

MÉTODOS PARA REDUCIR LA DIMENSIÓN: ANÁLISIS DE COMPONENTES PRINCIPALES

Autor: Juana María Alonso Revenga

OBJETIVOS Y COMPETENCIAS A ALCANZAR.

- Extraer las Componentes Principales de una matriz de datos.
- Conocer las reglas que justifiquen el número de Componentes Principales a retener.
- Interpretar, cuando ello sea posible, las Componentes Principales.
- Obtener la proyección de las observaciones sobre cada componente principal
- Conocer la utilidad que tienen las Componentes Principales para utilizarse como entrada en otras técnicas estadísticas.
- Conocer las hipótesis del modelo Factorial.
- Conocer el grado de satisfacción obtenido desde el punto de vista estadístico del Análisis Factorial llevado a cabo.
- Conocer los fundamentos de la rotación de los Factores originales.
- Proyectar las observaciones sobre el plano Factorial.
- Representar gráficamente tanto las variables como las observaciones sobre el espacio Factorial obtenido

.1.1. INTRODUCCIÓN.

Cuando se recoge la información de una muestra de datos, lo más frecuente es tomar el mayor número posible de variables. Sin embargo, si tomamos demasiadas variables **es difícil visualizar relaciones** entre ellas.

Otro problema que se presenta es la **fuerte correlación** que muchas veces se presenta entre las variables: si tomamos demasiadas variables (cosa que en general sucede cuando no se sabe demasiado sobre los datos o sólo se tiene ánimo exploratorio), lo normal es que estén relacionadas o que midan lo mismo bajo distintos puntos de vista. Por ejemplo, en estudios médicos, la presión sanguínea a la salida del corazón y a la salida de los pulmones están fuertemente relacionadas.

Se hace necesario, pues, reducir el número de variables. Es importante resaltar el hecho de que el concepto de **mayor información se relaciona con el de mayor variabilidad o varianza**. Cuanto mayor sea la variabilidad de los datos (varianza) se considera que existe mayor información, lo cual está relacionado con el concepto de entropía.

Estas técnicas fueron inicialmente desarrolladas por Pearson a finales del siglo XIX y posteriormente fueron estudiadas por Hotelling en los años 30 del siglo XX. Sin embargo, hasta la aparición de los ordenadores no se empezaron a popularizar.

Para estudiar las relaciones que se presentan entre p variables correlacionadas (que miden información común) se puede transformar el conjunto original de variables en otro (que no tenga repetición o redundancia en la información) **llamado conjunto de componentes principales**. Las nuevas variables son combinaciones lineales de las anteriores y se van construyendo según el orden de importancia en cuanto a la variabilidad total que recogen de la muestra. De modo ideal, se buscan $m < p$ variables que sean combinaciones lineales de las p originales y que estén incorreladas, recogiendo la mayor parte de la información o variabilidad de los datos.

Si las variables originales están incorreladas de partida, entonces no tiene sentido realizar un análisis de componentes principales.

El análisis de componentes principales es una técnica matemática que no requiere la suposición de normalidad multivariante de los datos, aunque si esto último se cumple se puede dar una interpretación más profunda de dichos componentes. **El análisis factorial tiene por objeto explicar si un conjunto de variables observadas por un pequeño número de variables latentes, o no observadas**, que llamaremos factores. El análisis factorial surge impulsado por el interés de Karl Pearson y Charles Spearman en comprender las dimensiones de la inteligencia humana en los años 30, y muchos de sus avances se han producido en el área de la psicometría.

El análisis factorial está relacionado con los componentes principales, pero existen ciertas diferencias. En primer lugar, **las componentes principales se construyen para explicar las varianzas, mientras que los factores se construyen para explicar las covarianzas o correlaciones entre las variables**. En segundo lugar, **componentes principales es un herramienta descriptiva, mientras que el análisis factorial presupone un modelo estadístico formal** de generación de la muestra dada.

El Análisis de Componentes Principales, es en realidad un caso particular del Análisis Factorial, y que por motivos de tiempo es el que estudiaremos en el máster

1.2. ANÁLISIS DE COMPONENTES PRINCIPALES.

La idea básica del método de Análisis de Componentes Principales es describir la variación de un conjunto de datos multivariante en términos de un conjunto incorrelado de variables, cada una de las cuales es una combinación lineal particular de las variables originales. Estas nuevas variables, denominadas **Componentes Principales**, se obtienen en orden decreciente a su importancia, es decir:

- **la primera Componente Principal recoge la máxima información** (variación) de los datos originales,
- **la segunda Componente se elige de forma que recoja la mayor cantidad de información, que no haya sido recogida por la primera Componente, es decir, recogerá información de las variables iniciales que sean incorreladas con la primera Componente Principal, y así sucesivamente.**

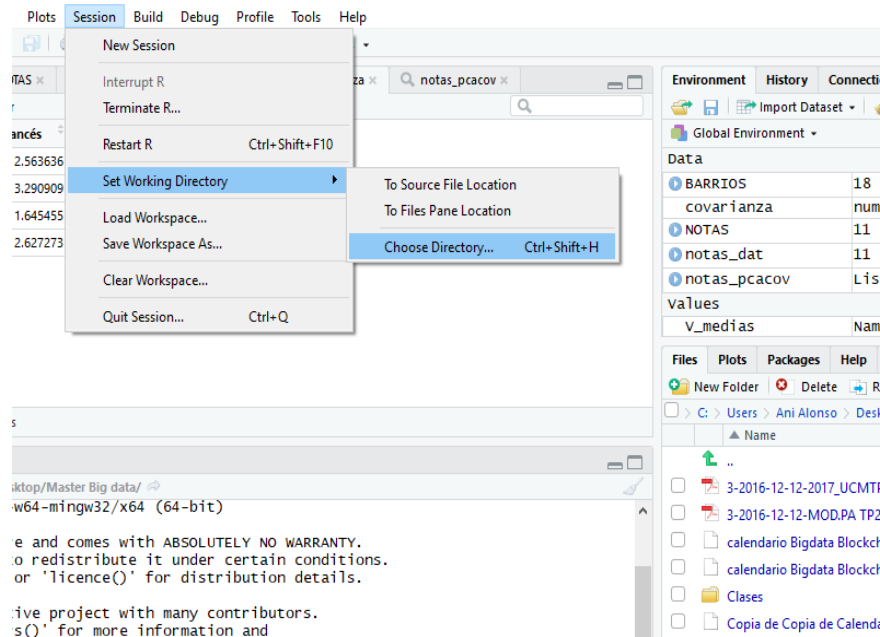
El objetivo de este análisis es ver si las primeras Componentes Principales recogen la mayor parte de la variación de los datos originales. Si esto es así, dichas Componentes se pueden utilizar para resumir los datos con la mínima pérdida de información. Esto dará lugar a importantes simplificaciones en los análisis posteriores que se puedan llevar a cabo. Al examinar este nuevo conjunto de variables se pueden obtener respuestas a preguntas importantes dentro del análisis de datos como: existencia de datos atípicos, multicolinealidad entre las variables, agrupación de observaciones en subgrupos y comprobar hipótesis de distribución normal multivariante

