

• Single Value Decomposition `svd()`

$L_2$  Euclidean norm:  $\|x\| = (\sum x_i^2)^{1/2}$   
 $L_1$  Manhattan norm:  $\|x\| = \sum |x_i|$   
 matrix norm:  $\|A\| = \{\text{trace}(A^T A)\}^{1/2}$

$X = UDV^T$ ,  $U$  n x p with orth col,  $V$  p x p with orth col

$D = \begin{pmatrix} d_1 & & \\ & \ddots & \\ 0 & & d_p \end{pmatrix}$   $d_1 \geq d_2 \geq \dots \geq d_p = 0$

$X = \sum d_j u_j v_j^T$  where  $u_j, v_j$  col of  $U, V$ .

$X^* = O_1 X O_2$ ,  $O_1$  n x n,  $O_2$  p x p  $\Rightarrow$  `svd(X) = svd(X*)  
 $\|X\|_2 = d_1$   $\|X\|_F = (d_1^2 + \dots + d_p^2)^{1/2}$`

• PCA: reduce dimensions

find  $r < p$  that has same structure.

PCA define  $y_i = Ax_i$   $A$  r x p matrix `princomp()` **PCs are orthogonal**

$X = UDV^T$ ,  $X^* = UD^*V^T$ ,  $D^* = \begin{pmatrix} d_1 & & \\ & \ddots & \\ & & d_r \end{pmatrix}$  **Comp =  $\sum$  loadings  $\cdot x_i$**   
 Let  $\tilde{X} = X - \bar{X}$ ,  $\tilde{X} = UDV^T$ ,  $S = \frac{1}{n-1} VD^*V^T$  **PCs are uncorrelated**

Define  $Y = \tilde{X}V = (\tilde{y}_1 \dots \tilde{y}_p)$  **pc-scores: column of  $Y$  ( $UD$  or  $\tilde{X}V$ )  $\ell_j \cdot y_j$**   
**pc-loadings: vector  $v_1 \dots v_p$**   
 first pc: overall strength  
 second pc: contrast

Andrew Curve (identify outliers)  
 $\phi_1 = \frac{x}{\sqrt{2}}$   $\phi_2 = \sin(2\pi t)$   $\phi_3 = \cos(2\pi t)$   $\phi_4 = \sin(4\pi t)$   $\phi_5 = \cos(4\pi t)$   $\int_0^1 \phi_1 \phi_2 = 0$

• ICA replace assumption  $\text{Cov}(Y_i) = I$  with  $Y_i$  independent  
 $X_i = \mu + AY_i$   $Y_i$  independent  $E(Y_i) = 0$   $\text{Cov}(Y_i) = I$   $A$  is mixing matrix.

two steps ① prewhitening  $L$  ② component extraction  $W$   
 - to estimate  $A$ ,  $\hat{A} = (WL)^{-1} = L^{-1}W^T$ ,  $y_i = \hat{A}^{-1}(x_i - \bar{x})$

pre-white: find  $L$  that  $\text{Cov}(L(x_i - \mu)) = \text{Cov}(LX_i) = I$   
 comp extract: choose  $W$  that comp of  $Y_i = Wz_i$  independent (practically impossible)

Entropy -  $H(X) = -\int_{-\infty}^{\infty} f(x) \ln(f(x)) dx = -E(\ln f(x))$

**fastICA**  $X$ : pre-whitened centered  $W$ : estimated orthogonal  
 $K$ : pre-whitened matrix  $A$ : estimated transpose of  $A$

• Factor Analysis represent  $X$  in terms of unobserved factors.

$C = V\Lambda V^T$ ,  $\Lambda$  diagonal  $\lambda_1 \geq \lambda_2 \dots \lambda_r > 0$   
 $C = (\lambda_1^2 v_1 \dots \lambda_r^2 v_r \dots)^T = LL^T$   $L$  is  $p \times r$

Factor Analysis Model: `factanal()`  
 $X = \mu + LF + \epsilon = \mu + \sum F_k l_k + \epsilon$ ,  $F_k$ : factors,  $l_k$ : loadings

$\text{Cov}(F) = I$ ,  $\text{Cov}(\epsilon) = \Psi = \begin{pmatrix} \psi_1 & & \\ & \ddots & \\ & & \psi_p \end{pmatrix}$ ,  $\psi_i > 0$ ,  $C = LL^T + \Psi$   
 $L$  is loading matrix (not unique)

$Y = DX + a \Rightarrow Y = \underbrace{D\mu}_{\text{mean}} + a + \underbrace{DL}_{\text{loadings}} F + D\epsilon$

