

MAT-269: Sesión 18, Método de Componentes Principales I

Felipe Osorio

fosorios.mat.utfsm.cl

Departamento de Matemática, UTFSM



Objetivo:

Reemplazar observaciones p -dimensionales por k combinaciones lineales de las variables donde k es mucho más pequeño que p . La elección de k debe explicar razonablemente la proporción de dispersión total $\text{tr } S$.

Definición 1:

Sea \mathbf{x} un vector aleatorio con media $\boldsymbol{\mu}$ y dispersión $\boldsymbol{\Sigma}$. Considere $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_p)$ una matriz ortogonal, tal que

$$\mathbf{T}^\top \boldsymbol{\Sigma} \mathbf{T} = \boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p),$$

donde $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$. Sea

$$\mathbf{y} = \mathbf{T}^\top (\mathbf{x} - \boldsymbol{\mu}),$$

es decir $y_j = \mathbf{t}_j^\top (\mathbf{x} - \boldsymbol{\mu})$, es llamado el j -ésimo componente principal de \mathbf{x} para $j = 1, \dots, p$ y $\mathbf{z}_j = y_j / \lambda_j^{1/2}$ es llamado el j -ésimo componente principal estandarizado.



Resultado 1:

Los y_j son no correlacionados y $\text{var}(y_j) = \lambda_j$. En efecto,

$$\begin{aligned}\text{Cov}(\mathbf{y}) &= \text{Cov}(\mathbf{T}^\top (\mathbf{x} - \boldsymbol{\mu})) = \mathbf{T}^\top \text{Cov}(\mathbf{x})\mathbf{T} \\ &= \mathbf{T}^\top \boldsymbol{\Sigma}\mathbf{T} = \boldsymbol{\Lambda}.\end{aligned}$$

Además, se tiene que

$$\text{Cov}(\mathbf{z}) = \text{Cov}(\boldsymbol{\Lambda}^{-1/2}\mathbf{y}) = \boldsymbol{\Lambda}^{-1/2}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^{-1/2} = \mathbf{I}.$$



Propiedad 1:

Los componentes principales $y_j = \mathbf{t}_j^\top (\mathbf{x} - \boldsymbol{\mu})$, $j = 1, \dots, p$ tienen las siguientes propiedades:

- (i) Para cualquier vector \mathbf{a}_1 de largo $\|\mathbf{a}_1\| = 1$, $\text{var}(\mathbf{a}_1^\top \mathbf{x})$ alcanza su valor máximo λ_1 cuando $\mathbf{a}_1 = \mathbf{t}_1$.
- (ii) Para cualquier vector unitario tal que $\mathbf{a}_j^\top \mathbf{t}_i = 0$ ($i = 1, \dots, j - 1$) $\text{var}(\mathbf{a}_j^\top \mathbf{x})$ alcanza su valor máximo λ_j cuando $\mathbf{a}_j = \mathbf{t}_j$.
- (iii) $\sum_{j=1}^p \text{var}(y_j) = \sum_{j=1}^p \text{var}(x_j) = \text{tr } \boldsymbol{\Sigma}$.



Previo:

Considere $\mathbf{T}^\top \mathbf{A} \mathbf{T} = \mathbf{\Lambda}$ con \mathbf{T} matriz ortogonal y $\mathbf{\Lambda}$ matriz diagonal. Sea

$$\mathbf{x} = \mathbf{T} \mathbf{y} = y_1 \mathbf{t}_1 + \cdots + y_p \mathbf{t}_p.$$

Entonces

$$\frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\mathbf{y}^\top \mathbf{T}^\top \mathbf{A} \mathbf{T} \mathbf{y}}{\mathbf{y}^\top \mathbf{y}} = \frac{\mathbf{y}^\top \mathbf{\Lambda} \mathbf{y}}{\mathbf{y}^\top \mathbf{y}} = \frac{\sum_{i=1}^p \lambda_i y_i^2}{\mathbf{y}^\top \mathbf{y}},$$

como $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p > 0$, tenemos

$$\frac{\sum_{i=1}^p \lambda_i y_i^2}{\mathbf{y}^\top \mathbf{y}} \leq \lambda_1 \frac{\sum_{i=1}^p \mathbf{y}^\top \mathbf{y}}{=} \lambda_1,$$

con la igualdad sólo si $y_1 = 1, y_2 = 0, \dots, y_p = 0$, es decir para $\mathbf{x} = \mathbf{t}_1$.



