

Introduction

→ sample mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

→ variance

$$S_{11} = \frac{1}{n-1} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$$

→ pairwise covariance

$$S_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

→ sample correlation coefficient

$$r_{jk} = \frac{S_{jk}}{\sqrt{S_{jj}} \sqrt{S_{kk}}}$$

Eigenvalues - Eigenvectors

$$q(\lambda) = |\Sigma - \lambda I| = 0 \rightarrow \lambda$$

singular $\rightarrow a_i \pi_i = \sum_j \pi_j$ right eigenvector of Σ to a_i

standardized if $\pi_i^T \pi_i = 1$

$$|\Sigma - aI| = |C| |\Sigma - aC^{-1}C| |C^{-1}| = |C \Sigma C^{-1} - aI|$$

$$\pi_i \Sigma = \pi_i a \rightarrow \pi_i C^{-1} \Sigma C = \pi_i a \rightarrow C \pi_i$$

Spectral Theorem

→ Every symmetric $p \times p$ matrix Σ

$$\Sigma = \mathcal{N} \mathcal{A} \mathcal{N}^T = \sum_{i=1}^p \pi_i a_i \pi_i^T$$

$$\mathcal{N}^T = \mathcal{N}^{-1}$$

$$\mathcal{N} = (\pi_1 \dots \pi_p)$$

$$\mathcal{A} = \text{diag}(a_1, \dots, a_n)$$

Multivariate Case

$$E(x) = [E(x_i)]$$

$$E(Ax+B) = AE(x) + B$$

$$\text{Cov}(x, y) = E((x_i - E(x_i))(y_i - E(y_i)))$$

$$\text{Cov}(x, y) = [E((x - \mu_x)(y - \mu_y)^T)] = E(xy^T) - \mu_x \mu_y^T$$

Multivariate Normal Distribution

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

pD^2 without $\frac{1}{2}!$

$$x - \mu \sim N_p(0, \Sigma)$$

$$Ax \sim N_p(A\mu, A\Sigma A^T)$$

$$pD^2 \stackrel{!}{=} c^2 \rightarrow \text{ellipse}$$

$$pD^2 \sim \chi_p^2$$

Wishart Distribution

→ symmetric $p \times p$ matrix ω

→ if ω can be represented as:

$$\omega = x^T x \quad | \quad x \sim N_p(0, \Sigma)$$

we can use the notation:

$$\omega \sim W_p(\Sigma, n)$$

→ Properties:

$$\rightarrow A^T \omega A \sim W_q(A^T \Sigma A, n)$$

$$\rightarrow \Sigma^{-\frac{p}{2}} \omega \Sigma^{-\frac{1}{2}} \sim W(1, n)$$

$$\rightarrow \frac{a^T \omega a}{a^T \Sigma a} \sim \chi_n^2 \quad (a \neq 0)$$

Central Limit Theorem

→ n independent obs \Rightarrow for large $n-p$ approximately:

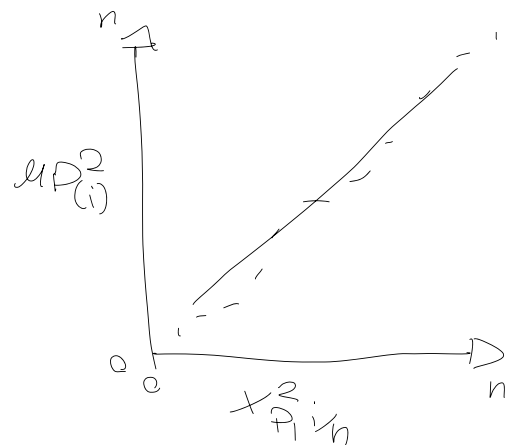
$$\sqrt{n}(\bar{x} - \mu) \sim \mathcal{N}(0, \Sigma)$$

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

\Rightarrow approximate normal distribution

Tests for multivariate normality - χ^2 plot

- (1) for every obs $\Rightarrow \text{LD}_i^2(x_i | \bar{x}, S)$
- (2) sort LD's: $\text{LD}_{(1)}^2 \leq \dots \leq \text{LD}_{(n)}^2$
- (3) compute quantiles i/n of χ_p^2 : $\chi_{p, i/n}^2$
- (4) Graph: if linear trend \Rightarrow multivar. normal



Transformation to Normality

→ "change of data scale"

→ examples:

Count $y \rightarrow \sqrt{y}$ Ratios $\rightarrow \log$ Correlations \rightarrow Fisher

→ Power Transformation x^λ

→ for appropriate λ : histogram or QQ-Plot

→ Adaptation: Box-Cox

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln x & \lambda = 0 \end{cases} \quad \left| \lambda \in \mathbb{R}, x > 0 \right.$$

→ only good approx. and not normality can be achieved

Tests, confidence regions

$$H_0: \mu = \mu_0, H_1: \mu \neq \mu_0, x_1, \dots, x_n \sim N_p(\mu, \Sigma)$$

→ Hotelling's T^2

$$T^2 = \underbrace{(\bar{x} - \mu_0)^T \left(\frac{1}{n} S\right)^{-1} (\bar{x} - \mu_0)}_{\text{MD}^2(\bar{x}, \mu_0, \frac{1}{n} S)} \sim \frac{(n-1)p}{(n-p)} F_{p, n-p}$$

→ significance level α

$$\mathcal{L} = P \left[T^2 > \frac{(n-1)p}{(n-p)} F_{p, n-p, 1-\alpha} \right]$$

→ confidence region

→ ellipsoid determined by all μ for which:

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

Clustering Analysis

Distance measures

$$L_1: d(i, j) = \|x_i - x_j\|_1 = \sum_{k=1}^p |x_{ik} - x_{jk}|$$

$$L_2: d(i, j) = \|x_i - x_j\|_2 = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2}$$

$$D = [(d_{ij})]$$

$$[i \setminus j]$$

Hierarchical clustering measures (5)

$$\max_{i \in C_k, j \in C_l} d(i, j)$$

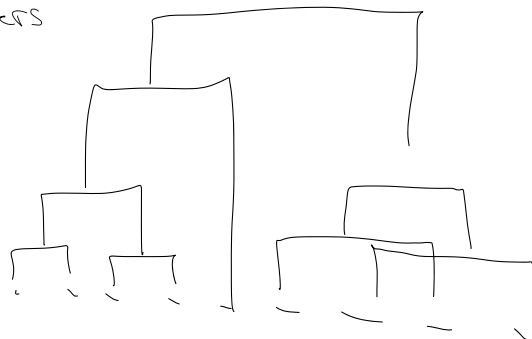
$$\min_{i \in C_k, j \in C_l} d(i, j)$$

$$\frac{1}{n_k n_l} \sum_{i \in C_k, j \in C_l} d(i, j)$$

$$\|\bar{x}(C_k) - \bar{x}(C_l)\|_2$$

$$\frac{\|\bar{x}(C_k) - \bar{x}(C_l)\|^2}{\frac{1}{n_k} + \frac{1}{n_l}}$$

... increase in variance when merging two clusters



Partitioning measures

K-Means

$$T = W(C) + B(C)$$

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i, j \in C_k} d(i, j)^2$$

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i \in C_k, j \notin C_k} d(i, j)^2$$

→ Cluster centers (μ) to minimize

$$W(C) = \sum_{k=1}^K n_k \sum_{i \in C_k} \|x_i - \bar{x}_k\|_2^2 \quad \leftarrow \text{K-means}$$

Model-based clustering

π_1, \dots, π_k ... mixing coefficients ($\sum_{i=1}^k \pi_i = 1$)

→ estimate μ, Σ, π using $E \& M$

$\Sigma: p \times p \rightarrow$ large samples difficult $\Rightarrow \sum_k = \frac{2}{k}$

$$k_1, \dots, k_n = a \mid a \in \mathbb{R}$$

$$k_1, \dots, k_n \neq a \mid a \in \mathbb{R}$$

Fuzzy clustering - Fuzzy k-Means

$$\min \rightarrow \sum_{i=1}^n \sum_{k=1}^k u_{ik}^2 \|x_i - m_k\|^2$$

$$\text{with: } m_{kj} = \frac{\sum_{i=1}^n u_{ik}^2 x_{ij}}{\sum_{i=1}^n u_{ik}^2}$$

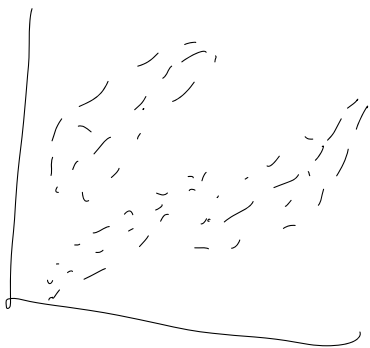
→ k given

→ u_{ik} to minimize

→ $\hat{=}$ same also as k-means

u_{ik} ... membership coeff.

m_{kj} ... weighted cluster center



→ D instead of x can be used

→ clusters spherically shaped

$$B_k = \sum_{k=1}^K \|\bar{x}_k - \bar{x}\|^2 \quad \text{Heterogeneity}$$

$$W_k = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \bar{x}_k\|^2 \quad \text{Homogeneity}$$

Validity measures (4) \Rightarrow ~~determine optimal k~~

$$H_k = \ln \frac{B_k}{W_k}$$

$$CH_k = \frac{B_k}{\frac{W_k}{n-k}}$$

\hookrightarrow identifies point where higher k does not result in better clustering "elbow point"

\hookrightarrow higher values $\hat{=}$ better clustering

avg. silh width

$$s = \frac{1}{n} \sum_{i=1}^n s_i$$

\rightarrow higher $s \hat{=}$ better classification

$$s_i = \frac{d_{i,c} - d_{i,k}}{\max(d_{i,c}, d_{i,k})}$$

$\in [-1, 1]$... -1 false clustering, 1 perfect clusters

$$d_{i,c} = \min(d_{i,c_j}) \rightarrow d_{i,c_j} = \frac{1}{n_j} \sum_{j \in C_j} d(i,j)^2 \dots \text{avg. dissimilarity to other clusters}$$

$$d_{i,k} = \frac{1}{n_k - 1} \sum_{i,j \in C_k} d(i,j)^2 \dots \text{avg. dissimilarity to same cluster}$$

\rightarrow smallest k that maximizes s distance to min. expectation

$\hat{=}$ exp under sample n of appropri. reference dists.