- Estimate  $L, \Psi$

① Iterative principal factor - until convergence  
 $\hat{L}\hat{L}^T + \hat{\Psi} \approx S = \frac{1}{n-1} \sum (x_i - \bar{x})(x_i - \bar{x})^T$

② Maximum likelihood estimation  
 Compute  $\mathcal{L}(L, \Psi)$ , test  $H_0: C = LL^T + \Psi$   
 Likelihood-Ratio test  
 large p-value: r factor is appropriate

- Scaling with rotation - easier to interpret ( $\ell_{kj} = 0$  for most  $j$ )

Orthogonal - LO where  $O$  orth  
 Oblique - LO where  $O$  invertible (all factors correlated)  
 $\hookrightarrow X = \mu + [LO][O^{-1}F] + \epsilon$ ,  $\text{Cov}(F^*) = [O^T O]^{-1}$

Varimax - max variance of squared elements of LO for  $O$  orth  
 Promax - power up varimax loading - result in non-orth loadings

Cluster Analysis - identify similar groups of observations

① Hierarchical clustering: use distance `hclust()`

Start with  $n$  clusters, then combine clusters

distance matrix  $D = \begin{pmatrix} 0 & d(x_1, x_2) & \dots \\ d(x_2, x_1) & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ d(x_n, x_1) & \dots & \dots & 0 \end{pmatrix}$

$d_s \leq d_a \leq d_c$   
 single linkage  $d_s(U, V) = \min\{d(x_i, x_j)\}$  produce less compact clusters  
 complete linkage  $d_c(U, V) = \max\{d(x_i, x_j)\}$  favour more compact clusters  
 average linkage  $d_a(U, V) = \text{average}\{d(x_i, x_j)\}$  compromise between extremes

$d(x, y) = d(y, x)$ ,  $d(x, y) \leq d(x, z) + d(z, y)$ ,  $d(x, y) = 0$  iff  $x = y$   
 $d(U, V) \leq d(U, W) + d(W, V)$

② Model-based clustering: assume a function

Mixture model:  $f(x) = \lambda_1 f_1(x; \theta_1) + \dots + \lambda_k f_k(x; \theta_k)$   
 $\lambda_i$  unknown,  $\sum \lambda_i = 1$ ,  $\theta_i$  unknown  
 Ideally  $f_i(x) f_j(x) \approx 0$  for  $i \neq j$

Example:  $f_j(x) = \frac{1}{(2\pi)^k \det(C_j)} \exp(-\frac{1}{2}(x - \mu_j)^T C_j^{-1}(x - \mu_j))$

k-means  $\sum \min\{\|x_i - \mu_1\|^2, \dots, \|x_i - \mu_k\|^2\}$  `kmeans()`  
 $\mu_i$  is the center of cluster

MLE estimation  $\mathcal{L} = \sum \ln\{\sum \lambda_j f_j(x_i; \theta_j)\}$   
 EM algorithm compute  $E_{(\theta, \lambda)}(\Delta_{ij} | X_i = x_i) = \hat{\delta}_{ij}(\theta, \lambda)$   
 assign observation  $i$  to cluster  $j$  if  $\hat{\delta}_{ij} > \hat{\delta}_{il}$  for all  $l \neq j$   
 ③  $\hat{\delta}_{ij} = \frac{\lambda_j f_j(x_i; \hat{\theta}_j)}{\sum_l \lambda_l f_l(x_i; \hat{\theta}_l)}$  ④  $\lambda_j = \frac{1}{n} \sum \delta_{ij}$  with  $\{\hat{\theta}_j\}$  updated



# Basics

$x_1 \dots x_n$  observations.  $x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{in} \end{bmatrix}$  p-dimension

