

# ACP

DANIEL MARTÍN GARCÍA (versión de Pablo Flores Vidal)

Dic 2024

## 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, disponer de un amplio número de variables puede ocasionar diferentes problemas. En este contexto complejo y multidimensional, surge la necesidad de técnicas que puedan simplificar y resumir la información de manera eficiente. Aquí es donde el Análisis de Componentes Principales (ACP) juega un papel fundamental.

El ACP es una técnica estadística (fuertemente anclada en las matemáticas) que busca reducir la dimensionalidad de un conjunto de datos, manteniendo al mismo tiempo la mayor cantidad de información ( $\approx$  variabilidad) posible. En un mundo donde los datos masivos son comunes, la posibilidad de simplificar la información sin perder sus características esenciales se convierte en una herramienta poderosa.

Uno de los principales problemas que el ACP ayuda a resolver es la multicolinealidad, es decir, la correlación alta entre diversas variables. Esto puede complicar la interpretación de los datos y afectar la eficacia de otros métodos de análisis. Transformando los datos en componentes principales, el ACP elimina este problema, facilitando la interpretación y el análisis.

Además, el ACP permite una visualización más clara de los datos. Al reducir la cantidad de variables a unos pocos componentes principales, se pueden representar gráficamente los datos en un espacio de baja dimensión. Esto es invaluable para la comprensión y el descubrimiento de patrones y relaciones en conjuntos de datos complejos.

La versatilidad del ACP es otra de sus fortalezas. No se limita a un campo de estudio específico, sino que se aplica en a cualquier disciplina donde quepa la ciencia de Datos, como la biología, economía, psicología, y demás. Su capacidad para detectar estructuras en los datos ha llevado a avances significativos en campos tan diversos como la genómica o las finanzas.

El ACP también contribuye a la mejora de la calidad de los datos. Identificar y eliminar las variables redundantes mejora la precisión de los modelos, permitiendo un análisis más ajustado a la realidad. Además, a mayor dimensionalidad de los datos, mayor es la capacidad computacional que se precisa. Por lo tanto, reducir la dimensionalidad mediante ACP no solo optimiza los recursos computacionales sino que también facilita el manejo y procesamiento de grandes conjuntos de datos.

En conclusión, el Análisis de Componentes Principales es más que una simple técnica de reducción de dimensionalidad. Es una herramienta multifacética que responde a los desafíos de la era del Big Data, permitiendo una interpretación de los datos más sencilla, eliminando la multicolinealidad, facilitando la visualización, y contribuyendo a la calidad y eficiencia de los análisis. Su aplicabilidad en diversos campos y su capacidad para desentrañar patrones ocultos en datos complejos subrayan su importancia y justifican su uso continuo en el análisis de datos moderno.

## 2 Análisis de Componentes Principales

### 2.1 ¿En qué consiste el ACP?

El Análisis de Componentes Principales es una técnica estadística que se aplica a un conjunto de datos compuesto por variables que están potencialmente correlacionadas entre sí. El objetivo central del ACP es transformar estas variables originales en un nuevo conjunto de variables, que son linealmente independientes y, por lo tanto, no correlacionadas entre sí, si la estructura de los datos lo permite. Estas nuevas variables se denominan componentes principales.

La transformación se lleva a cabo mediante una combinación lineal de las variables originales, donde se seleccionan los coeficientes de manera que el primer componente principal capture la mayor parte de la variabilidad presente en los datos originales. Esto significa que esta primera componente resume, en la medida de lo posible, la información contenida en las variables originales.

El segundo componente principal se calcula de manera similar, pero está restringido a ser ortogonal al primer componente, es decir, no correlacionado con él. Este segundo componente explica la siguiente mayor parte de la variabilidad que no ha sido capturada por el primer componente.

Este proceso continúa de manera secuencial, con cada componente principal subsiguiente explicando una proporción decreciente de la variabilidad restante, y siempre siendo ortogonal a todos los componentes principales anteriores. La transformación se realiza de tal manera que se puede representar la estructura compleja y multidimensional de los datos originales en un número reducido de dimensiones, sin perder significativamente la información esencial.

El resultado es una representación simplificada y estructurada de los datos originales, donde los componentes principales capturan y resumen la información crítica, facilitando así la interpretación y el análisis posterior. La eficacia de esta transformación depende de la naturaleza de las relaciones entre las variables originales y de la variabilidad inherente en los datos, y el cumplimiento de los requisitos que mencionaremos a continuación garantiza una aplicación más precisa y significativa del ACP.

A continuación, en la Figura 1, se muestra (a la izquierda) una nube de puntos que representa la relación entre las variables originales  $(X, Y)$ , altamente correlacionadas. Por otro lado, a la derecha se presenta el resultado del ACP: dos variables ortogonales que han sido derivadas de las variables originales  $(X, Y)$ . Aunque la dimensionalidad no se ha reducido (sigue siendo un espacio bidimensional), lo notable aquí es que la correlación entre las variables originales ha sido eliminada. Las nuevas variables, ahora ortogonales entre sí, ofrecen una vista descorrelacionada de los datos,

preservando la estructura subyacente pero en un marco donde las variables no están linealmente relacionadas.

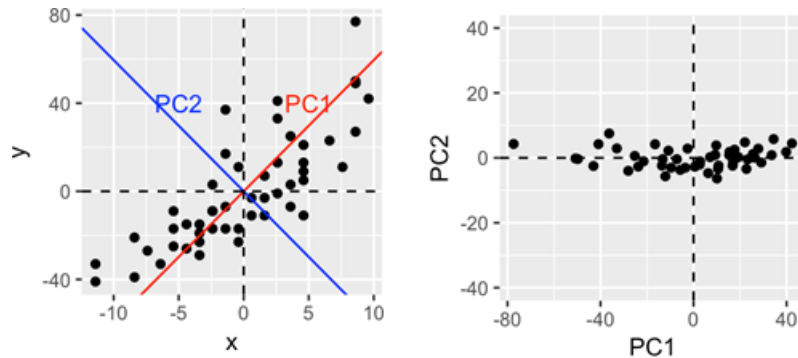


Figure 1: Representación de las Componentes.

## 2.2 Diferencias con el Análisis Factorial

El Análisis de Componentes Principales (ACP) y el Análisis Factorial (AF) son técnicas de reducción de la dimensionalidad, pero tienen propósitos y enfoques distintos, lo que puede generar confusión.

### 1. Objetivo principal

ACP: Busca transformar un conjunto de variables originales en un nuevo conjunto de variables no correlacionadas llamadas componentes principales. Estas componentes explican la máxima varianza posible del conjunto de datos. Su meta es descriptiva, resumir los datos de forma más manejable.

AF: Intenta identificar factores latentes (no observables) que explican las correlaciones entre las variables originales. Su meta es explicar las relaciones subyacentes en los datos, idealmente en términos de constructos teóricos.

### 2. Modelo matemático

ACP: No asume un modelo subyacente. Es puramente algorítmico y se basa en la diagonalización de la matriz de covarianzas o correlaciones. Se busca una combinación lineal de las variables que maximice la varianza explicada.

AF: Utiliza un modelo subyacente que supone que las variables observadas están influenciadas por factores comunes (latentes) y factores únicos (error). Se basa en la ecuación:

$$\mathbf{X} = \Lambda \mathbf{F} + \epsilon$$

donde:  $\mathbf{X}$ : Variables observadas,  $\Lambda$ : Matriz de cargas factoriales (relación entre factores y variables),  $\mathbf{F}$ : Factores latentes,  $\epsilon$ : Error específico o varianza única.

### 3. Varianza explicada

ACP: Explica toda la varianza de las variables originales (tanto común como única). Por esto, los componentes principales suelen ser útiles para comprensión y visualización de datos.

AF: Solo intenta explicar la varianza común entre las variables, ignorando la varianza única. Esto hace que sea más útil para identificar patrones o constructos subyacentes.

#### 4. Método de extracción

ACP: Se basa en la descomposición en valores propios (eigenvalues) de la matriz de correlación o covarianza.

AF: Emplea métodos como máxima verosimilitud, mínimos cuadrados ponderados o métodos iterativos para estimar las cargas factoriales y los factores comunes.

#### 5. Interpretación

ACP: Los componentes principales son combinaciones lineales de las variables originales, y no tienen una interpretación teórica más allá de representar la dirección de mayor varianza en los datos.

AF: Los factores latentes tienen una interpretación más teórica y suelen representar constructos subyacentes (por ejemplo, "inteligencia general" en psicometría o "satisfacción del cliente" en estudios de mercado).

#### Resumen:

ACP: Ideal para resumir datos y reducir la dimensionalidad, explicando varianza sin enfocarse en factores latentes.

AF: Diseñado para descubrir relaciones subyacentes y latentes en los datos, basado en un modelo teórico.

### 2.3 ¿Qué requisitos precisan los datos originales?

Para que un conjunto de datos sea susceptible de análisis mediante el ACP, es vital cumplir con ciertos requisitos:

- **Escala de las Variables:** Las variables de entrada deben estar en la misma escala o ser estandarizadas, para que cada una contribuya de manera equitativa al análisis.
- **Ausencia de Valores Atípicos:** Los datos extremos pueden dificultar el análisis, y su manejo adecuado es fundamental.
- **Datos perdidos:** La ausencia de datos puede complicar el análisis, por lo que pueden ser necesarias técnicas de imputación.
- **Distribución de las Variables Originales:** Aunque no es necesario que las variables originales tengan una distribución normal, esto puede mejorar la interpretabilidad de los resultados.

**Nota:** Es fundamental destacar que el ACP presupone correlaciones entre las variables originales. Aunque técnicamente es posible llevar a cabo el análisis en ausencia de correlaciones, el resultado sería trivial y no proporcionaría ninguna comprensión nueva o valiosa acerca de la estructura de los datos. En otras palabras, sin correlaciones entre las variables, el ACP no obtendría ningún resultado útil.