$$\text{Gap}_n(k) = E_n \{ \log(\tilde{W}_k) \} - \log(\bar{W}_k)$$

$$\tilde{W}_k = \sum_{k=1}^K \frac{1}{2n_k} \sum_{i,j \in C_k} d(i,j)^2$$

\Rightarrow smallest k, so that $\text{Gap}_n(k) \leq 1$ std. error away from 1. local max

$$\text{Gap}_n(k) = E_n \{ \log(\hat{W}_k) \} - \log(\tilde{W}_k)$$

Gap indicates how good cluster is compared to random

$$y = X\beta + \Sigma$$

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \Sigma$$

↓

$$y_1 = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p} + \Sigma_1$$

⋮

$$y_n = \beta_0 + \beta_1 x_{n1} + \dots + \beta_p x_{np} + \Sigma_n$$

$$\rightarrow y = X\beta + \Sigma$$

$$y \dots n \times 1$$

$$x \dots n \times (p+1)$$

$$\beta = (p+1) \times 1$$

$$\Sigma \dots n \times 1$$

$$E(\Sigma) = 0 \quad \text{cov}(\Sigma) = \sigma^2 I_n$$

$$LS = (y - X\beta)^T (y - X\beta) \rightarrow \min \Rightarrow$$

$$\hat{\beta} = (X^T X)^{-1} X^T y \rightarrow \hat{y} = X\hat{\beta}$$

$$X : p+1 \leq n \dots \text{not full rank} \Rightarrow H = X(X^T X)^{-1} X^T$$

$$\hat{y} = Hy \rightarrow y - \hat{y} = \tilde{\Sigma} = (I - H)y$$

$$E(\tilde{\Sigma}) = 0 \quad \text{cov}(\tilde{\Sigma}) = \sigma^2 (I - H)$$

$$E(\hat{\beta}) = \beta \quad \text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$$s^2 = \frac{\sum \tilde{\Sigma}^T \tilde{\Sigma}}{n - p - 1}$$

$$\left. \begin{aligned}
 y_{1j} &= \beta_{0j} + \beta_{1j}x_{11} + \dots + \beta_{pj}x_{p1} + \varepsilon_{1j} \\
 \vdots \\
 y_{nj} &= \beta_{0j} + \beta_{1j}x_{1n} + \dots + \beta_{pj}x_{pn} + \varepsilon_{nj}
 \end{aligned} \right\} Y = XB + \varepsilon$$

$Y \dots n \times m$
 $X \dots n \times (p+1)$
 $B \dots (p+1) \times m$
 $\varepsilon \dots n \times m$

$$LS = (Y - XB)^T (Y - XB) \rightarrow \hat{B} = (X^T X)^{-1} X^T Y$$

$$\hat{Y} = X \hat{B}$$

$$\text{Cov}(\varepsilon_j, \varepsilon_k) = \sigma^2 \delta_{jk} \quad E(\varepsilon_j) = 0$$

$$Y \sim N_n(XB, \sigma^2 I) \rightarrow \alpha \text{ for } \beta_j \text{ and } \sigma^2$$

$$\lambda = \left(\frac{|SR|}{|S_{R_0}|} \right)^{\frac{n}{2}} \quad \text{Wilks Lambda} \quad SR = \frac{RRT}{n-p-1} \quad R = Y - \hat{Y}$$

$$SR_0 = \frac{R_0 R_0^T}{n-q-1} \quad R_0 = Y^* - \hat{Y}_0$$

$$q < p$$

→ approximates for large n 's to $-2 \log \lambda \sim \chi^2_{m(p-q)}$

$$\rightarrow (F(x; T, G)) = \lim_{\varepsilon \rightarrow 0} \frac{T((1-\varepsilon)G + \varepsilon x) - T(G)}{\varepsilon} \quad \leftarrow \text{rate of change}$$

$$\rightarrow \text{max bias}(m, T, x) = \sup_{\tilde{x}} \|T(\tilde{x}) - T(x)\|$$

$$\rightarrow \Sigma_n^2(T, x) = \min \left\{ \frac{m}{n}; \text{max bias}(m, T, x) = \infty \right\}$$

$$\rightarrow \text{efficiency} \approx \frac{1}{\text{var}} \quad (+ \text{Fisher Information})$$

\rightarrow μ -estimator of $n \rightarrow x, \dots, n+(p+1)$

$$\hat{\beta} = \text{argmin}_{\beta} \sum_{i=1}^n r_i(\beta)^2 \quad \rightarrow \quad \hat{\beta} = \text{argmin}_{\beta} \sum_{i=1}^n \rho(r_i(\beta))$$

$$\rightarrow \rho(r) = r^2 \dots \text{LS}$$

$$\rho(r) = |r| \dots \text{L1}$$

\rightarrow scale invariance:

$$\hat{\beta} = \text{argmin}_{\beta} \sum_{i=1}^n \rho\left(\frac{r_i(\beta)}{\hat{\sigma}}\right)$$

$\hat{\sigma} \dots$ robust scale estimator of residuals

\rightarrow μ -estimating equations

\rightarrow enhances robustness by getting rid of dependence on data scale

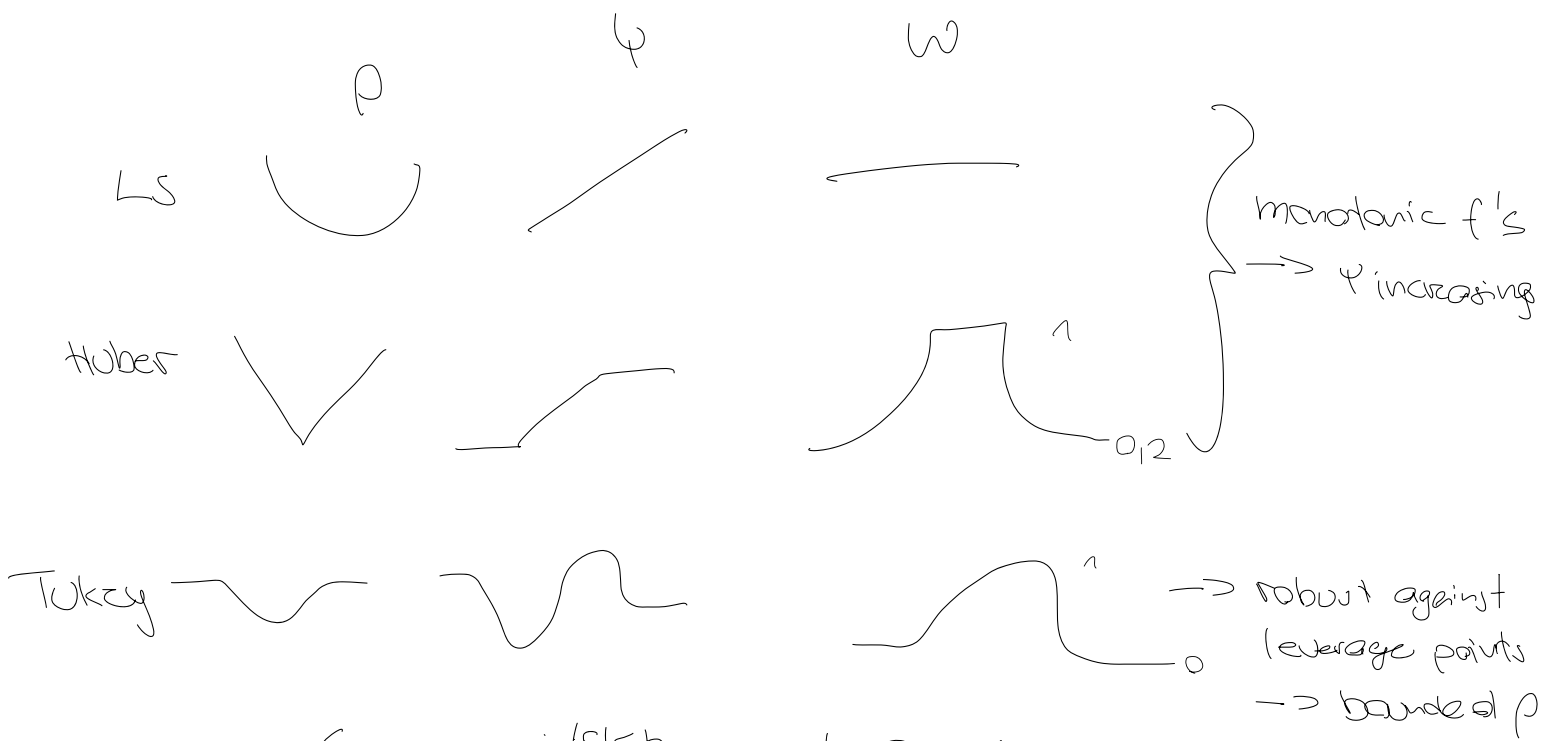
$$\sum_{i=1}^n \psi\left(\frac{r_i(\beta)}{\hat{\sigma}}\right) x_i = 0 \quad \rightarrow \quad \psi(r) \dots \text{naive equations}$$

$\psi = \rho'$

$$\omega(r) = \psi(r)/r \quad \rightarrow \quad \sum_{i=1}^n r \omega\left(\frac{r_i(\beta)}{\hat{\sigma}}\right) x_i$$

$$\omega_i = \omega(r_i(\beta)/\hat{\sigma}) \quad \rightarrow \quad \sum_{i=1}^n \omega_i (y_i - x_i^T \beta) x_i = 0$$

$r_i = r_i(\beta_m)$



$$\rho_b(r) = \begin{cases} r^2 & : |r| \leq b \\ b \operatorname{sign}(r) & : |r| > b \end{cases}$$

$b=0 \rightarrow L_1$
 $b=\infty \rightarrow LS$

$$\rho_{BK}(r) = \begin{cases} \frac{1}{k} \left(\frac{r}{k} \right)^2 \left(3 - 3 \left(\frac{r}{k} \right)^2 + \left(\frac{r}{k} \right)^4 \right) & : |r| \leq k \\ 1 & : |r| > k \end{cases}$$

$k \dots$ tuning parameter
 $k \rightarrow \infty \dots LS$
 \rightarrow efficiency vs. robust

