

### 3.6.1 Modelo lineal general multivariado

La distinción entre los modelos lineales multivariados y los modelos univariados es, como su nombre lo señala, que el modelo multivariado involucra más de una variable dependiente o respuesta.

Considérese que las observaciones multivariadas  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ , conforman un conjunto de observaciones independientes de una población normal  $p$ -variante; es decir,  $\mathbf{Y}_\alpha \sim N_p(\mathbf{X}_\alpha\boldsymbol{\beta}, \boldsymbol{\Sigma})$ , para  $\alpha = 1, \dots, n$ . Los vectores  $\mathbf{X}_\alpha$  de tamaño  $(1 \times q)$  son conocidos. Tanto la matriz  $\boldsymbol{\Sigma}_{p \times p}$ , como la matriz  $\boldsymbol{\beta}_{q \times p}$  son desconocidas. Los  $\mathbf{Y}_\alpha$  corresponden a las variables respuesta en un modelo de regresión (dependientes), mientras que las  $\mathbf{X}_\alpha$  son las variables regresoras o explicativas. En tales condiciones los vectores se pueden relacionar a través de un *modelo lineal general multivariado*, tal como el siguiente:

$$\begin{pmatrix} y_{11} \cdots y_{1p} \\ y_{21} \cdots y_{2p} \\ \vdots \cdots \vdots \\ y_{n1} \cdots y_{np} \end{pmatrix} = \begin{pmatrix} x_{11} \cdots x_{1q} \\ x_{21} \cdots x_{2q} \\ \vdots \cdots \vdots \\ x_{n1} \cdots x_{nq} \end{pmatrix} \begin{pmatrix} \beta_{11} \cdots \beta_{1p} \\ \beta_{21} \cdots \beta_{2p} \\ \vdots \cdots \vdots \\ \beta_{q1} \cdots \beta_{qp} \end{pmatrix} + \begin{pmatrix} \varepsilon_{11} \cdots \varepsilon_{1p} \\ \varepsilon_{21} \cdots \varepsilon_{2p} \\ \vdots \cdots \vdots \\ \varepsilon_{n1} \cdots \varepsilon_{np} \end{pmatrix}$$

En forma condensada, el modelo lineal multivariado anterior se escribe de la manera siguiente:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$$

La matriz  $\mathbf{X}$  conforma, en la mayoría de los casos, la matriz de diseño o la matriz de variables regresoras,  $\boldsymbol{\beta}$  es la matriz de parámetros desconocidos y la matriz aleatoria  $\mathbf{E}$  contiene los errores.

Para los propósitos de este texto, se propone, estima e infiere sobre los modelos ligados una estructura de una y dos vías de clasificación, mediante una conformación adecuada de la matriz de diseño  $\mathbf{X}$  y de la matriz de parámetros  $\boldsymbol{\beta}$ . Además, se extiende el análisis de perfiles, de medidas repetidas y de curvas de crecimiento, para el caso de varias poblaciones multivariadas.

Tal como en el modelo lineal clásico ( $q = 1$ ), los estimadores de máxima verosimilitud para  $\boldsymbol{\beta}$  y  $\boldsymbol{\Sigma}$  son:

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left( \sum_{\alpha=1}^n \mathbf{X}'_\alpha \mathbf{X}_\alpha \right)^{-1} \left( \sum_{\alpha=1}^n \mathbf{X}'_\alpha \mathbf{Y}_\alpha \right) \\ \hat{\boldsymbol{\Sigma}} &= \frac{1}{n} \sum_{\alpha=1}^n (\mathbf{Y}_\alpha - \mathbf{X}_\alpha \hat{\boldsymbol{\beta}})(\mathbf{Y}_\alpha - \mathbf{X}_\alpha \hat{\boldsymbol{\beta}})' \end{aligned} \quad (3.36)$$

### 3.6.3 Análisis de varianza multivariado

Desde un punto de vista práctico, el análisis de varianza multivariado es una técnica con la cual se puede verificar la igualdad de los vectores de medias ligados a varias poblaciones multivariadas.