- mean  $\mu = E(X) = \begin{pmatrix} E(x_1) \\ \vdots \\ E(x_p) \end{pmatrix}$   $E(AX+b) = AE(X)+b$
- covariance matrix  $C = Cov(X) = E(X-\mu)(X-\mu)^T = \begin{pmatrix} \sigma_1^2 & \dots & \sigma_{1p} \\ \vdots & \ddots & \vdots \\ \sigma_{p1} & \dots & \sigma_p^2 \end{pmatrix}$   $\sigma_{jk} = Cov(X_j, X_k)$   
 $Cov(AX+b) = ACA^T$   $\sigma_{j^2} = Var(X_j)$   
 $C$  non singular  $\Rightarrow (X-\mu)^T C^{-1} (X-\mu) \sim \chi^2(p)$   $\rightarrow a, b, c, s, s$  multinomial
- Estimations:  $\hat{\mu} = \bar{x} = \frac{1}{n} \sum x_i$   $\hat{C} = S = \frac{1}{n-1} \sum (x_i - \bar{x})(x_i - \bar{x})^T$
- correlation matrix  $R = D^{-1/2} C D^{-1/2}$ ,  $D = diag(C)$   $corr()$   $corr(x_1, x_2) = \frac{Cov(x_1, x_2)}{\sqrt{Var(x_1)Var(x_2)}}$
- concentration matrix:  $K = C^{-1}$  solve(corr())  
 $K$  determines dependency  $k_{ij} \neq 0 \Rightarrow$  dependence  $V(x_1, x_2) = V(x_1) + V(x_2) \pm 2Cov(x_1, x_2)$   
 $\Rightarrow f(x) = \frac{1}{(2\pi)^{p/2} |det(C)|^{1/2}} \exp(-\frac{1}{2}(x-\mu)^T C^{-1} (x-\mu))$

Multivariate Normal  $X \sim N_p(\mu, C)$ ,  $Y = AX+b \Rightarrow Y \sim N_r(A\mu+b, ACA^T)$   
 $X \sim N_p(\mu, C) \Rightarrow a^T X \sim N(a^T \mu, a^T C a)$

Conditional  $X_1 | X_2 = x_2 \sim N_r(\mu_{1|2}, C_{11|2})$   
 $\mu_{1|2} = \mu_1 + C_{12} C_{22}^{-1} (x_2 - \mu_2)$   
 $C_{11|2} = C_{11} - C_{12} C_{22}^{-1} C_{21}$

## Assess Normality

Mahalanobis distance:  $d_i = (x_i - \bar{x})^T S^{-1} (x_i - \bar{x})$  if  $x_i^2$  then normal  
 normal  $\Rightarrow$  straight line

Shapiro-Wilk test: if normal  $corr(x_i, \bar{x}) \approx 1$ , if not  $corr(x_i, \bar{x}) < 1$   
 $H_0$ : is normal  $< 0.05$  reject

Kurtosis  $K(X) = \frac{E((X-\mu)^4)}{Var(X)^2}$   $K(X) > 1$ ,  $K(aX+b) = K(X)$   
 $X \sim N(\mu, \sigma^2) \Rightarrow K(X) = 3$

$|K(X) - 3|$  measures non-normality

Scatterplot matrix pairs() Histogram hist()  
 $p$  increase, effectiveness decrease. Optimal bandwidth  $\textcircled{1} 3.49 \times SD \times n^{-1/5}$   
 $\textcircled{2} 2 \times IQR \times n^{-1/5}$

Curse of dimension:  $x_1 \dots x_n$  in  $\mathbb{R}^p$  as  $p$  increase, points become sparse.  
 Bi-plot: PCI vs PC2  
 vector indicate how  $x_i$  correlated with first two PCs  
 $cos(\theta) = corr(x_i, x_j)$

AIC =  $-2 \times \max \log \text{likelihood} + 2 \times$  number of parameter  
 BIC =  $-2 \times \max \log \text{likelihood} + \ln(n) \times$  number of parameter  
 choose model that minimize AIC and BIC

## Generalized Linear Model

Take  $H(x)$  cts.  $0 < H(x) < 1$ ,  $H(x)$  strictly increasing  
 Define  $P(G=1|X=x) = H(\beta_0 + x^T \beta)$   $\ln\left(\frac{P(G=1|X=x)}{P(G=0|X=x)}\right) = \beta_0 + x^T \beta$   
 $P(G=0|X=x) = 1 - H(\beta_0 + x^T \beta)$

Then  $H^{-1}(P(G=1|X=x)) = \beta_0 + x^T \beta$   
 link function

Link example Logit  $H^{-1}(x) = \ln(x/(1-x))$   
 Probit  $H^{-1}(x) = \Phi^{-1}(x)$  inverse normal  
 Log-log  $H^{-1}(x) = -\ln(-\ln(x))$

MLE:  $L(\beta_0, \beta) = \sum \{g_i(\beta_0 + x_i^T \beta) - \ln(1 + \exp(\beta_0 + x_i^T \beta))\}$

It satisfy  $\sum \{g_i - \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)}\} = 0$  ( $\sum x_i = 0$ )

- LDA  $k=2$  is logistic regression  
 LDA  $k=2$ ,  $P(G=1|X=x) = \frac{\lambda_1 f_1 / \lambda_0 f_0}{1 + \lambda_1 f_1 / \lambda_0 f_0}$   $\frac{\lambda_1 f_1}{\lambda_0 f_0} = \exp(\beta_0 + x^T \beta)$   
 $\beta = C^{-1}(\mu_1 - \mu_0)$   $\beta_0 = \frac{1}{2}(\mu_0^T C^{-1} \mu_0 - \mu_1^T C^{-1} \mu_1) + \ln(\lambda_1 / \lambda_0)$