Computation

- IRWLS:**
- (1) $r_i = r_i(\hat{\beta}_m)$ in iter = 0 : $\hat{\beta}_m = \hat{\beta}_0 \rightarrow m \dots$ approx. at iteration m
 - (2) $w_i = W(r_i/\hat{\sigma}) \dots$ update weights
 - (3) $\hat{\beta}_{m+1} \dots$ update β , reiterate $\rightarrow \hat{\beta}_0$ must be robust
- \rightarrow no leverage points: L_1 as initial $\hat{\beta}$ and $\hat{\sigma}$ computed as robust scale of r_i 's (e.g. MAD)
- \rightarrow leverage points: $\hat{\beta}_0$ cannot use previous res. scale \Rightarrow

Scale Estimators

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \hat{\sigma}(\tau(\beta))$$

Based on $\hat{\sigma}$ we differentiate:

→ Estimators using Scales based on ordered values

→ LMS: $\hat{\sigma}(\tau) = \operatorname{med}(|r_i|)$

→ LTS: $\hat{\sigma}(\tau) = \left(\frac{1}{h} \sum_{i=1}^h |r_i|^2 \right)^{\frac{1}{2}} \quad | \quad h \in [n/2, n]$

→ S-Estimator

Regression estimators β with σ given by:

→ Scale M-estimators

M-estimator of scale: Solution σ of

$$\sigma \rightarrow \frac{1}{n} \sum_{i=1}^n \rho\left(\frac{r_i}{\sigma}\right) = \sigma \quad | \quad \rho \in (0, \rho(\infty))$$

$$\sigma = 1, \rho(z) = z^2 \Rightarrow \text{RLSE}$$

$$\sigma = 0.5, \rho(z) = 1(|z| > 1) \Rightarrow \sigma = \operatorname{med}(|r_i|)$$

$$\text{with } w_\sigma(z) = \frac{\rho'(z)}{z^2} \rightarrow \hat{\sigma}^2 = \frac{1}{n \hat{\sigma}} \sum_{i=1}^n w_i r_i^2$$

$$\text{bp: } \min(\hat{\sigma}, 1 - \hat{\sigma})$$

... can be solved iteratively

→ computed using IRWLS using initial approximation

→ achieve max breakdown but low efficiency

M-estimators

$\hat{\beta}_0$ from S-estimator → IRWLS with σ (from M-scale ($\rho = \rho_{S,1}$))

→ BP of $\hat{\beta}_0$ and efficiency γ (recommended 0.95)

Affine equivariance

$$f(Ax_1 + b, \dots, Ax_n + b) = A f(x_1, \dots, x_n) + B$$

$$C(Ax_1 + b, \dots, Ax_n + b) = A C(x_1, \dots, x_n) A^T$$

$\mathcal{H} < \mathcal{D}$

n data points for which determinant of empirical covariance matrix minimal

$$h = \frac{n}{2} \rightarrow \text{max Breakdown point} \\ \rightarrow \text{low efficiency}$$

t ... mean of n obs.

C ... given by cov matrix with smallest det, scaled to obtain consistency for normal distribution

Multivariate S -estimator

$$\text{Residuals} \Rightarrow \text{small MD's} : MD^2 := (x - t)^T C^{-1} (x - t)$$

\rightarrow Using h -scale: $\min : \check{\sigma}(MD^2(x_1, t, C), \dots, MD^2(x_n, t, C)) \quad |C| = 1$

\rightarrow low asymptotic efficiency

Multivariate M -estimator

high $\sum \psi_i$ high efficiency \Rightarrow Affine equivariant, bounded IF

\rightarrow using ADW's

multivariate outlier detection

- 1) MCD for target C
- 2) $MD^2 := (x - t)^T C^{-1} (x - t)$
- 3) $MD^2 = \chi_{p+1, 0.975}^2$

$$h_{ii} = \frac{1}{n-1} MD^2 + \frac{1}{n} \rightarrow MD^2 \text{ non robust}$$

$$h_{ii} = 2 \cdot \frac{p+1}{n} \text{ identification of (vg. points)}$$

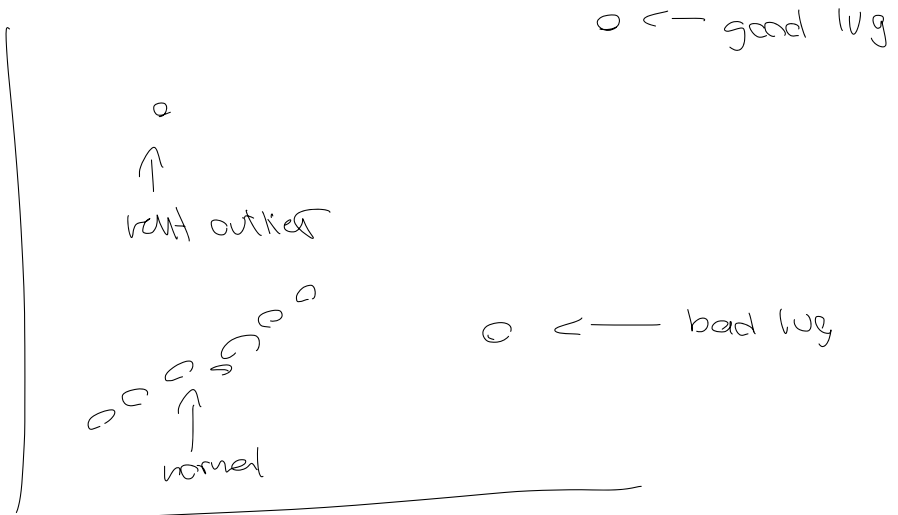
$$\text{tr}(H) = p+1 \rightarrow \text{avg}(h_{ii}) = \frac{p+1}{n}$$

$$\rightarrow \text{use robust } MD^2 := (x - t)^T C^{-1} (x - t)$$

$$MD^2 = \chi_{p+1, 0.975}^2$$

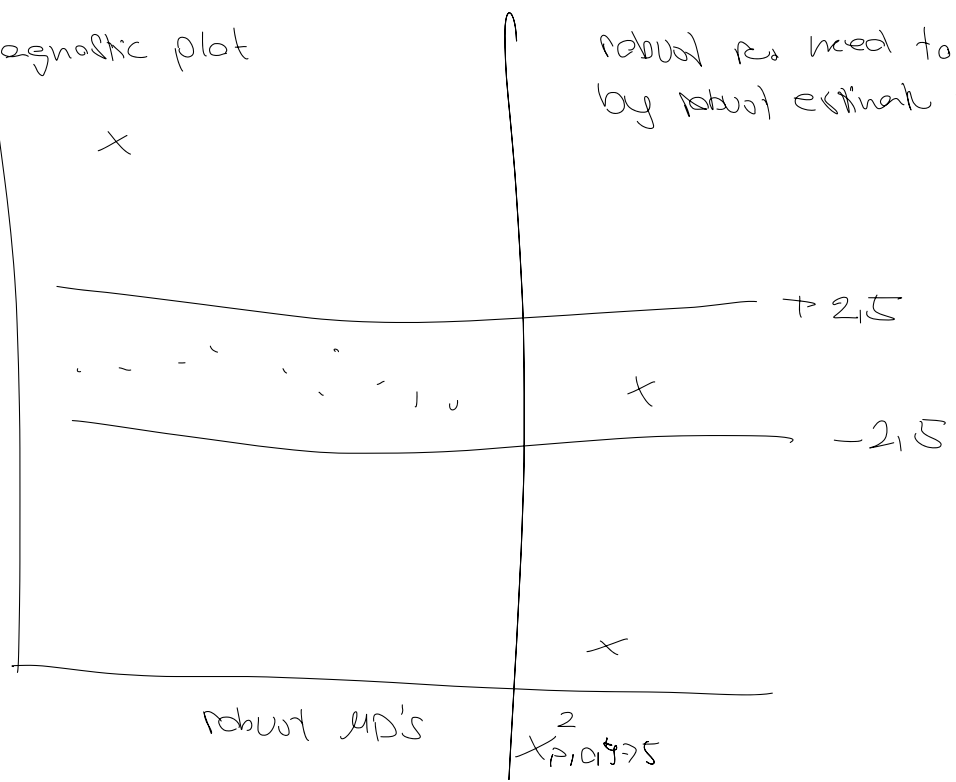
\rightarrow not possible to decide between good and bad (vg. using h_{ii})

\rightarrow only x -values considered in h_{ii}



Regression diagnostic plot

standardized robust residuals



robust res need to be scaled by robust estimate of res. scale

Robust multivariate regression

$$Y = XB + E$$

→ robustly estimate B and Σ

→ multivar. S-estimator: MD's based on r_i and C

→ minimize $\hat{\sigma}^2(r_i^T C^{-1} r_i, \dots)$ using ρ -scale ($\rho(1) = 1$)

PCA

$$z = N^T(x - \mu)$$

$$x \dots p \times 1$$

$$\mu \dots p \times p \text{ with } (\mu_1, \dots, \mu_p)$$

$$z \dots p \times 1 \dots \text{Principal components}$$

$$\text{Var}(z_i) = \mu_i^T \Sigma \mu_i$$

$$\mu_i \mu_i^T = 1, \mu_i \mu_j^T = 0$$

$$\mu^{-1} = \mu^T$$

→ Calculation of Principal Components

→ Maximise Var under restrictions ⇒ Lagrange optimization

$$\Phi_1 = \mu_1^T \Sigma \mu_1 - a_1 (\mu_1^T \mu_1 - 1)$$

→ $\Sigma \mu_1 = a_1 \mu_1 \dots \mu_1$ eigenvector of Σ to eigenvalue a_1

$$\text{Var}(z_1) = \mu_1^T \Sigma \mu_1 = \mu_1^T a_1 \mu_1 = a_1$$

of z_1 's eigenvalues, the largest one is chosen for PC₁ since we want to maximize $\text{Var}(z_1)$

$$\Phi_2 = \mu_2^T \Sigma \mu_2 - a_2 (\mu_2^T \mu_2 - 1) - b \mu_2^T \mu_1 \dots \mu_2^T \mu_1 = 0 \text{ for uncorrelatedness}$$

→ $\Sigma \mu_2 = a_2 \mu_2$ → μ_2 eigenvector to a_2 which is PC₂ (second largest PC)

→ PC solution

$$\Sigma = N A N^T \rightarrow \Sigma = N^T A N \dots \text{based on decamp of } \Sigma$$