Muchas hipótesis en el campo multivariado pueden expresarse como las hipótesis concernientes al análisis de regresión esquematizado anteriormente. Dentro de este estilo, se presenta la técnica del análisis de varianza para arreglos de una y dos vías de clasificación.

### 3.6.4 Modelos de una vía de clasificación

Considérese que  $Y_{ij}$  es una observación de una población  $N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$  con  $i = 1, \dots, q$ , y  $j = 1, \dots, n_i$ . Los datos se pueden visualizar de la siguiente forma

Población	Muestra	Media muestral
1	$Y_{11}, Y_{12}, \dots, Y_{1n_1}$	$\bar{Y}_{1\bullet}$
2	$Y_{21}, Y_{22}, \dots, Y_{2n_2}$	$\bar{Y}_{2\bullet}$
$\vdots$	$\vdots$	$\vdots$
$q$	$Y_{q1}, Y_{q2}, \dots, Y_{qn_q}$	$\bar{Y}_{q\bullet}$

La media  $\bar{Y}_{i\bullet}$  en cada muestra se obtiene mediante

$$\bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{n_i} Y_{i\bullet}, \text{ para } i = 1, \dots, q.$$

La media general  $\bar{Y}_{\bullet\bullet}$  se obtiene de

$$\bar{Y}_{\bullet\bullet} = \frac{1}{N} \sum_{i=1}^q \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{N} \sum_{i=1}^q \bar{Y}_{i\bullet}$$

con  $N = \sum_{i=1}^q n_i$ , el número total de observaciones.

El modelo que relaciona las observaciones con los parámetros  $\boldsymbol{\mu}_i$  es de la forma

$$Y_{ij} = \boldsymbol{\mu}_i + \mathbf{E}_{ij}, \text{ con } \mathbf{E}_{ij} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma}), \text{ para } i = 1, \dots, q \text{ y } j = 1, \dots, n_i.$$

$$\begin{pmatrix} Y'_{11} \\ Y'_{12} \\ \vdots \\ Y'_{1n_1} \\ \vdots \\ Y'_{q1} \\ Y'_{q2} \\ \vdots \\ Y'_{qn_q} \end{pmatrix} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_q} \end{pmatrix} \begin{pmatrix} \mu'_1 \\ \mu'_2 \\ \vdots \\ \mu'_q \end{pmatrix} + \begin{pmatrix} \varepsilon'_{11} \\ \varepsilon'_{12} \\ \vdots \\ \varepsilon'_{1n_1} \\ \vdots \\ \varepsilon'_{q1} \\ \varepsilon'_{q2} \\ \vdots \\ \varepsilon'_{qn_q} \end{pmatrix}$$

$$\mathbf{Y} = \bigoplus_{i=1}^q \mathbf{1}_{n_i} \mu_i + \mathbf{E}.$$

La hipótesis a verificar es la igualdad de los vectores de medias de las  $q$ -poblaciones; es decir,

$$H_0 : \mu_1 = \cdots = \mu_q. \quad (3.41)$$

La hipótesis planteada en (3.41) se puede escribir en la forma

$$\begin{aligned} H_0 : \mu_1 - \mu_q = \mu_2 - \mu_q = \cdots = \mu_{q-1} - \mu_q = \mathbf{0} \\ : \begin{pmatrix} 1 & 0 & \cdots & 0 & -1 \\ 0 & 1 & \cdots & 0 & -1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -1 \end{pmatrix} \begin{pmatrix} \mu'_1 \\ \mu'_2 \\ \vdots \\ \mu'_q \end{pmatrix} = \mathbf{0}. \end{aligned}$$

La ecuación (3.39) se utiliza para contrastar esta hipótesis. La región de rechazo a un nivel de significación  $\alpha$  es

$$\Lambda = \frac{|\mathbf{E}|}{|\mathbf{E} + \mathbf{H}|} = \frac{|N\widehat{\Sigma}|}{|N\widehat{\Sigma}_\omega|} < \Lambda_{(\alpha, p, \nu_H, \nu_E)} \quad (3.42)$$

donde  $\nu_H = q - 1$  son los grados de libertad para la hipótesis,  $\nu_E = N - q$  son los grados de libertad del error ( $N = \sum_{i=1}^q n_i$ ).

