (c) only the conditional mean of X_1 is changed as X_2 and X_3 are independent.

In part (b) the variance is reduced by a factor depending on the conditioning (but not the value) of both X_1 and X_3 since they are correlated. In part (c), only the variance of X_1 is affected and not the covariance or variance of X_2 .

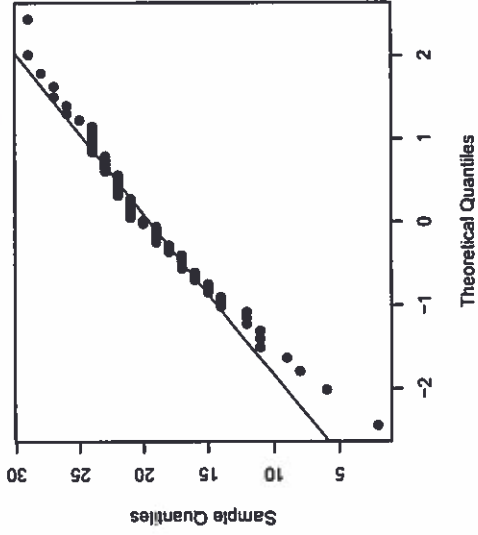
Q 9. Using the data in Table 4.6 from the book (T4-6), check for univariate and multivariate normality for the females. Comment.

The normal QQ plots look reasonably straight with some mild curvature in Supp and Leader but nothing exaggerated. However, the multivariate chisquared qq plot definitely has some stronger curvature indicating multivariate normality might be violated.

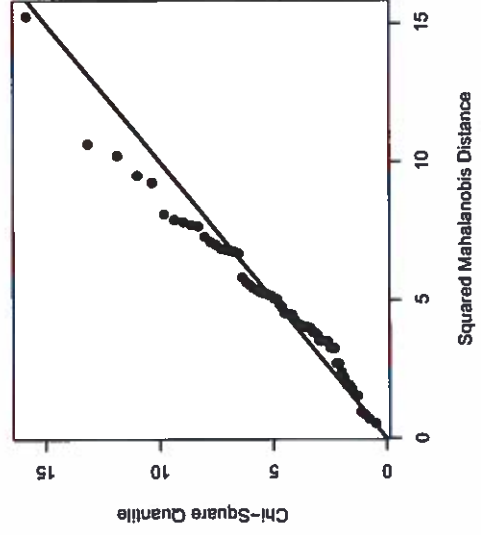
(As a side note, anderson-darling tests applied to the marginal distributions give a p-value of about 0.02 for Supp and 0.001 for Leader, and Royston's test for Multivariate Normality gives a p-value of about 0.01.

```
# -----  
#  
# Import functions  
#  
# -----  
  
library(mvnTest) # For the Royston test & joint chisq plot  
  
# -----  
#  
# Import data and draw plots  
#  
# -----  
  
psych <- read.table(file = "datasets/Wichern_data/T4-6.dat")  
names(psych) <- c("Indep", "Supp", "Benev", "Conform", "Leader", "Gender", "Socio")  
  
is.female <- which(psych$Gender == 2)  
psych.f <- psych[is.female, ]  
  
dev.new(width = 15, height = 10)  
par(mfrow = c(2, 3))  
  
for(j in c("Indep", "Supp", "Benev", "Conform", "Leader")){  
  qqnorm(psych.f[,j], main = paste("Normal QQ plot for", j), pch = 19)  
  qqline(psych.f[,j])  
}  
  
R.test(psych.f[,c("Indep", "Supp", "Benev", "Conform", "Leader")], qqplot = TRUE)
```

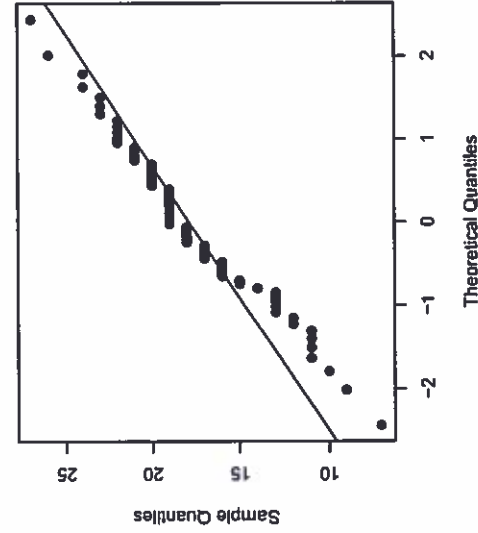
Normal QQ plot for Benev



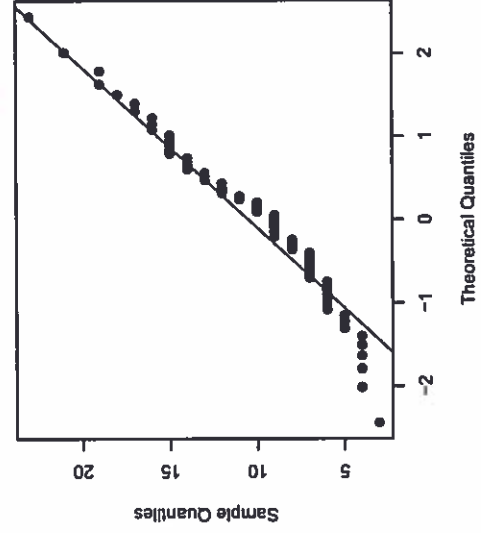
Chi-Square Q-Q Plot



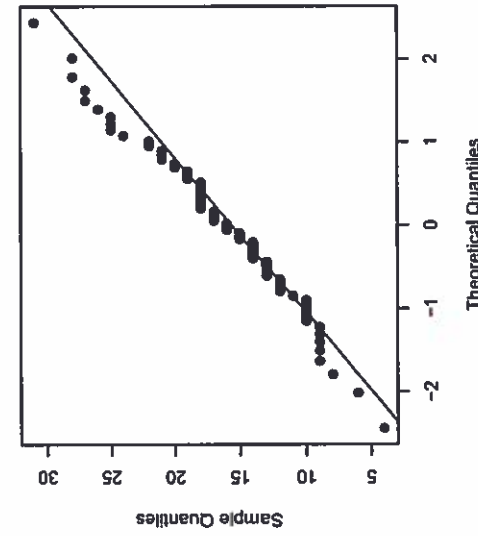
Normal QQ plot for Supp



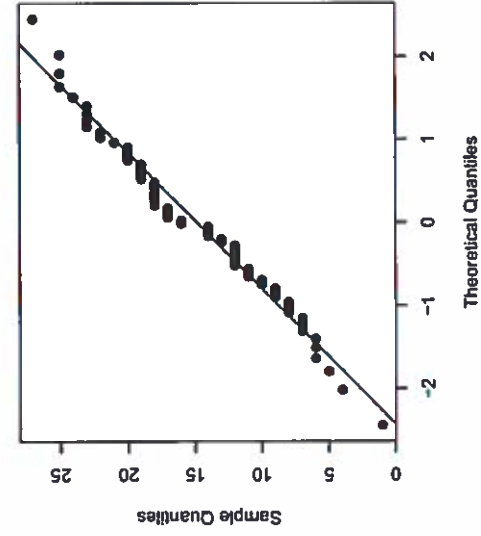
Normal QQ plot for Leader



Normal QQ plot for Indep



Normal QQ plot for Conform



Squared Mahalanobis Distance



Q 10. Let $\mathbf{X}_1 \sim N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$ and $\mathbf{X}_2 \sim N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$. Given independent random samples of sizes n_1 and n_2 respectively, obtain the MLE's of $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, and $\boldsymbol{\Sigma}$.

The likelihood is

$$\begin{aligned} \ell(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma} | \mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}, \mathbf{x}_{21}, \dots, \mathbf{x}_{2n_2}) \\ &= \left[\prod_{j=1}^{n_1} \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_{1j} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{1j} - \boldsymbol{\mu}_1)} \right] \left[\prod_{j=1}^{n_2} \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_{2j} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{2j} - \boldsymbol{\mu}_2)} \right] \\ &= \frac{1}{(2\pi)^{(n_1+n_2)p/2} |\boldsymbol{\Sigma}|^{(n_1+n_2)/2}} e^{-\frac{1}{2}[\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{1j} - \boldsymbol{\mu}_1) + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{2j} - \boldsymbol{\mu}_2)]} \end{aligned}$$

By a theorem from class we have

$$\begin{aligned} &= \frac{1}{(2\pi)^{(n_1+n_2)p/2} |\boldsymbol{\Sigma}|^{(n_1+n_2)/2}} \times \\ &e^{-\frac{1}{2}[\text{tr}(\boldsymbol{\Sigma}^{-1} \{ \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T \}) + n_1(\bar{\mathbf{x}}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_1 - \boldsymbol{\mu}_1) + n_2(\bar{\mathbf{x}}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_2 - \boldsymbol{\mu}_2)]} \\ &= \left[\frac{1}{(2\pi)^{(n_1+n_2)p/2}} \right] \left[\frac{1}{|\boldsymbol{\Sigma}|^{(n_1+n_2)/2}} e^{-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \{ \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T \})} \right] \times \\ &\quad \left[e^{-\frac{1}{2} n_1 (\bar{\mathbf{x}}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_1 - \boldsymbol{\mu}_1)} \right] \left[e^{-\frac{1}{2} n_2 (\bar{\mathbf{x}}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}}_2 - \boldsymbol{\mu}_2)} \right] \end{aligned}$$

All 4 products are strictly positive, the first being a constant and thus ignorable. The third and fourth are clearly maximized (uniquely since $\boldsymbol{\Sigma}$ is positive definite and thus $\boldsymbol{\Sigma}^{-1}$ is positive definite) when the quadratic form is 0, which occurs if $\hat{\boldsymbol{\mu}}_1 = \bar{\mathbf{x}}_1$ and $\hat{\boldsymbol{\mu}}_2 = \bar{\mathbf{x}}_2$. The second term is maximized by applying the theorem from class with

$$\begin{aligned} b &= 1/2(n_1 + n_2) \\ B &= \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T \end{aligned}$$

yielding

$$\begin{aligned} \hat{\boldsymbol{\Sigma}} &= \frac{1}{n_1 + n_2} \left[\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T \right] \\ &= \frac{1}{n_1 + n_2} [(n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2] \end{aligned}$$

Q 11. Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- (a) Find the mle of k - a scalar- where $\boldsymbol{\Sigma} = k\boldsymbol{\Sigma}_0$ and both $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_0$ are known.
 (b) Find the mle of k - a scalar- where $\boldsymbol{\mu} = k\boldsymbol{\mu}_0$ and both $\boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma}$ are known.

