

• \hat{S} = sample covariance matrix, $\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$, diagonal is sample variance
 • estimated covariance $S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$, diagonal is sample variance
 • corr: $R_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$, $\sigma_i = \sqrt{\lambda_i}$; $R = D^{-\frac{1}{2}} S D^{-\frac{1}{2}}$, $D^{-\frac{1}{2}} = \text{diag}(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_p})$
 • Euclidean dist: $d(x, y) = \sqrt{\sum_{k=1}^p (x_k - y_k)^2}$ \rightarrow est. of R , where trace R
 est. mean: $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$ (both est. are unbiased) $\Rightarrow D^{-\frac{1}{2}} C D^{-\frac{1}{2}}$
 Manhattan dist. $d_m(x_i, x_j) = \sum_{k=1}^p |x_{ik} - x_{jk}|$ dimensionless by $\frac{1}{s_k}$
 $N(\mu, \sigma^2)$, $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]$; $X \sim N(\mu, C) \Rightarrow f(x) = f(x_1, \dots, x_p) = \frac{1}{(2\pi)^{p/2} |C|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(x - \mu)^T C^{-1} (x - \mu)\right]$ positive definite
 Two approaches to check multi. normality: ① Mahalanobis distance / ② Ad hoc.
 $X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \sim p-r$, $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}$, $C = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix}$, $K = C^{-1}$ (concentration matrix)
 $x_1 | x_2 = x_3 \sim N_r(\mu_{1|2}, C_{1|2})$ where $\mu_{1|2} = \mu_1 + C_{12} C_{22}^{-1} (\mu_2 - \mu_1)$, $C_{1|2} = K_{11} - K_{12} K_{22}^{-1} K_{21}$ conditional mean
 $C_{1|2} = C_{11} - C_{12} C_{22}^{-1} C_{21}$ conditional covariance
 Graphical Lasso: K "sparse", maximize $\ln(|K|) - \text{trace}(K) - \lambda \|K\|_1$.
 scatterplots (2D), ① find interesting 2D projections / ② Grand Tour Method
 \Rightarrow (P) pairs! downside: miss out on higher dim structure.
 None-interesting projection = projections that are normal.
 Kurtosis, X, μ, σ^2 : $K(X) = E[(X - \mu)^4] / \sigma^4$, $K(aX + b) = K(X)$
 If $X \sim N(\mu, \sigma^2)$, $K(X) = 3$, long tailed dist'n. $K > 3$, multi-modal $K(X) > 3$
 If $X \sim N(\mu, \sigma^2)$, $K(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4 / (C^T S^2)$, find a to optimize $K(a)$
 $K(a) = \frac{1}{n} \sum_{i=1}^n (a^T x_i - a^T \bar{x})^4 / (C^T S^2)$, Density Estimation: $f_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$
 Types of kernel: ① Rectangle kernel, $K(x) = \frac{1}{2}$ for $|x| \leq 1$
 ② kernel $K(x) = -|x|$ for $|x| \leq 1$ ③ Gaussian kernel $K(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$
 ④ Epanechnikov kernel $K(x) = \frac{3}{4}(1-x^2)$, $|x| \leq 1$.
 Kernel density est for large P is problematic \Rightarrow use projection pursuit
 PCA approaches: ① Low rank matrix approximation / ② Projection pursuit \Rightarrow maximize
 assume $p \leq n$, $X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$ rank(X) = p , approx X by a matrix with variance
 rank $< p$. $\hat{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} = \sum_{i=1}^p \sum_{j=1}^p |X_{ij} - M_{ij}|^2$
 find \hat{X}^* to minimize $\|X - \hat{X}\|_F^2 = \sum_{i=1}^p \sum_{j=1}^p |X_{ij} - M_{ij}|^2$
 Singular Value Decomposition (SVD)
 Write $X = UDV^T$, U, V have orthonormal columns, $U = (u_1, \dots, u_p)$, $V = (v_1, \dots, v_p)$
 Now $D^* = \begin{pmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & d_p \end{pmatrix}$, $X^* = U D^* V^T$, rank(X^*) = $\min(\text{rank}(U), \text{rank}(D^*), \text{rank}(V)) = \text{rank}(D^*) = r$
 $X^* V = U D^* V^T V = U D^* = (d_1 u_1, \dots, d_r u_r, 0, \dots, 0)$
 From SVD \Rightarrow PCA: $\hat{X} = U D V^T$, $S = \frac{1}{n-1} \sum_{i=1}^n \sum_{j=1}^n (x_i - \bar{x})(x_j - \bar{x})^T = \frac{1}{n-1} V D U^T U D V^T = V \left(\frac{D^2}{n-1} \right) V^T$
