

MAT-269: Regresión Multivariada

Felipe Osorio

fosorios.mat.utfsm.cl

Departamento de Matemática, UTFSM



Modelo de regresión multivariado

En esta sección extendemos el modelo de regresión lineal considerando que ahora se dispone de k variables de respuesta. El **modelo lineal multivariado** asume la forma

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U},$$

donde \mathbf{Y} y \mathbf{U} son matrices aleatorias $n \times k$, \mathbf{X} es matriz de diseño $n \times p$ y \mathbf{B} es matriz de **coeficientes de regresión** $p \times k$. Asumiremos que $\text{rg}(\mathbf{X}) = p$ y $n \geq p + k$.

Además supondremos que las filas de la matriz de **disturbios aleatorios** son **independientes** $N_p(\mathbf{0}, \Sigma)$. Es decir,

$$\mathbf{U} \sim N_{n,k}(\mathbf{0}, \mathbf{I}_n \otimes \Sigma),$$

o análogamente

$$\mathbf{Y} \sim N_{n,k}(\mathbf{X}\mathbf{B}, \mathbf{I}_n \otimes \Sigma).$$



Resultado 1

Si $Y \sim N_{n,k}(XB, I_n \otimes \Sigma)$ y $n \geq k + p$ los **estimadores máximo verosímiles** de B y Σ son

$$\begin{aligned}\hat{B} &= (X^T X)^{-1} X^T Y, \\ \hat{\Sigma} &= \frac{1}{n} (Y - X\hat{B})^T (Y - X\hat{B})\end{aligned}$$

Además $(\hat{B}, \hat{\Sigma})$ es suficiente para (B, Σ) .



Demostración:

En efecto, como $\mathbf{Y} \sim N_{n,k}(\mathbf{XB}, \mathbf{I}_n \otimes \Sigma)$, la función de densidad conjunta de \mathbf{Y} es

$$f(\mathbf{Y}) = (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{Y} - \mathbf{XB}) \Sigma^{-1} (\mathbf{Y} - \mathbf{XB})^\top \right\},$$

ignorando términos que no involucran $\theta = (\mathbf{B}, \Sigma)$, tenemos

$$\ell_n(\mathbf{B}, \Sigma) = -\frac{n}{2} \log |\Sigma| - \frac{1}{2} \text{tr} \Sigma^{-1} (\mathbf{Y} - \mathbf{XB})^\top (\mathbf{Y} - \mathbf{XB}).$$

Diferenciando con relación a \mathbf{B} y Σ obtenemos

$$d_B \ell_n(\mathbf{B}, \Sigma) = -\frac{1}{2} \text{tr} \Sigma^{-1} d_B Q(\mathbf{B})$$

$$d_\Sigma \ell_n(\mathbf{B}, \Sigma) = -\frac{n}{2} \text{tr} \Sigma^{-1} d\Sigma + \frac{1}{2} \text{tr} \Sigma^{-1} (d\Sigma) \Sigma^{-1} Q(\mathbf{B})$$

donde $Q(\mathbf{B}) = (\mathbf{Y} - \mathbf{XB})^\top (\mathbf{Y} - \mathbf{XB})$.



En efecto,

$$d_B Q(B) = -(dB)^T X^T (Y - XB) - (Y - XB)^T X dB,$$

de este modo

$$d_B \ell_n(B, \Sigma) = \frac{1}{2} \operatorname{tr} \Sigma^{-1} \{ (dB)^T X^T (Y - XB) + (Y - XB)^T X dB \},$$

recordando que $\operatorname{tr} A = \operatorname{tr} A^T$ y como Σ es simétrica

$$d_B \ell_n(B, \Sigma) = \operatorname{tr} \Sigma^{-1} (Y - XB)^T X dB,$$

el diferencial es cero si y sólo si

$$X^T (Y - XB) = 0,$$

es decir \hat{B} es solución del sistema de ecuaciones

$$X^T X \hat{B} = X^T Y,$$

como $\operatorname{rg}(X) = p$, tenemos $\hat{B} = (X^T X)^{-1} X^T Y$.