→ Properties

$$N = (\mu_1, \dots, \mu_p) \quad A = \text{Diag}(a_1, \dots, a_p)$$

$$\hookrightarrow \text{Cor}(z_i, z_j) = 0$$

not scale invariant

$$E(z) = N^T [E(x - \mu)] \rightarrow \text{if centered } \mu = 0$$

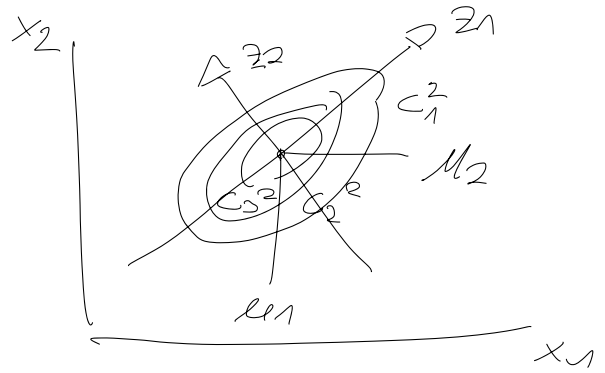
$$\text{Cov}(z) = N^T \text{Cov}(x - \mu) N = N^T \Sigma N = A \rightarrow \text{different } z \text{ uncorrelated}$$

$$\text{Cov}(x, z) = N A$$

$$\text{Cor}(x, z) = 1 = (\text{Diag}(\Sigma))^{-\frac{1}{2}} N A \frac{1}{2} \rightarrow \sigma_{ij}^{-\frac{1}{2}}$$

Graphical representation

$$(x-\mu)^T \Sigma^{-1} (x-\mu) := c^2$$



$$c^2 = (x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T M A^{-1} M^T (x-\mu)$$

$$= [M^T (x-\mu)]^T A^{-1} M^T (x-\mu) = z^T A^{-1} z = \sum_{i=1}^p \frac{z_i^2}{a_i}$$

→ components of z represent main axes of ellipsoids

↳ eigenvectors give direction

→ PC1 along largest expansion of ellipsoid and so on

↳ eigenvalues give strength

due to scaling
→ expansion = $\sqrt{a_i}$

estimation in real world conditions

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\Sigma} = S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\hat{z} = (x - \hat{\mu} \mathbf{1}^T)^T \hat{S}^{-1} \hat{\mu} \mathbf{1}^T \rightarrow \hat{S}^{-1} \hat{S} \hat{\mu} \mathbf{1}^T = \hat{A} \hat{\mu} \mathbf{1}^T$$

$\mathbf{1} \dots (\mathbf{1}_1 \dots \mathbf{1}_n)^T \dots$ needed for dimensionality

$$\hat{A} = (\hat{\lambda}_1 \dots \hat{\lambda}_p)$$

$$\text{cov}(x, z) = \hat{\lambda} = (\text{Diag}(S))^{-\frac{1}{2}} \hat{S} \hat{A}^{-\frac{1}{2}}$$

Number of relevant PC's

→ total variance: $\text{trace}(\hat{A})$... sum of eigenvalues

(1) Statistical tests

H_0 : last $p-k$ PC's contain some variance

→ start with $k=0$ until H_0 can't be rejected

$$\left(n - \frac{2p+1}{6}\right) (p-k) \ln\left(\frac{m_a}{m_g}\right) \sim \chi^2_{(p-k+2)(p-k-1)/2}$$

$$m_a = \frac{\hat{a}_{k+1} + \dots + \hat{a}_p}{p-k}$$

$$m_g = \sqrt[p-k]{\hat{a}_{k+1} \dots \hat{a}_p}$$

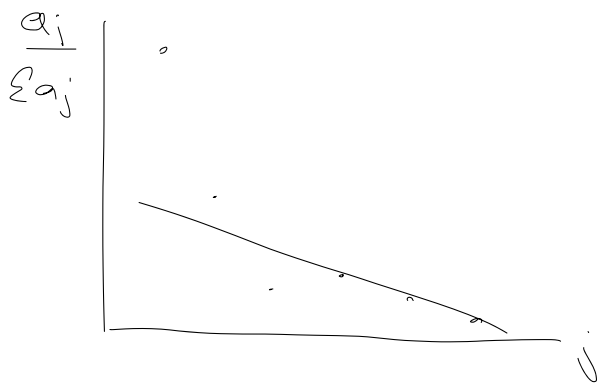
(2) Rules of thumb

$$\rightarrow \frac{\sum_{j=1}^k \hat{a}_j}{\sum_{i=1}^p \hat{a}_i} \geq \gamma_1$$

$\gamma_1 \dots \gamma_{p-1} \geq \gamma_0$

→ exclude PC's with var lower average (1 if standardized)

(3) Scree Graph



(4) Squared correlation coefficient between x and z

$$\lambda_{ij}^2 = \frac{r_{ij}^2 a_j}{s_{ii}}$$

if low for last few PC's \Rightarrow exclude

(5) resampling procedures (e.g.: Jackknife)

→ estimate uncertainty (e.g. at explained variance by k -PC's)

SVD

$$X = UDV^T$$

$X \dots n \times p$

$D \dots n \times p$ $d_{ii} \geq 0$ singular values $i = \text{rank}(X)$

$U \dots n \times n$, orthogonal eigenvectors of XX^T in cols

$V^T \dots p \times p$, orthogonal eigenvectors of $X^T X$ in cols

$$Z = (X - 1X^T) \Gamma \rightarrow \text{mean centered: } \tilde{Z} = X \Gamma \rightarrow X = Z \Gamma^T$$

$$\rightarrow X = UDV^T = Z \Gamma^T$$

$(x-u)^T(x-u) = x^T x$ since centered

$$S = \frac{1}{n-1} X^T X = \hat{M} \hat{A} \hat{M}^T \rightarrow M \text{ matrix with normalized eigenvectors of } S \text{ which in this case is } U \text{ (of } X^T X) \rightarrow \hat{M} \equiv U$$

$\hat{M} = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T$

$$\Rightarrow X = UDV^T = ZU^T \rightarrow Z = UD$$

$$\|X\|_F = \sqrt{\sum_{i=1}^n \|x_i\|_2^2} = \sqrt{\sum_{i=1}^n \sum_{j=1}^p x_{ij}^2}$$

← measure of total power of matrix

$X \dots n \times p$
 $V_m \dots n \times m$
 $B \dots p \times m$

$XV = UB = Z$ transformation matrix
 $\text{rank}(B) \leq m \leq \text{rank}(X)$ $BB^T = I$

→ "maximising variance of projected data"

$$V_m = \underset{B}{\text{argmax}} \|XB\|_F^2$$

mapping of x onto lower dim space $n \times p \Rightarrow n \times m$

→ "minimising RSS between $q.f.$ and projection"
 ↳ maximising power of projection $\hat{=}$ maximising variance

$$X = XVV^T = XV_m V_m^T + \Sigma$$

$$V_m = \underset{B}{\text{argmin}} \|X - \tilde{X}\|_F^2$$

$$\tilde{X} = XBB^T$$

↳ minimizing F of $XB \hat{=}$ minimizing power difference

Biplots

$$X \approx X_{(2)} = UDU^T = (u_1 \ u_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix} \begin{pmatrix} v_1^T \\ v_2^T \end{pmatrix}$$

$$X_{(2)} = GH^T$$

$$G = (u_1 \ u_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^{c-1}$$

→ gives coord's of data points in reduced space

$$H = (v_1 \ v_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^c$$

→ determines direction and magnitude of arrows indicating eq. vars

$0 \leq c \leq 1 \dots$ determines distribution of singular values to G and H

Obtained for $c=1$ (+ rescaling):

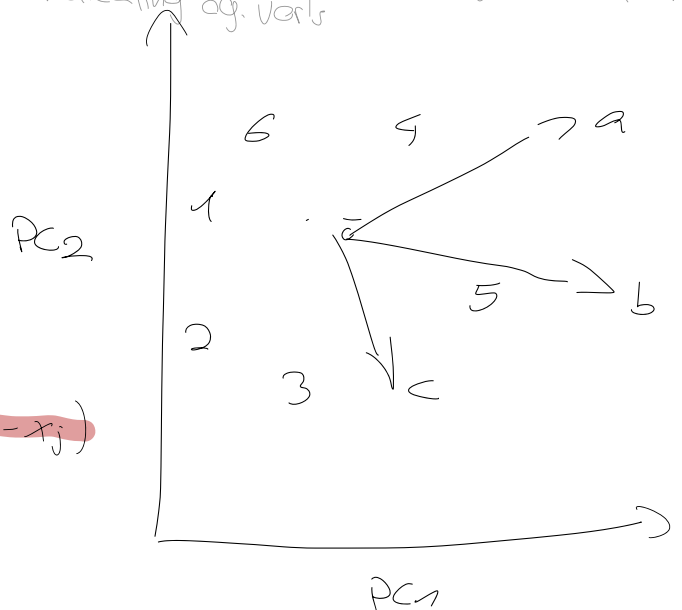
$$\rightarrow g_i^T h_j \approx x_{ij}$$

$$\rightarrow HH^T \approx S = \frac{1}{n-1} X^T X$$

$$\rightarrow \|h_j\|^2 \approx \text{Var}(x_j)$$

$$\rightarrow \cos(h_i, h_j) \approx r_{ij}$$

$$\rightarrow \|g_i - g_j\|^2 \approx MD^2 = (x_i - x_j)^T S^{-1} (x_i - x_j)$$



Diagnostics

$$SD_i = \left(\sum_{j=1}^k \frac{z_{ij}^2}{\hat{q}_j} \right)^{\frac{1}{2}}$$

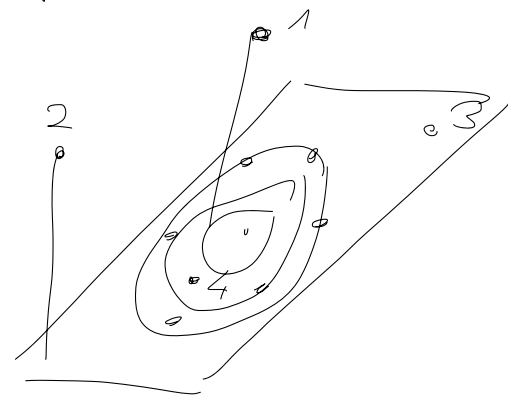
... only k relevant PCs considered

\hat{q}_j = MD score vector to PCA center with respect to covariance matrix A

$$OD_i = \|x_i - \hat{\mu}_k\|_2$$

$\hat{\mu}_k$ = $\| \cdot \|_2$ of observation to projection into space of first k PC's