EJEMPLO:

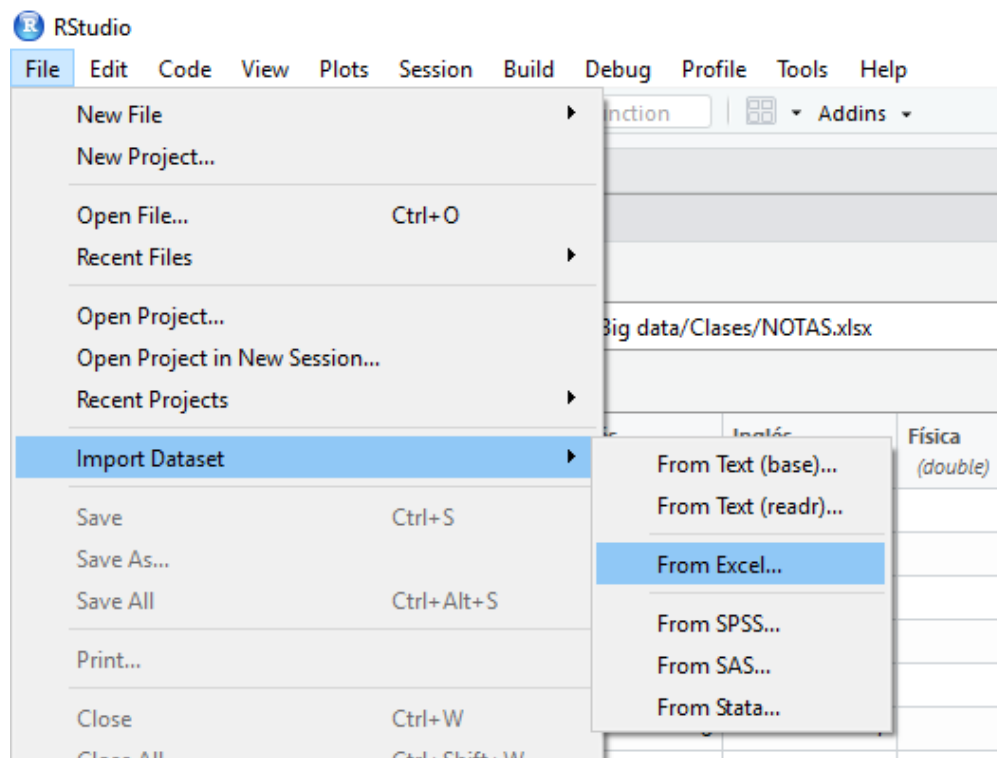
Se desea analizar las características de un grupo de alumnos con el objetivo de asignarles un índice que les oriente en cuanto a sus facultades para proseguir sus estudios. Para ello, se consideran las calificaciones obtenidas en Matemáticas, Francés, Inglés y Física. Los datos de partida son los siguientes:

| Obs | Mats | Francés | Inglés | Física |
|-----|------|---------|--------|--------|
| 1 | 1 | 4 | 5 | 3 |
| 2 | 5 | 5 | 4 | 4 |
| 3 | 6 | 10 | 9 | 6 |
| 4 | 3 | 6 | 6 | 4 |
| 5 | 2 | 4 | 6 | 1 |
| 6 | 6 | 8 | 7 | 8 |
| 7 | 6 | 6 | 6 | 7 |
| 8 | 3 | 5 | 8 | 4 |
| 9 | 8 | 5 | 5 | 8 |
| 10 | 2 | 5 | 8 | 4 |
| 11 | 9 | 7 | 7 | 9 |

Para comenzar la sesión en R elegimos la carpeta en donde vamos a crear nuestros conjuntos de datos



Importamos los datos desde Excel



A continuación cargamos las librerías que vamos a necesitar.

```
install.packages("readxl")
```

```
install.packages("knitr")
```

```
install.packages("ggplot2")
```

```
install.packages("factoextra")
```

```
install.packages("pastecs")
install.packages("FactoMineR")
install.packages("corrplot")
```

```
library(readxl)
library(ggplot2)
library(factoextra)
library(pastecs)
library(knitr)
library(FactoMineR)
library(corrplot)
```

1.2.1. ¿CÓMO SE OBTIENEN LAS COMPONENTES PRINCIPALES?

Desde un punto de vista algebraico las Componentes Principales son una combinación particular de p variables aleatorias. Geométricamente, estas combinaciones lineales representan la selección de un nuevo sistema de coordenadas obtenido mediante la rotación del sistema original. Los nuevos ejes representan las direcciones con máxima variabilidad y producen una descripción más simple y ordenada de la estructura de la matriz de Varianzas-Covarianzas.

Sea la variable $X = (X_1, X_2, \dots, X_p)^T$ cuyo vector de medias es

$$\mu_X = E[X] = (\mu_1, \mu_2, \dots, \mu_p)^T$$

y matriz de varianzas-covarianzas $\Sigma = E[(X - \mu_X)(X - \mu_X)^T]$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdot & \cdot & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \sigma_{p1} & \sigma_{p2} & \cdot & \cdot & \sigma_{pp} \end{pmatrix}$$

Luego $V(X_i) = \sigma_{ii} = E[(X_i - \mu_i)^2]$ y $COV(X_i, X_j) = \sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)]$.

Supongamos que disponemos de un conjunto de datos, observaciones de la variable X , dispuesto de forma matricial (con n observaciones y p variables):

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdot & \cdot & x_{1p} \\ x_{21} & x_{22} & \cdot & \cdot & x_{2p} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & \cdot & \cdot & x_{np} \end{pmatrix}$$

El estimador del vector de medias será $\bar{X} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)^T$ donde

$$\bar{x}_i = \frac{\sum_{k=1}^n x_{ik}}{n}$$

y el estimador de la matriz de varianzas-covarianzas

$$S = \begin{pmatrix} S_{11} & S_{12} & \cdot & \cdot & S_{1p} \\ S_{21} & S_{22} & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ S_{p1} & S_{p2} & \cdot & \cdot & S_{pp} \end{pmatrix}$$

$$\text{Donde } S_{ii} = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2}{n} \quad \text{y} \quad S_{ij} = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{n}.$$

A partir de la matriz S se puede calcular la matriz de correlaciones muestrales

$$R = \begin{pmatrix} 1 & r_{12} & \cdot & \cdot & r_{1p} \\ r_{21} & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ r_{p1} & r_{p2} & \cdot & \cdot & 1 \end{pmatrix}$$

$$\text{Donde } r_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}} \sqrt{S_{jj}}}.$$

Estos estimadores se pueden calcular utilizando notación matricial, como indicamos a continuación.

