

Usando las primeras dos columnas de g en la página 38, la primera y segunda componentes principales para los datos de billetes del banco suizo son

$$y_1 = 0.044x_1 + 0.112x_2 + 0.139x_3 + 0.768x_4 + 0.202x_5 - 0.579x_6$$

$$y_2 = 0.011x_1 + 0.071x_2 + 0.066x_3 - 0.563x_4 + 0.659x_5 - 0.489x_6$$

x_1 = longitud del billete.

x_2 = ancho del billete (a la izquierda).

x_3 = ancho del billete (a la derecha).

x_4 = distancia de la figura en el billete al borde inferior del billete

x_5 = distancia de la figura en el billete al borde superior del billete.

x_6 = longitud de la diagonal del billete

De forma que la primera componente y_1 , esencialmente corresponde a la diferencia entre el ancho del borde inferior del billete y la longitud de la diagonal

La segunda componente principal y_2 , corresponde a la diferencia entre el ancho del borde superior del billete y la suma de: el ancho del borde inferior y la longitud de la diagonal en los billetes.

Una medida que nos dice qué tan bien explican la variabilidad las primeras q componentes principales, está dada por

$$\psi_q = \frac{\sum_{j=1}^q \lambda_j}{\sum_{j=1}^p \lambda_j} = \frac{\sum_{j=1}^q \text{VAR}(Y_j)}{\sum_{j=1}^p \text{VAR}(Y_j)}$$

La figura D muestra una gráfica de las parejas ordenadas $(i, \frac{\lambda_i}{\sum_{j=1}^p \lambda_j})$; $i=1, 2, \dots, 6$.

Podemos darnos cuenta de que la 1ª componente principal explica un 66% de la variabilidad y las primeras dos componentes explican 88%. En algunos libros se sugiere seleccionar $i_0 \in \{1, 2, \dots, 6\}$ en donde la gráfica de

La figura ⁽¹⁾ D presenta un "codo" o se dobla. Este valor de i_0 es el número de componentes principales que se usarán para representar a los datos originales X .

El siguiente teorema resultará de interés

TEOREMA Sea X un vector aleatorio de dimensión $p \times 1$ tal que $E(X) = \mu$ y $VAR(X) = \Sigma$. Sea Y el vector de componentes principales de X , entonces

$$(1) \quad COV(X, Y) = \Pi \Lambda,$$

donde Π es la matriz de dimensiones $p \times p$ cuyas columnas son los vectores propios de Σ (en la descomposición de Jordan de Σ) y Λ es una matriz diagonal (de dimensiones $p \times p$) que tiene los valores propios de Σ , $\lambda_1, \dots, \lambda_p$ como elementos de la diagonal.

(1) Scree Plot

- (2) La correlación $\rho_{X_i Y_j}$ entre la variable X_i y la componente principal Y_j está dada por

$$\rho_{X_i Y_j} = r_{ij} \left(\frac{\lambda_j}{\sigma_{X_i}^2} \right)^{1/2}, \quad \dots \text{ (rho)}$$

donde $\sigma_{X_i}^2 = \text{VAR}(X_i) = \Sigma_{ii}$ la entrada i, i de la matriz Σ y r_{ij} es la entrada i, j de la matriz Γ .

Sea \mathcal{D}_i la columna i en Γ , notemos que $\sum_{j=1}^p \lambda_j r_{ij}^2 = \mathcal{D}_i' \Lambda \mathcal{D}_i = \text{elemento } (i, i) \text{ de la matriz } \Gamma \Lambda \Gamma' = \Sigma$, entonces tenemos

$$\begin{aligned} (\star) \left\{ \sum_{j=1}^p \rho_{X_i Y_j}^2 \right. &= \sum_{j=1}^p r_{ij}^2 \frac{\lambda_j}{\Sigma_{ii}} = \frac{\sum_{j=1}^p r_{ij}^2 \lambda_j}{\Sigma_{ii}} \\ &= \frac{\Sigma_{ii}}{\Sigma_{ii}} = 1. \end{aligned}$$

La anterior propiedad nos inspira a interpretar la "correlación" $\rho_{X_i Y_j}^2$ como "la proporción de

de la varianza de la variable X_i explicada por la j -ésima componente principal Y_j .

Si siguiendo estas ideas, el porcentaje de la varianza de X_i explicado por las primeras q componentes principales Y_1, \dots, Y_q está dado por $\sum_{j=1}^q \rho_{X_i Y_j}^2$.

Como todas las cantidades mencionadas son poblacionales, necesitamos definir como calcular la correlación en la ecuación (rho) para unos datos x_{i1}, \dots, x_{in} , esta correlación muestral está dada por

$$r_{X_i Y_j} = g_{ij} \left(\frac{l_j}{\hat{\Sigma}_{ii}} \right)^{1/2}.$$

Notemos que podemos argumentar como en (★) pero con cantidades muestrales:

$$\sum_{j=1}^p l_j g_{ij}^2 = g_i' \mathbf{I} g_i = \text{elemento } (i,i) \text{ de la matriz } \mathbf{g} \mathbf{I} \mathbf{g}' = \hat{\Sigma} \quad (g_i = \text{columna } i \text{ en } \mathbf{g}),$$