Modelo de regresión multivariado

Por otro lado,

$$d_{\Sigma} Q(\mathbf{B}) = -\frac{1}{2} \operatorname{tr} \Sigma^{-1} (n \Sigma - Q(\mathbf{B})) \Sigma^{-1} d \Sigma,$$

y por tanto,

$$\hat{\Sigma} = \frac{1}{n} Q(\hat{\mathbf{B}}).$$

Note además que

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\mathbf{B})^{\top} (\mathbf{Y} - \mathbf{X}\mathbf{B}) &= (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} + \mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B})^{\top} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} + \mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B}) \\ &= Q(\hat{\mathbf{B}}) + (\hat{\mathbf{B}} - \mathbf{B})^{\top} \mathbf{X}^{\top} \mathbf{X} (\hat{\mathbf{B}} - \mathbf{B}) \end{aligned}$$

(pues $\mathbf{X}^{\top} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) = 0$), de ahí que

$$\begin{aligned} L(\mathbf{B}, \Sigma) &= |\Sigma|^{-n/2} \exp \left\{ -\frac{1}{2} \operatorname{tr} \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})^{\top} (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right\} \\ &= |\Sigma|^{-n/2} \exp \left\{ -\frac{1}{2} \operatorname{tr} \Sigma^{-1} (Q(\hat{\mathbf{B}}) + (\hat{\mathbf{B}} - \mathbf{B})^{\top} \mathbf{X}^{\top} \mathbf{X} (\hat{\mathbf{B}} - \mathbf{B})) \right\} \\ &= |\Sigma|^{-n/2} \exp \left\{ -\frac{n}{2} \operatorname{tr} \hat{\Sigma} \Sigma^{-1} + \frac{1}{2} \operatorname{tr} \Sigma^{-1} (\hat{\mathbf{B}} - \mathbf{B})^{\top} \mathbf{X}^{\top} \mathbf{X} (\hat{\mathbf{B}} - \mathbf{B}) \right\}, \end{aligned}$$

es decir $(\hat{\mathbf{B}}, \hat{\Sigma})$ es suficiente para (\mathbf{B}, Σ) .



Resultado 2

Si $\mathbf{Y} \sim N_{n,k}(\mathbf{X}\mathbf{B}, \mathbf{I}_n \otimes \mathbf{\Sigma})$ los estimadores máximo verosímiles $\hat{\mathbf{B}}$ y $\hat{\mathbf{\Sigma}}$ son independientemente distribuidos como:

$$\begin{aligned}\hat{\mathbf{B}} &\sim N_{q,k}(\mathbf{B}, (\mathbf{X}^\top \mathbf{X})^{-1} \otimes \mathbf{\Sigma}), \\ n\hat{\mathbf{\Sigma}} &\sim W_k(n-p, \mathbf{\Sigma}).\end{aligned}$$

Demostración:

Sea $\mathbf{E} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$, con

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top.$$

Sea $\mathbf{M} = \mathbf{I} - \mathbf{H}$, entonces

$$\mathbf{M}\mathbf{X} = (\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}, \quad \mathbf{M}^2 = \mathbf{M}.$$

De este modo,

$$\mathbf{E}^\top \mathbf{E} = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{Y}^\top \mathbf{M}\mathbf{Y} = n\hat{\mathbf{\Sigma}}.$$



Modelo de regresión multivariado

Considere la matriz

$$\begin{pmatrix} \hat{B} \\ E \end{pmatrix} = \begin{pmatrix} (X^T X)^{-1} X^T \\ M \end{pmatrix} Y$$

luego $(\hat{B}^T, E^T)^T$ es normal con media

$$E \begin{pmatrix} \hat{B} \\ E \end{pmatrix} = \begin{pmatrix} (X^T X)^{-1} X^T \\ M \end{pmatrix} E(Y) = \begin{pmatrix} \hat{B} \\ E \end{pmatrix} = \begin{pmatrix} (X^T X)^{-1} X^T \\ M \end{pmatrix} X B = \begin{pmatrix} B \\ 0 \end{pmatrix}$$

mientras que la covarianza es:

$$\begin{aligned} & \left(\begin{pmatrix} (X^T X)^{-1} X^T \\ M \end{pmatrix} \otimes I_p \right) (I_n \otimes \Sigma) \left((X(X^T X)^{-1}, M^T) \otimes I_p \right) \\ &= \begin{pmatrix} (X^T X)^{-1} X^T \\ M \end{pmatrix} (X(X^T X)^{-1}, M^T) \otimes \Sigma \\ &= \begin{pmatrix} (X^T X)^{-1} & (X^T X)^{-1} X^T M^T \\ M X (X^T X)^{-1} & M M^T \end{pmatrix} \otimes \Sigma \end{aligned}$$



