

Método del Análisis Factorial

1. *Introducción*
2. *Metodología Estadística*
 - 2.1 *Pruebas estadísticas preliminares*
 - 2.2 *El modelo del análisis factorial*
 - 2.3 *Métodos de estimación*
 - 2.4 *Rotación de factores*
 - 2.5 *Interpretación de los factores*
 - 2.6 *Puntuaciones o Scores.*
3. *Ejemplo de aplicación*

1. Introducción

El AF y el ACP son técnicas para examinar interdependencia de variables. El objetivo del ACP es explicar la mayor parte de la variabilidad total de un conjunto de variables con el menor número de Comp. posibles. En el AF, los Fact. son seleccionados para explicar las interrelaciones entre variables. El AF es un método que permite construir un modelo para explicar la correlación existente entre un conjunto de variables, en términos de otro conjunto de menor número de variables denominadas factores. En el AF las variables originales juegan el rol de variables dependientes, cuya interrelación se desea que sean explicadas con la selección de un conjunto de factores (comunes y únicos) que no son observables y que deberán ser hallados.

El AF tiene como objetivo explicar la estructura causal que origina las relaciones entre un conjunto de variables, así como la variación específica de cada una de ellas.

El AF, presupone la existencia de un conjunto de variables **subyacentes** (factores latentes) desconociendo cuántas y cuáles son, pero que deberán ser halladas en base de la estructura de correlación o variabilidad y de la explicación que puede aportar cada una de ellas sobre las variables originales.

- El ACP es una técnica de reducción de datos que se sitúa en el campo de la Estadística Descriptiva, mientras que el AF implica la elaboración de un modelo que requiere la formulación de hipótesis estadísticas y la aplicación de métodos de inferencia.

2.1 Pruebas Estadísticas Preliminares

1) Prueba de los coeficientes de correlación

Examinando la matriz de correlaciones se puede evidenciar la existencia de correlaciones significativas. Para corroborar se realiza la prueba estadística de significación de cada uno de los coeficientes de correlación:

Formulación de hipótesis: $H_p: \rho = 0$

$H_a: \rho \neq 0$

Cálculo de la prueba estadística: $t_C = \frac{r - \rho}{S_r}$ donde: $S_r = \sqrt{\frac{1 - r^2}{n - 2}}$

Decisión estadística: Se acepta H_p , si $t_{(\frac{\alpha}{2}, n-2)} \leq t_C \leq t_{(1-\frac{\alpha}{2}, n-2)}$

2.1 Pruebas Estadísticas Preliminares

2) Prueba de esfericidad de Barlett

Permite probar si existe una intercorrelación significativa entre las variables originales.

Así, para p variables se puede usar la matriz de correlación R_p cuyos elementos de la diagonal son 1s y los que están fuera de la diagonal son coeficientes que miden la intercorrelación entre cada par de variables. Si todos estos coeficientes son nulos (no existe correlaciones entre las p variables), entonces la matriz R_p será igual a la identidad, con lo que su determinante será igual a la unidad.

Formulación de las hipótesis:

$$H_p : |R_p| = 1$$

$$H_a : |R_p| \neq 1$$

Estadística de Barlett:

$$\chi_c^2 = - \left[n - 1 - \frac{1}{6}(2p + 5) \right] \ln |R|$$

Decisión estadística. Se rechaza H_p , si $\chi_c^2 \geq \chi^2_{(1-\alpha)(p^2-p)}$. Si se acepta H_p , entonces las variables no están correlacionadas y por lo tanto no tiene sentido aplicar el análisis de componentes principales.

2.1 Pruebas Estadísticas Preliminares

3. Con el Software R se hará la prueba de normalidad p-variada de Shapiro.

H0: Las p variables tienen distribución normal p-variada.

NOTA: Desde un punto de vista estadístico, se pueden obviar los supuestos de normalidad, homocedasticidad y linealidad siendo conscientes que su incumplimiento produce una disminución en las correlaciones observadas.. En realidad sólo es necesaria la normalidad cuando se aplica una prueba estadística a la significación de los factores; sin embargo raramente se utilizan estas pruebas. Es deseable que haya cierto grado de multicolinealidad, porque uno de los objetivos es identificar series de variables interrelacionadas.