Más explícitamente

$$\begin{aligned}
 \mathbf{H} &= \sum_{i=1}^q n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^T \\
 \mathbf{E} &= \sum_{i=1}^q \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet}) (Y_{ij} - \bar{Y}_{i\bullet})^T \\
 \mathbf{E} + \mathbf{H} &= \sum_{i=1}^q \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet}) (Y_{ij} - \bar{Y}_{\bullet\bullet})^T. \tag{3.45}
 \end{aligned}$$

Esta escritura de  $\mathbf{E}$  permite encontrar un estimador insesgado de  $\Sigma$ . De esta manera:

$$\mathbf{E} = (n_1 - 1)\mathbf{S}_1 + \cdots + (n_q - 1)\mathbf{S}_q = \sum_{i=1}^q (n_i - 1)\mathbf{S}_i, \tag{3.46}$$

donde  $\mathbf{S}_i$  es la matriz de covarianzas de la  $i$ -ésima muestra. Así, la matriz de varianzas y covarianzas estimada, puesto que las poblaciones se han considerado con igual matriz de covarianzas, es:

$$\mathbf{S}_p = \frac{1}{\sum_{i=1}^q (n_i - 1)} \mathbf{E} = \frac{\sum_{i=1}^q (n_i - 1)\mathbf{S}_i}{\sum_{i=1}^q (n_i - 1)}.$$

Es inmediato que para  $p = 1$  (caso univariado), la razón de máxima verosimilitud se reduce a la conocida estadística  $F$ ; así, se rechaza  $H_0$  si:

$$\frac{\sum_i n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2}{\sum_{i=1}^q \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2} \left( \frac{N - q}{q - 1} \right) > F_{(\alpha, q-1, N-q)}$$

La distribución exacta de  $\Lambda$  ha sido obtenida para algunos casos especiales, la tabla 3.8 los resume.

Para muestras de tamaño grande se tiene la estadística de Bartlett

$$\begin{aligned}
 V &= - \left( N - 1 - \frac{(p + q)}{2} \right) \ln \Lambda \\
 &= - \left( N - 1 - \frac{(p + q)}{2} \right) \ln \left( \frac{|\mathbf{E}|}{|\mathbf{E} + \mathbf{H}|} \right),
 \end{aligned}$$

la cual tiene aproximadamente una distribución Ji-cuadrado con  $p(q-1)$  grados de libertad. Se rechaza  $H_0$  para valores de  $V$  mayores que  $\chi^2_{(\alpha, p(q-1))}$ .

Tabla 3.8: Relación entre las estadísticas  $\Lambda$  y  $F$

No. Variables	No. Grupos	Transformación	Distribución $F$
$p = 1$	$q \geq 2$	$\left(\frac{1-\Lambda}{\Lambda}\right) \left(\frac{N-q}{q-1}\right)$	$F_{(q-1, N-q)}$
$p = 2$	$q \geq 2$	$\left(\frac{1-\Lambda^{1/2}}{\Lambda^{1/2}}\right) \left(\frac{N-q-1}{q-1}\right)$	$F_{(2(q-1), 2(N-q-1))}$
$p \geq 1$	$q = 2$	$\left(\frac{1-\Lambda}{\Lambda}\right) \left(\frac{N-p-1}{p}\right)$	$F_{(p, N-p-1)}$
$p \geq 1$	$q = 3$	$\left(\frac{1-\Lambda^{1/2}}{\Lambda^{1/2}}\right) \left(\frac{N-p-2}{p}\right)$	$F_{(2p, 2(N-p-2))}$

**Ejemplo 3.6.1.** Con los siguientes datos se quiere establecer si tres métodos de enseñanza producen el mismo rendimiento promedio en matemáticas y escritura en niños de características similares.