En efecto,

$$\text{Cov} \left(\text{vec} \begin{pmatrix} \hat{\mathbf{B}} \\ \mathbf{E} \end{pmatrix} \right) = \begin{pmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{M} \end{pmatrix} \otimes \boldsymbol{\Sigma}.$$

De ahí que

$$\hat{\mathbf{B}} \sim \text{N}(\mathbf{B}, (\mathbf{X}^\top \mathbf{X})^{-1} \otimes \boldsymbol{\Sigma}), \quad \mathbf{E} \sim \text{N}(\mathbf{0}, \mathbf{M} \otimes \boldsymbol{\Sigma}),$$

luego $\hat{\mathbf{B}}$ y \mathbf{E} son independientes. La parte final del resultado sigue de notar que $n\hat{\boldsymbol{\Sigma}}$ puede ser escrito como $\mathbf{E}^\top \mathbf{E}$.



Ejemplo: Datos de dializadores (Vonesh y Carter, 1987)¹

Datos de un estudio para evaluar las características de **ultrafiltración** in vivo de un **grupo de dializadores**.

Los dializadores se evaluaron en **tres centros** y cada uno de ellos utilizó un tipo **diferente** de sistema de administración de dializado.

Los datos corresponden a la **tasa de ultrafiltración de los 4 dializadores** Y_1, Y_2, Y_3, Y_4 .

Este conjunto de datos ha sido usado en varios artículos científicos y bajo el **supuesto de normalidad**. Además, este supuesto no es rechazado por el test de Mardia (1974), basado en las medidas de sesgo y kurtosis multivariadas.

¹Biometrics **43**, 617-628.



Ejemplo: Datos de dializadores

```
# Carga biblioteca 'heavy' y datos de ejemplo  
> library(heavy)  
> data(dialyzer)
```

```
> dialyzer  
      y1    y2    y3    y4 centre  
1    600  1026  1470  1890      1  
2    516   930  1380  1770      1  
3    480   900  1380  1860      1  
4    528   930  1410  1872      1  
5    540   978  1410  1920      1  
6    564   996  1422  1920      1  
7    564  1062  1500  1980      1  
8    492   900  1392  1860      1  
9    516   960  1380  1800      1  
10   528   930  1356  1860      1  
11   564  1020  1380  1884      1
```

...

```
38  480   780  1140  1710      3  
39  540   840  1200  1650      3  
40  780   780  1290  1680      3
```



Ejemplo: Datos de dializadores

Vamos a considerar un modelo de regresión multivariada con matriz de respuestas

$$\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3, \mathbf{Y}_4),$$

y matriz de diseño

$$\mathbf{X} = \begin{pmatrix} \mathbf{1}_{17} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{12} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1}_{11} \end{pmatrix}.$$

Por tanto tenemos que $\mathbf{B} \in \mathbb{R}^{3 \times 4}$ y $\mathbf{\Sigma}$ es matriz simétrica y definida positiva 4×4 .



Ejemplo: Datos de dializadores

```
# Ajuste de un modelo de regresión lineal multivariado
# (bajo errores normales)
> fm <- heavyLm(cbind(y1,y2,y3,y4) ~ -1 + centre, data = dialyzer,
+             family = normal())

# Salida con estimación de los coeficientes de regresión
> fm
Call:
heavyLm(formula = cbind(y1, y2, y3, y4) ~ -1 + centre, data = dialyzer,
        family = normal())
Converged in 1 iterations

Coefficients:
           y1           y2           y3           y4
centre1  541.0588  973.4118 1404.3529 1873.4118
centre2  472.5000  830.5000 1230.5000 1653.0000
centre3  591.8182  850.9091 1276.3636 1655.4545

Degrees of freedom: 40 total; 37 residual
```



Ejemplo: Datos de dializadores

```
# Salida de función 'summary'
> summary(fm)
Multivariate regression under heavy-tailed distributions
Data: dialyzer; Family: normal()

Coefficients:
      y1      y2      y3      y4
centre1 541.0588 973.4118 1404.3529 1873.4118
centre2 472.5000 830.5000 1230.5000 1653.0000
centre3 591.8182 850.9091 1276.3636 1655.4545

Scatter matrix estimate:
      y1      y2      y3      y4
y1 4402.5848
y2 380.6878 2337.0659
y3 867.8210 2712.4908 6160.8440
y4 506.2871 1987.2839 4808.4436 7134.5452

Degrees of freedom: 40 total; 37 residual
Log-likelihood: -761.7318 on 22 degrees of freedom
```



