

①

Multivariate Random Vector : Some Basic Concepts

Multivariate random vector : $\tilde{x} = (x_1, \dots, x_k)'$

Joint distribution function
 $x_i \text{ s.r.v. on } (\mathbb{R}, \mathcal{F}_t, \mathbb{P})$

$$F_{x_1, \dots, x_k}(x_1, \dots, x_k) = P(x_1 \leq x_1, \dots, x_k \leq x_k)$$

$\forall (x_1, \dots, x_k) \in \mathbb{R}^k$

For discrete multivariate distributions

$$F_{x_1, \dots, x_k}(x_1, \dots, x_k) = \sum_{i_1 \leq x_1} \dots \sum_{i_k \leq x_k} P(x_1 = i_1, \dots, x_k = i_k)$$

$$P(x_1 = x_1, \dots, x_k = x_k) \rightarrow \text{p.m.f.}$$

$\forall (x_1, \dots, x_k) \in \text{Support of mult dist.}$

Marginal p.m.f. of x_i

$$P(x_i = x_i) = \sum_{x_2} \dots \sum_{x_k} P(x_1 = x_1, x_2 = x_2, \dots, x_k = x_k)$$

say marginal for x_i ($i \in \{1, \dots, k\}$)

Joint marginal of x_i, x_j (or any subset)

$$\sum_{\text{except } x_i, x_j} \dots \sum P(x_1 = x_1, \dots, x_k = x_k)$$

For continuous multivariate dist's

j.t. d.f.

$$F_{x_1, \dots, x_K}(x_1, \dots, x_K) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_K} f_{x_1, \dots, x_K}(u_1, \dots, u_K) \frac{K}{\prod_i} du_i$$

$$f_{x_1, \dots, x_K}(x_1, \dots, x_K) : \text{j.t p.d.f.}$$

Marginal p.d.f. of x_i :

$$f_{x_i}(x_i) = \int_{-\infty}^{x_i} \dots \int_{-\infty}^{x_i} f(u_1, \dots, u_K) \frac{K}{\prod_{j=1}^K} du_j$$

j.t marginal p.d.f. of (x_i, x_j) (or any subset)

$$f_{x_i, x_j}(x_i, x_j) = \int_{-\infty}^{x_i} \dots \int_{-\infty}^{x_j} f(u_1, \dots, u_K) \frac{K}{\prod_{l=1}^{K-2}} du_l$$

$K-2$ fold, except over x_i & x_j

Expectation vector of \underline{x}

$$\mathbb{E}(\underline{x}) = \begin{pmatrix} \mathbb{E}x_1 \\ \vdots \\ \mathbb{E}x_K \end{pmatrix} = \underline{\mu}$$

Let $\underline{y} = A \underline{x} + \underline{b}$; A : matrix of constants
 \underline{b} : vector of constants

$$\mathbb{E}(\underline{y}) = A \mathbb{E}(\underline{x}) + \underline{b}$$

$$\mathbb{E}(\underline{x}' \underline{x}) = \underline{x}' \mathbb{E}(\underline{x})$$

$\underline{x} \in \mathbb{R}^K$

3

Note:
More

Generally, for a random matrix

$$Z_{\mathbb{P}X_q} = ((Z_{ij})) \quad Z_{ij} \text{ are r.v.s. on } (\Omega, \mathcal{F}, P)$$

b.d.f. of Z is jt p.d.f of Z_{ij} 's ($i = 1 \dots p$
 $j = 1 \dots q$)

$$E(z) = ((Ez_{ij}))$$

$$f: \mathbb{Z}^m \rightarrow \mathbb{B} \oplus \mathbb{C}^{q \times n} + \mathbb{D}^{m \times n}$$

$$\text{Then } E z = B E(z) c + D$$

Covariance matrix of \tilde{x}

$$\text{Var}(\underline{x}) = E((\underline{x} - E(\underline{x}))(\underline{x} - E(\underline{x}))')$$

$$\Sigma = E(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})'$$

$(i,j)^{th}$ element of Σ is

$$\sigma_{ij} = E(x_i - \mu_i)(x_j - \mu_j)$$

$$\text{i.e } \tau_{ij} = \text{cov}(x_i, x_j)$$

$$\sigma_{ii} = E(x_i - \mu_i)^2 = \text{Var}(x_i)$$

Note : (i) Total variation in $\tilde{x} = \text{tr } \sum = \sum_{i=1}^k A_{ii}$

(ii) Generalized variance of $\tilde{X} = |\Sigma|$

(4)

Let $\underline{x} \sim k_{x_1}$ & $\underline{y} \sim p_{x_1}$ be two random vectors

$$E(\underline{x}) = \underline{\mu}_x \text{ & } E(\underline{y}) = \underline{\mu}_y$$

$$\text{Cov}(\underline{x}, \underline{y}) = E(\underline{x} - \underline{\mu}_x)(\underline{y} - \underline{\mu}_y)'$$

Note: $\text{Cov}(\underline{x}) = \Sigma$ is always a sym matrix

Result: Characterization of a covariance matrix

Any $p \times p$ real sym matrix Σ is a covariance matrix iff Σ is positive semi definite (i.e. $\underline{\alpha}' \Sigma \underline{\alpha} \geq 0 \forall \underline{\alpha} \in \mathbb{R}^p$).

Pf: Suppose Σ is cov matrix of \underline{x} , then

$$\forall \underline{\alpha} \in \mathbb{R}^p$$

$$\begin{aligned} 0 \leq V(\underline{\alpha}' \underline{x}) &= E((\underline{\alpha}' \underline{x} - \underline{\alpha}' \underline{\mu})^2) ; E(\underline{x}) = \underline{\mu} \\ &= E(\underline{\alpha}' (\underline{x} - \underline{\mu}))^2 \\ &= E(\underline{\alpha}' (\underline{x} - \underline{\mu})) (\underline{\alpha}' (\underline{x} - \underline{\mu}))' \\ &= E(\underline{\alpha}' (\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})' \underline{\alpha}) \\ &= \underline{\alpha}' \Sigma \underline{\alpha} \end{aligned}$$

$$\Rightarrow \underline{\alpha}' \Sigma \underline{\alpha} \geq 0 \quad \forall \underline{\alpha} \in \mathbb{R}^p$$

$\Rightarrow \Sigma$ is p.s.d.

Alternatively, suppose Σ is p.s.d. with rank r ($\leq p$).

Then, $\Sigma = CC'$; C is $p \times r$ matrix of rank r

Let \underline{y} be $r \times 1$ vector of indep r.v.s \Rightarrow

$$E(\underline{y}) = \underline{0} \text{ and } \text{cov}(\underline{y}) = I_r$$

Transform $\underline{y} \rightarrow \underline{x} = C\underline{y}$

$$\text{then } E(\underline{x}) = C E(\underline{y}) = \underline{0}$$

$$\begin{aligned} \text{cov}(\underline{x}) &= E(\underline{x}\underline{x}') = E(C\underline{y})(C\underline{y})' \\ &= C E(\underline{y}\underline{y}') C' = C'C = \Sigma \\ &= \underline{\underline{\Sigma}} \end{aligned}$$

$\Rightarrow \Sigma$ is a covariance matrix

Note: $\underline{x} \Rightarrow E(\underline{x}) = \underline{u}$; $\text{cov}(\underline{x}) = \Sigma$

Transform $\underline{x} \rightarrow \underline{y} = A\underline{x} + \underline{b}$

$$\begin{aligned} \text{cov}(\underline{y}) &= E((\underline{y} - E(\underline{y}))(\underline{y} - E(\underline{y}))') \\ &= E(A\underline{x} + \underline{b} - (A\underline{u} + \underline{b}))(A\underline{x} + \underline{b} - (A\underline{x} + \underline{b}))' \\ &= E(A\underline{x} - A\underline{u})(A\underline{x} - A\underline{u})' \\ &= A E(\underline{x} - \underline{u})(\underline{x} - \underline{u})' A' = A \Sigma A' \end{aligned}$$

(6)

Sp. case: Suppose $V = \text{diag}(\tau_{11}, \dots, \tau_{KK})$
with $\tau_{ii} > 0$

Take $A = (V^{1/2})^{-1}$

$$\underline{x} \rightarrow \underline{y} = A \underline{x}$$

$$\text{Cov}(\underline{y}) = (V^{1/2})^{-1} \sum (V^{1/2})^{-1} = \Sigma_y$$

$$(i,j)^{\text{th}} \text{ entry of } \Sigma_y = \frac{\tau_{ij}}{(\tau_{ii} \tau_{jj})^{1/2}} = \rho_{ij}$$

ρ_{ij} is $\text{corr}^*(x_i, x_j)$

i.e. Σ_y = corr* matrix of \underline{x}

$$\text{i.e. } \text{corr}^*(\underline{x}) = (V^{1/2})^{-1} \sum (V^{1/2})^{-1} = \rho.$$

Note: If Cov matrix of \underline{x} is not p.d., then

w.p. 1, components of \underline{x} are linearly related.

Pf. If $\Sigma \not\succ 0$, then \exists an $\alpha \in \mathbb{R}^p$ ($\alpha \neq 0$) \Rightarrow

$$0 = \alpha' \Sigma \alpha = V(\alpha' \underline{x})$$

$$\text{i.e. } P(\alpha' \underline{x} = \alpha' \underline{\mu}) = 1$$

$$\text{i.e. } P(\alpha' (\underline{x} - \underline{\mu}) = 0) = 1$$

$$\text{i.e. } \sum_i \alpha_i (x_i - \mu_i) = 0 \text{ w.p. 1 for not all } \alpha_i = 0$$

i.e. w.p 1 x_i 's are linearly related.

Note: Partitions of covariance matrix

$$\underset{k \times 1}{\tilde{X}} \text{ is } \Rightarrow E(\tilde{X}) = \underline{\mu} \text{ & } \text{cov}(\tilde{X}) = \Sigma$$

Consider the partition

$$\tilde{X} = \begin{pmatrix} \tilde{X}^{(1)} \\ \vdots \\ \tilde{X}^{(2)} \end{pmatrix}_{k-p+1}; \underline{\mu} = \begin{pmatrix} \underline{\mu}^{(1)} \\ \vdots \\ \underline{\mu}^{(2)} \end{pmatrix}$$

$$\text{cov}(\tilde{X}) = E(\tilde{X} - \underline{\mu})(\tilde{X} - \underline{\mu})'$$

$$= E \left((\tilde{X}^{(1)} - \underline{\mu}^{(1)}) (\tilde{X}^{(1)} - \underline{\mu}^{(1)})' \quad (\tilde{X}^{(1)} - \underline{\mu}^{(1)}) (\tilde{X}^{(2)} - \underline{\mu}^{(2)})' \\ \vdots \qquad \qquad \qquad \vdots \qquad \qquad \vdots \qquad \qquad \vdots \\ (\tilde{X}^{(2)} - \underline{\mu}^{(2)}) (\tilde{X}^{(1)} - \underline{\mu}^{(1)})' \quad (\tilde{X}^{(2)} - \underline{\mu}^{(2)}) (\tilde{X}^{(2)} - \underline{\mu}^{(2)})' \right)$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

$$\Sigma_{11} = \text{cov}(\tilde{X}^{(1)}) ; \Sigma_{22} = \text{cov}(\tilde{X}^{(2)})$$

$$\Sigma_{12} = \text{cov}(\tilde{X}^{(1)}, \tilde{X}^{(2)}) = \Sigma_{21}'$$

If elements of $\tilde{X}^{(1)}$ are indep of the elements of $\tilde{X}^{(2)}$
 Then $\Sigma_{12} = 0$, i.e. $\Sigma = \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$

Converse is NOT true in general.

Random sampling from multivariate pop's

Multivariate popⁿ with mean vector μ and covariance Σ

μ & Σ usually unknown

of unknowns $p + \frac{p(p+1)}{2}$ (for p dimensional pop^p)

Let $\tilde{x}_1 = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{p1} \end{pmatrix}, \dots, \tilde{x}_n = \begin{pmatrix} x_{1n} \\ x_{2n} \\ \vdots \\ x_{pn} \end{pmatrix}$ be n random samples from the popⁿ (i.i.d)

Random sample matrix

$$X_{p \times n} = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ x_{21} & \cdots & x_{2n} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pn} \end{pmatrix}$$

Observation matrix / data matrix

$$\begin{aligned} X_{p \times n} &= \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ x_{21} & \cdots & x_{2n} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pn} \end{pmatrix} \\ &= (\tilde{x}_1, \dots, \tilde{x}_n) = \begin{pmatrix} \tilde{y}_1' \\ \vdots \\ \tilde{y}_p' \end{pmatrix} \end{aligned}$$

$$\tilde{x}_j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{pj} \end{pmatrix} \quad \begin{matrix} j^{\text{th}} \text{ obsn vector} \\ j = 1(1)n \end{matrix}$$

$$\tilde{y}_i = (y_{i1}, \dots, y_{in}) \quad \begin{matrix} i^{\text{th}} \text{ variable obsns} \\ i = 1(1)p \end{matrix}$$

$\bar{\vec{x}}$: observed sample mean vector

$$\bar{\vec{x}} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix} = \frac{1}{n} \begin{pmatrix} y_1' \mathbf{1}_n \\ \vdots \\ y_p' \mathbf{1}_n \end{pmatrix} = \frac{1}{n} \begin{pmatrix} \sum_{j=1}^n x_{1j} \\ \vdots \\ \sum_{j=1}^n x_{pj} \end{pmatrix} = \frac{1}{n} \vec{x} \mathbf{1}_n$$

Observed sample variance covariance matrix:

$$S_n = \frac{1}{n} \begin{pmatrix} \sum_{j=1}^n (x_{1j} - \bar{x}_1)^2 & \cdots & \sum_{j=1}^n (x_{1j} - \bar{x}_1)(x_{pj} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n (x_{2j} - \bar{x}_2)^2 & \cdots & \sum_{j=1}^n (x_{2j} - \bar{x}_2)(x_{pj} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n (x_{pj} - \bar{x}_p)^2 & \cdots & \vdots \end{pmatrix}$$

$$\text{or } S_{n-1} = \frac{n}{n-1} S_n$$

Generalized sample variance: $|S_n|$ (or $|S_{n-1}|$)

Total sample variation: $\text{tr } S_n$ (or $\text{tr } S_{n-1}$)

Note that $(n-1)S_{n-1} = nS_n$

$$= \begin{pmatrix} \sum x_{1j}^2 & \sum x_{1j}x_{2j} & \cdots & \sum x_{1j}x_{pj} \\ \vdots & \ddots & \ddots & \vdots \\ \sum x_{2j}^2 & \cdots & \sum x_{2j}x_{pj} & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \sum x_{pj}^2 \end{pmatrix}$$

$$= n \begin{pmatrix} \bar{x}_1^2 & \bar{x}_1 \bar{x}_2 & \cdots & \bar{x}_1 \bar{x}_p \\ \bar{x}_2^2 & \cdots & \bar{x}_2 \bar{x}_p & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \bar{x}_p^2 & \cdots & \bar{x}_p \bar{x}_p & \vdots \end{pmatrix}$$

$$\begin{aligned}
 n S_n &= \mathbf{x} \mathbf{x}' - n \bar{\mathbf{x}} \bar{\mathbf{x}}' \\
 &= \mathbf{x} \mathbf{x}' - n \left(\frac{1}{n} \mathbf{x} \mathbf{1}_{-n} \right) \left(\frac{1}{n} \mathbf{x} \mathbf{1}_{-n} \right)' \\
 &= \mathbf{x} \mathbf{x}' - \frac{1}{n} \mathbf{x} \mathbf{1}_{-n} \mathbf{1}_{-n}' \mathbf{x}' \\
 &= \mathbf{x} \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n' \right) \mathbf{x}' = (n-1) S_{n-1}
 \end{aligned}$$

Alternatively,

$$\begin{aligned}
 n S_n &= \mathbf{x} \mathbf{x}' - n \bar{\mathbf{x}} \bar{\mathbf{x}}' \\
 &= (\underline{x}_1; \dots; \underline{x}_n) \begin{pmatrix} \underline{x}_1' \\ \vdots \\ \underline{x}_n' \end{pmatrix} - n \bar{\mathbf{x}} \bar{\mathbf{x}}' \\
 &= \sum \underline{x}_j \underline{x}_j' - n \bar{\mathbf{x}} \bar{\mathbf{x}}' \\
 &= \sum_{j=1}^n (\underline{x}_j - \bar{\mathbf{x}})(\underline{x}_j - \bar{\mathbf{x}})' = (n-1) S_{n-1}
 \end{aligned}$$

Note: Corresponding random matrices

$$\begin{aligned}
 n S_n &= (n-1) S_{n-1} \\
 &= \mathbf{X} \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n' \right) \mathbf{X}' \\
 &= \sum (\underline{x}_j - \bar{\mathbf{x}})(\underline{x}_j - \bar{\mathbf{x}})'
 \end{aligned}$$

Multivariate Normal distⁿ

Defⁿ: Let $\underline{\underline{x}}_{p \times 1}$ be a random vector with $E(\underline{\underline{x}}) = \underline{\underline{\mu}}$ and $Cov(\underline{\underline{x}}) = \Sigma$. We say that $\underline{\underline{x}}$ has a multivariate normal distⁿ ($\underline{\underline{x}} \sim N_p(\underline{\underline{\mu}}, \Sigma)$) iff $\forall \underline{\alpha} \in \mathbb{R}^p$ ($\underline{\alpha} \neq \underline{0}$), $\underline{\alpha}' \underline{\underline{x}}$ has univariate normal distⁿ

Some results

If $\underline{\underline{x}} \sim N_p(\underline{\underline{\mu}}, \Sigma)$, then

- (i) $\underline{\underline{x}} - \underline{\underline{\mu}} \sim N_p(\underline{0}, \Sigma)$
- (ii) each component of $\underline{\underline{x}} \sim N_1$
- (iii) Any subvector of $\underline{\underline{x}}$ has mult normal
- (iv) $A \underline{\underline{x}} + \underline{b} \sim N_q(A\underline{\underline{\mu}} + \underline{b}, A\Sigma A')$

(v) If $\Sigma > 0$, then p.d.f. of $\underline{\underline{x}}$ is

$$f_{\underline{\underline{x}}}(\underline{\underline{x}}) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2}(\underline{\underline{x}} - \underline{\underline{\mu}})' \Sigma^{-1} (\underline{\underline{x}} - \underline{\underline{\mu}}))$$

(vi) If $\Sigma \neq 0$, \nexists p.d.f. of $\underline{\underline{x}}$

(v) If $\underline{\underline{x}} = \begin{pmatrix} \underline{x}^{(1)} \\ \vdots \\ \underline{x}^{(2)} \end{pmatrix}_{p \times q \times 1}$, $\underline{\underline{\mu}} = \begin{pmatrix} \underline{\mu}^{(1)} \\ \vdots \\ \underline{\mu}^{(2)} \end{pmatrix}$, $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

$\underline{x}^{(1)}, \underline{x}^{(2)}$ indep iff $\Sigma_{12} = 0$

(12)

Random sampling from multivariate normal

$\tilde{X}_1, \dots, \tilde{X}_n$ r. s. from $N_p(\underline{\mu}, \Sigma)$, $\Sigma > 0$

Unbiased estimators:

$$\tilde{\bar{X}} = \frac{1}{n} \sum_{i=1}^n \tilde{X}_i ; E \tilde{\bar{X}} = \underline{\mu}$$

$$S_{n-1} = \frac{1}{n-1} \sum_{i=1}^n (\tilde{X}_i - \tilde{\bar{X}})(\tilde{X}_i - \tilde{\bar{X}})'$$

$$E S_{n-1} = \frac{1}{n-1} E \left(\sum \tilde{X}_i \tilde{X}_i' - n \tilde{\bar{X}} \tilde{\bar{X}}' \right)$$

$$= \frac{1}{n-1} \left(\sum_i E \tilde{X}_i \tilde{X}_i' - n E \tilde{\bar{X}} \tilde{\bar{X}}' \right)$$

$$= \frac{1}{n-1} \left(n(\Sigma + \underline{\mu} \underline{\mu}') - n(\Sigma_{1/n} + \underline{\mu} \underline{\mu}') \right)$$

$$\text{Thus } \underline{\mu} \stackrel{?}{=} \tilde{\bar{X}} \quad \& \quad \Sigma \stackrel{?}{=} S_{n-1} \quad \begin{array}{l} \text{(This is true for random)} \\ \text{sampling from any} \\ \text{multivariate pop^n} \end{array}$$

Maximum likelihood estimators

$\tilde{\bar{X}}$ is MLE of $\underline{\mu}$

$S_n = \frac{1}{n} \sum (\tilde{X}_i - \tilde{\bar{X}})(\tilde{X}_i - \tilde{\bar{X}})'$ is MLE of Σ

DTSFⁿ:

$$\tilde{\bar{X}} \sim N_p(\underline{\mu}, \Sigma_{1/n})$$

$$n S_n = (n-1) S_{n-1} \sim W_p(n-1, \Sigma)$$

→ a Wishart distⁿ of order p with
d. f. $n-1$ and associated covariance
matrix Σ

Also $\tilde{\bar{X}}$ & S_{n-1} (or S_n) are independently distributed.

Principal Component Analysis (PCA)

$$\underline{\underline{X}} = (X_1, \dots, X_p)'$$

$$\Sigma = \text{Cov}(\underline{\underline{X}}) = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{pmatrix} = ((\sigma_{ij}))$$

$$\text{Total variation in } \underline{\underline{X}} = \text{tr } \Sigma = \sum_{i=1}^p \sigma_{ii}$$

Aim of PCA

Replace $\underline{\underline{X}}$ with $\underline{\underline{Y}} = (Y_1, \dots, Y_p)'$, where Y_i 's are linear combinations of X_i 's \Rightarrow

(i) Y_i 's are uncorrelated, i.e.

$$\text{Cov}(Y_i, Y_j) = 0 \quad \forall i \neq j$$

(ii) Total variation in $\underline{\underline{X}}$

$$= \text{total variation in } \underline{\underline{Y}}$$

& (iii) total variation in $(Y_1, \dots, Y_k)'$

$$\approx \text{total variation of } \underline{\underline{X}} ; k \ll p$$

Note: (iii) may not always be possible

Major uses of PCA

- (i) Data dimension reduction ~~& visualization~~
- (ii) Outlier mining
- (iii) Detection of clusters in the multidimensional data cloud
- (iv) Projection & visualization
- (v) ranking of multidimensional data
- (vi) Checking for multi-variate normality

Defⁿ: Principal components

PCs are uncorrelated linear combinations y_1, \dots, y_p (derived from x_1, \dots, x_p) whose variances are in decreasing order, with y_1 explaining the maximum variability, y_2 explaining the second maximum variability and so on.

i.e. the 1st PC, y_1 , is the linear combination $\underline{l}' \underline{x}$ that maximizes $V(\underline{l}' \underline{x})$ subject to $\underline{l}' \underline{l} = 1$

The 2nd PC, y_2 , is the linear combination $\underline{l}_2' \underline{x}$ that maximizes $V(\underline{l}' \underline{x})$ subject to $\underline{l}' \underline{l} = 1$ and $\text{cov}(\underline{l}_1' \underline{x}, \underline{l}' \underline{x}) = 0$

The ith PC, y_i , is the linear combination $\underline{l}_i' \underline{x}$ that maximizes $V(\underline{l}' \underline{x})$ subject to $\underline{l}' \underline{l} = 1$ and $\text{cov}(\underline{l}_k' \underline{x}, \underline{l}' \underline{x}) = 0 \quad \forall k < i$.

Derivation of PCs

Let Σ be the covariance matrix associated with the random vector \underline{x} . The eigen value - eigen vector pairs of Σ are $(\lambda_1, \underline{e}_1), \dots, (\lambda_p, \underline{e}_p)$, where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

The ith PC is given by

$$y_i = \underline{e}_i' \underline{x} \quad i = 1, \dots, p$$

With $V(y_i) = \lambda_i$ & $\text{cov}(y_i, y_j) = 0 \quad \forall i \neq j$

Pf: Consider any $\underline{\lambda}' \underline{x}$, $\underline{\lambda} \in \mathbb{R}^p$; $\Rightarrow \underline{\lambda}' \underline{\lambda} = 1$

$$V(\underline{\lambda}' \underline{x}) = \underline{\lambda}' V(\underline{x}) \underline{\lambda}$$

$$= \underline{\lambda}' \sum \underline{\lambda}$$

Consider the eigenvalue-eigen vector decomposition of Σ as

$$\Sigma = P D_\lambda P'$$

$$D_\lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$$

$$\Rightarrow V(\underline{\lambda}' \underline{x}) = \underline{\lambda}' P D_\lambda P' \underline{\lambda}$$

$$= \underline{\beta}' D_\lambda \underline{\beta}$$

$$= \sum_{i=1}^p \beta_i^2 \lambda_i$$

where $\underline{\beta} = P' \underline{\lambda} \in \mathbb{R}^p$

$$\underline{\lambda} = P \underline{\beta}$$

$$\text{Now, } \underline{\lambda}' \underline{\lambda} = 1 \Rightarrow \underline{\beta}' P' P \underline{\beta} = \underline{\beta}' \underline{\beta} = 1$$

Thus

$$\max_{\underline{\lambda} \rightarrow \underline{\lambda}' \underline{\lambda} = 1} V(\underline{\lambda}' \underline{x}) = \max_{\underline{\beta} \rightarrow \underline{\beta}' \underline{\beta} = 1} \underline{\beta}' D_\lambda \underline{\beta} = \max_{\underline{\beta} \rightarrow \underline{\beta}' \underline{\beta} = 1} \sum_{i=1}^p \lambda_i \beta_i^2$$

Realize that,

$$\sum_{i=1}^p \lambda_i \beta_i^2 \leq \lambda_1 \sum \beta_i^2 = \lambda_1$$

$$\text{i.e. } V(\underline{\lambda}' \underline{x}) \leq \lambda_1$$

and

$$V(\underline{e}_1' \underline{x}) = \underline{e}_1' \sum \underline{e}_i$$

$$= \underline{e}_1' P D_\lambda P' \underline{e}_1$$

$$= \underline{e}_1' (\underline{e}_1 : \dots : \underline{e}_p) D_\lambda \begin{pmatrix} \underline{e}_1' \\ \vdots \\ \underline{e}_p' \end{pmatrix} \underline{e}_1$$

$$\text{i.e. } V(\underline{e}' \underline{x}) = (1, 0, \dots, 0) D_\lambda \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$= \lambda_1 = \max_{\underline{l} \geq \underline{l}' \underline{l} = 1} V(\underline{l}' \underline{x})$$

$\Rightarrow y_1 = \underline{e}' \underline{x}$ is the 1st PC

For the second PC, we need $y_2 = \underline{l}_2' \underline{x} \Rightarrow y_2$ is uncorrelated with y_1 and $V(\underline{l}_2' \underline{x})$ to be the maximum under this constraint.

$$\text{Now, } \text{cov}(\underline{l}' \underline{x}, \underline{e}' \underline{x})$$

$$= E(\underline{l}'(\underline{x} - \underline{\mu})) (\underline{e}'(\underline{x} - \underline{\mu}))'$$

$$= \underline{l}' \sum \underline{e}_i$$

$$= \underline{l}' \left(\sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i' \right) \underline{e}_1$$

$$= \underline{l}' (\lambda_1 \underline{e}_1) = \lambda_1 \underline{l}' \underline{e}_1$$

Thus $\text{cov}(\underline{l}' \underline{x}, \underline{e}' \underline{x}) = 0 \Leftrightarrow \underline{l} \perp \underline{e}_1$

Thus $\max_{\underline{l} \geq \underline{l}' \underline{l} = 1} V(\underline{l}' \underline{x})$ subject to $\text{cov}(\underline{l}' \underline{x}, y_1) = 0$ is

equivalent to $\max_{\substack{\underline{l} \geq \underline{l} \perp \underline{e}_1 \\ \underline{l}' \underline{l} = 1}} V(\underline{l}' \underline{x})$

$$V(\underline{l}' \underline{x}) = \underline{l}' P D_\lambda P'$$

$$= \underline{l}' (\underline{e}_1 : \dots : \underline{e}_p) \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} \underline{e}_1' \\ \vdots \\ \underline{e}_p' \end{pmatrix}$$

$\underline{\lambda} \perp \underline{e}_1$

$$\begin{aligned}\underline{\lambda}' P D \lambda P' \underline{\lambda} &= (0, b_1, \dots, b_p) \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} 0 \\ b_1 \\ \vdots \\ b_p \end{pmatrix} \\ &= \sum_{i=2}^p b_i^2 \lambda_i\end{aligned}$$

$$\left. \begin{aligned} &\leq \lambda_2 \\ \text{i.e. } \underline{\lambda} &= P \underline{b} \\ \underline{\lambda}' \underline{\lambda} &= 1 \Rightarrow \underline{b}' \underline{b} = 1 \quad \text{i.e. } \sum_2^p b_i^2 = 1 \end{aligned} \right)$$

$$\begin{aligned}\text{And } V(\underline{e}_2' \underline{x}) &= \underline{e}_2' \sum \underline{e}_2 \\ &= \underline{e}_2' (\underline{e}_1 : \dots : \underline{e}_p) D \lambda \begin{pmatrix} \underline{e}_1' \\ \vdots \\ \underline{e}_p' \end{pmatrix} \underline{e}_2 \\ &= (0, 1, 0, \dots, 0) D \lambda \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \underline{e}_2 \\ &= \lambda_2 \\ &= \max_{\substack{\underline{\lambda} \geq \underline{e}_1 \\ \underline{\lambda}' \underline{\lambda} = 1}} V(\underline{\lambda}' \underline{x})\end{aligned}$$

$\Rightarrow y_2 = \underline{e}_2' \underline{x}$ is the 2nd PC

For the $(k+1)^{\text{th}}$ PC ($k+1 \leq p$), we need to

maximize $V(\underline{\lambda}' \underline{x}) \Rightarrow \text{Cor}(\underline{\lambda}' \underline{x}, \underline{y}_i) = 0 \forall i \leq k$

i.e. maximize $V(\underline{\lambda}' \underline{x}) \Rightarrow \underline{\lambda} \perp \underline{e}_i \forall i \leq k$

$\nexists \underline{\lambda} \perp e_1, \dots, e_k$

$$V(\underline{\lambda}' \underline{x}) = \underline{\lambda}' (e_1 : \dots : e_k : e_{k+1} : \dots : e_p) D_\lambda \begin{pmatrix} e'_1 \\ \vdots \\ e'_k \\ e'_{k+1} \\ \vdots \\ e'_p \end{pmatrix} \underline{\lambda}$$

$$= (0, \dots, 0, c_{k+1}, \dots, c_p) D_\lambda \begin{pmatrix} 0 \\ \vdots \\ 0 \\ c_{k+1} \\ \vdots \\ c_p \end{pmatrix}$$

$$= \sum_{i=k+1}^p c_i \lambda_i$$

$$\underline{c} = P' \underline{\lambda} \quad \text{with} \quad \underline{c} = (0, \dots, 0, c_{k+1}, \dots, c_p)'$$

$$\underline{\lambda}' \underline{\lambda} = 1 \Rightarrow \underline{c}' \underline{c} = 1 \quad \text{i.e.} \quad \sum_{k+1}^p c_i^2 = 1$$

Thus $\nexists \underline{\lambda} \ni e_1, \dots, e_k$

$$V(\underline{\lambda}' \underline{x}) = \sum_{k+1}^p c_i^2 \lambda_i \leq \lambda_{k+1}$$

$$\text{and } \lambda_{k+1} = V(e_{k+1}' \underline{x}) (= e_{k+1}' P D_\lambda P' e_{k+1})$$

$$\Rightarrow V(e_{k+1}' \underline{x}) = \lambda_{k+1} = \max_{\underline{\lambda} \ni \underline{\lambda} \perp e_1, \dots, e_k} V(\underline{\lambda}' \underline{x})$$

$$\wedge \underline{\lambda}' \underline{\lambda} = 1$$

$$\Rightarrow \cancel{\lambda}_{k+1} \in P C$$

$$y_{k+1} = e_{k+1}' \underline{x}$$

Note 1 : y_1, \dots, y_p are \Rightarrow

$$\text{Cov}(\underline{y}) = D_\lambda$$

$$y_i = e'_i \underline{x}$$

$$\text{Cov}(y_i, y_j) = \text{Cov}(e'_i \underline{x}, e'_j \underline{x})$$

$$= e'_i \sum e'_j$$

$$= e'_i \left(\sum_{k=1}^p \lambda_k e_k e'_k \right) e'_j$$

$$= \begin{cases} \lambda_i, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$$

$$\Rightarrow \text{Cov}(\underline{y}) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \lambda_p \end{pmatrix} = \Sigma_y$$

Note 2 : Total variation in $\underline{x} = \text{Total variation in } \underline{y}$

$$\text{Total variation in } \underline{x} = \text{tr } \Sigma$$

$$= \text{tr}(P D_\lambda P')$$

$$= \text{tr}(D_\lambda P' P)$$

$$= \text{tr } D_\lambda = \sum_{i=1}^p \lambda_i$$

$$\text{Total variation in } \underline{y} = \text{tr}(\text{Cov}(\underline{y}))$$

$$= \text{tr } D_\lambda = \sum_{i=1}^p \lambda_i$$

Note 3 : If $\underline{x} \sim N_p(\underline{\mu}, \Sigma)$, then

$$(i) \quad y_i = e_i' \underline{x} \quad (\forall i=1(1)p) \sim N_1$$

$$(ii) \quad \underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix} = \begin{pmatrix} e_1' \underline{x} \\ \vdots \\ e_p' \underline{x} \end{pmatrix} = \begin{pmatrix} e_1' \\ \vdots \\ e_p' \end{pmatrix} \underline{x} \sim$$

$$\text{i.e. } \underline{y} = P' \underline{x} \sim N_p(P' \underline{\mu}, P' \Sigma P)$$

$$P' \Sigma P = P'(P D_A P') P = D_A$$

(iii) y_1, \dots, y_p are independent N_1

($\text{cov} = 0$ for joint normal \Rightarrow independence)

Note 4: $\text{Cov}(x_i, y_j) = ?$

$$y_j = e_j' \underline{x}$$

$$\underline{y} = P' \underline{x}$$

$$\underline{x} = P \underline{y} = (e_1 : \dots : e_p) \underline{y}$$

$$= \begin{pmatrix} e_{11} & \dots & e_{1p} \\ e_{21} & \dots & e_{2p} \\ \vdots & & \vdots \\ e_{p1} & \dots & e_{pp} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix}$$

$$\text{i.e. } x_i = \sum_{k=1}^p e_{ik} y_k$$

$$\text{Cov}(x_i, y_j) = \text{Cov}\left(\sum_{k=1}^p e_{ik} y_k, y_j\right) = e_{ij} \lambda_j$$

$$\text{Corr}^n(x_i, y_j) = e_{ij} \sqrt{\frac{\lambda_j}{\lambda_i}}$$

e_{ij} is the i^{th} element (row) of e_j

$$V(X_i) = \sigma_{ii} \quad & V(Y_j) = \lambda_j$$

$$\Rightarrow \text{Corrl}^*(X_i, Y_j) = \frac{\rho_{ij} \lambda_j}{(\sigma_{ii} \lambda_j)^{1/2}} = \rho_{ij} \sqrt{\frac{\lambda_j}{\sigma_{ii}}}$$

Note 5: Principal Components derived from correlation matrix

Consider the standardized variables

$$X_i \rightarrow Z_i = \frac{X_i - \bar{x}_i}{\sqrt{\sigma_{ii}}}, \quad i = 1, \dots, p; \quad \sigma_{ii} > 0$$

define

$$V = \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$$

$$\underline{z} = (z_1, \dots, z_p)' = V^{-1/2} (\underline{x} - \underline{\mu})$$

$$E \underline{z} = \underline{0}; \quad \text{cov}(\underline{z}) = V^{-1/2} \Sigma V^{-1/2} = I_p$$

Let (r_i, f_i) ($i = 1, \dots, p$), $r_1 \geq r_2 \dots \geq r_p \geq 0$ be the eigen value - eigen vector (o.n.) pairs of P matrix, then

(a) i.e. PC derived from P is

$$Y_i = f_i' \underline{z} = f_i' V^{-1/2} (\underline{x} - \underline{\mu})$$

$$(b) \text{ total variation in } \underline{y} = \sum_{i=1}^p V(Y_i) = \sum_{i=1}^p V(z_i) = p$$

↑
total variation in \underline{z}

$$(c) \rho_{z_i, Y_j} = f_{ij} \sqrt{r_j} \quad ; \quad i, j = 1, \dots, p$$

Remark: (i) Proportion of total variation in \tilde{X} explained by the K^{th} PC (based on Σ) is

$$\lambda_K / \sum_i \lambda_i$$

(ii) Proportion of total variation in \tilde{X} explained by the first K PCs (based on Σ) is

$$\sum_{i=1}^K \lambda_i / \sum_{i=1}^p \lambda_i$$

(iii) If $\sum_{i=1}^K \lambda_i \approx \sum_{i=1}^p \lambda_i$ for some small K , then

It is possible to have a meaningful data dimension reduction as it is enough in such a situation to consider only (y_1, \dots, y_K) .