(a) The likelihood as a function of k is

$$\begin{aligned} L(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}_0) &= \frac{1}{(2\pi)^{np/2} |k\boldsymbol{\Sigma}_0|^{n/2}} e^{-\frac{1}{2} [\sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T [k\boldsymbol{\Sigma}_0]^{-1} (\mathbf{x}_j - \boldsymbol{\mu})]} \\ &= \left[\frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}_0|^{n/2}} \right] \left[\frac{1}{k^{np/2}} e^{-\frac{1}{2k} [\frac{1}{2} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x}_j - \boldsymbol{\mu})]} \right] \end{aligned}$$

With log-likelihood

$$\begin{aligned} \ell(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}_0) &= \ln L(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}_0) \\ &= \text{constant} - \frac{np}{2} \ln k - \frac{1}{2k} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) \end{aligned}$$

with derivative

$$\frac{d\ell}{dk} = -\frac{np}{2k} + \frac{1}{2k^2} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x}_j - \boldsymbol{\mu})$$

setting to 0 and solving yields

$$\widehat{k} = \frac{1}{np} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x}_j - \boldsymbol{\mu})$$

which is a (local) maximum because

$$\begin{aligned} \left. \frac{d^2\ell}{dk^2} \right|_{k=\widehat{k}} &= \frac{np}{2\widehat{k}^2} - \frac{1}{\widehat{k}^3} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) \\ &= -\frac{np^2}{2\widehat{k}} < 0 \end{aligned}$$

$(\widehat{k} - 2)^2 = 22$

(b)

$$L(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} e^{-\frac{1}{2} [\sum_{j=1}^n (\mathbf{x}_j - k\boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_j - k\boldsymbol{\mu}_0)]}$$

With log-likelihood

$$\begin{aligned} \ell(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) &= L(k|\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) \\ &= \text{constant} - \frac{1}{2} \left[\sum_{j=1}^n (\mathbf{x}_j - k\boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_j - k\boldsymbol{\mu}_0) \right] \\ &= \text{constant} - \frac{1}{2} \left[\sum_{j=1}^n (\mathbf{x}_j^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_j - 2k\boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_j + k^2 \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0) \right] \\ &= \text{constant} - \frac{1}{2} \left[\sum_{j=1}^n (-2k\boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_j + k^2 \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0) \right] \\ &= \text{constant} - \frac{1}{2} \left[nk^2 \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 - 2 \sum_{j=1}^n k\boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_j \right] \end{aligned}$$

The derivative is

$$\frac{d\ell}{dk} = -nk\mu_0^T \Sigma^{-1} \mu_0 + \sum_{j=1}^n \mu_0 \Sigma^{-1} \mathbf{x}_j$$

setting to 0 and solving yields

$$\begin{aligned} \hat{k} &= \frac{\sum_{j=1}^n \mu_0 \Sigma^{-1} \mathbf{x}_j}{n\mu_0^T \Sigma^{-1} \mu_0} \\ &= \frac{\mu_0 \Sigma^{-1} \bar{\mathbf{x}}}{\mu_0^T \Sigma^{-1} \mu_0} \end{aligned}$$

which is a global maximum because

$$\frac{d^2\ell}{dk^2} = -n\mu_0^T \Sigma^{-1} \mu_0 < 0$$

provided $\mu_0 \neq \mathbf{0}$. (Of course if $\mu_0 = \mathbf{0}$ then estimating the scale constant k is not possible.)



MATH 755 HW 3

Due 10/04/2018

EX 5.67 / 1. Johnson & Wichern 5.1

2. Johnson & Wichern 5.18 (a)

Let $X = \begin{bmatrix} 2 & 10 \\ 8 & 9 \\ 6 & 9 \\ 8 & 10 \end{bmatrix}$ Test $H_0: \mu = [7 \ 11]$
 a) Find T^2 for testing $\mu = [7 \ 11]$
 b) Specify the dist of T^2 in a) $\mu = [7 \ 11]$
 c) Use a) and b) at $\alpha = 0.05$. ? conclusion

```
college <- read.table(file = "Wichern_data/T5-2.dat")
```

a) Test: $\mu = \begin{bmatrix} 500 \\ 30 \end{bmatrix}$ and $H_1: \mu \neq \begin{bmatrix} 500 \\ 30 \end{bmatrix}$ at $\alpha = 0.05$ b) Length and direction of the axes of the 95% CI for μ .

3. Johnson & Wichern 5.19

c) Q&Q plot

```
lumber <- read.table(file = "Wichern_data/T5-11.dat")
names(lumber) <- c("stiff", "bend")
```

a) Confidence ellipse for $\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ $\mu_1 = E(X_1)$
 $\mu_2 = E(X_2)$

not finished drawing standard ellipse

4. Johnson & Wichern 6.24

```
egypt <- read.table(file = "Wichern_data/T6-13.dat")
names(egypt) <- c("MaxBreath", "BasHeight", "BasLength", "NasHeight", "TimePeriod")
```

5. Find a dataset online which contains (i) two factors and (ii) at least two continuous responses.

- (a). Run a one-way MANOVA using only one factor. State your conclusions in plain language.
- (b). Run a two-way MANOVA with interactions. State your conclusions in plain language. Refit a two-way MANOVA without interactions if the interactions are not significant. State your conclusions in plain language.
- (c). One assumption we made in MANOVA is that the variance-covariance matrices of each group of residuals are equal. Does this assumption hold in your dataset?

1
d



MAT 755 Midterm Fall 2018

Due: October 16, 2018

Instructions: This is a take-home exam. It will weight 15% of your final grade. You should work on all the questions independently. If you submit your work electronically (Blackboard, email), please name your file as "lastname_midterm_F18".

1. Show that

$$Q = \begin{pmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{pmatrix}$$

is an orthogonal matrix.

2. Using the matrix

$$Q = \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix}$$

(a). Calculate QQ^T and obtain its eigenvalues and eigenvectors.

(b). Calculate Q^TQ and obtain its eigenvalues and eigenvectors. Check that the nonzero eigenvalues are the same as those in part (a).

(c). Obtain the singular value decomposition of Q .

3. Consider an arbitrary $n \times p$ matrix Q .

(a). Show that Q^TQ is a symmetric $p \times p$ matrix.

(b). Show that Q^TQ is necessarily nonnegative definite.

$X_1 = \text{lather}$

$X_2 = \text{mildness}$

2 way 8.

way 1: $\bar{x}_1 = \begin{bmatrix} 8 \\ 4 \end{bmatrix} S_1$

way 2: $\bar{x}_2 = \begin{bmatrix} 10 \\ 3 \end{bmatrix} S_2$

4. Bars of soap are manufactured in each of two ways. The characteristics $X_1 = \text{lather}$ and $X_2 = \text{mildness}$ are measured. The summary statistics for 61 and 61 bars produced by methods 1 and 2, respectively, are

$$\bar{x}_1 = [8, 4]^T, \bar{x}_2 = [10, 3]^T, S_1 = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}, S_2 = \begin{pmatrix} 7 & 2 \\ 2 & 7 \end{pmatrix}.$$

(a). Test the hypothesis that the population means $\mu_1 = \mu_2$ at 10% significance level, using the pooled covariance matrix to estimate the population covariance matrix Σ and assuming that both populations are with multivariate normal distribution.

(b). Find the 90% Bonferroni simultaneous confidence interval for the individual mean difference.

5. Wichern and Johnson 7.24 on page 425

```
bull <- read.table(file="Wichern_data/T1-10.dat", header=F)
names(bull) <- c("Breed", "SalePr", "YrHgt", "FtFrBody", "PrctFFB", "Frame", "BkFat", "SaleHt", "SaleWt")
```

6. A researcher randomly assigns 33 subjects to one of three groups.

- Group 1 receives technical dietary information interactively from an on-line website.
- Group 2 receives the same information from a nurse practitioner.
- Group 3 receives the information from a video tape made by the same nurse practitioner.

Each subject then made three ratings: difficulty, usefulness, and importance of the information in the presentation. The researcher is interested in whether or not there is a difference in the modes of presentation. (Data: presentation.sav)

(a). Test the equality of the group means.

(b). What assumptions do you make in part (a)? Assess the validity of the assumptions.

(c). Construct and comment the 95% Bonferroni confidence intervals for differences in mean components for the three responses for each pair of populations.

17 Show that

$$Q = \begin{pmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{pmatrix} \text{ is an orthogonal matrix}$$

* We have Q is an orthogonal matrix iff

$$\begin{cases} \textcircled{1} Q \text{ is a square matrix (done)} \\ \textcircled{2} Q^T Q = Q Q^T = I \end{cases}$$

$$Q^T Q = \begin{pmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{pmatrix} \begin{pmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{pmatrix} = I$$

$$Q Q^T = \begin{pmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{pmatrix} \begin{pmatrix} \frac{5}{13} & -\frac{12}{13} \\ \frac{12}{13} & \frac{5}{13} \end{pmatrix} = I$$

} \Rightarrow $\textcircled{2}$ has been done

$\textcircled{1}$ and $\textcircled{2} \Rightarrow Q$ is an orthogonal matrix

2) Using the matrix $Q = \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix}$

1) Calculate QQ^T and obtain its eigenvalues and eigenvectors

2) Calculate Q^TQ and obtain its eigenvalues and eigenvectors.

Check that the nonzero eigenvalues are the same as those in part a)

3) Obtain the singular value decomposition of Q

$$1) \quad QQ^T = \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix} \begin{pmatrix} 4 & 3 \\ 8 & 6 \\ 8 & 9 \end{pmatrix} = \begin{pmatrix} 144 & 132 \\ 132 & 126 \end{pmatrix} = A$$

Characteristic polynomial

$$\det(A - \lambda I) = \det \begin{pmatrix} 144 - \lambda & 132 \\ 132 & 126 - \lambda \end{pmatrix} = (144 - \lambda)(126 - \lambda) - 132^2 =$$

$$= \lambda^2 - 270\lambda + 720$$

$$\Rightarrow \text{Eigenvalues } \lambda_{1,2} = 135 \pm 3\sqrt{1245}$$

Use the matrix

$$Q = \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix}$$

a) Calculate QQ^T and obtain its eigenvalues and eigenvectors.

$$* QQ^T = \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix} \begin{pmatrix} 4 & 3 \\ 8 & 6 \\ 8 & 9 \end{pmatrix} = \begin{pmatrix} 144 & 132 \\ 132 & 126 \end{pmatrix}$$

* Find eigenvalues:

• Characteristic polynomial

$$\det(QQ^T - \lambda I) = \det \begin{pmatrix} 144 - \lambda & 132 \\ 132 & 126 - \lambda \end{pmatrix} = (144 - \lambda)(126 - \lambda) - 132^2 = \lambda^2 - 270\lambda + 720$$

$$\text{then } \lambda_{1,2} = 3(45 \pm \sqrt{1945}) \Leftrightarrow \begin{cases} \lambda_1 = 267.306 \\ \lambda_2 = 2.694 \end{cases}$$

* Find eigenvectors:

• Find eigenvectors associated with $\lambda_1 = 3(45 + \sqrt{1945}) = 267.306$

the eigenvector $y = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$ is a solution of $QQ^T y = \lambda_1 y \Leftrightarrow (QQ^T - \lambda_1 I) y = 0$

$$(QQ^T - \lambda_1 I) = \begin{pmatrix} 144 - 3 \cdot 45 - 3 \cdot \sqrt{1945} & 132 & | & 0 \\ 132 & 126 - 3 \cdot 45 - 3 \cdot \sqrt{1945} & | & 0 \end{pmatrix}$$

$$R_1 \rightarrow \frac{R_1}{3} \Leftrightarrow \begin{pmatrix} 48 - 45 - \sqrt{1945} & 44 & | & 0 \\ 44 & 42 - 45 - \sqrt{1945} & | & 0 \end{pmatrix}$$

$$R_2 \rightarrow \frac{R_2}{3} \Leftrightarrow \begin{pmatrix} 3 - \sqrt{1945} & 44 & | & 0 \\ 44 & -3 - \sqrt{1945} & | & 0 \end{pmatrix}$$

$$\Rightarrow \begin{cases} R_1 \Rightarrow \begin{cases} u_2 = \frac{-(3 - \sqrt{1945})u_1}{44} \\ u_1^2 + u_2^2 = 1 \end{cases} \Rightarrow \begin{cases} u_1 = -0.7307 \\ u_2 = -0.6826 \end{cases} \end{cases}$$

