

CAPITULO IV

4. ANÁLISIS MULTIVARIADO DE LA POBLACIÓN INVESTIGADA

4.1. INTRODUCCIÓN

Luego de realizar el análisis de cada una de las variables tratadas en el capítulo anterior, en este capítulo se utilizará técnicas multivariadas que nos permitirá analizar las interacciones que tienen entre sí las variables y su correspondiente influencia en las unidades respectivas.

Para poder realizar este análisis utilizaremos en este capítulo cuatro técnicas idóneas:

1. Correlación lineal
2. Análisis Bivariado
3. Tablas de Contingencia
4. Componentes Principales
5. Correlación Canónica

La primera técnica estadística multivariada nos indicara si las variables analizadas se encuentran o no relacionadas linealmente entre sí; la segunda técnica nos permitirá realizar el análisis de las frecuencias obtenidas entre las diferentes categorías de las dos variables seleccionadas; con la tercera técnica se determinara si existe independencia o no entre las variables de estudio; por medio de la cuarta técnica nos permitirá concluir si las variables de estudio pueden ser expresadas a través de factores que agrupen características de estas variables y reduzcan considerablemente el número de variables de estudio; y por medio de la quinta técnica nos permitirá conocer si existe algún tipo de asociación entre 2 grupos de variables.

Para este análisis multivariado se ha decidido eliminar las variables de NACIONALIDAD (IP5) y LENGUA (IP6) de la información personal de los trabajadores y servidores públicos del MEC, se decidió esto en base de que más del 99% de la información proporcionada por estos trabajadores poseen nacionalidad ecuatoriana y hablan el idioma español, y no serán aleatorias en la práctica.

4.2 Método de imputación de datos

Antes de comenzar aplicar las técnicas estadísticas multivariada, se encontró que existe evidencia de que los datos obtenidos en el Censo del Magisterio no fueron completos por falta de respuesta o no declaración por parte del informante, para lo cual se procedió a imputar los datos de tal forma que si la variable tiene menos del 10% de ausencia de datos, se aplicó una técnica de simulación denominada Método de la Transformada Inversa, que consiste en generar números aleatorios (0-1) y observar en que intervalo de la distribución de frecuencia acumulada de cada variable encaja este número, y de esta manera asignar la codificación correspondiente al intervalo donde encajó dicho número, para lo cual se procedió hacer un software que permita realizar esta técnica debido a la gran

cantidad de información con la que se trabaja en cada variable de la base de datos del censo del Magisterio Fiscal en el litoral ecuatoriano; y si la variable tiene más del 10% de ausencia de información, se le asignó el valor de 0, en donde este valor representa para nuestro análisis información no proporcionada por el trabajador del MEC, en las variables que cumplan con este eventualidad.

4.3 Definiciones de estadística multivariada

El análisis estadístico puede obtener conclusiones significativas al hacer uso de las técnicas multivariadas; el presente estudio utilizará estas técnicas para determinar si existe algún efecto en la interacción de las variables, para esto se necesita algunas definiciones importantes que se detallarán a continuación.

4.3.1 Vector aleatorio

Sean X_1, X_2, \dots, X_p , p variables aleatorias o características sujetas a investigación, entonces se procede a definir un vector p variado $\mathbf{X} \in \mathbf{R}^p$, el cual está compuesto por p variables, esto es:

$$\mathbf{X}^t = [X_1 \ X_2 \ \dots \ X_p]$$

4.3.2 Matriz de datos

Se define a la matriz \mathbf{X} como el arreglo compuesto por p filas y n columnas, en donde cada elemento de este arreglo está determinado por x_{ij} , que representa la i -ésima variable o característica sujeta a investigación y la j -ésima unidad de investigación.

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \cdots & x_{pn} \end{bmatrix} = [X_1 \ X_2 \ \dots \ X_n]$$

Donde $X_1 \ X_2 \ \dots \ X_n$ representa una muestra de tamaño n tomada de una población p -variadas de tamaño N .

4.3.3 Vector de medias

Sea $\mathbf{X}^t = [X_1 \ X_2 \ \dots \ X_p]$, un vector p variado, se define al vector de medias como el valor esperado de \mathbf{X} .

$$\mathbf{m} = E[\mathbf{X}] = \begin{bmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_p] \end{bmatrix} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \vdots \\ \mathbf{m}_p \end{bmatrix}$$

4.3.4 Matriz de varianzas y covarianzas

Sea $\mathbf{X}^t = [X_1 \ X_2 \ \dots \ X_p]$, un vector p variado, se define la matriz de varianzas y covarianzas como el valor esperado del producto de la diferencia del vector p-variado y su vector de medias con su traspuesta.