Ideally: $K \ll p$!

Remark: $r_{ij} \propto P_{x_i, y_j}$ and hence the magnitude of r_{ij} indicates how important is x_i to y_j , this helps in interpreting y_j .

Remark: In case the units of variables are different or the variables have widely varying ranges, we should obtain PCs from standardized variables i.e. Work with P instead of Σ to get PCs.

Examples

$$(1) \quad \Sigma = \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

Solving $\det(\Sigma - \lambda I) = 0$ gives $\lambda_1 = 5.83$, $\lambda_2 = 2.0$, $\lambda_3 = 0.17$

$$\underline{e}_1' = (0.383, -0.924, 0)$$

$$\underline{e}_2' = (0, 0, 1)$$

$$\underline{e}_3' = (0.924, 0.383, 0)$$

PCs

$$Y_1 = .383 X_1 - .924 X_2$$

$$Y_2 = X_3$$

$$Y_3 = .924 X_1 + .383 X_2$$

1st PC explains $\frac{5.83}{8} \times 100 = 73\%$ of total variation

1st 2 PCs explain $\frac{5.83 + 2.0}{8} \times 100 = 98\%$ of total variation

$$\rho_{X_1, Y_1} = .383 \sqrt{\frac{5.83}{1}} \approx .925 \quad \left(\rho_{X_i, Y_j} = e_{ij} \sqrt{\frac{\lambda_j}{\sigma_{ii}}} \right)$$

$$\rho_{X_2, Y_1} = (-.924) \sqrt{\frac{5.83}{5}} \approx -.998$$

Example 2

$$\Sigma = \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix} \quad x_1, x_2 \rightarrow \text{widely varying in spread}$$

$$\lambda_1 = 100.16, \lambda_2 = 0.84$$

$$e_1' = (.04, .999); e_2' = (.999, -0.04)$$

$$y_1 = .04x_1 + .999x_2; y_2 = .999x_1 - .04x_2$$

y_1 explains $\frac{100.16}{101} \times 100 \approx 99\%$ of total variation

$$\rho_{x_1, y_1} = 0.04 \sqrt{\frac{100.16}{1}} \approx 0.4 \quad \left. \begin{array}{l} x_2 \text{ is much more imp} \\ \text{to } y_1 \text{ than } x_1 \text{ is to } y_1 \end{array} \right\}$$

$$\rho_{x_2, y_1} = .999 \sqrt{\frac{100.16}{100}} \approx .99 \quad \left. \begin{array}{l} \text{Var}(x_2) \text{ dominating!} \\ \dots \end{array} \right\}$$

Consider standardized variables

$$z_1 = \frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \quad z_2 = \frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}}$$

$$\Sigma_z = P_x = \begin{pmatrix} 1 & .4 \\ .4 & 1 \end{pmatrix}$$

$$z_1 = 1.4; z_2 = 0.6$$

$$f_1' = (.707, .707); f_2' = (.707, -.707)$$

$$y_1 = .707 z_1 + .707 z_2$$

$$y_2 = .707 z_1 - .707 z_2$$

$$\rho_{z_1, y_1} = .707 \sqrt{\frac{1.4}{1}} = \rho_{z_2, y_1}$$

z_1 & z_2 are equally important

y_1 explains $\frac{1.4}{2} \times 100 = 70\%$ of total variation

Principal Components from data matrix \mathbf{x}

observed data matrix

$$\mathbf{x}_{p \times n} = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ x_{21} & \cdots & x_{2n} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pn} \end{pmatrix} \quad \begin{array}{l} x_{11}, \dots, x_{nn} \text{ } n \text{ obsns} \\ \text{from multivariate} \\ \text{pop}^n \text{ with mean} \\ \text{vector } \mathbf{\bar{x}} \text{ and cov mat } \Sigma \end{array}$$

Sample mean vector : $\hat{\mathbf{x}} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix} = \frac{1}{n} \mathbf{x} \mathbf{1}_n'$

Sample variance covariance matrix \mathbf{x} :

$$\begin{aligned} S_{n-1} &= \frac{1}{n-1} \mathbf{x} \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n' \right) \mathbf{x}' = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j - \hat{\mathbf{x}})(\mathbf{x}_j - \hat{\mathbf{x}})' \\ &= \frac{1}{n-1} \begin{pmatrix} \sum (\mathbf{x}_{1j} - \bar{x}_1)^2 & \sum (\mathbf{x}_{1j} - \bar{x}_1)(\mathbf{x}_{2j} - \bar{x}_2) & \cdots & \sum (\mathbf{x}_{1j} - \bar{x}_1)(\mathbf{x}_{pj} - \bar{x}_p) \\ \vdots & \vdots & \ddots & \vdots \\ \sum (\mathbf{x}_{pj} - \bar{x}_p)^2 & \sum (\mathbf{x}_{pj} - \bar{x}_p)(\mathbf{x}_{1j} - \bar{x}_1) & \cdots & \sum (\mathbf{x}_{pj} - \bar{x}_p)(\mathbf{x}_{pj} - \bar{x}_p) \end{pmatrix} \end{aligned}$$

From the observed data matrix \mathbf{x} , our objective is to construct uncorrelated linear combinations of the measured characteristics that account for as much of the sample variation as possible.

Suppose $\underline{l}' \hat{\mathbf{x}}$ is sample linear combination

We have $\underline{l}' \hat{\mathbf{x}}_1, \dots, \underline{l}' \hat{\mathbf{x}}_n$: n transformed obsns

Sample mean of $(\underline{l}' \hat{\mathbf{x}}_1, \dots, \underline{l}' \hat{\mathbf{x}}_n)$:

$$\begin{aligned} &\frac{1}{n} \sum_{j=1}^n \underline{l}' \hat{\mathbf{x}}_j \\ &= \frac{1}{n} \sum_{j=1}^n (l_1 x_{1j} + l_2 x_{2j} + \cdots + l_p x_{pj}) \\ &= l_1 \bar{x}_1 + l_2 \bar{x}_2 + \cdots + l_p \bar{x}_p = \underline{l}' \hat{\mathbf{x}} \end{aligned}$$

Sample variance of $(\underline{\lambda}' \underline{x}_1, \dots, \underline{\lambda}' \underline{x}_n)$:

$$\begin{aligned}
 & \frac{1}{n-1} \left((\underline{\lambda}' \underline{x}_1 - \underline{\lambda}' \bar{\underline{x}})^2 + \dots + (\underline{\lambda}' \underline{x}_n - \underline{\lambda}' \bar{\underline{x}})^2 \right) \\
 &= \frac{1}{n-1} \left((\underline{\lambda}' (\underline{x}_1 - \bar{\underline{x}}))^2 + \dots + (\underline{\lambda}' (\underline{x}_n - \bar{\underline{x}}))^2 \right) \\
 &= \frac{1}{n-1} \left((\underline{\lambda}' (\underline{x}_1 - \bar{\underline{x}})) (\underline{\lambda}' (\underline{x}_1 - \bar{\underline{x}}))' + \dots + (\underline{\lambda}' (\underline{x}_n - \bar{\underline{x}})) (\underline{\lambda}' (\underline{x}_n - \bar{\underline{x}}))' \right) \\
 &= \frac{1}{n-1} \left(\underline{\lambda}' (\underline{x}_1 - \bar{\underline{x}}) (\underline{x}_1 - \bar{\underline{x}})' \underline{\lambda} + \dots + \underline{\lambda}' (\underline{x}_n - \bar{\underline{x}}) (\underline{x}_n - \bar{\underline{x}})' \underline{\lambda} \right) \\
 &= \underline{\lambda}' \left(\frac{1}{n-1} \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}}) (\underline{x}_j - \bar{\underline{x}})' \right) \underline{\lambda} \\
 &= \underline{\lambda}' S_{n-1} \underline{\lambda}
 \end{aligned}$$

Further, suppose $\underline{\lambda}_i' \underline{x}$ & $\underline{\lambda}_j' \underline{x}$ be 2 linear combinations,
sample covariance betⁿ n pairs

$$\begin{aligned}
 & (\underline{\lambda}_i' \underline{x}_1, \underline{\lambda}_j' \underline{x}_1), \dots, (\underline{\lambda}_i' \underline{x}_n, \underline{\lambda}_j' \underline{x}_n) \\
 & \frac{1}{n-1} \left((\underline{\lambda}_i' \underline{x}_1 - \underline{\lambda}_i' \bar{\underline{x}}) (\underline{\lambda}_j' \underline{x}_1 - \underline{\lambda}_j' \bar{\underline{x}})' + \dots + (\underline{\lambda}_i' \underline{x}_n - \underline{\lambda}_i' \bar{\underline{x}}) (\underline{\lambda}_j' \underline{x}_n - \underline{\lambda}_j' \bar{\underline{x}})' \right) \\
 &= \frac{1}{n-1} \sum_{k=1}^n (\underline{\lambda}_i' (\underline{x}_k - \bar{\underline{x}}) (\underline{x}_k - \bar{\underline{x}})' \underline{\lambda}_j') \\
 &= \underline{\lambda}_i' \left(\frac{1}{n-1} \sum_{k=1}^n (\underline{x}_k - \bar{\underline{x}}) (\underline{x}_k - \bar{\underline{x}})' \right) \underline{\lambda}_j \\
 &= \underline{\lambda}_i' S_{n-1} \underline{\lambda}_j
 \end{aligned}$$

Defⁿ: Sample Principal Components

The sample PCs are uncorrelated linear combinations of observation vectors, have sample variances in decreasing order. The 1st sample PC is the lin comb

$\underline{\hat{b}}^T \underline{x}_j$ which maximises the sample variance of $\underline{\hat{b}}^T \underline{x}_j$ subject to $\underline{\hat{b}}^T \underline{\hat{b}} = 1$ (i.e. sample variance of $(\underline{\hat{b}}^T \underline{x}_1, \dots, \underline{\hat{b}}^T \underline{x}_n)$)

$$= \max_{\underline{\hat{b}} \rightarrow \underline{\hat{b}}^T \underline{\hat{b}} = 1} \underline{\hat{b}}^T S_{n-1} \underline{\hat{b}}$$

The 2nd sample PC is the linear comb $\underline{\hat{b}}_2^T \underline{x}_j \rightarrow$ it

maximizes sample variance of $\underline{\hat{b}}_2^T \underline{x}_j \rightarrow \underline{\hat{b}}^T \underline{\hat{b}} = 1$ and

zero sample covariance for the pairs $(\underline{\hat{b}}_1^T \underline{x}_j, \underline{\hat{b}}_2^T \underline{x}_j); j=1, \dots, n$

The ith sample PC is the linear comb $\underline{\hat{b}}_i^T \underline{x}_j \rightarrow$ it

maximizes the sample variance of $\underline{\hat{b}}_i^T \underline{x}_j$ and having

zero sample covariance for all pairs $(\underline{\hat{b}}_i^T \underline{x}_j, \underline{\hat{b}}_k^T \underline{x}_j)$

$\forall k < i : j=1, \dots, n$

S_{n-1} : Sample variance covariance matrix

$(\hat{\lambda}_1, \hat{e}_1), \dots, (\hat{\lambda}_p, \hat{e}_p)$: eigenvalue - eigenvector pairs of S_{n-1}

$\max_{\underline{\hat{b}} \rightarrow \underline{\hat{b}}^T \underline{\hat{b}} = 1} \underline{\hat{b}}^T S_{n-1} \underline{\hat{b}} = \hat{\lambda}_1$ = sample variance of

$$(\hat{e}_1^T \underline{x}_1, \dots, \hat{e}_1^T \underline{x}_n)$$

$$= \hat{e}_1^T S_{n-1} \hat{e}_1$$

$$= \hat{e}_1^T \left(\sum_j \hat{\lambda}_j \hat{e}_j \hat{e}_j^T \right) \hat{e}_1 = \hat{\lambda}_1$$

1st Sample PC is $\hat{e}_1' \tilde{x}_j = \hat{y}_{1(j)}$

For 2nd sample PC we need to find

$$\max_{\substack{\hat{e} \geq \hat{e}_1 \\ \hat{e}' \hat{e} = 1}} \hat{e}' S_{n-1} \hat{e} = \hat{\lambda}_2 = \text{sample variance of } (\hat{e}_2' \tilde{x}_1, \dots, \hat{e}_2' \tilde{x}_n)$$

2nd Sample PC is $\hat{e}_2' \tilde{x}_j = \hat{y}_{2(j)}$

k^{th} ($k \leq p$) Sample PC is $\hat{e}_k' \tilde{x}_j = \hat{y}_{k(j)}$

Note i) Sample variance of $\hat{y}_k = \hat{\lambda}_k$

ii) Sample covariance of \hat{y}_i & $\hat{y}_k = 0 \quad \forall i \neq k$

(Sample covariance of $(\hat{e}_i' \tilde{x}_j, \hat{e}_k' \tilde{x}_j)$ pairs

$$= \hat{e}_i' S_{n-1} \hat{e}_k$$

$$= \hat{e}_i' \left(\sum_{j=1}^n \hat{\lambda}_j \hat{e}_j \hat{e}_j' \right) \hat{e}_k = 0$$

iii) total sample variation = $\sum_{i=1}^p \hat{\lambda}_i$

= total sample variation in sample PCs

$$= \sum_{i=1}^p \hat{\lambda}_i$$

$$\text{iv) } r_{x_i, \hat{y}_k} = \hat{e}_{ik} \sqrt{\frac{\hat{\lambda}_k}{\hat{\lambda}_{kk}}}$$

Uses of PCA

Sample PCs $\hat{y}_i = \hat{e}_i' \tilde{x}; i=1(1)n$

$$\tilde{x}_1 \rightarrow (\hat{y}_1^{(1)}, \hat{y}_2^{(1)}, \dots, \hat{y}_p^{(1)})$$

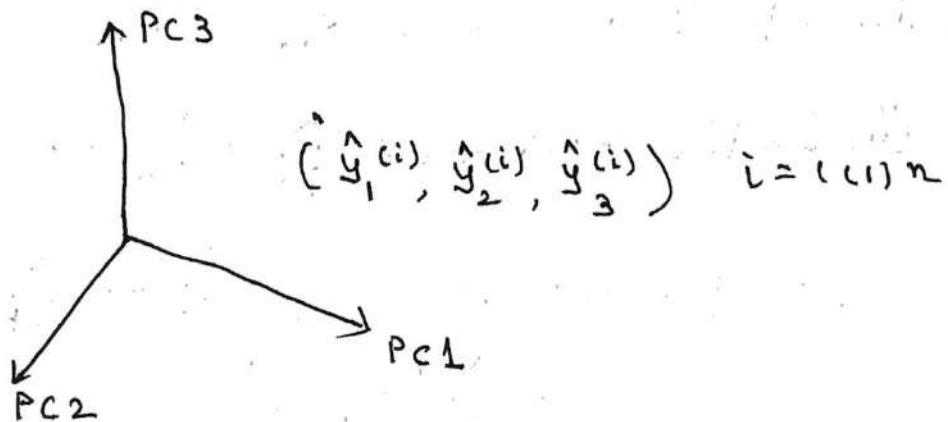
$$\tilde{x}_2 \rightarrow (\hat{y}_1^{(2)}, \hat{y}_2^{(2)}, \dots, \hat{y}_p^{(2)})$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$$

$$\tilde{x}_n \rightarrow (\hat{y}_1^{(n)}, \hat{y}_2^{(n)}, \dots, \hat{y}_p^{(n)})$$

(i) Data projection & visualization

Projection onto PC plane



- Meaningful when $(\hat{y}_1, \hat{y}_2, \hat{y}_3)$ collectively explains a "significant" amount of total sample variation
- a similar 2d plot on $(PC1, PC2)$ plane, provided the first 2 PC collectively captures "significant" proportion of total sample variation

(ii) outlier detection : from 2d/3d plot (PC plane)

(iii) cluster detection : from above plots

(iv) Data dimension reduction

through $(\hat{y}_1, \dots, \hat{y}_k)$ for $k < p$

(v) Ranking of cases

$$\left. \begin{array}{l} \underline{x}_1 \rightarrow \hat{y}_1^{(1)} \\ \underline{x}_2 \rightarrow \hat{y}_1^{(2)} \\ \vdots \\ \underline{x}_n \rightarrow \hat{y}_1^{(n)} \end{array} \right\} \begin{array}{l} \text{rank based on } \hat{y}_1^{(i)} \\ \text{coeffs of variables (original) play} \\ \text{a crucial role} \end{array}$$

(vi) Variable clustering

$$r_{\hat{y}_k, x_i} = \hat{\epsilon}_{ik} \sqrt{\frac{\lambda_k}{\lambda_{ii}}}$$

- Represent each of the original p variables by $(\hat{r}_{\hat{y}_1, x_i}, \hat{r}_{\hat{y}_2, x_i}, \hat{r}_{\hat{y}_3, x_i})^T$ $i = 1, \dots, p$
- Identify variable clusters from this plot

(vii) Checking for multivariate normality

- Check for univariate normality of each of the p PCs
- For i th PC use $(\hat{y}_1^{(1)}, \dots, \hat{y}_1^{(n)})$ to check for univariate normality
- If any of the p checks for univariate normality fails, N_p assumption will be rejected

Cluster Analysis

Grouping of objects (with multidimensional feature vectors) based on self similarities ("closeness").

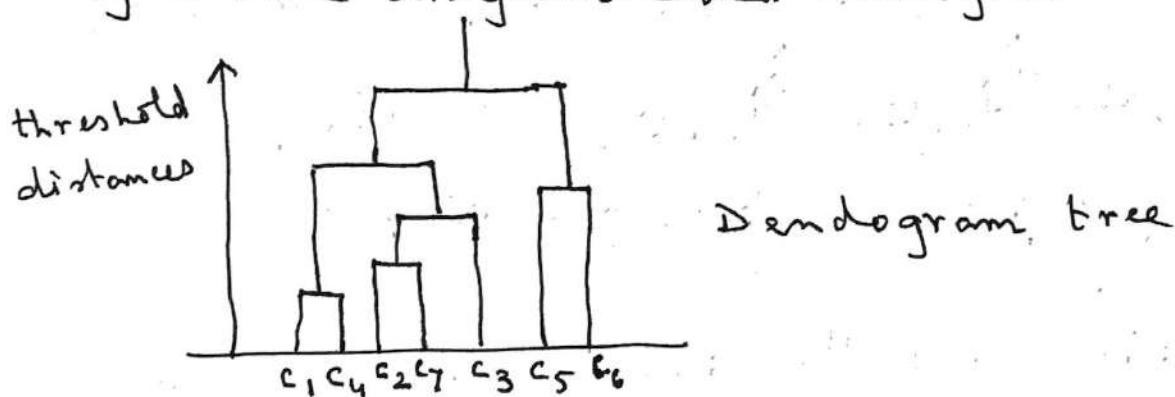
Two approaches

- Hierarchical clustering
- non-hierarchical clustering

Hierarchical clustering

In HCA the observation vectors (cases) are grouped together on the basis of mutual distances in such a manner so as to ensure hierarchy in cluster formation.

HCA output is usually visualized through a hierarchical tree, which is a nested set of partitions represented by a tree diagram called dendrogram.



Characteristics

- (i) Sectioning a tree at a particular level produces a partition into g disjoint groups
- (ii) If two groups are chosen from different partitions (i.e. the results of partitioning at different levels) then either the groups are disjoint or one group contains the other

- (iii) a numerical value is associated with each partition up/down the tree where branches join. This is a measure of distance between 2 merged clusters.
- (iv) We get different trees corresponding to different distance measures ("bet" cases).

Two algorithms for HCA

- Agglomerative clustering algorithm
- Divisive clustering algorithm

Agglomerative clustering algorithm

- Operates with successive merger of objects
- Starts with n clusters, each containing a single data point.
- At each stage merges the 2 most similar groups to form a new cluster, thus reducing the number of clusters by 1.
- The process of merger of group of objects continue till all subgroups are fused together to form a single cluster

Divisive clustering algorithm

- Operates by successively splitting groups of objects
- Starts with a single cluster having all the n objects
- Initial group split into 2 clusters \Rightarrow the distance bet["] the split groups ~~are~~ is maximum, i.e. the subgroups are as far as possible.

- Process of splitting continues till there are n clusters, each having a single object (leaf node)
- DCA is computationally inefficient.

Note: Result from agglomerative (or divisive) clustering algorithm displayed through dendrogram.

Steps in agglomerative hierarchical algorithm

Step I : n clusters, each having a single object and an $n \times n$ symmetric matrix of distances ($D = ((d_{ij}))$)

Step II : Search the distance matrix for the nearest (most similar) pair of clusters. Let the distance bet" the most similar clusters, say U and V , be d_{UV} .

Step III : Merge $U \& V$ to form cluster (U, V) , at a merger (fusion) level d_{UV} .

Update the distance matrix by

(i) deleting the rows & col^ms corresponding to clusters U & V

and (ii) adding a row & col^m giving the distance bet" (U, V) and the remaining clusters (existing ones other than $U \& V$)..

Step IV : Repeat (II) & (III) $n-1$ times, recording the identity of the clusters that are merged and the merger levels (i.e the

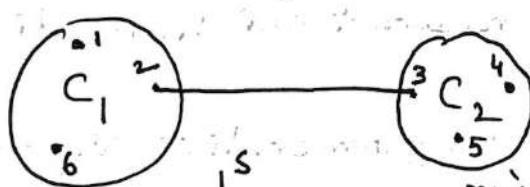
distances at which mergers take place).

Step 2 : Construct the dendrogram tree from the information on mergers and level of mergers.

Remark: step ii requires calculation of distance betⁿ clusters (for distance matrix updation).

The following are commonly used distance measures betⁿ clusters of objects

(i) Single link distance (minimum distance or nearest neighbor distance)



$$d^s_{C_1, C_2} = \min_{\substack{i \in C_1 \\ j \in C_2}} d_{ij}$$

(ii) Complete link distance (maximum distance or farthest distance)

$$d^c_{C_1, C_2} = \max_{\substack{i \in C_1 \\ j \in C_2}} d_{ij}$$

(iii) Average link distance

N_{C_1} : # of objects in C_1

N_{C_2} : # of objects in C_2

$$d^A_{C_1, C_2} = \frac{1}{N_{C_1} N_{C_2}} \sum_{\substack{i \in C_1 \\ j \in C_2}} d_{ij}$$

(iv) Centroid distance : Distance betⁿ cluster means

(v) Median distance : Distance betⁿ cluster medians

Remark : Single linkage agglomerative algorithm

In step ii, (distance matrix update stage), we compute the distance betⁿ (U, V) and any existing cluster W as

$$d_{(U,V),W} = \min(d_{UW}, d_{VW})$$

↑ ↑
present in prev round

d_{UW} & d_{VW} are nearest neighbor distances

Remark : Complete linkage agglomerative algorithm

In step ii, compute distance betⁿ (U, V) and any existing W as

$$d_{(U,W),W} = \max(d_{UW}, d_{VW})$$

Remark : Average linkage agglomerative algorithm

In step ii, compute distance betⁿ (U, V) and W computed as

$$d_{(U,V),W} = \frac{\sum_{i \in (U,V)} \sum_{k \in W} d_{ik}}{N_{(U,V)} N_W}$$

d_{ik} : distance betⁿ object i in the cluster (U, V)
and object k in cluster W

$N_{(U,V)}$: # of objects in (U, V) & N_W : # of objects in W

Example: Single linkage clustering

Five objects c_1, c_2, c_3, c_4, c_5

Distance matrix

$$D = \begin{pmatrix} c_1 & 0 & & & \\ c_2 & 9 & 0 & & \\ c_3 & 3 & 7 & 0 & \\ c_4 & 6 & 5 & 9 & 0 \\ c_5 & 11 & 10 & 2 & 8 & 0 \end{pmatrix}$$

$c_1 \ c_2 \ c_3 \ c_4 \ c_5$

$$d_{c_3 c_5} = \min d_{ij} = 2$$

\Rightarrow Merge c_3 & c_5 to form cluster (c_3, c_5) at merger level 2 - total 4 clusters

Row col^m deletion containing c_3 & c_5

$$\begin{pmatrix} 0 & & & & \\ 9 & 0 & & & \\ c_3 & 3 & 7 & 9 & - \\ c_5 & 6 & 5 & 0 & 1 \\ 11 & 10 & 2 & 8 & 0 \end{pmatrix}$$

Computation of $d_{(c_3, c_5), c_i}$ for $i = 1, 2, 4$

$$\begin{aligned} d_{(c_3, c_5) c_1} &= \min(d_{c_1 c_3}, d_{c_1 c_5}) \\ &= \min(3, 11) = 3 \end{aligned}$$

$$d_{(c_3, c_5) c_2} = \min(d_{c_2 c_3}, d_{c_2 c_5}) = \min(7, 10) = 7$$

$$d_{(c_3, c_5) c_4} = \min(d_{c_3 c_4}, d_{c_4 c_5}) = \min(9, 8) = 8.$$

Updated distance matrix x

(C_3, C_5)	0				
C_1	3	0			
C_2	7	9	0		
C_4	8	6	5	0	

$$d_{(C_3, C_5)C_1} \text{ is min}$$

Merge (C_3, C_5) with C_1 to form cluster

(C_1, C_3, C_5) at merger level 3 - total 3 clusters

Row column deletion step

(C_3, C_5)	0	-	-	-	-
C_1	3	0	-	-	-
C_2	7	9	0		
C_4	8	6	5	0	

Calculate distance bet " (C_3, C_5, C_1) & C_2 and
 (C_3, C_5, C_1) & C_4

$$\begin{aligned} d_{(C_1, C_3, C_5)C_2} &= \min(d_{C_1 C_2}, d_{(C_3, C_5)C_2}) \\ &= \min(9, 7) = 7 \end{aligned}$$

$$\begin{aligned} d_{(C_1, C_3, C_5)C_4} &= \min(d_{C_1 C_4}, d_{(C_3, C_5)C_4}) \\ &= \min(6, 8) = 6 \end{aligned}$$

Update distance matrix

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{ccc} 0 & & \\ & & \\ & & \end{array} \right) \\ c_2 & \left(\begin{array}{cc} 7 & 0 \\ & \end{array} \right) \\ c_4 & \left(\begin{array}{ccc} 6 & 5 & 0 \\ & \boxed{5} & \\ & & \end{array} \right) \end{matrix}$$

d_{c_2, c_4} is min

Merge c_2 & c_4 to form cluster (c_2, c_4) at merger level 5 - total 2 clusters (c_1, c_3, c_5) & (c_2, c_4)

Row column deletion

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc|c} 0 & & \\ & & | \\ c_2 & 7 & 0 & \\ & & \hline & & \end{array} \right) \\ c_4 & \left(\begin{array}{cc|c} 6 & 5 & 0 \\ & & \hline & & \end{array} \right) \end{matrix}$$

Distance matrix update

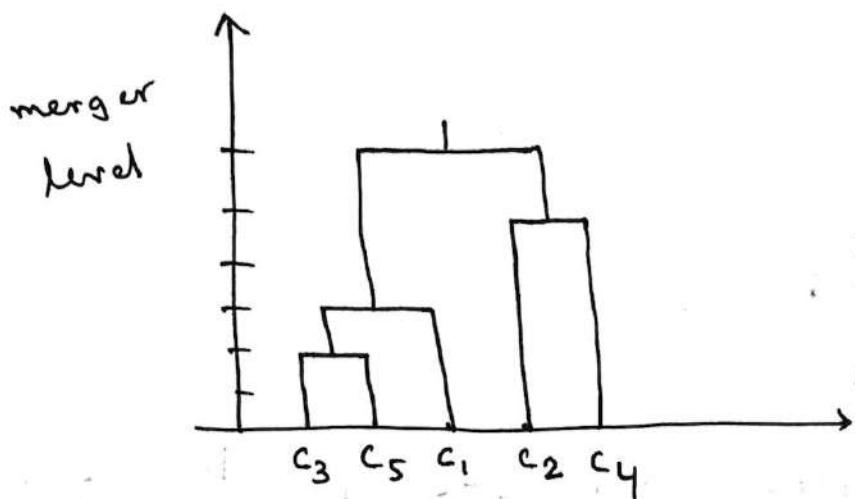
$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc} 0 & \\ & \end{array} \right) \\ (c_2, c_4) & \left(\begin{array}{cc} ? & 0 \end{array} \right) \end{matrix}$$

$$\begin{aligned} d_{(c_1, c_3, c_5)(c_2, c_4)} &= \min(d_{(c_1, c_3, c_5)c_2}, d_{(c_1, c_3, c_5)c_4}) \\ &= \min(7, 6) = 6 \end{aligned}$$

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc} 0 & \\ & \end{array} \right) \\ (c_2, c_4) & \left(\begin{array}{cc} 6 & 0 \end{array} \right) \end{matrix}$$

\Rightarrow Merge (c_1, c_3, c_5) & (c_2, c_4) at merger level 6
- single cluster

Dendrogram



Measures of dissimilarity

$$\mathbf{x}_{p \times n} = (x_1 : \dots : x_n)$$

$\mathbf{x} \rightarrow D$: dissimilarity matrix
 $p \times p$

$$D = ((d_{rs})) \quad \Rightarrow$$

$$(i) d_{rs} \geq 0 \quad \forall r, s$$

$$(ii) d_{rr} = 0 \quad \forall r$$

$$(iii) d_{rs} = d_{sr} \quad \forall r, s$$

Remark: (i) "symmetry cond" may not always be satisfied

$$D \rightarrow (D + D')/2$$

(ii) If in addition to the cond's (i)-(iii),
 d_{rs} also satisfies triangle inequality

$$d_{rs} \leq d_{rt} + d_{ts} \quad \forall r, s, t$$

then the dissimilarity measure is a metric and
we use the term "distance"

Some commonly used dissimilarity measures for numeric values

x_i, x_j : $p \times 1$ feature vector

$$x_i = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{pi} \end{pmatrix}$$

(a) Euclidean distance

$$d_e = \left(\sum_{i=1}^p (x_{1i} - x_{2i})^2 \right)^{1/2}$$

(b) Absolute value distance or City-block distance

$$d_{cb} = \sum_{i=1}^p |x_{1i} - x_{2i}|$$

(c) Chebyshov distance

$$d_{ch} = \max_i |x_{li} - x_{rj}|$$

(d) Pearson distance

$$d_p = \left(\sum_{i=1}^p (x_{li} - x_{rj})^2 / s_i^2 \right)^{1/2}$$

s_i^2 : sample variance for i^{th} variable

(e) Quadratic distance

$$d_q = \left((\underline{x}_i - \underline{x}_j)' Q (\underline{x}_i - \underline{x}_j) \right)^{1/2}$$

Q is p.d.

e.g. $Q = S^{-1}$: S = within group covariance matrix

case of feature vector with ordinal values

Realizable values are ordered set.

e.g. Academic grades: A, B, C, D, E, F

opinion ~~Opinion~~: strongly disagree, disagree, neutral, agree, strongly agree

One approach is to transform M (say) ordinal values

to

$$\frac{j - \frac{1}{2}}{M}; j = 1(1)M$$

in the prescribed order of their original values.

Treat the transformed values as quantitative numeric values.

Case of nominal (categorical, unordered) / ordinal values

Transform to binary variables set.

Nominal variable with K states

→ represent as K binary variables

Nominal variable at m^{th} state

$$(0, \dots, \underset{m^{\text{th}}}{\uparrow} 1, \dots, 0)$$

Sly for ordinal

Similarity measures for binary variables

Let \tilde{x} & \tilde{y} be 2 vectors of binary variables.

Define

a : # of occurrences of $x_i = 1$ and $y_i = 1$

b : # of occurrences of $x_i = 0$ & $y_i = 1$

c : $x_i = 1$ & $y_i = 0$

d : $x_i = 0$ & $y_i = 0$

$a + b + c + d = p$, # of variables.

Some similarity measures:

(a) Simple matching coefficient:

$$\delta_{sm} = \frac{a+d}{a+b+c+d}$$

(b) Russell & Rao

$$\delta_{RR} = \frac{a}{a+b+c+d}$$

(c) Jaccard

$$\delta_J = \frac{a}{a+b+c}$$

(d) Czekanowski:

$$\delta_{CZ} = \frac{2a}{2a+b+c}$$

Note:

dissimilarity measure

$$\text{e.g. } d_{xy} = 1 - \delta_{xy} \text{ for } S_m$$

Mixed type: dimensions have both numerical / categorical

Transform all to binary.

Example

	Height*	Weight [#]	Eye color	Hair color	Handedness	Gender
Ind 1	68	140	green	blond	R	F
Ind 2	73	185	brown	brown	R	M
Ind 3	67	165	blue	blond	R	M
Ind 4	64	120	brown	brown	R	F
Ind 5	76	210	brown	brown	L	M

* : in inch

: in lb

Define, $x_{1i} = \begin{cases} 1, & ht \geq 72 \text{ in} \\ 0, & \text{o/w.} \end{cases}$ $x_{2i} = \begin{cases} 1, & wt \geq 150 \text{ lb} \\ 0, & \text{o/w} \end{cases}$

$$x_{3i} = \begin{cases} 1, & \text{brown eyes} \\ 0, & \text{o/w} \end{cases} \quad x_{4i} = \begin{cases} 1, & \text{blond hair} \\ 0, & \text{o/w} \end{cases}$$

$$x_{5i} = \begin{cases} 1, & R \\ 0, & L \end{cases} \quad x_{6i} = \begin{cases} 1, & \text{female} \\ 0, & \text{male} \end{cases}$$

	x_1	x_2	x_3	x_4	x_5	x_6
Ind 1	0	0	0	1	1	1
Ind 2	1	1	1	0	1	0

 \rightarrow

	Ind 1		
Ind 2	1	0	
	1 (a)	2 (b)	
0	3 (c)	0 (d)	

$$s_{dm}^{(1,2)} = \frac{1}{6}$$

Combine $\neq (i, j)$ $i \neq j$ pairs and get similarity matrix

$$S = \begin{pmatrix} 1 & 1 & & & & \\ 2 & \frac{1}{6} & 1 & & & \\ 3 & - & - & 1 & & \\ 4 & - & - & - & 1 & \\ 5 & - & - & - & - & 1 \\ 6 & - & - & - & - & - \end{pmatrix}$$

Distance betⁿ groups of objects

(A) Methods based on feature vectors

Already discussed (single linkage, complete linkage, average linkage, centroid distance, median distance)

(B) Methods based on probabilistic distance.

Distance measures use complete information about the structure of groups (classes) provided by the class conditional densities. The distance measure, say J , satisfies

$$(i) \quad J = 0 \quad \text{if} \quad p(x|\pi_1) = p(x|\pi_2)$$

$$(ii) \quad J \geq 0$$

(iii) J attains the maximum value when the classes are disjoint

Commonly used probabilistic distances

(I) Bhattacharya distance

$$J_B = -\log \left(\int \{p(x|\pi_1) p(x|\pi_2)\}^{1/2} dx \right)$$

(II) Chernoff distance

$$J_C = -\log \left(\int \{p(x|\pi_1) [p(x|\pi_2)]^{\lambda}\}^{1-\lambda} dx \right)$$

$\lambda \in [0, 1]$

(III) Divergence distance

$$J_D = \int (p(x|\pi_1) - p(x|\pi_2)) \log \left(\frac{p(x|\pi_1)}{p(x|\pi_2)} \right) dx$$

(iv) Patrick - Fischer distance

$$J_{PF} = \left(\int \left\{ p_1 p(\underline{x} | \pi_1) - p_2 p(\underline{x} | \pi_2) \right\}^2 d\underline{x} \right)^{1/2}$$

p_1 & p_2 are prior probabilities of the classes.