 Sample cov. matrix (corr matrix if scaled)
 Columns of V are eigenvectors v_1, \dots, v_p of S with eigenvalues the diagonal elements of $D^2/(n-1) \Rightarrow d_1/(n-1), \dots, d_r/(n-1)$
 Principal Components: $y = \hat{X} V = \hat{X} (v_1, \dots, v_p) = (\hat{x}_1 v_1, \dots, \hat{x}_p v_p) \Rightarrow$ PC scores
 Prop. pursuit approach: x_1, \dots, x_n with sample varn S on correlation matrix R
 find a with $\|a\|^2 = a^T a = 1$ max. sample variance of $[a^T x_i]$:
 $\sum_{i=1}^n (a^T x_i - a^T \bar{x})^2 = a^T S a$, $S = V \Lambda V^T$, $\Lambda = D^2 / (n-1)$, $a^T S a = a^T V \Lambda V^T a = \sum_{k=1}^p \lambda_k (v_k^T a)^2$, $\Lambda = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda_p \end{pmatrix}$, max. at $a = v_1 / \|v_1\|$. R better than S
 $\sum_{k=1}^p \lambda_k (v_k^T a)^2$, $\Lambda = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda_p \end{pmatrix}$, max. at $a = v_1 / \|v_1\|$. R better than S
 $\max a^T S a$ s.t. $\|a\|^2 = 1 \Rightarrow a^T v_1 = 0 \Rightarrow a = v_2$ etc. dominates.
 note that $\sum_{k=1}^p (v_k^T x_i - v_k^T \bar{x})^2 = V^T S V = \lambda_k$ biplot (Important)
 PCs' variances, PCA \Rightarrow transforms original data to obtain uncorrelated variables.
 PC scores: columns of $\hat{X} V = U D = (u_1, \dots, u_p) \begin{pmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & d_p \end{pmatrix} = (du_1, \dots, du_p)$
 PC loadings: columns of $V = (v_1, \dots, v_p)$
 Multidimensional scaling: $D = (d_{ij})$, $d = (x_i - x_j)$, Recover (d_{ij}) from $B = X X^T = \begin{pmatrix} b_{11} & \dots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nn} \end{pmatrix}$
 distance matrix: $D = (d_{ij})$, $d = (x_i - x_j)$, Recover (d_{ij}) from $B = X X^T = \begin{pmatrix} b_{11} & \dots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nn} \end{pmatrix}$
 $d_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij}$, $SVD \Rightarrow X = U \Lambda^{1/2} V^T$, $X X^T = U \Lambda U^T$, $U = (u_1, \dots, u_p)$
 optimal k dim rep. of pts ($\lambda_1^{\frac{1}{2}} u_1, \lambda_2^{\frac{1}{2}} u_2, \dots, \lambda_k^{\frac{1}{2}} u_k$)
 how to get low dim rep for data? $B = U \Lambda U^T$, $\Lambda = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \lambda_k \end{pmatrix}$, $\lambda_1 \geq \lambda_2 \geq \dots$
 if $k=2$, plot $\lambda_1^{\frac{1}{2}} u_1$ v.s. $\lambda_2^{\frac{1}{2}} u_2$
 Andrew curves (idea: rep. each x_i as a function of $[0, 1]$) $g_i(t) = \sum_{k=1}^p x_{ik} \phi_k(t)$
 Need to choose ϕ_1, \dots, ϕ_p s.t. $g_i(t) = g_j(t)$ for $t \Rightarrow x_i = x_j$
 take $\phi_i(t) = \frac{1}{\sqrt{2}} \cos(i\pi t)$, $\phi_1(t) = \cos(2\pi t)$, $\phi_2(t) = \cos(4\pi t)$, \dots
 Practice: Q1. ① joint dist'n of multi-normal: "section". ② dist'n of $x_i + x_j$
 $x_3 + x_4 + x_5$, μ is the sum, variance $\sigma^2 = a^T \Lambda a$, a is vector. ③ conditioned of x_i , given $x_3 = -1$, know joint mean $\mu_{1|5}$, C_{15} , $\text{Corr}(x_1, x_5) = \rho_{15} = \frac{\mu_{1|5} - \mu_1}{\sqrt{1 - \rho_{15}^2}}$
 $\text{Var}(x_1 | x_5) = \mu_{1|5}^2 + \rho_{15}^2 \text{Var}(x_1)$, variance $\sigma_1^2(1 - \rho_{15}^2)$
 Q2: $R = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \rho & \dots & \rho \end{pmatrix}$ p-variate, x_1, \dots, x_n . Show one of PC has equal loadings, i.e. coeff. of all p variables are equal.
 Loadings are eigenvectors of correlation, their sum equal to 1
 if one of PC has equal loadings then $\vec{V} = (1, 1, \dots, 1)^T$ for example
 is an eigenvector of \hat{R} , verify $\hat{R} \vec{V} = \lambda \vec{V}$ for some λ .
 $\hat{R} \vec{V} = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \rho & \dots & \rho \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = (1 + (p-1)\rho) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

