

Análisis de Componentes Principales (Primera Parte)

Mathias Bourel

IMERL - Facultad de Ingeniería, Universidad de la República, Uruguay

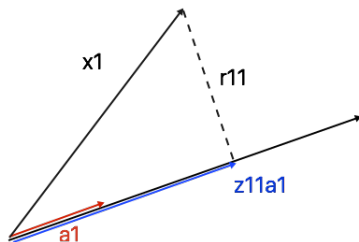
May 2, 2019

Enfoque geométrico

Consideramos la matriz de datos $X \in \mathcal{M}_{n \times p}$ **centrada**, es decir que la media de cada columna es 0. Queremos encontrar un subespacio de dimensión menor que p que represente de manera adecuada los datos. Más precisamente queremos encontrar un subespacio de dimensión menor que p tal que cuando proyectamos los individuos sobre él, la estructura se distorciona lo menos posible.

Consideremos una recta por el origen (subespacio de dimensión 1) generada por un vector $a_1 \in \mathbb{R}^p$ unitario. Si consideramos un individuo x_i su proyección sobre el subespacio generado por a_1 es

$$z_{i1} a_1 = \frac{\mathbf{x}_i' a_1}{\|a_1\|^2} a_1 = \mathbf{x}_i' a_1 a_1 = a_1' \mathbf{x}_i a_1$$



Enfoque geométrico

Si queremos minimizar $\sum_{i=1}^n r_{i1}^2 = \sum_{i=1}^n \|\mathbf{x}_i - z_{i1}\mathbf{a}_1\|^2 = \sum_{i=1}^n (\mathbf{x}_i - z_{i1}\mathbf{a}_1)'(\mathbf{x}_i - z_{i1}\mathbf{a}_1)$ observamos que por el teorema de Pitágoras

$$\mathbf{x}'_i \mathbf{x}_i = z_{i1}^2 + r_{i1}^2$$

y entonces

$$\sum_{i=1}^n \mathbf{x}'_i \mathbf{x}_i = \sum_{i=1}^n z_{i1}^2 + \sum_{i=1}^n r_{i1}^2$$

Como el término $\sum_{i=1}^n \mathbf{x}'_i \mathbf{x}_i$ es constante, minimizar $\sum_{i=1}^n r_{i1}^2$ equivale a maximizar $\sum_{i=1}^n z_{i1}^2$ que no es otra cosa que la varianza muestral **de los datos proyectados** dado que los datos son centrados. En efecto

$$\sum_{i=1}^n z_{i1} = \sum_{i=1}^n a'_1 \mathbf{x}_i = a'_1 \left(\sum_{i=1}^n \mathbf{x}_i \right) = a'_1 \bar{\mathbf{x}} = 0$$

Objetivos

- Reducir el número de variables sin perder (demasiada) información: al proyectar los n individuos sobre un espacio de dimensión l con $l < p$ tal que la dispersión en el espacio proyectado sea máxima.
- Simplificar la descripción del conjunto de datos. Analizar la estructura y relación de las observaciones y de las variables.
- Las componentes principales deben tener varianza máxima (mayor información relacionado con mayor variabilidad).

Para eso:

- Cada componente principal es una combinación lineal de las variables originales.

$$\text{Probabilidad : } z_j = a_{j1}x_1 + a_{j2}x_2 + \cdots + a_{jp}x_p = \quad \forall j = 1, \dots, l, \quad l < p$$

$$\text{Estadística : } z_j = a_{j1}x_1 + a_{j2}x_2 + \cdots + a_{jp}x_p = \mathbf{X}a_j = \begin{pmatrix} \mathbf{x}'_1 a_j \\ \mathbf{x}'_2 a_j \\ \vdots \\ \mathbf{x}'_n a_j \end{pmatrix} = \begin{pmatrix} z_{1j} \\ z_{2j} \\ \vdots \\ z_{nj} \end{pmatrix}$$

- Las componentes principales son no correlacionadas dos a dos, y de esta manera eliminamos información repetida:

x_1, \dots, x_p correladas $\rightarrow z_1, \dots, z_l$ **incorreladas**

Cálculo de las componentes

- 1 Vamos a imponer que $\|a'_j\| = 1 \quad \forall j = 1, \dots, p$
- 2 Vamos a buscar a_1 tal que z_1 tenga la mayor varianza y $\|a_1\| = 1$.
- 3 Vamos a buscar a_2 tal que z_2 sea incorrelada con z_1 , con varianza menor que z_1 y $\|a_2\| = 1$.
- 4 Vamos a buscar a_3 tal que z_3 sea incorrelada con z_1 y z_2 , con varianza menor que z_1 y z_2 y $\|a_3\| = 1$.
- 5 ...

Cálculo de las componentes

Sea Σ la matriz de covarianzas de \mathbf{X} . Habitualmente se usa la matriz de correlaciones ya que se estandariza los datos (cada columna tiene media cero y desvío 1).