Remark: Probabilistic distances are computed under specific distributional setups.

Example : Gtr I objects - $N_p(\underline{\mu}_1, \Sigma_1)$ $\Sigma_i > 0$

Gtr II objects - $N_p(\underline{\mu}_2, \Sigma_2)$

Divergence distance calculations

$$\begin{aligned} J_D &= \int \left((2\pi)^{-p/2} |\Sigma_1|^{-1/2} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1)\right) \right. \\ &\quad \left. - (2\pi)^{-p/2} |\Sigma_2|^{-1/2} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2)\right) \right) \\ &\quad \left(\log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} + \log \left(\exp\left(-\frac{1}{2}((\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2))\right) \right) \right) d\underline{x} \\ &= \log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \int (p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2)) d\underline{x} \\ &\quad + \int \left(-\frac{1}{2}(\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) \right) \left(p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2) \right) d\underline{x} \\ &\quad + \int \left(\frac{1}{2}(\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \right) \left(p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2) \right) d\underline{x} \end{aligned}$$

$$= (1 - 1) \log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}}$$

$$+ \left(-\frac{1}{2}\right) \int (\underline{x}' \Sigma_1^{-1} \underline{x} + \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - 2 \underline{\mu}_1' \Sigma_1^{-1} \underline{x}) \\ (\Pr(\underline{x} | \pi_1) - \Pr(\underline{x} | \pi_2)) d\underline{x}$$

$$+ \left(\frac{1}{2}\right) \int (\underline{x}' \Sigma_2^{-1} \underline{x} + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - 2 \underline{\mu}_2' \Sigma_2^{-1} \underline{x}) \\ (\Pr(\underline{x} | \pi_1) - \Pr(\underline{x} | \pi_2)) d\underline{x}$$

$$= -\frac{1}{2} \left(\left\{ E_{\pi_1} (\underline{x}' \Sigma_1^{-1} \underline{x}) + \cancel{\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1} - 2 \underline{\mu}_1' \Sigma_1^{-1} E_{\pi_1} (\underline{x}) \right\} \right. \\ \left. - \left\{ E_{\pi_2} (\underline{x}' \Sigma_1^{-1} \underline{x}) + \cancel{\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1} - 2 \underline{\mu}_1' \Sigma_1^{-1} E_{\pi_2} (\underline{x}) \right\} \right)$$

$$+ \frac{1}{2} \left(\left\{ E_{\pi_1} (\underline{x}' \Sigma_2^{-1} \underline{x}) + \cancel{\underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2} - 2 \underline{\mu}_2' \Sigma_2^{-1} E_{\pi_1} (\underline{x}) \right\} \right.$$

$$\left. - \left\{ E_{\pi_2} (\underline{x}' \Sigma_2^{-1} \underline{x}) + \cancel{\underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2} - 2 \underline{\mu}_2' \Sigma_2^{-1} E_{\pi_2} (\underline{x}) \right\} \right)$$

— *

Notice that

$$(i) E_{\pi_i} (\underline{x}) = \underline{\mu}_i ; i = 1, 2$$

and we need to calculate terms of type

$$E_{\pi_i} (\underline{x}' \Sigma_j^{-1} \underline{x}) ; i, j = 1, 2 \\ = ?$$

If $\underline{x} \sim N_p(\underline{\mu}, \Sigma)$, then

$$\begin{aligned}
 E(\underline{x}' \Sigma^{-1} \underline{x}) &= E(\text{tr}(\underline{x}' \Sigma^{-1} \underline{x})) \\
 &= E(\text{tr}(\Sigma^{-1} \underline{x} \underline{x}')) \\
 &= \text{tr}(E(\Sigma^{-1} \underline{x} \underline{x}')) \\
 &= \text{tr}(\Sigma^{-1} E(\underline{x} \underline{x}')) \\
 &= \text{tr}(\Sigma^{-1} (\Sigma + \underline{\mu} \underline{\mu}')) \\
 &= \text{tr} \Sigma^{-1} \Sigma + \text{tr}(\Sigma^{-1} \underline{\mu} \underline{\mu}') \\
 &= \text{tr} I_p + \text{tr}(\underline{\mu}' \Sigma^{-1} \underline{\mu}) \\
 &= \text{tr} I_p + \underline{\mu}' \Sigma^{-1} \underline{\mu} \quad - (**)
 \end{aligned}$$

In general,

$$\text{From } (*), E(\underline{x}' A \underline{x}) = \text{tr} A \Sigma + \underline{\mu}' A \underline{\mu}$$

$$\begin{aligned}
 J_D &= (\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 \\
 &\quad - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2) \\
 &\quad - \frac{1}{2} \left(E_{\pi_1}(\underline{x}' \Sigma_1^{-1} \underline{x}) - E_{\pi_2}(\underline{x}' \Sigma_1^{-1} \underline{x}) \right. \\
 &\quad \left. - E_{\pi_1}(\underline{x}' \Sigma_2^{-1} \underline{x}) + E_{\pi_2}(\underline{x}' \Sigma_2^{-1} \underline{x}) \right)
 \end{aligned}$$

$$\begin{aligned}
&= \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right) \\
&\quad - \frac{1}{2} \left(\left(\text{tr} \Sigma_1^{-1} \Sigma_1 + \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 \right) - \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 \right) \right. \\
&\quad \left. - \left(\text{tr} \Sigma_2^{-1} \Sigma_1 + \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right) + \left(\text{tr} \Sigma_2^{-1} \Sigma_2 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 \right) \right) \\
&= \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - \frac{1}{2} \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \frac{1}{2} \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 - \frac{1}{2} \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 \right) \\
&\quad + \left(-\frac{1}{2} \text{tr} I_p + \frac{1}{2} \text{tr} \Sigma_1^{-1} \Sigma_2 + \frac{1}{2} \text{tr} \Sigma_2^{-1} \Sigma_1 - \frac{1}{2} \text{tr} I_p \right) \\
&= \left(\frac{1}{2} \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \frac{1}{2} \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 \right) \\
&\quad + \frac{1}{2} \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \text{tr} \Sigma_2^{-1} \Sigma_1 - 2 \text{tr} I_p \right) \\
&= \frac{1}{2} \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 + \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - 2 \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - 2 \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right) \\
&\quad + \frac{1}{2} \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \text{tr} \Sigma_2^{-1} \Sigma_1 - 2 \text{tr} I_p \right) \\
&= \frac{1}{2} (\underline{\mu}_2 - \underline{\mu}_1)' (\Sigma_1^{-1} + \Sigma_2^{-1}) (\underline{\mu}_2 - \underline{\mu}_1) \\
&\quad + \frac{1}{2} \text{tr} (\Sigma_1^{-1} \Sigma_2 + \Sigma_2^{-1} \Sigma_1 - 2 I_p).
\end{aligned}$$

Note: If $\Sigma_1 = \Sigma_2$, then

$$J_D = (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma^{-1} (\underline{\mu}_2 - \underline{\mu}_1)$$

→ square of Mahalanobis distance

Note: For $N_p(\underline{\mu}_1, \Sigma_1) \& N_p(\underline{\mu}_2, \Sigma_2)$ setup.

Chernoff distance

$$J_C = \frac{1}{2} \lambda(1-\lambda) (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma_{\lambda}^{-1} (\underline{\mu}_2 - \underline{\mu}_1) + \frac{1}{2} \log \left(\frac{|\Sigma_{\lambda}|}{|\Sigma_1|^{1-\lambda} |\Sigma_2|^{\lambda}} \right)$$

$$\Sigma_{\lambda} = (1-\lambda) \Sigma_1 + \lambda \Sigma_2 ; \lambda \in [0, 1]$$

Note: If $\Sigma_1 = \Sigma_2$ i.e. $N_p(\underline{\mu}_1, \Sigma) \& N_p(\underline{\mu}_2, \Sigma)$.

$$J_D = 8 J_B = (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma^{-1} (\underline{\mu}_2 - \underline{\mu}_1)$$

Note: For $N_p(\underline{\mu}_1, \Sigma_1) \& N_p(\underline{\mu}_2, \Sigma_2)$ setup

$$J_{PF} = (2\pi)^{-p/2} \left(|2\Sigma_1|^{-1/2} + |2\Sigma_2|^{-1/2} - 2 |\Sigma_1 + \Sigma_2|^{-1/2} \right)$$

Patrick-Fischer

$$\exp \left(-\frac{1}{2} (\underline{\mu}_2 - \underline{\mu}_1)' (\Sigma_1 + \Sigma_2)^{-1} (\underline{\mu}_2 - \underline{\mu}_1)' \right)$$

Non-hierarchical clustering methods

Features

- no distance/dissimilarity matrix calculations
- no hierarchy of clusters
- starts from either
 - (i) a random initial partition of items into groups
 - or (ii) an initial set of randomly selected seed points, which will form the nuclei of clusters.

K-means clustering or iterative relocation method

K-means method is an iterative algorithm that assigns each item to the cluster having the nearest centroid (mean).

Steps for K-means algorithm

- (I) Partition the items into K initial clusters.
- (II) Reassign items to the cluster whose centroid is nearest (in Euclidean sense mostly). Recalculate the centroids for the cluster receiving the new item and for the cluster losing that item.
- (III) Repeat step (II) till no more reassignment can take place.

Note: Rather than starting with a partition of items into K initial groups in step (I), we can also specify

K initial seed points (centroids) and then proceed to step (ii) after a walk through the data.

Remark: Let's try to understand the logic behind K-means method !!

Let there be N objects to be put into K clusters.

Suppose $c(i)=k$ denote that object i in cluster k

$$k=1(1)K; i=1(1)N$$

A natural "loss function" (criterion) for clustering :

$$W(c) = \frac{1}{2} \sum_{k=1}^{K\#} \sum_{c(i)=k} \sum_{c(i')=k} d_{ii'} - (1)$$

$$\text{where, } d_{ii'} = d(x_i, x_{i'})$$

$W(c)$ is within cluster point scatter, a measure of compactness of clusters. It characterizes the extent to which observations assigned, under a partition, to same clusters tend to be close to one another.

Further, consider the total point scatter

$$T = \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N d_{ii'}$$

$$\text{i.e. } T = \frac{1}{2} \sum_{k=1}^{K\#} \sum_{c(i)=k} \left(\sum_{c(i')=k} d_{ii'} + \sum_{c(i') \neq k} d_{ii'} \right)$$

$$\text{i.e. } T = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i')=k} d_{ii'} + \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i') \neq k} d_{ii'}$$

$$\text{i.e. } T = W(c) + B(c)$$

$$B(c) = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i') \neq k} d_{ii'} \quad \text{is between-}$$

cluster point scatter. A measure of separation of clusters; large when different clusters are far apart.

Note:

$W(c)$ & $B(c)$ depends on cluster assignment
 T is indep of cluster assignment

Note: It is natural to maximize $B(c)$ or equivalently minimize $W(c)$ over all \Rightarrow possible cluster assignments

Note: The combinatorial optimization is not feasible even for moderate N & K .

The number of distinct cases (assignments) is

$$S(N, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^N$$

$$\text{e.g } S(10, 4) = 34,105$$

$$S(19, 4) \approx 10^{10}$$

Feasible strategies are based on iterative greedy descent - K-means algorithm is one such algorithm.

K-means algorithm

Iterative descent clustering algorithm

Let

$$d(\tilde{x}_i, \tilde{x}_{i'}) = d_{ii'} = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|\tilde{x}_i - \tilde{x}_{i'}\|^2$$

Within-cluster point scatter

$$\begin{aligned} W(c) &= \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i')=k} \|\tilde{x}_i - \tilde{x}_{i'}\|^2 \\ &= \sum_{k=1}^K N_k \sum_{c(i)=k} \|\tilde{x}_i - \bar{\tilde{x}}_k\|^2 \end{aligned}$$

where $N_k = \sum_{i=1}^N I(c(i)=k)$

$\bar{\tilde{x}}_k$: mean of of k^{th} cluster

Objective is thus to get $c^* \ni$

$$c^* = \underset{c}{\operatorname{argmin}} \sum_{k=1}^K N_k \sum_{c(i)=k} \|\tilde{x}_i - \bar{\tilde{x}}_k\|^2$$

and we can obtain c^* by enlarged optimization problem

$$\min_{c, \tilde{m}_1, \dots, \tilde{m}_K} \sum_{k=1}^K N_k \sum_{c(i)=k} \|\tilde{x}_i - \tilde{m}_k\|^2 \quad (2)$$

as for any set of S observations

$$\tilde{\tilde{x}}_S = \underset{m}{\operatorname{argmin}} \sum_{i=1}^S \|\tilde{x}_i - \tilde{m}\|^2$$

K-mean algorithm is an alternating optimization procedure for minimization of (2)

K-means clustering algorithm

Step I. For a given cluster assignment C , the total within cluster variability:

$$\sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \bar{m}_k\|^2$$

is minimized w.r.t. $\bar{m}_1, \dots, \bar{m}_K$ by taking means of the currently assigned clusters

Step II: Given a current set of cluster means,

$$\sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \bar{m}_k\|^2 \quad \text{i.e. } \bar{m}_k = \bar{x}_k$$

is minimized by assigning each obsn to the closest cluster mean

$$\text{i.e. } C(i) = \underset{1 \leq k \leq K}{\operatorname{argmin}} \|x_i - \bar{m}_k\|^2$$

Step I & II are iterated until the assignments do not change.

Remark: The result may represent a suboptimal local minimum.

One should start the algorithm with many initial random choices and choose the solⁿ having minimum objective function value.

Example : K-means

Cases	Variables	
	x_1	x_2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

$$K = 2$$

Step I : Arbitrary partition (A, B) (C, D)

(cluster centroids)

Cluster	\bar{x}_1	\bar{x}_2
(A, B)	2	2
(C, D)	-1	-2

Step II : Compute distances from cluster centroids and
reassign each item to the nearest group

(If an item is moved from the initial configuration,
the cluster centroid must be updated before we
can proceed further)

$$\underline{A} \quad d^2(A, (A, B)) = 3^2 + 1^2 = 10$$

$$d^2(A, (C, D)) = 6^2 + 5^2 = 61$$

$\Rightarrow A$ is closer to (A, B) than (C, D)

\Rightarrow No reassignment required.

$$\underline{B} \quad d^2(B, (A, B)) = 3^2 + 1^2 = 10$$

$$d^2(B, (C, D)) = 0^2 + 3^2 = 9$$

\Rightarrow Reassign B to (c, D) to form (B, c, D)

Cluster centroid updation

Cluster	Coordinates of centroid	
	\bar{x}_1	\bar{x}_2
A	5	3
(B, c, D)	-1	-1

$$d^2(A, A) = 0 ; d^2(A, (B, c, D)) = 52$$

$$d^2(B, A) = 40 ; d^2(B, (B, c, D)) = 4$$

$$d^2(c, A) = 41, d^2(c, (B, c, D)) = 5$$

$$d^2(D, A) = 89, d^2(D, (B, c, D)) = 5$$

\Rightarrow no reassignment is necessary and relocation iteration stops.

(A), (B, c, D) is the final clustering.

Comparing different cluster partitions

Objective is to have a criterion function for comparing different cluster partitions of the data and hence to choose the best cluster partition.

Let the N data points be given by $\tilde{x}_1, \dots, \tilde{x}_N$.

The sample variance covariance matrix $\hat{\Sigma}$ is given by

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\tilde{x}_i - \tilde{m})(\tilde{x}_i - \tilde{m})'$$

$$\text{where } \tilde{m} = \frac{1}{N} \sum_{i=1}^N \tilde{x}_i = \bar{\tilde{x}}$$

Let there be K clusters and define

$$z_{ji} = \begin{cases} 1, & \text{if } \tilde{x}_i \in \text{cluster } j \\ 0, & \text{o/w} \end{cases}$$

$$\tilde{m}_j = \frac{1}{N_j} \sum_{i=1}^N z_{ji} \tilde{x}_i (= \bar{\tilde{x}}_j) \quad \text{mean of cluster } j$$

$$n_j = \sum_{i=1}^N z_{ji} \quad \# \text{ of } \tilde{x}_i \text{'s in cluster } j$$

The within-cluster sum of squares and cross product (SSCP) scatter matrix is

$$S_W = \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{x}_i - \tilde{m}_j)(\tilde{x}_i - \tilde{m}_j)'$$

(The pooled within cluster scatter matrix over K clusters)

Note that

$$\begin{aligned}
 \sum &= \frac{1}{N} \sum_{i=1}^N (\underline{x}_i - \underline{m})(\underline{x}_i - \underline{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m})(\underline{x}_i - \underline{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\overline{\underline{x}_i - \underline{m}_j} + \overline{\underline{m}_j - \underline{m}})(\underline{x}_i - \underline{m}_j + \underline{m}_j - \underline{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \left((\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)' \right. \\
 &\quad \left. + (\underline{m}_j - \underline{m})(\underline{m}_j - \underline{m})' \right. \\
 &\quad \left. + (\cancel{\underline{x}_i - \underline{m}_j})(\cancel{\underline{m}_j - \underline{m}}) \right. \\
 &\quad \left. + (\cancel{\underline{m}_j - \underline{m}})(\cancel{\underline{x}_i - \underline{m}_j})' \right)
 \end{aligned}$$

i.e. $\sum = \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)'$

$$+ \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{m}_j - \underline{m})(\underline{m}_j - \underline{m})'$$

(e.g. $\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{m}_j - \underline{m})(\underline{x}_i - \underline{m}_j)'$

$$= \frac{1}{N} \sum_{j=1}^K (\underline{m}_j - \underline{m}) \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m}_j)'$$

$$= \frac{1}{N} \sum_{j=1}^K (\underline{m}_j - \underline{m}) (n_j \underline{m}_j - n_j \underline{m}_j)' = 0$$

So $\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m}_j)(\underline{m}_j - \underline{m})' = 0$

$$\text{i.e. } \sum = S_W + \frac{1}{N} \sum_{j=1}^K n_j (\underline{m}_j - \underline{m})(\underline{m}_j - \underline{m})'$$

$$= S_W + S_B$$

S_B : between-cluster sum of squares and cross product
 (SSCP) A scatter matrix
 (indicates the scatter of the cluster means
 about total grand mean)

Popular optimum clustering criteria are based on
 univariate (scalar) functions of above matrices, S_W, S_B, \sum

e.g. (a) Minimization of $\text{tr}(S_W)$

$$\begin{aligned}\text{tr}(S_W) &= \text{tr} \left(\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)' \right) \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \text{tr}((\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)') \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \text{tr}((\underline{x}_i - \underline{m}_j)'(\underline{x}_i - \underline{m}_j)) \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \|\underline{x}_i - \underline{m}_j\|^2\end{aligned}$$

$$\text{i.e. } \text{tr}(S_W) = \frac{1}{N} \sum_{j=1}^K S_j$$

$$S_j = \sum_{i=1}^n z_{ji} \|\underline{x}_i - \underline{m}_j\|^2$$

The within-cluster sum of squares for j^{th} cluster

Thus minimization of $\text{tr}(S_W)$ is same as
minimization of total within cluster sum of squares
about the K centroids (equiv to max $\text{tr} S_B$)

$$(b) \text{ Min } \frac{\|S_W\|}{\|\hat{\Sigma}\|}$$

function of $\hat{\Sigma}$ or Σ has component of rows & columns added together

$$(c) \text{ Min } \text{tr}(\hat{\Sigma}^{-1} S_W)$$

function of $\hat{\Sigma}$ will be constant with $\hat{\Sigma}$ & S_W

$$(d) \text{ Max } \text{tr}(S_W^{-1} S_B)$$

constant with S_W & S_B as S_W^{-1} is proportional to S_W & S_B is multiple of identity matrix

thus, $S_W^{-1} S_B$ is constant w.r.t. S_W & S_B

thus, $\text{tr}(S_W^{-1} S_B)$ is constant w.r.t. S_W & S_B

$\text{tr}(S_W^{-1} S_B) = \text{tr}(S_B S_W^{-1})$ (commutative property)

$\text{tr}(S_B S_W^{-1}) = \text{tr}(S_B)$ (multiple of identity matrix)

$\text{tr}(S_B) = \text{tr}(S_B^2)$ (square of matrix)

$\text{tr}(S_B^2) = \sum_{i=1}^n \lambda_i^2$ (eigenvalues)

$\lambda_i^2 = \lambda_i \lambda_i$ (square of eigenvalue)

$\lambda_i \lambda_i = \lambda_i^2$ (square of eigenvalue)

$\text{tr}(S_B^2) = \sum_{i=1}^n \lambda_i^2$ (square of eigenvalues)

Density estimation

Objective is to find

$$f(\mathbf{x}) - \text{density of feature vector}$$

Approaches

- (i) non-parametric approach
- (ii) parametric approach

Non-parametric density estimation methods

(A) Histogram method

For one-dimensional data :

$$\hat{p}(x) = \frac{n_j}{\left(\sum_{j=1}^N n_j \right) dx}$$

Where,
 n_j : # of samples in the histogram cell of width dx

the contains the pt x

N : # of cells in the histogram

dx : width of the cell

For multidimensional feature vector :

$$\hat{p}(\mathbf{x}) = \frac{n_j}{\left(\sum_{j=1}^N n_j \right) dv}$$

dv : volume of the j^{th} bin

(B) K-nearest neighbor method

Note that if X is a.r.v. (cont type),

$$P(x \leq X \leq x + \Delta x) = F(x + \Delta x) - F(x)$$

$$\lim_{\Delta x \rightarrow 0} \frac{F(x + \Delta x) - F(x)}{\Delta x} = p(x) \quad \text{p.d.f}$$

$$\text{i.e. } F(x + \Delta x) - F(x) = P(x \leq X \leq x + \Delta x)$$

$$\approx p(x) \Delta x \quad \text{for small } \Delta x$$

For multivariate set up \underline{X} with p.d.f $p(\underline{x})$

$P(\underline{X}$ will fall in a given region C , centered at say \underline{x})

$$= \int_C p(\underline{x}) d\underline{x} \approx V(C) p(\underline{x}), \text{ if we assume } C \ni$$

volume of $V(C)$ is small and
 $p(\underline{x})$ does not vary appreciably
within region C

$$\text{let } \theta = V(C) p(\underline{x})$$

Realize that θ can also be approximated by the proportion of samples falling within C

$$\text{i.e. } \theta \approx \frac{K}{n};$$

Where ; K , the number of samples falling within C
out of total n samples.

$$\text{i.e. } \frac{K}{n} \approx p(\underline{x}) V$$

$$\Rightarrow p(\underline{x}) = \frac{K}{n V}$$

K-nearest neighbor approach fixes K and then

determines the volume V which contains k samples centered at the point \underline{x} .

If \underline{x}_k is the k^{th} nearest neighbor pt to \underline{x} , then C may be taken to be sphere centered at \underline{x} with radius $\|\underline{x} - \underline{x}_k\|$. The volume of such a sphere in p dimension is $2^p \pi^{p/2} / p \Gamma_{p/2}$

Remark : This approach is different from the histogram approach wherein bin size is fixed.

(C) Kernel methods (Parzen methods)

Consider a 1-dimensional sample, x_1, \dots, x_n

An estimate of cumulative distⁿ F^n at x is

$$\hat{F}(x) = \frac{\# \text{ of observations} \leq x}{n}$$

estimate of p.d.f. at x :

$$\hat{f}(x) = \frac{\hat{F}(x+h) - \hat{F}(x-h)}{2h}$$

→ proportion of obsns falling within the interval $[x-h, x+h] / 2h$

i.e. Using a Kernel f^n (rectangular kernel)

$$K(z) = \begin{cases} \frac{1}{2}, & |z| \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$\hat{p}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

$\leftarrow \rightarrow$
 $= \frac{1}{2} (\# \text{ of observations with } h \text{ distance from } x)$

i.e. pts within h distance from x contribute $\frac{1}{2}$

to the density and pts outside this distance
contribute 0.

Remark: 'h' is referred to as spread or smoothing parameter (or band width)

Remark: Examples of popular univariate Kernel f^n s.

(i) ~~Rect~~ Rectangular : $K(z) = \begin{cases} \frac{1}{2}, & |z| \leq 1 \\ 0, & \text{otherwise} \end{cases}$

(ii) Triangular : $K(z) = \begin{cases} 1 - |z|, & \text{for } |z| \leq 1 \\ 0, & \text{otherwise} \end{cases}$

(iii) Gaussian : $K(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2} \neq z$

(iv) Bi-Weight/quartic : $K(z) = \begin{cases} \frac{15}{16} (1-z^2)^2, & |z| \leq 1 \\ 0, & \text{otherwise} \end{cases}$

(v) Bartlett-Epanechnikov:

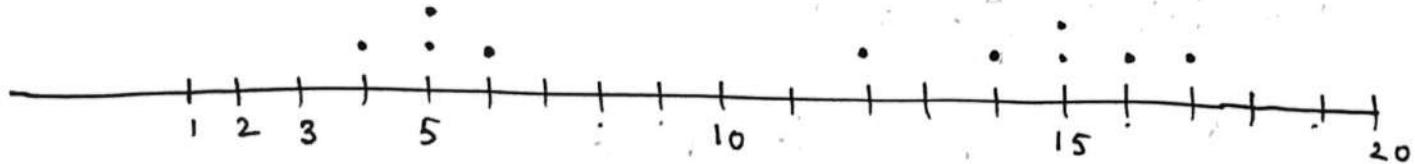
$$K(z) = \begin{cases} \frac{3}{4} (1 - \frac{z^2}{5}) / \sqrt{5}, & |z| \leq \sqrt{5} \\ 0, & \text{otherwise} \end{cases}$$

Example: $X = \{4, 5, 5, 6, 12, 14, 15, 15, 16, 17\}$

$n = 10$ samples

(a) Knn density estimate with $K=4$

$$\hat{p}(3) = \frac{4}{10} (V_4(3))^{-1} \quad \left(\hat{p}(z) = \frac{K}{n} V^{-1} \right)$$



$$\text{For } V_4(3), r = 3 \Rightarrow V_4(3) = 2r = 6$$

$$\hat{p}(3) = \frac{4}{10 \times 6} = \frac{1}{15}$$

$$\hat{p}(10) = \frac{4}{10} (V_4(10))^{-1}$$

$$\text{For } V_4(10), r = 5 \Rightarrow V_4(10) = 10$$

$$\Rightarrow \hat{p}(10) = \frac{4}{10 \times 10} = \frac{1}{25}$$

So $\hat{p}(15) = \frac{4}{10 \times 2} = \frac{1}{5} \quad (r=1)$

Kernel density estimate with rectangular kernel

with $h = 4$ bandwidth

$$\hat{p}(x) = \frac{1}{n h} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right); \quad K(z) = \begin{cases} \frac{1}{2}, & |z| \leq 1 \\ 0, & |z| > 1 \end{cases}$$

$$\text{i.e. } \hat{p}(x) = \frac{1}{n h} \sum_{i=1}^n \left(\frac{1}{2} I_{(|x-x_i| \leq h)} \right)$$

e.g.

$$\hat{p}(3) = \frac{1}{10 \times 4} \left(\sum_{i=1}^{10} \frac{1}{2} I_{(|3-x_i| \leq 4)} \right)$$

$$\text{i.e. } \hat{p}(3) = \frac{1}{40} \left(\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + 0 + 0 + \dots + 0 \right)$$

$$\text{i.e. } \hat{p}(3) = \frac{2}{40} = \frac{1}{20}$$

$$\hat{p}(10) = \frac{1}{10 \times 4} \left(\sum_{i=1}^{10} \frac{1}{2} I_{(|10-x_i| \leq 4)} \right)$$

$$\text{i.e. } \hat{p}(10) = \frac{1}{40} \times \frac{3}{2} = \frac{3}{80}$$

$$\hat{p}(15) = \frac{1}{40} \times \left(6 \times \frac{1}{2} \right) = \frac{3}{40}$$

- - - - -

Multivariate Kernel density estimate

Approach I : Assume independence of the component variables and estimate univariate Kernel density estimates for the components and get

$$\hat{p}(\underline{x}) = \prod_{i=1}^p \hat{p}_i(x_i)$$

Approach II : Generalization of univariate approach for multivariate case.

$$\hat{p}(\underline{x}) = \frac{1}{n h^p} \sum_{i=1}^n K\left(\frac{\underline{x} - \underline{x}_i}{h}\right)$$

$$\text{i.e. } \hat{p}(\underline{x}) = \frac{1}{n h^p} \sum_{i=1}^n K\left(\frac{x_1 - x_{i1}}{h}, \dots, \frac{x_p - x_{ip}}{h}\right)$$

Remark : A more general form is using different bandwidths

$$\hat{p}(\underline{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\prod_j h_j} K\left(\frac{x_1 - x_{i1}}{h_1}, \dots, \frac{x_p - x_{ip}}{h_p}\right)$$

Remark : A simple approach is to use a product Kernel

$$\hat{p}(\underline{x}) = \frac{1}{n h^p} \sum_{i=1}^n \left(\prod_{j=1}^p \hat{K}\left(\frac{x_j - x_{ij}}{h_j}\right) \right)$$

$$\text{or } \frac{1}{n} \sum_{i=1}^n \left(\prod_{j=1}^p \frac{1}{h_j} \hat{K}\left(\frac{x_j - x_{ij}}{h_j}\right) \right)$$

using diff bandwidth

$$\text{or } \frac{1}{n} \sum_{i=1}^n \left(\prod_{j=1}^p \frac{1}{h_j} \tilde{k}_j \left(\frac{x_j - x_{ij}}{h_j} \right) \right)$$



\tilde{K} (or \tilde{k}_j) is taken as one of the univariate kernels discussed earlier

Remark: Alternatively, one can use a genuine multivariable kernel

e.g. multivariate Gaussian Kernel

$$K(\underline{y}) = \frac{1}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2} \underline{y}' \underline{y}\right)$$

multivariate Epanechnikov Kernel, multivariate quartic Kernel are other choices of mult Kernel.

Epanechnikov Kernel

$$K(\underline{y}) = \begin{cases} (1 - \underline{y}' \underline{y}) (p+2)/2 c_p & \text{for } |\underline{y}| \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$c_p = \pi^{p/2} / \Gamma(p/2 + 1) = 2 \pi^{p/2} / p \Gamma(p/2)$$

Parametric density estimation

Most commonly used assumption: multivariate Gaussian or a mixture of multivariate Gaussian

Multivariate Gaussian:

$\underline{x}_1, \dots, \underline{x}_n$ realizations of $N_p(\underline{\mu}, \Sigma)$; $\Sigma > 0$

Use $\underline{x}_1, \dots, \underline{x}_n$ to find $\hat{f}(\underline{x})$ $\underline{x} \in \mathbb{R}^p$.

$$f(\underline{x}) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu})\right)$$

$\Theta = (\underline{\mu}, \Sigma)$ set of unknown parameters

Likelihood f^n

$$L(\Theta) = \prod_{j=1}^n f(\underline{x}_j)$$

$$L(\Theta) = (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu})\right)$$

Note that

$$\begin{aligned} & \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu}) \\ &= \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}} + \bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}} + \bar{\underline{x}} - \underline{\mu}) \\ &= \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}})' \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}}) + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\ &\quad + 2 \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}})' \cancel{\Sigma^{-1} (\bar{\underline{x}} - \underline{\mu})} \end{aligned}$$

$$\begin{aligned}
 &= \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}})' \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}}) + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &\Rightarrow = \text{tr} \left(\sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}})' \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}}) \right) + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &= \sum_{j=1}^n \text{tr} (\underline{x}_j - \bar{\underline{x}})' \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}}) + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &= \sum_{j=1}^n \text{tr} \Sigma^{-1} (\underline{x}_j - \bar{\underline{x}})(\underline{x}_j - \bar{\underline{x}})' + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &= \text{tr} \Sigma^{-1} \sum_{j=1}^n (\underline{x}_j - \bar{\underline{x}})(\underline{x}_j - \bar{\underline{x}})' + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &= \text{tr} \Sigma^{-1} (n-1)A + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \\
 &= \text{tr} \Sigma^{-1} A + n(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu})
 \end{aligned}$$

$$\Rightarrow L(\theta) = (2\pi)^{-n/2} |\Sigma|^{-n/2} \exp \left(-\frac{1}{2} \text{tr} \Sigma^{-1} A - \frac{n}{2} (\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \right) \quad —(*)$$

Note that for a fixed $\Sigma > 0$

$$(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu}) \geq 0 \text{ and } 0 \text{ only for } \bar{\underline{x}} = \underline{\mu}$$

$L(\underline{\mu}, \Sigma)$ for a fixed $\Sigma (> 0)$ is max if exponent is max w.r.t $\underline{\mu}$
 i.e. If $(\bar{\underline{x}} - \underline{\mu})' \Sigma^{-1} (\bar{\underline{x}} - \underline{\mu})$ is min w.r.t $\underline{\mu}$

i.e. If $\underline{\mu} = \bar{\underline{x}}$ \leftarrow indep of the fixed level of Σ

$$\Rightarrow \hat{\underline{\mu}}_{MLE} = \bar{\underline{x}}$$

The log-likelihood at $\hat{\mu} = \bar{x}$

$$l(\hat{\mu}, \Sigma) = \log L(\hat{\mu}, \Sigma)$$

$$= -\frac{n}{2} \log 2\pi + \frac{n}{2} \log |\Sigma'| - \frac{1}{2} \text{tr} \Sigma' A - (*)$$

maximization of $(*)$ w.r.t. Σ is equiv to maximization

$$\text{of } \frac{n}{2} \log |\Sigma'| - \frac{1}{2} \text{tr} \Sigma' A$$

$$= \frac{n}{2} \log |\Sigma' A| - \frac{1}{2} \text{tr} \Sigma' A - \frac{n}{2} \log |A|$$

i.e. maximization of

$$\frac{n}{2} \log |\Sigma' A| - \frac{1}{2} \text{tr} \Sigma' A - (**)$$

Let $\lambda_1, \dots, \lambda_p$ be eigenvalues of $\Sigma' A$

$$(**) = \frac{n}{2} \log \prod_{j=1}^p \lambda_j - \frac{1}{2} \sum_{j=1}^p \lambda_j$$

$$= \frac{n}{2} \sum_{j=1}^p \log \lambda_j - \frac{1}{2} \sum \lambda_j$$

$$= \frac{1}{2} \sum_{j=1}^p (n \log \lambda_j - \lambda_j) - (***)$$

(***) is maximized w.r.t λ_j at $\lambda_j = n$ & n

$$\text{i.e. } \Sigma' A = P n I_p P' = n I_p.$$

$$\Rightarrow \Sigma' = n A^{-1} \quad \text{i.e. } \Sigma = \frac{1}{n} A \text{ maximizes likelihood w.r.t } \Sigma$$

$$\Rightarrow \hat{\Sigma}_{MLE} = \frac{1}{n} \sum_{j=1}^n (\underline{x}_j - \hat{\underline{x}})(\underline{x}_j - \hat{\underline{x}})'$$

Based on $\tilde{x}_1, \dots, \tilde{x}_n$ obtain

$$\hat{\mu} = \bar{\tilde{x}}$$
 and $\hat{\Sigma} = \frac{1}{n} A = S_n$

Estimate density as

$$\hat{f}(\tilde{x}) = (2\pi)^{-p/2} |S_n|^{-1/2} \exp\left(-\frac{1}{2}(\tilde{x} - \bar{\tilde{x}})' S_n^{-1} (\tilde{x} - \bar{\tilde{x}})\right)$$

Mixture Normal setup.