Expresiones matriciales.

Los estadísticos anteriores, \bar{X} , S , R , se pueden escribir en función de la matriz de datos, X , según las siguientes expresiones matriciales:

Vector de medias muestrales:

$$\bar{\mathbf{X}} = \frac{1}{n} \mathbf{X}' \mathbf{1}_n$$

Matriz de Covarianzas muestral:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n' \right) \mathbf{X}$$

Matriz de Correlaciones muestrales:

$$\mathbf{R} = \mathbf{D}^{-1/2} \mathbf{S} \mathbf{D}^{-1/2}$$

donde : $\mathbf{D} = \begin{pmatrix} s_{11} & 0 & \dots & 0 \\ 0 & s_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & s_{pp} \end{pmatrix}$

Existe la siguiente relación entre ellas:

$$\Rightarrow \mathbf{S} = \mathbf{D}^{1/2} \mathbf{R} \mathbf{D}^{1/2}$$

Luego, partiendo bien de la matriz de Correlaciones o de la matriz de Varianzas-Covarianzas estimadas, se hallará su descomposición, en función de sus valores propios y la matriz formada por sus autovectores correspondientes.

$$\mathbf{S} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^{-1} \quad \text{o bien} \quad \mathbf{R} = \mathbf{P}_R \mathbf{\Lambda}_R \mathbf{P}_R^{-1}.$$

Donde \mathbf{P} es la matriz de contiene los autovectores y $\mathbf{\Lambda}$ una matriz diagonal formada por los autovalores. Según se ha explicado en el apartado anterior, la primera Componente Principal, será aquella variable cuyo vector director coincide con el del autovector correspondiente al mayor autovalor de la matriz de partida, sea la de Correlaciones o la de Varianzas-Covarianzas. **Designaremos a partir de este momento por e_α a los autovectores. Así las Componentes principales se construyen como combinación lineal de las variables originales estandarizadas con coeficientes determinados por los autovectores.**

$$CP_j = e_{1j} X_1^* + e_{2j} X_2^* + \dots + e_{pj} X_p^*$$

Para realizar un análisis de componentes principales con R seguiremos los siguientes pasos:

Utilizaremos las librerías Factoextra y factominer, que se resumen en los siguientes esquemas.

FactoMineR & factoextra

Analyzing & Visualizing Multivariate Data

FactoMineR

factoextra

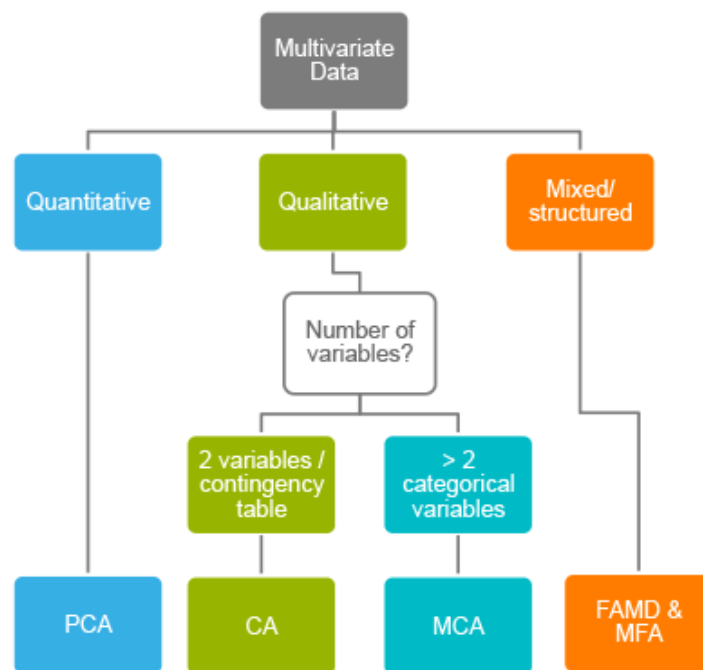


- **Performs** PCA, (M)CA, FAMD, MFA, HCPC & more
- **Provides** the coordinates, the quality of representation and the contribution of individuals & variables
- **Predicts** the results for supplementary individuals & variables
- **Produces** ggplot2-based elegant data visualization and **facilitates** the interpretation
- **Creates** human readable outputs
- **Simplifies** cluster analysis and visualization

Dentro de la librería factoextra se encuentra el análisis de Componentes Principales, PCA.

Principal Component Methods

Methods to Summarize & Visualize Multivariate Data



- PCA: Principal Component Analysis
- (M) CA: (Multiple) Correspondence Analysis
- FAMD: Factor Analysis of Mixed Data
- MFA: Multiple Factor Analysis

Además existen funciones que nos permiten visualizar mediante gráficos el resultado obtenido con el PCA.

| Functions | Description |
|-------------------------------|---|
| fviz_eig (or fviz_eigenvalue) | Visualize eigenvalues. |
| fviz_pca | Graph of PCA results. |
| fviz_cos2 | Visualize element cos2. ¹ |
| fviz_contrib | Visualize element contributions. ² |

También dispone de funciones que nos permiten extraer información del análisis de PCA realizado.

| Functions | Description |
|-----------------|--------------------------------------|
| get_eigenvalue | Access to the dimension eigenvalues. |
| get_pca | Access to PCA outputs. |
| facto_summarize | Summarize the analysis. |

Continuando con el **EJEMPLO** anterior, sobre el fichero con las notas de 11 alumnos que habíamos importado de Excel.

Primero importamos los datos desde Excel y creamos un **Data Frame** con los datos de las variables:

```
# Importar La base de datos de NOTAS
NOTAS<-read_excel("C:/Users/reven/OneDrive/Desktop/Master Big
data/Clases/ACP y Factorial/NOTAS.xlsx")
datos<- as.data.frame(NOTAS)
rownames(datos)<-datos[,1]
notas_dat<-datos[, -1]
```

Es importante crear así el Data Frame para tener la primera columna alfanumérica con el identificador del individuo que más tarde utilizaremos para las representaciones gráficas. Además creamos la matriz notas_dat con solo los valores numéricos de las variables necesaria como entrada para la **función PCA**

Calculamos los descriptivos unidimensionales

```
#Descriptivos
```

```
tabla1<-stat.desc(notas_dat,basic=FALSE)
knitr::kable(tabla1, digits =2,caption = "Estadísticos descriptivos")
```

Estadísticos descriptivos

| | Mats | Francés | Inglés | Física |
|--------------|------|---------|--------|--------|
| median | 5.00 | 5.00 | 6.00 | 4.00 |
| mean | 4.64 | 5.91 | 6.45 | 5.27 |
| SE.mean | 0.79 | 0.55 | 0.45 | 0.75 |
| CI.mean.0.95 | 1.76 | 1.22 | 1.01 | 1.68 |
| var | 6.85 | 3.29 | 2.27 | 6.22 |
| std.dev | 2.62 | 1.81 | 1.51 | 2.49 |
| coef.var | 0.56 | 0.31 | 0.23 | 0.47 |

Recordamos que siempre es conveniente realizar el análisis de Componentes principales sobre la matriz de correlaciones en lugar de sobre la matriz de varianzas-covarianzas.