Example $p=2$, $K = \begin{pmatrix} 1 & p \\ p & 1 \end{pmatrix}$: eigenvalues solve $\det(K - \lambda I) = 0$

$$\lambda^2 - 2\lambda + 1 - \hat{p}^2 = 0 \Rightarrow \lambda_1 = 1 + |\hat{p}|, \lambda_2 = 1 - |\hat{p}|$$

eigenvectors (loadings) = $v_1 = \left(\frac{\sqrt{2}}{\text{sgn}(\hat{p})\sqrt{2}} \right)$ $v_2 = \left(\frac{\sqrt{2}}{-\text{sgn}(\hat{p})\sqrt{2}} \right)$

Joint distribution of $(X_1, X_2, X_5)'$ is 3-variate normal with mean $(\mu_1, \mu_2, \mu_5)'$. Covariance matrix $C_{15} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 5 \end{pmatrix}$. The distribution of $X_1 + X_2 + X_5$ is normal with $\mu = \mu_1 + \dots + \mu_5$, variance is $\sigma^2 = ((1)(1)(1))C((1)(1)(1))^T = 140$.

Conditional distribution: X_1 given $X_5 = -1$? X_1, \dots, X_5 ? $C_{15} = \begin{pmatrix} 55 & 6 \\ -6 & 60 \end{pmatrix}$

$\text{Cor}(X_1, X_5) = \rho_{15} = -6/\sqrt{55 \cdot 60}$. The conditional dist'n of X_1 given $X_5 = -1$ is normal with mean $\mu_{1|5} = \mu_1 + \rho_{15} \frac{\sigma_5}{\sigma_1} (\bar{x} - \mu_5) = 1 - \frac{6}{55} x$ and variance is $\sigma_{1|5}^2 = (1 - \rho_{15}^2) = 54.4$. Plug in $x = -1$, $\mu_{1|5} = 1.1$.

Correlation matrix? $K = C^{-1}$, linkage if not equal to 0.

$= \begin{pmatrix} 1 & \varphi & \cdot \\ \varphi & 1 & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$ show PC's have equal loadings. The loadings are eigenvectors of the correlation, standardized so that their sums of squares is 1. If one of the PCs has equal loadings then (e.g.) vector $v = (1, 1, \dots, 1)^T$ is an eigenvector of R . We just need to verify $Rv = \lambda v$ for some λ : $Rv = \begin{pmatrix} 1 & \varphi & \cdot \\ \varphi & 1 & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = ((1+\varphi)\varphi) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

b). Sesp p is even loadings of 1st 2 PCs are $(p^{\frac{1}{2}}, \dots, p^{\frac{1}{2}})^T$ & $(p^{\frac{1}{2}}, \dots, p^{\frac{1}{2}}, -p^{\frac{1}{2}})^T$ the square roots of the eigenvalues $\lambda_1^{\pm} \geq \dots \geq \lambda_r^{\pm}$ and ψ to be a diagonal matrix s.t. so that the diagonal elements of $LL^T + \psi$ are 1s. The same can be done with the sample covariance matrix making appropriate modifications.