• Find eigenvector of QQ^T associated with $\lambda_2 = 3 \cdot 45 - 3\sqrt{1945}$

$$QQ^T - \lambda_2 I = 0 \Leftrightarrow \begin{pmatrix} 144 - 3 \cdot 45 + 3\sqrt{1945} & 132 & | & 0 \\ 132 & 126 - 3 \cdot 45 + 3\sqrt{1945} & | & 0 \end{pmatrix}$$

$$R_1 \rightarrow \frac{R_1}{3} \Leftrightarrow \begin{pmatrix} 48 - 45 + \sqrt{1945} & 44 & | & 0 \\ 44 & 42 - 45 + \sqrt{1945} & | & 0 \end{pmatrix}$$

$$\Rightarrow \begin{bmatrix} 3 + \sqrt{1945} & 44 & | & 0 \\ 44 & -3 + \sqrt{1945} & | & 0 \end{bmatrix} \Rightarrow \begin{cases} u_2 = \frac{-(3 + \sqrt{1945})}{44} \\ u_1^2 + u_2^2 = 1 \end{cases} \Rightarrow \begin{cases} u_1 = 0.68263 \\ u_2 = -0.730762 \end{cases}$$



1
= 0

Calculate $Q^T Q$ and obtain its eigenvalues and eigenvectors.
 Check that the nonzero eigenvalues are the same as those in part a.

$$Q^T Q = \begin{pmatrix} 4 & 3 \\ 8 & 6 \\ 8 & 9 \end{pmatrix} \begin{pmatrix} 4 & 8 & 8 \\ 3 & 6 & 9 \end{pmatrix} = \begin{pmatrix} 25 & 50 & 59 \\ 50 & 100 & 118 \\ 59 & 118 & 145 \end{pmatrix}$$

* Find eigenvalues of $Q^T Q$

• Characteristic polynomial

$$\det(Q^T Q - \lambda I) = \det \begin{pmatrix} 25-\lambda & 50 & 59 \\ 50 & 100-\lambda & 118 \\ 59 & 118 & 145-\lambda \end{pmatrix} \xrightarrow{R_2 - 2R_1} \det \begin{pmatrix} 25-\lambda & 2\lambda & 59 \\ 50 & -\lambda & 118 \\ 59 & 118 & 145-\lambda \end{pmatrix}$$

$$\xrightarrow{R_2 - 2R_1 + R_1} \det \begin{pmatrix} 25-\lambda & 2\lambda & 59 \\ 125-\lambda & 0 & 295 \\ 59 & 0 & 145-\lambda \end{pmatrix} = -2\lambda \det \begin{pmatrix} 125-\lambda & 295 \\ 59 & 145-\lambda \end{pmatrix}$$

$$= -2\lambda (\lambda^2 - 270\lambda + 720)$$

$$= \det(Q Q^T - \lambda I) \text{ from a)}$$

⇒ the nonzero eigenvalues are the same as those in part a

Some have eigenvalues of $Q^T Q$ are

$$\begin{cases} \lambda_1 = 3(45 + \sqrt{1945}) = 267.306 \\ \lambda_2 = 3(45 - \sqrt{1945}) = 2.6935 \\ \lambda_3 = 0 \end{cases}$$

* Find eigenvector of $Q^T Q$ associated with $\lambda_1 = 3(45 + \sqrt{1945})$.

the eigenvector $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$ satisfies $Q^T Q \mathbf{x} = \lambda_1 \mathbf{x}$
 $\Leftrightarrow (Q^T Q - \lambda_1 I) \mathbf{x} = \mathbf{0}$

$$\Rightarrow \begin{pmatrix} -242.306 & 50 & 59 \\ 50 & -167.306 & 118 \\ 59 & 118 & -112.306 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \Rightarrow \text{unit eigenvector} \begin{pmatrix} -0.304 \\ -0.608 \\ 0.733 \end{pmatrix}$$

* Find eigenvector of $Q^T Q$ associated with $\lambda_2 = 3(45 - \sqrt{1945})$

$$\begin{pmatrix} 22.3065 & 50 & 59 \\ 50 & 97.3065 & 118 \\ 59 & 118 & 142.3065 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \Rightarrow \text{unit eigenvector} \begin{pmatrix} -0.3279 \\ -0.6559 \\ 0.6798 \end{pmatrix}$$

* Find eigenvector associated with $\lambda_3 = 0$

$$\begin{pmatrix} 25 & 50 & 59 \\ 50 & 100 & 118 \\ 59 & 118 & 145 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \iff \begin{pmatrix} 25 & 50 & 59 & | & 0 \\ 0 & 0 & 0 & | & 0 \\ 9 & 18 & 27 & | & 0 \end{pmatrix} \rightarrow \begin{pmatrix} 25 & 50 & 59 & | & 0 \\ 0 & 0 & 0 & | & 0 \\ 9 & 18 & 27 & | & 0 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \begin{pmatrix} -2v_2 \\ v_2 \\ 0 \end{pmatrix} \Rightarrow \text{corresponding unit vector } \begin{pmatrix} -2/\sqrt{5} \\ 1/\sqrt{5} \\ 0 \end{pmatrix}$$

c) Find SVD of Q

We have the singular values $\Rightarrow \sigma_1 = \sqrt{267.36} = 16.3511$ $\sigma_2 = \sqrt{2.6935} = 1.6411$

We have the SVD of Q is $Q = U \Sigma V^T$ where U and V are orthogonal matrices
 Σ : diagonal matrix

So then:

$$\textcircled{1} Q Q^T = (U \Sigma V^T)(V \Sigma U^T) = U \Sigma \underbrace{V^T V}_{=I} \Sigma U^T$$

$\Rightarrow u_i$ is eigen vector associative with eigenvalue of $Q Q^T$

$$\textcircled{2} Q^T Q = (V \Sigma U^T)(U \Sigma V^T) = V \Sigma \underbrace{U^T U}_{=I} \Sigma V^T = V \Sigma \Sigma V^T$$

$\Rightarrow v_i$ is eigen vector associative with eigenvalue of $Q^T Q$.

\Rightarrow We have the SVD of Q is:

$$Q = U \Sigma V^T \text{ where } U = \begin{pmatrix} -0.7307 & 0.6826 \\ -0.6826 & -0.7307 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 16.3511 & 0 \\ 0 & 1.6411 \end{pmatrix}$$

$$V = \begin{pmatrix} -0.304 & -0.3279 \\ -0.608 & -0.6559 \\ -0.733 & 0.6798 \end{pmatrix}$$

Consider arbitrary $n \times p$ matrix Q

a) Show that $Q^T Q$ is a symmetric $p \times p$ matrix

b) Show that $Q^T Q$ is necessarily nonnegative definite

a)

* First, we have that $\underbrace{Q^T}_{p \times n} \underbrace{Q}_{n \times p}$ is a $p \times p$ matrix
square

* Second, we have that:

$$\text{Put } A = Q^T Q$$

$$\Rightarrow A^T = (Q^T Q)^T = Q^T Q^{TT} = Q^T Q = A$$

$\Rightarrow A = Q^T Q$ is a square matrix
and $A = A^T$
 $\Rightarrow A$ is symmetric

b) Since $Q^T Q$ is symmetric, we have that $x^T Q^T Q x$ is well defined and we want to prove that $x^T Q^T Q x \geq 0, \forall x \neq 0, \mathbb{R}^p$

* Put $y = Qx$, then we have $y^T y = (Qx)^T (Qx) = x^T Q^T Q x$ (1)

• Since $y^T y = \sum_{i=1}^n y_i^2 \geq 0$ when $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$

equality happens when $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$

when x is a solution of $x = 0$

$$(1) + (2) \Rightarrow x^T Q^T Q x \geq 0, \forall x \neq 0$$

equality happens when $Qx = 0$

$\Rightarrow Q^T Q$ is necessarily nonnegative definite

1

2

3

MAT 755 HW4

Your names

Due: Tuesday, Oct. 23

1. Johnson & Wichern 8.10

```
stocks <- read.table("Wichern_data/T8-4.dat",  
  col.names=c("JPM", "Citibank", "Wells", "RoyalDutchShell", "ExxonMobil"))
```

2. Johnson & Wichern 8.22

```
bulls<-read.table("Wichern_data/T1-10.dat")  
names(bulls) <- c("Breed", "SalePr", "YrHgt", "FtFrBody", "PrctFFB", "Frame", "BkFat", "SaleHt", "SaleWt")  
bulls<-bulls[,3:9]
```

3.

Let $\mathbf{X} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\mu} = [1, 2]^T$ and

$$\boldsymbol{\Sigma} = \begin{bmatrix} 5/2 & 1 \\ 1 & 5/2 \end{bmatrix}$$

- Find the PCA factor loadings
- Find total variance and proportion of variance explained by each principal component
- Find the correlations between the principal components Y_i to the raw variables X_j
- Compute the principal components for the point $\mathbf{x} = [1 + \sqrt{2}, 2 - \sqrt{2}]^T$.

1

2

3

4

3) Let $X \sim (\mu, \Sigma)$ $\mu = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\Sigma = \begin{bmatrix} 5/2 & 1 \\ 1 & 5/2 \end{bmatrix}$

a) Find the PCA factor loadings

b) Find total variance and proportion of variance explained by each principal component

c) Find the correlation between the principal component Y_i to the raw variables X_j

d) Compute the principal components for the point $x = \begin{bmatrix} 1 + \sqrt{2} \\ 2 - \sqrt{2} \end{bmatrix}$

a) * We first want to find eigenvalues λ_1, λ_2 and eigenvector \vec{a}_1 and \vec{a}_2 of Σ .

• Characteristic polynomial:

$$\det(\Sigma - \lambda I) = \det \begin{pmatrix} 5/2 - \lambda & 1 \\ 1 & 5/2 - \lambda \end{pmatrix} = (5/2 - \lambda)^2 - 1 = \lambda^2 - 5\lambda + \frac{21}{4}$$

$$\det(\Sigma - \lambda I) = 0 \Rightarrow \lambda_1 = \frac{7}{2} \quad \lambda_2 = \frac{3}{2}$$

• Find eigenvector $\vec{a}_1 = \begin{pmatrix} a_{11} \\ a_{21} \end{pmatrix}$, \vec{a}_1 is a solution of $\Sigma \vec{a}_1 = \lambda_1 \vec{a}_1 \Rightarrow (\Sigma - \lambda_1 I) \vec{a}_1 = \vec{0}$.

$$[\Sigma - \lambda_1 I \mid \vec{0}] \Leftrightarrow \left[\begin{array}{cc|c} 5/2 - 7/2 & 1 & 0 \\ 1 & 5/2 - 7/2 & 0 \end{array} \right] \Rightarrow \left[\begin{array}{cc|c} -1 & 1 & 0 \\ 1 & -1 & 0 \end{array} \right] \Rightarrow \vec{a}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

• Find eigenvector \vec{a}_2 associative with $\lambda_2 = 3/2$

$$[\Sigma - \lambda_2 I \mid \vec{0}] = \left[\begin{array}{cc|c} 5/2 - 3/2 & 1 & 0 \\ 1 & 5/2 - 3/2 & 0 \end{array} \right] \Rightarrow \left[\begin{array}{cc|c} 1 & 1 & 0 \\ 1 & 1 & 0 \end{array} \right] \Rightarrow \vec{a}_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$$

\Rightarrow The PCA factor loadings are $\vec{a}_1 = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$ and $\vec{a}_2 = \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix}$

b) * The total variance $\lambda_1 + \lambda_2 = \frac{7}{2} + \frac{3}{2} = 5$

* The proportion of variance explained by the first principal component, W_1 is

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{\frac{7}{2}}{\frac{10}{2}} = \frac{7}{10}$$

* The proportion of variance explained by the second pc, W_2 is

$$\frac{\lambda_2}{\lambda_1 + \lambda_2} = \frac{\frac{3}{2}}{\frac{10}{2}} = \frac{3}{10}$$

→ Find the correlation between the principal component Y_i and raw variable X_j

We have $\begin{bmatrix} \vec{a}_1 \\ \vec{a}_2 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$ $\lambda_1 = \frac{7}{2}$ $\lambda_2 = \frac{3}{2}$ $\Sigma = \begin{bmatrix} 5/2 & 1 \\ 1 & 5/2 \end{bmatrix}$.