- 1: big OD_i , low SD_i → vertical outlier
- 2: big OD_i , big SD_i → bad leverage point
- 3: low OD_i , big SD_i → good leverage point
- 4: low OD_i , low SD_i → good point



→ outlier detection:

$$SD_i \approx \sqrt{\frac{2}{k} \chi_{k, 0.975}^2}$$

$$OD_i > \left(\text{median}_i \left(OD_i^{\frac{2}{3}} \right) + \text{MAD} \left(OD_i^{\frac{2}{3}} \right) z_{0.975} \right)^{\frac{3}{2}} \dots OD_i^{\frac{2}{3}} \text{ closer to normality}$$

$$\text{MAD} = 1.483 \cdot \text{median}_i \left(|y_i - \text{median}_j(y_j)| \right) \quad z_{0.975} \dots 0.975 \text{ quantile of } N(0,1)$$

Factor Analysis

$$y_i = \frac{x_i - E(x_i)}{\sqrt{\text{Var}(x_i)}} \quad | \quad i=1, \dots, p \quad \dots \text{data standardized}$$

$$y = (y_1, \dots, y_p)^T \quad p \times 1$$

$$f = (f_1, \dots, f_k) \quad k < p \quad \text{factors}$$

$$y = \Lambda f + e$$

$$\Lambda = [\Lambda_{ij}] \quad \dots \quad p \times k \quad \text{loadings matrix}$$

$$e = (e_1, \dots, e_p)^T \quad p \times 1 \quad \text{unique factor}$$

$$E(f) = 0 = E(e)$$

$$\text{Cov}(e_i, e_j) = 0 = \text{Cov}(e_i, e_j) \quad i \neq j \quad \rightarrow \quad \text{Cov}(e) = \Psi = \text{Diag}(\psi_{11}, \dots, \psi_{pp})$$

$$\text{Var}(f_i) = 1$$

$$\rho = \text{Corr}(x) = \overset{\text{data standardized}}{\text{Cov}(y)} = \text{Cov}(\Lambda f + e) =$$

$$= \Lambda \text{Cov}(f) \Lambda^T + \underbrace{\Lambda \text{Cov}(f|e)}_0 + \underbrace{\text{Cov}(e|f) \Lambda^T}_0 + \text{Cov}(e) = \Lambda \Phi \Lambda^T + \Psi$$

$$\text{factors uncorrelated} \Rightarrow \Phi = I \Rightarrow \rho = \Lambda \Lambda^T + \Psi$$

$$\rho_{\text{red}} = \rho - \Psi = \Lambda \Lambda^T = \begin{bmatrix} \lambda_{11}^2 & \dots & \lambda_{1p}^2 \\ \lambda_{21}^2 & \dots & \lambda_{2p}^2 \\ \vdots & \ddots & \vdots \\ \lambda_{k1}^2 & \dots & \lambda_{kp}^2 \end{bmatrix} \rightarrow \lambda_{i1}^2 = 1 - \psi_{ii} = \sum_{j=1}^k f_{ij}^2 \quad \dots \text{ communalities}$$

$f_{ij}^2 \Rightarrow$ % variance explained by factors

Non-uniqueness

\Rightarrow factors loadings not unique \Rightarrow additional constraints for uniqueness:

$$\Lambda^T \Psi^{-1} \Lambda \quad \text{or} \quad \Lambda \Lambda^T \quad \dots \text{diagonal}$$

parameters

upper bound for k : $s > 0 \rightarrow$ fewer parameters than corr. matrix

$$s = \text{corr. matrix} - (\text{factor model} - \text{restrictions}) = \frac{1}{2}(p-k)^2 - \frac{1}{2}(p+k)$$

Parameter estimation using PFA (Principal factor analysis)

→ communalities

(1) highest correlation coefficient

$$\max_{i \neq j} |\hat{\rho}_{ij}^1| \rightarrow \text{may indicate strongest lin. relationship}$$

(2) Squared multiple correlation coefficient

$$\hat{\rho}_{1,2,\dots,p}^2 = 1 - \frac{1}{\hat{\rho}_{ii}} \rightarrow \text{var. proportion of } i\text{th var explained by other var's}$$

(3) Iterative estimation

- 1) fix k (criterion from PCA on # components)
- 2) initialize \hat{k}_i^2 using (1) or (2)
- 3) replace $\hat{\rho}_{ii}^1$ with \hat{k}_i^2
- 4) estimate $\hat{\Lambda}$ (see below)
- 5) re-estimate $\hat{k}_i^2 = \sum_{j=1}^k \hat{\rho}_{ij}^2 \quad | i=1, \dots, p$
- 6) repeat until stable

→ loadings

$$\hat{\rho}_{red} = \hat{\Lambda} \hat{\Lambda}^T = \hat{\rho} - \hat{\Psi} = \hat{\Lambda} \hat{\Lambda}^T = \underbrace{\hat{\Lambda}_{1:k} \hat{\Lambda}_{1:k}^T}_{\text{model}} + \sum_{i=k+1}^p a_i \gamma_i \gamma_i^T \Rightarrow \hat{\Lambda} = \hat{\Lambda}_{1:k} \hat{\Lambda}_{1:k}^{\frac{1}{2}}$$

→ Update uniqueness

$$\hat{\Psi}_{ii} = 1 - \sum_{j=1}^k \hat{\rho}_{ij}^2 \rightarrow \text{valid if } \hat{\Psi}_{ii} > 0$$

Rotation

$$\sum_{s: j=1}^k \sum_{i=1}^p (\lambda_{is} \lambda_{ij})^2 \rightarrow \min \Rightarrow \text{loadings matrix simpler: many small absolute values, only a few large ones}$$

→ orthogonal rotation

orthogonal $k \times k$ matrix $T \rightarrow \hat{\Lambda} = \Lambda T$

→ factors uncorrelated (90° to each other)

→ do not change k

→ $(k_i^2)^2 = \text{const}$

→ sum over all variables λ const. \Rightarrow maximise one term $\hat{=}$ minimize the other

$\hat{=}$ maintains overall structure of factor solution in terms of variance

(1) $\text{QMAX} = \sum_{i=1}^p \sum_{j=1}^k \hat{\lambda}_{ij}^4 \Rightarrow$ one dominant factor

(2) $\text{VLEX} = \frac{1}{p} \sum_{j=1}^k \sum_{i=1}^p \left(\frac{\hat{\lambda}_{ij}}{k_i} \right)^4 - \sum_{j=1}^k \left[\sum_{i=1}^p \left(\frac{\hat{\lambda}_{ij}}{k_i} \right)^2 \right]^2 \Rightarrow$ max $\hat{=}$ absolute large and small loadings

→ normalize with communalities, that dominate

→ oblique rotation

→ ~~factors no longer uncorrelated~~ $\Rightarrow \phi$ needs to be considered

→ $\hat{\Lambda} = \Lambda T$... Two bigger orthogonal

→ $P_{\text{red}} = \hat{\Lambda} \text{Cov}(\hat{f}) \hat{\Lambda}^T \rightarrow \hat{f} = T^{-1} f \Rightarrow T$ must be invertible

→ maximises variance of squared loadings of each factor

(1) $\text{QRELW} = \sum_{s: j=1}^k \sum_{i=1}^p \hat{\lambda}_{is}^2 \hat{\lambda}_{ij}^2 \Rightarrow$ discourages multiple large loadings

(2) $\text{OBRILW} = \sum_{s: j=1}^k \left(\sum_{i=1}^p \hat{\lambda}_{is}^2 \hat{\lambda}_{ij}^2 - \frac{\gamma}{p} \sum_{i=1}^p \hat{\lambda}_{is}^2 \sum_{i=1}^p \hat{\lambda}_{ij}^2 \right) \Rightarrow$ higher (low) loadings but not medium

→ $\gamma = 0 \rightarrow \text{QRELW}$

→ $\gamma = 1 \rightarrow \text{OBRILW}$

Estimation of factor scores

→ Bartlett

$y = \Lambda f + e$... heteroscedastic (cov(e) not equal) $\Rightarrow \cdot \Psi^{-\frac{1}{2}}$ (weights)

$\Rightarrow \Psi^{-\frac{1}{2}} y = \Psi^{-\frac{1}{2}} \Lambda f + \Psi^{-\frac{1}{2}} e \rightarrow \Psi^{-\frac{1}{2}} \text{Cov}(e) \Psi^{-\frac{1}{2}} = \Psi^{-\frac{1}{2}} \Psi \Psi^{-\frac{1}{2}} = I \Rightarrow$ homoscedastic

$\Rightarrow \hat{f} = (\underbrace{\Lambda^T}_{\beta^T} \underbrace{\Psi^{-1}}_{\Psi^{-1}} \underbrace{\Lambda}_{\Lambda})^{-1} \underbrace{\Lambda^T}_{\beta^T} \underbrace{\Psi^{-1}}_{\Psi^{-1}} \underbrace{y}_{y}$... used when aim is to reproduce factors as close to original as possible

$\Rightarrow \hat{F} = Y \Psi^{-1} \Lambda (\Lambda^T \Psi^{-1} \Lambda)^{-1}$

→ Regression

$f = B y + \rho$ regresses f on y

$\rightarrow \hat{B} = f y^T (y y^T)^{-1} \Rightarrow \hat{f} = \hat{B} y = f y^T (y y^T)^{-1} y$

$\rightarrow \hat{f} = f (\Lambda f + e)^T (y y^T)^{-1} y = \Lambda f f^T + \underbrace{e f^T}_{\text{Cov}(e, f) = 0} (y y^T)^{-1} y = \Lambda f f^T (y^T y)^{-1} y$

→ sample version: $\hat{F} = \frac{\hat{F} \hat{F}^T}{n-1} \frac{n-1}{(Y^T Y)} Y =$

$= Y R^{-1} \Lambda \hat{\Phi} \rightarrow Y R^{-1} \Lambda$... used for orthogonal factors
 \uparrow $\underbrace{\quad}_{= I}$ for orthogonal
 correlation matrix of observed variables (y)