2.1 Pruebas Estadísticas Preliminares

3) Ajuste del Modelo

Para evaluar la adecuación el conjunto de datos de la muestra al análisis factorial, se utiliza el estadístico propuesto por Kaiser-Meyer-Olkin (KMO), definido como:

$$KMO = \frac{\sum \sum_{h \neq j} r_{jh}^2}{\sum \sum_{h \neq j} r_{jh}^2 + \sum \sum_{h \neq j} a_{jh}^2}$$

donde r_{jh} son los coeficiente de correlación simple entre las variables
 a_{jh} son los coeficiente de correlación parcial entre las variables

En el caso que existiera una adecuación de los datos, los coeficientes de correlación parcial serán pequeños y por consiguiente el KMO estará próximo a 1. El KMO varía entre 0 y 1. Para valores de KMO menores a 0.5 se considerará no aceptable la aplicación del análisis factorial al conjunto de datos.

2.1 Pruebas Estadísticas Preliminares

4) Bondad de ajuste de cada variable

También se propone una medida de adecuación para cada una de las variables. Los índices de adecuación es:

$$MSA_j = \frac{\sum_{h \neq j} r_{jh}^2}{\sum_{h \neq j} r_{jh}^2 + \sum_{h \neq j} a_{jh}^2}$$

Si MSA_j se próximo a 1, indicaría que la variable j es adecuada para el análisis factorial. Por el contrario, para valores de MSA_j menores a 0.5 se considera a la variable no aceptable para el AF y puede ser eliminada para el análisis.

2.2. El Modelo del Análisis Factorial

Modelo de Análisis Factorial

$$X_1 = l_{11}F_1 + l_{12}F_2 + l_{13}F_3 + \cdots + l_{1m}F_m + e_1$$

$$X_2 = l_{21}F_1 + l_{22}F_2 + l_{23}F_3 + \cdots + l_{2m}F_m + e_2$$

.....

$$X_p = l_{p1}F_1 + l_{p2}F_2 + l_{p3}F_3 + \cdots + l_{pm}F_m + e_p$$

Donde las X_i , $i = 1, \dots, p$ son las p variables observables que están tipificadas (tienen media cero y varianza 1), las F_i , $i = 1, \dots, m$ son m factores comunes ($m < p$) y las e_i , $i = 1, \dots, p$ son los factores únicos o específicos de cada variable X_i . Las l_{jh} son las cargas factoriales (son los pesos del factor común h en la variable j).

Los factores comunes y los específicos no son observables.

2.2. El Modelo del Análisis Factorial expresado en forma matricial

El modelo en forma matricial

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ \vdots \\ X_p \end{bmatrix} = \begin{bmatrix} l_{11} & l_{12} & l_{13} & \cdots & l_{1m} \\ l_{21} & l_{22} & l_{23} & \cdots & l_{2m} \\ l_{31} & l_{32} & l_{33} & \cdots & l_{3m} \\ \vdots & & & & \\ l_{p1} & l_{p2} & l_{p3} & \cdots & l_{pm} \end{bmatrix} \begin{bmatrix} F_1 \\ F_2 \\ F_3 \\ \vdots \\ F_m \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_p \end{bmatrix}$$

El modelo en forma matricial condensada

$$X = Lf + e$$

2.2. El Modelo del Análisis Factorial: hipótesis o supuestos

Para realizar inferencias con el modelo anterior se debe formular supuestos estadísticos acerca de los factores comunes y los factores únicos.