Ejemplo: Datos de dializadores

Es decir, tenemos que

$$\hat{\mathbf{B}} = \begin{pmatrix} 541.059 & 973.412 & 1404.353 & 1873.412 \\ 472.500 & 830.500 & 1230.500 & 1653.000 \\ 591.818 & 850.909 & 1276.364 & 1655.455 \end{pmatrix}.$$

y

$$\hat{\mathbf{\Sigma}} = \begin{pmatrix} 4402.585 & 380.688 & 867.821 & 506.287 \\ 380.688 & 2337.066 & 2712.491 & 1987.284 \\ 867.821 & 2712.491 & 6160.844 & 4808.444 \\ 506.287 & 1987.284 & 4808.444 & 7134.545 \end{pmatrix}.$$

Además, $\ell_n(\hat{\mathbf{B}}, \hat{\mathbf{\Sigma}}) = -761.732$.



Ejemplo: Datos de dializadores

```
# Alternativamente, podemos usar la función 'lm' de R
> f0 <- lm(cbind(y1,y2,y3,y4) ~ -1 + centre, data = dialyzer)

# Salida:
> f0

Call:
lm(formula = cbind(y1, y2, y3, y4) ~ -1 + centre, data = dialyzer)

Coefficients:
          y1          y2          y3          y4
centre1  541.1   973.4  1404.4  1873.4
centre2  472.5   830.5  1230.5  1653.0
centre3  591.8   850.9  1276.4  1655.5
```



El objetivo de esta sección es estimar B sujeto a restricciones del tipo:

$$AB = C,$$

donde A es matriz $r \times p$ de rango r y C es matriz $t \times k$. Sabemos que A puede ser particionada como:

$$A = (A_r, A_s),$$

donde A_r es no singular. De este modo,

$$AB = (A_r, A_s) \begin{pmatrix} B_r \\ B_s \end{pmatrix} = A_r B_r + A_s B_s = C,$$

es decir,

$$B_r = A_r^{-1}(C - A_s B_s).$$



Substituyendo en el modelo, tenemos

$$\begin{aligned} Y &= XB + U = (X_r, X_s) \begin{pmatrix} B_r \\ B_s \end{pmatrix} + U, \\ &= X_r B_r + X_s B_s + U, \\ &= X_r A_r^{-1} (C - A_s B_s) + X_s B_s + U, \\ &= X_r A_r^{-1} C + (X_s - X_r A_r^{-1} A_s) B_s + U, \end{aligned}$$

que puede ser escrito como:

$$Y_R = X_R B_s + U, \quad U \sim N(0, I_n \otimes \Sigma),$$

con

$$Y_R = Y - X_r A_r^{-1} C, \quad X_R = X_s - X_r A_r^{-1} A_s.$$



De este modo,

$$\begin{aligned}\tilde{\mathbf{B}}_s &= (\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \mathbf{X}_R^\top \mathbf{Y}, \\ \tilde{\mathbf{B}}_r &= \mathbf{A}_r^{-1} (\mathbf{C} - \mathbf{A}_s \tilde{\mathbf{B}}_s)\end{aligned}$$

Además, como $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ sigue que

$$\mathbf{Y}_R \sim \mathcal{N}(\mathbf{X}_R \mathbf{B}_s, \mathbf{I}_n \otimes \boldsymbol{\Sigma}),$$

y por tanto,

$$\tilde{\mathbf{B}}_s \sim \mathcal{N}(\mathbf{B}_s, (\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \otimes \boldsymbol{\Sigma}).$$

Como

$$\begin{aligned}\tilde{\mathbf{B}} &= \begin{pmatrix} \tilde{\mathbf{B}}_r \\ \tilde{\mathbf{B}}_s \end{pmatrix} = \begin{pmatrix} \mathbf{A}_r^{-1} (\mathbf{C} - \mathbf{A}_s \tilde{\mathbf{B}}_s) \\ \tilde{\mathbf{B}}_s \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A}_r^{-1} \mathbf{C} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \tilde{\mathbf{B}}_s.\end{aligned}$$