Es éste un problema de análisis de varianza multivariado, con  $p = 2$  que corresponde a los puntajes en matemáticas y escritura por estudiante. El número



Tabla 3.9: Datos de rendimiento bajo tres métodos de enseñanza

Método 1	(69) (75)	(69) (70)	(71) (73)	(78) (82)	(79) (81)	(73) (75)
Método 2	(69) (70)	(68) (74)	(75) (80)	(78) (85)	(68) (68)	(63) (74)
	(63) (66)	(71) (76)	(72) (78)	(71) (73)	(70) (73)	(56) (83)
Método 3	(72) (79)	(64) (65)	(74) (74)	(74) (74)	(72) (75)	(82) (84)
	(76) (76)	(68) (65)	(78) (79)	(70) (71)	(60) (61)	(69) (68)

Fuente: Freund, Litell & Spector (1986)

de poblaciones es  $q = 3$ ; es decir, las tres metodologías. Los resultados del experimento se muestran en la tabla 3.9.

Se hará el análisis de varianza univariado (ANDEVA); es decir, para cada una de las dos variables, y el análisis de varianza multivariado (ANAVAMU) que se sugiere en este capítulo.

Las tablas 3.10 y 3.11 corresponden al análisis de varianza para cada una de las variables en forma separada.

Tabla 3.10: ANDEVA para matemáticas

F. de Variación	G. L.	S. C.	C. M.	Valor $F$	$Pr > F$
Métodos	2	60.6051	30.3025	0.91	0.4143
Error	28	932.8788	33.3171		
Total	30	993.4839			

De los resultados mostrados en la tabla 3.10 se puede afirmar que las metodologías no producen rendimientos promedios diferentes en matemáticas, en esta clase de niños.

Una conclusión similar se puede extraer de la tabla 3.11 para la variable escritura.

Tabla 3.11: ANDEVA para escritura

F. de Variación	G. L.	S. C.	C. M.	Valor $F$	$Pr > F$
Métodos	2	49.7359	24.8679	0.56	0.5776
Error	28	1243.9416	44.4265		
Total	30	1293.6775			

$$Y_{ij} = \mu + \mu_i + \varepsilon_{ij} \quad \text{con } i = 1, 2 \text{ y } 3 \quad j = 1, \dots, n_i.$$

En este caso  $n_1 = 6$ ,  $n_2 = 14$ ,  $n_3 = 11$  y  $N = 31$ .

Mediante la hipótesis nula se afirma que los métodos producen un rendimiento en promedio igual en matemáticas y en escritura; es decir,

$$H_0 : \mu_1 = \mu_2 = \mu_3.$$

Las matrices de sumas de cuadrados (covariabilidad) dentro y entre tratamientos se obtienen aplicando (3.45)

$$\mathbf{E} = \begin{pmatrix} 932.87879 & 1018.6818 \\ 1018.6818 & 1243.9416 \end{pmatrix} \quad \mathbf{H} = \begin{pmatrix} 60.6050 & 31.5117 \\ 31.5117 & 49.7358 \end{pmatrix}.$$

El valor del lambda de Wilks es

$$\Lambda = \frac{|\mathbf{E}|}{|\mathbf{E} + \mathbf{H}|} = \frac{\begin{vmatrix} 932.8788 & 1018.6818 \\ 1018.6818 & 1243.9416 \end{vmatrix}}{\begin{vmatrix} 993.4838 & 1050.1935 \\ 1050.1935 & 1293.6774 \end{vmatrix}} = 0.6731.$$

De la tabla 3.8 y como  $p = 2$  y  $q = 3$ , se puede utilizar la estadística

$$\begin{aligned} \left( \frac{1 - \Lambda^{1/2}}{\Lambda^{1/2}} \right) \left( \frac{n - q - 1}{q - 1} \right) &= \left( \frac{1 - \sqrt{0.6731}}{\sqrt{0.6731}} \right) \left( \frac{31 - 3 - 1}{3 - 1} \right) \\ &= 2.954851. \end{aligned}$$

El valor anterior comparado con  $F_{(5\%, 2(3-1), 2(31-3-1))} = F_{(5\%, 4, 54)} \approx 2.5$  (tabla C.8), permite afirmar que el puntaje promedio no es el mismo para las tres metodologías. ¡El resultado no es el mismo que se obtuvo con los análisis de varianzas univariados!