### 2.4 Pasos del Proceso

Supongamos que tenemos un conjunto de datos con  $n$  observaciones y  $p$  variables donde  $X_{ij}$ , es el valor de la observación  $i$ -ésima, en la variable  $j$ -ésima, y queremos reducirlo a  $k$  componentes principales.

1. **Estandarización de los Datos:** Si las escalas de las variables originales difieren, deben ser estandarizadas:

$$Z_{ij} = \frac{X_{ij} - \bar{X}_j}{s_j}, i \in n, j \in p,$$

donde  $\bar{X}_j$  es la media de la variable  $j$ -ésima, y  $s_j$  es su desviación estándar.

2. **Cálculo de la Matriz de Covarianzas  $\Sigma$ :** Si los datos originales están en la misma escala, (poco habitual), podremos utilizar la matriz de covarianzas, la cual, es una matriz cuadrada  $p \times p$  que contiene las covarianzas entre cada par de variables originales.

La entrada en la  $i$ -ésima fila y  $j$ -ésima columna de la matriz de covarianzas  $\Sigma$  es la covarianza entre  $X_i$  y  $X_j$ , denotada como  $\text{Cov}(X_i, X_j)$ . La covarianza entre dos variables se calcula mediante la fórmula:

$$\text{Cov}(X_i, X_j) = \frac{1}{n-1} \sum_{k=1}^n (X_{i,k} - \bar{X}_i)(X_{j,k} - \bar{X}_j)$$

Donde:

- $n$  es el número de observaciones.
- $X_{i,k}$  es la  $k$ -ésima observación de la variable  $X_i$ .
- $\bar{X}_i$  es la media de la variable  $X_i$ .

La matriz de covarianzas  $\Sigma$  es una matriz simétrica de  $p \times p$ , donde cada elemento  $\Sigma_{ij}$  representa la covarianza entre  $X_i$  y  $X_j$ .

En forma matricial, la matriz de covarianza  $\Sigma$  se expresa como:

$$\Sigma = \begin{bmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Cov}(X_p, X_p) \end{bmatrix}$$

Es importante señalar que al trabajar con datos muestrales, se utiliza el factor  $\frac{1}{n-1}$  para calcular la matriz de covarianza muestral. Si se tiene toda la población, se usaría  $\frac{1}{n}$  en su lugar.

3. **Cálculo de la Matriz de Correlaciones:** Si las variables originales no están medidas en la misma escala, entonces, utilizaremos la matriz de correlaciones, la cual, se calcula a partir de los datos estandarizados y vuelve a ser una matriz cuadrada  $p \times p$  que contiene las correlaciones entre cada par de variables estandarizadas. Se calcula como:

$$R = \begin{bmatrix} \rho(Z_1, Z_1) & \rho(Z_1, Z_2) & \dots & \rho(Z_1, Z_p) \\ \rho(Z_2, Z_1) & \rho(Z_2, Z_2) & \dots & \rho(Z_2, Z_p) \\ \vdots & \vdots & \ddots & \vdots \\ \rho(Z_p, Z_1) & \rho(Z_p, Z_2) & \dots & \rho(Z_p, Z_p) \end{bmatrix}$$

En esta notación:

- $\rho(Z_i, Z_j)$  representa la correlación entre  $Z_i$  y  $Z_j$ .
- $R$  es la matriz de correlaciones.
- La matriz es simétrica y cada elemento  $R_{ij}$  representa la correlación entre  $Z_i$  y  $Z_j$ .

Ten en cuenta que el coeficiente de correlación se indica comúnmente con la letra griega  $\rho$  (rho).

4. **Determinación de los Valores Propios:** A partir de la matriz de covarianza (o de correlaciones), se calculan los valores propios, resolviendo la ecuación característica  $\det(\Sigma - \lambda I) = 0$ , donde  $\Sigma$  es la matriz de covarianzas o la matriz de correlaciones,  $I$  es la matriz identidad, matriz de dimensión  $p \times p$  con unos en la diagonal y ceros en todas las demás posiciones; y  $\lambda$  es el vector que contiene el conjunto de valores propios
5. **Determinación de los  $k$  primeros Valores Propios:** Los autovalores se ordenan en orden descendente, y se eligen los primeros  $k$ , que representarán los componentes principales. Para elegir estos  $k$  primeros, hay que tener en cuenta que el valor de cada uno de los autovalores  $\lambda_i$  para  $i = 1, \dots, p$ , es proporcional a la variabilidad explicada por las componentes asociadas. Por ejemplo, estando los autovalores ordenados, tal que  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ , la componente asociada al primer autovalor  $\lambda_1$ , explica un  $\frac{\lambda_1}{\sum_{i=1}^p \lambda_i} * 100$  por ciento de la variabilidad total de las variables originales.
6. **Determinación de los vectores propios asociados:** Una vez que tenemos los valores propios  $\lambda$ , podemos encontrar los autovectores asociados, resolviendo la siguiente ecuación lineal para cada  $\lambda$ :  $(\Sigma - \lambda I)w = 0$ . Donde  $w$  es el vector propio que estamos buscando.
7. **Construcción de la Matriz de Transformación:** Los primeros  $k$  vectores propios forman la matriz de transformación. Es decir,  $W = [w_1, w_2, \dots, w_k]$ , donde  $w_i$ , para  $i = 1, \dots, k$  son los primeros  $k$  vectores propios.
8. **Cálculo de los Componentes Principales:** Una vez que se ha obtenido la matriz de transformación  $W$ , las  $s$  componentes principales  $CP$  se pueden calcular como la proyección de los datos originales  $X$  o los datos estandarizados  $Z$  en el espacio de las nuevas dimensiones (componentes principales). Esto se realiza mediante una multiplicación matricial entre los datos y la matriz de transformación:

$$CP = XW \quad \text{ó} \quad CP = ZW$$

donde  $CP$  es una matriz de dimensiones  $n \times k$  que contiene los componentes principales de los datos. Cada columna de  $CP$  representa un componente principal, y cada fila de  $CP$  representa los componentes principales para una observación individual en el conjunto de datos original. Este paso produce las coordenadas de las observaciones en el nuevo espacio definido por los componentes principales.

**Nota:** Incluso aunque las variables se encuentren en la misma escala, siempre es aconsejable estandarizar los datos, por los siguientes motivos:

- **Influencia de Variables:** Si no estandarizas, las variables con mayor variabilidad tendrán más peso en la formación de los componentes principales. Esto puede ser problemático si la variabilidad es alta debido a ruido o errores de medición, y no porque la variable sea particularmente importante para el fenómeno estudiado.
- **Comparabilidad:** Si planeas comparar los resultados de tu ACP con otros estudios que también utilizan ACP pero en diferentes escalas, la estandarización podría ser esencial para hacer que los resultados sean comparables.

## 2.5 Selección del Número de Componentes

Una de las decisiones fundamentales en el Análisis de Componentes Principales es la selección del número de componentes principales, denotado como  $k$ , que se deben retener. Las componentes principales resultantes son linealmente independientes y ortogonales entre sí, y proporcionan un resumen compacto y eficiente de la variabilidad en los datos originales.

La elección de  $k$  es esencial para garantizar que el ACP capture una cantidad significativa de la variabilidad presente en los datos originales sin incluir componentes innecesarios. Una técnica común para determinar  $k$  es utilizar el *método del codo*, que se basa en el gráfico de la varianza explicada.

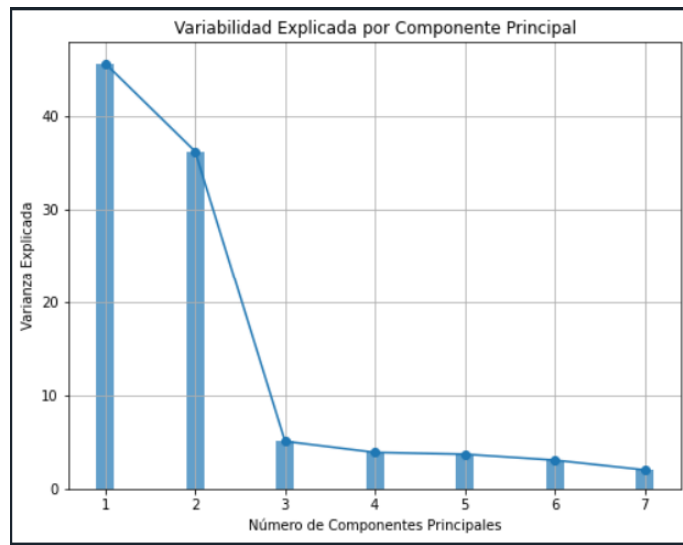


Figure 2: Representación de la variabilidad explicada

En este gráfico, el eje  $x$  representa el número de componentes y el eje  $y$  representa la proporción de varianza explicada. Buscamos el punto en el que agregar componentes adicionales ya no proporciona una ganancia significativa en la varianza explicada. Este punto se asemeja a la forma de un 'codo' en la curva, y marca el número óptimo de componentes a retener. En otras palabras, seleccionamos  $k$  de manera que capturemos la mayor parte de la variabilidad en los datos originales sin introducir redundancias innecesarias.

Este proceso de selección de  $k$  es crucial para obtener un modelo de ACP efectivo que reduzca la dimensionalidad de los datos de manera informativa y eficiente.

