

CAPÍTULO IV

IV. ANÁLISIS ESTADÍSTICO MULTIVARIADO

4.1 INTRODUCCIÓN

Una vez concluido el análisis univariado, es necesario realizar un estudio más elaborado, en la que se pueda combinar varias variables en estudio, de tal manera que permita examinar la interacción y la influencia existente entre las variables.

Las técnicas multivariada a utilizar en esta investigación son:

- a. Correlación lineal
- b. Tabla de contingencia
- c. Componentes principales

4.2 DEFINICIONES BÁSICAS

Las técnicas de análisis multivariado están organizados con base en un esquema que divide los procedimientos en **métodos de dependencia**, éstos especifican una o más variables como si se hubiesen pronosticado mediante un conjunto de variables independientes (el análisis de varianza es un ejemplo de este tipo de análisis), y **métodos de interdependencia**, éstas tratan de explicar la interrelación entre todas las variables tomadas como un conjunto que interesa al investigador (el análisis de componentes principales es un ejemplo de este tipo de procedimiento).

Para poder utilizar estas técnicas es necesario representar la información disponible por medio de una tabla de datos rectangular que corresponden a las observaciones disponibles (n filas), y p columnas, correspondientes al número de características medidas.

		TABLA XCIII MATRIZ DE DATOS					
		Variables					
		1	2	3	...	P	
Individuos	1	X_{11}	X_{12}	X_{13}	...	X_{1p}	
	2	X_{21}	X_{22}	X_{23}	...	X_{2p}	
	3	X_{31}	X_{32}	X_{33}	...	X_{3p}	
	⋮	⋮	⋮	⋮	⋮	⋮	
	n	X_{n1}	X_{n2}	X_{n3}	...	X_{np}	

4.2.1 Tabla de contingencia

En muchas ocasiones, una muestra (de tamaño n) tomada de una población pueden clasificarse de acuerdo con dos criterios diferentes, por lo tanto, es interesante saber si los métodos de clasificación son estadísticamente independientes. Supóngase que el primer criterio tiene r niveles, y que el segundo tiene c niveles. Sea f_{ij} la frecuencia observada para el nivel i del primer método de clasificación y el nivel j del segundo criterio. En general, los datos aparecerán como se muestran en la tabla XCIV, usualmente se la conoce como **tabla de contingencia**.

TABLA XCIV DISEÑO DE UNA TABLA DE CONTINGENCIA						
		COLUMNAS				
		1	2	3	...	c
RENGLONES	1	f_{11}	f_{12}	f_{13}	...	f_{1c}
	2	f_{21}	f_{22}	f_{23}	...	f_{2c}
	3	f_{31}	f_{32}	f_{33}	...	f_{3c}
	⋮	⋮	⋮	⋮	⋮	⋮
	r	f_{r1}	f_{r2}	f_{r3}	...	f_{rc}

El interés recae en probar la hipótesis de que los métodos de clasificación renglón – columnas son independientes. Si se rechaza esta hipótesis, entonces se concluye que existe alguna relación entre los dos criterios de clasificación.

El procedimiento de prueba puede obtenerse por medio de un estadístico válido para n grande. En general tenemos, su p_{ij} es la probabilidad de que un elemento seleccionado al azar caiga en la ij -ésima celda, dado que las dos clasificaciones son independientes. Entonces $p_{ij} = u_i v_j$, donde u_i es la probabilidad de que un elemento seleccionado al azar pertenezca a la renglón de la clase i , y v_j es la probabilidad de que un elemento pertenezca a la columna de la clase j . Ahora, si se supone independencia, los estimadores de u_i y v_j son:

$$\hat{u}_i = \frac{1}{n} \sum_{j=1}^c f_{ij}$$

$$\hat{v}_j = \frac{1}{n} \sum_{i=1}^r f_{ij}$$

Por lo tanto, la frecuencia esperada de cada celda es

$$e_{ij} = n \hat{u}_i \hat{v}_j = \frac{1}{n} \sum_{j=1}^c f_{ij} \sum_{i=1}^r f_{ij}$$

Entonces, para n grande, el estadístico

$$\chi_o^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(f_{ij} - e_{ij})^2}{e_{ij}}$$

tiene una distribución aproximada ji-cuadrada con $(r-1)(c-1)$ grados de libertad. Por consiguiente, la hipótesis de independencia debe rechazarse si el valor observado del estadístico de prueba χ_o^2 es mayor que $\chi_{\alpha, (r-1)(c-1)}^2$.