distance matrix

height 6
1 2 3 4 5 6
1 2 3 4 5 6
2 3 4 5 6
3 4 5 6
4 5 6
5 6
6 7 8 9
7 8 9
8 9
9 10
10 11
11 12
12 13
13 14
14 15
15 16
16 17
17 18
18 19
19 20
20 21
21 22
22 23
23 24
24 25
25 26
26 27
27 28
28 29
29 30
30 31
31 32
32 33
33 34
34 35
35 36
36 37
37 38
38 39
39 40
40 41
41 42
42 43
43 44
44 45
45 46
46 47
47 48
48 49
49 50
50 51
51 52
52 53
53 54
54 55
55 56
56 57
57 58
58 59
59 60
60 61
61 62
62 63
63 64
64 65
65 66
66 67
67 68
68 69
69 70
70 71
71 72
72 73
73 74
74 75
75 76
76 77
77 78
78 79
79 80
80 81
81 82
82 83
83 84
84 85
85 86
86 87
87 88
88 89
89 90
90 91
91 92
92 93
93 94
94 95
95 96
96 97
97 98
98 99
99 100
single linkage
complete linkage

HW1 ~~(X₁, X₂)~~ given X₂ = 2, X₄ = 3, X₅ = -1
conditional distn ~~(X₁, X₂)~~ given X₂ = 2, X₄ = 3, X₅ = -1
 $\mu = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$, $C_{12}(2 \times 2)$, $C_{22} = (2 \times 3)$, $C_{12} = (2 \times 3)$
 $M_{12} = \mu_1 + C_{12} C_{22}^{-1} (X_2 - \mu_2)$, $C_{12}^{-1} = C_{11} - C_{12} C_{22}^{-1} C_{21} = (2 \times 2)$

Shapiro-Wilk's test: p-value < 0.05 not normal.

HW2
(a) If $g(t) = \frac{x_{11}}{\sqrt{p}} + x_{12} \sin(2\pi t) + x_{13} \cos(2\pi t) + \dots$
Show that $\int_0^P [g_i(t) - g_j(t)] dt = \sum_{k=1}^p (x_{ik} - x_{jk})^2$
Write $g_i(t) = \sum_{k=1}^p x_{ik} \phi_k(t)$ where $\phi_k(t) = \frac{1}{\sqrt{p}}, \phi_2, \dots, \phi_p$ are sines & cosines. Note that
 $\int_0^P \phi_k^2(t) dt = \frac{1}{p}$ & $\int_0^P \phi_k(t) dt = 0 \quad \forall k, l$
and $\int_0^P [g_i(t) - g_j(t)] dt = \int_0^P \sum_{k=1}^p \sum_{l=1}^p (x_{ik} - x_{jk})(x_{il} - x_{jl}) \phi_k(t) \phi_l(t) dt$
 $= \sum_{k=1}^p (x_{ik} - x_{jk}) \int_0^P \phi_k^2(t) dt = \frac{1}{p} \sum_{k=1}^p (x_{ik} - x_{jk})^2$

HW3
Sps S symmetric positive definite $p \times p$ with $S = V \Lambda V^T$ where Λ is diagonal with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. V is an orthogonal matrix with columns x_1, x_2, \dots, x_p . (a) Sps approximate S by $\psi + LL^T$ where $\text{diag } L = (\lambda_1^{\frac{1}{2}} v_1, \dots, \lambda_r^{\frac{1}{2}} v_r)$, ψ is diagonal with elements equal to those of S. If l_{ij} is $(i-j)$ th element of L, $\psi_{ii} = ?$

$s_{ii} = \psi_{ii} + \sum_{j=1}^r l_{ij}^2$, $\psi_{ii} = s_{ii} - \sum_{j=1}^r l_{ij}^2$

(b) $D = \psi + LL^T - S$. $d_{ij} = \sum_{k=1}^p d_{ijk}^2 \leq \lambda_{n+1}^2 + \dots + \lambda_p^2$. The diag of D are 0 by definition; if we define $D^* = LL^T - S$ then D & D^* have the same off-diag elements and so $\sum_{i=1}^p \sum_{j=1}^p d_{ij}^2 \leq \sum_{i=1}^p \sum_{j=1}^p (d_{ij}^*)^2$. Now $D^* = LL^T - S = V \Lambda^* V^T - V \Lambda V^T$ where Λ^* is $p \times p$ diagonal with $\lambda_1, \dots, \lambda_r, 0, \dots, 0$. Thus $D^* = -\sum_{j=r+1}^p \lambda_j V_j V_j^T$. The i-th row of D^* is $d_i^T = -\sum_{j=r+1}^p \lambda_j V_j$. and $d_i^T d_i = \sum_{j=r+1}^p \lambda_j V_j^2$. Therefore $\sum_i d_i^T d_i = \sum_{i=1}^p \lambda_i^2 \sum_{j=r+1}^p V_j^2 = \lambda_1^2 + \dots + \lambda_p^2$

(c) The result in (b) says $LL^T + \psi$ is a good approx to S if $\lambda_1^2 + \dots + \lambda_p^2$ is small. choose r s.t. $\frac{\lambda_1^2 + \dots + \lambda_p^2}{\lambda_1^2 + \dots + \lambda_p^2} \rightarrow 1$.

one vs. two factor model

Sometimes the loadings are not easy/obvious to interpret. P-value higher \rightarrow fit better.