## 2.6 Relación entre las Componentes Principales y las variables

En el contexto del Análisis de Componentes Principales, es fundamental comprender la relación entre las componentes principales y las variables originales. Esta relación se basa en los siguientes conceptos clave:

- $\lambda_i$ : Representa el autovalor asociado a la  $i$ -ésima componente principal ( $CP_i$ ) obtenido durante el ACP. Los autovalores indican cuánta varianza (variabilidad) se explica por cada  $CP_i$ . Un autovalor alto ( $\lambda_i$  grande) implica que  $CP_i$  retiene una gran cantidad de información o variabilidad de los datos originales, mientras que un autovalor bajo ( $\lambda_i$  pequeño) indica que  $CP_i$  retiene menos información.
- $e_{ij}$ : Representa el coeficiente ubicado en la posición  $ij$  del vector propio asociado a la  $CP_i$ . Este coeficiente indica la contribución relativa de la variable original  $X_j$  en la construcción de la  $CP_i$ . En otras palabras,  $e_{ij}$  muestra cuánto peso tiene la variable  $X_j$  en la definición de los valores de  $CP_i$ .
- $Cov(CP_i, X_j)$ . La covarianza entre la  $i$ -ésima componente principal ( $CP_i$ ) y la  $j$ -ésima variable original ( $X_j$ ) se calcula como  $\lambda_i \cdot e_{ij}$ . Esta covarianza nos proporciona información sobre la relación lineal entre  $CP_i$  y  $X_j$ . Si el valor resultante es positivo, significa que ambas variables tienden a aumentar o disminuir juntas (correlación positiva). Si es negativo, indica que una variable tiende a aumentar cuando la otra disminuye (correlación negativa). Un valor cercano a cero sugiere que no hay una fuerte relación lineal entre  $CP_i$  y  $X_j$ .
- $Corr(CP_i, X_j)$  La correlación entre una componente principal  $CP_i$  y una variable  $X_j$  se calcula como

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

### 2.6.1 Contribuciones de las Componentes Principales a la Variabilidad Explicada de las Variables Originales

Para evaluar si las nuevas variables generadas por el ACP capturan eficazmente la información contenida en las variables originales, y para determinar qué variables se explican mejor (o peor) mediante estas nuevas componentes, utilizamos una medida conocida como ‘cosenos al cuadrado’, que se calculan como:

$$Cos_{ij}^2 = Corr^2(CP_i, X_j).$$

Los cosenos al cuadrado, a pesar de su nombre, no están relacionados directamente con la función coseno, sino que se derivan de la geometría subyacente del ACP. Estos valores representan las correlaciones al cuadrado entre cada variable original y las componentes principales. Expresan la proporción de la varianza de cada variable original que se explica mediante cada componente principal.



En otras palabras, los cosenos al cuadrado nos indican cuánto de la variabilidad de una variable original está contenida en cada componente principal. Un valor cercano a 1 implica que la variable se explica en gran medida por la componente, mientras que un valor cercano a 0 indica que la componente tiene poco poder para explicar la variabilidad de esa variable.

El uso de los cosenos al cuadrado es esencial para evaluar la capacidad del ACP para resumir y retener información relevante de las variables originales, lo que nos permite identificar qué variables son más influyentes en la definición de cada componente principal.

### 2.6.2 Contribuciones de las Variables a las Componentes Principales

Para comprender cómo cada variable original contribuye a la construcción de las CP en el Análisis de Componentes Principales, podemos calcular su contribución específica. La contribución de una variable en la creación de una CP se calcula mediante la fórmula:

$$Contribucion_{ji} = Cos_{ij}^2 \sqrt{\lambda_i} \text{ para todo } i=1,\dots,k; j=1,\dots,p.$$

Esta fórmula nos permite cuantificar la influencia de cada variable en la definición de una CP. Al evaluar las contribuciones de las variables, podemos identificar cuáles son las más relevantes para cada CP y comprender mejor cómo se estructuran las CP en función de las variables originales.

Estas contribuciones están en términos relativos, por lo que se acostumbra a normalizar estas contribuciones entre  $[0, 1]$  o entre  $[0, 100]$  para que sea más fácilmente comparable la contribución de una variable con respecto al resto en la construcción de una CP.

## 3 Ejemplo con Python

Para comprender mejor cómo funciona el Análisis de Componentes Principales (ACP) en la práctica, consideremos un ejemplo concreto que involucra las notas de 120 alumnos en siete asignaturas diferentes: Matemáticas, Física, Lengua, Inglés, Historia, Literatura y Economía. Lo que crea un conjunto de datos multidimensional. A priori, es claro que estas notas pueden mostrar correlación entre sí. Nuestro objetivo es utilizar el ACP para simplificar esta información compleja (reduciendo la dimensionalidad) y obtener una visión más clara de las tendencias subyacentes en el desempeño académico de los estudiantes (creando las componentes principales).

A medida que avanzamos en este ejemplo, exploraremos paso a paso cómo aplicar el ACP a estos datos utilizando Python, desde la selección del número adecuado de componentes principales hasta la interpretación de los resultados. A lo largo del proceso, veremos cómo el ACP nos permite reducir la dimensionalidad de nuestros datos y facilitar la toma de decisiones informadas en el entorno educativo.

El fichero con el que vamos a trabajar se llama ‘NOTAS.xlsx’ y a continuación, en la Figura 3 se muestran los primeros 15 registros.

	A	B	C	D	E	F	G	H	I
1	Alumno	Matematicas	Fisica	Lengua	Ingles	Historia	Literatura	Economia	EXTRA_ESC
2	1	8,21	6,58	6,50	7,42	9,27	8,04	5,80	Musi-Dep
3	2	5,37	3,93	6,04	6,07	6,19	7,32	4,32	Deporte
4	3	4,38	2,99	6,43	6,25	5,40	6,74	6,24	Deporte
5	4	7,64	5,81	8,02	8,07	7,28	9,03	6,34	Musi-Dep
6	5	5,39	5,76	5,21	5,68	5,44	5,78	6,47	Deporte
7	6	3,67	4,44	5,57	6,77	6,12	5,04	1,81	Deporte
8	7	4,89	4,82	6,50	6,37	7,69	6,79	5,36	Deporte
9	8	3,32	3,34	5,93	4,49	6,19	5,75	3,80	Deporte
10	9	4,21	4,43	6,06	8,17	6,53	8,44	3,24	Deporte
11	10	5,09	4,50	4,97	6,50	4,25	5,82	6,79	Deporte
12	11	6,02	5,36	4,67	5,78	6,00	6,69	6,99	Música
13	12	4,01	2,67	7,00	8,08	6,52	6,58	4,02	Deporte
14	13	6,60	6,35	5,49	6,98	7,28	6,06	5,27	Música
15	14	7,75	5,69	4,36	6,54	4,41	2,40	6,62	Musi-Dep

Figure 3: Calificaciones de los 15 primeros alumnos.

Lo primero que vamos a necesitar para empezar a trabajar, es tener todas las librerías necesarias cargadas en el entorno "prompt" de Anaconda, mediante el código "pip install nombre\_librería" como se muestra a continuación.

```
C:\Windows\system32\cmd.exe
Microsoft Windows [Versión 10.0.19045.3324]
(c) Microsoft Corporation. Todos los derechos reservados.
(base) C:\Users\Adria>pip install nombre_libreria
```

Figure 4: Instalar librerías en entorno prompt.

Una vez están cargadas las librerías en el sistema, debemos llamarlas en nuestro script. A continuación se muestran todas las librerías utilizadas para realizar un ACP completo junto con una breve descripción de su utilidad.

```
1 import os # Proporciona funciones para interactuar con el sistema operativo
2 import pandas as pd # Manipulación y análisis de datos tabulares (filas y
  columnas).
3 import numpy as np # Operaciones numéricas y matriciales.
4 import seaborn as sns # Visualización estadística de datos.
5 import matplotlib.pyplot as plt # Creación de gráficos y visualizaciones.
6
7 # Matplotlib es una herramienta versátil para crear gráficos desde cero,
8 # mientras que Seaborn simplifica la creación de gráficos estadísticos.
9
10 from sklearn.decomposition import PCA # Implementación del Análisis de
  Componentes Principales (PCA).
11 from sklearn.preprocessing import StandardScaler # Estandarización de datos
  para análisis estadísticos.
```

Listing 1: Importación de librerías en Python.

Definimos nuestro entorno de trabajo utilizando la función "chdir" de la librería "os".

```
1 os.chdir('G:/Mi unidad/Apuntes titulo propio1/Datos') #Definimos nuestro
  entorno de trabajo.
```

Listing 2: Deinición entorno de trabajo.

Una vez cargadas las librerías y definido el entorno de trabajo, procedemos a la lectura en formato DataFrame de la base de datos, la cual, esta en formato tabla o tabulado. Para ello hacemos uso de la función *read\_excel* de la librería *pandas*. Le asignamos el nombre de *notas*. Indicamos que la primera variable (*Alumno*) es el índice del DataFrame y posteriormente eliminamos dicha variable junto con la última (*EXTRA.ESC*). Quedándonos así únicamente con las variables numéricas objeto del ACP.

```
1 # Cargar un archivo Excel llamado 'notas.xlsx' en un DataFrame llamado
  notas.
2 notas = pd.read_excel('notas.xlsx')
3
4 # Establecer la columna 'Alumno' como índice del DataFrame notas y
  eliminarla.
5 notas = notas.set_index('Alumno', drop=True)
6
7 # Quita 'EXTRA.ESC' y crea versión con sólo 2 decimales del data frame
8 notas = notas.drop('EXTRA.ESC', axis=1)
9 notas_2 = notas.round(decimals = 2)
10 notas_2.head()
```

Listing 3: Carga y preparación de los datos.

Antes de empezar con el ACP es conveniente realizar un análisis descriptivo de las variables, para observar si hay alguna anomalía en los datos, como datos fuera de rango, atípicos o perdidos.

```
1
2 # Cálculo de los estadísticos descriptivos.
3 # Genera una lista con los nombres de las variables.
4 variables = list(notas)
5
6 # Calcula los estadísticas descriptivas para cada variable y crea un
  DataFrame con los resultados.
7 estadisticos = pd.DataFrame({
8     'Mínimo': notas[variables].min(),
9     'Percentil 25': notas[variables].quantile(0.25),
10    'Mediana': notas[variables].median(),
11    'Percentil 75': notas[variables].quantile(0.75),
12    'Media': notas[variables].mean(),
13    'Máximo': notas[variables].max(),
14    'Desviación Estándar': notas[variables].std(),
15    'Varianza': notas[variables].var(),
```

```

16     'Coeficiente de Variación': (notas[variables].std() / notas[variables].
17     mean()),
17     'Datos Perdidos': notas[variables].isna().sum() # Cuenta los valores
18     NaN por variable.
18 })

```

Listing 4: Estudio descriptivo de los datos.