4.2.2 Componentes principales

Es una técnica multivariada de interdependencia, que estudia p variables de investigación, a través de las cuales se generarán k variables latentes, $k < p$, que contengan aproximadamente tanta información como las p variables originales, donde los objetivos generales de este análisis son la reducción de los datos y la interpretación.

Sea $\mathbf{X}^t = [X_1, X_2, \dots, X_p]$ el vector aleatorio p -variado con media \mathbf{m} y matriz de varianza y covarianza (S), supongamos además que los valores propios de S son $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_p \geq 0$ y sus respectivos vectores propios ortonormales $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \dots, \mathbf{b}_p$; definamos p variables no observadas denotadas por Y_1, Y_2, \dots, Y_p como una combinación lineal de $X_1, X_2, X_3, \dots, X_p$, tal que:

$$\begin{aligned} Y_1 &= \beta_{11}X_1 + \beta_{12}X_2 + \dots + \beta_{1p}X_p \\ Y_2 &= \beta_{21}X_1 + \beta_{22}X_2 + \dots + \beta_{2p}X_p \\ &\vdots \\ Y_p &= \beta_{p1}X_1 + \beta_{p2}X_2 + \dots + \beta_{pp}X_p \end{aligned}$$

En síntesis, tenemos:

$$\begin{aligned} Y_i &= \beta_{i1}X_1 + \beta_{i2}X_2 + \dots + \beta_{ip}X_p = \mathbf{b}_i^t \mathbf{X} & \mathbf{b}_i, \mathbf{X} \in \mathbb{R}^p \\ &= [\mathbf{b}_1 \quad \mathbf{b}_2 \quad \dots \quad \mathbf{b}_p] \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} \end{aligned}$$

De resultados podríamos saber que:

$$E[Y_i] = \mathbf{b}_i^t \mathbf{m} \quad i = 1, 2, \dots, p$$

$$\text{Var}(Y_i) = \mathbf{b}_i^t \mathbf{S} \mathbf{b}_i = \lambda_i$$

$$\text{Cov}(Y_i, Y_j) = \mathbf{b}_i^t \mathbf{S} \mathbf{b}_j = 0 \quad i \neq j$$

Las componentes principales de \mathbf{X} son aquellas combinaciones lineales Y_1, Y_2, \dots, Y_p que son no correlacionadas entre sí y cuyas varianzas son tan grandes como sea posible.

Primera componente principal = combinación lineal $\mathbf{b}_1^t \mathbf{X}$ que maximiza

$$\text{Var}(\mathbf{b}_1^t \mathbf{X}) \text{ sujeto a:}$$

$$\langle \mathbf{b}_1, \mathbf{b}_1 \rangle = 1$$

Segunda componente principal = combinación lineal $\mathbf{b}_2^t \mathbf{X}$ que maximiza

$$\text{Var}(\mathbf{b}_2^t \mathbf{X}) \text{ sujeto a:}$$

$$\langle \mathbf{b}_2, \mathbf{b}_2 \rangle = 1 \text{ y}$$

$$\text{Cov}(\mathbf{b}_1^t \mathbf{X}, \mathbf{b}_2^t \mathbf{X}) = 0$$

$$\text{Var}(\mathbf{b}_2^t \mathbf{X}) > \text{Var}(\mathbf{b}_1^t \mathbf{X})$$

l -ésima componente principal = combinación lineal $\mathbf{b}_l^t \mathbf{X}$ que maximiza

$$\text{Var}(\mathbf{b}_l^t \mathbf{X}) \text{ sujeto a:}$$

$$\langle \mathbf{b}_l, \mathbf{b}_l \rangle = 1$$

$$\text{Cov}(\mathbf{b}_k^t \mathbf{X}, \mathbf{b}_l^t \mathbf{X}) = 0 ; k < l$$

$$\text{Var}(\mathbf{b}_l^t \mathbf{X}) > \text{Var}(\mathbf{b}_k^t \mathbf{X})$$