Para comprobar si las variables están relacionadas de forma lineal entre ellas calculamos la matriz de correlaciones y la guardamos en la matriz R

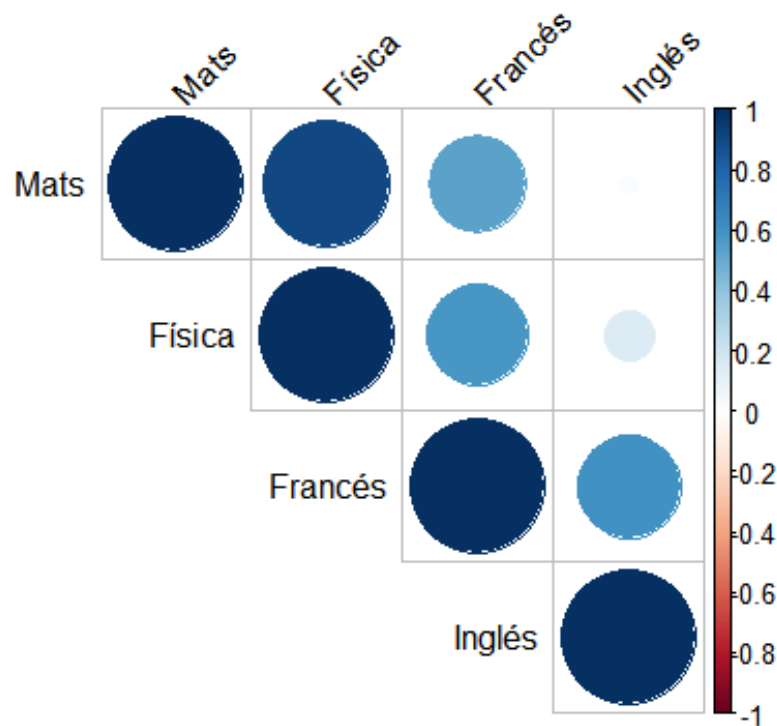
```
#Matriz correlaciones
```

```
R<-cor(notas_dat, method="pearson")
knitr::kable(R, digits =2,caption = "Correlaciones")
```

| | Mats | Francés | Inglés | Física |
|---------|------|---------|--------|--------|
| Mats | 1.00 | 0.54 | 0.02 | 0.91 |
| Francés | 0.54 | 1.00 | 0.60 | 0.58 |
| Inglés | 0.02 | 0.60 | 1.00 | 0.15 |
| Física | 0.91 | 0.58 | 0.15 | 1.00 |

Cuando tenemos muchas variables es muy útil poder analizar las correlaciones mediante una salida gráfica. La siguiente sintaxis representa las correlaciones según su valor en el gráfico anterior.

```
corrplot(R, type="upper", order="hclust",tl.col="black", tl.srt=45)
```



La sintaxis básica del procedimiento de R para realizar un análisis de componentes principales

```
#=====
# Estimamos PCA
#=====
library(FactoMineR)
fit<-PCA(notas_dat,scale.unit=TRUE,ncp=4,graph=TRUE)
```

La opción `scale.unit=TRUE`, significa que calculamos los autovalores de la matriz de correlaciones, porque las variables son estandarizadas

$$X_{ij}^* = \frac{X_{ij} - \bar{X}_j}{\sqrt{S_{jj}}}$$

Recordamos que las variables estandarizadas tienen media cero y varianza 1. Así la suma de las varianzas de estas nuevas variables es igual al número de variables p . `ncp=4` nos indica el número de componentes que se van a calcular, normalmente un número inferior al de variables que tenemos.

El resultado del análisis de componentes principales se guardará en el fichero de **datos fit**. Si queremos ver el contenido de este fichero haremos

```
head(fit)
```

Resultados del PCA guardados en `fit`:

1 "\$eig" [Autovalores](#)

2 "\$var" [Resultados para las variables](#), que son los siguientes:

3 "\$var\$coord": Coordenadas de las variables en las Componentes"

4 "\$var\$cor" : Correlaciones entre las variables y las Componentes

5 "\$var\$cos2" "cosenos al cuadrado de las variables"

6 "\$var\$contrib" "contributions of the variables"

7 "\$ind" [Resultados para los individuos](#), que son los siguientes:

8 "\$ind\$coord" "coord. for the individuals"

9 "\$ind\$cos2" "cos2 for the individuals"

10 "\$ind\$contrib" "contributions of the individuals"

11 "\$call" "summary statistics" : [Medidas estadísticas](#)

12 "\$call\$centre" "mean of the variables"

13 "\$call\$ecart.type" "standard error of the variables"

Veamos cada uno de los elementos desglosados

1 "\$eig" [Autovalores](#)

`$`eig``

| | eigenvalue | percentage of variance | cumulative percentage of variance |
|--------|------------|------------------------|-----------------------------------|
| comp 1 | 2.48576274 | 62.144068 | 62.14407 |
| comp 2 | 1.19280496 | 29.820124 | 91.96419 |
| comp 3 | 0.23856665 | 5.964166 | 97.92836 |
| comp 4 | 0.08286565 | 2.071641 | 100.00000 |

Los autovalores suman p (diagonal de la matriz R), es decir el número de variables.

Los autovectores asociados a dichos autovalores nos dan los coeficientes de la combinación lineal con la que se construyen las componentes.

`svdv`

| | [,1] | [,2] | [,3] | [,4] |
|------|-----------|------------|------------|------------|
| [1,] | 0.5487816 | -0.4158181 | 0.1001327 | 0.7182670 |
| [2,] | 0.5382892 | 0.3426941 | -0.7625266 | -0.1065772 |
| [3,] | 0.2834952 | 0.7829499 | 0.5292340 | 0.1628840 |
| [4,] | 0.5733358 | -0.3108784 | 0.3583825 | -0.6679840 |

Así la Componente 1 será:

$$CP_1 = 0.548X_1^* + 0.538X_2^* + 0.283X_3^* + 0.573X_4^*$$

1.2.3. DETERMINACIÓN DEL NÚMERO DE COMPONENTES PRINCIPALES

Como se indicó con anterioridad uno de los objetivos del Análisis de Componentes Principales consiste en detectar la verdadera dimensión del problema. Esto provoca la necesidad de tomar una decisión ¿cuántas Componentes debemos retener?.