Como podemos observar en la Figura 5 los datos no muestran ningún indicio de anomalías, contamos con todos los datos, y estos, están en los rangos esperados.

Índice	Mínimo	Percentil 25	Mediana	Percentil 75	Media	Máximo	Desviación Estándar	Varianza	Coeficiente de Variación	Datos Perdidos
Matematicas	1.65588	4.33087	5.26635	6.2365	5.30219	8.85618	1.46979	2.16028	0.277204	0
Fisica	9.211836	3.55401	4.50995	5.64319	4.63173	8.3716	1.51536	2.2963	0.327168	0
Lengua	2.96587	5.14468	6.04806	6.62505	5.97155	8.57129	1.27483	1.6252	0.213484	0
Ingles	3.27163	5.77339	6.7487	8.03096	6.80625	10.6239	1.47323	2.17041	0.216452	0
Historia	2.32689	5.55147	6.26843	7.13108	6.35473	9.41208	1.37323	1.88575	0.216095	0
Literatura	2.4029	5.61109	6.58599	7.5491	6.61477	10.4562	1.45661	2.1217	0.220205	0
Economia	1.64023	4.04546	5.07997	6.14481	5.1205	8.77921	1.50773	2.27324	0.294449	0

Figure 5: Descriptivos de las variables originales.

También es importante analizar las relaciones entre las variables, ya que si estas mostrasen fuerte independencia, no tendría sentido realizar un ACP. Para ello calculamos la matriz de correlaciones.

```

1  # Calcula la matriz de correlación entre las variables del DataFrame 'notas
2  # '.
2  R = notas.corr()
3
4  # Crea una nueva figura de dimensiones 10x8 para el gráfico.
5  plt.figure(figsize=(10, 8))
6
7  # Genera un mapa de calor (heatmap) de la matriz de correlación 'R'
8  # utilizando Seaborn.
8  # 'annot=True' agrega los valores de correlación en las celdas.
9  # 'cmap' establece el esquema de colores (en este caso, 'coolwarm' para
10 # colores fríos y cálidos).
10 # 'fmt' controla el formato de los números en las celdas ('.2f' para dos
11 # decimales).
11 # 'linewidths' establece el ancho de las líneas que separan las celdas.
12 sns.heatmap(R, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5)

```

Listing 5: Matriz de correlaciones y representación.

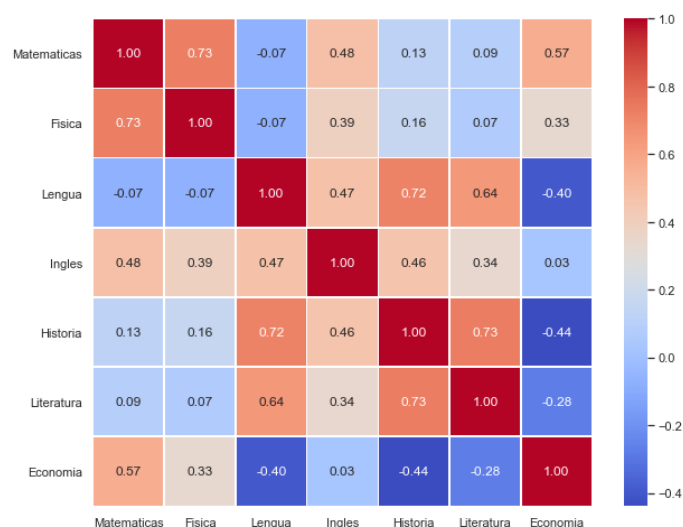


Figure 6: Gráfico de correlaciones de las variables originales.

Observemos que las variables con mayor correlación positiva son Lengua e Historia, Matemáticas y Física y también Historia y Literatura. Por otro lado, se observa una destacable correlación negativa entre las notas de Economía con las asignaturas de Lengua e Historia.

Ahora si, una vez cargados y validado los datos, y observado las *dependencias* entre las variables, estamos en disposición de empezar el ACP. Para ello, lo primero, estandarizamos los datos, que pese a estar en la misma escala y no ser imprescindible para realizar el análisis, siempre es aconsejable.

```

1 # Estandarizamos los datos:
2 notas_std = StandardScaler().fit_transform(notas)

```

Listing 6: Estandarización de los datos.

Para llevar a cabo nuestro ACP vamos a hacer uso del objeto PCA de la librería scikit-learn. Este objeto proporciona una variedad de resultados y métodos que permiten explorar y comprender los componentes principales de un conjunto de datos, así como reducir su dimensionalidad. A continuación, se presenta una lista de las opciones de configuración, resultados y métodos disponibles en un objeto PCA, junto con una breve descripción de cada uno de ellos. Estos recursos son esenciales para realizar análisis y transformaciones de datos en el espacio de los componentes principales.

### Opciones de Configuración:

- **n\_components:** Número de componentes principales a calcular y mantener.
- **whiten:** Indica si queremos que la matriz de componentes principales tenga varianza unitaria.
  - **whiten=False:** (valor predeterminado), las varianzas de las componentes no se estandarizan, lo que significa que pueden tener varianzas diferentes.
  - **whiten=True:** Esto puede ser útil para simplificar la interpretación de las componentes y cuando se utilizan en algoritmos de clasificación o regresión.

- **svd\_solver**: El algoritmo que PCA utilizará para calcular los componentes principales. La elección del **svd\_solver** dependerá de la cantidad de datos y componentes que estés manejando.
  1. **'auto'** (Opción por defecto): Seleccionará automáticamente el método más adecuado en función del tamaño de los datos y **n\_components**.
  2. **'randomized'**: Suele ser una buena opción cuando tienes un conjunto de datos grande y solo necesitas un subconjunto de componentes principales.
  3. **'full'**: Es adecuado para conjuntos de datos más pequeños o cuando necesitas calcular todos los componentes principales.
  4. **'arnold'**: Resulta útil cuando necesitas un número específico de componentes principales y deseas un equilibrio entre eficiencia computacional y precisión.

### Resultados:

- **components\_**: Devuelve los vectores de los componentes principales.
- **explained\_variance\_**: Devuelve los autovalores asociados a la matriz de correlaciones o covarianzas, los cuales, están ordenados en orden descendente.
- **explained\_variance\_ratio\_**: Devuelve la proporción de la varianza total explicada por cada componente principal.
- **mean\_**: Devuelve la media de cada variable original antes de la estandarización (si no se realizó).
- **n\_components\_**: Devuelve el número de componentes principales que se mantuvieron.
- **noise\_variance\_**: Proporciona una medida de cuánta información se ha perdido al reducir la dimensionalidad de los datos con PCA (solo disponible si se especifica **whiten=True**).
- **n\_features\_**: Devuelve el número de variables originales en los datos de entrada.

### Métodos:

- **fit(X)**: Ajusta el modelo PCA a los datos de entrada **X**.
- **fit\_transform(X)**: Ajusta el modelo PCA a los datos de entrada **X** y luego transforma **X** en el espacio de los componentes principales.
- **transform(X)**: Transforma los datos de entrada **X** en el espacio de los componentes principales, utilizando el modelo previamente ajustado.
- **inverse\_transform(X\_reduced)**: Realiza la inversa de la transformación PCA para llevar los datos desde el espacio de los componentes principales de nuevo al espacio original.

Procedemos a realizar el ACP a nuestro ejemplo mediante el siguiente código:

```

1      # Crea una instancia de Análisis de Componentes Principales (ACP):
2      # - Utilizamos PCA(n_components=7) para crear un objeto PCA que realizará
      un análisis de componentes principales.
3      # - Establecemos n_components en 7 para retener el maximo de las
      componentes principales (maximo= numero de variables).
4      pca = PCA(n_components=7)
5
6      # Aplicar el Análisis de Componentes Principales (ACP) a los datos
      estandarizados:
7      # - Usamos pca.fit(notas_estandarizadas) para ajustar el modelo de ACP a
      los datos estandarizados.
8      fit = pca.fit(notas_estandarizadas)

```

Listing 7: Aplicación del ACP.

A continuación, procedemos a calcular los autovalores asociados a las componentes principales (valores  $\lambda$  que resuelven la ecuación  $\det(\Sigma - \lambda I) = 0$ ), la varianza explicada y la varianza acumulada de las mismas. Así como su representación en una tabla.

```

1      # Obtener los autovalores asociados a cada componente principal.
2      autovalores = fit.explained_variance_
3
4      # Obtener la varianza explicada por cada componente principal como un
      porcentaje de la varianza total.
5      var_explicada = fit.explained_variance_ratio_
6
7      # Calcular la varianza explicada acumulada a medida que se agregan cada
      componente principal.
8      var_acumulada = np.cumsum(var_explicada)
9
10     # Crear un DataFrame de pandas con los datos anteriores y establecer índice
      .
11     data = {'Autovalores': autovalores, 'Variabilidad Explicada': var_explicada
      , 'Variabilidad Acumulada': var_acumulada}
12     tabla = pd.DataFrame(data, index=['Componente {}'.format(i) for i in range
      (1, fit.n_components_+1)])

```

Listing 8: Cálculo de los autovalores/variabilidad explicada.

Como podemos observar en la Figura 7, usando las dos primeras componentes conseguiremos explicar el 74.9% de la variabilidad total de las variables originales. La decisión de cuántas componentes seleccionar es subjetiva y se debe decidir teniendo en cuenta dos factores, por un lado, la cantidad de variabilidad que deseamos explicar y por otro, la dimensión que queremos conseguir. En este ejemplo, consideramos que añadiendo una tercera componente no se consigue *demasiado* beneficio en cuanto a la cantidad de variabilidad explicada. No obstante, esta decisión podría ser diferente, si por ejemplo, quisiéramos explicar un mínimo superior de variabilidad.

Índice	Autovalores	Variabilidad Explicada	Variabilidad Acumulada
Componente 1	2.94116	41.6664	41.6664
Componente 2	2.34795	33.2626	74.929
Componente 3	0.574517	8.13899	83.068
Componente 4	0.561622	7.95632	91.0243
Componente 5	0.265562	3.76214	94.7865
Componente 6	0.220849	3.1287	97.9152
Componente 7	0.147164	2.08482	100

Figure 7: Gráfico de correlaciones de las variables originales.