Se desea obtener la mayor proporción del total de varianza de la población, explicada por las k -ésima componentes principales, donde el valor individual de su aporte está dado por:

$$\frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} \quad \text{para } k = 1, 2, \dots, p$$

4.3 ANÁLISIS DE CORRELACIÓN LINEAL APLICADAS A LAS CARACTERÍSTICAS DE ESTUDIO

En la tabla XCV se muestra la matriz de correlación entre las variables de estudio para el año lectivo 1998-99, el propósito es determinar que variables están relacionadas la una con la otra, es decir que dependientes son entre sí. La mayoría de los coeficientes son bajas (valor absoluto menor a 0.30), tomaremos como valores significativos a los coeficientes cuyos valores absolutos sean mayores a 0.30, es decir tienen una correlación aceptable en el intervalo [0.3 ; 0.6) y correlación alta [0.6 ; 1.0), en la matriz de correlación (tabla XCV) los valores con **negritas** tienen correlaciones significativas. La más notable relación de dependencia lineal está entre el número de aulas que hay en la institución educativa (X_{10}) y el número de alumnos, sean estos promovidos (X_{12}), no promovidos (X_{13}) y desertores (X_{14}), en este caso los estimadores de coeficientes de correlación son iguales a 0.742, 0.636, 0.573 respectivamente.

La variable X_5 (sexo) está correlacionada significativamente con X_{12} y X_{13} , que corresponden respectivamente al número de estudiantes promovidos y reprobados.

Definición de las variables de la tabla XCV:

X_1 : Cantón	X_8 : Tenencia del edificio
X_2 : Parroquia	X_9 : Tipo de construcción
X_3 : Zona (urbana o rural)	X_{10} : Número de aulas
X_4 : Jornada	X_{11} : Personal del colegio
X_5 : Clasificación por sexo	X_{12} : Número de estud. promovidos
X_6 : Tipo de establecimiento	X_{13} : Número de estud. no promovidos
X_7 : Propietario del edificio	X_{14} : Número de estudiantes desertores

En la tabla XCV, observamos que entre la clasificación por sexo del colegio (X_5) y los estudiantes reprobados de la institución (X_{13}), tienen un coeficiente de correlación significativamente alto (-0.401), de esto podemos concluir que posee una relación inversa entre la variables, es decir, los alumnos de los colegios fiscales masculinos tienden a tener más alumnos reprobados que en los femeninos y mixtos. La variable X_5 también está correlacionada significativamente con X_{10} (número de aulas en los colegios), con un coeficiente de -0.382 , esta tiene una dependencia inversa, es decir, los colegios fiscales masculinos tienen más aulas que las demás instituciones educativas fiscales en la provincia del Guayas.

Con respecto a las variables X_{12} (número de estudiantes aprobados) y X_{13} (número de estudiantes reprobados) están correlacionadas significativamente (0.784), tienen una dependencia directa entre ella, es decir, los alumnos aprobados y reprobados aumentan en la mismas proporción en los colegios fiscales de la provincia del Guayas.

4.4 ANÁLISIS DE INDEPENDENCIA DE LAS CARACTERÍSTICAS EN ESTUDIO UTILIZANDO TABLA DE CONTINGENCIA

En esta sección determinaremos que características interaccionan, lo que desaseamos probar si existe independencia entre dos variables aleatorias. En la tabla XCVI se muestra un resumen de todas las pruebas de independencias realizadas para las variables (valor p de la prueba), los valores resaltados en negrillas indican los casos en que se rechazó la hipótesis de independencia (H_0).

H_0 : La variable X_i es independiente de la variable X_j ; para $i \neq j$
 Vs.
 H_a : $\neg H_0$

Veamos algunos casos:

X_6 (Tipo de establecimiento) vs. X_5 (Sexo): El tipo del establecimiento del nivel medio depende de su sexo, la mayor cantidad de instituciones son mixtas e instituciones regulares (informática, ciencias humanísticas, etc) llegando a obtener 85.6% del total provincial.