Demostración:

Note que $\text{var}(\mathbf{a}_1^\top \mathbf{x}) = \mathbf{a}_1^\top \text{Cov}(\mathbf{x}) \mathbf{a}_1 = \mathbf{a}_1^\top \boldsymbol{\Sigma} \mathbf{a}_1$, lo que permite mostrar (i).

Para verificar (ii), como $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_p$ son mutuamente ortogonales, ellos forman una base para \mathbb{R}^p . Tenemos que \mathbf{a}_j puede ser expresado como

$$\mathbf{a}_j = c_j \mathbf{t}_j + c_{j+1} \mathbf{t}_{j+1} + \dots + c_p \mathbf{t}_p,$$

además $\mathbf{a}_j \mathbf{a}_j = 1 = \sum_{r=j}^p c_r^2 \mathbf{t}_r^\top \mathbf{t}_r = \sum_{r=j}^p c_r^2$, y

$$\begin{aligned} \text{var}(\mathbf{a}_j^\top \mathbf{x}) &= \mathbf{a}_j^\top \left(\sum_{r=j}^p c_r \boldsymbol{\Sigma} \mathbf{t}_r \right) = \mathbf{a}_j^\top \sum_{r=j}^p c_r \lambda_r \mathbf{t}_r = \left(\sum_{r=j}^p c_r \mathbf{t}_r \right)^\top \left(\sum_{r=j}^p c_r \lambda_r \mathbf{t}_r \right) \\ &= \sum_{r=j}^p c_r^2 \lambda_r \geq \lambda_j \sum_{r=j}^p c_r^2 = \lambda_j, \end{aligned}$$

con la igualdad si y sólo si $\mathbf{a}_j = \mathbf{t}_j$. De este modo, $\text{var}(\mathbf{a}_j^\top \mathbf{x})$ es maximizada en $\mathbf{a}_j = \mathbf{t}_j$.

Finalmente, (iii) sigue, pues:

$$\sum_{j=1}^p \text{var}(y_j) = \text{tr } \boldsymbol{\Lambda} = \text{tr } \boldsymbol{\Lambda} \mathbf{T}^\top \mathbf{T} = \text{tr } \mathbf{T} \boldsymbol{\Lambda} \mathbf{T}^\top = \text{tr } \boldsymbol{\Sigma}$$



Interpretación:

- ▶ (i) y (ii) indica que la componente y_1 es la combinación lineal normalizada de los elementos de $\mathbf{x} - \boldsymbol{\mu}$ con varianza máxima λ_1 .
- ▶ Ahora $\mathbf{a}^\top \mathbf{x}$ es no correlacionado con y_1 , esto es

$$\begin{aligned}\text{Cov}(\mathbf{a}^\top \mathbf{x}, \mathbf{t}_1^\top (\mathbf{x} - \boldsymbol{\mu})) &= \text{Cov}(\mathbf{a}^\top \mathbf{x}, \mathbf{t}_1^\top \mathbf{x}) = \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{t}_1 \\ &= \lambda \mathbf{a}^\top \mathbf{t}_1 = 0,\end{aligned}$$

si y sólo si $\mathbf{a} \perp \mathbf{t}_1$. De este modo, desde (ii) y_2 es la combinación lineal no correlacionada con y_1 con varianza máxima ($= \lambda_2$).

- ▶ En general y_j es la combinación lineal no correlacionada con y_1, \dots, y_{j-1} con varianza máxima λ_j .