Then $\rho_{Y_1, X_1} = \frac{a_{11} \sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = a_{11} = \frac{1/\sqrt{2} \sqrt{7/2}}{\sqrt{5/2}} = \frac{\sqrt{7}}{\sqrt{10}}$

$\rho_{Y_1, X_2} = \frac{a_{12} \sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{1/\sqrt{2} \sqrt{7/2}}{\sqrt{5/2}} = \frac{\sqrt{7}}{\sqrt{10}}$

$\rho_{Y_2, X_1} = \frac{a_{21} \sqrt{\lambda_2}}{\sqrt{\sigma_{11}}} = \frac{1/\sqrt{2} \sqrt{3/2}}{\sqrt{5/2}} = \frac{\sqrt{3}}{\sqrt{10}}$

$\rho_{Y_2, X_2} = \frac{a_{22} \sqrt{\lambda_2}}{\sqrt{\sigma_{22}}} = \frac{-1/\sqrt{2} \sqrt{3/2}}{\sqrt{5/2}} = -\frac{\sqrt{3}}{\sqrt{10}}$

→ Compute the principal component for the point $\mathbf{x} = \begin{bmatrix} 1+\sqrt{2} \\ 2-\sqrt{2} \end{bmatrix}$

$$\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1+\sqrt{2} \\ 2-\sqrt{2} \end{bmatrix} = \begin{bmatrix} 3/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix}$$

5.17 In example 3.1, we consider the matrix $A = \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}$ and asserted, among other things, that its 2 norm is approximately 2.9208. Using the SVD, work out (on paper) the exact values of $\sigma_{\min}(A)$ and $\sigma_{\max}(A)$ for this matrix.

* Let $r = \#$ of nonzero singular value of A

then since $\det A \neq 0 \Rightarrow r = \text{rank } A = 2$

* By a property of SVD, r nonzero singular values of $A =$ square root of nonzero eigenvalues of A^*A

• Consider $A^*A = \begin{bmatrix} 1 & 0 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 8 \end{bmatrix}$

then characteristic function

$$\det \begin{bmatrix} 1-\lambda & 2 \\ 2 & 8-\lambda \end{bmatrix} = (1-\lambda)(8-\lambda) - 4 = \lambda^2 - 9\lambda + 4$$

$$\Rightarrow \text{Let } \lambda^2 - 9\lambda + 4 = 0 \text{ then } \lambda_{1,2} = \frac{9 \pm \sqrt{81-16}}{2} = \frac{9 \pm \sqrt{65}}{2} > 0$$

Then since $r = 2$, then 2 nonsingular values of A are

$$\sigma_{\max}(A) = \sqrt{\frac{9+\sqrt{65}}{2}} = 2.9208$$

$$\sigma_{\min}(A) = \sqrt{\frac{9-\sqrt{65}}{2}} = 0.6847$$





5.97 Consider the matrix $A = \begin{bmatrix} -2 & 11 \\ -10 & 5 \end{bmatrix}$

a) Determine, on paper, a real SVD of A in the form $A = U \Sigma V^T$, find the one that has the minimal number of minus signs in U and V

* We will find U, V by $\begin{cases} AA^T = U \Sigma V^T V \Sigma^T U^T = U \Sigma \Sigma^T U^T & (1) \Rightarrow \lambda_1^2 u_1 + \lambda_2^2 u_2 \\ AV = U \Sigma \Rightarrow V = A^{-1} U \Sigma & (2) \end{cases}$

* Find U :

$$AA^T = \begin{bmatrix} -2 & 11 \\ -10 & 5 \end{bmatrix} \begin{bmatrix} -2 & -10 \\ 11 & 5 \end{bmatrix} = \begin{bmatrix} 125 & 75 \\ 75 & 125 \end{bmatrix}$$

$$\det[AA^T - \lambda I] = \det \begin{bmatrix} 125 - \lambda & 75 \\ 75 & 125 - \lambda \end{bmatrix} = (125 - \lambda)^2 - 75^2$$

u_i is eigenvector associated with eigenvalue of AA^T
= (singular value of A)²

Let $\det[AA^T - \lambda I] = 0$, then $\lambda_1 = 200$ $\lambda_2 = 125 - 75 = 50$

By a theorem, the nonzero singular values of A = square root eigenvalues of AA^T

$$\sigma_1 = \sqrt{\lambda_1} = \sqrt{200} = 10\sqrt{2}$$

$$\sigma_2 = \sqrt{\lambda_2} = \sqrt{50} = 5\sqrt{2}$$

$$\Sigma = \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} = \begin{bmatrix} 14.142 & 0 \\ 0 & 7.0710 \end{bmatrix}$$

Now find the eigenvector of AA^T associated with $\lambda_1 = 200$

$$[AA^T - \lambda_1 I] = \begin{bmatrix} -75 & 75 \\ 75 & -75 \end{bmatrix} \Rightarrow \begin{bmatrix} -1 & 1 & 0 \\ 1 & -1 & 0 \end{bmatrix} \Rightarrow \vec{u}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

Now find the eigenvector of AA^T associated with $\lambda_2 = 50$

$$[AA^T - \lambda_2 I] = \begin{bmatrix} 75 & 75 \\ 75 & 75 \end{bmatrix} \Rightarrow \vec{u}_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$$

$$\Rightarrow U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

Now we find V by (2)

$$\vec{v}_1 = \sigma_1^{-1} A^{-1} \vec{u}_1 = \frac{1}{10\sqrt{2}} \frac{1}{100} \begin{bmatrix} 5 & -11 \\ 10 & -2 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \frac{10\sqrt{2}}{100} \begin{bmatrix} \frac{5-11}{\sqrt{2}} \\ \frac{10-2}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{-6}{10} \\ \frac{8}{10} \end{bmatrix} = \begin{bmatrix} -\frac{3}{5} \\ \frac{4}{5} \end{bmatrix}$$

$$\vec{v}_2 = \sigma_2^{-1} A^{-1} \vec{u}_2 = \frac{1}{5\sqrt{2}} \frac{1}{100} \begin{bmatrix} 5 & -11 \\ 10 & -2 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} = \frac{5\sqrt{2}}{100} \begin{bmatrix} \frac{5+11}{\sqrt{2}} \\ \frac{10+2}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{+16}{20} \\ \frac{12}{20} \end{bmatrix} = \begin{bmatrix} \frac{4}{5} \\ \frac{3}{5} \end{bmatrix}$$

$$V = \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}$$

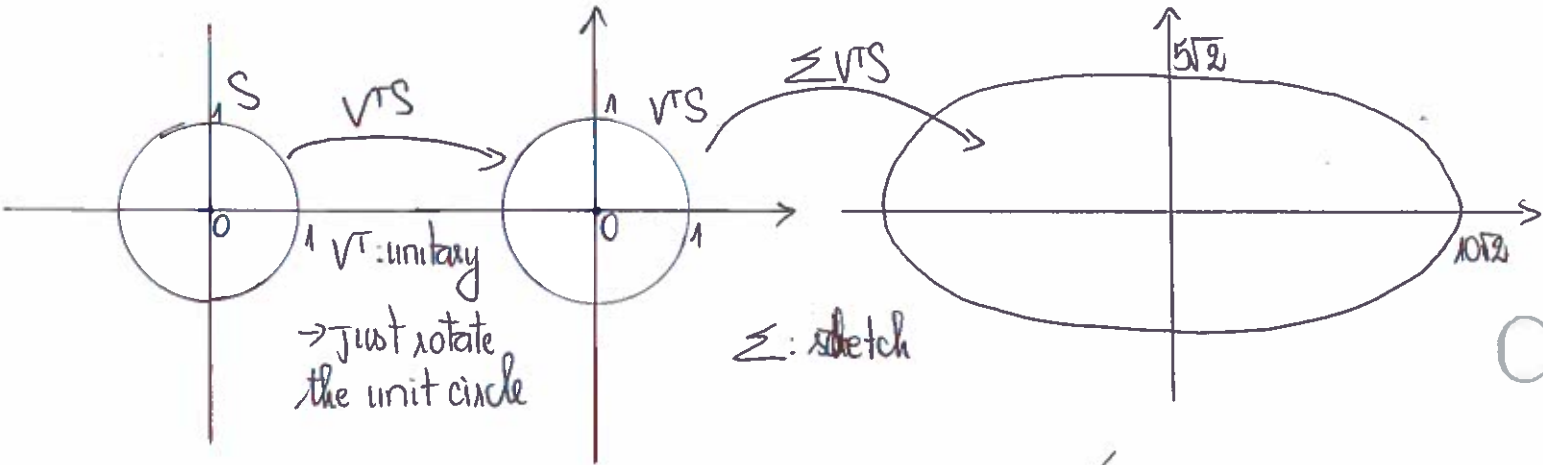
b7 List the singular values, left singular vectors and right singular vectors of A.
 Draw a careful, labeled picture of the unit ball in \mathbb{R}^2 and its image under A, together with the singular vectors, with the coordinates of their vertices marked.

From a7 singular values: $10\sqrt{2}$, $5\sqrt{2}$

left singular vectors: $\begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$ $\begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$ $A = U\Sigma V^T$

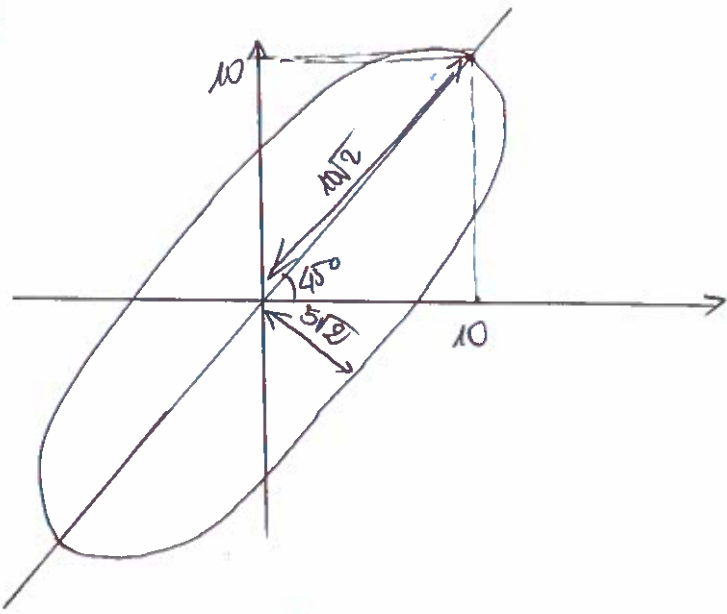
right singular vectors: $\begin{pmatrix} -3/5 \\ 4/5 \end{pmatrix}$ $\begin{pmatrix} 4/5 \\ 3/5 \end{pmatrix}$ \neq

So when we put $S =$ unit ball in \mathbb{R}^2 . $AS = U\Sigma V^T S$



$$U = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}$$

U : rotate $\frac{\pi}{4} = 45^\circ$ counter-clockwise



e) What are the 1-, 2-, ∞ -, Frobenius norms of A.