- 1 Como las variables originales tienen media cero entonces el vector $z_1 = \mathbf{X}a_1$ tiene también media cero y su varianza es $\text{Var}(z_1) = \frac{1}{n}z_1'z_1 = \frac{1}{n}a_1'\mathbf{X}'\mathbf{X}a_1 = a_1'\Sigma a_1$. Para maximizar $\text{Var}(z_1)$ de manera que $\|a_1\| = 1$:

$$L(a_1) = \overbrace{a_1'\Sigma a_1}^{\text{var}(z_1)} - \lambda(a_1'a_1 - 1)$$

$$\frac{\partial L(a_1)}{\partial a_1} = 0 \Rightarrow 2\Sigma a_1 - 2\lambda a_1 = 0$$

$$\Rightarrow (\Sigma - \lambda I)a_1 = 0 \Rightarrow \det(\Sigma - \lambda I) = 0 \text{ para } a_1 \neq 0$$

$$\Rightarrow \lambda \text{ es valor propio de } \Sigma \text{ asociado al vector propio } a_1$$

Recordar que Σ es diagonalizable en una base ortonormal pues es simétrica.

Cálculo de las componentes

Al ser la matriz de covarianzas Σ semidefinida positiva y de tamaño $p \times p$, consideramos $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ los valores propios de Σ .

$$\text{Var}(z_1) = \text{Var}(\mathbf{X}a_1) = a_1' \Sigma a_1 = a_1' \lambda_1 a_1 = \lambda_1 a_1' a_1 = \lambda_1$$

Para maximizar la varianza, tomo entonces el mayor valor propio λ_1 de Σ y el correspondiente vector propio $a_1' = (a_{11}, a_{12}, \dots, a_{1p})'$ (normalizado) y entonces

$$z_1 = \mathbf{X}a_1 = a_{11}x_1 + a_{12}x_2 + \dots + a_{1p}x_p$$

es la combinación lineal de los x_1, \dots, x_p con la mayor varianza.

2 Queremos ahora encontrar $z_2 = \mathbf{X}a_2$ tal que $\begin{cases} \text{Cov}(z_2, z_1) = 0 \\ \|a_2\| = 1 \end{cases}$

$$0 = \text{Cov}(z_2, z_1) = a_2' \Sigma a_1 = a_2' \lambda_1 a_1 \Leftrightarrow a_2' a_1 = 0$$

Maximizamos entonces la varianza de z_2 de manera que $\|a_2\| = 1$ y que $a_2' a_1 = 0$.

$$L(a_2) = \overbrace{a_2' \Sigma a_2}^{\text{Var}(z_2)} - \lambda(a_2' a_2 - 1) - \delta a_2' a_1$$

$$\frac{\partial L(a_2)}{\partial a_2} = 0 \Rightarrow 2\Sigma a_2 - 2\lambda a_2 - \delta a_1 = 0$$

Multiplicando por a_1' se tiene

$$2a_1' \Sigma a_2 - \delta = 0 \Rightarrow \delta = 2a_1' \Sigma a_2 = 2a_2' \Sigma a_1 = 0$$

Cálculo de las componentes

$$\frac{\partial L(a_2)}{\partial a_2} = 0 \Leftrightarrow 2\mathbf{\Sigma}a_2 - 2\lambda a_2 = 0 \Leftrightarrow (\mathbf{\Sigma} - \lambda I)a_2 = 0$$

Elijo entonces λ el 2do mayor valor propio de $\mathbf{\Sigma}$ con vector propio asociado a_2 .

Repetimos este procedimiento p veces, obteniendo los vectores a_1, a_2, \dots, a_p y se obtiene

una matriz ortogonal $A = \begin{pmatrix} | & | & \dots & | \\ a_1 & a_2 & \dots & a_p \\ | & | & \dots & | \end{pmatrix}$

Como $z_1 = \mathbf{X}a_1$, $z_p = \mathbf{X}a_p$, entonces:

Relación entre las viejas y las nuevas variables

- Observar que se puede escribir (poniendo las características en filas):

$$\begin{pmatrix} - & z_1 & - \\ - & z_2 & - \\ & \vdots & \\ - & z_p & - \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} & \dots & a_{p1} \\ a_{12} & a_{22} & \dots & a_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1p} & a_{2p} & \dots & a_{pp} \end{pmatrix} \begin{pmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_p \end{pmatrix}$$

$$Z' = A'X'$$

- O si no:

$$\begin{pmatrix} | & | & \dots & | \\ z_1 & z_2 & \dots & z_p \\ | & | & \dots & | \end{pmatrix} = \begin{pmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_p \\ | & | & \dots & | \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{pmatrix}$$

$$Z = XA$$

A las columnas de Z se le llaman *componentes principales* de X .