LDA vs QDA: CV leave one out, 10-fold \rightarrow train & test, but when predicting, no CV. rate slightly higher. (model for QDA misclassification is more complicated than LDA & the variability in QDA is more complicated than LDA)

Estimating additional parameters may lead to WTA having a higher misclassification rate.

$$\text{HW4: single-linkage: } \min d(A, B) \cdot d(A, C) = \min \{ \min_{a \in A} d(a, b), \min_{a \in A} d(a, c) \}$$

$$= \min_{a \in A, b \in B \cup C} d(a, b)$$

$$= d(A, B \cup C)$$

The other equality follows since $\min(x, y) = (x+y)/2 - |x-y|/2$

complete-linkage is similar.

Show $d(A, C) \leq d(A, B) + d(B, C)$ complete-linkage.

disjoint clusters A, B, C.

For, $a, b, c \in A, B, C$ $d(a, c) \leq d(a, b) + d(b, c)$

$$\max_{a \in A, b \in B, c \in C} d(a, c) = d(A, C) \leq \max_{a, b, c} (d(a, b) + d(b, c))$$

$$\text{likewise } \max_{a, b, c} (d(a, b) + d(b, c)) \leq \max_{a \in A} d(a, b) + \max_{b \in B} d(b, c) \\ = d(A, B), d(B, C)$$

This won't work for single-linkage.

• If EM cannot separate ~~#~~ k-groups after enough iterations, it is possible that using mixture model is not good. The initial clusters provided by k-means are not clear for EM (local extrema) to isolate the groups.

QDA vs LDA CLASSIFICATIONS

LDA assumes same var-cov. matrices between input variables of classes. The assumption is important for classification stage of the analysis. If the matrices substantially differ, observations will tend to be assigned to the class where variability is greater. To overcome the problem, QDA which is a modification of LDA allowing for the above heterogeneity of classes' cov matrices, "more effective", so clusters have extremely different shapes. If similar shape, LDA > QDA b/c QDA has more misclassifications b/c it counts more variabilities.

PCA vs ICA

Additionally, VS FA! FA, PCA & ICA are all related, three of them seek basis vectors that the data is projected against, such that you maximize [insert criteria here] think of the basis vectors as just encapsulating linear combinations.

Criteria?

Second-order criteria:
In PCA, you are finding basis vectors that explain the variance of your data. The first (i.e. highest ranked) basis vector is going to be the one that best fits all the variance from your data. The second one also has this criteria, but must be orthogonal to the first, and so on & so forth. (Turns out those basis vectors for PCA are nothing but the eigenvectors of your data's covariance matrix).

In FA, you are finding basis vectors that explain the variance of your data. The first (i.e. highest ranked) basis vector is going to be the one that best fits all the variance from your data. The second one also has this criteria, but must be orthogonal to the first, and so on & so forth. (Turns out those basis vectors for PCA are nothing but the eigenvectors of your data's covariance matrix).

Higher order Criteria's.

In ICA, you are again finding basis vectors, but this time, you want basis vectors that give a result, such that

this resulting vector is one of the independent components of the original data. You can do this by maximization of the absolute value of normalized kurtosis — a 4th order statistic. That is, you project your data on the same vector, and measure the kurtosis of the result. You change your basis vector a little (usually through gradient ascent), and then measure the kurtosis again, etc, etc. Eventually you will happen onto a basis that gives you a result that has the highest possible kurtosis, and this is your independent component. In graph if the data is like

orthogonal as well.

PCA#1 & PCA#2

but ICA would be like

independent

Another version of difference b/w FA & PCA

PCA involves extracting linear composites of observed variables.

Factor analysis is based on a formal model of predicting observed from theoretical latent variables factors.

1. Run FA if you assume or wish to test a theoretical model for of latent factors causing observed variables.

2. Run PCA if you want to simply reduce your correlated obs variables to a smaller set of important independent composite variables.