Correlation analysis

Multiple correlation analysis

→ dependency of feature x on p -dim. feature y

→ minimize $RSE = E(x - a_0 - a^T y)^2$

$$\begin{matrix} x \dots 1 \times 1 \\ y \dots 1 \times p \end{matrix}$$

Linear prediction function

→ $a_0 + a^T y$

$$\rightarrow a_0 = \mu_x - a^T \mu_y \quad a = \Sigma_{yy}^{-1} \sigma_{yx}$$

$$\rightarrow \mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{xx} & \sigma_{yx}^T \\ \sigma_{yx} & \Sigma_{yy} \end{pmatrix}$$

$$\rightarrow RSE = \sigma_{xx} - \sigma_{yx}^T \Sigma_{yy}^{-1} \sigma_{yx}$$

→ linear prediction function of x with minimal RSE and max corr with x

$$\text{corr}(x, a_0 + a^T y) = \frac{\sigma_{yx}^T \Sigma_{yy}^{-1} \sigma_{yx}}{\sigma_{xx}}$$

multiple correlation coefficient:

→ correlation between x and best linear predictor function

$$r_{x,y} = \sqrt{\frac{\sigma_{yx}^T \Sigma_{yy}^{-1} \sigma_{yx}}{\sigma_{xx}}}$$

→ $r_{x,y}^2$ multiple coefficient of determination

→ indicates how well feature x is explained by properties $y = (y_1, \dots, y_p)^T$

Test for hypothesis multiple correlation is zero $\stackrel{!}{=} \text{bivariate correlation zero}$

$$H_0: r_{x,y} = 0$$

assume normal distribution, significance (α) R , n , # observations in sample

$$F = \frac{(n-1-p) r_{x,y}^2}{p(1-r_{x,y}^2)} > F_{p, n-1-p, 1-\alpha}$$

Canonical correlation

$$x \dots 1 \times p \quad y \dots 1 \times q \quad p \leq q$$

- linear dependence between two group of variables \Rightarrow subspace (one corr. coeff does not suffice)
- k -th canonical pair of variables

$$y_k = a_k^T x = e_k^T \Sigma_{11}^{-\frac{1}{2}} x \quad m_k = b_k^T y = f_k^T \Sigma_{22}^{-\frac{1}{2}} y \quad \dots \text{linear combinations}$$

with:

$$\Sigma_{11} = E[(x - \mu_1)(x - \mu_1)^T]$$

$$\Sigma_{22} = E[(y - \mu_2)(y - \mu_2)^T]$$

$$\Sigma_{12} = \Sigma_{21}^T = E[(x - \mu_1)(y - \mu_2)^T]$$

- maximise $\text{Corr}(y_k, m_k) = \rho_k$ over all linear comb. uncorrelated with last $1, 2, \dots, k-1$ canonical variables

- eigenvalues and eigenvectors

$$\rho_1^2 \geq \dots \geq \rho_p^2 \quad a_i \text{ of } \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-\frac{1}{2}} \rightarrow \text{respective } f_1, e_1, \dots, e_p$$

→ $\geq \dots$ eigenvalues sorted in descending corr. magnitude

$$\rho_1, \dots, \rho_p \quad p \text{ largest } a_i \text{ of } \Sigma_{22}^{-\frac{1}{2}} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \rightarrow \text{respective } f_1, \dots, f_p$$

→ Every f_i proportional to $\Sigma_{22}^{-\frac{1}{2}} \Sigma_{21} \Sigma_{11}^{-1} e_i$

→ standard: same order for a_i 's as above

- solution is derived from eigenvalues of combined covariance matrices, thus an analysis of covariance structures between the sets

→ Proof: using Cauchy-Schwarz inequality and reformulation as Eigenvalue Problem

\Rightarrow eigenvalues of combined matrix represent directions in which corr. is maximised

- Properties

$$\text{Var}(y_k) = \text{Var}(m_k) = 1$$

$$\text{Cov}(y_k, y_l) = \text{Cov}(m_k, m_l) = \text{Cov}(y_k, m_l) = 0 \quad k \neq l$$

$$\Rightarrow \text{Cov}\begin{pmatrix} y \\ m \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \begin{matrix} f \dots (f_1, \dots, f_p)^T \\ m \dots (m_1, \dots, m_p)^T \\ p \dots \text{Diag}(p_1, \dots, p_p) \end{matrix}$$

$$\rho_1 = 1 ?$$

→ at least one $x_i = y_j$

→ indicates perfect linear relationship between linear comb. of sets X and Y

→ implications: linear-comb. of these variables completely correlated → lin. dependency between X and Y , no conclusions on relationships of other variables can be made if
 $x = (1, 0, \dots, 0)$ $y = (1, 0, \dots, 0)$

canonical cor. coeff. invariant to transformations

$$x^* = U^T x + U \quad y^* = V^T y + V \quad U \dots p \times p \quad V \dots q \times q$$

$$U \dots 1 \times p \quad V \dots 1 \times q$$

→ can. corr. between of $x, y =$ between x^*, y^*

$$\rightarrow a_i^* = U^{-1} a_i \quad b_i^* = V^{-1} b_i$$

Tests

→ Likelihood ratio test $H_0: \Sigma_{12} = 0$

$$\lambda^{2/n} = |I - S_{22}^{-1} S_{21} S_{11}^{-1} S_{12}| = \prod_{i=1}^p (1 - r_i^2) \rightarrow \lambda(q, n-1-p, p) \text{ Wilks' distr.}$$

↑
sample canonical correlation coefficient

→ Bartlett's approximation

→ for large n instead of 1

$$-\left[n - \frac{1}{2}(p+q+3) \right] \ln \prod_{i=1}^p (1 - r_i^2) \sim \chi_{pq}^2$$

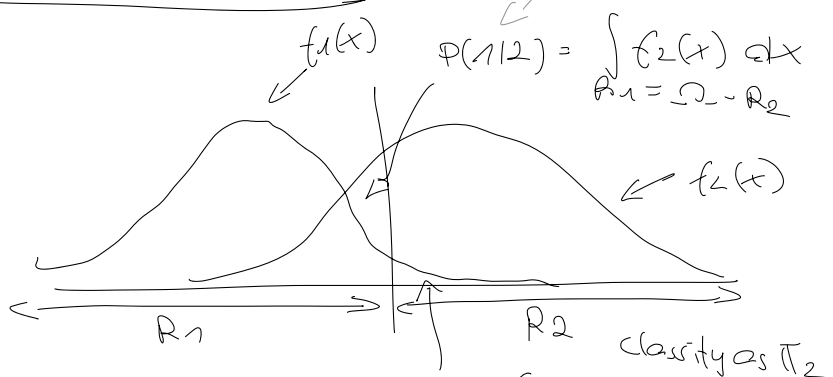
→ Test only s can. cor. are $\neq 0$

$$-\left[n - \frac{1}{2}(p+q+3) \right] \ln \prod_{i=s+1}^p (1 - r_i^2) \sim \chi_{(p-s)(q-s)}^2$$

Discriminant Analysis

$$\rightarrow P(A|B) = P(B|A) \cdot P(A)$$

Satz von Bayes: prob's updated based on new evidence



$\mathcal{R}_x \dots$ space of observations to which we assign objects Π_x

$\Pi_x \dots$ "real" population of x

$$f(2|1) = P(x \in R_2 | \Pi_1) \dots$$

x is falsely classified as Π_2

class. as Π_1 $f(2|1) = \int_{R_2} f_1(x) dx$
 $R_2 = \Omega - R_1$
 class. as Π_2 $f(1|2) = \int_{R_1} f_2(x) dx$
 $R_1 = \Omega - R_2$

\rightarrow Application of Bayes theorem

$$P(\Pi_i|x) = \frac{P(x|\Pi_i)P(\Pi_i)}{\sum_j P(x|\Pi_j)P(\Pi_j)} \rightarrow P(x \in \Pi_i) = P(x \in R_i | \Pi_i)P(\Pi_i)$$

\rightarrow Prior probability: $p_i \dots$ prop that obj. comes from $\Pi_i \dots P(\Pi_i)$

$$\rightarrow E_{cl} : c(2|1)P(2|1)p_1 + c(1|2)P(1|2)p_2$$

\rightarrow class. rule to minimize E_{cl}

$$R_1, \text{ obs which: } \frac{f_1(x)}{f_2(x)} \geq \frac{c(1|2)p_1}{c(2|1)p_2} \quad R_2 : <$$

$$\rightarrow TP_{cl} : p_1 P(2|1) + p_2 P(1|2)$$

$$g=2$$

$$\Sigma_1 = \Sigma_2 = \Sigma$$

$$P(x \in \Pi_1) = P(x \in R_1 | \Pi_1)P(\Pi_1)$$

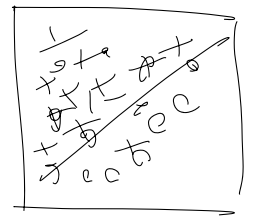
\rightarrow minimizing E_{cl} for $\Pi_1, \Pi_2 \sim N(\mu_i, \Sigma)$

$$(1) x_0 \text{ to } \Pi_1 : (p_1 - p_2)^T \Sigma^{-1} x_0 - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln \frac{c(1|2)p_1}{c(2|1)p_2}$$

\rightarrow linear in x

\rightarrow sample version

$$S_{pooled} = \frac{1}{n_1 + n_2 - 2} \sum_{l=1}^2 \sum_{j=1}^{n_l} (x_{lj} - \bar{x}_l)(x_{lj} - \bar{x}_l)^T$$



→ unbiased estimate of Σ , $\hat{\Sigma}$ cov. estimate of group data

(1b) $x_0 \text{ to } \Pi_1: (\bar{x}_1 - \bar{x}_2)^T \text{Spooled } x_0 - \frac{1}{2} (\bar{x}_1 - \bar{x}_2)^T \hat{\Sigma}^{-1} (\bar{x}_1 + \bar{x}_2) \geq \ln \frac{C(1|2)P_1}{C(2|1)P_2}$

→ simplification

→ sums and scalars; if $\ln \left(\frac{C(1|2)P_1}{C(2|1)P_2} \right) = 1 \rightarrow \ln(u) = 0$, then:

$\hat{y}_1 = \underbrace{(\bar{x}_1 - \bar{x}_2)^T}_{a^T} \text{Spooled}^{-1} x = a^T x$

$\hat{m}_1 = \frac{1}{2} (\bar{x}_1 - \bar{x}_2)^T \text{Spooled}^{-1} (\bar{x}_1 + \bar{x}_2) = \frac{1}{2} (y_1 + y_2)$