$$\mathbf{S} = E[(\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^t]$$

$$\mathbf{S} = \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_p) \\ \text{cov}(X_2, X_1) & \text{cov}(X_2, X_2) & \cdots & \text{cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_p, X_1) & \text{cov}(X_p, X_2) & \cdots & \text{cov}(X_p, X_p) \end{bmatrix}$$

$$\mathbf{S} = \begin{bmatrix} \mathbf{s}_{11} & \mathbf{s}_{12} & \cdots & \mathbf{s}_{1p} \\ \mathbf{s}_{21} & \mathbf{s}_{22} & \cdots & \mathbf{s}_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{s}_{p1} & \mathbf{s}_{p2} & \cdots & \mathbf{s}_{pp} \end{bmatrix}$$

donde σ_{ii} es la varianza de la i -ésima variable y σ_{ij} es la covarianza entre la i -ésima y j -ésima variable, además $\sigma_{ij} = \sigma_{ji}$, por lo tanto Σ es simétrica y diagonalizable ortogonalmente.

$$i = 1, 2, \dots, p$$

$$j = 1, 2, \dots, p$$

4.3.5 Matriz de correlación

Sea S la matriz de varianza y covarianza de un vector aleatorio \mathbf{X} p -variado, se define $\mathbf{V}^{1/2}$ como una matriz diagonal.

$$\mathbf{V}^{1/2} = \begin{bmatrix} \sqrt{s_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{s_{22}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{s_{pp}} \end{bmatrix}$$

donde $\sqrt{s_{ii}}$ es la desviación estándar de la variable aleatoria \mathbf{X}_i , entonces se define a la matriz de correlación como:

$$\mathbf{r} = (\mathbf{V}^{1/2})^{-1} \boldsymbol{\Sigma} (\mathbf{V}^{1/2})^{-1} = \begin{bmatrix} \mathbf{s}_{11}/\mathbf{s}_1\mathbf{s}_1 & \mathbf{s}_{12}/\mathbf{s}_1\mathbf{s}_2 & \cdots & \mathbf{s}_{1p}/\mathbf{s}_1\mathbf{s}_p \\ \mathbf{s}_{21}/\mathbf{s}_2\mathbf{s}_1 & \mathbf{s}_{22}/\mathbf{s}_2\mathbf{s}_2 & \cdots & \mathbf{s}_{2p}/\mathbf{s}_2\mathbf{s}_p \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{s}_{p1}/\mathbf{s}_p\mathbf{s}_1 & \mathbf{s}_{p2}/\mathbf{s}_p\mathbf{s}_2 & \cdots & \mathbf{s}_{pp}/\mathbf{s}_p\mathbf{s}_p \end{bmatrix}$$

$$\mathbf{r} = \begin{bmatrix} 1 & \mathbf{r}_{12} & \cdots & \mathbf{r}_{1p} \\ \mathbf{r}_{21} & 1 & \cdots & \mathbf{r}_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{r}_{p1} & \mathbf{r}_{p2} & \cdots & 1 \end{bmatrix}$$

donde ρ_{ij} es el coeficiente de correlación entre las variables \mathbf{X}_i y \mathbf{X}_j

$$i = 1, 2, \dots, p$$

$$j = 1, 2, \dots, p$$

4.3.6 Componentes principales

Componentes principales es una técnica estadística multivariada que permite la reducción de datos, algebraicamente es una combinación lineal de p variables aleatorias X_1, X_2, \dots, X_p , geoméricamente componentes principales representa la elección de un nuevo sistema de coordenadas obtenidas al rotar el sistema original con el nuevo ejes de coordenadas formado por cada

componente. Los nuevos ejes representan la dirección de máxima variabilidad obtenida.

Sea $\mathbf{X}^t = [X_1 \ X_2 \ \dots \ X_p]$ un vector p variado, y cada una de las variables que lo componen son variables aleatorias observables y no necesariamente normales; sea Σ la matriz de varianza y covarianza del vector p variado \mathbf{X} , y sea $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ los valores propios correspondientes a Σ .

