

Resumen AED

1. PCA. Buscamos transformar p variables en q combinaciones lineales ortogonales.
 $X' = (x_1, x_2, \dots, x_p)$ $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. El resumo componente principal.

$$y_i = e_i' X = e_{i1} x_1 + e_{i2} x_2 + \dots + e_{ip} x_p$$

En general, $y = P' X$ $\text{cov}(y) = \Sigma_y = P' \Sigma_x P$

$$\Sigma_x = P \Lambda P'$$
 , luego $\Sigma_y = \text{diag}(\lambda_i)$

$$\text{var}(y_i) = \lambda_i$$

$$\text{cov}(y_i, y_k) = 0$$

luego, $\sigma_{11} + \dots + \sigma_{pp} = \sum \sigma_{ii} = \lambda_1 + \dots + \lambda_p = \sum \lambda_i = \text{var}(y_i)$

→ Proporción total de varianza: $\text{prop}(y_k) = \frac{\lambda_k}{\sum \lambda_i}$

→ correlación: $r_{y_i, x_k} = \frac{\text{cov}(y_i, x_k)}{\sqrt{\lambda_i} \sqrt{\sigma_{kk}}} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$

Para interpretación sacar p y Diap 22 para ejemplo

1.1 Población multivariada Normal $X \sim N_p(\mu, \Sigma)$

→ Elipsoide con centro en μ y con ejes $u_i = \sqrt{\lambda_i} e_i$

• Cuando las varianzas son muy diferentes, entonces estandarizamos

$$Z_i = \frac{x_i - \mu_i}{\sqrt{\sigma_{ii}}}, \quad Z = V^{-1/2} (X - \mu), \quad V^{-1/2} = \text{diag}\left(\frac{1}{\sqrt{\sigma_{ii}}}\right)$$

$$E(Z) = 0 \text{ y } \Sigma_Z = R \Rightarrow \tilde{y} = \tilde{e}_i' Z, \quad \sum \text{var}(Z_i) = \text{tr}(R) = p$$

Ej. Men's back data, diapositiva 32

1.2 Sample principal components. $\hat{\Sigma}$ estimador de Σ , $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$.

$\hat{\Sigma}$ tiene λ es de $\hat{\Sigma}$

$$\hat{y}_i = \hat{e}_i' X = \hat{e}_{i1} x_1 + \dots + \hat{e}_{ip} x_p, \quad \text{var}(\hat{y}_i) = \hat{\lambda}_i, \quad \text{cov}(\hat{y}_i, \hat{y}_k) = 0$$

• $\text{tr}(\hat{\Sigma}) = \sum \hat{\sigma}_{ii} = \sum \hat{\lambda}_i$ prop-var: $\frac{\hat{\lambda}_i}{\sum \hat{\lambda}_k}$

correlación: $r_{\hat{y}_i, x_k} = \frac{\sqrt{\hat{\lambda}_i}}{\sqrt{\hat{\sigma}_{kk}}} \hat{e}_{ik}$ • Si x 's son estandar, $r_{\hat{y}_i, x_k} = \sqrt{\hat{\lambda}_i} \hat{e}_{ik}$

* las componentes principales son el vector direccional del elipsoide.

Ej. Otro ejemplo de interpretación, diapositiva 59.

2. FA: Buscamos un modelo. $X \rightarrow \text{obs con } n \times p \Sigma$. Factores comunes F_i . Erros E_i

→ Modelo: $X - \mu = L \times F + E$ $E(F) = 0$ $\text{cov}(F) = I$

$$E(E) = 0$$

$$\text{cov}(E) = \Psi = \begin{bmatrix} \psi_1 & 0 \\ 0 & \psi_p \end{bmatrix}$$

$$\text{cov}(F, E) = 0$$

$$\text{cov}(X) = \Sigma = LL' + \Psi$$

$$\text{cov}(X, F) = L$$

$$\text{var}(X_i) = \hat{e}_{i1}^2 + \dots + \hat{e}_{im}^2 + \psi_i, \quad \text{cov}(X_i, X_k) = \hat{e}_{i1} \hat{e}_{k1} + \dots + \hat{e}_{im} \hat{e}_{km}, \quad \text{cov}(X_i, F_j) = \hat{e}_{ij}$$

• Comunalidad: $h_i^2 = \sum \hat{e}_{ik}^2$, $\hat{\sigma}_{ii} = h_i^2 + \psi_i$

* Si $m=p$, entonces $\Sigma = LL'$, es decir, Ψ es nula.