- Como $Var(z_1) = \lambda_1$, $Var(z_2) = \lambda_2$, ..., $Var(z_p) = \lambda_p$ y son incorreladas:

$$\Sigma_Z = Var(Z) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_p \end{pmatrix} \underbrace{=}_{Z=XA} A' Var(\mathbf{X}) A.$$

Entonces:

$$\Sigma = A \Sigma_Z A'$$

Porcentajes de variabilidad

$$\sum_{i=1}^p \text{Var}(z_i) = \sum_{i=1}^p \lambda_i = \text{tr}(\Sigma_Z) = \text{tr}(A' \Sigma_X A) = \text{tr}(\Sigma_X A A') = \text{tr}(\Sigma_X)$$

Porcentaje de variabilidad de la variable i :

$$\frac{\text{Var}(z_i)}{\sum_{i=1}^p \text{Var}(z_i)} = \frac{\lambda_i}{\sum_{i=1}^p \lambda_i} \quad \left(\text{con matriz correlaciones } \frac{\lambda_i}{p} \right)$$

Porcentaje de variabilidad de las m primeras variables i :

$$\sum_{j=1}^m \frac{\lambda_j}{\sum_{i=1}^p \lambda_i} \quad \text{donde } m < p$$

Nos quedamos con un número mucho menor de componentes que recogen un porcentaje amplio de la variabilidad total (fijada por el usuario). En general no se elige más de 3.

Interpretación geométrica

- Cada eje de \mathbb{R}^p representa una de las p variables.
- Supongamos que tenemos N individuos, y nos focalizamos en el individuo n , entonces las coordenadas de \mathbf{x}'_n son los datos de las p variables para este individuo.
- $\mathbf{z}'_n = \mathbf{x}'_n \mathbf{A} = \mathbf{P}(\mathbf{x}_n)$ son las coordenadas del individuo \mathbf{x}'_n en el nuevo sistema de referencia determinado por las componentes principales.
- Podemos entonces pensar que “proyectamos” la nube de la población dada por \mathbf{X} sobre un subespacio de dimensión la cantidad de componentes principales que retendremos.

Correlación entre las nuevas y las viejas variables

Como

$$\text{Cov}(z_j, x_i) = \text{Cov}\left(z_j, \sum_{k=1}^p a_{ik} z_k\right) = a_{ij} \text{Var}(z_j) = \lambda_j a_{ij}$$

entonces la correlación, si x_i está estandarizada es:

$$\text{Cor}(z_j, x_i) = \frac{\lambda_j a_{ij}}{\sqrt{\lambda_j}}$$

Consideraciones finales

- 1 Se calculan las componentes principales sobre variables originales estandarizadas (media 0 y varianza 1). Tomo entonces las componentes principales sobre la matriz de correlaciones y se le da la misma importancia a todas las variables.
- 2 Si las variables x_1, \dots, x_p ya son incorreladas, entonces no tiene sentido hacer componentes principales. Si se hace se obtiene las mismas variables ordenadas de mayor a menor varianza. Para ver eso se hace el test de esfericidad de Bartlett (package psych) o el indice de Kayser-Meyer-Olkin (KMO).
- 3 Si Σ tiene un valor propio con multiplicidad mayor que 1 se toma vectores propios ortogonales en el subespacio propio correspondiente.
- 4 Se conservan en general dos o tres componentes.

\mathbb{R}^p

$x_1 \leftrightarrow X_1$
 $x_2 \leftrightarrow X_2$
 $x_j \leftrightarrow X_j$
 $x_p \leftrightarrow X_p$

x_{i1} x_{i2} x_{ij} x_{ip}

x_i

$x_i' \leftrightarrow P(x_i)$

Proyector
 individuos.
 y las variables

$a_1 \leftrightarrow z_1$
 $a_2 \leftrightarrow z_2$
 $a_3 \leftrightarrow z_3$

z_{i1} z_{i2} z_{i3}

$\mathbb{R}^3 = \mathbb{R}^q$

$q = 3 \ll p$

Espacio Factorial.

$P(x_i) = \begin{pmatrix} x_i' \cdot a_1 \\ x_i' \cdot a_2 \\ x_i' \cdot a_3 \\ \vdots \\ x_i' \cdot a_p \end{pmatrix}$

$P(x_i)' = x_i' A = \begin{pmatrix} c_{i1} & c_{i2} & c_{i3} \\ z_{i1} & z_{i2} & z_{i3} \\ \vdots & \vdots & \vdots \\ z_{ip} \end{pmatrix}$

$= \begin{pmatrix} x_i \end{pmatrix} \begin{pmatrix} a_1 & a_2 & \dots & a_p \end{pmatrix} = \begin{pmatrix} z_{i1} \\ z_{i2} \\ \vdots \\ z_{ip} \end{pmatrix}$

$X = \begin{pmatrix} x_1 & x_2 & \dots & x_p \\ z_1 & z_2 & \dots & z_p \\ \vdots & \vdots & \vdots & \vdots \\ z_n & z_n & \dots & z_n \end{pmatrix} \quad n$

Referencias

- 1 Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani, An Introduction to Statistical Learning with Applications in R, Springer, 2013.
- 2 Daniel Peña, Análisis Multivariante, Mac Graw Hill, 2002.