

y en tal caso $E(\tilde{X}) = 0$ porque

$$\begin{aligned} E(\tilde{X}) &= E(X - E(X)) = E(X) - E(E(X)) \\ &= E(X) - E(X) = 0 \end{aligned}$$

Vamos a escribir las relaciones (PC1), (PC2), (PC3), etc... que definen las componentes principales de un vector aleatorio X usando Σ y μ (la matriz de varianzas-covarianzas y el vector de medias) que son características de F , la distribución de X .

$$\mu = E(X)$$

de dim. $p \times p$ $\longrightarrow \Sigma = \text{VAR}(X) = \Gamma \Lambda \Gamma^T$
 simétrica y
 positivo def.

descomposición de
Jordan

$$\Lambda = \text{DIAG}(\lambda_1, \dots, \lambda_p)$$

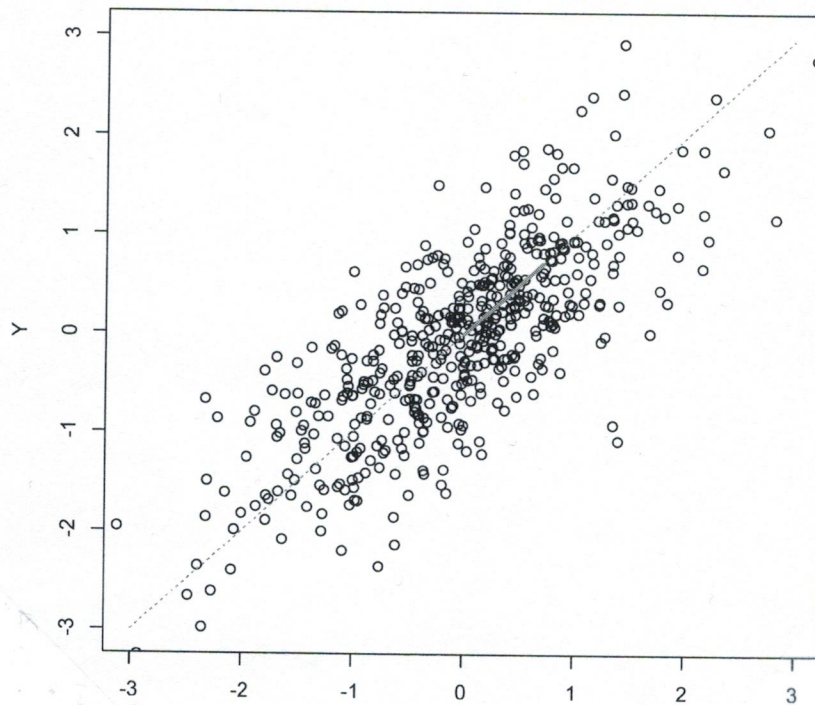
$$\Gamma = (\gamma_1^T, \gamma_2^T, \dots, \gamma_p^T)$$

Las relaciones (PC1), (PC2), (PC3), ...
 se pueden escribir en forma matricial como

$$Y = \Gamma^T (X - \mu) \dots \text{(CompPrin)}$$

Ejemplo: Supóngase $X \sim N_2(0, \Sigma); \Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$
con $\rho > 0$

<< FIGURA B >>
Normal sample, n = 500



Para encontrar los valores propios de Σ tenemos que resolver la ecuación

$$\begin{vmatrix} 1-\lambda & \rho \\ \rho & 1-\lambda \end{vmatrix} = 0 \quad (1-\lambda)^2 - \rho^2 = 0$$

cuyas soluciones son $\beta_1 = 1+\rho$ y $\beta_2 = 1-\rho$, de donde

$$\Lambda = \begin{pmatrix} 1+\rho & 0 \\ 0 & 1-\rho \end{pmatrix}.$$

El vector propio correspondiente a $\lambda_1 = 1+\rho$ se obtiene de

$$\begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = (1+\rho) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \text{o bien}$$

$$\left. \begin{aligned} x_1 + \rho x_2 &= (1+\rho) x_1 \\ \rho x_1 + x_2 &= (1+\rho) x_2 \end{aligned} \right\} \begin{aligned} &\text{ambas implican} \\ &x_1 = x_2 \end{aligned}$$

Todos los vectores $x = \begin{pmatrix} x_1 \\ x_1 \end{pmatrix}$ son propios para Σ , el vector x renormalizado asociado a λ_1 es

$$y_1 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}.$$

La figura muestra y_1 en color rojo (línea en rojo, gruesa) su dirección coincide con la dirección

en la que los datos tienen su mayor variabilidad (línea punteada en rojo). El segundo vector renormalizado (asociado a $\lambda_2 = 1 - \rho$) es

$$\mathbf{x}_2^u = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{pmatrix},$$

de forma que

$$\Gamma = (\mathbf{x}_1^u, \mathbf{x}_2^u) = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

que satisface

$$\Sigma = \Gamma \Lambda \Gamma'$$

Nota: Si se ajustara una regresión lineal a los datos (en donde esto tenga sentido), en general la dirección del vector propio asociado a λ_1 y la pendiente de la línea recta ajustada vía mínimos cuadrados serán diferentes. La razón es que la estimación vía mínimos cuadrados tiene por objetivo minimizar distancias (errores) ~~verticales~~, este objetivo es

diferente al que se tiene al maximizar una forma cuadrática seleccionando un vector propio de la matriz de varianzas-covarianzas de los datos.

La transformación de componentes principales en la ecuación (CompPrin), para el vector aleatorio \mathbf{X} queda como

$$\mathbf{Y} = \mathbf{T}'(\mathbf{X} - \boldsymbol{\mu}) = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \mathbf{X},$$

es decir

$$\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} X_1 + X_2 \\ X_1 - X_2 \end{pmatrix}.$$

La primera componente principal es

$$Y_1 = \frac{1}{\sqrt{2}} (X_1 + X_2).$$

La segunda componente principal es

$$Y_2 = \frac{1}{\sqrt{2}} (X_1 - X_2).$$