2.1 Rotación: Sea T una matriz $m \times m$ ortogonal. $X - \mu = L^* F^* + E$

$$L^* = LT$$

$$F^* = T' F$$

$$E(F^*) = 0$$

$$\text{cov}(F^*) = I$$

En general, $F \approx F^* \quad L \neq L^* \quad \Sigma = (L^*)(L^*)' + \Psi$

2.2 Método de Estimación \rightarrow PCA: Eigen-descomposición de Σ . MLE: estimador de máxima verosimilitud.

2.2.1 PCA: Sacamos autovalores y autovectores de Σ . $\lambda_1 \geq \dots \geq \lambda_p$

$$\tilde{L} = [\sqrt{\lambda_1} \hat{e}_1 \dots \sqrt{\lambda_m} \hat{e}_m] \quad \hat{\Psi} = \Sigma - \tilde{L} \tilde{L}^T, \quad \hat{\Psi}_{ii} = \Sigma_{ii} - \sum \hat{L}_{ij}^2 = S_{ii} - h_i^2$$

$$\bullet \text{ SS}(\Sigma - (\tilde{L} \tilde{L}^T + \hat{\Psi})) \leq \lambda_{m+1}^2 + \dots + \lambda_p^2, \quad \text{SS}(A) = \text{tr}(AA^T)$$

En general,

proporción total de la varianza para factor j :

$$\left\{ \begin{array}{ll} \frac{\hat{\lambda}_j}{S_{11} + \dots + S_{pp}} & \text{para } \Sigma \\ \frac{\hat{\lambda}_j}{p} & \text{para } R \end{array} \right.$$

Estimación de factores:

- Algoritmo:
1. Obtener estimación inicial de $\Psi \rightarrow$ SMC de X_i y las otras. El SMC es la diagonal de $R^{-1} \rightarrow \Psi^*$.
 2. Encontrar $\hat{\lambda}_i^*$ y \hat{e}_i^* de $R - \text{diag}(\Psi_i^*)$. wego.
 $L^* = (\sqrt{\lambda_1^*} \hat{e}_1^*, \dots, \sqrt{\lambda_m^*} \hat{e}_m^*)$
 3. $\Psi_i^* = 1 - h_i^{*2} = 1 - \sum \hat{L}_{ij}^{*2}$
 4. Repetir 2 y 3 hasta obtener convergencia.

2.2.2 MLE: Obtenemos el estimador de máxima verosimilitud

$$L(\mu, \Sigma) = (2\pi)^{-\frac{(n-1)p}{2}} |\Sigma|^{-\frac{(n-1)}{2}} \exp\left[-\frac{1}{2} (Z^T (\Sigma^{-1} (Z(X_j - \bar{x})(X_j - \bar{x})^T))\right] \\ \times (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)}$$

Proporción total de la varianza para factores $j \rightarrow \frac{\hat{L}_{1j}^2 + \hat{L}_{2j}^2 + \dots + \hat{L}_{pj}^2}{S_{11} + S_{22} + \dots + S_{pp}}$

Con $Z = V^{-1/2} (X - \mu)$ $\hat{\beta} = V^{-1/2} \hat{\Sigma} V^{-1/2}$ $\hat{\Psi}_2 = V^{-1/2} \hat{\Psi} V^{-1/2}$
 $L_2 = V^{-1/2} L$ $\hat{\beta} = \hat{L}_2 \hat{L}_2^T + \hat{\Psi}_2$ $\hat{h}_i^2 = \hat{L}_{i1}^2 + \dots + \hat{L}_{im}^2$

proporción total de la varianza para factor $j \rightarrow \frac{\hat{L}_{1j}^2 + \dots + \hat{L}_{pj}^2}{p}$

factores $\hat{L} = V^{-1/2} \hat{L}_2$ $\hat{\Phi} = V^{-1/2} \hat{\Phi}_2 V^{-1/2}$