X_7 (Propietario) vs. X_8 (Tenencia): El propietario del edificio del establecimiento del nivel medio depende de su tenencia, la mayor cantidad de instituciones son propiedad del estado y cuya tenencia es propia, cuyo porcentaje es del 72.6% del total provincial.

En la tabla XCVI, se muestran todos los valores p de la prueba de independencia utilizando tablas de contingencias, tales que los indicados con **negritas** indican que hay dependencia entre la variable del i-ésimo renglón, con la j-ésima columna.

X_6 (Tipo de establecimiento) vs. X_9 (Tipo de construcción): El tipo de colegio fiscal es independiente del tipo de construcción del mismo, la mayoría de estas instituciones educativas del nivel medio de la provincia del Guayas son regulares y están construidas de hormigón armado, teniendo como porcentaje 71.8% del total provincial.

X_{12} (Estudiantes promovidos) vs. X_{13} (estudiantes reprobados)

H_0 : El número de estudiantes promovidos es independiente de los estudiantes reprobados

Vs.

H_a : $\neg H_0$

Estadístico de prueba χ^2	Grados de libertad	Valor p
191.475	30	1.92E-25

En este caso, el valor p es suficientemente pequeño, nos permite concluir que existe evidencia estadística para rechazar la hipótesis nula (H_0), es decir, las variables X_{12} y X_{13} son dependientes. La tabla XCVII muestra la tabla de contingencia asociada a la prueba de independencia.

Tabla XCVII
Tabla de contingencia para la prueba de
independencia entre las variables X_{12} vs X_{13}

			X_{13} : Estudiantes no promovidos						TOTAL	
			[0, 43)	[43, 86)	[86, 129)	[129, 172)	[172, 215)	[215, 258)	[258, 301)	f. j e. j %
X_{12}: Estudiantes promovidos	[0, 358)	f _{1j} e _{1j} %	160 126.5 59.5	10 20.3 3.7	1 11.4 0.4	0 4.4 0.0	0 1.3 0.0	0 5.1 0.0	0 1.9 0.0	171 171 63.6
	[358, 716)	f _{2j} e _{2j} %	28 35.5 10.4	9 5.7 3.3	8 3.2 3.0	2 1.2 0.7	0 0.4 0.0	1 1.4 0.4	0 0.5 0.0	48 48 17.8
	[716, 1074)	f _{3j} e _{3j} %	8 20.7 3.0	7 3.3 2.6	7 1.9 2.6	2 0.7 0.7	1 0.2 0.4	2 0.8 0.7	1 0.3 0.4	28 28 10.4
	[1074, 1432)	f _{4j} e _{4j} %	0 5.2 0.0	3 0.8 1.1	1 0.5 0.4	2 0.2 0.7	0 0.1 0.0	1 0.2 0.4	0 0.1 0.0	7 7 2.6
	[1432, 1790)	f _{5j} e _{5j} %	2 7.4 0.7	1 1.2 0.4	1 0.7 0.4	1 0.3 0.4	1 0.1 0.4	3 0.3 1.1	1 0.1 0.4	10 10 3.7
	[1790, 2148)	f _{6j} e _{6j} %	1 3.7 0.4	2 0.6 0.7	0 0.3 0.0	0 0.1 0.0	0 0.0 0.0	1 0.1 0.4	1 0.1 0.4	5 5 1.9
	TOTAL	f _{i.} e _{i.} %	199 199 74.0	32 32 11.9	18 18 6.7	7 7 2.6	2 2 0.7	8 8 3.0	3 3 1.1	269 269 100

Fuente: Dirección nacional de estudios (Base de dato año lect. 1998–99)

Analizando la tabla XCVII, observamos que en 59.5% de colegios fiscales de la provincia del Guayas, los estudiantes aprobaron entre [0 ; 358) y reprobaron [0; 43). Hay un pequeño porcentaje de instituciones educativas del nivel medio fiscales que tienen una gran cantidad de alumnos matriculados que representan 0.4% del total provincial, del cual aprobaron entre [1790;2148) y reprobaron entre [258 ; 301) estudiantes.