$$\|A\|_1 = \max_{1 \leq j \leq 2} \|a_j\|_1 = \max\{13, 16\} = 16 \quad \checkmark$$

$$\|A\|_2 = \|\Sigma\|_2 = \sigma_1 = 10\sqrt{2} = 14.1421 \quad \checkmark$$

$$\|A\|_\infty = \max_{1 \leq i \leq 2} \|r_i^*\|_1 = \max\{13, 15\} = 15 \quad \checkmark$$

$$\|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2} = \sqrt{200 + 50} = \sqrt{250} = 5\sqrt{10} = 15.8114 \quad \checkmark$$

d) Find A^{-1} via the SVD

$$A = U \Sigma V^T$$

$$A^{-1} = (U \Sigma V^T)^{-1} = (V^T)^{-1} \Sigma^{-1} U^{-1} = V \Sigma^{-1} U^T = \frac{1}{10} \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{10\sqrt{2}} & 0 \\ 0 & \frac{1}{5\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{20} & 0 \\ 0 & \frac{1}{10} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{20} & \frac{1}{20} \\ \frac{1}{10} & -\frac{1}{10} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{20} & -\frac{11}{100} \\ \frac{1}{10} & -\frac{1}{50} \end{bmatrix} = \begin{bmatrix} 0.05 & -0.11 \\ 0.1 & -0.02 \end{bmatrix} \quad \checkmark$$

e) Find the eigenvalues of A

$$\det(A - \lambda I) = \det \begin{pmatrix} -2-\lambda & 11 \\ -10 & 5-\lambda \end{pmatrix} = -(2+\lambda)(5-\lambda) + 110$$

$$= \lambda^2 - 3\lambda + 100$$

$$\lambda_{1,2} = \frac{3 \pm \sqrt{9-400}}{2} = \frac{3 \pm \sqrt{391}i}{2} \quad \checkmark$$

$$\lambda_1 = 1.5 + 9.8869i$$

$$\lambda_2 = 1.5 - 9.8869i$$

f) Verify that $\det A = \lambda_1 \lambda_2$ $|\det A| = \sigma_1 \sigma_2$

$$\det A = -10 + 110 = 100 \quad \checkmark$$

$$\text{from quadratic equation } \lambda^2 - 3\lambda + 100 \Rightarrow \lambda_1 \lambda_2 = 100$$

$$\sigma_1 \sigma_2 = 10\sqrt{2} \cdot 5\sqrt{2} = 100 \quad \checkmark$$

g) Area of the ellipse = $\pi \sigma_1 \sigma_2 = 100\pi$ \checkmark

* 5.47 Suppose $A \in \mathbb{C}^{m \times m}$ has an SVD $A = U \Sigma V^*$

Find an eigenvalue decomposition $H = X \Lambda X^{-1}$ of the $2m \times 2m$ Hermitian

$$\text{matrix } H = \begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix}$$

* Note that since $A = U \Sigma V^* \Rightarrow AV = U \Sigma$
 $A^* = V \Sigma^* U^* \Rightarrow A^* U = V \Sigma^*$ ✓

So we have

$$\begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix} \begin{pmatrix} -U & U \\ V & V \end{pmatrix} = \begin{pmatrix} AV & AV \\ -A^*U & A^*U \end{pmatrix} = \begin{pmatrix} U \Sigma & U \Sigma \\ -V \Sigma^* & V \Sigma^* \end{pmatrix} =$$

$$\begin{pmatrix} -U & U \\ V & V \end{pmatrix} \begin{pmatrix} -\Sigma & 0 \\ 0 & \Sigma \end{pmatrix} \checkmark$$

This is NOT H

This means we have an eigenvalue decomposition for H is $H = X \Lambda X^{-1}$

$$\text{where } X = \begin{pmatrix} -U & U \\ V & V \end{pmatrix} \quad \Lambda = \begin{pmatrix} -\Sigma & 0 \\ 0 & \Sigma \end{pmatrix} \checkmark$$

* Note that (in case $A \in \mathbb{R}^{n \times m} \Rightarrow H$ is symmetric. ?

* We can also have an eigenvalue decomposition where X is orthonormal. $H = Y \Lambda Y^{-1}, Y$

$$Y = \frac{1}{\sqrt{2}} X$$

9.5
10

5.4 Suppose $A \in \mathbb{C}^{m \times m}$ has an SVD $A = U \Sigma V^*$

Find an eigenvalue decomposition $H = X \Lambda X^{-1}$ of the $2m \times 2m$ Hermitian matrix

$$H = \begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix}$$

* We want to find eigenvalues and associated eigenvectors of H , which means we want to find λ and X such that:

$$\begin{aligned}
 \begin{matrix} H X = \lambda X \\ \begin{matrix} 2m \times 2m & 2m \times 1 & 2m \times 1 \end{matrix} \end{matrix} & \Rightarrow \begin{bmatrix} 0_{m \times m} & A^*_{m \times m} \\ A_{m \times m} & 0_{m \times m} \end{bmatrix} \begin{bmatrix} \vec{x}_{1 \times m} \\ \vec{x}_{2 \times m} \end{bmatrix} = \lambda \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \end{bmatrix} \\
 X = \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \end{pmatrix} & \Rightarrow \begin{bmatrix} A^* \vec{x}_2 \\ A \vec{x}_1 \end{bmatrix} = \begin{bmatrix} \lambda \vec{x}_1 \\ \lambda \vec{x}_2 \end{bmatrix} \quad (1) \\
 & \hspace{15em} (2)
 \end{aligned}$$

• (1) $\Rightarrow A A^* \vec{x}_2 = \lambda A \vec{x}_2 \stackrel{(2)}{=} \lambda^2 \vec{x}_2$
 $\Rightarrow \vec{x}_2$ is a left singular vector of A associated with singular value λ

• (2) $\Rightarrow A^* A \vec{x}_1 = \lambda A^* \vec{x}_1 \stackrel{(1)}{=} \lambda^2 \vec{x}_1$
 $\Rightarrow \vec{x}_1$ is a left singular vector of A associated with singular value λ

Thus we see that:



Verify that my answer for 5.1 is correct

```
A=[1 2; 0 2]
[U,S,V]=svd(A)
sigma_max=S(1,1)
sigma_min=S(2,2)
result
% A =
%      1      2
%      0      2
% U =
%      0.7497  -0.6618
%      0.6618   0.7497
% S =
%      2.9208         0
%         0      0.6847
% V =
%      0.2567  -0.9665
%      0.9665   0.2567
%
% sigma_max =
%      2.9208
% sigma_min =
%      0.6847
```

'Test digit 5.mat'
 training set azip
 test set testzip
 dtest ← correct result.

* Compare:
 mismatch_index = find(dtest == TestDigit 25)
 (2007 - 122) / 2007 = 0.9392 ← % success = 93%
 5 subspace ⇒ % success = 90%

* Consider unsuccessful cases.
 • sample Id = randi([1 2007], 100, 1) -

Verify that my answer for 5.3 a,c,d, and e are correct

```
A=[-2 11; -10 5]
[U,S,V]=svd(A)
%5.3a
% A =
%      -2      11
%     -10       5
%
% U =
%     -0.7071  -0.7071
%     -0.7071   0.7071
%
% S =
%     14.1421         0
%         0      7.0711
%
% V =
%      0.6000  -0.8000
%     -0.8000  -0.6000
```

plot(sample Id, dtest(sample Id), 0)
 hold on
 plot(sample Id, test digit 25(sample Id), 'x')
 image(testzip(6881))

```
%5.3c
norm(A, 1)
norm(A, 2)
norm(A, inf)
```



```
norm(A, 'fro')
```

```
% ans =  
%      16  
%  
% ans =  
%      14.1421  
% ans =  
%      15  
% ans =  
%      15.8114
```

```
%%5.3d-----
```

```
Ainverse=inv(A)
```

```
% Ainverse =  
%      0.0500   -0.1100  
%      0.1000   -0.0200
```

```
%%5.3e-----
```

```
[lambda]=eig(A)
```

```
lambda1=lambda(1)
```

```
lambda2=lambda(2)
```

```
% lambda =  
%      1.5000 + 9.8869i  
%      1.5000 - 9.8869i
```

```
% lambda1 =  
%      1.5000 + 9.8869i
```

```
% lambda2 =  
%      1.5000 - 9.8869i
```

MAT755_HW3

Tran Le

September 29, 2018

48/50

Johnson and Wichen 5.1

a. Evaluate T^2 , for testing $\mu^T = [1 \ 2]$, using the data

$$X = \begin{bmatrix} 2 & 12 \\ 8 & 9 \\ 6 & 9 \\ 8 & 10 \end{bmatrix}$$

```
rm(list = ls())
X=matrix(c(2,12,8,9,6,9,8,10), nrow=4, ncol=2, byrow=TRUE)
n <- dim(X)[1] # number of row of X (#observations)
p <- dim(X)[2] # number of columns (#variables)

Xbar<-apply(X,2,mean)
Xcov <- cov(X)

mu.o<-c(7, 11)
T2 <- n * mahalanois(x = Xbar, center = mu.o, cov = Xcov)
T2
## [1] 13.63636
```

b. Specific the distribution of T^2 for the situation in (a).

We have that T^2 has $((n-1) * p) / (n-p) F_{p, n-p, \alpha}$

c. Using (a) and (b), test H_0 at the $\alpha = 0.05$. What conclusion do you reach?

From (a) and (b), we want to compare the value of $T_{observed}$ and the critical value

```
alpha <- 0.05
crit.val <- p * (n-1) / (n-p) * qf(1-alpha, p, n-p)
crit.val
## [1] 57
```

So we have that $T_{observed} < \text{critical value}$, fail to reject H_0 . In conclusion, at level $\alpha = 0.05$, we fail to reject that $\mu^T = [1 \ 2]$.

2. Johnson & Wichern 5.18a

Use the college test data in Table 5.2 (see example 5.5). At the level $\alpha = 0.05$. Test the hypotheses

$$\begin{cases} H_0: \mu^T = [500, 50, 30] \\ H_1: \mu^T \neq [500, 50, 30] \end{cases}$$

```
getwd() #get working directory
## [1] "C:/Users/tran1/Desktop/MAT755/Homework/HW3"

college <- read.table("C:/Users/tran1/Desktop/MAT755/Homework/HW3/T5-2.dat")
head(college)

##   V1 V2 V3
## 1 468 41 26
## 2 428 39 26
## 3 514 53 21
## 4 547 67 33
## 5 614 61 27
## 6 501 67 29

n <- dim(college)[1] # number of observations
p <- dim(college)[2] # number of variables

alpha <- 0.05
col.xbar <- sapply(college, mean)
col.cov <- cov(college)

mu0 <- c(500, 50, 30)

T2 <- n * mahalnobis(x = col.xbar, center = mu0, cov = col.cov)
T2
## [1] 223.3102

crit.val <- p * (n-1) / (n-p) * qf(1-alpha, p, n-p)
crit.val
## [1] 8.333483
```

So we have that $T^2 > \text{crit.val}$, we reject H_0 , at level $\alpha = 0.05$, it is reasonable to believe that the group of students represented by the scores in Table 5.2 is scoring differently than the average scores for thousands of college students over the last 10 years.