Realizamos la representación de la variabilidad explicada, para observar gráficamente como esta se estabiliza una vez seleccionadas las dos primeras componentes.

```

1  # Representacion de la variabilidad explicada:
2  #Creamos la función plot_varianza_explicada
3  #Args:
4  #var_explicada (array): Un array que contiene el porcentaje de
5  #varianza explicada
6  #por cada componente principal. Generalmente calculado como
7  #var_explicada = fit.explained_variance_ratio_ * 100.
8  #n_components (int): El número total de componentes principales.
9  #Generalmente calculado como fit.n_components.
10
11 def plot_varianza_explicada(var_explicada, n_components):
12     # Crear un rango de números de componentes principales de 1 a
13     # n_components
14     num_componentes_range = np.arange(1, n_components + 1)
15
16     # Crear una figura de dimensiones 8x6
17     plt.figure(figsize=(8, 6))
18
19     # Trazar la varianza explicada en función del número de componentes
20     # principales
21     plt.plot(num_componentes_range, var_explicada, marker='o')
22
23     # Etiquetas de los ejes x e y
24     plt.xlabel('Número de Componentes Principales')
25     plt.ylabel('Varianza Explicada')
26
27     # Título del gráfico
28     plt.title('Variabilidad Explicada por Componente Principal')
29
30     # Establecer las marcas en el eje x para que coincidan con el número de
31     # componentes
32     plt.xticks(num_componentes_range)
33
34     # Mostrar una cuadrícula en el gráfico

```



```

31 plt.grid(True)
32
33 # Agregar barras debajo de cada punto para representar el porcentaje de
    variabilidad explicada
34 # - 'width': Ancho de las barras de la barra. En este caso, se
    establece en 0.2 unidades.
35 # - 'align': Alineación de las barras con respecto a los puntos en el
    eje x.
36 # 'center' significa que las barras estarán centradas debajo de los
    puntos.
37 # - 'alpha': Transparencia de las barras. Un valor de 0.7 significa que
    las barras son 70% transparentes.
38 plt.bar(num_componentes_range, var_explicada, width=0.2, align='center',
    , alpha=0.7)
39
40 # Mostrar el gráfico
41 plt.show()
42
43 plot_varianza_explicada(var_explicada, fit.n_components_)

```

Listing 9: Representación de la variabilidad explicada por Componentes Principales.

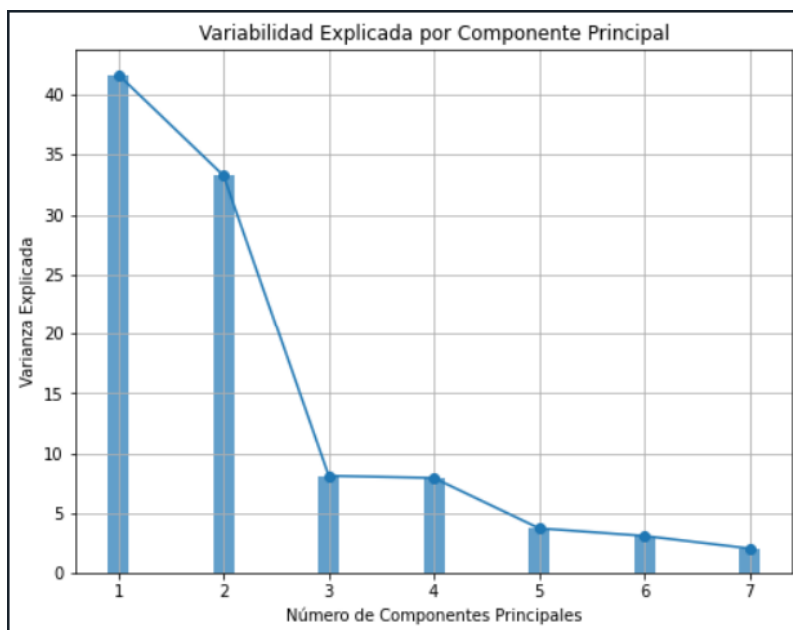


Figure 8: Varianza explicada.

Podemos observar que la variabilidad explicada por cada componente, decae y se estabiliza a partir de la tercera. Por lo tanto, optamos por retener únicamente dos de las siete componentes principales, las cuales explican conjuntamente el 74.91% de la varianza de las variables originales. Los

dos autovalores correspondientes a estas componentes seleccionadas son 2.94 y 2.35, representando el 41.7% y el 33.3% de la variabilidad, respectivamente.

Una vez hemos decidido quedarnos con las dos primeras componentes, procedemos a calcularlas. Para ello necesitamos conocer la matriz de transformación formada por los autovectores o vectores propios asociados a los autovalores seleccionados.

```

1 # Crea una instancia de ACP con las dos primeras componentes que nos
  # interesan y aplicar a los datos.
2 pca = PCA(n_components=2)
3 fit = pca.fit(notas_estandarizadas)
4
5 # Obtener los autovectores asociados a cada componente principal y
  # transponerlos.
6 autovectores = pd.DataFrame(pca.components_.T,
7                               columns = ['Autovector {}'.format(i) for i in
8                                           range(1, fit.n_components_+1)],
7                               index = ['{}_z'.format(variable) for variable
8                                           in variables])

```

Listing 10: Cálculo de los autovalores/variabilidad explicada.

Transponemos los autovectores para tenerlos por columnas e indicamos que las filas están representadas por las variables originales tipificadas  $Z$ .

Índice	Autovector 1	Autovector 2
Matematicas_z	-0.117466	-0.602636
Fisica_z	-0.128305	-0.53538
Lengua_z	-0.496457	0.162045
Ingles_z	-0.388664	-0.293545
Historia_z	-0.534553	0.0568696
Literatura_z	-0.482633	0.0657072
Economia_z	0.231363	-0.479813

Figure 9: Matriz de transformación  $W = [w_1, w_2]$ .

Una vez tenemos nuestros datos tipificados  $Z$  y nuestra matriz de transformación  $W$  procedemos a calcular las componentes principales  $CP$  tales que  $CP = ZW$ , es decir,

$$CP_1 = -0.117Mat_z - 0.128Fis_z - 0.496Len_z - 0.388Ing_z - 0.534His_z - 0.482Lit_z + 0.231Eco_z$$

$$CP_2 = -0.602Mat_z - 0.535Fis_z + 0.162Len_z - 0.293Ing_z + 0.056His_z + 0.065Lit_z - 0.479Eco_z.$$

```

1 # Crea una instancia de ACP con las dos primeras componentes que nos
  # interesan y aplicar a los datos.
2 pca = PCA(n_components=2)
3 fit = pca.fit(notas_estandarizadas)

```

```

4
5 # Obtener los autovectores asociados a cada componente principal y
   transponerlos.
6 autovectores = pd.DataFrame(pca.components_.T,
7                               columns = ['Autovector {}'.format(i) for i in
                                           range(1, fit.n_components_+1)],
8                               index = ['{}_z'.format(variable) for variable
                                         in variables])
9
10 resultados_pca = pd.DataFrame(fit.transform(notas_estandarizadas),
11                               columns=['Componente {}'.format(i) for i in
                                           range(1, fit.n_components_+1)],
12                               index=notas_estandarizadas.index)
13
14
15 notas_z_cp = pd.concat([notas_estandarizadas, resultados_pca], axis=1)
16 variables_cp = notas_z_cp.columns

```

Listing 11: Cálculo de las componentes principales.

Se muestra a continuación, en la Figura 10, la base de datos estandarizadas (5 primeras observaciones) con las dos componentes seleccionadas indexadas.

Alumno	Matematicas_z	Fisica_z	Lengua_z	Ingles_z	Historia_z	Literatura_z	Economia_z	Componente 1	Componente 2
1	1.98479	1.2881	0.412892	0.4168	2.13099	0.979363	0.454951	-2.27193	-1.97392
2	0.0463853	-0.462778	0.0538108	-0.501058	-0.119002	0.483532	-0.531895	-0.0708595	0.655825
3	-0.627619	-1.08649	0.357885	-0.381331	-0.699891	0.0874642	0.743655	0.687632	0.738973
4	1.59566	0.782447	1.61198	0.862701	0.675615	1.6637	0.809081	-2.40033	-1.61301
5	0.0571131	0.747611	-0.598737	-0.764955	-0.665394	-0.578819	0.901758	1.33561	-0.815696

Figure 10: Primeras 5 observaciones con las Componentes Principales seleccionadas.

Este es el momento en el que podríamos *eliminar* nuestras variables originales y trabajar con las componentes principales calculadas, por ejemplo, en objetivos más *avanzados* en el análisis de datos como es el caso de la *predicción*. Hemos entonces, reducido la dimensionalidad de nuestra base de datos original, sustituyendo un total de 7 variables por tan solo dos componentes, las cuales, combinadas, explican aproximadamente un 75% de la variabilidad original.

Vamos ahora a estudiar cómo de fuerte (o débil) es la relación entre las componentes principales seleccionadas y cada una de las variables originales. Para ello calculamos la correlación entre las variables y las componentes retenidas.

```

1
2 # Cálculo de las correlaciones entre las variables originales y las
   componentes seleccionadas.
3 # Guardamos el nombre de las variables del archivo conjunto (variables y
   componentes).
4 variables_cp = notas_z_cp.columns
5

```

```

6 # Calculamos las correlaciones y seleccionamos las que nos interesan (
  variables contra componentes).
7 correlacion = pd.DataFrame(np.corrcoef(notas_estandarizadas.T,
  resultados_pca.T),
8                               index = variables_cp, columns = variables_cp)
9
10 n_variables = fit.n_features_
11 correlaciones_notas_con_cp = correlacion.iloc[:fit.n_features_, fit.
  n_features_:]

```

Listing 12: Correlación entre variables y componentes.