Así, $\tilde{\mathbf{B}}$ sigue una distribución normal con

$$\mathbb{E}(\tilde{\mathbf{B}}) = \begin{pmatrix} \mathbf{A}_r^{-1} \mathbf{C} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \mathbb{E}(\tilde{\mathbf{B}}_s) = \begin{pmatrix} \mathbf{B}_r \\ \mathbf{B}_s \end{pmatrix}$$

y

$$\text{vec } \tilde{\mathbf{B}} = \text{vec} \begin{pmatrix} \mathbf{A}_r^{-1} \mathbf{C} \\ \mathbf{0} \end{pmatrix} + \left(\mathbf{I} \otimes \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \right) \text{vec } \tilde{\mathbf{B}}_r,$$

de donde sigue que

$$\begin{aligned} \text{Cov}(\text{vec } \tilde{\mathbf{B}}) &= \left(\mathbf{I} \otimes \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \right) \text{Cov}(\text{vec } \tilde{\mathbf{B}}_s) \left(\mathbf{I} \otimes \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \right)^{\top} \\ &= (\mathbf{X}_R^{\top} \mathbf{X}_R)^{-1} \otimes \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix} \Sigma \begin{pmatrix} -\mathbf{A}_r^{-1} \mathbf{A}_s \\ \mathbf{I} \end{pmatrix}^{\top} \end{aligned}$$



Por otro lado,

$$\begin{aligned} Y_R - X_R \tilde{B}_s &= Y - X_r A_r^{-1} C - (X_s - X_r A_r^{-1} A_s) \tilde{B}_s \\ &= Y - X_r A_r^{-1} (C - A_r^1 A_s \tilde{B}_s) - X_s \tilde{B}_s \\ &= Y - X \tilde{B}, \end{aligned}$$

lo que lleva a

$$\begin{aligned} \tilde{\Sigma} &= \frac{1}{n} (Y_R - X_R \tilde{B}_s)^\top (Y_R - X_R \tilde{B}_s) \\ &= \frac{1}{n} (Y - X \tilde{B})^\top (Y - X \tilde{B}) \\ &= \frac{1}{n} Q(\tilde{B}). \end{aligned}$$



Considere

$$Q(\hat{B}) = (Y - X\hat{B})^\top (Y - X\hat{B}) = R$$

$$Q(\tilde{B}) = (Y - X\tilde{B})^\top (Y - X\tilde{B}) = S$$

Cuando $H_0 : AB = C$ es verdadera, R y $H = S - R$ son independientemente distribuidos $W_k(n - p, \Sigma)$ y $W_k(r, \Sigma)$, respectivamente.

Además, podemos escribir

$$H = (A\hat{B} - C)^\top [A(X^\top X)^{-1}A^\top]^{-1}(A\hat{B} - C),$$

como $n - k \geq p$, ambos R y S son definidas positivas con probabilidad 1.



Sea $L(\mathbf{B}, \mathbf{\Sigma})$ la función de verosimilitud para las filas de \mathbf{Y} , el test de razón de verosimilitudes para $H_0 : \mathbf{AB} = \mathbf{C}$, es

$$\Lambda = \frac{L(\tilde{\mathbf{B}}, \tilde{\mathbf{\Sigma}})}{L(\hat{\mathbf{B}}, \hat{\mathbf{\Sigma}})} = \frac{|\tilde{\mathbf{\Sigma}}|^{-n/2}}{|\hat{\mathbf{\Sigma}}|^{-n/2}},$$

de este modo

$$T = \Lambda^{2/n} = \frac{|\hat{\mathbf{\Sigma}}|}{|\tilde{\mathbf{\Sigma}}|} = \frac{|\mathbf{R}|}{|\mathbf{S}|} = \frac{|\mathbf{R}|}{|\mathbf{R} + \mathbf{H}|} = |\mathbf{I} - \mathbf{V}|,$$

donde $\mathbf{V} = \mathbf{S}^{-1/2} \mathbf{H} \mathbf{S}^{-1/2}$. Cuando H_0 es verdadera $T \sim \Lambda(k, r, n - p)$ y por el principio de razón de verosimilitudes, rechazamos $H_0 : \mathbf{AB} = \mathbf{C}$ si T es muy pequeño, es decir, si $|\mathbf{S}|$ es mucho mayor que $|\mathbf{R}|$.