One of the most widely used assumption

$$f(\tilde{x}) = \sum_{j=1}^g \pi_j f(\tilde{x} | \theta_j)$$

g : # of mixing densities

π_j : mixing proportion for j^{th} group/component

$f(\tilde{x} | \theta_j)$: density for j^{th} component in the mixture

j^{th} component $N_p(\mu_j, \Sigma_j)$; $j=1 \dots g$, $\Sigma_j > 0$

$$\theta_j = (\mu_j, \Sigma_j)$$

$$\Phi = (\pi_1, \dots, \pi_g, \mu_1, \Sigma_1, \dots, \mu_g, \Sigma_g)$$

unknown parameters

Likelihood f^n

$$L(\Phi) = \prod_{i=1}^n \left(\sum_{j=1}^g \pi_j f(x_i | \theta_j) \right)$$

Remark: (x_1, \dots, x_n) is incomplete data

use E-M algorithm

\underline{x} : incomplete data without class labels

$$\underline{y}' = (\underline{x}', \underline{z}') \quad \text{complete data}$$

where, $\underline{z} = \text{indicator vector of length } g \text{ with } 1 \text{ at the } k^{\text{th}} \text{ position}$

$$\begin{aligned} p(\underline{y}' | \Phi) &= p(\underline{x}, \underline{z} | \Phi) = \frac{p(\underline{x}, \underline{z}, \Phi)}{p(\underline{z}, \Phi)} \cdot \frac{p(\underline{z}, \Phi)}{p(\Phi)} \\ &= p(\underline{x} | \underline{z}, \Phi) p(\underline{z} | \Phi) \\ &= p(\underline{x} | \theta_k) \pi_k \end{aligned}$$

$$\therefore p(\underline{y}' | \Phi) = (p(\underline{x} | \theta_1) \pi_1)^{z_1} \cdots (p(\underline{x} | \theta_k) \pi_k)^{z_k} \cdots (p(\underline{x} | \theta_g) \pi_g)^{z_g}$$

$$\text{Let } z_j = \begin{cases} 1, & \text{if } j=k \\ 0, & \text{otherwise} \end{cases}$$

$$p(\underline{y}' | \Phi) = \prod_{j=1}^g (p(\underline{x} | \theta_j) \pi_j)^{z_j}$$

Consider for simplicity $g=2$

$$\underline{z} = (z_1, z_2) \quad z_1 = 1 \text{ if } \underline{x} \text{ correspond to component 1}$$

$$\text{Let } \pi_2 = \pi, \pi_1 = 1 - \pi \quad z_2 = 1 \text{ if } \underline{x} \text{ correspond to component 2}$$

$$p(\underline{y}' | \Phi) = (p(\underline{x} | \theta_1) \pi_1)^{z_1} (p(\underline{x} | \theta_2) \pi_2)^{z_2}$$

$$\therefore p(\underline{y}' | \Phi) = (p(\underline{x} | \theta_1) (1-\pi))^{z_1} (p(\underline{x} | \theta_2) \pi)^{z_2}$$

$$p(\underline{y}_1, \dots, \underline{y}_n | \Phi) = \prod_{i=1}^n \prod_{j=1}^2 (p(\underline{x}_i | \theta_j) \pi_j)^{z_{ij}}$$

$$= \prod_{i=1}^n \left(\{p(\underline{x}_i | \theta_1) (1-\pi)\}^{z_{1i}} \{p(\underline{x}_i | \theta_2) \pi\}^{z_{2i}} \right)$$

Note: For a general 'g',

$$g(y_1, \dots, y_n | \Phi) = \prod_{i=1}^n \left(\frac{g}{\pi} \left(p(x_i | \theta_j) \pi_j \right)^{z_{ji}} \right).$$

log likelihood f^n

$$\ell(\Phi) = \log g(y_1, \dots, y_n | \Phi)$$

$$= \sum_{i=1}^n \log \left(\left\{ p(x_i | \theta_1) (1-\pi) \right\}^{z_{1i}} \left\{ p(x_i | \theta_2) \pi \right\}^{z_{2i}} \right)$$

$$\text{i.e } \ell(\Phi) = \sum_{i=1}^n \left(z_{1i} \log (p(x_i | \theta_1) (1-\pi) + z_{2i} \log (p(x_i | \theta_2) \pi) \right)$$

$$= \sum_{i=1}^n \left(z_{1i} \log (p(x_i | \theta_1)) + z_{2i} \log (p(x_i | \theta_2)) \right) + \sum_{i=1}^n (z_{1i} \log (1-\pi) + z_{2i} \log \pi)$$

Note: For a general 'g'

$$\ell(\Phi) = \sum_{i=1}^n \left(\sum_{j=1}^g z_{ji} \log (p(x_i | \theta_j)) \right) + \sum_{i=1}^n \left(\sum_{j=1}^g z_{ji} \log \pi_j \right).$$

Remark:

Note that if (z_{1i}, z_{2i}) is known $\forall i$, the MLE is simple

$$\begin{cases} \underline{\mu}_1 \rightarrow \bar{x}_1 \\ \Sigma_1 \rightarrow S_1 \end{cases} \text{ from all } \bar{x}_i \Rightarrow z_{1i} = 1$$

$$\begin{cases} \underline{\mu}_2 \rightarrow \bar{x}_2 \\ \Sigma_2 \rightarrow S_2 \end{cases} \text{ from all } \bar{x}_i \Rightarrow z_{2i} = 1$$

But \bar{x}_i 's are unknown!

E - M algorithm steps

E - step :

$$E(z_{ji} | \underline{x}_i, \underline{\Phi}^{(m)}) = P(z_{ji} = 1 | \underline{x}_i, \underline{\Phi}^{(m)}) = \omega_{ji}$$

ω_{ji} : prob that $\underline{x}_i \in$ group j given current estimates $\underline{\Phi}^{(m)}$

$$\omega_{ji} = \frac{\pi_j^{(m)} p(\underline{x}_i | \theta_j)}{\pi_1^{(m)} p(\underline{x}_i | \theta_1) + \pi_2^{(m)} p(\underline{x}_i | \theta_2)}$$

Form the f^n

$$\begin{aligned} Q(\underline{\Phi}, \underline{\Phi}^{(m)}) &= \sum_{i=1}^n \left(\omega_{1i} \log(p(\underline{x}_i | \theta_1)) + \omega_{2i} \log(p(\underline{x}_i | \theta_2)) \right) \\ &\quad + \sum_{i=1}^n \sum_j \omega_{ji} \log \pi_j \end{aligned}$$

Note that $Q(\underline{\Phi}, \underline{\Phi}^{(m)}) = E(\log g(y_1, \dots, y_n) | \underline{x}_1, \dots, \underline{x}_n, \underline{\Phi})$

M-step : Maximize Q w.r.t. π_i & θ_i

Maximization of Q w.r.t. π_i subject to $\sum \pi_j = 1$

$$\tilde{Q} = Q - \lambda (\sum \pi_j - 1)$$

$$\frac{\partial \tilde{Q}}{\partial \pi_j} = \frac{\partial}{\partial \pi_j} \left(\sum_{i=1}^n \sum_j \omega_{ji} \log \pi_j \right) - \frac{\partial}{\partial \pi_j} (\lambda (\sum \pi_j - 1))$$

$$= \sum_{i=1}^n \frac{\omega_{ji}}{\pi_j} - \lambda = 0$$

$$\Rightarrow \lambda = \sum_{j=1}^n \omega_{ji} / \pi_j \quad i.e. \lambda \pi_j = \sum_{i=1}^n \omega_{ji}$$

$$\lambda \pi_j = \sum_i w_{ji}$$

$$\lambda \sum_j \pi_j = \sum_j \sum_i w_{ji} = \sum_{i=1}^n \left(\sum_j w_{ji} \right)$$

i.e. $\lambda = n$

$$\therefore \hat{\pi}_j = \frac{1}{n} \sum_{i=1}^n w_{ji}$$

Also for $\hat{\mu}_j$

$$\hat{\mu}_j = \frac{\sum_{i=1}^n w_{ji} \bar{x}_i}{\sum_{i=1}^n w_{ji}} = \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{ji} \bar{x}_i$$

$$\therefore \hat{\Sigma}_j = \frac{\sum_{i=1}^n w_{ji} (\bar{x}_i - \hat{\mu}_j)(\bar{x}_i - \hat{\mu}_j)'}{\sum_{i=1}^n w_{ji}}$$

$$\text{i.e. } \hat{\Sigma}_j = \frac{1}{n \hat{\pi}_j} \sum_{i=1}^n w_{ji} (\bar{x}_i - \hat{\mu}_j)(\bar{x}_i - \hat{\mu}_j)'$$

The E-M algorithm alternates betⁿ E-step of estimating w_{ji} and M-step of calculating $\hat{\pi}_j$, $\hat{\mu}_j$ on $\hat{\Sigma}_j$, given w_{ji} .

The iteration continues till convergence of likelihood.

Example : $p=1$; $q=2$

Comp 1 : $N(\mu_1, \sigma_1^2)$; $\theta = (\mu_1, \sigma_1^2)$

Comp 2 : $N(\mu_2, \sigma_2^2)$; $\theta = (\mu_2, \sigma_2^2)$

$$\bar{\Phi} = (\pi, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)$$

$$Q(\bar{\Phi}, \bar{\Phi}^{(m)}) = \sum_{i=1}^n \left(w_{1i} \log p(x_i | \theta_1) + w_{2i} \log p(x_i | \theta_2) \right)$$

E-step :

$$+ \sum_{i=1}^n \left(\sum_{j=1}^2 w_{ji} \log \pi_j \right) - (*)$$

$$w_{1i} = \frac{\pi_1^{(m)} p(x_i | \theta_1^{(m)})}{\pi_1^{(m)} p(x_i | \theta_1^{(m)}) + (1-\pi_1^{(m)}) p(x_i | \theta_2^{(m)})}$$

Starting r.h.s. can be obtained from cluster analysis output.

M-step : $\hat{\pi}_j = \frac{\sum_{i=1}^n w_{ji}}{n}; j = 1, 2$

and $\frac{\partial Q}{\partial \mu_j} = \sum_{i=1}^n w_{ji} \frac{\partial}{\partial \mu_j} \left(-\frac{1}{2} \log 2\pi \sigma_j^2 - \frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 \right)$

$$= \sum_{i=1}^n w_{ji} \left(\frac{1}{2\sigma_j^2} (x_i - \mu_j) \right) = 0$$

i.e. $\sum_{i=1}^n w_{ji} x_i = \mu_j \sum_{i=1}^n w_{ji}$

$$\Rightarrow \hat{\mu}_j = \frac{\sum_i w_{ji} x_i}{\sum_i w_{ji}} = \frac{\sum_i w_{ji} x_i}{n \hat{\pi}_j}$$

$$\begin{aligned}\frac{\partial \ell}{\partial \sigma_j^2} &= \sum_{i=1}^n w_{ji} \frac{\partial}{\partial \sigma_j^2} \left(-\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_j^2 - \frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 \right) \\ &= \sum_{i=1}^n w_{ji} \left(-\frac{1}{2\sigma_j^2} + \frac{1}{2(\sigma_j^2)^2} (x_i - \mu_j)^2 \right)\end{aligned}$$

$$\left. \begin{array}{l} \frac{\partial \ell}{\partial \mu_j} = 0 \\ \frac{\partial \ell}{\partial \sigma_j^2} = 0 \end{array} \right\} \Rightarrow \begin{array}{l} \hat{\mu}_j = \frac{1}{\sum_i w_{ji}} \sum_i w_{ji} x_i \\ \hat{\sigma}_j^2 = \frac{1}{\sum_i w_{ji}} \sum_i w_{ji} (x_i - \hat{\mu}_j)^2 \end{array}$$

Start with initial $(\pi_1^{(0)}, \theta_1^{(0)}, \theta_2^{(0)}) \rightarrow$ obtain

$(w_{1i}, w_{2i}) \quad i = 1 \dots n \rightarrow$ obtain $(\hat{\pi}_1, \hat{\mu}_1, \hat{\sigma}_1^2,$

$\hat{\mu}_2, \hat{\sigma}_2^2) = (\pi_1^{(1)}, \theta_1^{(1)}, \theta_2^{(1)}) \rightarrow$ alternate

bet "E-step & M-step till convergence of the likelihood.

Association Rule Mining

Market basket analysis: to find associations between different itemsets that customers place in shopping basket

Applications: cross marketing, catalog / floor design, design attractive packs, web log analysis for e-commerce,

Rule generation:

Antecedent \Rightarrow Consequent (support, confidence)

Data structure:

$T = \{t_1, \dots, t_n\}$ set of transactions

Each t_k is an itemset

$I = \{i_1, \dots, i_m\}$

Typical data

Cid	A1	A2	A3	...	Ak
c1	0	0	1	0	1
c2	1	0	0	1	0
⋮	⋮	⋮	⋮	⋮	⋮
cn	-	-	-	-	-

Aim :

- Find frequent patterns, i.e. associations among sets of items in T

- Represent these relationships as association rules of the form

$X \Rightarrow Y$ (support, confidence)

Important definitions:

Support count : # of occurrences of an itemset in the database

$$\sigma(\{\text{itemset}\})$$

Support : Fraction of transactions containing the itemset

$$S(\text{itemset}) = \frac{\sigma(\text{itemset})}{|T|}$$

Frequent itemset : An itemset whose support \geq a threshold, minsup

Confidence : A measure of how often B appears in
 [Rule: $(A \Rightarrow B)$]
 transactions containing A.

% of transactions containing A which also contains B

$$C(A \Rightarrow B) = \frac{s(A, B)}{s(A)} \quad (\text{est of conditional prob})_{B|A}$$

Example:

Tid	Items
T1	bread, egg, peanut-butter
T2	bread, peanut-butter
T3	bread, milk, peanut-butter
T4	beer, bread
T5	beer, milk

Rule	Support	Confidence
$\text{bread} \Rightarrow \text{peanut-butter}$	$\frac{3}{5} = 0.6$	$\frac{3}{4} = 0.75$
$\text{peanut-butter} \Rightarrow \text{bread}$	$\frac{3}{5} = 0.6$	$\frac{3}{3} = 1.0$
$\text{beer} \Rightarrow \text{bread}$	$\frac{1}{5} = 0.2$	$\frac{1}{2} = 0.5$
$\text{peanut-butter} \Rightarrow \text{egg}$	$\frac{1}{5} = 0.2$	$\frac{1}{3} = 0.33$
$\text{egg} \Rightarrow \text{peanut-butter}$	$\frac{1}{5} = 0.2$	$\frac{1}{1} = 1$
$\text{egg} \Rightarrow \text{milk}$	0	0

ARM task: Given a set of transactions, the goal of ARM is to find all rules \Rightarrow

- (i) support $\geq \text{min sup}$
- & (ii) confidence $\geq \text{min conf}$

(min sup, min conf) : fixed apriori

Brute force approach

- List all possible association rules
- Compute support & confidence for each rule
- prune as per threshold

Approach is computationally prohibitive

The Apriori Algorithm

Agrawal & Srikant: "Fast algorithms for mining association rules in large databases", Int Conf on VLDB, 1994.

2-Step ARM of Apriori algorithm

S1: Generate all frequent itemsets with support $\geq \text{min sup}$

S2: Generate association rules using these frequent itemsets

Anti-monotonicity property of Apriori algorithm

Downward closure property

- Any subset of a frequent itemset is frequent

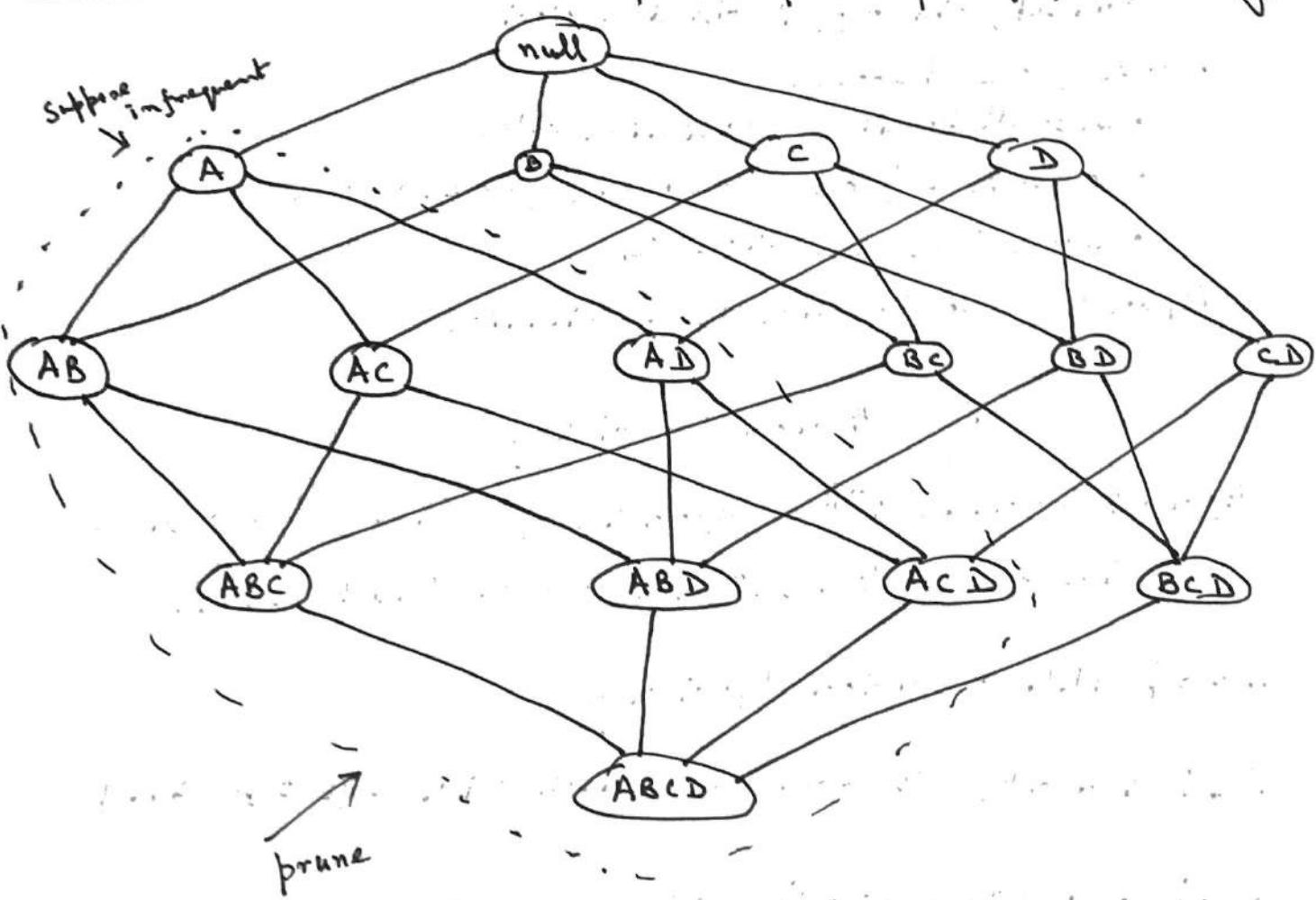
i.e. $\forall X \subseteq Y$

$$S(X) \geq S(Y) \quad \forall X, Y$$

\Rightarrow If an itemset is not frequent, none of its supersets can be frequent (\Rightarrow prune all such sets)

\Rightarrow If an itemset is not frequent, there is no need to explore its supersets.

Example: Itemset lattice \rightarrow apriori principle of pruning



- Suppose B is found to be infrequent (support < min sup)
 - \Rightarrow all itemsets with B will also be infrequent
 - \Rightarrow prune all such branches. (i.e. all its supersets)
 - i.e. exclude itemset AB, BC, BD, ABC, ABD, BCD, ABCD.
- suppose AB is found to be infrequent
 - \Rightarrow all itemsets with AB will also be infrequent
 - \Rightarrow prune all such branches (i.e. all supersets)
 - i.e. exclude itemsets, ABC, ABD, ABCD.

Apriori Step-1 in pseudo codes

- $K = 1$
- Generate frequent itemsets of length 1
- Prune itemsets of higher orders (i.e. supersets), if necessary
- Generate itemsets of length $K+1$ from frequent itemsets of length K
- Compute the support of new candidate itemsets w.r.t. min sup. Prune if necessary.
- $K = K+1$
- Repeat until no frequent itemsets are found.

Generation (step) of Itemsets for next level

Let L_K denote the frequent itemsets at level K and C_K denote the set of all candidates at level K

- Items in L_{K-1} are listed in an order

Step 1: Self joining $L_{K-1} * L_{K-1}$

i.e. joining of 2 items from L_{K-1}

$\{p.\text{item}_1, p.\text{item}_2, \dots, p.\text{item}_{K-1}\}$
 $\wedge \{q.\text{item}_1, q.\text{item}_2, \dots, q.\text{item}_{K-1}\}$ under the
 He itemset given order

insert into C_K , $p.\text{item}_1 p.\text{item}_2 \dots p.\text{item}_{K-2} p.\text{item}_{K-1} q.\text{item}_{K-1}$

Step 2: Pruning of C_K set

Itemsets c in C_K and

$(K-1)$ subsets s of c (under the given order)

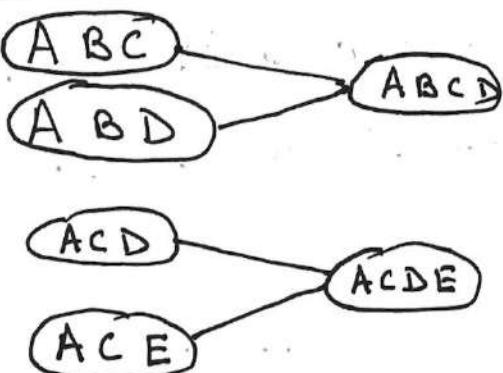
If (s is not in L_{K-1}) delete c from C_K

Justification is anti-monotone property

Example

$$L_3 = \{ABC, ABD, ACD, ACE, BCD\}$$

Self joining step: $L_3 * L_3$



Pruning step:

Consider $A C D E \rightarrow$ subsets $A C D, A C E, C D E, A D E$

$C D E / A D E$ is not in L_3

\Rightarrow Prune $A C D E$ from C_4

Consider $A B C D \rightarrow$ subsets $A B C, A B D, A C D, B C D$
all subsets in L_3

\Rightarrow Pruning not reqd for $A B C D$

\Rightarrow Pass $A B C D$ to C_4

Example

$T_1 : \{A, C, D\}$

$T_2 : \{B, C, E\}$

$T_3 : \{A, B, C, E\}$

$T_4 : \{B, E\}$

C_1		Freq. Itemset		C_2
item	support count	L_1		
A	2	A	2	AB
B	3	B	3	Ac
C	3	C	3	AE
D	-1- · · · prune	E	3	Bc
E	3			BE
				CE

2nd scan

Itemset	Support Count	L_2	C_3	L_3
- AB	1			
- AC	2			
- AE	+			
BC	2			
BE	3			
CE	2			

all other joins pruned.

$S(BC\bar{E}) = 2$

Under C_3 , the only Itemset $BC\bar{E}$ has support count 2, i.e. $\geq \text{minsup}$

Step 2 of apriori algorithm

Generate association rules using frequent itemsets

Given any frequent itemset L :

- find all non-empty subsets F of L
- output each rule $F \Rightarrow \{L-F\}$ that satisfies the threshold on confidence (min conf)

Example: Let $L = \{A, B, C\}$ is a frequent itemset

Candidate rules are

$$\begin{array}{lll} AB \Rightarrow C; & AC \Rightarrow B; & BC \Rightarrow A \\ A \Rightarrow BC; & B \Rightarrow AC; & C \Rightarrow AB \end{array}$$

In general, there are $2^{|L|} - 2$ candidate rules.

Remark: Efficiency in rule generation

confidence of rules generated from the same itemset have anti-monotone property

$$c(ABC \Rightarrow D) \geq c(AB \Rightarrow CD) \geq c(A \Rightarrow BCD)$$

say $c(ABC \Rightarrow D) \geq c(AC \Rightarrow BD) \geq c(c \Rightarrow ABD)$

[Consider, e.g., $ABC \Rightarrow D$ & $AB \Rightarrow CD$

$$AB \cdot C \bar{A} \bar{B} \bar{C}$$

$$S(AB) \geq S(ABC)$$

$$\Rightarrow \frac{1}{S(ABC)} \leq \frac{1}{S(AB)}$$

$$\Rightarrow \frac{S(ABCD)}{S(ABC)} \geq \frac{S(ABCD)}{S(AB)}$$

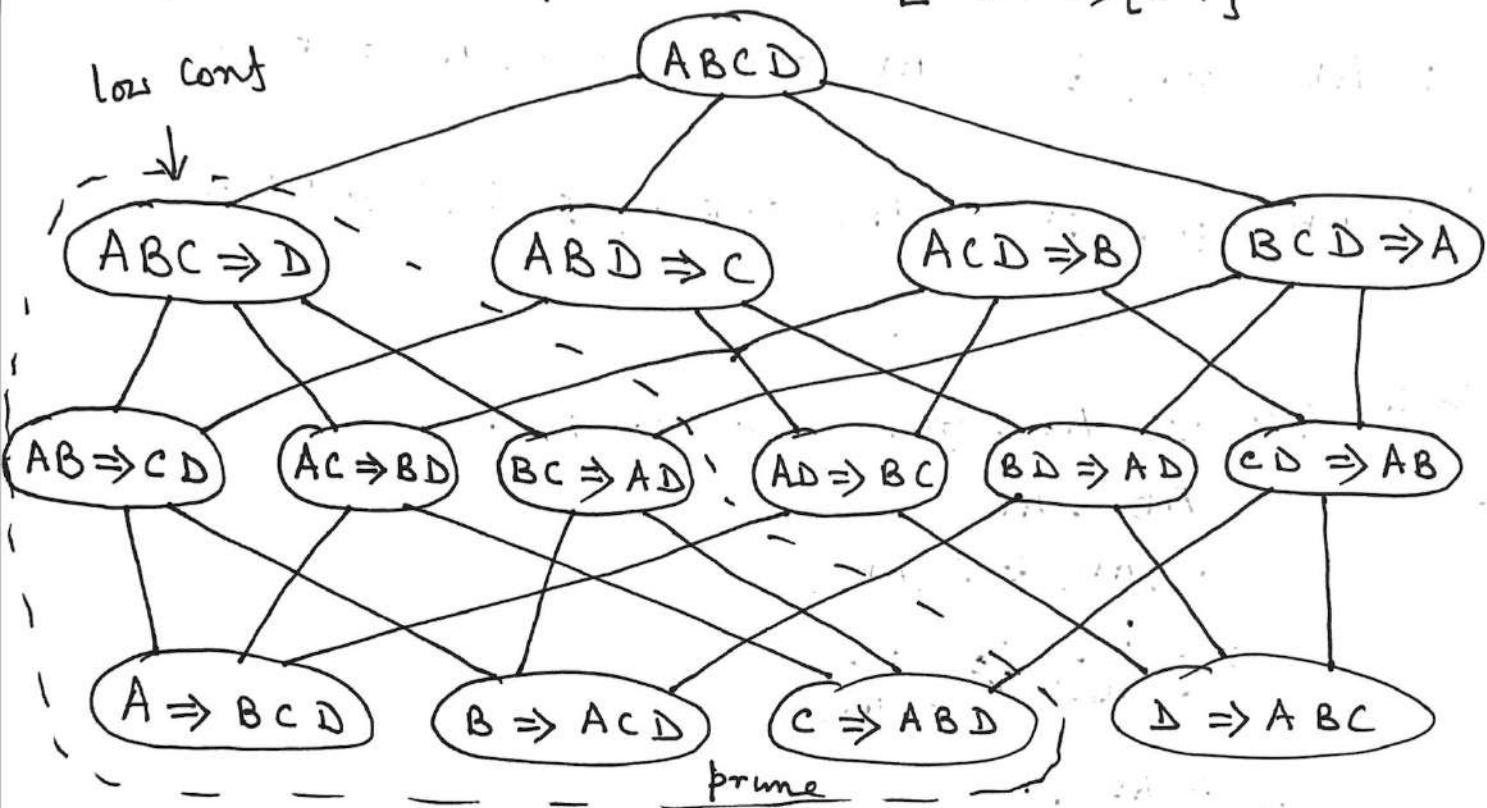
$$\text{i.e. } c(ABC \Rightarrow D) \geq c(AB \Rightarrow CD)$$

We apply this property to prune rule generation sequence.

Example of pruning using anti-monotone property

Suppose $ABCD$ is a frequent itemset and look at possible rules for ARM

$$L \rightarrow F \Rightarrow \{L - F\}$$



~~Suppose~~ Suppose the rule $ABC \Rightarrow D$ is low conf
i.e. $C(ABC \Rightarrow D) < \text{min conf}$

Then by anti-monotone property, all sub rules below it in the lattice will be $< \text{min conf}$ and can be pruned, i.e.

$$AB \Rightarrow CD, AC \Rightarrow BD, BC \Rightarrow AD$$

$$A \Rightarrow BCD, B \Rightarrow ACD \& C \Rightarrow ABD$$

are pruned.

Example

QD

$$T_1 : \{A, C, D\}$$

$$T_2 : \{B, C, E\}$$

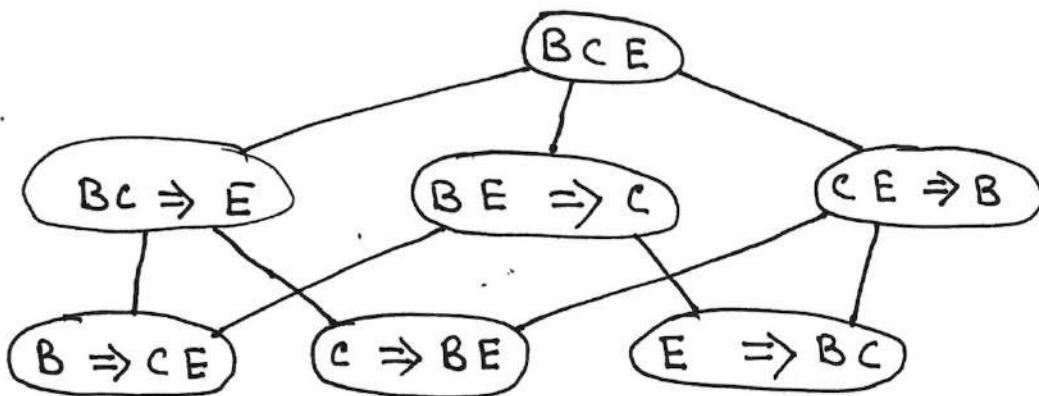
$$T_3 : \{A, B, C, E\}$$

$$T_4 : \{B, E\}$$

Frequent itemsets with $\text{min sup} \geq 0.5$ are
 $s(\cdot) = 2$

$$BCE, AC, BC, BE, CE \\ s(\cdot) = 2 \quad s(\cdot) = 2 \quad s(\cdot) = 3 \quad s(\cdot) = 2$$

Consider rules for BCE itemset



$$C(BC \Rightarrow E) = \frac{s(BC E)}{s(BC)} = \frac{2}{2} = 1$$

$$\Rightarrow BC \Rightarrow E (.5, 1)$$

$$C(BE \Rightarrow C) = \frac{s(BC E)}{s(BE)} = \frac{2}{3}$$

if min Conf = 0.75, then $BE \Rightarrow C$

along with $B \Rightarrow CE$ & $E \Rightarrow BC$
are prunned.

$$C(CE \Rightarrow B) = \frac{s(BC E)}{s(CE)} = \frac{2}{2} = 1$$

$$C(C \Rightarrow BE) = \frac{s(BC E)}{s(C)} = \frac{2}{3} \quad \text{pruned at min Conf 0.75}$$

Discrimination & Classification

Discriminant analysis: optimal way to separate heterogeneous populations

Classification: classification of new observation vector in one of the possible populations (classes) using discriminant function.

Examples / Applications

Medical diagnostics, alert systems for adverse events, credit card fraud detection, loan classification,

Data structure

Learning sample

$$\mathcal{L} = \{(x_1, j_1), \dots, (x_n, j_n)\}$$

n preclassified examples

$$\tilde{x}_i \in \mathbb{X}, \text{ feature space } \forall i = 1(1)n$$

$$j_i \in \{1, 2, \dots, J\} \quad \forall i = 1(1)n$$

J : possible number of classes.

$$\mathcal{C} = \{1, 2, \dots, J\} - \text{set of classes}$$

aim : Given $\tilde{x} \in \mathbb{X}$; to find a systematic way of predicting class membership
 i.e. assigns one of the classes in $\{1, 2, \dots, J\}$ to $\tilde{x} \in \mathbb{X}$.

Def": Classifier

A classifier or a classification rule is a function $d(\cdot)$ defined on $\mathcal{X} \ni x$, $d(x)$ is one equal to one of the numbers $1, 2, \dots, J$.

Note: An alternate way to look at a classifier is through partitions

$$\text{Let } A_j = \{x : d(x) = j\}$$

$$A_1, \dots, A_J \text{ are } \Rightarrow$$

$$A_i \cap A_j = \emptyset \quad \forall i \neq j$$

$$\& \bigcup_i A_i = \mathcal{X}$$

Alt def": A classifier is a partition of \mathcal{X} into J disjoint sets $A_1, \dots, A_J \ni x \in A_j$ the assigned class membership is j .

Fisher Linear Discriminant Function (FLDF)

Consider a 2-class problem

i.e. 2 pop's π_1 & π_2

$$\underline{x} | \pi_1 \sim \underline{\mu}_1, \Sigma ; \underline{\mu}_1 = E(\underline{x} | \pi_1), \Sigma = \text{cov}(\underline{x} | \pi_1)$$

$$\Sigma > 0;$$

$$\underline{x} | \pi_2 \sim \underline{\mu}_2, \Sigma ; \underline{\mu}_2 = E(\underline{x} | \pi_2), \Sigma = \text{cov}(\underline{x} | \pi_2)$$

- Aim :
- To find a f": $g(\cdot) \Rightarrow$ If \underline{x}_1 is from π_1 & \underline{x}_2 is from π_2 , then $g(\underline{x}_1) \neq g(\underline{x}_2)$ should look as "different as possible". Such a $g(\cdot)$ will be the constructed discriminant f".
 - Given a new obsn \underline{x} , use $g(\cdot)$ to classify it to one of π_1 & π_2 .