2.3. Rotación: Hay 2 tipos ortogonal y oblicua.

2.3.1 Ortogonal: T matriz ortogonal, $\Sigma = L^T L^T + \Psi$ $L^T = LT$
con las estimaciones es igual.

2.4. Factor Scores: $\hat{f}_j^* = T^T f_j$ $f_j = \hat{L}_1^T \hat{\Sigma}^{-1} (X_j - \bar{x}) = \hat{L}_1^T R^{-1} z_j$
 $z_j = \hat{\Sigma}^{-1/2} (X_j - \bar{x})$ $\hat{\beta} = \hat{L}_2 \hat{L}_2^T + \hat{\Psi}_2$

3. Clasificación y Discriminantes.

3.1 clasificación 2 poblaciones $X' = [X_1, \dots, X_p]$ $f_1(x)$ y $f_2(x)$. π_1 y π_2

Si $x_0 \in R_1 \Rightarrow x_0 \in \pi_1$ Si $x_0 \in R_2 \Rightarrow x_0 \in \pi_2$

$$P(2|1) = \int_{R_2} f_1(x) dx$$

$$P(1|2) = \int_{R_1} f_2(x) dx$$

costo de clasificación:

	π_1	π_2
π_1	0	$c(2 1)$
π_2	$c(1 2)$	0

clasificación, pred \rightarrow columnas
verdadero \rightarrow filas.

→ valor esperado del costo: $ECM = C(211)P(112)P_1 + C(112)P(112)P_2$
 * la idea es minimizar el costo, una buena clasificación tiene ECM pequeño

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \frac{C(112)}{C(211)} \cdot \frac{P_2}{P_1} \quad R_2: \frac{f_1(x)}{f_2(x)} < \frac{C(112)}{C(211)} \cdot \frac{P_2}{P_1}$$

→ Bayes: $\frac{P_1 f_1(x)}{P_1 f_1(x) + P_2 f_2(x)} \geq \frac{P_2 f_2(x)}{P_1 f_1(x) + P_2 f_2(x)} \Rightarrow \delta_1(x) \geq \delta_2(x)$

3.2 Clasificaciones de 2 poblaciones multivariadas normales:

• Cada una es normal. $f_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_i)' \Sigma^{-1} (x - \mu_i) \right]$

→ Con $\Sigma_1 = \Sigma_2 = \Sigma$. $\delta_1 (\mu_1 - \mu_2)' \Sigma^{-1} x_0 - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln \left(\frac{C(112)P_2}{C(211)P_1} \right)$

entonces $x_0 \in \Pi_1$. Si no, entonces $x_0 \in \Pi_2$

En general, $\delta_1 (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}_{pool}^{-1} x_0 - \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}_{pool}^{-1} (\bar{x}_1 + \bar{x}_2) \geq \ln \left(\frac{C(112)P_2}{C(211)P_1} \right)$

entonces $x_0 \in \Pi_1$. Si no, entonces $x_0 \in \Pi_2$

$$\hat{\Sigma}_{pool} = \frac{(n_1 - 1) \hat{\Sigma}_1 + (n_2 - 1) \hat{\Sigma}_2}{n_1 + n_2 - 2}$$

→ Fisher's Approach: $\hat{y}_0 = (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}_{pool}^{-1} x_0$
 $\hat{m} = \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}_{pool}^{-1} (\bar{x}_1 + \bar{x}_2)$

• Si $\hat{y}_0 \geq \hat{m} \Rightarrow x_0 \in \Pi_1$ • Si $\hat{y}_0 < \hat{m} \Rightarrow x_0 \in \Pi_2$

→ Con $\Sigma_1 \neq \Sigma_2$.