LDA estimate mean covariance  
 Logistic regression use maximum likelihood

# Supervised Learning

## Classification

Given data, find a rule  $\phi: X \rightarrow G$   
 $X$  is range of possible  $x$ ,  $G$  is set of all possible classes

- ① Logistic Regression model  $P(G=g|X=x) = \Psi_g(x; \beta)$  glm()

$P(G=1|X=x) = \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)}$   $G = \{0, 1\}$

$P(G=g|X=x) = \frac{\exp(x^T \beta_g)}{1 + \exp(x^T \beta_1) + \dots + \exp(x^T \beta_k)}$   $G = \{0, 1, \dots, k\}$

Optimal classification rule  $R_g = \{x: \lambda_g f_g(x) > \lambda_j f_j(x) \forall j \neq g\}$   
 $\hat{\phi}(x) = \arg \max_g \lambda_g f_g(x)$ :  $g = 1, \dots, k$   
 (=g if  $\lambda_g f_g > \lambda_j f_j$  for all  $j \neq g$ )

Posterior distribution  $P(G=g|X=x) = \frac{\lambda_g f_g(x)}{\lambda_1 f_1 + \dots + \lambda_k f_k}$

- Application: multivariate normal

$\hat{\phi}(x) = g$  if  $\ln\left(\frac{f_g(x)}{f_j(x)}\right) > \ln\left(\frac{\lambda_j}{\lambda_g}\right)$  for all  $j \neq g$ .

$\ln\left(\frac{f_g(x)}{f_j(x)}\right) = \{x^T C^{-1} \mu_g - \frac{1}{2} \mu_g^T C^{-1} \mu_g\} - \{x^T C^{-1} \mu_j - \frac{1}{2} \mu_j^T C^{-1} \mu_j\}$

- ② Discriminant Analysis

- LDA linear discriminant lda() class: predicted class. posterior: possibility estimate

$\hat{\mu}_g = \frac{\sum \sum I(g_i=g) x_i}{\sum \sum I(g_i=g)}$   $\hat{C} = \frac{1}{n-k} \sum (x_i - \hat{\mu}_{g_i})(x_i - \hat{\mu}_{g_i})^T$

$\hat{\lambda}_g = \frac{1}{n} \sum I(g_i=g)$   $\hat{\phi}(x) = \arg \max_g \hat{d}_g(x)$

discriminant scores:  $\hat{d}_g(x) = x^T \hat{C}^{-1} \hat{\mu}_g - \frac{1}{2} \hat{\mu}_g^T \hat{C}^{-1} \hat{\mu}_g + \ln(\hat{\lambda}_g)$

LDA assumes  $C_1 = C_2 = \dots = C_k = C$

- QDA quadratic discriminant qda()

$X|G=g \sim N_p(\mu_g, C_g)$   $\hat{\phi}(x) = \arg \max_g \hat{d}_g(x)$

discriminant score:  $\hat{d}_g(x) = \ln(\hat{\lambda}_g) - \frac{1}{2}(x - \mu_g)^T \hat{C}_g^{-1} (x - \mu_g) - \frac{1}{2} \ln(\det(\hat{C}_g))$

QDA assumes  $C_1 \dots C_k$  are arbitrary

- QDA allows for flexible boundaries  
 QDA estimate more parameters, variance of  $\hat{C}_1 \dots \hat{C}_k > \hat{C}$  for LDA  
 QDA typically has lower bias due to flexibility  
 note: bias-variance trade-off

- Cross Validation - Leave out  $m$  obs - training set - estimate rule  
 - remain  $n-m$  obs - test set - estimate error rate

- ③ Tree-based classification, use half spaces

Example  $B_j = [x_2 < 5] \times [x_1 > 3]$

Advantage: no assumptions - flexibility - any

Disadvantage: complex, depends on stopping rule

MANOVA: assess feasibility of LDA,  $H_0: \mu_1 = \mu_2 = \dots = \mu_k$  manova()

$S_T = \frac{\sum n_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^T}{S_B} + \frac{\sum_{j=1}^n \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)^T}{S_W}$

high  $F \rightarrow$  reject  $H_0$ , low  $F \rightarrow$  fail to reject  $H_0$

Test statistics: Wilks Lambda / Pillai  
 Hotelling-Lawley trace / Roy's maximal root