FLDF approach : change π_1 & π_2 into univariate pop

by transforming \underline{x} to $\underline{l}' \underline{x}$

$$\Rightarrow \underline{l}' \underline{x} | \pi_1 \sim \underline{l}' \underline{\mu}_1, \underline{l}' \Sigma \underline{l}$$

$$\& \underline{l}' \underline{x} | \pi_2 \sim \underline{l}' \underline{\mu}_2, \underline{l}' \Sigma \underline{l}$$

Note that

Separate out the 2 univariate pop as much as possible

\Leftrightarrow maximization of distance bet" the 2 pop's w.r.t. \underline{l}

Take the distance betⁿ the 2 univariate pop's as

$$\frac{(\underline{\lambda}' \underline{\mu}_1 - \underline{\lambda}' \underline{\mu}_2)^2}{\underline{\lambda}' \Sigma \underline{\lambda}} - (*)$$

Maximize (*) w.r.t. $\underline{\lambda}$

i.e. find $\underline{\lambda}^* = \arg \max_{\underline{\lambda}} \frac{(\underline{\lambda}' (\underline{\mu}_1 - \underline{\mu}_2))^2}{\underline{\lambda}' \Sigma \underline{\lambda}}$

Define $\underline{a}' = \underline{\lambda}' \Sigma^{\frac{1}{2}}$

Then $\frac{(\underline{\lambda}' (\underline{\mu}_1 - \underline{\mu}_2))^2}{\underline{\lambda}' \Sigma \underline{\lambda}}$

$$= \frac{(\underline{a}' \Sigma^{\frac{1}{2}} (\underline{\mu}_1 - \underline{\mu}_2))^2}{\underline{a}' \underline{a}} - (*^2)$$

By C-S inequality

$$(*)^2 \leq \frac{(\underline{a}' \underline{a}) (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)}{\underline{a}' \underline{a}}$$

$$= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)$$

→ square of Mahalanobis distance

Note that equality in the above is attained

at

$$\underline{a}' = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{\frac{1}{2}}$$

i.e. $\underline{\lambda}' \Sigma^{\frac{1}{2}} = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{\frac{1}{2}}$

$$\text{i.e. } \underline{\lambda}' = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1}$$

So we get the FLDF as

$$g(\underline{x}) = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \quad \text{as the discriminant f'}$$

Next task is to frame the classification rule based on FLDF

i.e. given a new obsn \underline{x}_0 , to assign it to π_1 or π_2

Realize that

$$E((\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} | \pi_1) = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{\mu}_1 = m_1, \text{ say}$$

$$\& E((\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} | \pi_2) = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{\mu}_2 = m_2, \text{ say}$$

$$\text{and } m_1 - m_2 = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)$$

$$\Rightarrow m_1 - m_2 \geq 0 ; \text{ equal to 0 only if } \underline{\mu}_1 = \underline{\mu}_2$$

$$\text{i.e. } m_1 \geq m_2$$

$$\text{let } y_0 = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x}_0$$

We will assign \underline{x}_0 to π_1 If $y_0 \geq 0$

y_0 is closer to m_1 than to m_2

i.e. If $y_0 \geq \frac{m_1 + m_2}{2}$; we assign \underline{x}_0 to π_1

i.e. if $y_0 \geq \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$

and assign \underline{x}_0 to π_2 if

$$y_0 < \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$$

Then the classification partition associated with FLD is (R_1, R_2)

$$R_1 = \{ \underline{x} : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \geq \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \}$$

$$R_2 = \{ \underline{x} : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} < \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \}$$

Remark: Usually $\underline{\mu}_1, \underline{\mu}_2, \Sigma$ are unknown. We use

the learning sample \mathcal{D} , of preclassified examples, to estimate $\underline{\mu}_1, \underline{\mu}_2, \Sigma$ and get

Sample analogue of $(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x}$ as

$$\underline{(\bar{x}_1 - \bar{x}_2)' S^{-1} x}$$

→ This is called the Fisher sample linear discriminant function.

Where, (\bar{x}_1, S_1) are sample mean vector and sample variance covariance matrix computed from pre-classified examples of π_1 class.

Similarly, $(\bar{x}_2, S_2) \rightarrow$ computed from pre-classified examples of π_2 class.

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

n_1 : # of pre-classified π_1 cases in \mathcal{L}

n_2 : # of pre-classified π_2 cases in \mathcal{L}

$$\hat{R}_1 = \left\{ \underline{x} : (\bar{x}_1 - \bar{x}_2)' \bar{S}^{-1} \underline{x} \geq \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \bar{S}^{-1} (\bar{x}_1 + \bar{x}_2) \right\}$$

$$\hat{R}_2 = \left\{ \underline{x} : (\bar{x}_1 - \bar{x}_2)' \bar{S}^{-1} \underline{x} < \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \bar{S}^{-1} (\bar{x}_1 + \bar{x}_2) \right\}$$

General Classification problem

Two-class problem

π_1 & π_2 : two classes

Suppose,

$\underline{x} | \pi_1$ has support \mathcal{X}_1 $\mathcal{X}_1 \cup \mathcal{X}_2 = \mathcal{X}$

$\underline{x} | \pi_2$ has support \mathcal{X}_2

If $\mathcal{X}_1 \cap \mathcal{X}_2 = \emptyset$, then there is no misclassification.

But usually $\mathcal{X}_1 \cap \mathcal{X}_2 \neq \emptyset$ and there is a chance

that observations can be misclassified.

Let

$c(i|j)$ denote the cost when an observation from π_j is misclassified as from π_i .

$$c(1|1) = c(2|2) = 0$$

Suppose the class conditional densities are

$$f_i(\underline{x}) ; i=1, 2$$

i.e. $\underline{x} | \pi_i$ has p.d.f. $f_i(\underline{x})$ $i=1, 2$

and the prior probabilities for π_1 and π_2 be p_1 and p_2 , resp.

Let $P(i|j)$: probability of misclassifying an observation from π_j into π_i .

Suppose the classification partition is $\{R_1, R_2\}$, (i.e. if $x \in R_i$, we classify it to π_i). Then

$$P(i|j) = \int_{R_i} f_j(x) dx$$

$$\text{e.g. } P(2|1) = \int_{R_2} f_1(x) dx$$

$$P(1|2) = \int_{R_1} f_2(x) dx$$

Also $P(\text{an obsn comes from } \pi_i \text{ and is misclassified to } \pi_j) = p_i P(j|i)$

Define

- Total Probability of misclassification (TPM).

$$= p_1 P(2|1) + p_2 P(1|2)$$

- Expected cost of misclassification (ECM)

$$= C(2|1) (p_1 P(2|1)) + C(1|2) (p_2 P(1|2))$$

Optimal Strategies

To find the partition $\{R_1^*, R_2^*\}$, say, \Rightarrow

TPM or ECM is minimized.

(I) TPM minimizing partition

$$\text{TPM} = p_2 \int_{R_1} f_2(x) dx + p_1 \int_{R_2} f_1(x) dx$$

$$= p_2 \int_{R_1} f_2(x) dx + p_1 \int_{x=R_1} f_1(x) dx$$

$$= p_2 \int_{R_1} f_2(x) dx + \int_{x=R_1} p_1 f_1(x) dx - \int_{R_1} p_1 f_1(x) dx$$

i.e.

$$\text{TPM} = \int_{R_1} (p_2 f_2(x) - p_1 f_1(x)) dx + \underbrace{\int_{x=R_1} p_1 f_1(x) dx}_{\text{indep of partition}}$$

Minimization of TPM w.r.t. $\{R_1, R_2\}$

\Leftrightarrow

$$\text{minimization of } \int_{R_1} (p_2 f_2(x) - p_1 f_1(x)) dx$$

Note that $\forall x \in \mathbb{X}; (p_2 f_2(x) - p_1 f_1(x))$ is either ≤ 0 or > 0 .

Thus, the TPM is minimized by choosing

$$R_1^* = \{ \underline{x} : b_2 f_2(\underline{x}) \leq b_1 f_1(\underline{x}) \}$$

i.e. $R_1^* = \left\{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq \frac{b_2}{b_1} \right\}$

& $R_2^* = \left\{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} < \frac{b_2}{b_1} \right\}$

In other words,

TPM minimizing classification rule is

assign \underline{x} to Π_1 if $\frac{f_1(\underline{x})}{f_2(\underline{x})} \geq \frac{b_2}{b_1}$

& to Π_2 if otherwise .

Remark: Consider testing of

$$H_0: \underline{x} \sim f_1(\underline{x}) \text{ ag } H_A: \underline{x} \sim f_2(\underline{x})$$

By Neyman-Pearson lemma, the MP(α) test has critical region

$$\omega = \left\{ \underline{x} : \frac{f_2(\underline{x})}{f_1(\underline{x})} > K \right\}$$

where K is \exists

$$P_{f_1} \left(\frac{f_2(\underline{x})}{f_1(\underline{x})} > K \right) = \alpha$$

Thus we observe that

$$R_2^* = \left\{ \underline{x} : \frac{f_2(\underline{x})}{f_1(\underline{x})} > \frac{p_1}{p_2} \right\}$$

Corresponds to the critical region of an MP test
of a fixed size (fixed by $\frac{p_1}{p_2}$)

$$P_{f_1} \left(\frac{f_2(\underline{x})}{f_1(\underline{x})} > \frac{p_1}{p_2} \right) = \alpha^*, \text{ say.}$$

Example :

		π_1		
		$f_1(\underline{x})$		
		1	2	3
x_1				
1		.1	.05	.15
2		.25	.2	.25

		π_2		
		$f_2(\underline{x})$		
		1	2	3
x_1				
1		.2	.2	.2
2		.2	.1	.1

$$p_1 = p_2 = \frac{1}{2}$$

TPM minimizing rule is

assign \underline{x} to π_1 if $f_1(\underline{x}) \geq f_2(\underline{x})$

$$\text{i.e. } R_1^* = \left\{ \underline{x} : f_1(\underline{x}) \geq f_2(\underline{x}) \right\}$$

$$\Rightarrow R_1^* = \{(1, 2), (2, 2), (3, 2)\}$$

$$R_2^* = \{(1, 1), (2, 1), (3, 1)\}$$

$$P(1|2) = \sum_{\underline{x} \in R_1^*} f_2(\underline{x})$$

$$= 0.2 + 0.1 + 0.1 = 0.4$$

$$P(2|1) = \sum_{\underline{x} \in R_2^*} f_1(\underline{x})$$

$$= 0.1 + 0.05 + 0.15 = 0.3$$

$$TPM(R_1^*, R_2^*) = \frac{1}{2} (P(1|2) + P(2|1))$$

$$= 0.5 (0.4 + 0.3)$$

Note: $\hat{p}_1 \neq \hat{p}_2$

If $\hat{p}_1 = 0.4$ & $\hat{p}_2 = 0.6$ Then

$$R_1^* = \{\underline{x} : 0.4 f_1(\underline{x}) \geq 0.6 f_2(\underline{x})\}$$

$$R_2^* = \{\underline{x} : 0.4 f_1(\underline{x}) < 0.6 f_2(\underline{x})\}$$

Example 2 :

$$\pi_1 \equiv N_p(\underline{\mu}_1, \Sigma)$$

$$\pi_2 \equiv N_p(\underline{\mu}_2, \Sigma); \Sigma > 0$$

Optimum TPM rule is

$$R_1^* = \left\{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq 1 \right\}$$

$$R_2^* = \left\{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} < 1 \right\}$$

Now,

$$\frac{f_1(\underline{x})}{f_2(\underline{x})} \geq 1.$$

\Leftrightarrow

$$(2\pi)^{-\frac{1}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_1)' \Sigma^{-1} (\underline{x} - \underline{\mu}_1)\right) \\ \geq (2\pi)^{-\frac{1}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_2)' \Sigma^{-1} (\underline{x} - \underline{\mu}_2)\right)$$

i.e. $(\underline{x} - \underline{\mu}_1)' \Sigma^{-1} (\underline{x} - \underline{\mu}_1) \leq (\underline{x} - \underline{\mu}_2)' \Sigma^{-1} (\underline{x} - \underline{\mu}_2)$

i.e. ~~$\underline{x}' \Sigma^{-1} \underline{x} + \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 - 2 \underline{\mu}_1' \Sigma^{-1} \underline{x}$~~

$$\leq \cancel{\underline{x}' \Sigma^{-1} \underline{x}} + \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 - 2 \underline{\mu}_2' \Sigma^{-1} \underline{x}$$

i.e. $-2 (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \leq \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1$

i.e. $-2 (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \leq \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 \\ + \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_2$

i.e. $-2 (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \leq -(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$

i.e. $(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$

Hence the TPM minimizing rule is

Assign \underline{x} to π_1 , if

$$(\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} \underline{x} \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$$

Remark: to π_2 if δ/κ

The above rule is same as FLD

Remark: if $\pi_1 = N_p(\underline{\mu}_1, \Sigma)$

$$\pi_2 = N_p(\underline{\mu}_2, \Sigma) \quad \Sigma > 0$$

then

$$R_1^* = \left\{ \underline{x} : (\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} \underline{x} \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) + \log\left(\frac{p_2}{p_1}\right) \right\}$$

$$R_2^* = \left\{ \underline{x} : (\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} \underline{x} < \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \bar{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) + \log\left(\frac{p_1}{p_2}\right) \right\}$$

Remark: 1, the learning sample is to be used to obtain estimated partition $(\hat{R}_1^*, \hat{R}_2^*)$

$$\hat{R}_1^* = \left\{ \underline{x} : (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' \bar{S}^{-1} \underline{x} \geq \frac{1}{2} (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' \bar{S}^{-1} (\bar{\underline{x}}_1 + \bar{\underline{x}}_2) + \log\left(\frac{p_2}{p_1}\right) \right\}$$

Remark:

$(\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' \bar{S}^{-1} \underline{x} - \frac{1}{2} (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' \bar{S}^{-1} (\bar{\underline{x}}_1 + \bar{\underline{x}}_2)$ is called Anderson's classification statistic

Optimum classification rules (ECM)

Recall that we have earlier derived optimum classification rule minimizing "Total Probability of Misclassification (TPM)". We now consider the problem of deriving optimum classification rule minimizing "Expected Cost of Misclassification (ECM)".

You may recall that for a 2-class problem, ECM is given by

$$ECM = p_1 P(2|1) C(2|1) + p_2 P(1|2) C(1|2)$$

p_1 & p_2 are prior probabilities of classes π_1 & π_2

$C(i|j)$: cost of misclassifying an obsn from π_j into π_i

$P(i|j)$: prob of misclassifying an obsn from π_j into π_i

$$P(1|2) = \int_{R_1} f_2(x) dx ; P(2|1) = \int_{R_2} f_1(x) dx$$

$f_i(x)$: class conditional density of class π_i

$\{R_1, R_2\}$ is the classification partition

$$ECM = p_1 C(2|1) \int_{R_2} f_1(x) dx + p_2 C(1|2) \int_{R_1} f_2(x) dx$$

$$= p_1 C(2|1) \int_{x-R_1} f_1(x) dx + p_2 C(1|2) \int_{R_1} f_2(x) dx$$

$$= \int_{R_1} \left(p_2 C(1|2) - f_2(x) - p_1 C(2|1) f_1(x) \right) dx + \left(p_1 C(2|1) \int_{x-R_1} f_1(x) dx \right) \leftarrow \text{indep of partition } \{R_1, R_2\}$$

The above ECM is minimized w.r.t. $\{R_1, R_2\}$ by choosing

$$R_1^* = \left\{ \underline{x} : b_2 C(1|2) f_2(\underline{x}) \leq b_1 C(2|1) f_1(\underline{x}) \right\}$$

$$R_2^* = \underline{x} - R_1^* = \left\{ \underline{x} : b_2 C(1|2) f_2(\underline{x}) > b_1 C(2|1) f_1(\underline{x}) \right\}$$

(logic is same as used for deriving opt TPM rule in class)

Remark: If $C(1|2) = C(2|1)$, then opt ECM rule is same as TPM rule.

Remark: The opt ECM partition $\{R_1^*, R_2^*\}$ correspond to a Most Powerful (MP) test for testing

$$H_0: \underline{x} \sim f_1(\underline{x}) \text{ ag } H_A: \underline{x} \sim f_2(\underline{x})$$

This MP test is of fixed size.

(Recall the logic given in class during discussion on connection bet" opt TPM partition and MP test using Neyman-Pearson lemma)

Example 1: $\pi_1: N_p(\underline{\mu}_1, \Sigma)$

$$\pi_2: N_p(\underline{\mu}_2, \Sigma) \quad \Sigma > 0$$

$$b_1 f_1(\underline{x}) C(2|1) \geq b_2 f_2(\underline{x}) C(1|2)$$

$$\Leftrightarrow \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq \frac{b_2 C(1|2)}{b_1 C(2|1)}$$

$$\Leftrightarrow \frac{(2\pi)^{-p/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2} (\underline{x} - \underline{\mu}_1)' \Sigma^{-1} (\underline{x} - \underline{\mu}_1))}{(2\pi)^{-p/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2} (\underline{x} - \underline{\mu}_2)' \Sigma^{-1} (\underline{x} - \underline{\mu}_2))} \geq \frac{b_2 C(1|2)}{b_1 C(2|1)}$$

$$\begin{aligned}
 & \Leftrightarrow -\frac{1}{2} \left((\underline{x} - \underline{\mu}_1)' \Sigma^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_2)' \Sigma^{-1} (\underline{x} - \underline{\mu}_2) \right) \\
 & \quad \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right). \\
 & \Leftrightarrow -\frac{1}{2} \left(\cancel{\underline{x}' \Sigma^{-1} \underline{x}} + \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 - 2 \underline{\mu}_1' \Sigma^{-1} \underline{x} - \cancel{\underline{x}' \Sigma^{-1} \underline{x}} - \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 + 2 \underline{\mu}_2' \Sigma^{-1} \underline{x} \right) \\
 & \quad \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \\
 & \Leftrightarrow -\frac{1}{2} \left(-2 \underbrace{(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x}}_{\text{FLDF}} - (\underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 + \underbrace{\underline{\mu}_2' \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_2}_{=0}) \right) \\
 & \quad \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \\
 & \Leftrightarrow -\frac{1}{2} \left(-2 (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} - (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right) \\
 & \quad \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right)
 \end{aligned}$$

i.e. $R_1^* = \left\{ \underline{x} : +\frac{1}{2} \left(2 (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} - (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right) \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \right\}$

$$R_2^* = \exists \epsilon \cdot R_1^*$$

Sp. Case : If $p_2 C(1|2) = p_1 C(2|1)$, then

$$\tilde{R}_1^* = \left\{ \underline{x} : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right\}$$

which is same as classification partition based on FLDF.

Remark : R_1^* given above estimated from learning sample \mathcal{L} of preclassified examples.

Example 2

$$\pi_1: N_p(\underline{\mu}_1, \Sigma_1)$$

$$\pi_2: N_p(\underline{\mu}_2, \Sigma_2); \Sigma_i > 0$$

$$R_1^* = \{ \underline{x} : p_1 f_1(\underline{x}) C(2|1) \geq p_2 f_2(\underline{x}) C(1|2) \}$$

Proceeding as in example 1, we get

$$p_1 f_1(\underline{x}) C(2|1) \geq p_2 f_2(\underline{x}) C(1|2)$$

\Leftrightarrow

$$-\frac{1}{2} \left((\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \right) - \frac{1}{2} \log \left(\frac{|I_1|}{|\Sigma_2|} \right) \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right)$$

$$\Leftrightarrow -\frac{1}{2} \left(\underline{x}' (\Sigma_1^{-1} - \Sigma_2^{-1}) \underline{x} \right) + (\underline{\mu}_1' \Sigma_1^{-1} - \underline{\mu}_2' \Sigma_2^{-1}) \underline{x} \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) + \frac{1}{2} \log \left(\frac{|I_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2)$$

The L.H.S. in above expression is
indep of \underline{x}

Called the quadratic discriminant f^n (for obvious reason).

Assignment rule:

For a new \underline{x}_0 ;

Assign \underline{x}_0 to π_1 if

$$-\frac{1}{2} \underline{x}_0' (\Sigma_1^{-1} - \Sigma_2^{-1}) \underline{x}_0 + (\underline{\mu}_1' \Sigma_1^{-1} - \underline{\mu}_2' \Sigma_2^{-1}) \underline{x}_0$$

$$\geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) + \frac{1}{2} \log \left(\frac{|I_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2)$$

Remark: Note that we have to use the learning sample \mathcal{L} to have sample counterpart $\{\hat{R}_1, \hat{R}_2\}$

If in \mathcal{L} there are n_1 preclassified obsns from Π_1 & n_2 preclassified cases from Π_2 ($n_1 + n_2 = n$ ← total preclassified cases in \mathcal{L}), then use n_1 obsns of preclassified Π_1 examples to calculate $\hat{\mu}_1 = \bar{x}_{(1)}$ & $\hat{\Sigma}_1 = S_{(1)}$.
 Similarly $\hat{\mu}_2 = \bar{x}_{(2)}$ & $\hat{\Sigma}_2 = S_{(2)}$ from n_2 preclassified cases from Π_2 .

Note that in this case no pooling is required as would be necessary in example for estimation of common covariance matrix.

Example 3 : Discrete distⁿ case

(III)

π_1 , p.m.f.

	x_1	0	1	2
x_2				
0		.1	.05	.15

	x_1	0	1	2
x_2				
1		.25	.2	.25

Prior prob

$$p_1 = 0.4$$

π_2 , p.m.f

	x_1	0	1	2
x_2				
0		.2	.2	.2

	x_1	0	1	2
x_2				
1		.2	.1	.1

Prior prob

$$p_2 = 0.6$$

Misclassification cost

$$C(1|2) = 5 \quad C(2|1) = 10$$

$$C(1|1) = C(2|2) = 0 \text{ (always)}$$

possible pairs $(0,0), (0,1), (1,0), (1,1), (2,0), (2,1)$.

Check the assignment

e.g. $(0,0) = \underline{x}$

$$p_1 f_1(\underline{x}) C(2|1) = 0.4 \times 0.1 \times 10 = 0.4$$

$$\& p_2 f_2(\underline{x}) C(1|2) = 0.6 \times 0.2 \times 5 = 0.6$$

Recall that

$$R_1^* = \{\underline{x} : p_1 f_1(\underline{x}) C(2|1) \geq p_2 f_2(\underline{x}) C(1|2)\}$$

$$\& R_2^* = \{\underline{x} : \text{,,} < \text{,,}\}$$

Hence assign rule for $(0,0)$ is π_2

Similarly assignment rules for all other possible (x_1, x_2) can be computed.

Remark: For all the 3 examples ECM (or TPM) of the ECM minimizing classification rules can be calculated once we calculate $P(1|2)$ & $P(2|1)$

e.g.: Consider "Example 1" of $\Pi_1: N_p(\underline{\mu}_1, \Sigma_1)$

$$\Pi_2: N_p(\underline{\mu}_2, \Sigma_2)$$

$$P(1|2) = P_{\Pi_2} \left(\tilde{x} \in R_1^* \right).$$

$$= P_{\Pi_2} \left((\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \tilde{x} \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) + \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right)$$

Note that under Π_2 :

$$\begin{aligned} z_1 = (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \tilde{x} &\sim N_1 \left((\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2, \right. \\ &\quad \left. \underbrace{(\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 - \underline{\mu}_2)}_{\Delta^2} \right) \\ \Rightarrow P(1|2) &= P_{\Pi_2} \left(\frac{z_1 - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2}{\Delta} \geq \left[\frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right. \right. \\ &\quad \left. \left. - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2 \right. \right. \\ &\quad \left. \left. + \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \right] \Delta^{-1} \right) \\ \text{let } y &= \frac{z_1 - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2}{\Delta} \end{aligned}$$

$$P(1|2) = P_{\Pi_2} \left(y \geq \frac{1}{2} \Delta + \frac{1}{\Delta} \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \right)$$

$$\Rightarrow P(1|2) = 1 - \bar{\Phi} \left(\frac{\Delta}{2} + \frac{1}{\Delta} \log \left(\frac{P_2 C(1|2)}{P_1 C(2|1)} \right) \right)$$

∴ $P(2|1)$ can be calculated

Note: TPM / ECM can be estimated by estimating Δ from learning set L .

Remark: Performance measure for classifier

Apparent error rate (APER)

- Calculate the "confusion matrix"

		Predicted class		
		π_1	π_2	
Actual class	π_1	n_{1c}	n_{1M}	$(n_{1c} + n_{1M} = n_1)$
	π_2	n_{2M}	n_{2c}	$(n_{2M} + n_{2c} = n_2)$

n_{1M} : # of cases from π_1 misclassified by

n_{1c} : # of - - - π_1 correctly classified

$$\text{APER} = \frac{n_{1M} + n_{2M}}{n_1 + n_2}$$

Classification in multiclass problem

Let $\pi_1, \pi_2, \dots, \pi_c$ be C classes with prior probabilities $p(\pi_1), p(\pi_2), \dots, p(\pi_c)$, respectively.

Bayes classifier

Assign \tilde{x} to π_j if the posterior prob

of the class π_j given the obsn \tilde{x} , i.e.,

$p(\pi_j | \tilde{x})$ is the highest over all classes

π_1, \dots, π_c

i.e. Assign \tilde{x} to π_j if

$$p(\pi_j | \tilde{x}) \geq p(\pi_k | \tilde{x}) ; k=1, \dots, c \\ k \neq j$$

$$\Rightarrow \frac{p(\pi_j) p(\tilde{x} | \pi_j)}{\sum_{i=1}^c p(\pi_i) p(\tilde{x} | \pi_i)} \geq \frac{p(\pi_k) p(\tilde{x} | \pi_k)}{\sum_{i=1}^c p(\pi_i) p(\tilde{x} | \pi_i)}$$

$$\Rightarrow p(\pi_j) p(\tilde{x} | \pi_j) \geq p(\pi_k) p(\tilde{x} | \pi_k)$$

\Rightarrow as in the 2-class problem (done in class), Bayes classifier for a $C (\geq 2)$ -class problem can be expressed in terms of prior ~~and~~ prob & class conditional densities

TPM minimizing classifier

TPM for a general c-class problem can be defined as

$$\text{TPM} = \sum_{i=1}^c p(\pi_i) P(\text{error} | \pi_i)$$

↑
 misclassification
 i.e. $P(\text{error} | \pi_i)$: prob of
 misclassifying a pattern
 from π_i

Let $\{R_1, R_2, \dots, R_c\}$ be the classification partition, then

$$P(\text{error} | \pi_i) = \int_{R_i^c} f(\underline{x} | \pi_i) d\underline{x}$$

$$\Rightarrow \text{TPM} = \sum_{i=1}^c \int_{x-R_i} f(\underline{x} | \pi_i) p(\pi_i) d\underline{x}$$

$$= \sum_{i=1}^c p(\pi_i) \int_{x-R_i} f(\underline{x} | \pi_i) d\underline{x}$$

$$= \sum_{i=1}^c p(\pi_i) \left(\int_{\underline{x}} f(\underline{x} | \pi_i) d\underline{x} - \int_{R_i} f(\underline{x} | \pi_i) d\underline{x} \right)$$

$$= \sum_{i=1}^c p(\pi_i) \left(1 - \int_{R_i} f(\underline{x} | \pi_i) d\underline{x} \right)$$

i.e.

$$\text{TPM} = 1 - \sum_{i=1}^C \int_{R_i} p(\pi_i) f(\underline{x} | \pi_i) d\underline{x}$$

\Rightarrow Minimizing the TPM w.r.t. partition $\{R_1, \dots, R_C\}$ is equivalent to maximizing

$$\sum_{i=1}^C \int_{R_i} p(\pi_i) f(\underline{x} | \pi_i) d\underline{x} \text{ w.r.t. } \{R_1, \dots, R_C\}$$

$\xleftarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}}$

This can be viewed as prob of correct classification

Maximization of the above is achieved by selecting R_i to be the region for which

$p(\pi_i) f(\underline{x} | \pi_i)$ is the largest among all classes

(i.e. among all $i=1 \dots C$)

i.e. the TPM minimizing classification rule

is

Assign \underline{x} to π_i if

$$p(\pi_i) f(\underline{x} | \pi_i) \geq p(\pi_k) f(\underline{x} | \pi_k)$$

Remark: As in the 2-class problem, the TPM minimizing rule is equivalent to the one obtained by maximizing the "posterior prob", i.e. the Bayes classifier

Minimum ECM classification rule

Let $p(\pi_i) = p_i \quad i = 1(1)c$

class conditional densities $f(\tilde{x} | \pi_i); i = 1(1)c$

$c(k|i)$: cost of misclassifying an item/pattern
from π_i to $\pi_k; k, i = 1(1)c$

$c(i|i) = 0 \quad i = 1(1)c$

$\{R_1, \dots, R_c\}$: classification partition

$$P(k|i) = P(\text{Classifying an item to } \pi_k | \pi_i) \\ = \int_{R_k} f(\tilde{x} | \pi_i) d\tilde{x}$$

$$P(i|i) = \int_{R_i} f(\tilde{x} | \pi_i) d\tilde{x} = \int_{\tilde{x} - \cup_{k \neq i} R_k} f(\tilde{x} | \pi_i) d\tilde{x} \\ = 1 - \sum_{\substack{k=1 \\ k \neq i}}^c P(k|i)$$

Note that the conditional ECM of \tilde{x} from π_1 ,
into π_2 or π_3 or ... or π_c is

$$\text{ECM}_1 = P(2|1)C(2|1) + P(3|1)C(3|1) + \dots + P(c|1)C(c|1) \\ = \sum_{k=2}^c P(k|1)C(k|1)$$

The above ECM, is with prob $p(\pi_i) = p_i$,

Thus for other conditional ECM's ECM_i

We can write similar expressions w.p. p_i
(as below)

\Rightarrow

$$\text{ECM} = p_1 \text{ECM}_1 + p_2 \text{ECM}_2 + \dots + p_c \text{ECM}_c$$

$$\text{i.e. } \text{ECM} = p_1 \left(\sum_{k=2}^c p(k|1) c(k|1) \right)$$

$$+ p_2 \left(\sum_{\substack{k=1 \\ k \neq 2}}^c p(k|2) c(k|2) \right) + \dots$$

$$\dots + p_c \left(\sum_{k=1}^{c-1} p(k|c) c(k|c) \right)$$

$$\text{i.e. } \text{ECM} = \sum_{i=1}^c p_i \sum_{\substack{k=1 \\ k \neq i}}^c p(k|i) c(k|i)$$

$$= \sum_{k=1}^c \sum_{\substack{i=1 \\ i \neq k}}^c p_i p(k|i) c(k|i)$$

$$= \sum_{k=1}^c \sum_{\substack{i=1 \\ i \neq k}}^c \int_{R_K} p_i c(k|i) f(x| \pi_i) dx$$

$$\text{i.e. } \text{ECM} = \sum_{K=1}^C \int \left(\sum_{\substack{i=1 \\ i \neq K}}^C p_i c(k|i) f(\underline{x}) \pi_i \right) d\underline{x}$$

$$= \sum_{K=1}^C \int_{R_K} h_K(\underline{x}) d\underline{x}$$

Note that

$$\sum_{K=1}^C \int_{R_K} (h_K(\underline{x}) - \min_j h_j(\underline{x})) d\underline{x} \geq 0$$

with equality only if

$$h_K(\underline{x}) = \min_j h_j(\underline{x}) \quad \underline{\text{if } \underline{x} \text{ in } R_K}$$

Thus the ECM minimizing classification rule
is

Assign \underline{x} to π_K if

$$\sum_{\substack{i=1 \\ i \neq K}}^C p_i c(k|i) f(\underline{x}) \pi_i$$

is minimum among all such C expressions

Remark: Under equal cost setup, i.e. all misclassification costs are same

we assign \underline{x} to π_K if

$$\sum_{\substack{i=1 \\ i \neq K}}^C p_i f(\underline{x}) \pi_i$$

is the smallest

$$\text{i.e. } \sum_{\substack{i=1 \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) \leq \sum_{\substack{i=1 \\ i \neq j \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) \quad (*)$$

Subtracting $\sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i)$ from both sides

of (*), we get

$$\begin{aligned} & \sum_{\substack{i=1 \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) - \sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i) \\ & \leq \sum_{\substack{i=1 \\ i \neq j}}^c p_i f(\underline{x} | \pi_i) - \sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i) \end{aligned}$$

$$\text{i.e. } p_j f(\underline{x} | \pi_j) \leq p_k f(\underline{x} | \pi_k) \quad \forall j \neq k$$

i.e. the rule is

assign \underline{x} to π_k if

$p_k f(\underline{x} | \pi_k)$ is largest

This is the same as TPM minimizing rule
(What we expect under equal misclassification cost).

Further, the above is (under equal cost)
equivalent to Bayes classifier.

Multiclass classification problem

Example 1: 3-class problem

Consider the following cost-prior table

		True membership			Cost matrix ↳ mis class →
		π_1	π_2	π_3	
prior prob	π_1	$c(1 1) = 0$	$c(1 2) = 500$	100	
	π_2	10	0	50	
	π_3	50	200	0	
prior prob		$p_1 = 0.05$	$p_2 = 0.60$	$p_3 = 0.35$	

Let \tilde{x}_0 be a new obsn \Rightarrow

$$f(\tilde{x}_0 | \pi_1) = 0.01 ; f(\tilde{x}_0 | \pi_2) = 0.85 \\ f(\tilde{x}_0 | \pi_3) = 2$$

ECM minimizing classifier

Compute $\sum_{\substack{i=1 \\ i \neq k}}^3 p_i c(k|i) f(\tilde{x}_0 | \pi_i) \quad \forall k = 1, 2, 3$

$$\begin{aligned} K=1 : \quad & \sum_{i=2}^3 p_i c(k|i) f(\tilde{x}_0 | \pi_i) \\ & = p_2 c(1|2) f(\tilde{x}_0 | \pi_2) + p_3 f(\tilde{x}_0 | \pi_3) c(1|3) \\ & = 325 \end{aligned}$$

$$\underline{K=2} : \sum_{\substack{i=1 \\ i \neq 2}}^3 p_i c(2|i) f(\tilde{x}_0 | \pi_i)$$

$$= p_1 c(2|1) f(\tilde{x}_0 | \pi_1) + p_3 c(2|3) f(\tilde{x}_0 | \pi_3)$$

$$= 35.06$$

$$\underline{K=3} : \sum_{i=1}^2 p_i c(3|i) f(\tilde{x}_0 | \pi_i) = 102.03$$

Since $\sum_{\substack{i=1 \\ i \neq 2}}^3 p_i c(2|i) f(\tilde{x}_0 | \pi_i)$ is the smallest

we assign \tilde{x}_0 to π_2

Note: Suppose in the same example, we take equal cost \Rightarrow TPM minimizing rule

$$\text{Calculate } p_1 f(\tilde{x}_0 | \pi_1) = 0.0005$$

$$p_2 f(\tilde{x}_0 | \pi_2) = 0.510$$

$$p_3 f(\tilde{x}_0 | \pi_3) = 0.7$$

Since $p_3 f(\tilde{x}_0 | \pi_3) > p_i f(\tilde{x}_0 | \pi_i) \quad i=1, 2$

\tilde{x}_0 is allocated to π_3

Bayes classifier would also have done the same assignment

Example 3 : Discrete populations

Consider 3 bivariate discrete populations with the following jt p.m.f.s :

		Π_1		Π_2			Π_3		
		x_1	x_2	1	2	1	2	1	2
1	1	0.5	0.2	0.2	0.1	0.3	0.4	0.25	0.25
	2	0.1	0.2			0.3	0.4	0.25	0.25

Prior probabilities are equal, i.e. $p(\Pi_i) = \frac{1}{3}; i=1,2,3$
 $(= p_i)$

Misclassification costs are

$$c(1|i) = 1; i=2,3$$

$$c(2|i) = 2; i=1,3$$

$$c(3|i) = 3; i=1,2$$

Let us denote by $f_i(x)$ to be class conditional prob for class $\Pi_i; i=1,2,3$.

