

Análisis de Componentes Principales (ACP)

Teoría y Práctica

Dr. Bartolo Villar

November 28, 2021

TICs, ENES-UNAM

Introducción

- El **Análisis de Componentes Principales (ACP)** es una técnica matemática de **aprendizaje no supervisado**.
- Consiste en explicar la estructura de **varianza-covarianza** de un conjunto de variables a través de **pocas combinaciones lineales** de las variables originales.
- Generalmente se utiliza para
 - reducción de la dimensionalidad de los datos,
 - interpretación,
 - sirve como un paso intermedio en investigación.
- Los k componentes principales reemplazan a las p variables originales.
- El ACP únicamente depende de la matriz de varianzas-covarianza, Σ , o en su defecto de la matriz de correlaciones, ρ .

Formulación matemática

Sea $\mathbf{X}' = (X_1, X_2, \dots, X_p)$ un vector aleatorio con matriz de covarianzas Σ a su vez con eigenvalores $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

Considere las siguientes combinaciones lineales

$$Y_1 = \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Y_2 = \mathbf{a}'_2 \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

$$\vdots \qquad \qquad \qquad \vdots$$

$$Y_p = \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

$$\text{var}(Y_i) = \mathbf{a}'_i \Sigma \mathbf{a}_i \qquad i = 1, 2, \dots, p$$

$$\text{covar}(Y_i, Y_k) = \mathbf{a}'_i \Sigma \mathbf{a}_k \qquad i, k = 1, 2, \dots, p$$

- Los componentes principales son aquellas **combinaciones lineales no correlacionadas** (ortogonales), $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ cuyas varianzas son las más grandes.
- El primer componente principal (CP1) es aquella combinación lineal con **máxima varianza**

$$\text{var}(\mathbf{Y}_1) = \mathbf{a}_1' \Sigma \mathbf{a}_1$$

- Note que la varianza puede incrementarse si se multiplica a \mathbf{a}_1 por una constante, para evitar esto se impone la siguiente restricción

$$\mathbf{a}_1' \mathbf{a}_1 = 1$$

$CP_1 \ Y_1$

Combinación lineal $\mathbf{a}'_1 \mathbf{X}$ que maximiza la $\text{var}(\mathbf{a}'_1 \mathbf{X})$ sujeto a que $\mathbf{a}'_1 \mathbf{a}_1 = 1$.

$CP_2 \ Y_2$

Combinación lineal $\mathbf{a}'_2 \mathbf{X}$ que maximiza la $\text{var}(\mathbf{a}'_2 \mathbf{X})$ sujeto a que $\mathbf{a}'_2 \mathbf{a}_2 = 1$, y la $\text{cov}(\mathbf{a}_1 \mathbf{X}, \mathbf{a}_2 \mathbf{X}) = 0$.

\vdots

$CP_i \ Y_i$

Combinación lineal $\mathbf{a}'_i \mathbf{X}$ que maximiza la $\text{var}(\mathbf{a}'_i \mathbf{X})$ sujeto a que $\mathbf{a}'_i \mathbf{a}_i = 1$, y la $\text{cov}(\mathbf{a}_i \mathbf{X}, \mathbf{a}_k \mathbf{X}) = 0$ para $i < k$.

Resultado

Sea Σ la matriz de covarianzas asociada con el vector aleatorio

$\mathbf{X}' = (X_1, X_2, \dots, X_p)$. Sea Σ tal que tenga pares de eigenvalores-eigenvectores $(\lambda_1, \mathbf{e}_1) \dots (\lambda_p, \mathbf{e}_p)$ con $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. Entonces el i -ésimo componente principal está dado por

$$\begin{aligned} Y_i &= \mathbf{e}_i' \mathbf{X} \\ &= e_{i1}X_1 + e_{i2}X_2 + \dots + e_{ip}X_p \end{aligned}$$

$$\text{var}(\mathbf{Y}_i) = \mathbf{e}_i' \Sigma \mathbf{e}_i = \lambda_i \quad i = 1, 2, \dots, p$$

$$\text{cov}(\mathbf{Y}_i, \mathbf{Y}_k) = \mathbf{e}_i' \Sigma \mathbf{e}_k = 0 \quad i \neq k.$$

Varianzas e eigenvalores

Resultado

Sea Σ la matriz de covarianzas asociada con el vector aleatorio

$\mathbf{X}' = (X_1, X_2, \dots, X_p)$. Sea Σ tal que tenga pares de eigenvalores-eigenvectores $(\lambda_1, \mathbf{e}_1) \dots (\lambda_p, \mathbf{e}_p)$ con $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$.