Como podemos observar en la Figura 11 la primera componente principal ( $CP_1$ ) muestra una correlación negativa significativa con las variables Lengua, Historia y Literatura. Esto sugiere que estas variables están fuertemente influenciadas por  $CP_1$  y, por lo tanto, están representadas de manera directa por esta componente. Por otro lado, la Componente Principal 2 ( $CP_2$ ) presenta una correlación notable solo con Matemáticas, Física y Economía, lo que indica que estas variables están predominantemente influenciadas por  $CP_2$ . Además, es interesante observar que la nota en inglés está correlacionada con ambas componentes principales, lo que sugiere cierta influencia mixta de ambas  $CP_1$  y  $CP_2$  en esta variable.

Índice	Componente 1	Componente 2
Matematicas_z	-0.200611	-0.919564
Fisica_z	-0.219122	-0.816938
Lengua_z	-0.847859	0.247265
Ingles_z	-0.663768	-0.447922
Historia_z	-0.912921	0.0867775
Literatura_z	-0.82425	0.100263
Economia_z	0.395127	-0.732148

Figure 11: Correlaciones entre variables y componentes.

Podemos comprobar que la correlación entre cada par (*Componente, variable*) se puede calcular como

$$Corr(CP_i, x_j) = w_{i,j} \frac{\sqrt{\lambda_i}}{s_j}, \text{ para } i=1,2 \text{ y para } j=1,\dots,7.$$

Por ejemplo, la correlación entre la asignatura de Física y la Componente 1 sería:

$$Corr(CP_1, x_2) = w_{1,2} \frac{\sqrt{\lambda_1}}{s_2} = -0.128 \frac{\sqrt{2.941}}{1} = -0.219.$$

Si multiplicamos estas correlaciones al cuadrado, conseguimos los valores conocidos como *cosenos al cuadrado* que representan la proporción de variabilidad que cada componente explica de cada variable.

```
1 cos2 = correlaciones_notas_con_cp **2
```

Listing 13: cos2.

En la figura 12 podemos observar, por ejemplo, que la variable Matemáticas es la que más explicada esta por el conjunto de las dos componentes, mientras que ingles es la que menos.

$$\text{Corr}(CP_1, \text{Mat})^2 + \text{Corr}(CP_2, \text{Mat})^2 = 0.04 + 0.84 = 0.885,$$

$$\text{Corr}(CP_1, \text{Ing})^2 + \text{Corr}(CP_2, \text{Ing})^2 = 0.44 + 0.20 = 0.64.$$

Índice	Componente 1	Componente 2
Matematicas_z	0.0402447	0.845598
Fisica_z	0.0480146	0.667387
Lengua_z	0.718864	0.0611399
Ingles_z	0.440588	0.200634
Historia_z	0.833424	0.00753034
Literatura_z	0.679388	0.0100527
Economia_z	0.156125	0.536041

Figure 12:  $\text{Cos}^2$

Aunque para comparar la variabilidad explicada por las componentes entre el conjunto de variables, es más *visual* el siguiente gráfico de barras, donde se representa la suma de los  $\text{cos}^2$  para cada variable.

```
1
2
3 def plot_cos2_bars(cos2):
4     """
5     Genera un gráfico de barras para representar la varianza explicada de
6     cada variable utilizando los cuadrados de las cargas (cos^2).
7
8     Args:
9         cos2 (pd.DataFrame): DataFrame que contiene los cuadrados de las
10        cargas de las variables en las componentes principales.
```

```

11     None
12     """
13     # Crea una figura de dimensiones 8x6 pulgadas para el gráfico
14     plt.figure(figsize=(8, 6))
15
16     # Crea un gráfico de barras para representar la varianza explicada por
17     # cada variable
18     sns.barplot(x=cos2.sum(axis=1), y=cos2.index, color="blue")
19
20     # Etiqueta los ejes
21     plt.xlabel('Suma de los cos^2')
22     plt.ylabel('Variables')
23
24     # Establece el título del gráfico
25     plt.title('Varianza Explicada de cada Variable por las Componentes
26               Principales')
27
28     # Muestra el gráfico
29     plt.show()
30 plot_cos2Bars(cos2)

```

Listing 14: Gráfico de barras para los  $\cos^2$ .

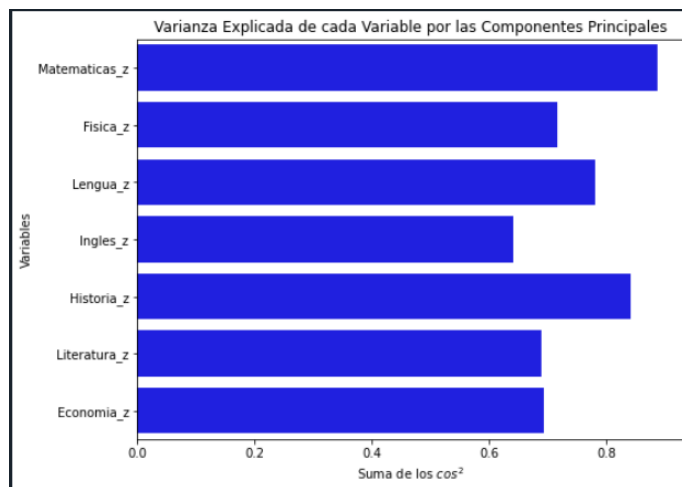


Figure 13: Gráfico de barras para la suma de los  $\cos^2$ .

Con la siguiente función *plot\_corr\_cos* generamos diferentes gráficos (en este caso particular solo 1), donde se representan las correlaciones entre cada variable y cada par de componentes principales (ejes). El color de cada vector en los gráficos indica el valor de los cosenos al cuadrado.

1

```

2 def plot_corr_cos(n_components, correlaciones_notas_con_cp):
3
4     """
5     Genera un graficos en los que se representa un vector por cada variable
6     , usando como ejes las compoenntes, la orientación y la longitud del
7     vector representa la correlación
8     entre cada variable y dos de las componentes. El color representa el
9     valor de la suma de los cosenos al cuadrado
10
11     Args:
12         n_components (int): Número entero que representa el número de
13         componentes principales seleccionadas.
14         correlaciones_notas_con_cp (DataFrame): DataFrame que contiene la
15         matriz de coorrelaciones entre variables y componentes
16
17     """
18     # Definir un mapa de color (cmap) sensible a las diferencias numéricas
19     cmap = plt.get_cmap('coolwarm') # Puedes ajustar el cmap según tus
20     preferencias
21
22     for i in range(n_components):
23         for j in range(i + 1, n_components): # Evitar pares duplicados
24             # Calcular la suma de los cosenos al cuadrado
25             sum_cos2 = correlaciones_notas_con_cp.iloc[:, i] ** 2 +
26                 correlaciones_notas_con_cp.iloc[:, j] ** 2
27
28             # Crear un nuevo gráfico para cada par de componentes
29             principales
30             plt.figure(figsize=(10, 10))
31
32             # Dibujar un círculo de radio 1
33             circle = plt.Circle((0, 0), 1, fill=False, color='b', linestyle
34                 ='dotted')
35
36             plt.gca().add_patch(circle)
37
38             # Dibujar vectores para cada variable con colores basados en la
39             suma de los cosenos al cuadrado
40             for k, var_name in enumerate(correlaciones_notas_con_cp.index):
41                 x = correlaciones_notas_con_cp.iloc[k, i] # Correlación en
42                 la primera dimensión
43                 y = correlaciones_notas_con_cp.iloc[k, j] # Correlación en
44                 la segunda dimensión
45
46                 # Seleccionar un color de acuerdo a la suma de los cosenos
47                 al cuadrado
48                 color = cmap(sum_cos2[k])
49
50                 # Dibujar el vector con el color seleccionado
51                 plt.quiver(0, 0, x, y, angles='xy', scale_units='xy', scale
52                     =1, color=color)

```

```

38         # Agregar el nombre de la variable junto a la flecha con el
39         mismo color
40         plt.text(x, y, var_name, color=color, fontsize=12, ha='
41             right', va='bottom')
42
43     # Dibujar líneas discontinuas que representen los ejes
44     plt.axhline(0, color='black', linestyle='--', linewidth=0.8)
45     plt.axvline(0, color='black', linestyle='--', linewidth=0.8)
46
47     # Etiquetar los ejes
48     plt.xlabel(f'Componente Principal {i + 1}')
49     plt.ylabel(f'Componente Principal {j + 1}')
50
51     # Establecer los límites del gráfico
52     plt.xlim(-1.1, 1.1)
53     plt.ylim(-1.1, 1.1)
54
55     # Agregar un mapa de color (colorbar) y su leyenda
56     sm = plt.cm.ScalarMappable(cmap=cmap)
57     sm.set_array([]) # Evita errores de escala
58     plt.colorbar(sm, orientation='vertical', label='cos^2') #
59     Agrega la leyenda
60
61     # Mostrar el gráfico
62     plt.grid()
63     plt.show()
64
65 plot_corr_cos(fit.n_components, correlaciones_notas_con_cp)

```

Listing 15: Grafico de las correlaciones entre variables y componentes.

Como se enunciaba anteriormente, en el gráfico de la Figura 14 se representa un vector por cada asignatura, siendo los ejes cada una de las componentes seleccionadas, de tal forma que la longitud del vector en cada uno de los ejes representa la correlación en dicha componente, el color del vector es la suma de estas correlaciones al cuadrado ( $\cos^2$ ). En este gráfico podemos visualizar lo ofrecido por la Figura 11, como las asignaturas de '*letras*' están fuertemente relacionadas con la segunda componente, las de '*ciencias*' con la primera e Inglés medianamente con ambas. Pero además podemos observar que las variables con el color *rojo* más intenso son Matemáticas e Historia, lo cual indica que son las variables mejor explicadas por el conjunto de componentes.