FA is a generalization of PCA which is based explicitly on ML. Like PCA, each data point is assumed to arise from sampling a pt in a subspace & then perturbing it with full-dimensional Gaussian noise. The difference is that FA allows the noise to have an arbitrary diagonal covariance matrix while PCA assumes the noise is spherical. In addition to estimating the subspace, FA estimates the noise covariance matrix.

Cluster analysis VS FA

Short: Cluster analysis is about grouping subjects, FA is about grouping variables.

Long: Clustering is asking how many groups the items can be split into. Factor analysis is asking whether the features of you have are explained as combinations of some smaller set of features.

You don't get any groupings out of a factor analysis, and you don't get any dimensionality reduction by doing clustering.

Cluster & Classification

Both split data into groups. With classification the groups (or classes) are specified before hand, with each training data instance belonging to a particular ~~data~~ class. It is the association b/w the instances features & the class they belong to that classification algorithms are supposed to learn. (Supervised learning)

With clustering the groups (or clusters) are based on the similarities of data instances

is used in training & the clustering algorithm is supposed to learn the grouping itself. (unsupervised learning)

Notes

Covariance C is non-negative definite $\alpha^T C \alpha \geq 0 \forall \alpha$

conditional distribution

$$X_1 | X_2 \sim N(\mu_1 + \rho \frac{X_2 - \mu_2}{\sigma_2}, \sigma_1^2(1 - \rho^2))$$

$$X \sim N_p(\mu, C) \quad \alpha^T X \sim N(\alpha^T \mu, \alpha^T C \alpha)$$

$$\alpha^T (X - \mu)^T C^{-1} (X - \mu) \sim \chi^2(p)$$

$$R \leftrightarrow C, \text{cor}_{ij} = \frac{\alpha_{ij}}{\sqrt{\alpha_{ii}\alpha_{jj}}}$$

Simple visualization of multivariate data (2D)-scatterplot

pros: - easy, more informative version of a correlation matrix; non-linear pairwise dependencies revealed.

cons: - miss out on higher dimensional structure.

projection pursuit: "interesting" = "non-normality"

4th order statistic Kurtosis: $K(X) = \frac{E((X - \mu)^4)}{\sigma^4}, K(X+b) \neq K(X)$

find a to maximize/minimize kurtosis

In PCA, λ 's are variances.

e.g. of PCA. x_1, \dots, x_n obs, p variables.

$X = (x_i^T)$ centred \tilde{X} , PCA transforms original data to obtain p uncorrelated variables

$$\tilde{X}V = UD = (u_1, \dots, u_p) \begin{pmatrix} d_1 & & \\ & \ddots & \\ & & d_p \end{pmatrix} \quad n \times p \quad p \times p \quad \text{diag}$$

$$= (d_1 u_1, \dots, d_p u_p), \quad d_1^2 + \dots + d_p^2 = \sum d_i^2 > 0$$

Covariance matrix

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{n-1} \tilde{X}^T \tilde{X} = VV^T \text{ where } \Lambda = \frac{\sum d_i^2}{n-1}$$

first PC maximizes $\alpha^T S \alpha$ s.t. $\alpha^T \alpha = 1 \Rightarrow \alpha = v_1, V = (v_1, \dots, v_p)$

k -th PC maximizes $\alpha^T S \alpha$ s.t. $\alpha^T v_j = 0 \& \alpha^T \alpha = 1$ for $j = 1, \dots, k-1$

PC scores: columns of $\tilde{X}V$; PC loadings columns of $V = (v_1, \dots, v_p)$

Naive PCA model

$$X = \mu + L E, \quad E \sim N_p(0, I) \quad (\text{with noise})$$

Factor analysis model $X = \mu + L E + \xi$

where ① $E \sim N_p(0, I)$ ② $\xi \sim N_p(0, \Psi)$ where $\Psi = \begin{pmatrix} \psi_{11} & & \\ & \ddots & \\ & & \psi_{pp} \end{pmatrix}$

③ Components of ξ are independent components of E

Note: Can also specify model in terms of $\text{Cov}(E)$, $\text{Cov}(\xi)$

and μ, E uncorrelated.

$$\Rightarrow X \sim N_p(\mu, LL^T + \Psi)$$

Why factor analysis?

① Unobserved factors explain relationships b/w variables

② Factor analysis is invariant (equiv) to changes in

scale for each variable.