→ $y_1 = (\bar{x}_1 - \bar{x}_2)^T \text{Spooled}^{-1} \bar{x}_1 = a^T \bar{x}_1$

→ $y_2 = (\bar{x}_1 - \bar{x}_2)^T \text{Spooled}^{-1} \bar{x}_2 = a^T \bar{x}_2$

} y results from linear comb. of observations

⇒ $x_0 \text{ to } \Pi_1: \hat{y}_1 \geq \hat{m}_1 \Rightarrow$ linear decision boundary in

⇒ for p -dim features → 1. dim variable y

$\Sigma_1 \neq \Sigma_2$

→ minimizing ECR for $\Pi_1, \Pi_2 \sim N(\mu_i, \Sigma_i)$

(2) $x_0 \text{ to } \Pi_1: -\frac{1}{2} x_0^T \underbrace{(\Sigma_1^{-1} - \Sigma_2^{-1})}_{\Sigma_0 \text{ alone}} x_0 + \underbrace{(\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1})}_{\mu_0^T \Sigma_0} x_0 - k \geq \ln \frac{C(2|1)P_1}{C(1|2)P_2}$

with $k = \frac{1}{2} \ln \left(\frac{|\Sigma_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\mu_1^T \Sigma_1^{-1} \mu_1 - \mu_2^T \Sigma_2^{-1} \mu_2)$

→ quadratic in x_0

$\mu_0^T \Sigma_0$ together

Evaluation

→ training set

(1) $AER = P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$

(2) $\frac{\# \text{ misclassified}}{\text{total}}$

→ no training set

(3) jackknife

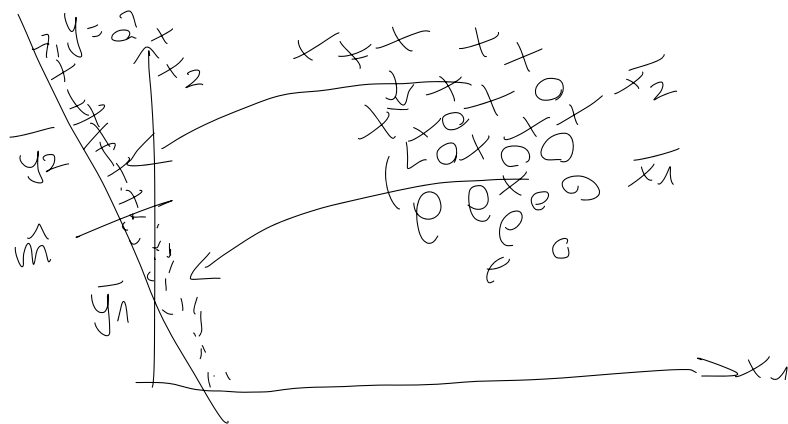
→ classify Π_1, Π_2 multiple times, omitting an observation as test case

estimated error rate: $\frac{\bar{n}_1 + \bar{n}_2}{n_1 + n_2}$, $\hat{P}(1|2) = \frac{\bar{n}_2}{n_2}$, $\hat{P}(2|1) = \frac{\bar{n}_1}{n_1}$

(4) k-fold cross validation

→ k parts, k-1 parts computed, last part for evaluation

Fisher



→ projected onto straight line → \hat{a} varied until max separation

$\frac{|\bar{y}_1 - \bar{y}_2|}{s_y} \rightarrow \max$ → pooled variance: $S_y^2 = \frac{1}{n_1 + n_2 - 2} \sum_{g=1}^2 \sum_{i \in \mathcal{I}_g} (x_{ig} - \bar{x}_g)^2$

$\hat{y} = \hat{a}^T x = (\bar{x}_1 - \bar{x}_2)^T S_{pooled}^{-1} x$ maximises $\frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(\hat{a}^T \bar{x}_1 - \hat{a}^T \bar{x}_2)^2}{\hat{a}^T S_{pooled} \hat{a}}$ over all \hat{a}
↑ pooled variance of y-values

→ maximum ratio: $D^2 = (\bar{x}_1 - \bar{x}_2)^T S_{pooled}^{-1} (\bar{x}_1 - \bar{x}_2)$... $2 \rho D^2$

→ transform multivariate obs to univariate, maximising separation

→ classification rule

x_0 to Π_1 : $\hat{y}_0 \geq \hat{m} := (\bar{x}_1 - \bar{x}_2)^T S_{pooled}^{-1} (\bar{x}_1 + \bar{x}_2)$... $\hat{y}_0 = \hat{y}$ with x_0

$$g > 2$$

process minimizing $E_{\mathcal{D}}$: x to Π_k if $\sum_{l=1, k \neq l}^g p_l f_l(x) c(k|l)$ minimal

$\rightarrow c$ equal for all $\Rightarrow x$ to Π_1 if $p_k f_k(x) > p_i f_i(x) \forall i \neq k$

$\rightarrow p_i, c, f$ must be known

$N(\mu_i, \Sigma_i)$ assumed

$\rightarrow \Sigma_1 = \Sigma_2 = \dots = \Sigma_n$

similar to (1) with μ_k instead of μ_1, μ_2

(3) $d_k(x) = \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \ln p_k$ largest $d_1(x) \dots d_g(x)$

\rightarrow pooled version (no training set)

$$S_{\text{pooled}} = \frac{1}{\sum_{i=1}^g n_i - g} \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)^T$$

(3b) x to Π_k if: $\hat{d}_k(x) = \bar{x}_k^T S_{\text{pooled}}^{-1} x - \frac{1}{2} \bar{x}_k^T S_{\text{pooled}}^{-1} \bar{x}_k + \ln p_k$ largest $\hat{d}_1(x) \dots \hat{d}_g(x)$

$\rightarrow \Sigma_1 \neq \Sigma_2 \neq \dots = \Sigma_n$

(4) x to Π_k if: $d_k^{\text{Q}}(x) = \underbrace{-\frac{1}{2} \ln |\Sigma_k|}_{\text{first term of } k} - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \ln p_k$ largest $d_1^{\text{Q}} \dots d_g^{\text{Q}}$

QD^2

\rightarrow estimated from training set

(4b) $\hat{d}_k^{\text{Q}}(x) = -\frac{1}{2} \ln |S_i| - \frac{1}{2} (x - \bar{x}_i)^T S_i^{-1} (x - \bar{x}_i) + \ln p_i$ largest $\hat{d}_1^{\text{Q}} \dots \hat{d}_g^{\text{Q}}$

$S_i \dots$ sample covariance matrix

Jackknife

$$\frac{\sum_{i=1}^g \frac{1}{n_i}}{\sum_{i=1}^g \frac{1}{n_i}}$$

Fisher

$$\rightarrow \Sigma_1, \dots, \Sigma_g = \Sigma$$

$$\sum_{i=1}^g p_i (\mu_{i,y} - \bar{\mu}_y)^2 \rightarrow \max$$

← weighted sum of squared distances of group - to the total mean

$$\sigma_y^2$$

↳ variance

$$\bar{\mu}_y = a^T \bar{\mu} = a^T \sum_{i=1}^g p_i \mu_i$$

$$\mu_{i,y} = E(y | x \in \pi_i) = a^T E(x | x \in \pi_i) = a^T \mu_i$$

$$\sigma_{y,i}^2 = \text{Var}(y | x \in \pi_i) = a^T \text{Cov}(x | x \in \pi_i) a = a^T \Sigma_i a \Rightarrow \sigma_y^2 = a^T \Sigma a$$

↑
 $\Sigma_1 = \dots = \Sigma_n$

$$\rightarrow \Sigma_1, \dots, \Sigma_g \neq \Sigma$$

→ resulting classification rule non-optimal ⇒ replace Σ by pooled variance

→ solution of maximization problem given by a_1, \dots, a_l of $W^{-1}B$

← here eigenvectors!

→ must be scaled: $a_i^T W a_i = 1 \quad | \quad i = 1, \dots, l$

→ $l = \min(g-1, p)$... $l = \#$ strictly positive eigenvalues ... $\#$ significant directions

max rank $W = g-1$ ← max rank $B = \#$ vars p

$$\Rightarrow \frac{a^T B a}{a^T W a} \quad (a \in \mathbb{R}^p, a \neq 0)$$

→ variation within groups: $W = \sum_{i=1}^g p_i \Sigma_i$

→ variation between groups: $B = \sum_{i=1}^g p_i (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T$

→ discriminant function: $y = a_j^T x$... projection of norm var x in direction a_j

→ classification rule

$$x \text{ to } \pi_k: \text{ if } d_k^F(x) = \sum_{j=1}^l (y_j - \mu_{j,y_j})^2 - 2 \log p_i \text{ smallest } d_1^F(x), \dots, d_g^F(x)$$

Euclidean distance in discriminant space

$$y_j = a_j^T x$$

$$\mu_{j,y_j} \stackrel{\wedge}{=} a_j^T \mu_i$$

→ also: $\sum_{j=1}^l (a_j^T (x - \mu_i))^2 - 2 \log p_i = (x - \mu_i)^T A A^T (x - \mu_i) - 2 \log p_i$

$\hat{=} \text{APD}^2$ in original space

Symbolic Data Analysis

$$Y_j: S \rightarrow B_j$$

$S = \{s_1, \dots, s_n\} \dots$ set of n units to be analyzed

$Y_1, \dots, Y_p \dots$ variables

$O_j \dots$ underlying domain of Y_j } $j = 1, \dots, p$
 $B_j \dots$ observation space of Y_j

$\rightarrow Y_j$ single valued variable $B_j \equiv O_j$

$\rightarrow Y_j$ multi valued variable $B_j = P(O_j)$

$\rightarrow Y_j$ interval variable (numeric) $Y_j: S \rightarrow B: Y_j(s_j) = [l_{ij}, u_{ij}], l_{ij} \leq u_{ij}$

$\rightarrow Y_j$ categorical model or histogram variable: $Y_j: B_j$ set of distributions on O_j

\rightarrow interval variable special case of histogram with $H_{Y_j(s_j)} = ([l_{ij}, u_{ij}], 1)$

Parametric Models

\rightarrow probabilistic models for interval variables

$$Y_j: S \rightarrow B: Y_j(s_i) = [l_{ij}, u_{ij}], l_{ij} \leq u_{ij} \quad \forall s_i$$

$$\begin{array}{cccc} & Y_1 & \dots & Y_p \\ s_1 & [l_{11}, u_{11}] & \dots & [l_{1p}, u_{1p}] \\ \vdots & \vdots & & \vdots \\ s_n & [l_{n1}, u_{n1}] & \dots & [l_{np}, u_{np}] \end{array}$$