Sean $\mathbf{Y}_1 = \mathbf{e}_1 \mathbf{X} \dots \mathbf{Y}_p = \mathbf{e}_p \mathbf{X}$ sus componentes principales. Entonces

$$\begin{aligned}\sigma_{11} + \dots + \sigma_{pp} &= \sum_{i=1}^p \text{var}(X_i) \\ &= \lambda_1 + \dots + \lambda_p \\ &= \sum_{i=1}^p \text{var}(Y_i)\end{aligned}$$

Varianza del k-ésimo componente

- Del resultado anterior podemos obtener la proporción de varianza explicada por el k-ésimo componente principal:

$$\left(\begin{array}{c} \text{proporción} \\ \text{de varianza del k-ésimo} \\ \text{componente} \end{array} \right) = \frac{\lambda_k}{\lambda_1 + \dots + \lambda_p} \quad \forall k = 1, \dots, p$$

- En la práctica, generalmente con los primeros 3 o 4 componentes principales se explica al menos el 80% de la varianza total de X .
- Por otra parte, cada componente de \mathbf{e} merece inspeccionarlos. La magnitud de e_{ik} mide la importancia de la k-ésima variable al i-ésimo componente

$$e_{ik} \propto \rho(Y_i, X_k)$$

Correlación entre Y_i y X_k

Resultado

Si $\mathbf{Y}_1 = \mathbf{e}_1 \mathbf{X} \dots \mathbf{Y}_p = \mathbf{e}_p \mathbf{X}$ son los componentes principales, con Σ siendo la matriz de covarianzas, entonces

$$\rho(Y_i, X_k) = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1, 2, \dots, p,$$

son los coeficientes de correlación entre Y_i y X_k .

Ejemplo

La matriz de covarianzas para $\mathbf{X} = (X_1, X_2, X_3)$ es

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Se puede verificar que

$$\lambda_1 = 5.828$$

$$\lambda_2 = 2$$

$$\lambda_3 = 0.172$$

$$\begin{aligned} \mathbf{e}'_1 &= [-0.383 \quad 0.924 \quad 0] \\ \mathbf{e}'_2 &= [0 \quad 0 \quad 1] \\ \mathbf{e}'_3 &= [0.924 \quad 0.383 \quad 0] \end{aligned}$$

Ejemplo...

Los componentes principales son:

$$Y_1 = \mathbf{e}_1' \mathbf{X} =$$

$$Y_2 = \mathbf{e}_2' \mathbf{X} =$$

$$Y_3 = \mathbf{e}_3' \mathbf{X} =$$

$$\text{var}(Y_1) = \text{var}(-0.383X_1 - 0.924X_2)$$

$$\text{cov}(Y_1, Y_2) = \text{cov}(-0.383X_1 - 0.924X_2, X_3)$$

$$\begin{aligned}\sigma_{1,1} + \sigma_{2,2} + \sigma_{3,3} &= 1 + 5 + 2 \\ &= \lambda_1 + \lambda_2 + \lambda_3 \\ &= 5.828 + 2 + 0.172\end{aligned}$$

Note que los dos primeros componentes principales explican el 98% de la varianza contenida en \mathbf{X} :

$$\frac{\lambda_1 + \lambda_2}{\lambda_1 + \lambda_2 + \lambda_3} = 0.98$$

$$\rho(Y_1, X_1) = \frac{e_{11}\sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} =$$

$$\rho(Y_1, X_2) = \frac{e_{12}\sqrt{\lambda_1}}{\sqrt{\sigma_{12}}} =$$

Representación gráfica del ACP

ACP con variables estandarizadas

$$Z_1 = \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}}$$

\vdots

$$Z_p = \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}}$$

En notación matricial

$$\mathbf{Z} = (\mathbf{V}^{1/2})^{-1}(\mathbf{X} - \boldsymbol{\mu}) \text{ donde } \mathbf{V}^{1/2} = \begin{bmatrix} \sqrt{\sigma_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{\sigma_{22}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\sigma_{pp}} \end{bmatrix}$$

Con $E(\mathbf{Z}) = \mathbf{0}$

$$\text{cov}(\mathbf{Z}) = (\mathbf{V}^{1/2})^{-1} \boldsymbol{\Sigma} (\mathbf{V}^{1/2})^{-1} = \boldsymbol{\rho}.$$