1. La esperanza de cada uno de los factores comunes es cero:

$$E(f) = 0$$

2. La matriz de covarianzas de los factores comunes es la matriz identidad, esto indica que la varianza de cada uno de los factores comunes es 1 y que no están correlacionados (ceros fuera de la diagonal):

$$E(f f') = I$$

3. La esperanza de cada uno de los factores únicos es cero:

$$E(e) = 0$$

2.2. El Modelo del Análisis Factorial: hipótesis o supuestos

4. La matriz de covarianza de los factores únicos o específicos es una matriz diagonal (pueden tener varianzas diferentes y no están correlacionadas porque fuera de la diagonal hay ceros):

$$E(ee') = \Omega$$

5. La matriz de covarianza entre los factores comunes y los factores específicos es nula (esto indica que los factores comunes y únicos no están correlacionados):

$$E(fe') = 0$$

2.2 Propiedades del Modelo

Como las variables X están tipificadas, su matriz de covarianzas es igual a la matriz de correlación poblacional R_p .

$$E(X'X) = R_p = \begin{bmatrix} 1 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{21} & 1 & \cdots & \rho_{2p} \\ & & \ddots & \\ \rho_{p1} & \rho_{p2} & \cdots & 1 \end{bmatrix}$$

Teniendo en cuenta las hipótesis anteriores R_p se puede descomponer:

$$R_p = LL' + \Omega$$

$$\begin{bmatrix} 1 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{21} & 1 & \cdots & \rho_{2p} \\ & & \ddots & \\ \rho_{p1} & \rho_{p2} & \cdots & 1 \end{bmatrix} = \begin{bmatrix} l_{11} & l_{12} & l_{13} & \cdots & l_{1m} \\ l_{21} & l_{22} & l_{23} & \cdots & l_{2m} \\ & & \vdots & & \\ l_{p1} & l_{p2} & l_{p3} & \cdots & l_{pm} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} & l_{31} & \cdots & l_{m1} \\ l_{12} & l_{22} & l_{32} & \cdots & l_{m2} \\ & & \vdots & & \\ l_{1m} & l_{2m} & l_{3m} & \cdots & l_{pm} \end{bmatrix} + \begin{bmatrix} \omega_1^2 & 0 & 0 & \cdots & 0 \\ 0 & \omega_2^2 & 0 & \cdots & 0 \\ & & \vdots & & \\ 0 & 0 & 0 & \cdots & \omega_p^2 \end{bmatrix}$$

2.2 Propiedades del Modelo

Propiedades del Modelo:

De lo anterior, la varianza de la variable tipificada X_1 se puede descomponer de la siguiente forma:

$$Var(X_1) = 1 = l_{11}^2 + l_{12}^2 + l_{13}^2 + \dots + l_{1m}^2 + \omega_1^2$$

Generalizando, la varianza de la variable tipificada X_j se puede descomponer de la siguiente forma:

$$Var(X_j) = 1 = l_{j1}^2 + l_{j2}^2 + l_{j3}^2 + \dots + l_{jm}^2 + \omega_j^2$$

Se conoce como la comunalidad a: $h_j^2 = l_{j1}^2 + l_{j2}^2 + l_{j3}^2 + \dots + l_{jm}^2$ y com especificidad a ω_j^2 . En conclusión:

$$Var(X_j) = 1 = h_j^2 + \omega_j^2$$

2.3 Métodos de Estimación

La estimación de los parámetros del modelo factorial implica el uso de métodos que conlleven a lo que se denomina la extracción de factores, determinando así el número de factores comunes a ser seleccionados en el análisis factorial. Existen varios métodos:

- 1) Método de los componentes principales (*método exploratorio*), cuando se busca la existencia y número de los factores
- 2) Método de máxima verosimilitud (*método confirmatorio*)
- 3) Método de mínimos cuadrados no ponderados
- 4) Método de mínimos cuadrados generalizados

2.4 Métodos de Rotación de Factores

La rotación de factores tiene como finalidad determinar fácilmente los factores comunes y la interpretación de sus interrelaciones con cada una de las variables originales. Los factores rotados obtenidos a partir de una solución inicial, presentarán una correlación alta (próxima a 1) con uno o grupo de variables originales y correlaciones bajas (próxima a 0) con el resto de variables. Así, se puede identificar rasgos o características comunes en un grupo de variables asociadas a un factor y dar una denominación a estas interrelaciones encontradas entre las variables del grupo.

- 1) Rotación ortogonal
- 2) Rotación oblicua.

