

Análisis de componentes ppales (ACP) o (PCA)

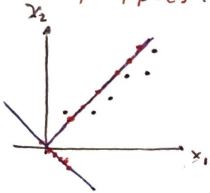
Thursday, March 31, 2022 6:52 PM

Análisis de componentes ppales (ACP) o (PCA)

Buscamos explicar la estructura de varianza-cov. de los datos y reducir dimensionalidad.

- Si hay p -vars, necesito p componentes pples xa capturar TODA la info.
- En general, se puede capturar gran parte de la info con $K \ll p$ comp. pples.
- La BD se convierte en n mediciones de K "vars".
- Normalm. se usa PCA antes de otros métodos (reg. lin)

Comp. pples: Las c.p son combinaciones lin de las p v.a's. x_1, \dots, x_p q' max la varianza.



Interp geométrica.

El c.p define una dir dada $a_i = (a_{i1}, \dots, a_{ip})$ y se busca encontrar las dirs q' explican la mayor Q de variabilidad de los datos.

- Aquí tenemos 2 vars y 2 dirs (1, 1)
- La projec° de los vals c. en esas dirs están como: ...
- Busco cuál me explica mejor los vals haciendo uso de la projec°.

Considera el vec aleat. $X = \begin{bmatrix} x_1 \\ x_p \end{bmatrix}$ con matriz de cov Σ y vals propios $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

Las comb. lin:

$$Y_1 = a_1' X = a_{11}x_1 + \dots + a_{1p}x_p$$

$$Y_p = a_p' X = a_{p1}x_1 + \dots + a_{pp}x_p$$

$$\text{var}(Y_i) = a_i' \Sigma a_i$$

$$\text{cov}(Y_i, Y_k) = a_i' \Sigma a_k$$

Las c.p son las comb. lin no correlacionadas (indep) q' maximizan la varianza.

C. P :
• Indep.
• Ortogonales entre sí.

Ej la 1era c.p = $\max \text{var}(Y_i) = a_i' \Sigma a_i$

(ACP sólo sirve para cuantitativas).

LAS COMPONENTES SON LOS VECES PROPIOS

Se tiene:

- 1era c.p = $a_1' X$ donde a_1 max $\text{var}(a_1' X)$

- i-ésima c.p = $a_i' X$ donde a_i max $\text{var}(a_i' X)$

Sujeto a $\text{cov}(a_i' X, a_k' X) = 0, k < i$
y a_i ortogonal a a_1, \dots, a_k .

Sea Σ la matriz de varianza-cov de X
Suponga q' Σ tiene los pares de vals, vect. propios:

$$(\lambda_1, e_1), \dots, (\lambda_p, e_p) \quad \lambda_1 \geq \dots \geq \lambda_p \geq 0.$$

* La i-ésima c.p es:

$$Y_i = e_i' X$$

$$\text{var}(Y_i) = e_i' \Sigma e_i = \lambda_i \quad i = 1, \dots, p$$

$$\text{cov}(Y_i, Y_k) = 0 \rightarrow \text{Las c.p son indep.}$$

* Si $Y_i = Y_k$ es como si fueran la misma.

Sea Σ la matriz de var-cov de X , con pares de vals y vecs propios:

$$(\lambda_1, e_1), \dots, (\lambda_p, e_p) \quad \lambda_1 \geq \dots \geq \lambda_p \geq 0.$$

Sean Y_1, \dots, Y_p las c.p. Entonces,

$$\begin{aligned} \sigma_{11} + \dots + \sigma_{pp} &= \text{var}(X_1) + \dots + \text{var}(X_p) = \Sigma \text{var}(X_i) \\ &= \lambda_1 + \dots + \lambda_p = \Sigma \text{var}(Y_i) \end{aligned}$$

Obs:

1. La varianza total poblacional se conserva.
2. La propor. de la variabilidad total poblacional explicada por la k -ésima C.P. es: $\frac{\lambda_k}{\lambda_1 + \dots + \lambda_p}$ \rightarrow la C.P. k explica este % de los datos.
3. En muchos casos, para las primeras C.P. esta propor. es del 0.5 - 0.9 (con las 4 primeras ya explico poco más del 70% de la info).
- (La α de C.P. depende del nivel de correla. de las vars.)
4. $e_i = [e_{ik}]$, $|e_{ik}|$ "mide" la importancia de la var k en la C.P. i .

Sea $Y_i = e_i' X$, $Y_p = e_p' X$ las C.P., entonces:

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{Xk}}} \rightarrow \text{correla. entre mi nueva } Y_i \text{ y mi og } X_k.$$

- ρ grande: X_k tiene mucho peso en Y_i (tiene info).
- $\rho > 0$: Si en X_k hay un val alto, espero un val alto en la C.P. Y_i .
- $\rho < 0$: Si X_k tiene val alto, espero q' C.P. Y_i tenga uno bajo.

Estandarizar?

Si no estandarizo, puedo estar dándole más del peso q' cuben a las vars con mucha variabilidad.

No obstante, deshacerme de esa mayor variabilidad, podría estar quitándome info.

$$\begin{aligned} \text{Se tienen las vars: } -z_1 &= \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} & Z &= (V^{1/2})^{-1} (X - \mu) \\ -z_p &= \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}} & & \text{matriz de desv. st.} \end{aligned}$$

$$E[Z] = 0; \text{ cov}(Z) = I_p$$

Los comp. ppls de Z se obtienen con los vals y vect propios de P .

La i -ésima C.P. de Z está dada por:

$$\begin{aligned} Y_i &= e_i' Z = e_i' (V^{1/2})^{-1} (X - \mu) \\ \sum_{i=1}^p \text{var}(Y_i) &= \sum \text{var}(Z_i) = 1 \cdot p = p. \end{aligned}$$

$$\rho_{Y_i, Z_k} = e_{ik} \sqrt{\lambda_i}$$

Prop. de variabilidad explicada por la C.P. k

$$= \frac{\lambda_k}{p} = \frac{\lambda_k}{\sum \text{var}(Y_i)}$$