\rightarrow Represent $Y_j(s_i)$ using:

• midpoint $c_{ij} = \frac{l_{ij} + u_{ij}}{2}$

• range $r_{ij} = u_{ij} - l_{ij}$

Gaussian model

\rightarrow assume: joint distribution of C 's and \log s of r 's is multivariate normal

$$(C, R^*) \sim N_{2p}(\mu, \Sigma)$$

$$\sum_{c_1, \dots, c_p, r_1, \dots, r_p}$$

with $\mu = [\mu_C^T, \mu_{R^*}^T]$ $\Sigma = \begin{bmatrix} \Sigma_{CC} & \Sigma_{CR^*} \\ \Sigma_{RC^*} & \Sigma_{R^*R^*} \end{bmatrix}$ $R^* := \ln(R)$

\rightarrow advantage: simple application of inference methods

Skew-Normal model

→ same assumption

→ more general model, less limitations

$$(C, R^T) \sim \text{SN}_{2p}(\xi, \Omega, \beta)$$

ξ ... location $1+p$

Ω ... scale, $p \times p$, symmetric positive definite

β ... shape (additional), $1+p$

↳ indicates skewness (e.g. $\beta > 0 \rightarrow$ longer right tail)

$$f(y; \xi, \Omega, \beta) = \frac{2}{\Omega} \phi_p(x - \xi; \Omega) \Phi(\beta^T \omega^{-1}(x - \xi)) \quad | x \in \mathbb{R}^p$$

ω ... $\text{Diag}(\sqrt{\Omega_{ii}})$

ϕ_p ... dens. of p -dim std. Gaussian vector

Φ ... Distr. func of std. norm var

↳ becomes 0.5 for $\beta = 0 \Rightarrow$ then $\text{SN}_{2p} = N$

Covariance Structures

→ C_{ij}, R_{ij} : two quantities related to one variable → dependencies?

→ Parametrization of covariance matrix which takes this into account

1) Non-restricted

$$\begin{bmatrix} \Sigma_{CC} & C_{CR} \\ \Sigma_{RC} & \Sigma_{RR} \end{bmatrix} \text{ no-restrictions}$$

2) y_j 's non correlated

$$\begin{bmatrix} \circ & \circ & \circ & \circ \\ \circ & \circ & \circ & \circ \\ \circ & \circ & \circ & \circ \\ \circ & \circ & \circ & \circ \end{bmatrix} \quad p \text{ } 2 \times 2 \text{ blocks}$$

Diag(a_{ii})

3) C's non correlated with r's

$$\begin{bmatrix} \Omega & 0 \\ 0 & \Omega \end{bmatrix} \quad \begin{array}{l} 2 \text{ } p \times p \text{ blocks} \\ \Sigma_{CR} = 0 = \Sigma_{RC} \end{array}$$

4) All C's and r's non correlated (incl. with self)

$$\begin{bmatrix} \circ & \circ \\ \circ & \circ \end{bmatrix} \quad \begin{array}{l} 2p \text{ blocks ... single real elements} \\ (3) \neq \text{Diag}(\Sigma_{CC}), \text{Diag}(\Sigma_{RR}) \end{array}$$

ML-estimation for interval data

→ Gaussian

empirical covariance matrix

$$\ln L(\mu, \Sigma) = -np \ln(2\pi) - \frac{n}{2} (\bar{X} - \mu)^T \Sigma^{-1} (\bar{X} - \mu) - \frac{n}{2} \ln |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} S)$$

← holds always due to $N(\mu, \Sigma)$

→ Σ^{-1} positive definite \Rightarrow ML estimate of μ is \bar{X} , reduces ML to:

$$\ln L(\mu, \Sigma) = \text{const.} - \frac{n}{2} \ln |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} S)$$

→ Σ structure \Rightarrow max L = separately maximising with respect to each block Σ

→ Skew-normal $\hat{=}$ Gaussian version

additional term

$$l = \text{constant} - \frac{n}{2} \ln |\Omega| - \frac{n}{2} \text{tr}(\Omega^{-1} V) + \sum_i (\rho_0(\Omega^T \Omega^{-1} (x_i - \eta)))$$

$$V = \frac{1}{n} \sum_i [(x_i - \eta)(x_i - \eta)^T] \approx S$$

$$\rho_0(x) = \ln(2\Phi(x))$$

Population covariance matrix since $\frac{1}{n}$ instead of $\frac{1}{n-1}$ (sample variance)

Robust Parameter estimation

→ Trimmed likelihood estimators

→ replace $\sum_{i=1}^n$ by a trimmed sum using TLE

→ remove obs. whose values would highly unlikely occur if fitted model true

→ Gaussian data: MCD and Weighted Trimmed Likelihood \Rightarrow same cov. estimators

LD form of TLE

→ TLE applied to all config's of Σ

Outlier Detection

- 1) Represent interval data as C and R
- 2) Assume $(C_i | R_i) \sim N(\mu, \Sigma)$, possibly restrict Σ
- 3) Robust parameter estimation $\rightarrow \hat{\mu}, \hat{\Sigma}$ (TLE)
- 4) robust MD's
- 5) interpret using EDA-graphics ($MD_i^2 > \chi_{d, 0.975}^2$)



Analysis of Variance

→ Comparison of means of ≥ 1 numerical variables in ≥ 2 populations

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k \quad H_1: \exists j_1 (\mu_{j_1} \neq \mu_1)$$

→ for interval valued variables

Each Y_j modelled by pair $(C_{j1}, R_j^*) \Rightarrow$ analysis of var: 2D-MANOVA of (C_{j1}, R_j^*)

→ MANOVA

→ > 1 variable considered (difference to ANOVA)

→ Compares $|T|$ of within-group matrix W to $|T|$ of global matrix T

$$W = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)(x_{ij} - \bar{x}_j)^T$$

$$T = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})(x_{ij} - \bar{x})^T$$

$x_{ij} \dots n \times k$

$\bar{x}_j, \bar{x}, \bar{x} \dots 1 \times p$

Similar to covariance formula

→ we consider p variables

→ Gaussian homoscedastic case:

→ for one variable \Rightarrow Wilk's Lambda

$$\lambda = \frac{|W|}{|T|} = \frac{|E_{act}|}{|E_{null}|}$$

$E_{null/act} \dots 2 \times 2$

→ General case

similar to end of chapter (3)

→ Likely ratio approach: maximise log-likelihood for H_0 and H_1

$$\lambda = \frac{L_{null}}{L_{act}} = \left(\frac{|E_{act}|}{|E_{null}|} \right)^{-\frac{n}{2}}$$

→ $-2 \ln \lambda$ follows asymptotically chi-square distribution

Discriminant Analysis

→ Gaussian

→ \forall config, estimate of optimum classification rule can be obtained with $\hat{\Sigma}$

→ direct generalization of classical linear and quadratic discriminant classification rules

→ Linear

$$Y = \operatorname{argmax}_g (\hat{\mu}_g^T \hat{\Sigma}^{-1} X - \frac{1}{2} \hat{\mu}_g^T \hat{\Sigma}^{-1} \hat{\mu}_g + \log(\hat{\pi}_g)) \quad g=2, \hat{\Sigma}_1 = \hat{\Sigma}_2$$

→ Quadratic

$$Y = \operatorname{argmax}_g \left(\hat{\mu}_g^T \hat{\Sigma}_g^{-1} X - \frac{1}{2} \hat{\mu}_g^T \hat{\Sigma}_g^{-1} \hat{\mu}_g + \log(\hat{\pi}_g) - \frac{1}{2} (\log(|\hat{\Sigma}_g|) + X^T \hat{\Sigma}_g^{-1} X) \right)$$

≈ LL-term

→ Skew-Normal

→ more complex due to skewness-parameter

→ alternative formulas based on configs

(1) groups differ only on μ

→ location

(2) groups differ on μ and Σ

→ location and scatter

(3) groups differ on μ and Σ and π_1 → general

→ location:

$$Y = \operatorname{argmax}_g \left(\hat{\xi}_g^T \hat{\Omega}^{-1} X - \frac{1}{2} \hat{\xi}_g^T \hat{\Omega}^{-1} \hat{\xi}_g + \log(\hat{\pi}_g) + \beta_0 (\hat{\xi}_g^T \hat{\Omega}^{-1} (X - \hat{\xi}_g)) \right)$$

Linear term

additional term

→ General

$$Y = \operatorname{argmax}_g \left(\hat{\xi}_g^T \hat{\Omega}^{-1} X - \frac{1}{2} \hat{\xi}_g^T \hat{\Omega}^{-1} \hat{\xi}_g + \log(\hat{\pi}_g) + \beta_0 (\hat{\xi}_g^T \hat{\Omega}^{-1} (X - \hat{\xi}_g)) - \frac{1}{2} (\log(|\hat{\Omega}_g|) + X^T \hat{\Omega}^{-1} X) \right)$$

GL-term

→ same between $\log p$ and $\log p$ - only difference: term β

Model based clustering

(1) finite mixture model $k \rightarrow \# \text{ clusters}$

$$f(x_i | \theta) = \sum_{l=1}^k \pi_l f_l(x_i | \theta_l)$$

data point x_i \hookrightarrow set of all parameters

$\pi_l \dots$ mixing proportions

$f_l(\dots)$ component density function

(2) ML estimation

$n \rightarrow \# \text{ data points}$

$$l(\theta; x) = \sum_{i=1}^n \ln f(x_i | \theta)$$

\rightarrow entire dataset constructed from mixture of several probability distributions, each corresponding to different cluster of model

(3) EM-Algorithm

\rightarrow avoid local optima, find max. likelihood estimate

(4) model selection

\rightarrow select model over $\#$ components k

\rightarrow Bayesian Information Criterion

$$BIC = -2 \ln(\hat{p}_1; x) + d_p \ln(n) \rightarrow \min BIC \hat{=} \text{optimal model}$$

\uparrow sample size
 $\underbrace{\hspace{10em}}_{\text{penalises worse model fit}} \quad \underbrace{\hspace{10em}}_{\text{penalises more parameters}}$
 $\hookrightarrow \# \text{ free parameters}$

(5) Model based Clustering using best model according to BIC