I would like to mimic your code to do the sample Hotelling T2-test using HotellingT2

```
library(ICSNP)
library(ellipse)
```

```
col.T2<-HotellingsT2(X=college, mu=mu0)
col.T2

##
## Hotelling's one sample T2-test
##
## data: college
## T.2 = 72.706, df1 = 3, df2 = 84, p-value < 2.2e-16
## alternative hypothesis: true location is not equal to c(500,50,30)
```

So we reject H_0

3. Johnson & Wichern 5.19 (lumber T5-11.dat)

Measurements of x_1 =stiffness and x_2 =beding strength for a sample of $n=30$ pieces of a particular grade of lumber are given in Table 5.11. The units are pounds/(inches)². Using the data in the table,

- Construct and sketch a 95% confidence ellise for the pair $[\mu_1, \mu_2]^T$ where $\mu_1 = E(X_1)$ and $\mu_2 = E(X_2)$.

```
#import the data
lumber <- read.table(file="C:/Users/tran1/Desktop/MAT755/Homework/HW3/T5-11.dat")
names(lumber) <- c("stiff", "bend")
head(lumber)

## stiff bend
## 1 1232 4175
## 2 1115 6652
## 3 2205 7612
## 4 1897 10914
## 5 1932 10850
## 6 1612 7627

lum.xbar<-sapply(lumber, mean)
lum.xbar

## stiff bend
## 1860.500 8354.133

lum.xbar[-1]

## bend
## 8354.133

lum.cov<-cov(lumber)
lum.cov
```

```

##          stiff      bend
## stiff 124054.7 361620.4
## bend  361620.4 3486333.2

lum.cov[-1,-1]

## [1] 3486333

alpha=0.05
ci.T2 <- rbind(lum.xbar - sqrt(p*(n-1)/(n-p)/n*qf(1-alpha, p, n-
p)*diag(lum.cov)),
              lum.xbar + sqrt(p*(n-1)/(n-p)/n*qf(1-alpha, p, n-
p)*diag(lum.cov)))
row.names(ci.T2)<-c("T2_L", "T2_U")
ci.T2

##          stiff      bend
## T2_L 1751.492 7776.253
## T2_U 1969.508 8932.014

```

Now we want to sketch the 95% confidence ellipse

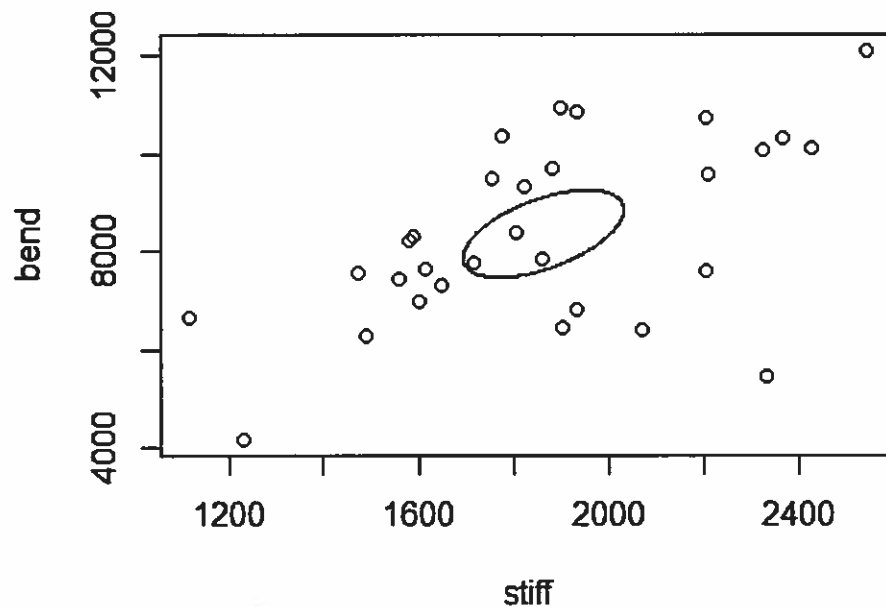
```

Ellipse <- function(covmat, centroid, csquare, resolution, plot = TRUE) {
angles <- seq(0, by = (2 * pi)/resolution, length = resolution)
sd <- covmat[1,2] / sqrt(covmat[1,1] * covmat[2,2])
projmat <- matrix(0,2,2)
projmat[1,1] <- sqrt(covmat[1,1] %**% (1+sd)/2)
projmat[1,2] <- -sqrt(covmat[1,1] %**% (1-sd)/2)
projmat[2,1] <- sqrt(covmat[2,2] %**% (1+sd)/2)
projmat[2,2] <- sqrt(covmat[2,2] %**% (1-sd)/2)
circle <- cbind(cos(angles), sin(angles))
Ellipse <- t(centroid + sqrt(csquare) * projmat %**% t(circle))
if (plot == TRUE) {lines(Ellipse)}
return(Ellipse)
}

cdellipse <- function (data, alpha=0.05, resolution=500)
{
xbar <- colMeans(data)
n <- dim(data)[1]
p <- dim(data)[2]
f <- qf(1-alpha, p, n-p)
csquare <- ((n-1)/n) * (p / (n-p)) * f
cat(csquare)
Ellipse <- Ellipse(cov(data), xbar, csquare, resolution)
}

old.par <- par(no.readonly = TRUE)
plot(lumber)
cdellipse(lumber, alpha = 0.05)

```



0.2306457

- b. Suppose $\mu_{10} = 2000$ and $\mu_{20} = 10000$ represent "typical values" for stiffness bending strength, respectively. Given the data in table 5.11 consistent with these values? Explain.

We have the given data is not consistent with these values since these values are not belong to the above 95% confidence ellipse.

c. ?

4. J& W 6.24

Four measurements were made of male Egyptian skulls for there different time periods. The measured variables are:

X_1 =maximum breath of skull (mm)

X_2 =basibregmatic heiht of skull (mm)

X_3 =basilaveplar length of skull (mm)

X_4 =nasal height of skull (mm)

- a. Construct a one way MANOVA of the Egyptian skull data.

```
egypt <- read.table(file="C:/Users/tran1/Desktop/MAT755/Homework/HW3/T6-13.dat")
```

```
names(egypt) <- c("MaxBreath", "BasHeight", "BasLength",
```

```
"NasHeight", "TimePeriod")
head(egypt)
```

```
## MaxBreath BasHeight BasLength NasHeight TimePeriod
## 1 131 138 89 49 1
## 2 125 131 92 48 1
## 3 131 132 99 50 1
## 4 119 132 96 44 1
## 5 136 143 100 54 1
## 6 138 137 89 56 1
```

```
attach(egypt)
dim(egypt)
```

```
## [1] 90 5
```

```
egypt$TimePeriod<-as.factor(egypt$TimePeriod)
egypt.manova<-
manova(cbind(MaxBreath,BasHeight,BasLength,NasHeight)~TimePeriod)
summary(egypt.manova, test="Wilks")
```

```
## Df Wilks approx F num Df den Df Pr(>F)
## TimePeriod 1 0.86076 3.4374 4 85 0.01182 *
## Residuals 88
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

interpret ...

```
by(egypt[, -5], egypt$TimePeriod, summary)
```

```
## egypt$TimePeriod: 1
## MaxBreath BasHeight BasLength NasHeight
## Min. :119.0 Min. :121.0 Min. : 89.00 Min. :44.00
## 1st Qu.:128.0 1st Qu.:131.2 1st Qu.: 95.00 1st Qu.:49.00
## Median :131.0 Median :134.0 Median :100.00 Median :50.00
## Mean :131.4 Mean :133.6 Mean : 99.17 Mean :50.53
## 3rd Qu.:134.8 3rd Qu.:136.0 3rd Qu.:102.75 3rd Qu.:53.00
## Max. :141.0 Max. :143.0 Max. :114.00 Max. :56.00
## -----
```

```
## egypt$TimePeriod: 2
## MaxBreath BasHeight BasLength NasHeight
## Min. :123.0 Min. :124.0 Min. : 90.00 Min. :45.00
## 1st Qu.:130.0 1st Qu.:129.2 1st Qu.: 97.00 1st Qu.:48.00
## Median :132.0 Median :133.0 Median : 98.50 Median :50.50
## Mean :132.4 Mean :132.7 Mean : 99.07 Mean :50.23
## 3rd Qu.:134.8 3rd Qu.:136.0 3rd Qu.:101.75 3rd Qu.:52.75
## Max. :148.0 Max. :145.0 Max. :107.00 Max. :56.00
## -----
```

```
## egypt$TimePeriod: 3
## MaxBreath BasHeight BasLength NasHeight
## Min. :126.0 Min. :123.0 Min. : 87.00 Min. :45.00
## 1st Qu.:132.2 1st Qu.:131.0 1st Qu.: 92.25 1st Qu.:48.25
```

```
## Median :136.0 Median :133.5 Median : 96.00 Median :50.00
## Mean :134.5 Mean :133.8 Mean : 96.03 Mean :50.57
## 3rd Qu.:137.0 3rd Qu.:137.0 3rd Qu.: 99.75 3rd Qu.:52.75
## Max. :140.0 Max. :145.0 Max. :106.00 Max. :60.00
```

- b. Construct 95% simultaneous confidence intervals to determine which mean components differ among the populations represented by the three time periods.

```
my.n <- nrow(egypt[, -5])
W <- (my.n - 1)*var(residuals(egypt.manova))
W

##           MaxBreath BasHeight BasLength NasHeight
## MaxBreath 1791.4500  183.5000  112.83333 293.11667
## BasHeight  183.5000 1944.3000  149.46667 178.23333
## BasLength  112.8333  149.4667 2196.02222 -10.98889
## NasHeight   293.1167  178.2333 -10.98889 842.20556

my.alpha<-0.05 # family confidence interval here
my.m <- 3 # number of groups
my.q <- 4 # number of variables

group.sample.sizes <- c(30,30,30)
for (k in 1:my.q) {
  pair.mean.diffs <- cbind( t(combn(my.m,2)),
    combn(tapply(egypt[,k],TimePeriod,mean),2,FUN=diff) )
  t.val <- qt(my.alpha/(my.q*my.m*(my.m-1)), df=my.n-my.m, lower=F)

  CI.L <- pair.mean.diffs[,3] -
    t.val*sqrt((diag(W)[k]/(my.n-my.m))*
      (1/group.sample.sizes[pair.mean.diffs[,1]] +
        1/group.sample.sizes[pair.mean.diffs[,2]])) )
  CI.U <- pair.mean.diffs[,3] +
    t.val*sqrt((diag(W)[k]/(my.n-my.m))*
      (1/group.sample.sizes[pair.mean.diffs[,1]] +
        1/group.sample.sizes[pair.mean.diffs[,2]])) )
  my.table.mat<-cbind(pair.mean.diffs,round(CI.L,3),round(CI.U,3),
    rep(k,times=nrow(pair.mean.diffs)) )

  my.table<-as.data.frame(my.table.mat)
  names(my.table)=c('grp1','grp2','diff.samp.means',
    'lower.CI','upper.CI','variable');
  print(my.table)
}

##  grp1 grp2 diff.samp.means lower.CI upper.CI variable
## 1    1    2             1.0   -2.448  4.448      1
## 2    1    3             3.1   -0.348  6.548      1
## 3    2    3             2.1   -1.348  5.548      1
##  grp1 grp2 diff.samp.means lower.CI upper.CI variable
## 1    1    2             -0.9   -4.492  2.692      2
```


##	2	1	3	0.2	-3.392	3.792	2
##	3	2	3	1.1	-2.492	4.692	2
##		grp1	grp2	diff.samp.means	lower.CI	upper.CI	variable
##	1	1	2	-0.100000	-3.918	3.718	3
##	2	1	3	-3.133333	-6.951	0.684	3
##	3	2	3	-3.033333	-6.851	0.784	3
##		grp1	grp2	diff.samp.means	lower.CI	upper.CI	variable
##	1	1	2	-0.300000	-2.664	2.064	4
##	2	1	3	0.033333	-2.331	2.398	4
##	3	2	3	0.333333	-2.031	2.698	4

c. Are the usual MANOVA assumptions realistic for these data?

I would like to compare the results we get when using MANOVA from part (a) and using 95% Bonferroni Confidence interval.

First, comparing the mean size of Maxbeareth changes throughout periods of time, we see that the mean Maxbeareth size was increasing. Then looking at the result when we use Bonferroni confidence interval, we see that the confidence interval for mean from period 1 to period 2 is $[-2.448, 4.448]$, we see that there is a chance that 0 belongs to this interval, however, there is more change that the different is a positive number, which means there was an increasing of Maxbreath mean size from period 1 to period 2.

Similarly, we see that there was a slightly change in the skull size during these periods of time. Further investments are needed to be done to have futher and more presice conclusion.

Test MANOVA assumption ...

5. Find a dataset online which contains (i) two factors and (ii) at least two continuous responses.

The data (see Micheal and Johnson problem 7.25)

Amitriptyline is prescribed by some physicians as an antidepressant. However, there are also conjectured side effects that seem to be related to the use of the drug: irregular heartbeat, abnormal blood pressures, and irregular waves on the electrocardiogram among other things. Data gathered on 17 patients who were admitted to the hospital after an amitriptyline overdose are given in Table 7.6. The two response variables are

Y1 = Total TCAD plasma level (TOT)

Y2 = Amount of amitriptyline present in TCAD plasma level(A.MI)

The five predictor variables are

Z1 = Gender: 1 iffemale, 0 if male (GEN)

Z2 = Amount of antidepressants taken at time of overdose (AMT)

Z3 = PR wave measurement (PR)

Z4 = Diastolic blood pressure (DIAP)

Z5 = QRSwavemeasurement(QRS)

(a). Run a one-way MANOVA using only one factor. State your conclusions in plain language.