Si $-\frac{1}{2} x_0' (\Sigma_1^{-1} - \Sigma_2^{-1}) x_0 + (\mu_2' \Sigma_1^{-1} - \mu_2' \Sigma_2^{-1}) x_0 - k \geq \ln \left(\frac{C(112)P_2}{C(211)P_1} \right)$

$\Rightarrow x_0 \in \Pi_1$, Si no, $x_0 \in \Pi_2$, $k = \frac{1}{2} \ln \left(\frac{|\Sigma_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\mu_1' \Sigma_1^{-1} \mu_1 - \mu_2' \Sigma_2^{-1} \mu_2)$

En general,

Si $-\frac{1}{2} x_0' (\hat{\Sigma}_1^{-1} - \hat{\Sigma}_2^{-1}) x_0 + (\bar{x}_1' \hat{\Sigma}_1^{-1} - \bar{x}_2' \hat{\Sigma}_2^{-1}) x_0 - k \geq \ln \left(\frac{C(112)P_2}{C(211)P_1} \right)$

$\Rightarrow x_0 \in \Pi_1$, Si no, $x_0 \in \Pi_2$, $k = \frac{1}{2} \ln \left(\frac{|\hat{\Sigma}_1|}{|\hat{\Sigma}_2|} \right) + \frac{1}{2} (\bar{x}_1' \hat{\Sigma}_1^{-1} \bar{x}_1 - \bar{x}_2' \hat{\Sigma}_2^{-1} \bar{x}_2)$

3.3 Evaluar la clasificación:

APER = $\frac{n_{11} + n_{22}}{n_1 + n_2}$

	Pred	
	Π_1	Π_2
Π_1	n_{11}	n_{12}
Π_2	n_{21}	n_{22}

$n_{1c} \rightarrow$ Elementos en Π_1 que son de Π_1

$n_{2c} \rightarrow$ Elementos en Π_2 que son de Π_2

$n_{1m} = n_1 - n_{1c}$ $n_{2m} = n_2 - n_{2c}$

3.4. Varias poblaciones: Para más de 2 categorías los costos siempre son = $d_i(x)$. $P(K|L) = \int_{R_K} f_L(x) dx$ $ECM = \left(\sum P(K|L) C(K|L) \right) \sum P_L$

→ clasificación con ECM. $x_0 \in \Pi_k$ si $P_k f_k(x) \geq P_i f_i(x) \forall i \neq k$

QDA: $d_i(x) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) + \ln P_i$

QDA-sample: $d_i(x) = -\frac{1}{2} \ln |\hat{\Sigma}_i| - \frac{1}{2} (x - \bar{x}_i)' \hat{\Sigma}_i^{-1} (x - \bar{x}_i) + \ln P_i$

LDA: $d_i(x) = \mu_i' \Sigma^{-1} x - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \ln P_i$

LDA-sample: $d_i(x) = \bar{x}_i' \hat{\Sigma}_{pool}^{-1} x - \frac{1}{2} \bar{x}_i' \hat{\Sigma}_{pool}^{-1} \bar{x}_i + \ln P_i$

$$\hat{\Sigma}_{pool} = \frac{(n_1 - 1) \hat{\Sigma}_1 + \dots + (n_p - 1) \hat{\Sigma}_p}{n_{11} + \dots + n_{pp} - p}$$

$x \in \Pi_k$ si $\max \{d_i(x)\} = d_k(x)$

4. Clustering: hay 2 métodos \Rightarrow k -means
jerárquico.

4.1 k -means. En el taller preparaval 3 también dice el paso a paso

Algoritmo: 1. Asigne una etiqueta random de 1 a k a cada observación.
Estos son clusters iniciales.

2. Itere hasta que dejen de cambiar los clusters:

a. Para cada k cluster saque un centroide, a partir de las medias según la etiqueta.

b. Cambie los clusters al que esté más cerca con distancia euclidiana.

Ej $n=3$. Diapositiva 17

4.2 Jerárquico: \Rightarrow complete. max
Single. min
Average. Avg

Algoritmo: 1. Empiece con n observaciones y las distancias.
Con la matriz de disimilitud.

2. For $i=n, n-1, \dots, 2$:

a. Fusione las 2 obs más cerca y saque las nuevas obs y sus distancias con las demás

b. saque la nueva matriz de disimilitud.

* la matriz de disimilitud se saca con distancia euclidiana, si son muchos datos se usa matriz de correlación R