1. El número de Componentes lo determina la proporción de variabilidad a explicar que se considere suficiente, como mínimo el 70%, pero si es posible en torno al 80-90%. En nuestro ejemplo de las notas los dos primeros autovalores explican el 91.96% de la variabilidad.

```
$`eig`
```

| | eigenvalue | percentage of variance | cumulative percentage of variance |
|--------|------------|------------------------|-----------------------------------|
| comp 1 | 2.48576274 | 62.144068 | 62.14407 |
| comp 2 | 1.19280496 | 29.820124 | 91.96419 |
| comp 3 | 0.23856665 | 5.964166 | 97.92836 |
| comp 4 | 0.08286565 | 2.071641 | 100.00000 |

2. **Scree Plot.** Representación gráfica de las proporciones de variabilidad explicada por cada Componente. Cuando los puntos de la gráfica se nivelan, el porcentaje de variabilidad explicado por estas componentes es próximo a cero como para poder despreciarlos. Se puede trazar una recta que aglutine en su entorno a los autovalores más pequeños y todos los que queden por encima corresponderían a las Componentes Principales retenidas.

```
# Scree plot
```

```
fviz_eig(fit, geom="line")+ theme_grey()
```

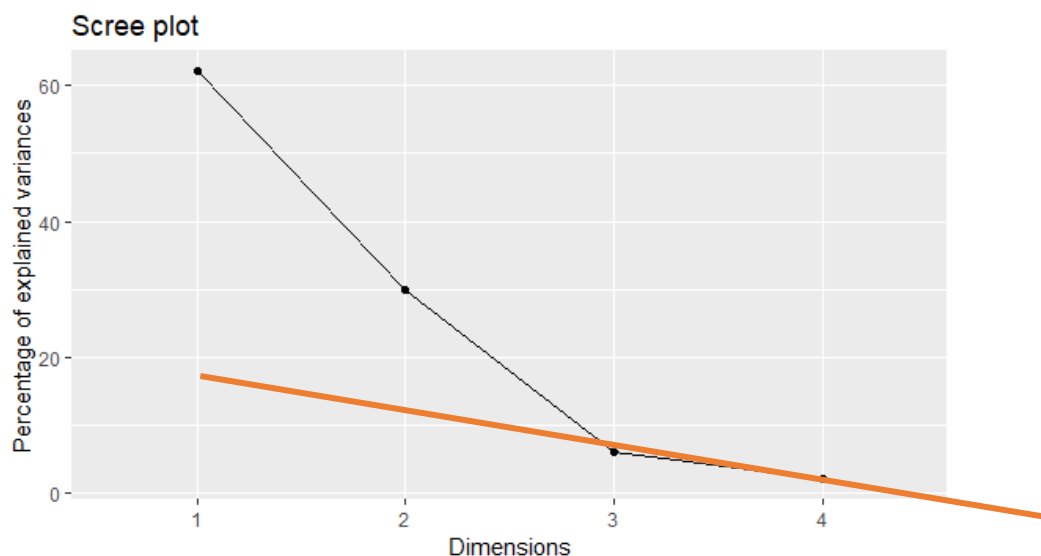
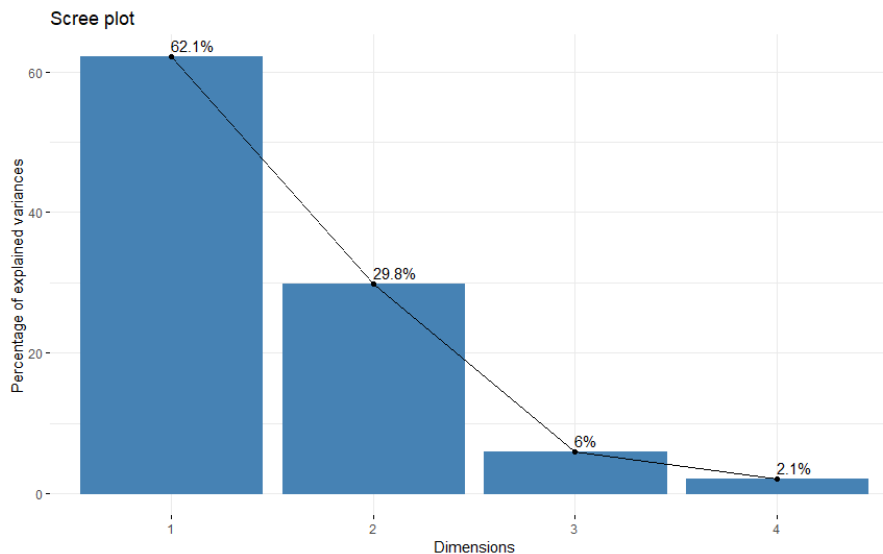


Figura 2.2. Ejemplo de Scree Plot

Parece adecuado retener las 2 primeras Componentes Principales.

Otra version del Scree Plot añade una barra a cada porcentaje.

```
# Scree plot
fviz_eig(fit,addlabels=TRUE)
```



Por lo que decidiríamos quedarnos con 2 de las cuatro componentes principales que explicarían el 91,96% de la varianza de las variables iniciales. Los dos autovalores asociados a las componentes seleccionadas valen 2.485 y 1.192, y en términos de porcentaje explican el 62.14% de la variabilidad y el 29.82%. Esto se lo indicamos en la opción `ncp=2`.

```
#=====
# Decidimos hacer el análisis con 2 componentes
#=====
fit<-PCA(notas_dat,scale.unit=TRUE, ncp=2, graph=TRUE)
```

1.2.3 RELACIÓN ENTRE LAS COMPONENTES PRINCIPALES Y LAS VARIABLES

En estadística, la relación entre variables numéricas se mide mediante la covarianza entre dichas variables. Las covarianzas entre cada componente principal y las variables vienen dadas por el producto de las coordenadas del vector propio que define el componente por el valor propio:

$$\text{Cov}(CP_i, X_j) = \lambda_i e_{ij}$$

La correlación entre una componente principal y una variable X es proporcional al coeficiente de esa variable en la definición del componente, y el coeficiente de

proporcionalidad es el cociente entre la desviación típica del componente y la desviación típica de la variable.