Interpretación:

- ▶ (iii) provee una técnica simple para decidir cual k debe ser escogido. La razón

$$r_k = \frac{\lambda_j}{\sum_{r=1}^p \lambda_r} = \frac{\text{var } y_j}{\text{tr } \Sigma},$$

mide la contribución de y_j a $\text{tr } \Sigma$ la **variación total** de \mathbf{x} . De este modo, podemos incorporar componentes sucesivos y_1, y_2, \dots y detener en y_k cuando r_k sea cercano a la unidad.

- ▶ Si se estandariza los x_j y se trabaja con $x_j^* = (x_j - \mu_j)/\sigma_{jj}^{1/2}$. Entonces $\text{Cov}(\mathbf{x}^*) = \mathbf{R}$ la matriz de correlación. Se puede extraer un nuevo conjunto de componentes principales,

$$\mathbf{y}^* = \mathbf{L}^\top \mathbf{x}^*,$$

donde \mathbf{L} es ortogonal y $\mathbf{L}^\top \mathbf{R} \mathbf{L}$ es diagonal. Sin embargo \mathbf{y}^* en general difiere de \mathbf{y} .



Componentes principales muestrales

En la práctica μ y Σ no son conocidos y deben ser estimados desde $\mathbf{x}_1, \dots, \mathbf{x}_n$. Sea

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}, \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top,$$

y suponga $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ y $\hat{\mathbf{T}} = (\hat{\mathbf{t}}_1, \dots, \hat{\mathbf{t}}_p)$ valores propios y matriz ortogonal de vectores propios de $\hat{\boldsymbol{\Sigma}}$.

Podemos definir el vector de componentes principales (scores) para cada observación

$$\mathbf{y}_i = \hat{\mathbf{T}}^\top (\mathbf{x}_i - \bar{\mathbf{x}}), \quad i = 1, \dots, n,$$

obteniendo la matriz de datos

$$\mathbf{Y}^\top = (\mathbf{y}_1, \dots, \mathbf{y}_n) = \hat{\mathbf{T}}^\top (\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}).$$



Frecuentemente se prefiere usar S el estimador insesgado de Σ en lugar de $\widehat{\Sigma}$ para definir las componentes principales. En este caso

$$S\widehat{t}_j = \frac{n}{n-1}\widehat{\Sigma}\widehat{t}_j = \left(\frac{n}{n-1}\widehat{\lambda}_j\right)\widehat{t}_j,$$

y los valores propios de S son $n\widehat{\lambda}_j/(n-1)$.

Si se utiliza $\widehat{\lambda}_j/\text{tr}\widehat{\Sigma}$ para determinar la magnitud relativa de un valor propio, entonces el factor de escala $n/(n-1)$ cancela y ambos enfoques (usando $\widehat{\Sigma}$ o S) son idénticos.



Resultado 2:

Sea $\mathbf{A} \in \mathbb{R}^{p \times k}$ y considere $\mathbf{y}_i(k) = \mathbf{g}(\mathbf{x}_i - \bar{\mathbf{x}})$, $i = 1, \dots, n$, para cualquier función $\mathbf{g} : \mathbb{R}^p \rightarrow \mathbb{R}^k$ tal que $\bar{\mathbf{y}}(k) = \mathbf{0}$. Si f es estrictamente creciente e invariante bajo transformaciones ortogonales. Entonces,

$$f\left\{\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{A}\mathbf{y}_i(k))(\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{A}\mathbf{y}_i(k))^\top\right\},$$

es minimizada con respecto a \mathbf{A} y \mathbf{g} cuando

$$\mathbf{A} = (\hat{\mathbf{t}}_1, \dots, \hat{\mathbf{t}}_k) = \hat{\mathbf{T}}_1, \quad \mathbf{g}(\mathbf{x} - \bar{\mathbf{x}}) = \hat{\mathbf{T}}_1^\top (\mathbf{x} - \bar{\mathbf{x}}).$$

Previo (Seber, 1984; pp. 177):

Sea $\mathbf{T} = (\mathbf{T}_1, \mathbf{T}_2)$ donde $\mathbf{T}_1 = (\mathbf{t}_1, \dots, \mathbf{t}_k)$ definida desde $\mathbf{T}^\top \Sigma \mathbf{T} = \mathbf{\Lambda}$. Entonces

$$f(\mathbf{\Delta}) = f(\text{Cov}(\mathbf{x} - \mathbf{A}\mathbf{y}(k))),$$

es minimizada cuando $\mathbf{A}\mathbf{y}(k) = \mathbf{T}_1 \mathbf{T}_1^\top \mathbf{x}$.