③ Relationship to graphical models:

for one-factor model: ~~exists~~ link exists b/w F & X_j if $f_{1j} \neq 0$.

for multi-factor model: X_i, X_j are conditional independent

given F, \dots, L (loading is not uniquely determined)

Both a feature & a bug.

$$\text{Generally, } X = \mu + L E + \xi, \quad \text{Cov}(E) = I, \quad \text{Cov}(\xi) = \Psi = \begin{pmatrix} \psi_{11} & & \\ & \ddots & \\ & & \psi_{pp} \end{pmatrix}$$

$$\text{Cov}(X) = LL^T + \Psi$$

$$\text{Cov}(E) = \begin{pmatrix} \psi_{11} & & \\ & \ddots & \\ & & \psi_{pp} \end{pmatrix}$$

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \approx LL^T + \Psi$$

$$\hat{\Psi} = \text{diag}(S - LL^T) \quad (\text{but assuming } S - LL^T \text{ diagonal and positive})$$

$$LL^T \approx S - \hat{\Psi} = V \Lambda V^T = (v_1, \dots, v_p) \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} \begin{pmatrix} v_1^T \\ \vdots \\ v_p^T \end{pmatrix}$$

Principal factor Method

Algorithm:

$$\textcircled{1} \text{ Set } \hat{\Psi} = 0 \quad \textcircled{2} \text{ set } S - \hat{\Psi} = V \Lambda V^T \& \hat{L} = (v_1, \dots, v_p)$$

$$\textcircled{3} \text{ Set } \hat{\Psi} = \text{diag}(S - \hat{L} \hat{L}^T) \quad \textcircled{4} \text{ Iterate } \textcircled{2} \text{ till convergence}$$

Hierarchical clustering methods

Start with n clusters each with 1 observation. Group together similar obs.

① Single ② complete ③ average linkage

General algorithm:

①

② Search D to find nearest pair of clusters (depends on approach) call these $U+V$

③ Merge $U+V$ into a single cluster $\tilde{V}^T \tilde{V}$ and redefine D (linkage)

④ Repeat ② & ③ until we have a single cluster

Dendrogram

k-means clustering, EM:

$f(x) = \lambda_1 f_1(x, \theta_1) + \dots + \lambda_k f_k(x, \theta_k)$
mixture of multivariate normals.

$$f_j(x, \mu_j, \Sigma_j) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp(-\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j))$$

Classification

Data problem: Given $(g_1, X_1), \dots, (g_n, X_n)$ as training data, determine a classification rule with small error rate.

Say R_k are boundaries of each class.

$$\text{error rate} = P(\text{error}) = \sum_{j=1}^k P(G=j)P(\text{error}|G=j)$$

$$= \sum_{j=1}^k \int_{R_j} P(X \notin R_j | G=j) = \sum_{j=1}^k \int_{R_j} 1 - \sum_{i \neq j} f_i(x) dx$$

$= 1 - \sum_{j=1}^k \int_{R_j} f_j(x) dx \quad \rightarrow \text{correct classification rate}$

Optimal classification rule: (Baye's Rule)

Classify x to group j if classifier

$$f_j(x) > f_i(x) \quad \forall i \neq j$$

$$\frac{f_j(x)}{f_i(x)} > \frac{\lambda_i}{\lambda_j} \quad \forall i \neq j$$

e.g. example from prac-final:

$$\int_0^1 \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \frac{1}{\pi} r dr d\theta = \int_0^1 \frac{\pi}{2} r dr = \frac{1}{2}$$

- naive error rate: $\frac{1}{n} \sum_{i=1}^n I(g_i \neq \hat{g}_i)$
- resubstitution error rate \downarrow biased downside.
- fix: CV (leave one out), estimate classifier from remaining obs. & look at error rate for obs. not used to est. the classifier.
- Two approaches (Optimal classification)

① Minimize $P(\text{error})$.

$$\begin{aligned} & \lambda_j f_j(x) > \lambda_i f_i(x) \\ \Rightarrow & \text{Look at conditional dist'n of } G \text{ given } X=x \\ \text{Bayes Thm } P(G=j|X=x) &= \frac{P(G=j) f_j(x)}{P(X=x)} \\ = \frac{P(G=j) f_j(x)}{P(X=x)} &= \frac{\lambda_j f_j(x)}{\lambda_1 f_1(x) + \dots + \lambda_K f_K(x)} \end{aligned}$$

max this value. (same as D?)

(But more information we get)