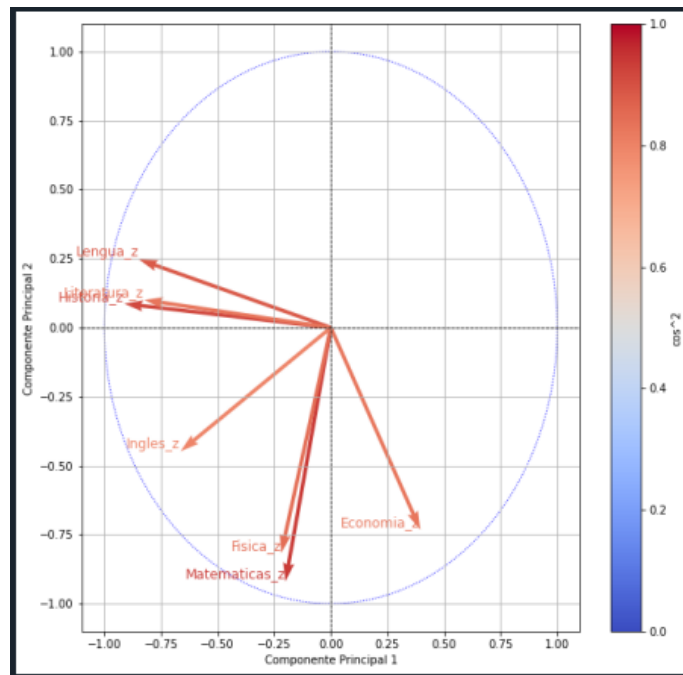


Figure 14: Correlaciones entre Variables y Componentes Principales.

Procedemos a continuación a calcular las contribuciones normalizadas de cada variable en la construcción de cada una de las componentes. Representamos los resultados mediante un *mapa de calor* donde el valor de las contribuciones dependerá de la intensidad del color azul utilizado.

```

1 def plot_contribuciones_proporcionales(cos2, autovalores, n_components):
2     """
3     Cacula las contribuciones de cada variable a las componentes
4     principales y
5     Genera un gráfico de mapa de calor con los datos
6     Args:
7         cos2 (DataFrame): DataFrame de los cuadrados de las cargas (cos^2).
8         autovalores (array): Array de los autovalores asociados a las
9         componentes principales.
10        n_components (int): Número de componentes principales seleccionadas
11        .
12    """
13    # Calcula las contribuciones multiplicando cos2 por la raíz cuadrada de
14    # los autovalores
15    contribuciones = cos2 * np.sqrt(autovalores)
16
17    # Inicializa una lista para las sumas de contribuciones
18    sumas_contribuciones = []

```

```

17 # Calcula la suma de las contribuciones para cada componente principal
18 for i in range(n_components):
19     nombre_componente = f'Componente {i + 1}'
20     suma_contribucion = np.sum(contribuciones[nombre_componente])
21     sumas_contribuciones.append(suma_contribucion)
22
23 # Calcula las contribuciones proporcionales dividiendo por las sumas de
    contribuciones
24 contribuciones_proporcionales = contribuciones.div(sumas_contribuciones
    , axis=1) * 100
25
26 # Crea una figura de dimensiones 8x8 pulgadas para el gráfico
27 plt.figure(figsize=(8, 8))
28
29 # Utiliza un mapa de calor (heatmap) para visualizar las contribuciones
    proporcionales
30 sns.heatmap(contribuciones_proporcionales, cmap='Blues', linewidths
    =0.5, annot=False)
31
32 # Etiqueta los ejes (puedes personalizar los nombres de las filas y
    columnas si es necesario)
33 plt.xlabel('Componentes Principales')
34 plt.ylabel('Variables')
35
36 # Establece el título del gráfico
37 plt.title('Contribuciones Proporcionales de las Variables en las
    Componentes Principales')
38
39 # Muestra el gráfico
40 plt.show()
41
42 # Devuelve los DataFrames de contribuciones y contribuciones
    proporcionales
43 return contribuciones_proporcionales
44
45 contribuciones_proporcionales = plot_contribuciones_proporcionales(cos2,
    autovalores, fit.n_components)

```

Listing 16: Contribuciones de las variables a la construcción de las CP.

En la Figura 15 volvemos a corroborar lo ya observado anteriormente, las variables de *letras*, participan fuertemente en la construcción de la primera Componente, las de *ciencias* en la segunda, e ingles medianamente en ambas.

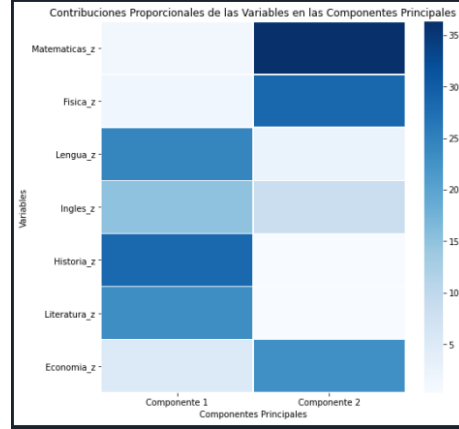


Figure 15: Contribuciones de las Variables en las Componentes Principales.

### 3.1 Coordenadas de los individuos en los nuevos ejes

Por último, vamos a calcular las nuevas coordenadas (o valores) de cada uno de los individuos en los nuevos ejes (Componentes Principales), para ello hacemos uso de nuestra matriz  $X$  o  $Z$  donde tenemos nuestras variables o variables tipificadas y de la matriz de transformación  $W$ . De tal manera que el valor de la  $i$  –ésima observación en la  $\alpha$  –ésima componente, se calcula como:

$$CP_{\alpha,i} = \sum_{j=1}^p z_{i,j} w_{\alpha,j} = \sum_{j=1}^p \frac{x_{i,j} - \bar{x}_j}{S_j} w_{\alpha,j}.$$

Como las variables están estandarizadas, es sencillo comprobar que la suma de los valores de los  $n$  individuos sobre cada Componente es nula,

$$\sum_{i=1}^n CP_{\alpha,i} = \sum_{i=1}^n \sum_{j=1}^p \frac{x_{i,j} - \bar{x}_j}{S_j} w_{\alpha,j} = \sum_{i=1}^n w_{\alpha,j} \sum_{j=1}^p \frac{x_{i,j} - \bar{x}_j}{S_j} = 0,$$

por lo que la media de las Componentes es cero y su varianza  $\lambda_{\alpha}$ .

Continuando con el ejemplo, tenemos en la siguiente figura, los valores de los 10 primeros alumnos en las Componentes Principales.

```

1
2 # Indexamos las componentes principales a la base de datos original.
3 notas_cp = pd.concat([notas, extra, resultados_pca], axis=1)

```

Listing 17: Base de datos original más CP.

Recordamos que los valores de las medias y desviaciones estándar de las variables originales son:

Alumno	Matematicas	Fisica	Lengua	Ingles	Historia	Literatura	Economia	EXTRA_ESC	Componente 1	Componente 2
1	8.20723	6.57551	6.49572	7.41773	9.26885	8.03536	5.80358	Musi-Dep	-2.27193	-1.97392
2	5.37009	3.93339	6.03986	6.07116	6.192	7.31615	4.32189	Deporte	-0.0708595	0.655825
3	4.38358	2.99218	6.42589	6.24681	5.39764	6.74164	6.23704	Deporte	0.687632	0.738973
4	7.63769	5.81247	8.01798	8.0719	7.27863	9.02802	6.33528	Musi-Dep	-2.40033	-1.61301
5	5.38579	5.7599	5.21145	5.684	5.44481	5.77518	6.47443	Deporte	1.33561	-0.815696
6	3.67495	4.4407	5.57038	6.76995	6.11794	5.03942	1.81466	Deporte	0.420651	1.66906
7	4.88547	4.8183	6.49907	6.37287	7.69025	6.79006	5.36362	Deporte	-0.616813	0.245222
8	3.3211	3.33974	5.93082	4.49104	6.1948	5.74971	3.79899	Deporte	1.04484	2.10859
9	4.21363	4.42519	6.05626	8.16895	6.52674	8.4445	3.2434	Deporte	-1.25451	0.949534
10	5.08802	4.49552	4.97355	6.50225	4.24672	5.8174	6.79072	Deporte	1.8463	-0.587592

Figure 16: Contribuciones de las Variables en las Componentes Principales.

Luego, el valor que toma el primer alumno en la Primera Componente, viene dado por:

$$CP_{11} = -\frac{8.2 - 5.30}{1.47}0.117 - \frac{6.57 - 4.63}{1.52}0.128 - \frac{6.49 - 5.97}{1.27}0.496 - \frac{7.41 - 6.81}{1.47}0.389 \\ - \frac{9.26 - 6.35}{1.37}0.534 - \frac{8.03 - 6.61}{1.46}0.483 + \frac{5.80 - 5.12}{1.46}0.231 = -2.27.$$

A continuación vamos a representar a los individuos utilizando como ejes cada par de Componentes principales (en este ejemplo solo 1).

```

1 def plot_pca_scatter(pca, datos_estandarizados, n_components):
2     """
3     Genera gráficos de dispersión de observaciones en pares de componentes
4     principales seleccionados.
5
6     Args:
7         pca (PCA): Objeto PCA previamente ajustado.
8         datos_estandarizados (pd.DataFrame): DataFrame de datos
9         estandarizados.
10        n_components (int): Número de componentes principales seleccionadas
11        .
12    """
13    # Representamos las observaciones en cada par de componentes
14    # seleccionadas
15    componentes_principales = pca.transform(datos_estandarizados)
16
17    for i in range(n_components):
18        for j in range(i + 1, n_components): # Evitar pares duplicados
19            # Calcular la suma de los valores al cuadrado para cada
20            # variable
21            # Crea un gráfico de dispersión de las observaciones en las dos
22            # primeras componentes principales
23            plt.figure(figsize=(8, 6)) # Ajusta la dimension de la figura
24            # si es necesario

```

```

19 plt.scatter(componentes_principales[:, i],
20             componentes_principales[:, j])
21
22 # Agrega etiquetas a las observaciones
23 etiquetas_de_observaciones = list(datos_estandarizados.index)
24
25 for k, label in enumerate(etiquetas_de_observaciones):
26     plt.annotate(label, (componentes_principales[k, i],
27                         componentes_principales[k, j]))
28
29 # Dibujar líneas discontinuas que representen los ejes
30 plt.axhline(0, color='black', linestyle='--', linewidth=0.8)
31 plt.axvline(0, color='black', linestyle='--', linewidth=0.8)
32
33 # Etiquetar los ejes
34 plt.xlabel(f'Componente Principal {i + 1}')
35 plt.ylabel(f'Componente Principal {j + 1}')
36
37 # Establece el título del gráfico
38 plt.title('Gráfico de Dispersión de Observaciones en PCA')
39
40 plt.show()
41
42 plot_pca_scatter(pca, notas_estandarizadas, fit.n_components)

```

Listing 18: Gráfico de dispersión de los individuos por cada par de CPs.