ECM classification rule is :

Assign \tilde{x} to Π_k if

$$\sum_{\substack{i=1 \\ i \neq k}}^3 p_i f_i(\tilde{x}) c(k|i) \text{ is smallest}$$

Based on the above find rule for all pairs \tilde{x}

pair (1,1)

K=1

$$\sum_{i=2}^3 p_i f_i(\underline{x}) \underset{\substack{\downarrow \\ \underline{x} = (1,1)}}{c(1|i)}$$

$$= p_2 f_2(\underline{x}) c(1|2) + p_3 f_3(\underline{x}) c(1|3) \\ = \frac{1}{3} (f_2(\underline{x}) + f_3(\underline{x}))$$

$$= \frac{1}{3} (0.2 + 0.25) = 0.15$$

K=2

$$\sum_{\substack{i=1 \\ i \neq 2}}^3 p_i f_i(\underline{x}) \underset{\substack{\downarrow \\ \underline{x} = (1,1)}}{c(2|i)}$$

$$= p_1 f_1(\underline{x}) c(2|1) + p_3 f_3(\underline{x}) c(2|3)$$

$$= \frac{2}{3} (f_1(\underline{x}) + f_3(\underline{x})) = \frac{2}{3} (0.5 + 0.25) \\ = 0.5$$

K=3

$$\sum_{i=1}^2 p_i f_i(\underline{x}) c(3|i)$$

$$= p_1 f_1(\underline{x}) c(3|1) + p_2 f_2(\underline{x}) c(3|2)$$

$$= \frac{3}{3} (f_1(\underline{x}) + f_2(\underline{x})) = 0.7$$

$$\sum_{\substack{i=1 \\ i \neq 1}}^3 p_i f_i(\underline{x}) c(1|i) \text{ is smallest}$$

\Rightarrow Assign $\underline{x} \overset{(1,1)}{\in} \Pi_1$

i.e. $(1,1) \in R_1$ partition

Pair $(1, 2)$

$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(1, 2) + f_3(1, 2)) = \frac{1}{3} (0.55) \leftarrow \text{smallest}$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(1, 2) + f_3(1, 2)) = \frac{2}{3} (0.35)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(1, 2) + f_2(1, 2)) = 0.4$$

$\underline{k=1}$ gives smallest

$$\Rightarrow (1, 2) \in R_1, \text{ partition}$$

Pair $(2, 1)$

$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(2, 1) + f_3(2, 1)) = \frac{1}{3} (0.35)$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(2, 1) + f_3(2, 1)) = \frac{2}{3} (0.45)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(2, 1) + f_2(2, 1)) = 0.3$$

$\underline{k=1}$ gives smallest

$$\Rightarrow (2, 1) \in R_1, \text{ partition}$$

Pair $(2, 2)$

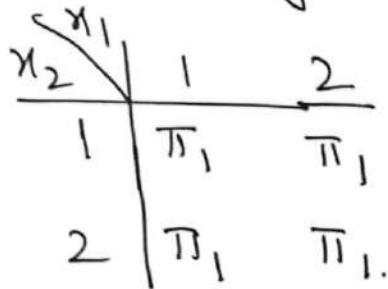
$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(2, 2) + f_3(2, 2)) = \frac{1}{3} (0.65)$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(2, 2) + f_3(2, 2)) = \frac{2}{3} (0.45)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(2, 2) + f_2(2, 2)) = 0.6$$

$$\Rightarrow (2, 2) \in R_1$$

\Rightarrow ECM minimizing classification rule is



Suppose we have a set \mathcal{D} of preclassified examples

$$\mathcal{D} = \{((1, 2), \pi_1), ((1, 1), \pi_1), ((1, 1), \bar{\pi}_1), ((1, 1), \pi_3), ((2, 2), \pi_2), ((2, 1), \pi_3), ((1, 2), \bar{\pi}_3), ((2, 2), \bar{\pi}_3) \}$$

If we apply ECM rule on \mathcal{D} , we get

	Coming from		
	π_1	π_2	π_3
$\bar{\pi}_1$	3	1	4
$\bar{\pi}_2$	0	0	0
$\bar{\pi}_3$	0	0	0

Misclassification rate would be $\frac{5}{8} !!$

↑
not good at all!

(i) Try to find the TPM opt rule

(ii) Try to calculate ECM/TPM of the two optimum rule

Example 2 : Multiclass Gaussian

$$\pi_i \equiv N_p(\mu_i, \Sigma_i)$$

Assume for simplicity that $C(k|i)$'s are all same

We know that under this equal cost, the rule is

Assign \tilde{x} to π_k if

$$p_k f(\tilde{x} | \pi_k) = \max_i p_i f(\tilde{x} | \pi_i)$$

$$\text{i.e. if } \log(p_k f(\tilde{x} | \pi_k)) = \max_i \log(p_i f(\tilde{x} | \pi_i))$$

$$\text{with } f(\tilde{x} | \pi_k) = (2\pi)^{-p/2} |\Sigma_k|^{-1/2} \exp\left(-\frac{1}{2}(\tilde{x} - \mu_k)' \Sigma_k^{-1} (\tilde{x} - \mu_k)\right)$$

Define quadratic discriminant score for the i th popn, π_i , as

$$S_i^{QDS}(\tilde{x}) = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\tilde{x} - \mu_i)' \Sigma_i^{-1} (\tilde{x} - \mu_i) + \log p_i$$

Classification Rule

Assign \tilde{x} to π_k if

$$S_K^{QDS}(\tilde{x}) = \max_i S_i^{QDS}(\tilde{x})$$

124

Note: With a given learning sample L of preclassified examples, we calculate the quadratic discriminant score for each class.

$$\hat{g}_i^{QDS}(x) = -\frac{1}{2} \log |\mathcal{S}_i| - \frac{1}{2} (x - \bar{x}_i)^T \mathcal{S}_i^{-1} (x - \bar{x}_i) + \log b_i$$

S_i : Sample variance-covariance matrix based on pre-classified examples corresponding to class Π_i in d

\bar{x}_i : sample mean vector based on π_i in d.

Remark: Sp case: suppose π_i is $N_p(\mu_i, \Sigma)$ $i=1(1)3$

$$\log p_i f_i(\underline{x}) = \underbrace{-\frac{1}{2} |\Sigma| - \frac{1}{2} \underline{x}' \Sigma^{-1} \underline{x}}_{\text{Same as pop's}} + \underline{\mu}_i' \Sigma^{-1} \underline{x} - \frac{1}{2} \underline{\mu}_i' \Sigma^{-1} \underline{\mu}_i + \log p_i$$

Define the linear discriminant score as

$$f_i^{LDS}(x) = \bar{u}_i' \bar{\Sigma}^{-1} x - \frac{1}{2} \bar{u}_i' \bar{\Sigma}^{-1} \bar{u}_i + \log p_i$$

Assignment rule is:

Assign \approx to π_k if

$$f_k^{\text{LDS}}(x) = \max_i f_i^{\text{LDS}}(x)$$

We have to use λ to get the following estimated LDS

$$\hat{f}_i^{\text{LDS}}(\underline{x}) = \underline{\bar{x}}_i' \hat{\Sigma}^{-1} \underline{x} - \frac{1}{2} \underline{\bar{x}}_i' \hat{\Sigma}^{-1} \underline{\bar{x}}_i + \log p_i$$

$\underline{\bar{x}}_i$ for $i=1(1)c$ calculated from pre-classified examples of π_i class in λ .

$$\lambda \left(\sum_{i=1}^c n_i - c \right) \hat{\Sigma} = \sum_{i=1}^c (n_i - 1) \hat{\Sigma}_i$$

$\hat{\Sigma}_i$ for $i=1(1)c$ calculated from pre-classified examples of π_i class in λ .
 (divisor $(n_i - 1)$ sample covariance matrix)

Estimated classification rule is :

Assign \underline{x} to π_K if

$$\hat{f}_K^{\text{LDS}}(\underline{x}) = \max_i \hat{f}_i^{\text{LDS}}(\underline{x})$$

Note: The above rule is equivalent to

Assign \underline{x} to π_K if

$$-\frac{1}{2} (\underline{x} - \underline{\bar{x}}_K)' \hat{\Sigma}^{-1} (\underline{x} - \underline{\bar{x}}_K) + \log p_K$$

i.e assign \underline{x} to p_K closest in terms of sqd Mahalanobis distance penalized by $\log p_K$.

Note: We may also need to estimate p_i 's from λ .

Nearest-neighbor classifier

It's a non-parametric classifier approach.

NN classifier uses observations in \mathcal{X} closest to the given \underline{x} , closest in the feature vector space.

Let $N_K(\underline{x})$ denote the neighborhood of \underline{x} defined by the K closest points \underline{x}_i 's in the training set.

Obtain a "majority vote rule" for the K nearest neighbors, i.e. pts within $N_K(\underline{x})$; voting w.r.t. their class membership. The class which is the winner (having maximum counts for ~~the~~ cases inside $N_K(\underline{x})$) is the assigned class.

Logistic discrimination model

Logistic regression based classification model

As opposed to the previous methods (TPM, ECM), this approach does not require strong parametric assumption regarding class conditional densities ($f(\tilde{x}|\pi_k)$).

Two-class problem

Let us first consider a 2-class classification problem

Suppose the 2 classes are π_1 & π_2

Logistic discrimination model assumes that the class conditional densities of the 2 classes satisfy

$$\log \left(\frac{f(\tilde{x}|\pi_1)}{f(\tilde{x}|\pi_2)} \right) = \beta_0 + \beta^T \tilde{x}$$

Where; $\tilde{x} \in \mathbb{R}^p$ is the feature vector

$f(\tilde{x}|\pi_i)$'s are class conditional densities

$\beta_0, \beta_{\tilde{x}}$: unknown deterministic constants

Note that

$\log \left(\frac{f(\underline{x} | \pi_1)}{f(\underline{x} | \pi_2)} \right)$ is referred to as the log-odds ratio.

Now

$$\log \left(\frac{f(\underline{x} | \pi_1)}{f(\underline{x} | \pi_2)} \right) = \beta_0 + \beta' \underline{x}$$

$$\Leftrightarrow \log \left(\frac{\frac{p(\pi_1 | \underline{x}) f(\underline{x})}{p(\pi_1)}}{\frac{p(\pi_2 | \underline{x}) f(\underline{x})}{p(\pi_2)}} \right) = \beta_0 + \beta' \underline{x}$$

$$\Leftrightarrow \log \left(\frac{p(\pi_2) p(\pi_1 | \underline{x})}{p(\pi_1) p(\pi_2 | \underline{x})} \right) = \beta_0 + \beta' \underline{x}$$

$$\Leftrightarrow \frac{p(\pi_2)}{p(\pi_1)} \cdot \frac{p(\pi_1 | \underline{x})}{p(\pi_2 | \underline{x})} = e^{\beta_0 + \beta' \underline{x}}$$

$$\Leftrightarrow \frac{p(\pi_1 | \underline{x})}{p(\pi_2 | \underline{x})} = \frac{p(\pi_1)}{p(\pi_2)} e^{\beta_0 + \beta' \underline{x}}$$

$$\Rightarrow \frac{1 - p(\pi_2 | \underline{x})}{p(\pi_2 | \underline{x})} = e^{\log \frac{p(\pi_1)}{p(\pi_2)} + \beta_0 + \beta' \underline{x}} \\ (\because p(\pi_1 | \underline{x}) = 1 - p(\pi_2 | \underline{x}))$$

$$\frac{1 - p(\pi_2 | \underline{x})}{p(\pi_2 | \underline{x})} = e^{\beta_0^* + \beta' \underline{x}} \\ (\beta_0^* = \beta_0 + \log \frac{p(\pi_1)}{p(\pi_2)})$$

$$\Rightarrow p(\pi_2 | \underline{x}) = \frac{1}{1 + e^{\beta_0^* + \beta' \underline{x}}}$$

and accordingly

$$p(\pi_1 | \underline{x}) = \frac{e^{\beta_0^* + \beta' \underline{x}}}{1 + e^{\beta_0^* + \beta' \underline{x}}}$$

Discrimination betⁿ the 2-classes depend on the odds-ratio $p(\pi_1 | \underline{x}) / p(\pi_2 | \underline{x})$

The class assignment rule is

Assign \underline{x} to π_1 if $\frac{p(\pi_1 | \underline{x})}{p(\pi_2 | \underline{x})} \geq 1$

& assign \underline{x} to π_2 if $\frac{p(\pi_1 | \underline{x})}{p(\pi_2 | \underline{x})} < 1$

i.e. the assignment rule is

Assign \tilde{x} to π_1 , If $\beta_0^* + \beta' \tilde{x} \geq 0$

& assign \tilde{x} to π_2 If $\beta_0^* + \beta' \tilde{x} < 0$

Note: The logistic discrimination is thus

based on the value of $\beta_0^* + \beta' \tilde{x}$.

Remark: Note that when we talk about classification in the above framework based on the value of " $\beta_0^* + \beta' \tilde{x}$ ", \tilde{x} is the feature vector observed and we need to estimate/compute β_0^* & β' based on learning set. We would address this point after we look at the multiclass framework of logistic discrimination.

Logistic discrimination: multiclass

Suppose that we have C classes π_1, \dots, π_C

Assume that log-likelihood ratios for any pair (s, c)
 $s = 1(1)C-1$
 of likelihoods is linear of the form

$$\log \left(\frac{f(\underline{x} | \pi_s)}{f(\underline{x} | \pi_c)} \right) = \beta_{s0} + \beta_s' \underline{x} ; s = 1(1)C-1$$

i.e. the $C-1$ log-odds specifies the model
 for discrimination

Note that

$$\log \left(\frac{f(\underline{x} | \pi_s)}{f(\underline{x} | \pi_c)} \right) = \beta_{s0} + \beta_s' \underline{x} ; s = 1(1)C-1$$

$$\Rightarrow \log \left(\frac{p(\pi_c)}{p(\pi_s)} \cdot \frac{p(\pi_s | \underline{x})}{p(\pi_c | \underline{x})} \right) = \beta_{s0} + \beta_s' \underline{x} ; s = 1(1)C-1$$

$$\Rightarrow \frac{p(\pi_c)}{p(\pi_s)} \cdot \frac{p(\pi_s | \underline{x})}{p(\pi_c | \underline{x})} = e^{\beta_{s0} + \beta_s' \underline{x}} ; s = 1(1)C-1$$

$$\Rightarrow \frac{p(\pi_s | \underline{x})}{p(\pi_c | \underline{x})} = e^{\beta_{s0}^* + \beta_s' \underline{x}} ; s = 1(1)C-1$$

$$(\beta_{s0}^* = \beta_{s0} + \log \left(\frac{p(\pi_s)}{p(\pi_c)} \right))$$

$$\Leftrightarrow \frac{1}{p(\pi_c | \underline{x})} \sum_{s=1}^{c-1} p(\pi_s | \underline{x}) = \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}$$

$$\Leftrightarrow \frac{1 - p(\pi_c | \underline{x})}{p(\pi_c | \underline{x})} = \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}$$

$$\left(\sum_{s=1}^{c-1} p(\pi_s | \underline{x}) = 1 \Rightarrow \sum_{s=1}^{c-1} p(\pi_s | \underline{x}) \right)$$

$$\Rightarrow p(\pi_c | \underline{x}) = \frac{1}{1 + \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}} = 1 - p(\pi_c | \underline{x})$$

$$\& p(\pi_s | \underline{x}) = p(\pi_c | \underline{x}) \frac{e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}}{e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}}$$

$$\text{i.e. } p(\pi_s | \underline{x}) = \frac{e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}}{1 + \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{s1}' \underline{x}}}.$$

$$s=1 \dots c-1$$

The multiclass logistic discrimination rule based on the above set of posterior of the classes is given by:

Assign \underline{x} to π_j if

$$p(\pi_j | \underline{x}) = \max_i p(\pi_i | \underline{x})$$

i.e Assign \underline{x} to π_j if ($j = 1(1)c-1$)

$$\left\{ \begin{array}{l} \beta_{j0}^* + \beta_{\cdot j}^{\top} \underline{x} = \max_{(i=1(1)c-1)} \{ \beta_{i0}^* + \beta_{\cdot i}^{\top} \underline{x} \} - (i) \\ \text{and } \beta_{j0}^* + \beta_{\cdot j}^{\top} \underline{x} > 0 \\ \text{otherwise assign to } \pi_c \end{array} \right. - (ii)$$

Note that (i) comes from the fact that

$$p(\pi_s | \underline{x}) = \frac{e^{\beta_{s0}^* + \beta_{\cdot s}^{\top} \underline{x}}}{1 + \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{\cdot s}^{\top} \underline{x}}} \quad \forall s = 1(1)c-1$$

$$\text{So } p(\pi_j | \underline{x}) = \max_{i=1(1)c-1} p(\pi_i | \underline{x})$$

(\Rightarrow)

condition (i)

Further (ii) comes from the fact that

$$p(\pi_j | \underline{x}) > p(\pi_c | \underline{x})$$

$j = 1(1)c-1$

$$(\Rightarrow) \quad \frac{e^{\beta_{j0}^* + \beta_{\cdot j}^{\top} \underline{x}}}{1 + \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{\cdot s}^{\top} \underline{x}}} > \frac{1}{1 + \sum_{s=1}^{c-1} e^{\beta_{s0}^* + \beta_{\cdot s}^{\top} \underline{x}}}$$

$$(\Rightarrow) \quad e^{\beta_{j0}^* + \beta_{\cdot j}^{\top} \underline{x}} > 1$$

\Rightarrow $\beta_{j0}^* + \beta_{\cdot j}^{\top} \underline{x} > 0$ If this does not happen then assign \underline{x} to π_c

In other words assignment rule is:

assign \underline{x} to π_j if ($j=1 \dots c-1$)

$$\begin{cases} p(\pi_j | \underline{x}) > p(\pi_k | \underline{x}) & k = 1 \dots c-1 \\ & k \neq j \\ & \Delta p(\pi_j | \underline{x}) > p(\pi_c | \underline{x}) \text{ otherwise assign} \\ & \underline{x} \text{ to } \pi_c \end{cases}$$

\Rightarrow assign \underline{x} to π_j ($j=1 \dots c-1$) if

$$\beta_{j0}^* + \beta_{rj}^T \underline{x} = \max_{s=1 \dots c-1} \{ \beta_{s0}^* + \beta_{rs}^T \underline{x} \} \quad \&$$

$$\beta_{j0}^* + \beta_{rj}^T \underline{x} > 0 \text{ otherwise assign} \\ \underline{x} \text{ to } \pi_c$$

—

Parameter estimation : 2 class problem

Consider a binary variable y with values 1 and 0. "1" associated with π_1 population and "0" associated with π_2 population.

Learning sample data : $(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_n, y_n)$

$$P(\pi_1 | \underline{x}) = \frac{e^{\beta_0^* + \beta_1' \underline{x}}}{1 + e^{\beta_0^* + \beta_1' \underline{x}}} = \frac{e^{\beta^*' \underline{x}^*}}{1 + e^{\beta^*' \underline{x}^*}}$$

$$\underline{x}^* = \begin{pmatrix} 1 \\ \underline{x} \end{pmatrix}; \quad \beta^* = \begin{pmatrix} \beta_0^* \\ \beta_1 \end{pmatrix}$$

$$P(\pi_2 | \underline{x}) = \frac{1}{1 + e^{\beta^*' \underline{x}^*}}$$

$$f_{y_i}(y_i) = P(y=y_i) = (P(\pi_1 | \underline{x}))^{y_i} (1 - P(\pi_1 | \underline{x}))^{1-y_i}$$

y_i is 0 or 1

Likelihood f^n

$$L(\beta^*) = \prod_{i=1}^n f_{y_i}(y_i)$$

$$L(\beta^*) = \prod_{i=1}^n (P(\pi_1 | \underline{x}_i))^{y_i} (1 - P(\pi_1 | \underline{x}_i))^{1-y_i}$$

Log likelihood

$$\log L = \sum_{i=1}^n \left(y_i \log P(\pi_1 | \underline{x}_i) + (1-y_i) \log (1 - P(\pi_1 | \underline{x}_i)) \right)$$

$$= \sum_{i=1}^n \left(y_i \log \left(\frac{P(\pi_1 | \underline{x}_i)}{1 - P(\pi_1 | \underline{x}_i)} \right) + \log (1 - P(\pi_1 | \underline{x}_i)) \right)$$

$$\log L = \sum_{i=1}^n \left(y_i \beta^*' \tilde{x}_i^* - \log (1 + e^{\beta^*' \tilde{x}_i^*}) \right)$$

$$\begin{aligned} \frac{\partial \log L}{\partial \beta^*} &= \sum_{i=1}^n \left(y_i \tilde{x}_i^* - (1 + e^{\beta^*' \tilde{x}_i^*})^{-1} e^{\beta^*' \tilde{x}_i^*} \tilde{x}_i^* \right) \\ &= \sum_{i=1}^n \tilde{x}_i^* \left(y_i - \frac{e^{\beta^*' \tilde{x}_i^*}}{1 + e^{\beta^*' \tilde{x}_i^*}} \right) \end{aligned}$$

$$\frac{\partial \log L}{\partial \beta^*} = 0 \Rightarrow \sum_{i=1}^n \tilde{x}_i^* \left(y_i - \frac{e^{\beta^*' \tilde{x}_i^*}}{1 + e^{\beta^*' \tilde{x}_i^*}} \right) = 0 \quad - (*)$$

$(*)$ is a system of $p+1$ nonlinear equations which needs to be solved to obtain the MLEs. An iteratively ~~Reweighted Least Squares~~ Reweighted Least Squares (IRLS) method is used for solving $(*)$ to get the MLEs.

CART: Classification & Regression Tree

CART, is a powerful tree based method which can be used for regression modelling and for classification tasks.

- CART is a non-parametric approach and does not require any distributional assumption.

Before we talk about construction of such trees, let us first look at how it works.

Classification Trees

Classification tree (or a decision tree) is an example of a multistage decision process.

Rather than using the complete set of features jointly to make a output decision, different subsets of features are used at different levels of the tree

Consider the following classification tree for a 3-class problem with 6-dimensional feature vector $\underline{x} = (x_1, x_2, x_3, x_4, x_5, x_6)'$.

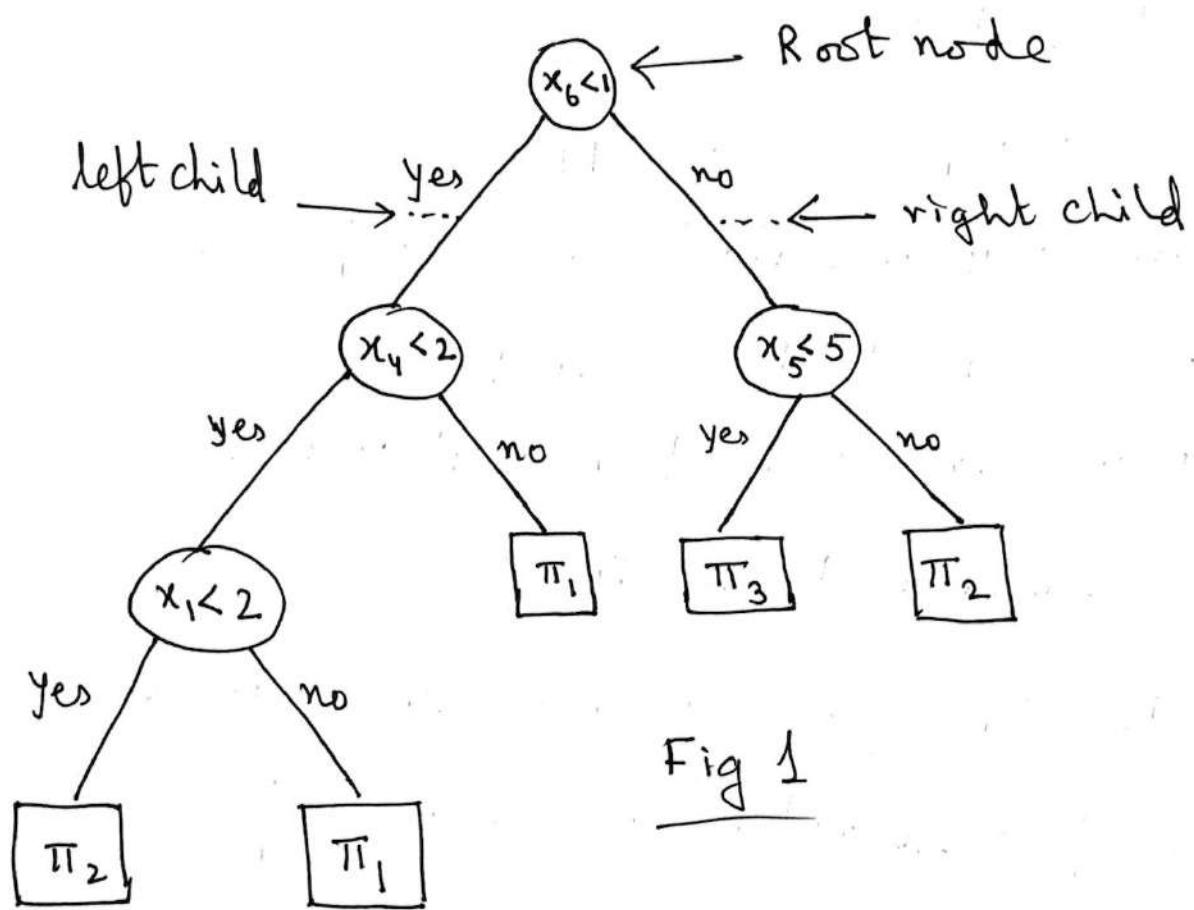


Fig 1

Explanation.

○ : These are internal nodes of the tree; associated with each internal node, there is a variable and a threshold

□ : These are called terminal node or leaf; there is a class label associated with each leaf

The top most internal node is the "root" of the tree

Remark: Construction of classification tree would revolve around

- (i) how to choose the variable associated with any internal node (including the root node)
- (ii) how to find the threshold associated with the variable at any internal node
- (iii) how to assign class label at the terminal nodes.
- (iv) when to declare a node ~~as~~ a terminal

How to use the classification tree?

Suppose we have a feature vector

$$\tilde{x} = (5, 3, 6, 2, 1, 3)' \text{ and}$$

wish to use classification tree of Fig 1 to classify the above feature vector.

S1: We begin with the root node; compare the value of x_6 (i.e. 3) with the threshold (which is 1) and find that threshold is exceeded and

hence take the right side path (called the ~~left~~
child of the node) right

S2: we arrive at the internal node $x_5 < 5$

We compare the value of the x_5 variable in the feature (i.e. 1) with the threshold at that node (i.e. 5) and find that the threshold is not exceeded; hence take the left side path (the left child).

S3: We ^{arrive} at ~~at~~ a terminal node with the last path direction. The terminal node that we have reached has a class label π_3 and hence the allotment of the given feature vector is π_3 , i.e. we assign $\underline{x} = (5, 3, 6, 2, 1, 3)'$ to π_3

The rule looks so simple!!
Isn't it??

Remark: The tree classifier in Fig 1 is an example of a binary decision tree.

Remark: Every internal node induces a partition of feature space. These partitions induced are hyperplanes parallel to the coordinate axes.

Remark: The classification partition is induced by the set of terminal nodes

Remark: For every internal node $t \in T$, there is a subspace $V(t) \in \mathcal{X}$. This $V(t)$ is the union of the subspaces of the terminal nodes that are its descendants (i.e. the terminal nodes below t in the tree structure).

Example: classification partitions

2-class problem: Π_1 & Π_2

2-dimensional feature vector

$$\underline{x} = (x_1, x_2)$$

Classification tree is given in Fig 2

Remark: The tree classifier in Fig 1 is an example of a "binary decision tree". More generally, the outcome of a decision at any node can also be more than 2.

Remark: Binary trees partition the feature space into 2 parts. The partitions induced (as in Fig 1) are hyperplanes parallel to the coordinate axes.

Let me try to illustrate the concept of partitions induced by such trees with help of a simple example

Example: 2-class problem: Π_1 & Π_2

2-dimensional feature vector

$$\underline{x} = (x_1, x_2)'$$

Let the constructed classification tree be given by Fig 2

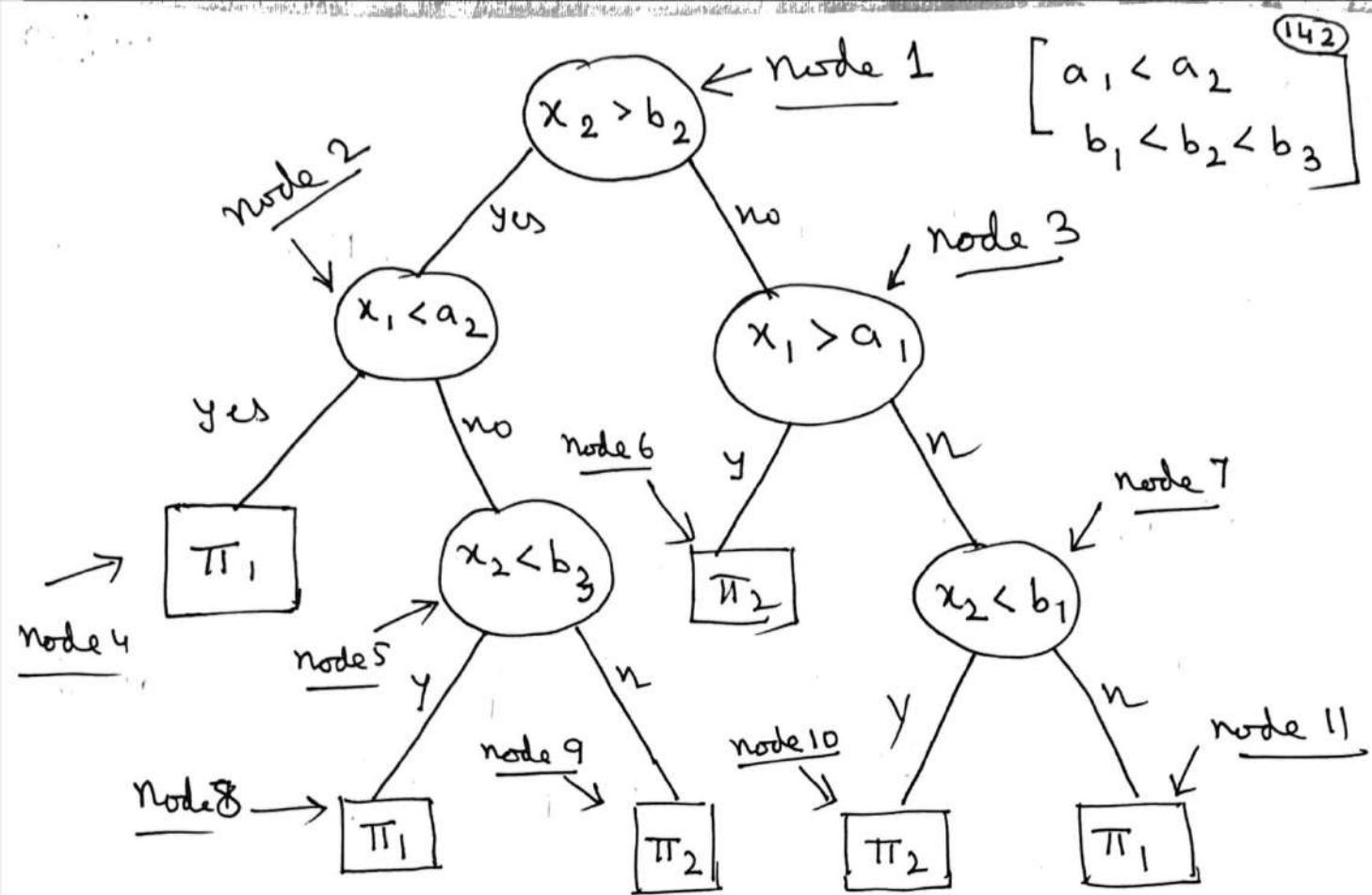


Fig 2.

To see how the above leads to partition of the feature space with class labels for partition regions follow the sequence of figures.