```
plasma <- read.table(file="C:/Users/tran1/Desktop/MAT755/Homework/HW3/T7-6.dat")
names(plasma) <- c("TOT", "AMI", "GEN", "AMT", "PR", "DIAP", "QRS")
attach(plasma)
head(plasma)
```

```
##   TOT  AMI  GEN  AMT  PR  DIAP  QRS
## 1 3389 3149   1 7500 220    0  140
## 2 1101  653   1 1975 200    0  100
## 3 1131  810   0 3600 205   60  111
## 4  596  448   1  675 160   60  120
## 5  896  844   1  750 185   70   83
## 6 1767 1450   1 2500 180   60   80
```

```
dim(plasma)
```

```
## [1] 17  7
```

```
plasma.manova <- manova(cbind(TOT, AMI) ~ QRS)
summary(plasma.manova, test="Wilks")
```

```
##           Df  Wilks approx F num Df den Df  Pr(>F)
## QRS          1 0.62363   4.2247     2    14 0.03668 *
## Residuals 15
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(b). Run a two-way MANOVA with interactions. State your conclusions in plain language. Refit a two-way MANOVA without interactions if the interactions are not significant. State your conclusions in plain language.

```
fit.inter <- manova(cbind(TOT, AMI) ~ PR * QRS)
#summary.aov(fit.inter) # univariate ANOVA tables
summary(fit.inter, test="Wilks")
```

```
##           Df  Wilks approx F num Df den Df  Pr(>F)
## PR          1 0.36468  10.4529     2    12 0.002352 **
## QRS          1 0.66291   3.0510     2    12 0.084863 .
## PR:QRS       1 0.57127   4.5028     2    12 0.034759 *
## Residuals 13
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

With level $\alpha = 0.05$, we see that the interaction is significant. The factor QRS alone is not significant, which means QRS wave measurement does not affect Total TCAD plasma level

(TOT) and amount of amitriptyline present in TCAD plasma level(A.MI). However, the interaction of QRS wave measurement and PR wave measurement does affect Total TCAD plasma level (TOT) and amount of amitriptyline present in TCAD plasma level(A.MI).

(c). One assumption we made in MANOVA is that the variance-covariance matrices of each group of residuals are equal. Does this assumption hold in your dataset?

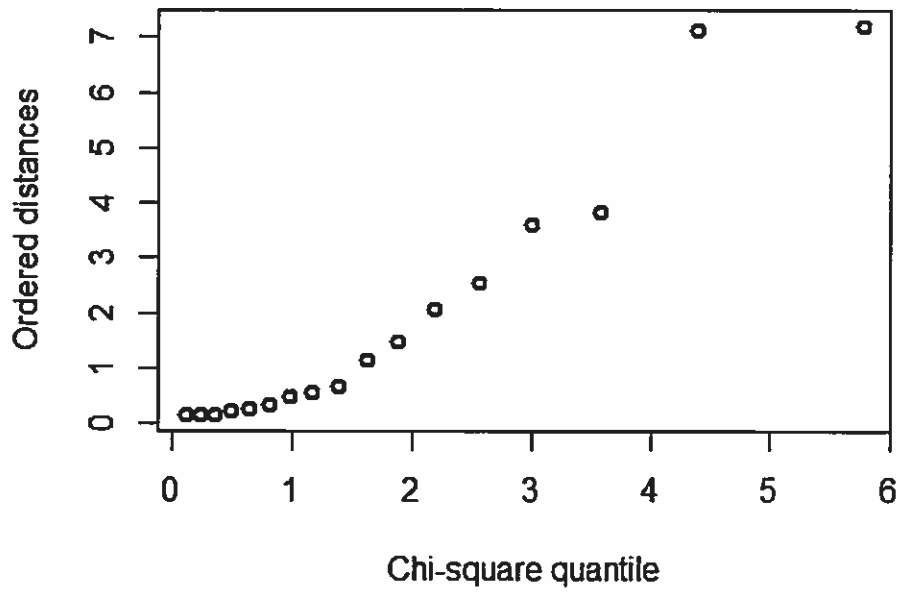
```
library(asbio)
Kullback(Y=plasma[,1:2], X=PR+QRS)

##
## Kullback test for equal covariance matrices
##      Chi* df P(Chi>Chi*)
## 214  NaN 42           NaN
```

Does box's m test work?

I could not explain why I have the NAN here, I have tried with all cases when X=PR, X=QRS, and X=PR+QRS, they all give the same result, NAN for P-value. So let's we try to check the condition whether or not the residuals follow multivariate normal probability with means equal zero.