Considere las siguientes combinaciones lineales:

$$Y_1 = \mathbf{a}_1^t \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Y_2 = \mathbf{a}_2^t \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

$$\vdots$$

$$Y_p = \mathbf{a}_p^t \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

Se define a las varianzas y covarianzas de las componentes principales como:

$$Var (Y_i) = \mathbf{a}_i^t \Sigma \mathbf{a}_i \quad i = 1, 2, \dots, p$$

$$Cov (Y_i, Y_j) = \mathbf{a}_i^t \Sigma \mathbf{a}_j \quad j = 1, 2, \dots, p$$

Se definen a Y_1, Y_2, \dots, Y_p como las componentes principales, no correlacionadas y ortonormales entre sí, y además $\text{Var}(Y_1) \geq \text{Var}(Y_2) \geq \dots \geq \text{Var}(Y_p) \geq 0$. Por lo cual, éstas deben cumplir con:

$$\begin{aligned} \|\mathbf{a}_i\| &= 1 && \text{para } i = 1, 2, \dots, p \\ \langle \mathbf{a}_i, \mathbf{a}_j \rangle &= 0 && \text{para } i \neq j \end{aligned}$$

donde $\|\mathbf{a}_i\|$ es la norma del vector \mathbf{a}_i y $\langle \mathbf{a}_i, \mathbf{a}_j \rangle$ es el producto interno entre los vectores \mathbf{a}_i y \mathbf{a}_j .

La primera componente principal es la combinación lineal $Y_1 = \mathbf{a}_1' \mathbf{X}$ de máxima varianza, esto es que maximiza la varianza de Y_1 , sujeta a que la norma del vector \mathbf{a}_1 sea unitaria.

La segunda componente principal es la combinación lineal $Y_2 = \mathbf{a}_2' \mathbf{X}$, sujeta a que la norma del vector \mathbf{a}_2 sea unitaria, la $\text{Cov}(Y_1, Y_2) = 0$ y a que $\text{Var}(Y_2) \leq \text{Var}(Y_1)$.

La i -ésima componente principal es la combinación lineal $Y_i = \mathbf{a}_i' \mathbf{X}$, sujeta a que la norma del vector \mathbf{a}_i sea unitaria, la $\text{Cov}(Y_{i-1}, Y_i) = 0$ y a que $\text{Var}(Y_i) \leq \text{Var}(Y_{i-1})$.

para $i = 2, 3, \dots, p$

Como resultados obtenemos que si S es la matriz de covarianzas asociada con el vector $\mathbf{X}^t = [X_1 \ X_2 \ \dots \ X_p]$, donde S tiene los pares de valores y vectores propios $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$, donde $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Entonces se puede probar que la i -ésima componente principal Y_i es igual a:

$$Y_i = \mathbf{e}_i^t \mathbf{X} = e_{i1}X_1 + e_{i2}X_2 + \dots + e_{ip}X_p \quad i = 1, 2, \dots, p$$

Y además:

$$Var(Y_i) = \mathbf{e}_i^t \Sigma \mathbf{e}_i = \lambda_i \quad i = 1, 2, \dots, p$$

$$Cov(Y_i, Y_j) = \mathbf{e}_i^t \Sigma \mathbf{e}_j = 0 \quad j = 1, 2, \dots, p$$

El porcentaje total de la varianza contenida por la i -ésima componente principal, está dada por:

$$\frac{I_i}{\sum_{i=1}^p I_i}$$

4.3.7 Correlación canónica

El análisis de correlación canónica es una técnica estadística multivariada que mide el grado de asociación entre 2 grupos de variables.

Sea \mathbf{X} un vector aleatorio p variado ($\mathbf{X} \in \mathbb{R}^p$) el cual podemos particionarlo en 2 grupos de vectores de tamaño q y $p-q$ respectivamente, de la siguiente forma:

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_q \\ \dots \\ X_{q+1} \\ \vdots \\ X_p \end{bmatrix} = \begin{bmatrix} \mathbf{X}^{(1)} \\ \dots \\ \mathbf{X}^{(2)} \end{bmatrix}$$

Donde el vector de medias correspondiente a $\mathbf{X}^{(1)}$ y $\mathbf{X}^{(2)}$ es la partición del vector de medias $\mathbf{X} \in \mathbb{R}^p$, que es el siguiente:

$$\mathbf{m} = E[\mathbf{X}] = \begin{bmatrix} \mathbf{m}_1 \\ \vdots \\ \mathbf{m}_q \\ \dots \\ \mathbf{m}_{q+1} \\ \vdots \\ \mathbf{m}_p \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}^{(1)} \\ \boldsymbol{\mu}^{(2)} \end{bmatrix}$$

Donde:

$$\mathbf{X}^{(1)} \text{ y } \mathbf{m}^{(1)} \in \mathbb{R}^q$$

$$\mathbf{X}^{(2)} \text{ y } \mathbf{m}^{(2)} \in \mathbb{R}^{p-q}$$

Las matrices de varianzas para $\mathbf{X}^{(1)}$ y $\mathbf{X}^{(2)}$ son Σ_{11} y Σ_{22} respectivamente y la matriz de covarianza de $\mathbf{X}^{(1)}$ con $\mathbf{X}^{(2)}$ es $\Sigma_{12} = \Sigma_{21}^t$.