$$\text{Corr}(CP_i, X_j) = \frac{\text{Cov}(CP_i, X_j)}{\sqrt{V(CP_i)V(X_j)}} = \frac{\lambda_i e_{ij}}{\sqrt{\lambda_i s_j^2}} = e_{ij} \frac{\sqrt{\lambda_i}}{s_j}$$

```
#Guardamos los estadísticos asociados a las variables en el objeto var
var<-get_pca_var(fit)
knitr::kable(var$cor, digits =3,caption = "Correlaciones de la CP con
las variables")
```

Correlaciones de la CP con las variables

| | Dim.1 | Dim.2 |
|---------|-------|--------|
| Mats | 0.865 | -0.454 |
| Francés | 0.849 | 0.374 |
| Inglés | 0.447 | 0.855 |
| Física | 0.904 | -0.340 |

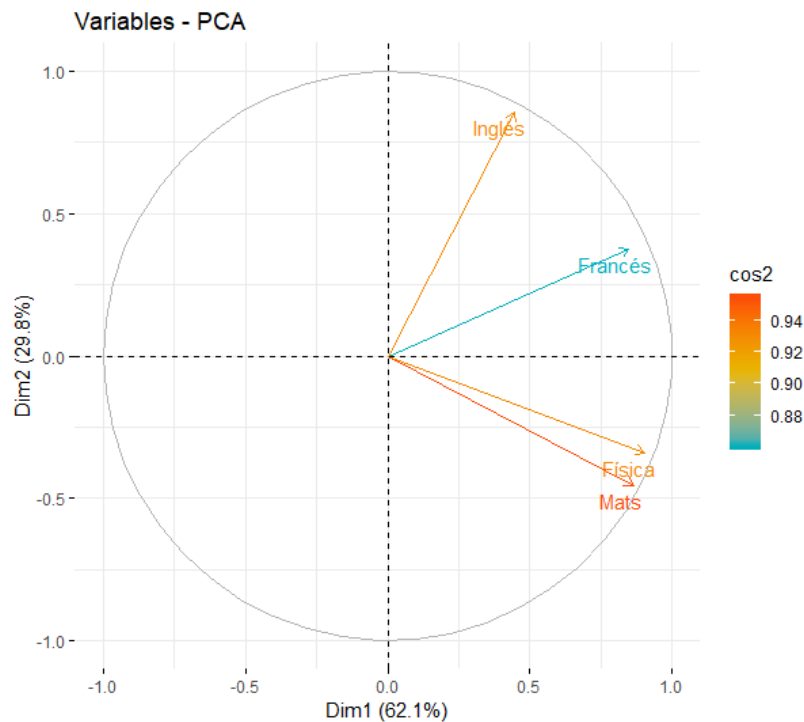
Por ejemplo la correlación entre la variable X_1 y la componente 1 es:

$$\text{Corr}(CP_1, X_1) = 0.548 \frac{\sqrt{2.486}}{\sqrt{6.854}} = 0.865$$

Observemos que la primera componente tiene una correlación grande y positiva con las variables Física, Matemáticas y Francés. Por esto estas variables están “representadas” de forma directa en la CP1. Mientras que la Componente 2 tiene correlación grande solo con Inglés (0.85) por lo que esta componente representa fundamentalmente a esta variable aunque también hay que tener en cuenta las correlaciones con las otras variables y el hecho de que tenga correlación negativa con Matemáticas y Física.

La siguiente representación gráfica nos puede ayudar a entender como están recogidas nuestras variables iniciales en las nuevas componentes, puesto que representa el coeficiente de correlación entre las variables y las nuevas componentes.

```
#=====
# Representación gráfica variables
#=====
fviz_pca_var(fit, col.var="cos2",gradient.cols = c("#00AFBB", "#E7B800",
"#FC4E07"),repel = TRUE)
```

Por ejemplo la primera componente recoge sobre todo los valores de Física y Matemáticas y algo menos Francés, por lo que podríamos identificar dicha componente como la que recoge las asignaturas de ciencias.

Los **cosenos al cuadrado** son **las correlaciones al cuadrado** que **expresan la proporción de la varianza de cada variable que es explicada por cada componente**

```
knitr::kable(var$cos2, digits = 2, caption = "Cosenos al cuadrado")
```

Cosenos al cuadrado

| | Dim.1 | Dim.2 |
|---------|-------|-------|
| Mats | 0.75 | 0.21 |
| Francés | 0.72 | 0.14 |
| Inglés | 0.20 | 0.73 |
| Física | 0.82 | 0.12 |

$$\text{Cos}^2(CP_1, X_1) = \text{Corr}(CP_1, X_1)^2 = 0.865^2 = 0.75$$

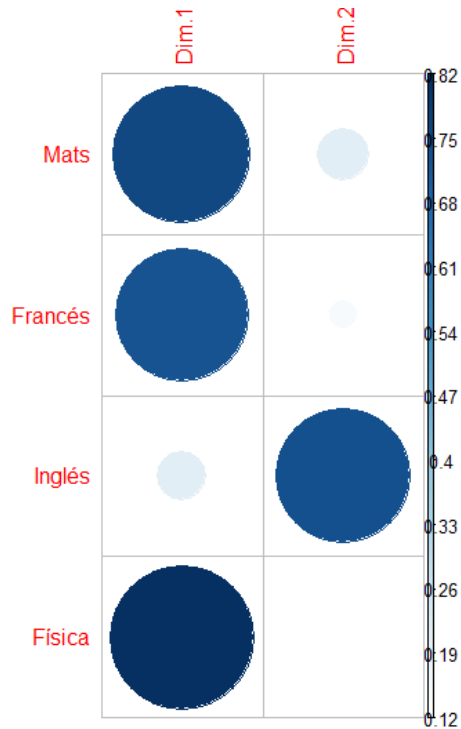
$$\text{Cos}^2(CP_2, X_1) = \text{Corr}(CP_2, X_1)^2 = -0.454^2 = 0.21$$

Luego la proporción de variabilidad de la variable Matemáticas que es explicada por las dos Componentes es

$$\text{Cos}^2(X_1) = 0.748 + 0.206 = 0.954$$

Gráficamente se muestra que la varianza de las variables Matemáticas, Física y Francés es explicada por la Componente 1 mientras que la Componente 2 explica Inglés

```
corrplot(var$cos2, is.cor=FALSE)
```



Mostramos el porcentaje de la varianza de las variables que es explicada por las dos Componentes en total

```
$var$contrib
```

| | Dim.1 | Dim.2 |
|---------|-----------|-----------|
| Mats | 30.116124 | 17.290471 |
| Francés | 28.975525 | 11.743927 |
| Inglés | 8.036952 | 61.301061 |
| Física | 32.871399 | 9.664541 |

1.2.4. COORDENADAS DE LOS INDIVIDUOS SOBRE LOS NUEVOS EJES.

Una vez seleccionadas las direcciones de los ejes que explican la máxima variabilidad, pasamos a situar las observaciones sobre los nuevos ejes luego tenemos que calcular sus nuevas coordenadas. Las coordenadas de los n puntos individuos sobre el eje correspondiente a \mathbf{e}_α son las n componentes del vector

$$\Psi_\alpha = \mathbf{X} \mathbf{e}_\alpha,$$

OBSERVACIÓN: A la hora de calcular las coordenadas, en realidad se emplea la matriz \mathbf{X} , en la que $x_{ij} - \bar{x}_j$, si el Análisis se realiza a partir de la Matriz de Varianzas-Covarianzas. Si el A.C.P. parte de la matriz de Correlaciones, para obtener las coordenadas de los individuos la matriz \mathbf{X} es aquella en que $\frac{x_{ij} - \bar{x}_j}{s_j}$.