```
chisplot <- function(x) {
  if (!is.matrix(x)) stop("x is not a matrix")
  ### determine dimensions
  n <- nrow(x)
  p <- ncol(x)
  xbar <- apply(x, 2, mean)
  S <- var(x)
  S <- solve(S)
  index <- (1:n)/(n+1)
  xcent <- t(t(x) - xbar)
  di <- apply(xcent, 1, function(x,S) x %*% S %*% x,S)
  quant <- qchisq(index,p)
  plot(quant, sort(di), ylab = "Ordered distances",
       xlab = "Chi-square quantile", lwd=2,pch=1)
}
chisplot(residuals(plasma.manova))
```



So from the QQ plot, we see that the residuals do not follow the multivariate normal probability distribution. So the results we got from Manova above may not be a good approach to work with this data set.

1

○

.

○

○

17 Consider the covariance matrix

$$\Sigma = \begin{bmatrix} 0.95 & 0.25 & 0.25 \\ 0.25 & 0.45 & 0.25 \\ 0.25 & 0.25 & 1.35 \end{bmatrix}$$

a) Find a factor model with $m=1$ factors based on Σ . Comment on the model.

b) Compute the communalities $h_i^2, i=1,2,3$, and interpret these quantities.

c) Calculate $\text{Corr}(X_i, F_j)$ for $i=1,2,3$ and $j=1$.

Which variable would carry the most weight in "naming" or "identifying" the common factor. Explain.

d) What proportion of the total variance is explained by the hidden factor.

a) We have the factor model is $X = LF + \mu + \epsilon$ and the covariance structure is

$$\Sigma = LL' + \Psi = \begin{bmatrix} l_{11} \\ l_{21} \\ l_{31} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} & l_{31} \end{bmatrix} + \begin{bmatrix} \psi_1 & & \\ & \psi_2 & \\ & & \psi_3 \end{bmatrix}$$

$$\Rightarrow 0.95 \stackrel{\textcircled{1}}{=} l_{11}^2 + \psi_1$$

$$0.25 \stackrel{\textcircled{2}}{=} l_{11} l_{21}$$

$$0.25 \stackrel{\textcircled{4}}{=} l_{11} l_{31}$$

$$0.45 \stackrel{\textcircled{3}}{=} l_{21}^2 + \psi_2$$

$$0.25 \stackrel{\textcircled{5}}{=} l_{21} l_{31}$$

$$1.35 \stackrel{\textcircled{6}}{=} l_{31}^2 + \psi_3$$

or

$$\bullet \textcircled{4} \text{ and } \textcircled{5} \Rightarrow l_{11} = l_{21}$$

$$l_{11} = l_{21} \text{ and } \textcircled{2} \Rightarrow l_{11} = l_{21} = \sqrt{0.25} = 0.5$$

$$l_{11} = l_{21} = -0.5$$

$$l_{11} = 0.5 \text{ and } \textcircled{4} \Rightarrow l_{31} = 0.5$$

$$l_{31} = -0.5$$

$$l_{11} = 0.5 \text{ and } \textcircled{1} \Rightarrow \psi_1 = 0.7$$

$$\psi_1 = 0.7$$

$$l_{21} = 0.5 \text{ and } \textcircled{3} \Rightarrow \psi_2 = 0.9$$

$$\psi_2 = 0.9$$

$$l_{31} = 0.5 \text{ and } \Rightarrow \psi_3 = 1.1$$

$$\psi_3 = 1.1$$

Since the variances are the same, we can use both of these models (this agrees with what we have known that the factor models are not unique)

In this case, variances are the same \Rightarrow can we use one of the two

$$X = \begin{bmatrix} 0.5 \\ 0.5 \\ 0.5 \end{bmatrix} F + \mu + \epsilon$$

*7 Compute the communalities $h_i^2, i=1,2,3$ and interpret these quantities.

$$h_1^2 = \lambda_{1,1}^2 = 0.25$$

$$h_2^2 = \lambda_{2,1}^2 = 0.25$$

$$h_3^2 = \lambda_{3,1}^2 = 0.25$$

We have that the communalities h_i represents the proportion of variance of X_i is explained by $m=1$ factor, which means that the higher the communality, the more informative X_i is

since $h_1^2 = h_2^2 = h_3^2 \Rightarrow$ the three variables X_1, X_2, X_3 share the same amount information with the rest of the variables.

\Rightarrow Cal $\text{cor}(X_i, F_j) \quad i=1,2,3, \quad j=1$. Which variables would carry the most weight?

$$\text{we have } \text{cov}(X_i, F) = L \Rightarrow \text{cov}(X_1, F) = 0.5 \quad \text{cor}(X_1, F) = \frac{\sigma_{X_1, F}}{\sigma_{X_1} \sigma_F} = \frac{0.5}{\sqrt{0.95} \cdot 1}$$

$$\text{cov}(X_2, F) = 0.5 \quad \text{cor}(X_2, F) = \frac{\sigma_{X_2, F}}{\sigma_{X_2} \sigma_F} = \frac{0.5}{\sqrt{0.45} \cdot 1}$$

$$\text{cov}(X_3, F) = 0.5 \quad \text{cor}(X_3, F) = \frac{\sigma_{X_3, F}}{\sigma_{X_3} \sigma_F} = \frac{0.5}{\sqrt{1.35} \cdot 1}$$

So we have that $\text{cor}(X_2, F) > \text{cor}(X_1, F) > \text{cor}(X_3, F)$

\Rightarrow The variable X_2 carries the most weight since it has the biggest correlation with F comparing with other variable

\Rightarrow What proportion of the total variance is explained by the hidden factor.

* We first need to find the eigenvalues $\lambda_i, i=1,3$ of Σ .

$$\det(\Sigma - \lambda I) = \det \begin{pmatrix} 0.95 - \lambda & 0.25 & 0.25 \\ 0.25 & 0.45 - \lambda & 0.25 \\ 0.25 & 0.25 & 1.35 - \lambda \end{pmatrix} = -\lambda^3 + 2.75\lambda^2 - 2.13\lambda + 0.4365$$

$$\Rightarrow \begin{cases} \lambda_1 = 1.5711 \\ \lambda_2 = 0.85331 \\ \lambda_3 = 0.32559 \end{cases}$$

\Rightarrow The total variance explained by the

by 1st factor is $\frac{\lambda_1}{\sum_{i=1}^3 \lambda_i} = \frac{1.5711}{0.95 + 0.45 + 1.35} = \frac{1.5711}{2.75} = 0.5713$

by 2nd factor is $\frac{0.85331}{2.75} = 0.3103$

by 3rd factor is $\frac{0.32559}{2.75} = 0.1184$

Wickham - Johnson - 9.6

Verify the following matrix identities.

a) $(I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} L = I - (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} L$

b) $(LL' + \Psi)^{-1} = \Psi^{-1} - \Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1}$

c) $L' (LL' + \Psi)^{-1} L = (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} L$

a) Put $B = L' \Psi^{-1} L$,

Problem a) \Leftrightarrow we need to prove that $(I + B)^{-1} B = I - (I + B)^{-1} B$ (*)

We multiply both sides of (a) by $(I + B)$, we have

$$(I + B)(I + B)^{-1} B = (I + B) [I - (I + B)^{-1} B]$$

$$\Leftrightarrow IB = (I + B)I - (I + B)(I + B)^{-1} B = (I + B) - I = B \quad (\text{correct equality})$$

so we have that (*) is correct \square a)

b) Prove that $(LL' + \Psi)^{-1} = \Psi^{-1} - \Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1}$

We multiply both sides of (b) by $(LL' + \Psi)$ we have:

• LHS $(LL' + \Psi) = (LL' + \Psi)^{-1} (LL' + \Psi) = I$

• RHS $(LL' + \Psi) = (\Psi^{-1} - \Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1}) (LL' + \Psi)$

$$= \Psi^{-1} LL' + \Psi^{-1} \Psi - \underbrace{\Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} LL'}_{\text{from (a)}} - \underbrace{\Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} \Psi}_{= I}$$

$$= \Psi^{-1} LL' + I - \underbrace{\Psi^{-1} L (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} LL'}_{\text{from (a)}} - I$$

$$= \Psi^{-1} LL' + I - \Psi^{-1} L [I - (I + L' \Psi^{-1} L)^{-1} L' \Psi^{-1} L] L'$$

$$= \Psi^{-1} LL' + I - \Psi^{-1} L I L' = I$$

So we have $(LHS)(LL' + \Psi) = RHS(LL' + \Psi)$

since $LL' + \Psi$ is invertible $\Rightarrow LHS = RHS \quad \square$ b

$$\Rightarrow \text{Prove that } L'(LL' + \Psi)^{-1} = (I + L'\Psi^{-1}L)^{-1}L'\Psi^{-1}$$

We first notice that Ψ , LL' and $I + L'\Psi^{-1}L$ are symmetric matrices.

Similarly we have that $(LL' + \Psi)^{-1} = \Psi^{-1} - \Psi^{-1}L(I + L'\Psi^{-1}L)^{-1}L'\Psi^{-1}$

$$\begin{aligned} \text{then we have } (LL' + \Psi)^{-1}L &= [\Psi^{-1} - \Psi^{-1}L(I + L'\Psi^{-1}L)^{-1}L'\Psi^{-1}]L \\ &= \Psi^{-1}L - \underbrace{\Psi^{-1}L(I + L'\Psi^{-1}L)^{-1}L'\Psi^{-1}L}_{\text{by (a)}} \end{aligned}$$

$$= \Psi^{-1}L(I + L'\Psi^{-1}L)^{-1}$$

take transpose both sides, note that Ψ , $L'L$ and $(I + L'\Psi^{-1}L)$ are symmetric, we have

$$L'(LL' + \Psi)^{-1} = ((I + L'\Psi^{-1}L)^{-1})'L'(\Psi^{-1})'$$

$$\Rightarrow L'(LL' + \Psi)^{-1} = (I + L'\Psi^{-1}L)^{-1}L'\Psi^{-1} \quad \square$$

Mat 755 Final exam

P1 - Classification rules are often evaluated in terms of the expected cost of misclassification $ECM = c(2|1) P(2|1) p_1 + c(1|2) P(1|2) p_2$

Q7 Show that the region R_1 and R_2 that minimize the ECM are

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

$$R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

* We have $P(2|1) = P(R_2 | \pi_1) =$ probability of misclassifying an object into R_2 while it belongs to π_1

$$= \int_{R_2} f_1(x) dx$$

Similarly $P(1|2) = \int_{R_1} f_2(x) dx$

* So we have

$$\begin{aligned} ECM &= c(2|1) P(2|1) p_1 + c(1|2) P(1|2) p_2 \\ &= c(2|1) \int_{R_2} f_1(x) dx p_1 + c(1|2) \int_{R_1} f_2(x) dx p_2 \end{aligned}$$

We note that since $\left\{ \begin{array}{l} \Omega = R_1 + R_2 \\ R_1 \cap R_2 = \emptyset \end{array} \right\} \Rightarrow 1 = \int_{R_1} f_1(x) dx + \int_{R_2} f_1(x) dx$
 $\Rightarrow \int_{R_2} f_1(x) dx = 1 - \int_{R_1} f_1(x) dx$

So then

$$ECM = c(2|1) \left[1 - \int_{R_1} f_1(x) dx \right] p_1 + c(1|2) \left[\int_{R_1} f_2(x) dx \right] p_2$$

$$= \underbrace{\int_{R_1} [c(1|2) f_2(x) p_2 - c(2|1) f_1(x) p_1] dx}_{(*)} + \underbrace{c(2|1) p_1}_{\text{constant}}$$

$\Rightarrow ECM$ is minimized when R_1 is chosen to be the region where $(*) \leq 0$

→ ECM is minimized when R_1 is chosen to be the region where

$$c(1|2) p_2 f_2(x) \leq c(2|1) p_1 f_1(x)$$

$$\Rightarrow R_1: \frac{f_1}{f_2} \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

And so since $\Omega = R_1 \cup R_2$, $R_1 \cap R_2 = \phi$

$$R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

b) A researcher has enough data to estimate the density function $\begin{cases} f_1(x) \\ f_2(x) \end{cases}$ associated with population π_1 and π_2 , respectively.

$$\text{Let } c(2|1) = 50, c(1|2) = 100$$

It is known that about 30% of all possible items belong to π_2 .

Give the minimum ECM rule for assigning a new item to one of the two populations

$$* c(2|1) = 50 \quad c(1|2) = 100$$

$$p_1 = 0.7 \quad p_2 = 0.3$$

$$\text{Then } \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) = \left(\frac{100}{50} \right) \left(\frac{0.3}{0.7} \right) = 2 * \frac{3}{7} = \frac{6}{7}$$

* So for each observation x , we need to compare $\frac{f_1(x)}{f_2(x)}$ and $\frac{6}{7}$.

assign the observation to group 1, R_1 , when $\frac{f_1(x)}{f_2(x)} \geq \frac{6}{7}$

assign the observation to group 2, R_2 , when $\frac{f_1(x)}{f_2(x)} < \frac{6}{7}$.

c) Measurements recorded on a new item yield the density values $f_1(x) = 0.3$ and $f_2(x) = 0.5$.

Given the preceding information, assign this item to π_1 or π_2 :

$$\frac{f_1(x)}{f_2(x)} = \frac{0.3}{0.5} = \frac{3}{5} < \frac{6}{7} = \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

so we assign the new item to π_2 , group 2.

* Problem 2:

Let $X = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n]$ be centered data

$$d_{ij}^2 = (\underline{x}_i - \underline{x}_j)^T (\underline{x}_i - \underline{x}_j)$$

$$B = X^T X$$

or Show that B has entries

$$b_{ij} = -\frac{1}{2} \left[d_{ij}^2 - \frac{2}{n} \sum_{i=1}^n d_{ij}^2 + \frac{1}{n^2} \sum_{i,j=1}^n d_{ij}^2 \right]$$

Note that $d_{ij}^2 = \underline{x}_i^T \underline{x}_i + \underline{x}_j^T \underline{x}_j - 2 \underline{x}_i^T \underline{x}_j \Rightarrow \underline{x}_i^T \underline{x}_j = -\frac{1}{2} (d_{ij}^2 - \underline{x}_i^T \underline{x}_i - \underline{x}_j^T \underline{x}_j)$ (1)

We also have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n d_{ij}^2 &= \frac{1}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i + \frac{1}{n} \sum_{i=1}^n \underline{x}_j^T \underline{x}_j - \frac{2}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_j \\ &= \underline{x}_i^T \underline{x}_i + \frac{1}{n} \sum_{i=1}^n \underline{x}_j^T \underline{x}_j = -\frac{2}{n} \sum_{i=1}^n \underline{x}_j^T \underline{x}_i = 0. \end{aligned} \quad (2)$$

$$\begin{aligned} \frac{1}{n^2} \sum_{i,j=1}^n d_{ij}^2 &= \frac{1}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i + \frac{1}{n} \underbrace{\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \underline{x}_j^T \underline{x}_j}_{= \frac{1}{n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \underline{x}_i^T \underline{x}_i} = \frac{2}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i \quad (3) \\ &= \frac{1}{n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \underline{x}_i^T \underline{x}_j \end{aligned}$$

* So now consider $B = X^T X$, then we have B will have entries:

$$b_{ij} = \underline{x}_i^T \underline{x}_j \stackrel{\text{by (1)}}{=} -\frac{1}{2} (d_{ij}^2 - \underline{x}_i^T \underline{x}_i - \underline{x}_j^T \underline{x}_j)$$

$$\stackrel{\text{by (2)}}{=} -\frac{1}{2} \left(d_{ij}^2 - \frac{1}{n} \sum_{i=1}^n d_{ij}^2 + \frac{1}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i - \frac{1}{n} \sum_{i=1}^n d_{ij}^2 + \frac{1}{n} \sum_{j=1}^n \underline{x}_j^T \underline{x}_j \right)$$

$$= -\frac{1}{2} \left(d_{ij}^2 - \frac{2}{n} \sum_{i=1}^n d_{ij}^2 + \frac{2}{n} \sum_{i=1}^n \underline{x}_i^T \underline{x}_i \right)$$

$$\stackrel{\text{by (3)}}{=} -\frac{1}{2} \left(d_{ij}^2 - \frac{2}{n} \sum_{i=1}^n d_{ij}^2 + \frac{2}{n^2} \sum_{i,j=1}^n d_{ij}^2 \right)$$

$$= -\frac{1}{2} \left(d_{ij}^2 - \frac{2}{n} \sum_{i=1}^n d_{ij}^2 + \frac{2}{n^2} \sum_{i,j=1}^n d_{ij}^2 \right) \quad \square \text{ or } \square$$

b7 If $B = V\Lambda V^T$ is the spectral decomposition of B .

Show that $X = \Lambda^{1/2} V^T$

* Note that since B is symmetric and since $B = X^T X \Rightarrow$ positive definite

$\Rightarrow B = V\Lambda V^T$ is well defined, where $\Lambda = \text{diag}(\lambda_i) = \begin{pmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{pmatrix}$

and $B = V\Lambda V^T$

$$= V\Lambda^{1/2} \Lambda^{1/2} V^T \quad \text{where } \Lambda^{1/2} = \text{diag}(\sqrt{\lambda_i})$$

$$= V(\Lambda^{1/2})^T \Lambda^{1/2} V^T \quad (\text{since } \Lambda^{1/2} \text{ is a diagonal matrix})$$

$$B = (\Lambda^{1/2} V^T)^T (\Lambda^{1/2} V^T)$$

Since $B = X^T X$

$$\} \Rightarrow X = \Lambda^{1/2} V^T \quad \square b$$