2.5 Interpretación de los factores

Correlación entre las componentes principales y las variables originales

La correlación entre la j-ésima variable original y la k-ésima componente principal, representa el grado de asociación entre ellas y su valor cuantificará la proporción de la variación total de la variable original j que es explicada por la componente k. Así se tiene:

$$r_{jh} = \frac{w_{hj} \sqrt{\lambda_h}}{\sqrt{\text{Var}(X_j)}} \quad \text{Si la variable } X_j \text{ esta tipificada: } r_{jh} = w_{hj} \sqrt{\lambda_h}$$

donde los w_{hj} son los coeficientes de los autovectores o cargas factoriales.

Estas correlaciones también representan la parte de variancia de cada variable que es explicada cada factor. Se cumple que la suma horizontal de los cuadrados de las cargas factoriales de una variable es igual a uno.

2.6 Puntuaciones o Scores

De acuerdo al modelo de Análisis Factorial, se hace la estimación de las puntuaciones o scores, que sirven para verificar la ortogonalidad de los factores y como datos de entrada para otros análisis estadísticos.

3. Ejemplo de aplicación

La base a datos que se utilizó corresponde a alpacas adultas en actividad reproductiva. Las alpacas adultas corresponden a animales de más de dos años de edad; edad en la cual son capaces de producir crías.

Los datos que se procesaron fueron recolectados durante la esquila del año 2019 en alpacas del plantel de reproductores de la SAIS Pachacutec en Junín.

Las variables se describen a continuación:

Descripción de las variables

| Variable | Descripción |
|--|---|
| Peso vivo=Peso | Medido en unidades de Kilos (Kg.). El registro se realiza momentos previos a la esquila con una balanza de plataforma. |
| Longitud de mecha =Mecha | Medido en Cm. El registro se realiza momentos posteriores a la esquila. Se utiliza una regla de 30cm, midiendo la longitud de la mecha de fibra en posición perpendicular al cuerpo del animal. |
| Diámetro promedio de fibra (DF) = Diametro | Medido en micras. Corresponde al promedio del diámetro medido de 1000 fibras. Se midieron 1000 fibras por animal. |
| Desviación estándar del DF = DesvE | Medido en micras. Corresponde a la desviación estándar de las 1000 fibras medidas por animal. |
| Coeficiente de variación del (DF) =CV | Medido en %. Corresponde al coeficiente de variación de las 1000 fibras medidas por animal. |
| Factor picazón = Picazon | Medido en %. Corresponde al porcentaje de fibras que presentan un diámetro de fibra mayor a 30.5 micras; de las 1000 fibras analizadas. |

4. Ejemplo de aplicación

Use el archivo **Distritos Peruanos** para hacer el Análisis Factorial.

Las variables son:

4. Ejemplo de aplicación

- NOASIST1: Porcentaje de población en el hogar de 6 a 11 años que no asiste a un centro educativo.
- NOASIST2: Porcentaje de población en el hogar de 2 a 17 años que no asiste a un centro educativo.
- ANALFT: Porcentaje de población en el hogar de 15 años a más que no saben leer ni escribir..
- EDUPRIM1: Porcentaje de población en el hogar de 15 años a más que tiene educación primaria completa.
- EDUSEC1: Porcentaje de población en el hogar de 18 años a más que tiene educación secundaria completa.
- EDUSUP1: Porcentaje de población en el hogar de 18 años a más que tiene educación superior no universitaria completa.
- ICASTELL: Porcentaje de población con idioma castellano como lengua aprendida desde la niñez.
- INATIVA: Porcentaje de población con idioma quechua, aymara u otra lengua nativa como lengua aprendida desde la niñez
- EDUYEARS: N° de años promedio de educación de la población en el hogar.
- EDU1564: N° de años promedio de educación de la población en el hogar de 15 y 64 años.
- EDU1599: N° de años promedio de educación de la población en el hogar de 15 años y más.
- EDUJEFE: N° de años promedio de educación del jefe del hogar.
- EDUCONY: N° de años promedio de educación del cónyuge del jefe del hogar