Como la matriz de datos de partida está centrada (hemos restado su media a cada componente), es sencillo comprobar que la suma de las coordenadas de los n individuos sobre cada eje es nula.

$$\sum_{i=1}^n \Psi_{\alpha i} = \sum_{i=1}^n \sum_{j=1}^p \frac{(x_{ij} - \bar{x}_j)}{s_j} e_{\alpha j} = \sum_{j=1}^p e_{\alpha j} \sum_{i=1}^n \frac{(x_{ij} - \bar{x}_j)}{s_j} = 0 \quad \forall \alpha$$

Por lo que la media de las n coordenadas es cero y su varianza será:

$$\Psi'_\alpha \Psi_\alpha = \mathbf{e}_\alpha' \mathbf{X}' \mathbf{X} \mathbf{e}_\alpha = \lambda_\alpha$$

Continuando con el **EJEMPLO** de las notas de los alumnos en las cuatro asignaturas, tenemos en el siguiente cuadro, las coordenadas de los 11 alumnos en las Componentes Principales.

| Dim.1 | Dim.2 |
|--------|--------|
| -2.228 | -0.268 |
| -0.994 | -1.411 |
| 2.250 | 1.875 |
| -0.728 | 0.209 |
| -2.294 | 0.372 |
| 1.716 | 0.128 |
| 0.655 | -0.683 |
| -0.645 | 1.101 |
| 0.827 | -1.889 |
| -0.865 | 1.267 |

| | |
|-------|--------|
| 2.305 | -0.701 |
|-------|--------|

Observemos, que la coordenada del primer alumno en la Primera Componente, viene dada por:

$$y_{11} = \frac{(1-4.63)}{2.61} \cdot 0.548 + \frac{(4-5.9)}{1.81} \cdot 0.538 + \frac{(5-6.45)}{1.51} \cdot 0.283 + \frac{(3-5.27)}{2.49} \cdot 0.573 = -2.228$$

Además, la media de las Componentes vale 0. Se puede comprobar que:

$$\frac{\sum_{i=1}^{11} cp_{i1}}{11} = 0$$

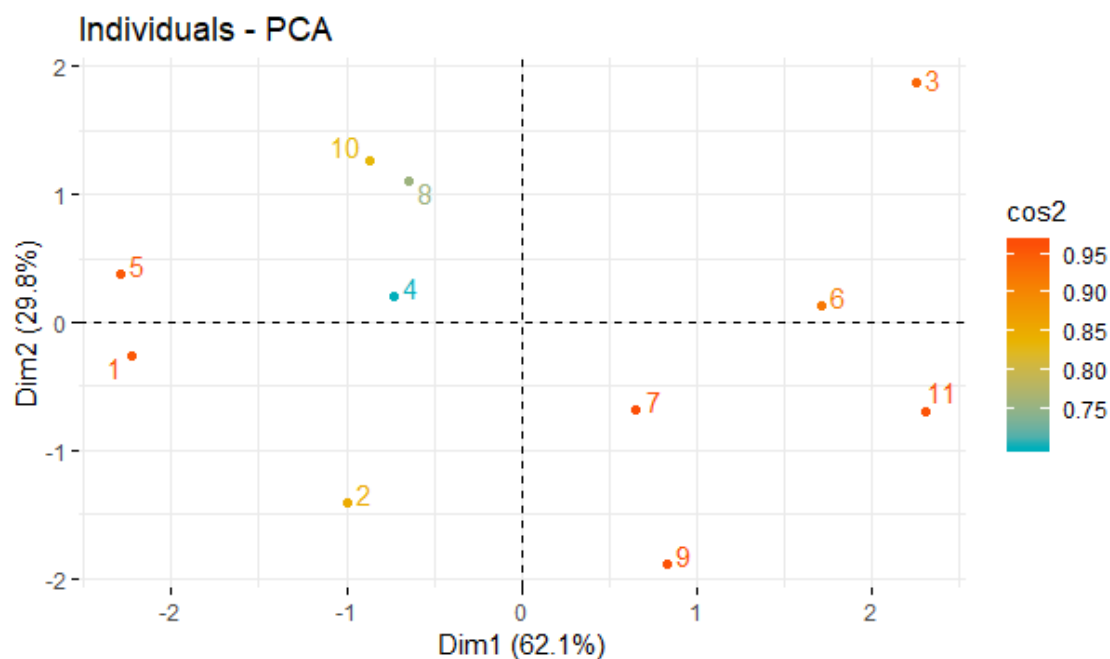
$$V(cp_1) = \frac{\sum_{i=1}^{11} cp_{i1}^2}{10} = \lambda_1 = 2.485.$$

Análogamente, para el resto de variables Ψ_2 , Ψ_3 y Ψ_4 .

En R además podemos representar los individuos en el plano de componentes, porque dichas representaciones pueden ofrecernos información adicional sobre el comportamiento de los individuos.