$$\Sigma = \begin{bmatrix} \mathbf{s}_{11} & \cdots & \mathbf{s}_{1q} & \mathbf{s}_{1,q+1} & \cdots & \mathbf{s}_{1p} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{s}_{q1} & \cdots & \mathbf{s}_{qq} & \mathbf{s}_{q,q+1} & \cdots & \mathbf{s}_{qp} \\ \hline \mathbf{s}_{q+1,1} & \cdots & \mathbf{s}_{q+1,q} & \mathbf{s}_{q+1,q+1} & \cdots & \mathbf{s}_{q+1,p} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{s}_{p1} & \cdots & \mathbf{s}_{pq} & \mathbf{s}_{p,q+1} & \cdots & \mathbf{s}_{pp} \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

Ahora consideremos las siguientes combinaciones lineales:

$$U = \mathbf{a}^t \mathbf{X}^{(1)}$$

$$V = \mathbf{b}^t \mathbf{X}^{(2)}$$

Y su respectiva varianza y covarianza son:

$$Var(U) = \mathbf{a}^t \Sigma_{12} \mathbf{a}$$

$$Var(V) = \mathbf{b}^t \Sigma_{12} \mathbf{b}$$

$$Cov(U, V) = \mathbf{a}^t \Sigma_{12} \mathbf{b}$$

Nosotros buscaremos coeficientes de \mathbf{a} y \mathbf{b} tal que:

$$Corr(U, V) = \frac{\mathbf{a}^t \Sigma_{12} \mathbf{b}}{\sqrt{\mathbf{a}^t \Sigma_{11} \mathbf{a}} \sqrt{\mathbf{b}^t \Sigma_{22} \mathbf{b}}}$$

Basándonos en esto, definimos:

El primer par de variables canónicas, es el par de combinaciones lineales U_1, V_1 que tiene varianza unitaria y que maximiza la correlación entre ambas.

El segundo par de variables canónicas, es el par de combinaciones lineales U_2, V_2 que tiene varianza unitaria y que maximiza la correlación entre ambas, y además en todos los casos no está correlacionada con el primer par de variables canónicas.

El k -ésimo par de variables canónicas, es el par de combinaciones lineales U_k, V_k que tiene varianza unitaria y que maximiza la correlación entre ambas, y además en todos los casos no está correlacionada con las $k-1$ pares de variables canónicas previas.

Se denomina a la correlación entre el k -ésimo par de variables canónicas la k -ésima correlación canónica.

Para encontrar los vectores \mathbf{a} y \mathbf{b} nos basamos en los siguientes resultados:

Suponga $q \leq p$ y que los vectores $\mathbf{X}^{(1)}$ y $\mathbf{X}^{(2)}$ tienen:

$$\text{Cov}(\mathbf{X}^{(1)}) = \Sigma_{11}$$

$$\text{Cov}(\mathbf{X}^{(2)}) = \Sigma_{22}$$

$$\text{Cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) = \Sigma_{12} = \Sigma_{21}^t$$

Los coeficientes de los vectores \mathbf{a} y \mathbf{b} , para la combinación lineal

$$U = \mathbf{a}^t \mathbf{X}^{(1)}$$

$$V = \mathbf{b}^t \mathbf{X}^{(2)}$$

$$\max \text{Corr}(U, V) = r_1^*$$

son:

a, b

Logrando el k-ésimo par de variables canónicas:

$$U_k = \mathbf{e}_k^t \mathbf{S}_{11}^{-1/2} \mathbf{X}^{(1)}$$

$$V_k = \mathbf{f}_k^t \mathbf{S}_{22}^{-1/2} \mathbf{X}^{(2)}$$

Con:

$$\text{Corr}(U_k, V_k) = r_k^*$$

Donde r_k^* representa la k-ésima correlación canónica entre las k-ésimas variables canónicas, y además $r_1^{*2}, r_2^{*2}, \dots, r_p^{*2}$ son los valores propios de la matriz resultado de la multiplicación de: $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$ y $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$ son los vectores propios asociados a ésta, y $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$ son los vectores propios de la matriz obtenida de la multiplicación de $\Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$.

Las variables canónicas tienen las siguientes propiedades:

$$\text{Var}(U_k) = \text{Var}(V_k) = 1$$

$$\text{Cov}(U_k, U_l) = \text{Cov}(U_l, U_k) = 0 \quad k \neq l$$

$$\text{Cov}(V_k, V_l) = \text{Cov}(V_l, V_k) = 0 \quad k \neq l$$

$$\text{Cov}(U_k, V_l) = \text{Cov}(U_l, V_k) = 0 \quad k \neq l$$

para $k, l = 1, 2, \dots, p$