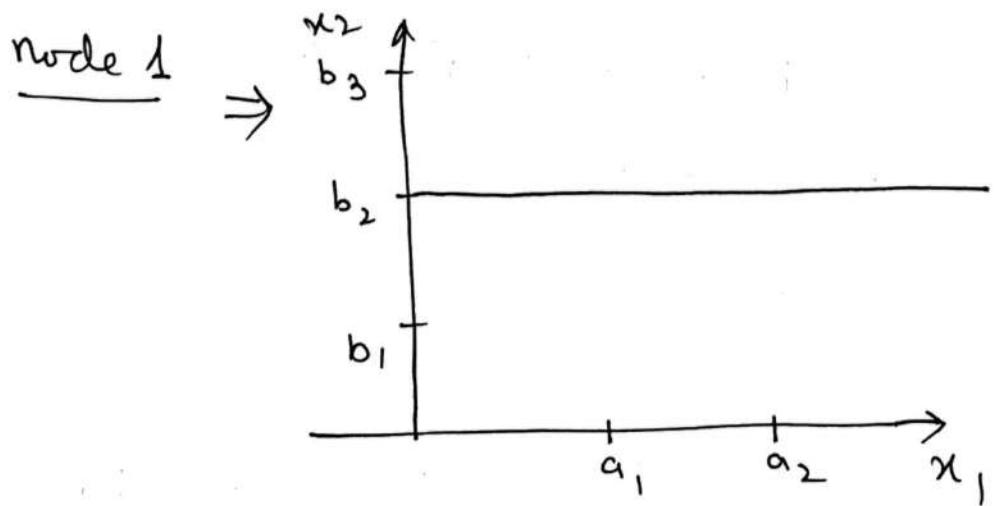


Fig 2.1

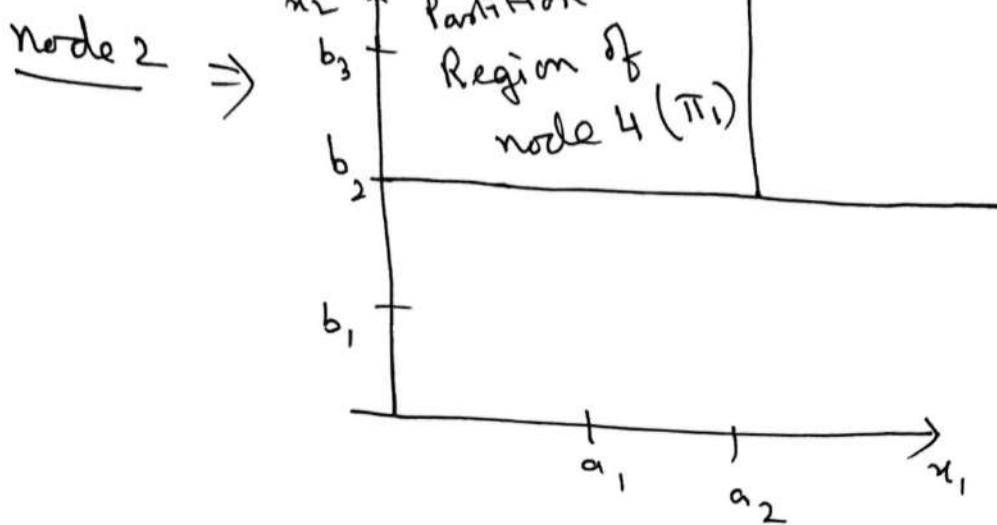


Fig 2.2

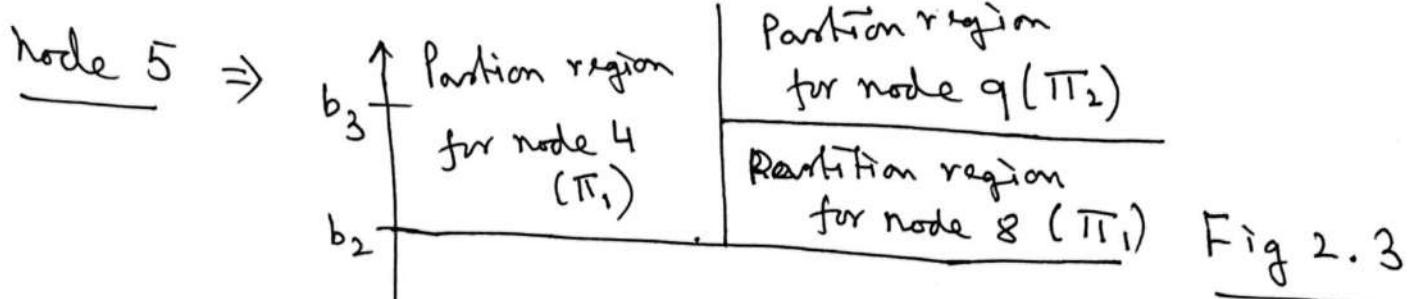


Fig 2.3

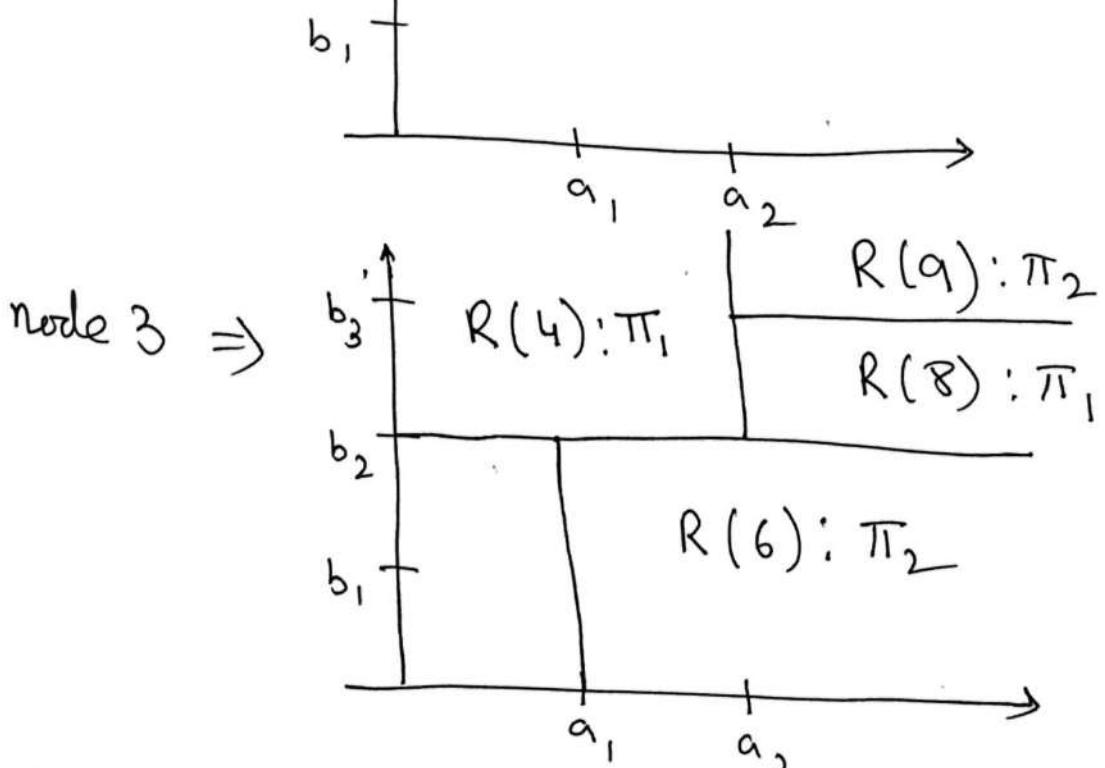


Fig 2.4

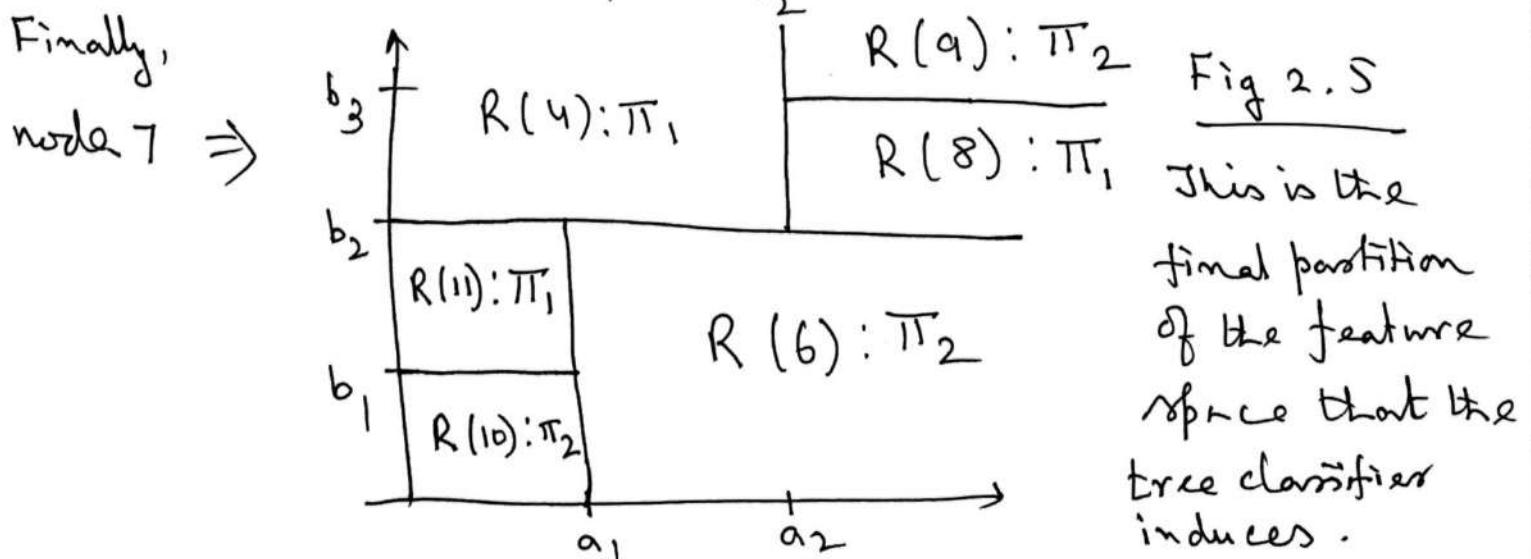
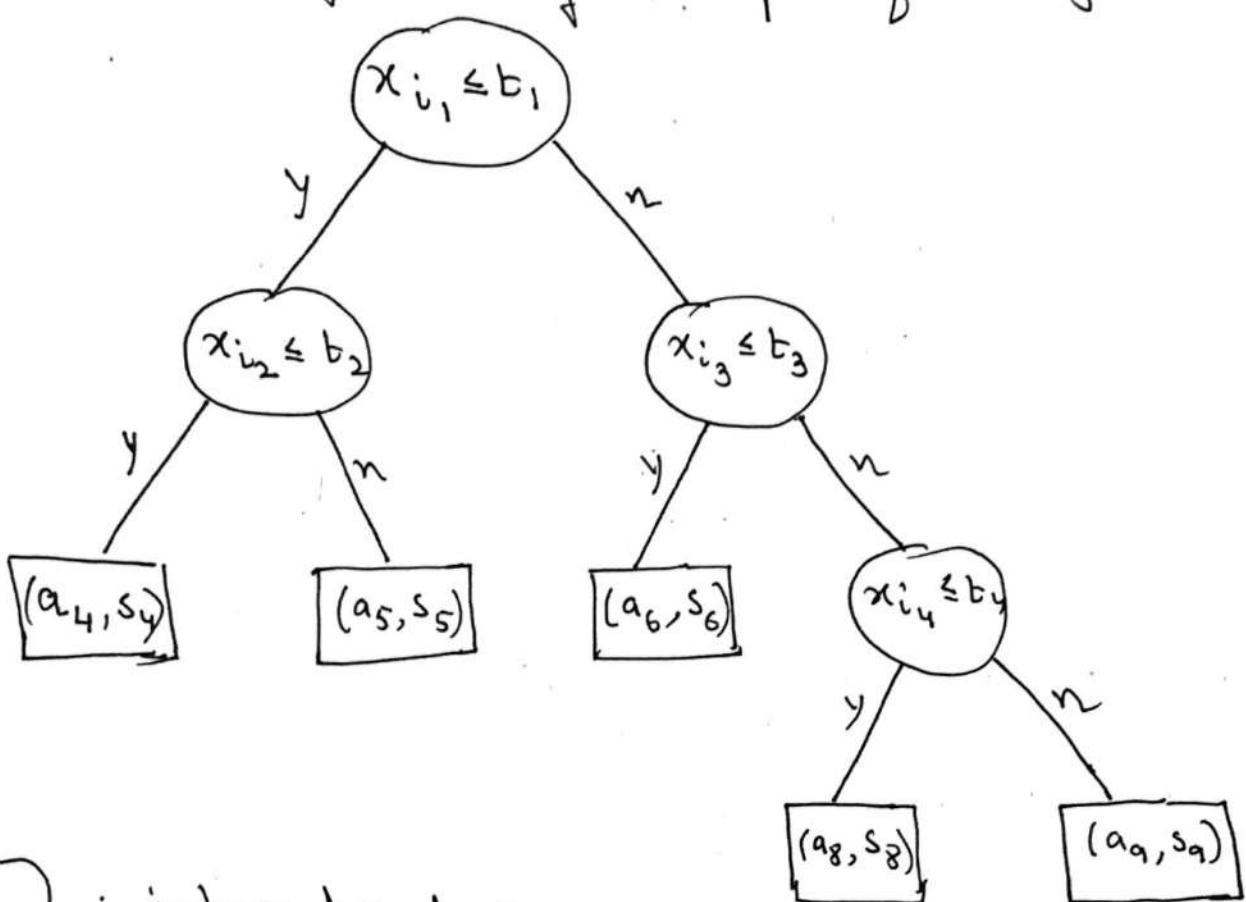


Fig 2.5

This is the final partition of the feature space that the tree classifier induces.

Regression Trees

Consider the following example of a regression tree



○ : internal nodes

: terminal nodes

i_j : j^{th} split variable

t_j : split point value

a_i : avg of all the response values of patterns in L
reaching terminal node i.

s_i : st.dev of \bar{x}_i

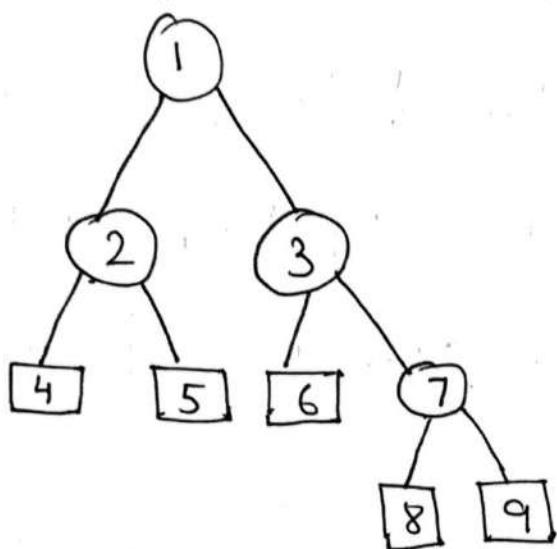
Remark: Regression trees look very much like
the same as classification trees, with
split variables & split points. The
terminal nodes are characterized by avg of

145

response variable values (in contrast to class labels as in classification trees) of patterns in the learning set \mathcal{L} that reaches a particular terminal node.

Remark: As in the classification trees, regression trees also induce a partition of the feature vector space. The predicted value of the response of a feature vector reaching a particular terminal node (i.e. belonging to the corresponding partition) would typically be the average of all the responses in the learning set (under sq error loss), whose feature vectors ~~of which~~ ~~partition~~ ~~partition~~ belongs to that partition.

Remark: Suppose the following is a classification or a regression tree



T : tree defined by the nodes $\{1, 2, \dots, 9\}$

\tilde{T} : set of terminal nodes $= \{4, 5, 6, 8, 9\}$

Partition induced by the tree is defined through the terminal nodes

$$U(4), U(5), U(6), U(8), U(9)$$

each $U(i) \in \mathcal{X} \leftarrow$ feature space

$$U(i) \cap U(j) = \emptyset$$

$$i \neq j ; i, j = 4, 5, 6, 8, 9$$

$$U(4) \cup U(5) \cup U(6) \cup U(8) \cup U(9) = \mathcal{X}$$

Towards tree construction

Some basic definitions

Tree: A tree is defined to be a set T of positive integers together with 2 functions $l(\cdot)$ and $r(\cdot)$ from T to $T \cup \{0\}$.

Each member of T corresponds to a node in the tree.

$l(t)$ and $r(t)$ are left node and right node values of t ($\forall t \in T$) \rightarrow

(i) $\forall t \in T$, either $l(t)$ and $r(t)$ are both 0 (in such a situation t is terminal) or $l(t)$ and $r(t)$ are both > 0 (internal node)

(ii) apart from the root node ($t=1$) there is a unique parent $s \in T \rightarrow \forall t \neq 1 \exists$ an $s \in T \rightarrow t = l(s)$ or $t = r(s)$.

Subtree

A subtree is a non-empty subset T_1 of T together with 2 fns $l_1(\cdot)$ and $r_1(\cdot) \Rightarrow$

$$l_1(t) = \begin{cases} l(t), & \text{if } l(t) \in T_1 \\ 0, & \text{of w} \end{cases}$$

$$\& \quad r_1(t) = \begin{cases} r(t), & \text{if } r(t) \in T_1 \\ 0, & \text{of w.} \end{cases}$$

and $T_1, l_1(\cdot)$ and $r_1(\cdot)$ form a tree.

Pruned subtree: A subtree that has same root as the original tree

Terminal node set:

\tilde{T} : set of terminal nodes.

$$\tilde{T} = \{t_1, t_2, \dots, t_M\}; t_i \in T$$

Partition induced by a tree T

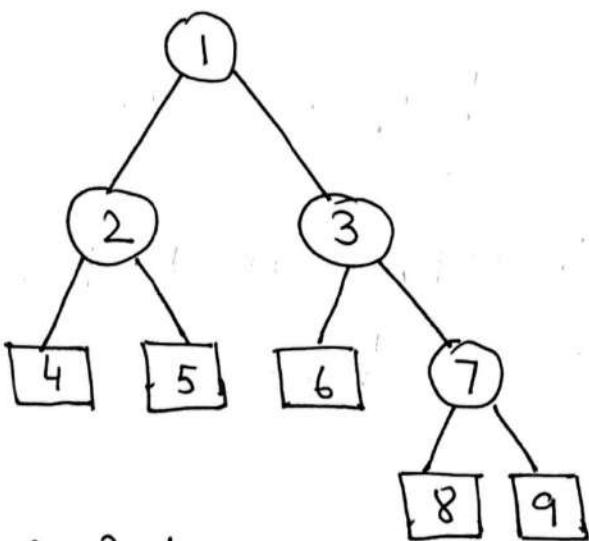
$\{U(t) : t \in \tilde{T}\}$: partition of feature space induced by T with terminal node set \tilde{T}

$$U(t) \cap U(s) = \emptyset \neq t \neq s; t, s \in \tilde{T}$$

$$\text{and } \bigcup_{t \in \tilde{T}} U(t) = \mathcal{X}$$

Note: Remember that associated with every non-terminal node t , there is a subspace $U(t)$ of \mathbb{X} , which is union of the subspaces of the terminal nodes that are its descendants.

Example :



$$T = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

$$\tilde{T} = \{4, 5, 6, 8, 9\}$$

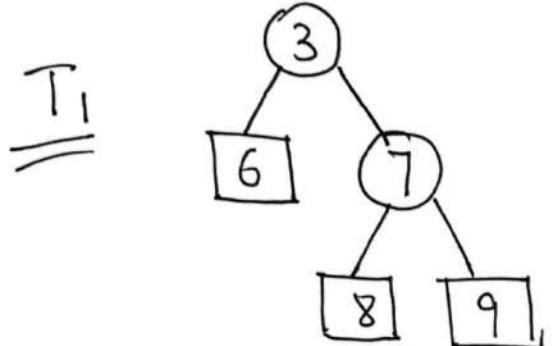
<u>t</u>	$l(t)$	$r(t)$	<u>t</u>	$l(t)$	$r(t)$
1	2	3	7	8	9
2	4	5	8	0	0
3	6	7	9	0	0
4	0	0			
5	0	0			
6	0	0			



enough to know the tree structure

Let

$$T_1 = \{3, 6, 7, 8, 9\}$$



$$\underline{T_1}$$

$$l_1(t) = \begin{cases} 1(t), & \text{if } l(t) \in T_1 \\ 0, & \text{otherwise} \end{cases}$$

$$r_1(t) = \begin{cases} r(t), & \text{if } r(t) \in T_1 \\ 0, & \text{otherwise} \end{cases}$$

$$\underline{t \quad l_1(t) \quad r_1(t)}$$

$$3 \quad 6 \quad 7$$

$$6 \quad 0 \quad 0$$

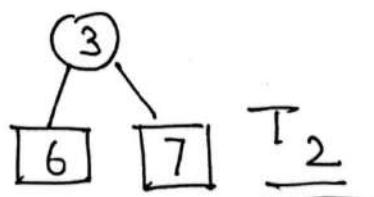
$$7 \quad 8 \quad 9$$

$$8 \quad 0 \quad 0$$

$$9 \quad 0 \quad 0$$

$$\underline{\quad x \quad \underline{\quad}}$$

$$T_2 = \{3, 6, 7\}$$



$$\underline{t \quad l_2(t) \quad r_2(t)}$$

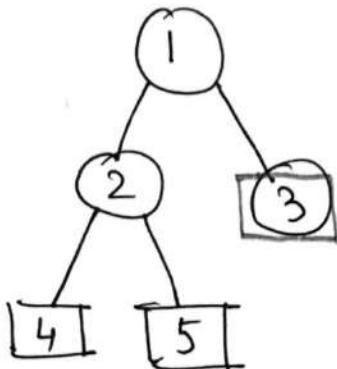
$$3 \quad 6 \quad 7$$

$$6 \quad 0 \quad 0$$

$$7 \quad 0 \quad 0$$

Also, take

$$T_3 = \{1, 2, 3, 4, 5\}$$



$$T_3 =$$

t	$l_3(t)$	$r_3(t)$
1	2	3
2	4	5
3	0	0
4	0	0
5	0	0

Remark : T_1 , T_2 , T_3 are all subtrees of T

T_1 is not a pruned subtree of T

T_2 is a pruned subtree of T_1 but is not a pruned subtree of T .

T_3 is a pruned subtree of T

Note : $\{2, 3, 4, 5\}$ is not a subtree, as 2 and 3 has no parent.

Classification Tree

(152)

Let us now look at various issues associated with construction of classification trees.

Let T be a tree with $\tilde{T} = \{t_1, \dots, t_M\}$ as the set of terminal nodes.

Let $\pi_{j(t)} \in \{\pi_1, \dots, \pi_c\}$ denote one of the class labels, that is associated with terminal node t .

A classification tree consists of $T, \tilde{T}, \{\pi_{j(t)}, t \in \tilde{T}\}$ and a partition $\{U(t) : t \in \tilde{T}\}$,

let the learning sample be $\mathcal{L} = \{(x_i, y_i) : i=1(\text{or})N\}$
(y_i 's are the class labels)

$\forall t \in T$, define

$$N(t) = \# \text{ of sample patterns in } \mathcal{L} \ni x_i \in U(t)$$

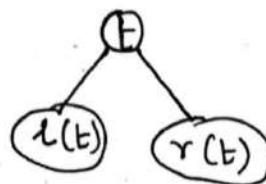
$$N_j(t) = \# \text{ of } x_i \in U(t) \text{ and } y_i = \pi_j$$

$$\text{clearly, } \sum_j N_j(t) = N(t) \quad \& \quad \sum_{t \in \tilde{T}} N(t) = N$$

$$p(t) = \frac{N(t)}{N}; \text{ estimate of } P(x \in U(t)) \text{ based on } \mathcal{L}$$

$$p(\pi_j | t) = \frac{N_j(t)}{N(t)}; \text{ estimate of } P(y = \pi_j | x \in U(t))$$

Recall that $\forall t \in T$



let $t_L = l(t)$ & $t_R = r(t)$

$$p_L^t = \frac{p(t_L)}{p(t)} ; \text{ estimate of } P(\tilde{x} \in U(t_L) | \tilde{x} \in U(t))$$

$$p_R^t = \frac{p(t_R)}{p(t)} ; \text{ estimate of } P(\tilde{x} \in U(t_R) | \tilde{x} \in U(t))$$

Rule for label assignment (class label assignment)

Assign label π_j to node t if

$$p(\pi_j | t) = \max_i p(\pi_i | t)$$

(i.e. majority voting inside the partition)

Splitting rules

Splitting rules at internal nodes are based on "node impurity" measures.

In general "node impurity" for any node $t \in T$ is defined as

$$\text{Imp}(t) = \phi(p(\pi_1 | t), p(\pi_2 | t), \dots, p(\pi_c | t))$$

Note:

Such a "node impurity" would be maximum if $p(\pi_j | t) = \frac{1}{c} \ \forall j$ (i.e. equal representation of all classes at t) and minimum if $\begin{cases} p(\pi_j | t) = 1 \text{ for some } j \\ p(\pi_i | t) = 0 \quad \forall i \neq j \end{cases} \begin{cases} \text{it's a pure node} \\ \text{if } j \neq i \end{cases}$

Examples of "node impurity" measures

(i) Gini Index (t) = $\sum_{\substack{i, j \\ i \neq j}} p(\pi_i | t) p(\pi_j | t)$

(If $c=2$ (two-class problem))

$$\text{Gini Index} (t) = \sum_{i=1}^2 \sum_{\substack{j=1 \\ j \neq i}}^2 p(\pi_i | t) p(\bar{\pi}_j | t)$$

$$= p(\pi_1 | t) p(\pi_2 | t) + p(\pi_2 | t) p(\bar{\pi}_1 | t)$$

$$= p(1-p) + (1-p)p \quad (p(\pi_1 | t) = p)$$

$$= 2p(1-p)$$

$$\xrightarrow{\max} \text{if } p = \frac{1}{2} \left(\text{i.e. } \frac{1}{c} \right)$$

$$\min \text{ if } p \text{ or } 1-p = 0/1$$

(ii) Misclassification error rate at node t

$$= \frac{1}{N(t)} \sum_{i: x_i \in V(t)} I(y_i \neq \pi_{j(t)})$$

$$\text{where } j(t) = \arg \max_i p(\pi_i | t)$$

$$= \frac{1}{N(t)} (N(t) - N_{j(t)}(t))$$

$$= 1 - \frac{N_{j(t)}(t)}{N(t)} = 1 - p(\pi_{j(t)} | t)$$

(iii) cross-entropy or deviance

$$- \sum_i p(\pi_i | t) \log p(\pi_i | t)$$

Remark: The above measures are for node impurity.
We can define tree impurity as

$$\text{Imp}(\mathcal{T}) = \sum_{t \in \mathcal{T}} p(t) \text{Imp}(t)$$

How to split a node?

This is by far the most important question!!

Consider a split using variable x_k at level l

$$\text{say, } S_k^l = \{x : x_k < l\}$$

We define a measure of "goodness of split" at node t as the change in impurity Δ . For the split S_k^l at t , this is

$$\Delta \text{Imp}(S_k^l, t) = \text{Imp}(t) - \left(p_L^t \text{Imp}(t_L) + p_R^t \text{Imp}(t_R) \right)$$

$$\left(\text{recall that } p_L^t = P(x \in \underbrace{U(t_L)}_{U(t)} | x \in U(t)) \right)$$

$$\& p_R^t = P(x \in U(t_R) | x \in U(t)).$$

We need to find $(k^*, l^*) \rightarrow \Delta \text{Imp}(S_k^l, t)$ is maximized over all $k=1 \dots p$ (dimension of feature vector) and l is allowed take one of a finite # of values within the range of possible values of the chosen feature variable.

The above "goodness of split" based approach is used to split nodes starting with the root node.

When to stop splitting?

We do not split a node if the change in the impurity due to split is less than a prescribed threshold

OR

Grow the tree till the terminal nodes are all pure (all ~~patterns~~ patterns belonging to the partition have same class label) and then apply pruning of the tree.

What is pruning? How to prune a classification tree?

There are various approaches of pruning. We discuss here 2 important approaches.

Pruning of a grown tree after applying splitting of nodes basically means cutting branches of the tree to get a subtree of the original tree.

Cost-complexity pruning

Let

$$r(t) = 1 - \max_i p(\pi_i | t)$$

↑
estimate of prob of misclassification at node t

define

$$R(t) = p(t) r(t)$$

Estimate of overall misclassification rate of the tree classifier is

$$R(T) = \sum_{t \in \tilde{T}} p(t) r(t) = \sum_{t \in \tilde{T}} R(t),$$

If α denotes the cost of complexity per terminal node then define

$$R_\alpha(t) = R(t) + \alpha \quad \begin{matrix} \text{as cost-complexity} \\ \text{criterion for node } t \end{matrix}$$

$$\& R_\alpha(T) = \sum_{t \in \tilde{T}} (R(t) + \alpha) \quad \begin{matrix} (\alpha: \text{tuning parameter}) \\ (\alpha \text{ trade-off parameter}) \end{matrix}$$

$$= \sum_{t \in \tilde{T}} R(t) + \alpha |\tilde{T}|$$

where $|\tilde{T}|$: cardinality of terminal node set \tilde{T}

Cost-complexity pruning approach : Starting from T_0 (a pure tree, i.e. a tree having all terminal nodes as pure), find the subtree T_α (for a fixed α) $\Rightarrow R_\alpha(T)$ is minimum.

Weakest link pruning approach

Let T_t be subtree with root t

& $\{t\}$ is subbranch of T_t consisting of a single node t

Let $R_\alpha(t) = R(t) + \alpha = r(t) \cdot p(t) + \alpha$
(as defined earlier)

& for T_t ; $R_\alpha(T_t) = R(T_t) + \alpha |\tilde{T}_t|$

Now $R_\alpha(T_t) < R_\alpha(\{t\})$ (i.e. T_t has smaller cost complexity than $\{t\}$).

if $R_\alpha(T_t) + \alpha |\tilde{T}_t| < R(t) + \alpha$

i.e. if $\alpha(|\tilde{T}_t| - 1) < R(t) - R(T_t)$

i.e. if $\alpha < \frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1} \quad (*)$

Note that as $\alpha \uparrow$ equality is achieved at $(*)$

and T_t and $\{t\}$ have same cost-complexity
and $\{t\}$ is preferred to T_t ; i.e. T_t can be pruned at t

We define $g(t) = \frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1}$

as the "strength of link" at node t

Find $t^* \Rightarrow g(t^*)$ is min and then
prune the tree at this "Weakest-link", i.e.
at t^* .

At the next step start with the pruned
subtree and find the weakest link and prune
again.

Continue to reach the root through pruning
to get the sequence of pruned subtrees

$$T_1 \supset T_2 \supset \dots \supset \{t\}$$

\uparrow org root

Under this approach, the pruned subtree
with $\min R_\alpha(T)$ is the best pruned tree.

Growing a regression tree

How to split nodes?

$$d = \{(x_i, y_i) : i=1 \text{ to } N\}$$

Consider a split variable j and a split pt t

$$R_1(j, t) = \{x : x_j < t\} \& R_2(j, t) = \{x : x_j \geq t\}$$

Set the criterion f^n as

$$\min_{j, t} \left[\min_{c_1} \sum_{i: x_i \in R_1(j, t)} (y_i - c_1)^2 + \min_{c_2} \sum_{i: x_i \in R_2(j, t)} (y_i - c_2)^2 \right] \quad (\#)$$

Note that for a fixed $j \& t$

$$\min \hat{c}_1 = \text{average } (y_i | x_i \in R_1(j, t))$$

$$\& \hat{c}_2 = \text{average } (y_i | x_i \in R_2(j, t))$$

Find (j^*, t^*) for the optimum split at a node \rightarrow

$$\sum_{i: x_i \in R_1(j, t)} (y_i - \hat{c}_1)^2 + \sum_{i: x_i \in R_2(j, t)} (y_i - \hat{c}_2)^2 \text{ is minimized}$$

here j varies over all $j = 1 \text{ to } p$ (feature vector dim)

& t in the grid of the j th variable

Growing a regression tree

How to split nodes?

$$\mathcal{D} = \{(x_i, y_i) : i=1 \dots N\}$$

Consider a split variable j and a split pt t

$$R_1(j, t) = \{x_i : x_j < t\} \& R_2(j, t) = \{x_i : x_j \geq t\}$$

Set the criterion f^* as

$$\min_{j, t} \left[\min_{c_1} \sum_{i: x_i \in R_1(j, t)} (y_i - c_1)^2 + \min_{c_2} \sum_{i: x_i \in R_2(j, t)} (y_i - c_2)^2 \right] \quad (*)$$

Note that for a fixed j &

$$\min \rightarrow \hat{c}_1 = \text{average } (y_i | x_i \in R_1(j, t))$$

$$\& \hat{c}_2 = \text{average } (y_i | x_i \in R_2(j, t))$$

Find (j^*, t^*) for the optimum split at a node \Rightarrow

$$\sum_{i: x_i \in R_1(j, t)} (y_i - \hat{c}_1)^2 + \sum_{i: x_i \in R_2(j, t)} (y_i - \hat{c}_2)^2 \text{ is minimized}$$

here j varies over all $j = 1 \dots p$ (feature vector dim)

& t in the grid of the j th variable

When to stop splitting?

(I) Split tree node only if the decrease in sum of squares residual due to split exceeds some pre-assigned threshold

OR

(II) Std dev of responses at terminal nodes is low (below a threshold)

OR

(III) Grow a large tree and stop splitting if the node size (# of patterns reaching that node) reaches a low threshold level and apply pruning.

Regression tree pruning

$$\forall t \in \tilde{T} \text{ define } \hat{\epsilon}_t = \frac{1}{N(t)} \sum_{x_i \in U(t)} y_i$$

$$N(t) : \# \{ x_i \in U(t) \}$$

Impurity measure:

Mean square error & obsn in $U(t)$ from λ

$$Q_t = \frac{1}{N(t)} \sum_{i: x_i \in U(t)} (y_i - \hat{\epsilon}_t)^2$$

Define cost complexity f^n as

$$C_\alpha(T) = \sum_{t \in \tilde{T}} N(t) Q_t + \alpha |\tilde{T}|$$

α : tuning parameter which governs the trade-off
bet tree size and goodness of fit

$|\tilde{T}|$: cardinality of \tilde{T} .

Pruning can be done based on $C_\alpha(T)$. For a fixed α , let T_α be a subtree of T obtained by pruning T . We try to find

$$T_\alpha \subset T \Rightarrow T_\alpha \text{ minimizes } C_\alpha(T)$$

for a fixed α .

(Note: larger the α smaller is the opt T_α)

Weakest link pruning: Collapse internal nodes (to prune) that produces the smallest per node increase in

$$\sum_{t \in \tilde{T}} N(t) Q(t) \text{ and continue till}$$

root node is reached.

Select the optimal from the pruned subtree sequence.

Two-classPerceptron learning rule for classificationTraining set $\mathcal{D} = ((\underline{x}_1, i_1), \dots, (\underline{x}_n, i_n))$

$$i_j \in \{1, 2\}, \quad \forall j = 1, \dots, n; \quad \underline{x}_i \in \mathcal{X}$$

Pops π_1 & π_2 .

A linear classifier

$$\underline{w}' \underline{x} + w_0 \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \underline{x} \in \begin{cases} \pi_1 \\ \pi_2 \end{cases}$$

$$\underline{w} = (w_0, w_1, \dots, w_p)' \quad \underline{x} = (1, x_1, \dots, x_p)'$$

or

$$\underline{w}' \underline{x} \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \underline{x} \in \begin{cases} \pi_1 \\ \pi_2 \end{cases}$$

Note: \underline{x} can also be $= (1, \phi_1(\underline{x}), \dots, \phi_D(\underline{x}))'$.

Define

$$\underline{y}_i' = (1, \underline{x}_i'), \text{ if } \underline{x}_i \in \pi_1,$$

$$\text{& } \underline{y}_i' = (-1, -\underline{x}_i'), \text{ if } \underline{x}_i \in \pi_2.$$

Ideally, we want a solution for $\underline{w} \Rightarrow$

$$\underline{w}' \underline{y} > 0 \text{ for as many samples as possible.}$$

If \exists a $\underline{w} \Rightarrow \underline{w}' \underline{y} > 0 \forall$ patterns then

the data are said to be 'linearly separable'.

$$\frac{19}{4} - 3 \frac{10}{4} \quad \frac{7}{4} + 2$$

If $\underline{v}' \underline{y}_i < 0$ then a misclassification occurs. (164)

Define the perceptron criterion function as

$$J_p(\underline{v}) = \sum (-\underline{v}' \underline{y}_i)$$

$\underline{y}_i \in \mathcal{Y} \rightarrow$ set of all misclassified patterns

- * $J_p(\underline{v})$ is proportional to the sum of distances of the misclassified samples to the decision boundary.

Objective is to find $\underline{v} \ni J_p(\underline{v})$ is minimised.
gradient-descent procedure used to solve $\min_{\underline{v}} J_p(\underline{v})$

$$\frac{\partial J_p(\underline{v})}{\partial \underline{v}} = \sum_{\underline{y}_i \in \mathcal{Y}} (-\underline{y}_i)$$

Iteration steps

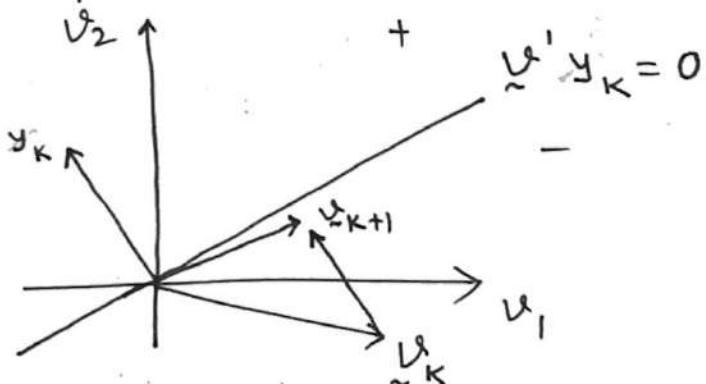
$$(i) \quad \underline{v}_{k+1} = \underline{v}_k + \rho \sum_{\underline{y}_i \in \mathcal{Y}} \underline{y}_i$$

batch update mode or many pattern adaptation

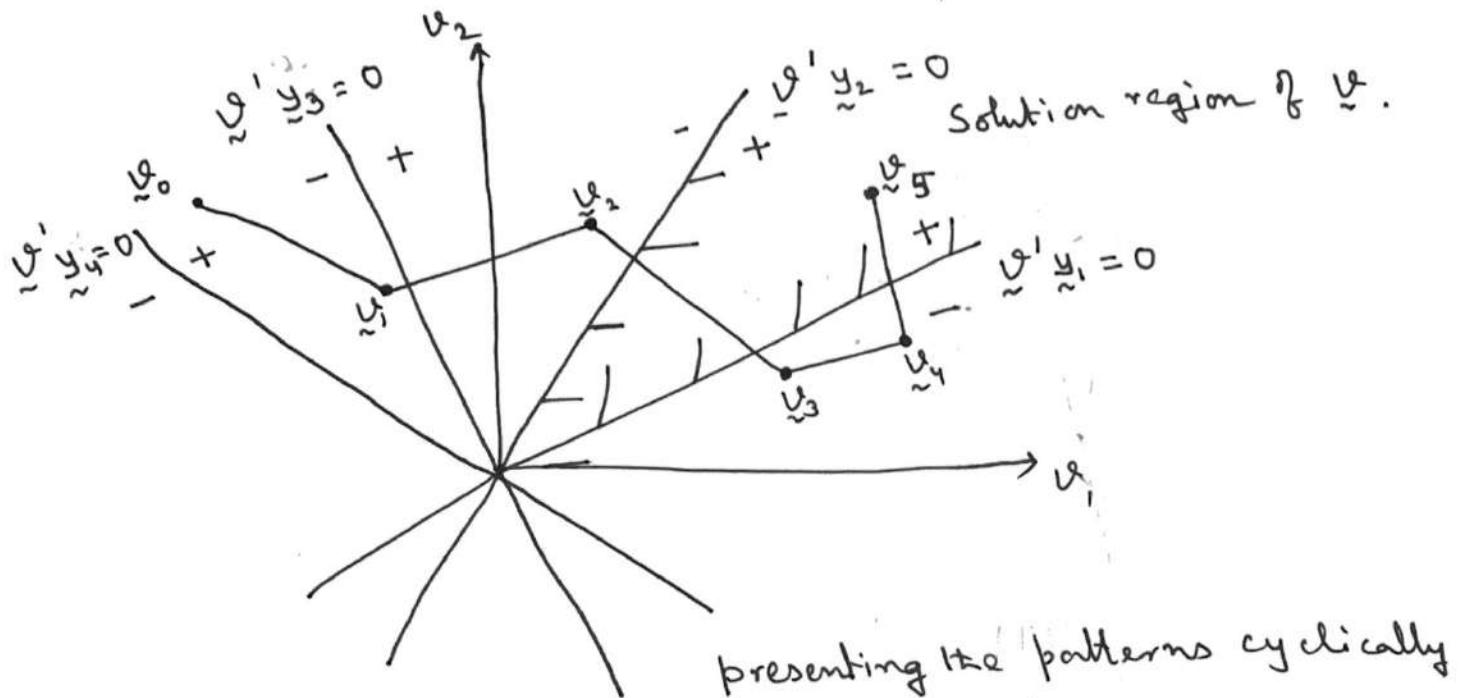
$$(ii) \quad \underline{v}_{k+1} = \underline{v}_k + \rho \underline{y}_i \quad \rho: \text{step size}$$

\underline{y}_i is a training pattern that is misclassified.

instantaneous update mode or single-pattern
weight update in weight space adaptation



example with 4 patterns on single-pattern adaption mode



presenting the patterns cyclically

$$\tilde{y}_1, \tilde{y}_2, \tilde{y}_3, \tilde{y}_4$$

Starting with \tilde{y}_0 , first updation for \tilde{y}_2 takes \tilde{v}_0 to \tilde{v}_1 ; next updation for \tilde{y}_3 ($\tilde{v}_1 \rightarrow \tilde{v}_2$); next updation for \tilde{y}_2 ($\tilde{v}_2 \rightarrow \tilde{v}_3$) ---.