Demostración:

Sea \mathbf{v} un vector aleatorio tomando el valor \mathbf{x}_i ($i = 1, \dots, n$) con probabilidad $\frac{1}{n}$, $\mathbf{w} = \mathbf{v} - \bar{\mathbf{x}}$ y $\mathbf{y}(k) = \mathbf{g}(\mathbf{w})$. Entonces $\boldsymbol{\eta} = \mathbf{E}(\mathbf{v}) = \sum_{i=1}^n \mathbf{x}_i \frac{1}{n} = \bar{\mathbf{x}}$, $\mathbf{E}(\mathbf{w}) = \mathbf{0}$, $\mathbf{E}(\mathbf{y}(k)) = \bar{\mathbf{y}}(k) = \mathbf{0}$ y

$$\begin{aligned}\text{Cov}(\mathbf{w}) &= \text{Cov}(\mathbf{v}) = \mathbf{E}\{(\mathbf{v} - \boldsymbol{\eta})(\mathbf{v} - \boldsymbol{\eta})^\top\} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\eta})(\mathbf{x}_i - \boldsymbol{\eta})^\top \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top = \hat{\boldsymbol{\Sigma}},\end{aligned}$$

y de ahí que

$$\begin{aligned}f\left\{\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{A}\mathbf{y}_i(k))(\mathbf{x}_i - \bar{\mathbf{x}} - \mathbf{A}\mathbf{y}_i(k))^\top\right\} \\ = f\left\{\mathbf{E}[(\mathbf{w} - \mathbf{A}\mathbf{y}(k))(\mathbf{w} - \mathbf{A}\mathbf{y}(k))^\top]\right\} = f(\text{Cov}(\mathbf{w} - \mathbf{A}\mathbf{y}(k))),\end{aligned}$$

que es minimizada para todos los vectores $\mathbf{y}(k)$ tales que $\mathbf{E}(\mathbf{y}(k)) = \mathbf{0}$ cuando $\mathbf{A}\mathbf{y}(k) = \hat{\mathbf{T}}_1 \hat{\mathbf{T}}_1^\top \mathbf{w}$ esta ecuación se satisface si hacemos $\mathbf{A} = \hat{\mathbf{T}}_1$ y $\mathbf{g}(\mathbf{w}) = \hat{\mathbf{T}}_1^\top \mathbf{w}$.



Calculando las componentes principales muestrales

En la práctica se recomienda obtener las componentes principales usando la descomposición valor singular (SVD) de la matriz de datos centrados,

$$\mathbf{Z} = \begin{pmatrix} \mathbf{z}_1^\top \\ \vdots \\ \mathbf{z}_n^\top \end{pmatrix},$$

con $\mathbf{z}_i = \mathbf{x}_i - \bar{\mathbf{x}}$, para $i = 1, \dots, n$. Es fácil notar que:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{Z}^\top \mathbf{Z}.$$

De este modo, obtener la SVD de $\mathbf{Z} = \mathbf{U} \mathbf{D} \mathbf{V}^\top$ con $\mathbf{U} \in \mathbb{R}^{n \times p}$ tal que $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$, $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$ de valores singulares ($d_1 \geq d_2 \geq \dots \geq d_p$). Es fácil notar que

$$d_1^2/(n-1), \dots, d_p^2/(n-1),$$

son los valores propios de \mathbf{S} , mientras que las columnas de \mathbf{V} corresponden a los vectores propios de \mathbf{S} . De este modo, los scores son calculados como sigue:

$$\mathbf{Y} = \mathbf{Z} \mathbf{V}.$$