Representación gráfica de los individuos

```
fviz_pca_ind(fit, col.ind = "cos2",
             gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),
             repel = TRUE)
```



REPRESENTACIÓN CONJUNTA: BILOT

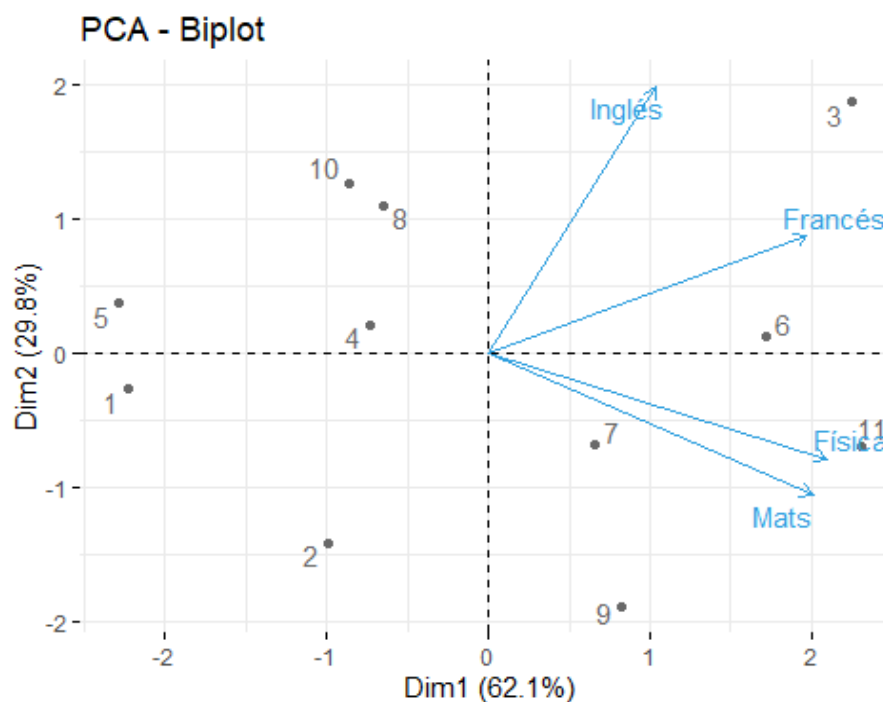
Los gráficos Biplot nos permiten representar de forma conjunta las variables y los individuos en el plano de las Componentes principales.

La construcción de este gráfico está basada en la descomposición en valores singulares (SVD) de una matriz de datos que tiene n filas (observaciones) y p columnas (variables) es $X = U L V^T$ donde U y V es una matriz

- ✓ **Cuidado:** Para representar de forma conjunta las variables y los individuos, la longitud de las flecha se ha modificado (no son las correlaciones) pero se mantienen los ángulos
- ✓ Los individuos próximos a una variable tienen valor alto de esa variable
- ✓ Los individuos en la misma dirección pero sentido opuesto tendrán un valor bajo de dicha variables.
- ✓ Ángulos agudos entre variables indican correlación positiva y cuanto menor ángulo mayor correlación. Si el ángulo es 90° las variables están incorreladas

Representación conjunta de los individuos y las variables en los planos de las CP

```
fviz_pca_biplot(fit, repel = TRUE,
  col.var = "#2E9FDF", # Variables color
  col.ind = "#696969") # Individuals color
```



1.2.5. COMENTARIOS FINALES AL ANÁLISIS DE COMPONENTES PRINCIPALES.

- El Análisis de Componentes Principales sobre Σ es apropiado en aquellos casos en donde todas las variables están medidas por lo menos en unidades comparables teniendo varianzas semejantes. Si las variables no se miden en las mismas unidades cualquier cambio en la escala de medición en alguna de las variables tendrá efecto sobre las Componentes Principales. Este cambio de escala puede invertir los papeles de las variables importantes.
- Se puede dar la circunstancia de que alguno de los autovalores (bien de la matriz de Varianzas-Covarianzas, bien de la de Correlaciones) sea 0. Esto significa que algunas de las variables originales serán linealmente dependientes de las otras. Así si retenemos m Componentes Principales, en esencia estamos suponiendo que hay $p-m$ autovalores iguales a 0 (realmente no significativamente diferentes de 0), y sus autovectores correspondientes definen $p-m$ restricciones lineales, linealmente independientes, sobre las variables observables. Las ecuaciones derivadas de las restricciones lineales (relaciones estructurales) proporcionan información acerca de cómo están relacionadas las variables entre sí. Nos ayudan a detectar relaciones de multicolinealidad entre un conjunto de variables predictoras en los problemas de regresión.
- Aunque cuando se usa el Análisis de Componentes Principales en la práctica el mayor interés suele centrarse en las Componentes con mayor varianza, hay ocasiones donde el examen de las dos últimas Componentes elegidas puede ser de utilidad. Por ejemplo estas Componentes pueden identificar observaciones atípicas que pueden provocar el incremento ficticio de la dimensión. Estos valores atípicos se pueden obtener representando las observaciones en el espacio generado por las dos últimas Componentes y observando la existencia de puntos aislados, los cuales serán identificados como atípicos.
- Por último, hay que tener en cuenta que las Componentes Principales no tienen que ser “variables interpretables”. El objetivo verdadero del A. C. P. no es conseguir variables interpretables, sino detectar la auténtica dimensión del problema.

1.2.5 SISTEMÁTICA DEL ANÁLISIS DE COMPONENTES PRINCIPALES.

En este apartado se presenta el esquema de los pasos a seguir para realizar un Análisis de Componentes Principales. Estos pasos son los siguientes:

- 1) Estudiar la matriz de Correlaciones para comprobar que existen muchas variables con correlaciones grandes, próximas a 1 en valor absoluto.
- 2) Obtención de autovalores y autovectores.
- 3) Determinación del número de Componentes a retener. Se pueden tener los siguientes criterios:

Como se diagonaliza la matriz de Correlaciones:

- autovalores mayores que 1, utilizar la gráfica de Catell, proporción de variabilidad explicada > 0.7 .
- 4) Interpretación, si procede, de las Componentes a través de los coeficientes de los autovectores respectivos y de las correlaciones con las variables
 - 5) Obtención de los valores que toman los individuos en cada Componente.
 - 6) Representar las variables en el espacio de las Componentes. Esto permitirá crear asociaciones entre variables.
 - 7) Representar los individuos en el espacio de las Componentes. Mediante dichas representaciones, se detectarán grupos de individuos, observaciones atípicas, etc.

Bibliografía.

- An Introduction of Applied Multivariate Analysis with R. Everitt B, Hothorn T. Ed Wiley. 2011. [Libro completo con explicaciones teóricas y ejemplos resueltos en R aunque con librerías básicas. Además aparecen otros temas relacionados como el análisis factorial](#)
- Principal Components Analysis. Jolliffe, I.T. Ed Springer 2002 . [Libro completo con explicaciones teóricas muy amplias y algunos ejemplos. El más citado por artículos de investigación.](#)
- Nuevos Métodos de Análisis Multivariante. Cuadras C.M. 2014 . [Libro completo con explicaciones teóricas y otras técnicas multivariantes](#)
- Practical Guide to Principal Component Methods in R. A. Kassambara. Ed STHDA.com. 2017 . [Libro completo que explica las librerías factominer y factoextra mediante ejemplos. Sin explicaciones teóricas](#)
- Package 'factoextra' . [Explicación del funcionamiento de la librería y la sintaxis detallada](#)
- Package 'factominer' . [Explicación del funcionamiento de la librería y la sintaxis detallada](#)
- <http://www.sthda.com/english/>