A solution of $J_p(\tilde{v}) = 0$. will be obtained for linearly separable patterns.

Note : An important variation of the above perceptron learning criterion is through introduction of a 'margin' $b > 0$.

A weight vector is updated whenever

$$\tilde{v}^T \tilde{y}_i \leq b$$

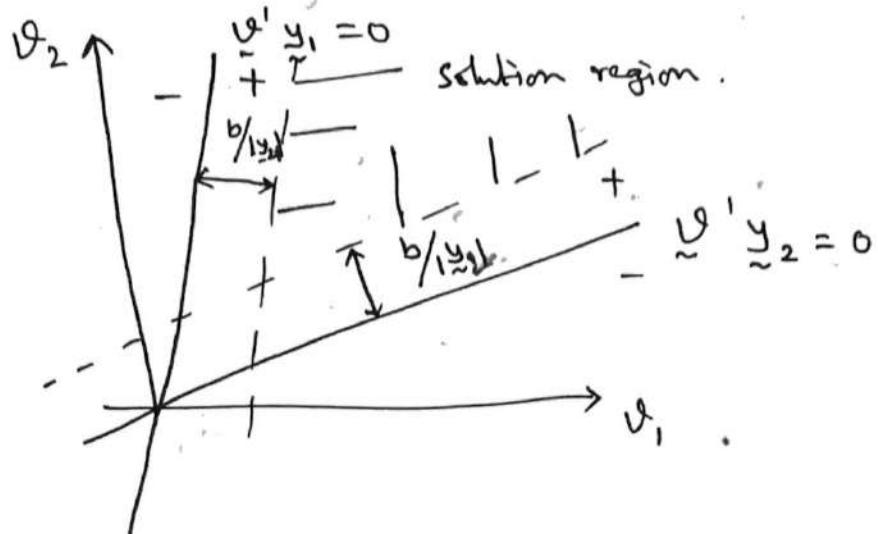
\Rightarrow the solution vector \tilde{v} must lie at a distance greater than (in term of feature space) $b/|\tilde{y}_i|$ from hyperplane $\tilde{v}^T \tilde{y}_i = 0$

Aim of the margin is to improve generalization

without the margin some points may lie too close to the separating boundary i.e. the separating hyperplane.

Maximal margin classifier \rightarrow Support Vector Machine (SVM) classifier.

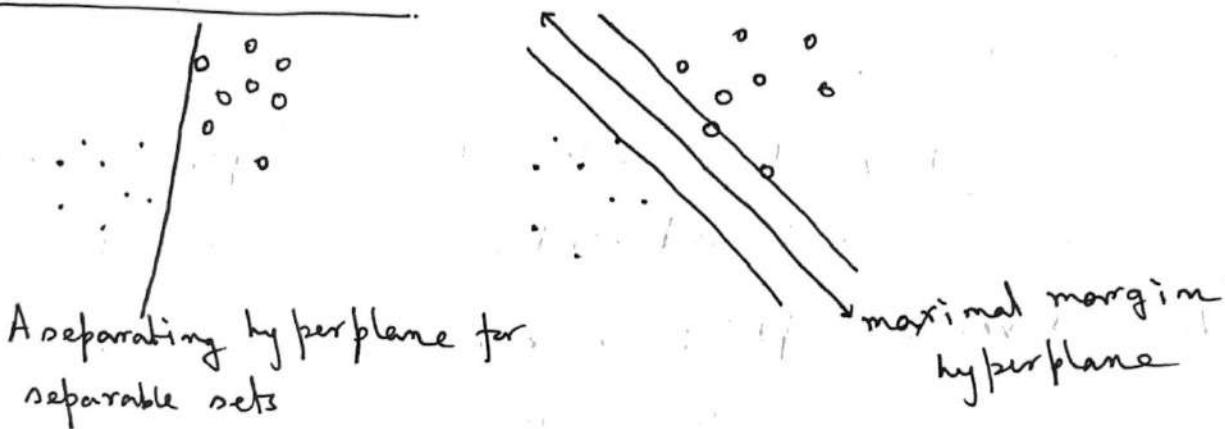
Solution region with a margin



SVM

Construction of 'best' separating hyperplane, i.e. maximal margin hyperplane

Basic SVM model:



Suppose the data is linearly separable

Two classes: $\pi_1 \text{ & } \pi_2 \rightarrow \text{labels } y_i = \pm 1$

$$\mathcal{D} = ((x_1, y_1), \dots, (x_n, y_n))$$

linear discriminant function

$$g(x) = w'x + w_0$$

Decision rule

$$\underline{w}' \underline{x} + w_0 \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \underline{x} \in \begin{cases} \Pi_1 \text{ with } y_i = +1 \\ \Pi_2 \text{ with } y_i = -1 \end{cases}$$

\Rightarrow Training pts are correctly classified if

$$y_i (\underline{w}' \underline{x}_i + w_0) > 0 \quad \forall i$$

Let $b > 0$ be a margin \Rightarrow we seek a solution \Rightarrow

$$y_i (\underline{w}' \underline{x}_i + w_0) \geq b$$

Perceptron learning algorithm yields a solution for which all points \underline{x}_i are at a distance greater than or equal to $b/|\underline{w}|$ from the separating hyperplane

Note: A scaling of b , w_0 and \underline{w} leaves $b/|\underline{w}|$ unaltered and

$$y_i (\underline{w}' \underline{x}_i + w_0) \geq b \text{ is still satisfied}$$

w.l.o.g. we take $b=1 \rightarrow$ canonical hyperplanes

$$CH: H_1: \underline{w}' \underline{x} + w_0 = +1$$

$$H_2: \underline{w}' \underline{x} + w_0 = -1$$

and

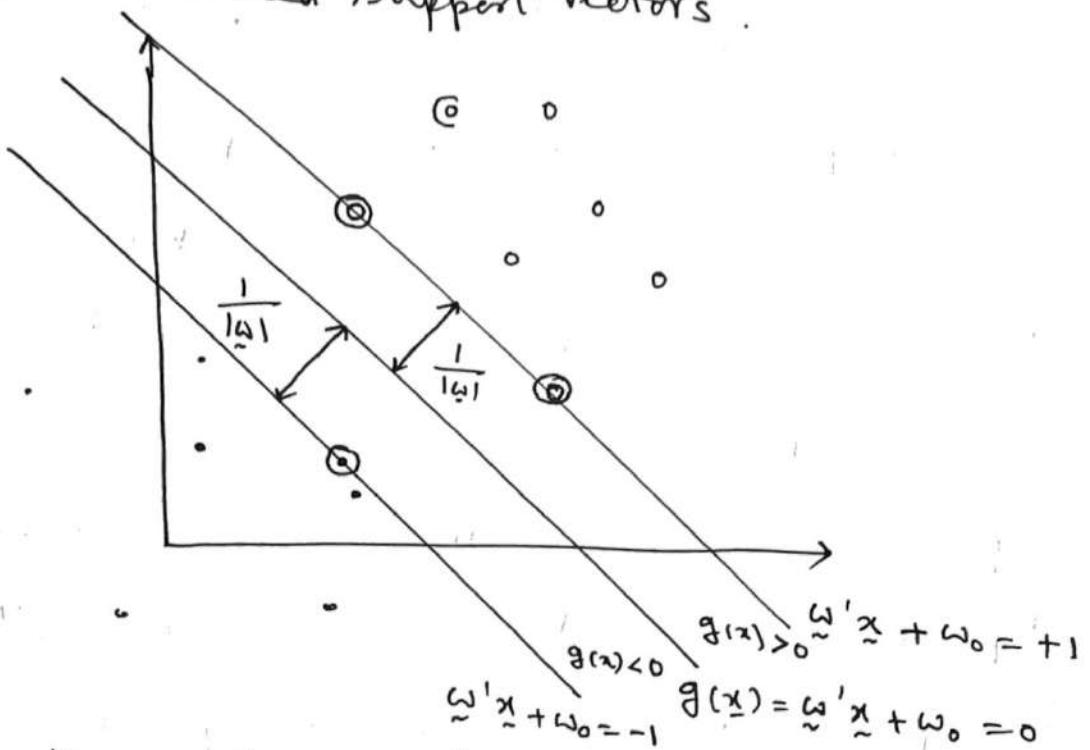
$$\underline{w}' \underline{x}_i + w_0 \geq +1 \quad \text{for } y_i = +1$$

$$\underline{w}' \underline{x}_i + w_0 \leq -1 \quad \text{for } y_i = -1$$

The distance of H_1 & H_2 from the separating hyperplane $g(\underline{x}) = \underline{w}' \underline{x} + w_0 = 0$ is $1/|\underline{w}|$ and is termed the margin bet^{canoncial} H_1/H_2 hyperplanes and the sep hyperplane

(168)

The points that lie on the canonical separating hyperplanes H_1 and H_2 are called 'Support vectors'.



Objective

Maximizing the margin subject to the constraint that the separating hyperplane separates the linearly separable groups i.e. $\text{Max } \frac{1}{\|w\|}$ i.e. $\text{Min } \|w\| \rightarrow \text{the constraint}$

$$C : y_i (\underline{w}' \underline{x}_i + w_0) \geq 1 \quad \forall i=1 \dots n$$

$$L_p = \frac{1}{2} \underline{w}' \underline{w} - \sum_{i=1}^n \alpha_i (y_i (\underline{w}' \underline{x}_i + w_0) - 1).$$

Where $\alpha_i \geq 0 ; i=1 \dots n$ are Lagrange multipliers.

(p+1) Primal parameter w_0, w_1, \dots, w_p .

Note (i) Solution obtained using quadratic programming

(ii) Case of linearly non-separable data solved using appropriate slack variables

(iii) Concepts can be extended for multi class problem

ANN

- Computational algorithm inspired by the parallelism of human brain
- Tries to mimic the cognitive ability of the human brain to identify patterns and to make decisions and forecasts based on past experience

How NN works?

- Non-linear models that can be trained to map the past and the future values of a time series
- extract hidden structure and relationship governing the data
- powerful pattern recognition capabilities

NN in the context of statistical methods

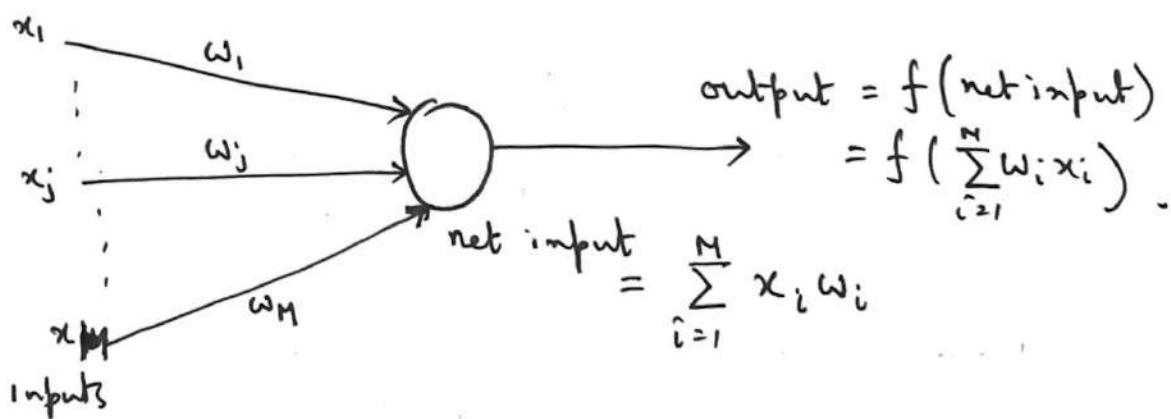
- Multivariate, non-linear, non-parametric inference technique
- completely data-driven and assumption free
- ability to model non-linear dynamics.

Some application areas

Stock market prediction, foreign exchange trading, bankruptcy prediction, volatility prediction, economic indicator forecasting, credit rating, prediction of international airline passenger traffic, nowcasting in aviation industry, ozone level prediction

Human brain: Hundred billion neurons - interconnected structure

Basic computational unit of ANN - a neuron - an artificial neuron



x_i : input value from the i^{th} input node

w_i : weight on the connection from the i^{th} input node to the PU.

PU - collects the inputs from all the inputs to get the net input value

$$\sum_{i=1}^M w_i x_i ; w_i \text{ denoting the relative importance of } i^{\text{th}} \text{ input}$$

PU processes the net input value and generates the output

$$\text{Out} = f\left(\sum_{i=1}^M w_i x_i\right)$$

↑ transfer function / squashing f^n .

A The simplest NN

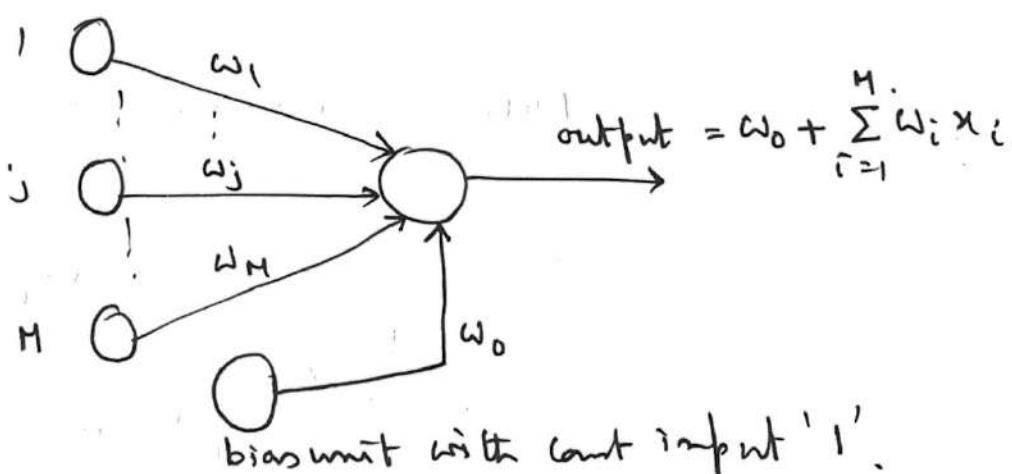
A single PU network.

MLR setup $\lambda = \{(x_i, y_i), \dots, (x_n, y_n)\}$.

Consider a transfer f^n as Identity

$$\text{output} = f\left(\sum_{i=1}^M w_i x_i\right) = \sum_{i=1}^M w_i x_i$$

MLR model without a constant term



Unknown parameters : weight connections (w_1, \dots, w_n)
or (w_0, w_1, \dots, w_n)

Optimum values of (w_0, w_1, \dots, w_M) can be obtained by minimizing residual S_{sq} .

Note: Injecting non-linearity into the model -

consider a non-linear transfer f^n .

e.g (i) A sigmoidal transfer f^n .

$$f(x) = \frac{1}{1 + e^{-x}}$$

(ii) A tanh transfer f^n

$$f(x) = \tanh(x).$$

For (i), output = $\frac{1}{1 + e^{-(w_0 + \sum_{i=1}^M w_i x_i)}} = g(x)$

A non-linear ^{model} of response in terms of the
the values of the feature vector components

Let (x_p, y_p) denote the p^{th} pattern in L .

Network output: $O_p = f(w_0 + \sum_{j=1}^M x_{pj} w_j) / f(\sum_{j=1}^M x_{pj} w_j)$.

Actual output: $y_p =$

Error for the p^{th} pattern : $y_p - O_p$

Error square : $E_p = \frac{1}{2}(y_p - O_p)^2 \leftarrow$ instantaneous error
for p^{th} pattern

Total error over L

$$E = \frac{1}{2} \sum_{p=1}^n E_p = \frac{1}{2} \sum_{p=1}^n (y_p - O_p)^2.$$

Training of the network : determination of optimum weight vector.

Modes of Learning : Instantaneous mode or single pattern
adaptation mode
or Batch mode or multiple pattern
adaptation mode.

Method for obtaining the weight updation : Gradient
descent
algorithm

Instantaneous mode :

$$\epsilon_p = \frac{1}{2} (y_p - o_p)^2$$

$$\epsilon_p = \frac{1}{2} (y_p - f(\sum_{j=1}^M w_j x_{pj}))^2$$

Gradient of the error ϵ_p w.r.t. the weights

$$\begin{aligned}\frac{\partial \epsilon_p}{\partial w_j} &= -(y_p - o_p) \cdot \frac{\partial f(\sum w_j x_{pj})}{\partial w_j} \\ &= -(y_p - o_p) \frac{\partial f(\sum w_j x_{pj})}{\partial (\sum w_j x_{pj})} \cdot \frac{\partial (\sum w_j x_{pj})}{\partial w_j} \\ &= -(y_p - o_p) f'(\sum w_j x_{pj}) x_{pj}.\end{aligned}$$

e.g. for sigmoidal transfer f^n .

$$f(\sum w_j x_{pj}) = \frac{1}{1 + e^{-\sum w_j x_{pj}}} = o_p$$

$$\left[f(x) = \frac{1}{1 + e^{-x}} ; f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) \right]$$

$$f'(\sum w_j x_{pj}) = o_p(1 - o_p)$$

$$\Rightarrow \frac{\partial \epsilon_p}{\partial w_j} = -(y_p - o_p) o_p (1 - o_p) x_{pj}$$

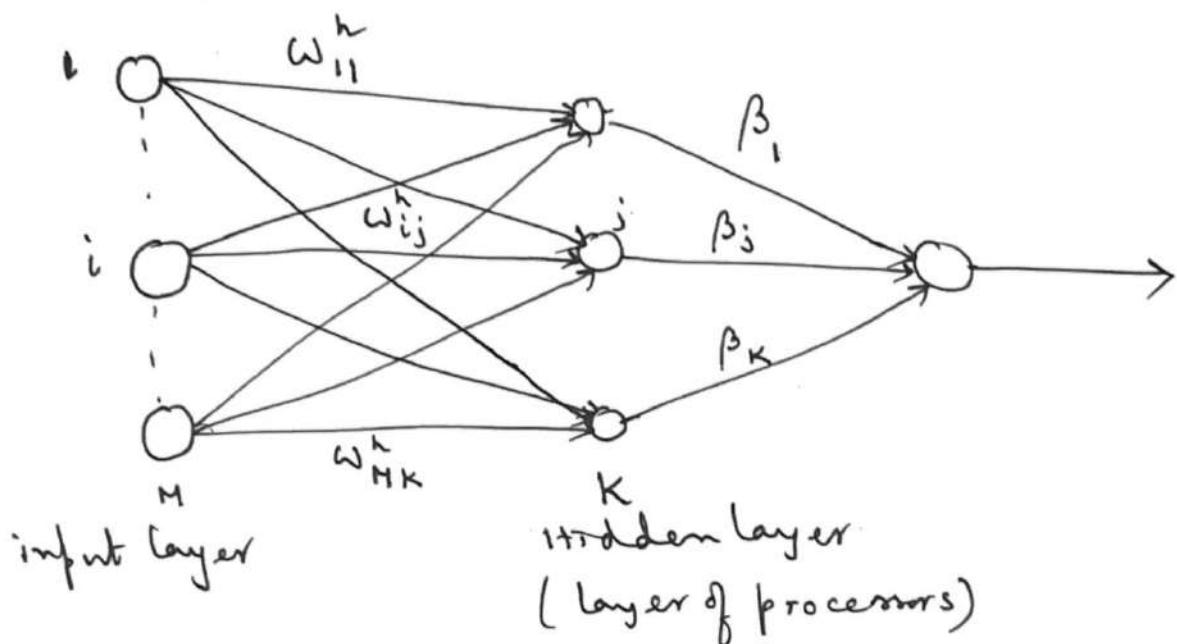
Gradient descent algorithm weight updation

$$\omega_j^{(k+1)} = \omega_j^{(k)} + \eta \left(-\frac{\partial E_p}{\partial \omega_j} \right) \Big|_{\omega_j^{(k)}}$$

$$\omega_j^{(k+1)} = \omega_j^{(k)} + \eta (y_p - o_p) o_p (1 - o_p) x_{pj}$$

Present the patterns at random and use the above equation for weight updation until stopping criterion is reached.

Multi-layer Perceptron (MLP)



M input, K hidden layer neurons and one output

M-K-1 MLP.

pth pattern $\rightarrow (x_p, y_p)$

$$\text{net}_{pj}^h = \sum_{i=1}^M w_{ij}^h x_{pi} \quad (\text{may have a bias unit})$$

net input value at jth neuron of hidden layer

output from the jth neuron

$$f(\text{net}_{pj}^h) = i_{pj}$$

i_{pj} : j^{th} input to the output unit

net input value at output

$$\sum_{j=1}^K \beta_j i_{pj} = \text{net}_p^o$$

Final network output

$$O_p = f(\text{net}_p^o) \quad \left(\begin{array}{l} \text{this 'f' can be diff} \\ \text{from the f used at} \\ \text{hidden layer neurons} \end{array} \right).$$

$$\text{i.e. } O_p = f(\text{net}_p^o)$$

$$= f\left(\sum_{j=1}^K \beta_j i_{pj}\right)$$

$$= f\left(\sum_{j=1}^K \beta_j f\left(\sum_{i=1}^M w_{ij}^h x_{pi}\right)\right).$$

$$\text{or } = g\left(\sum_{j=1}^K \beta_j f\left(\sum_{i=1}^M w_{ij}^h x_{pi}\right)\right).$$

Note: We can have different 'f' at different hidden layer neurons, in that case

$$O_p = g\left(\sum_{j=1}^K \beta_j f_j\left(\sum_{i=1}^M w_{ij}^h x_{pi}\right)\right).$$

Note: Suppose we consider sigmoid transfer f^* at all neurons (hidden and output) then

$$O_p = f\left(\sum_{j=1}^K \beta_j \left\{1 + \exp\left(-\sum_{i=1}^M w_{ij}^h x_{pi}\right)\right\}^{-1}\right)$$

$$= \left\{1 + \exp\left(-\sum_{j=1}^K \beta_j \left[1 + \exp\left(-\sum_{i=1}^M w_{ij}^h x_{pi}\right)\right]\right)\right\}^{-1}$$

$$\text{i.e. } O_p = f^*\left(x_{p1}, \dots, x_{pM}\right) \rightarrow \text{a non-linear model.}$$

Note: The network above is called a feed-forward network.

unknown parameters for an $M-K-1$ network: weight connections.

input layer \rightarrow hidden layer: $M \times K$ weight matrix

$$\begin{pmatrix} w_{11}^h & w_{12}^h & \dots & w_{1K}^h \\ \vdots & \vdots & \ddots & \vdots \\ w_{M1}^h & w_{M2}^h & \dots & w_{MK}^h \end{pmatrix}.$$

hidden layer \rightarrow output node: K -dimensional weight vector
 $(\beta_1, \dots, \beta_K)$.

Total # of unknown weights: $(MK + K)$.

Learning of the network: determination of weights.

Either instantaneous mode OR batch mode. Using
 Back-propagation algorithm (b-t of errors).

Consider for example learning on an instantaneous mode.

$$\epsilon_p = \frac{1}{2} (y_p - o_p)^2 \quad \leftarrow \text{instantaneous sq error for } p^{\text{th}} \text{ pattern}$$

$$\epsilon_p = \frac{1}{2} (y_p - f^o(\text{net}_p^o)) \quad \left| \begin{array}{l} o_p = f^o(\sum_j \beta_j f_j^o(\text{net}_j^o)) \\ \text{net}_p^o = \sum_{j=1}^K \beta_j i_{pj} \end{array} \right.$$

Gradient of ϵ_p w.r.t. output hidden to output layer wts.

$$\frac{\partial \epsilon_p}{\partial \beta_j^o} = -\frac{1}{2} (y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \beta_j^o}$$

$$= -\frac{1}{2} (y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \text{net}_p^o} \cdot \frac{\partial \text{net}_p^o}{\partial \beta_j^o}.$$

$$\text{net}_p^o = \sum_{j=1}^K \beta_j i_{pj}$$

$$= -\frac{1}{2} (y_p - o_p) f'^o(\text{net}_p^o) i_{pj}.$$

$$\text{i.e. } \frac{\partial \epsilon_p}{\partial \beta_j^o} = -\frac{1}{2} (y_p - o_p) o_p (1 - o_p) i_{pj}$$

for a sigmoid transfer f^o at the output node

Slay gradient of E_p w.r.t. input layer \rightarrow hidden layer, $\nabla_{w_{ij}^h}$

$$\frac{\partial E_p}{\partial w_{ij}^h} = \frac{\partial}{\partial w_{ij}^h} \frac{1}{2} (y_p - f^o(\sum_{j=1}^k \beta_j^o + \sum_{i=1}^m w_{ij}^h x_{pi}))^2$$

$$\begin{aligned} o_p &= f^o(\text{net}_p^o) \\ \text{net}_p^o &= \sum_{j=1}^k \beta_j^o + i_{pj} \\ i_{pj} &= f_j(\sum_{i=1}^m w_{ij}^h x_{pi}) \\ i_{pj} &= f_j(\text{net}_{pj}^h) \\ \text{net}_{pj}^h &= \sum_{i=1}^m w_{ij}^h x_{pi} \end{aligned}$$

$$= -(y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \text{net}_p^o} \cdot \frac{\partial \text{net}_p^o}{\partial i_{pj}} \cdot \frac{\partial i_{pj}}{\partial \text{net}_{pj}^h} \cdot \frac{\partial \text{net}_{pj}^h}{\partial w_{ij}^h}$$

$$= -(y_p - o_p) f'^o(\text{net}_p^o) \beta_j^o f'_j(\text{net}_{pj}^h) x_{pi}$$

Define, δ output layer delta as.

$$\delta_p^o = (y_p - o_p) f'^o(\text{net}_p^o)$$

and hidden layer delta as

$$\delta_{pj}^h = \delta_p^o \beta_j^o f'_j(\text{net}_{pj}^h)$$

Weight updation eq's

output layer weights

$$\Leftrightarrow \beta_j^o(k+1) = \beta_j^o(k) + \eta \delta_p^o i_{pj} \Big|_{\beta_j^o(k), w^{(k)}}$$

hidden layer weights

$$w_{ij}^h(k+1) = w_{ij}^h(k) + \eta \delta_{pj}^h x_{pi} \Big|_{\beta_j^o(k+1), w^{(k)}}$$

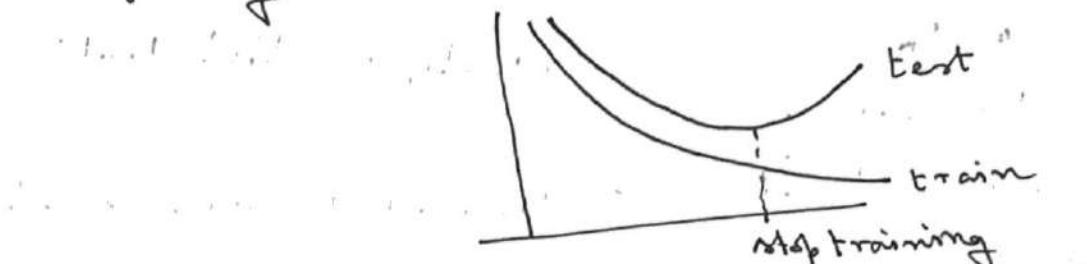
η : learning rate parameter

Iteration continues till stopping criterion is reached
Either low error below threshold OR change in wt vector is
lower than a threshold.

- Note: Information flows from input level to output level
 → error calculated on output layer → propagated back for adjustment of hidden to output layer wts.
 → with updated hidden to output layer wts error propagates backwards for weight updation of input to hidden layer wts.
 → Back-propagation of error.

Some important points: (after the discussion on variants)

- (i) Data splitting - training set and test set
- (ii) Preprocessing - normalization of inputs and output
 - Along channel & across channel
 - $\xrightarrow{\frac{* - \bar{x}_i}{Sd_{x_i}}}$
 - $\xrightarrow{\frac{* - \min}{\max - \min}}$
- (iii) Overfitting - MSE on training & test set



- Weight decay - add a penalty term for large weights

$$R(\theta) + \lambda J(\theta)$$

$$R(\theta) = \sum_{p=1}^n (y_p - \hat{y}_p)^2 \quad \theta: \text{set of weights}$$

$$(MK + K)$$

$$J(\theta) = \sum_{i,j} w_{ij}^2 + \sum_K \beta_K^2$$

λ : tuning parameter, $\lambda \geq 0$

larger λ tend to shrink weights towards 0

(approach similar to Ridge regression)

- Weight elimination

$$J(\theta) = \sum_{i,j} \frac{w_{ij}^2}{1 + w_{ij}^2} + \sum_K \frac{\beta_K^2}{1 + \beta_K^2}$$

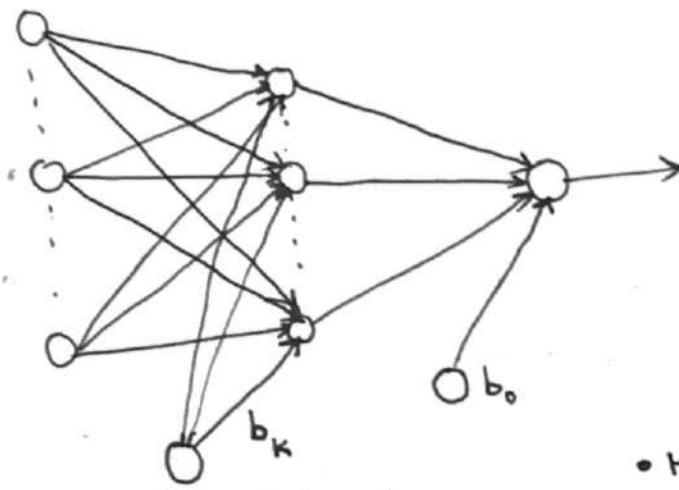
(iv) Number of hidden layers & units

One approach is to resort to experimentation

Systematic approach  Simulated annealing
Genetic algorithm

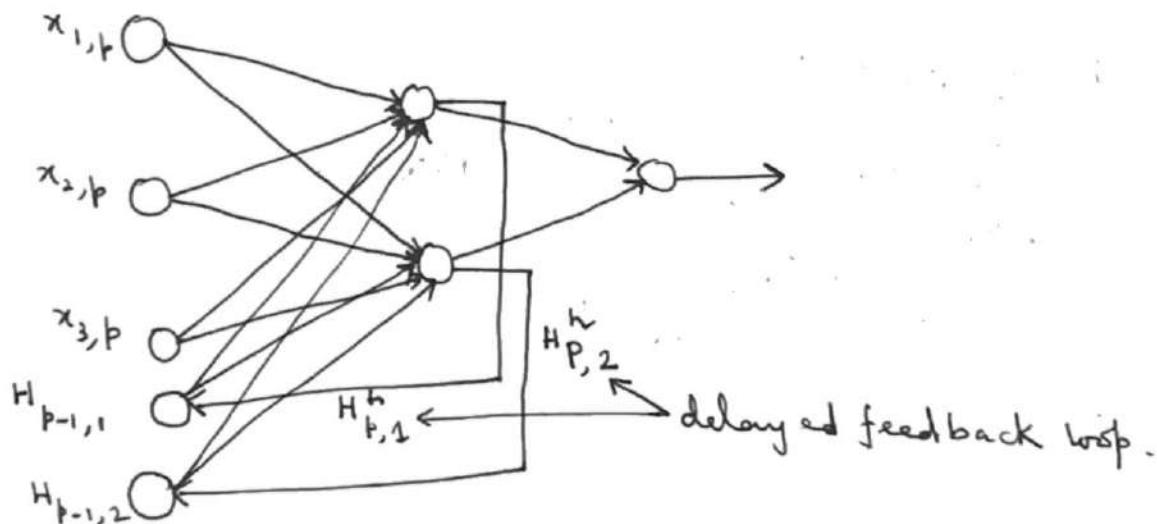
Variants of the feedforward MLP

- MLP with bias units



• Multiple output MLP

- ~~MLP~~ Recurrent MLP (more appropriate for time series set up)
3-2-1



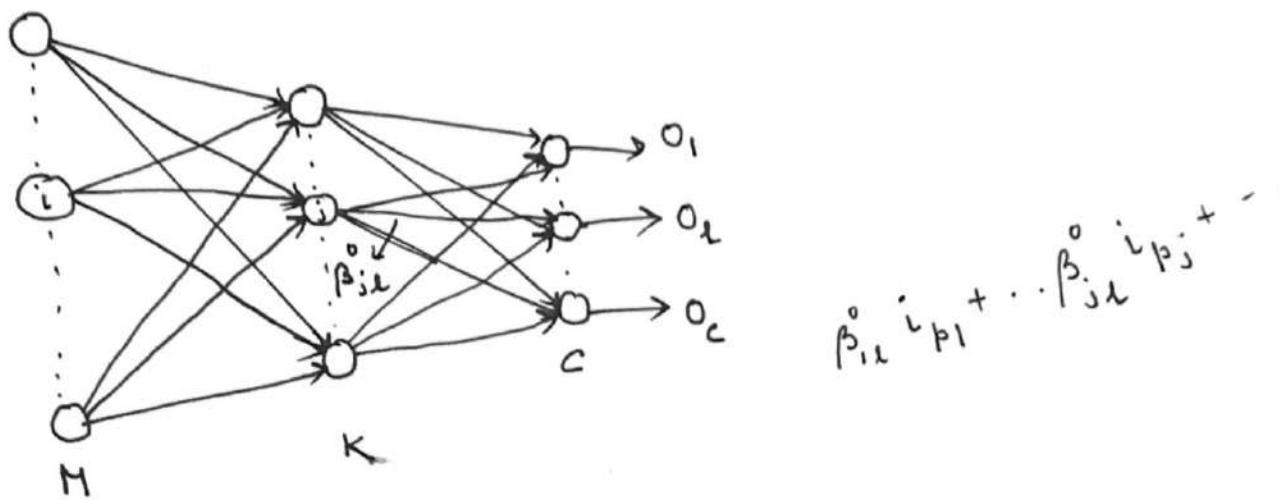
$$\hat{y}_p = \phi \left(\sum_{k=1}^2 \beta_k^0 \psi_k^h \left(\sum_{i=1}^3 x_{pi} w_{ij}^h + \sum_{k=1}^2 w_{kj}^f \psi_k^h \left(\sum_{i=1}^3 x_{p-1, i} w_{ij}^h \right) \right) \right)$$

← feedback input reaching hidden units.

ANN classification model

c-class problem

(179)



$$\text{net}_{p_j}^h = \sum_{i=1}^M x_{pi} w_{ij}^h$$

$$i_{pj} = f_j^h(\text{net}_{pj}^h) = f_j^h\left(\sum_{i=1}^M x_{pi} w_{ij}^h\right)$$

At l^{th} output unit

inputs are $i_{p_1}, i_{p_2}, \dots, i_{p_K}$
 Net input at l^{th} output $\sum_{j=1}^K \beta_{jl}^o i_{pj} = T_l$
 $\lambda = 1 \dots c$

output of l^{th} output unit

$$o_l = \frac{e^{i_p T_l}}{\sum_{k=1}^c e^{i_p T_k}} ; \quad l = 1 \dots c$$

$$\sum_{l=1}^c o_l = 1 ; \quad 0 \leq o_l \leq 1$$

Corresponding classifier is

$$G(x) = \underset{l}{\operatorname{argmax}} \quad o_l$$

Total error :

$$\sum_{i=1}^n \sum_{k=1}^c (y_{ik} - o_k)^2$$

or deviance $- \sum_{i=1}^n \sum_{k=1}^c y_{ik} \log o_k(x_i)$

(*) Discuss how training set is to be constructed.

This will specify what is

Steps for ANN model building

(180)

- (1) Creation of Data records
- (2) Splitting of data into training & test set
- (3) Preprocessing of inputs and output
- (4) Decide network architecture - # of inputs, # of hidden layers, # of neurons in hidden layer, transfer fns
- (5) Training of the network
- (6) Calculate error measure for training & test set