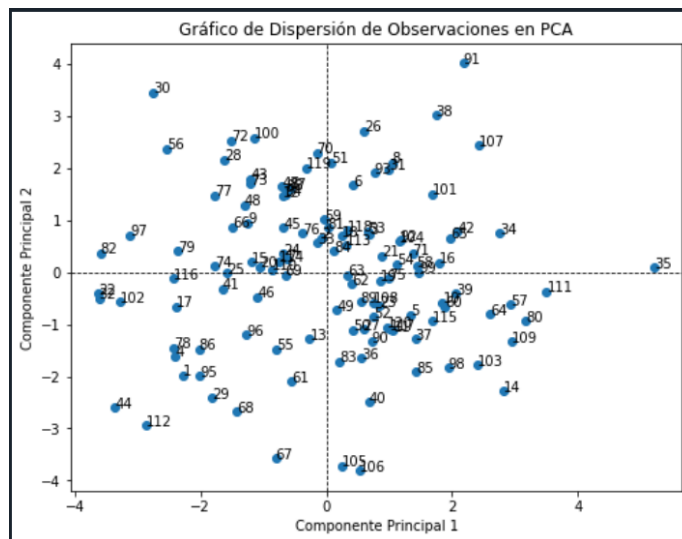


Figure 17: Dispersión de los individuos en los nuevos ejes.

Podemos modificar el código anterior para que el gráfico nos incluya un vector que represente

para cada una de las variables la correlación escalada con las componentes (ejes) del gráfico.

```
1
2
3 def plot_pca_scatter_with_vectors(pca, datos_estandarizados, n_components,
4 components_):
5     """
6     Genera gráficos de dispersión de observaciones en pares de componentes
7     principales seleccionados
8     con vectores de las correlaciones escaladas entre variables y
9     componentes
10
11     Args:
12         pca (PCA): Objeto PCA previamente ajustado.
13         datos_estandarizados (pd.DataFrame): DataFrame de datos
14         estandarizados.
15         n_components (int): Número de componentes principales seleccionadas
16         .
17         components_: Array con las componentes.
18     """
19     # Representamos las observaciones en cada par de componentes
20     seleccionadas
21     componentes_principales = pca.transform(datos_estandarizados)
22
23     for i in range(n_components):
24         for j in range(i + 1, n_components): # Evitar pares duplicados
25             # Calcular la suma de los valores al cuadrado para cada
26             variable
27             # Crea un gráfico de dispersión de las observaciones en las dos
28             primeras componentes principales
29             plt.figure(figsize=(8, 6)) # Ajusta el tamaño de la figura si
30             es necesario
31             plt.scatter(componentes_principales[:, i],
32                         componentes_principales[:, j])
33
34             # Añade etiquetas a las observaciones
35             etiquetas_de_observaciones = list(datos_estandarizados.index)
36
37             for k, label in enumerate(etiquetas_de_observaciones):
38                 plt.annotate(label, (componentes_principales[k, i],
39                                     componentes_principales[k, j]))
40
41             # Dibujar líneas discontinuas que representen los ejes
42             plt.axhline(0, color='black', linestyle='--', linewidth=0.8)
43             plt.axvline(0, color='black', linestyle='--', linewidth=0.8)
44
45             # Etiquetar los ejes
46             plt.xlabel(f'Componente Principal {i + 1}')
47             plt.ylabel(f'Componente Principal {j + 1}')
48
49             # Establece el título del gráfico
```

```

39 plt.title('Gráfico de Dispersión de Observaciones y variables
40           en PCA')
41
42 # A adimos vectores que representen las correlaciones
43   escaladas entre variables y componentes
44 coeff = np.transpose(fit.components_)
45 scaled_coeff = 8 * coeff #8 = escalado utilizado, ajustar en
46   función del ejemplo
47 for var_idx in range(scaled_coeff.shape[0]):
48     plt.arrow(0, 0, scaled_coeff[var_idx, i], scaled_coeff[
49         var_idx, j], color='red', alpha=0.5)
50     plt.text(scaled_coeff[var_idx, i], scaled_coeff[var_idx, j]
51             ],
52             notas_estandarizadas.columns[var_idx], color='red', ha
53             ='center', va='center')
54
55 plt.show()
56
57 plot_pca_scatter_with_vectors(pca, notas_estandarizadas, fit.n_components,
58                               fit.components_)

```

Listing 19: Gráfico de dispersión de los individuos y variables por cada par de CPs.

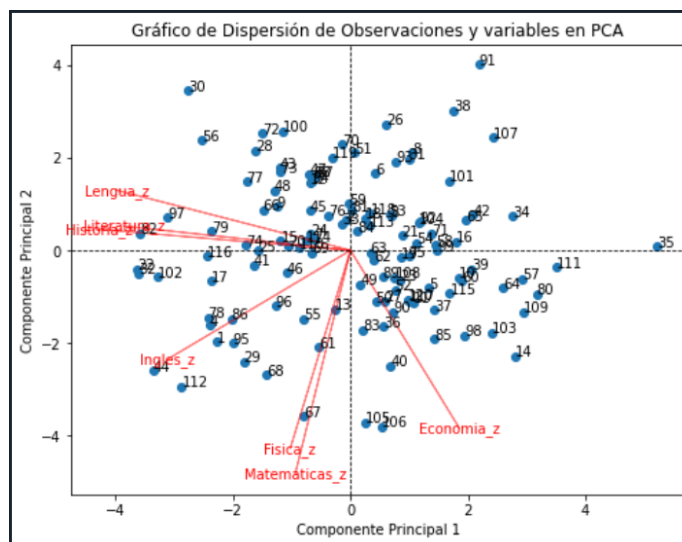


Figure 18: Dispersión de los individuos en los nuevos ejes.

Observemos, por ejemplo, que los alumnos 67 y 105 deben tener buenas notas en las asignaturas de *ciencias* mientras que en letras tienen notas sobre la media (su coordenada en la  $CP_1$  es casi cero). Destacamos también al alumno 35, que debe tener muy malas notas en las signaturas de *letras*.

### 3.2 Variables e individuos suplementarios

El ACP a menudo se utiliza para resumir y comprender las relaciones entre variables en un conjunto de datos. Sin embargo, en situaciones del mundo real, es común que surjan nuevos elementos o variables que no estaban presentes en el conjunto de datos original. Estos nuevos elementos se conocen como 'individuos suplementarios' cuando se refieren a observaciones, y 'variables suplementarias' cuando se refieren a características adicionales.

En nuestro ejemplo, en el que hemos realizado un ACP para comprender las notas escolares de los estudiantes, en un momento posterior podría incorporarse un nuevo estudiante, o decidir medir o incorporar una nueva variable que no estaba presente inicialmente en nuestro conjunto de datos. Estos nuevos elementos o variables deben ser incorporados en nuestro análisis de ACP existente para evaluar su relación con las componentes principales ya calculadas.

La inclusión de individuos y variables suplementarias en el ACP nos permite evaluar cómo se sitúan en relación con la estructura existente de las componentes principales y cómo contribuyen a la variabilidad general. Esto es esencial para mantener nuestro análisis actualizado y relevante en entornos en constante cambio.

Para aplicar esto a nuestro ejemplo particular vamos a utilizar los datos de 4 nuevos alumnos, los cuales están en el archivo *'NOTAS\_S.xlsx'*

	A	B	C	D	E	F	G	H	I
1	Alumno	Matematicas	Fisica	Lengua	Ingles	Historia	Literatura	Economia	EXTRA_ESC
2	121	6,8	5,7	5,6	7,6	8,3	7,7	5,8	Deporte
3	122	7,5	6,5	4,5	6,5	3,6	4,8	7,3	Música
4	123	4,7	3,2	3,5	5,4	4,6	5,6	5,3	Deporte
5	124	8,5	7,8	6,7	8,9	6,7	8,8	9,5	Musi-Dep

Figure 19: Individuos suplementarios.

Para tratar a estos nuevos estudiantes, debemos estandarizar sus variables, utilizando la media y la desviación estándar del conjunto de datos original. Posteriormente, los añadimos a nuestra base de datos estandarizada.

```
1 # Cargar un archivo Excel llamado 'notas.xlsx' en un DataFrame llamado
   notas.
2 notas_S = pd.read_excel('notas_S.xlsx')
3
4 # Establecer la columna 'Alumno' como índice del DataFrame notas y
   eliminarla.
5 notas_S = notas_S.set_index('Alumno', drop=True)
6
7 # Guarda la variable el índice y 'EXTRA_ESC' en un dataframe
8 extra_S = notas_S.iloc[:, [7]]
9
10 # Elimina la variable 'EXTRA_ESC' del DataFrame 'notas'.
11 notas_S = notas_S.drop(notas_S.columns[7], axis=1)
12
13 # Calcular la media y la desviación estándar de 'notas'
14 media_notas = notas_S.mean()
15 desviacion_estandar_notas = notas_S.std()
16
```



```

17 # Estandarizar 'notas_S' utilizando la media y la desviación estándar de '
    notas'
18 notas_S_estandarizadas = pd.DataFrame(((notas_S - media_notas) /
    desviacion_estandar_notas))
19
20 notas_S_estandarizadas.columns = ['{}_z'.format(variable) for variable in
    variables]
21
22 # Agregar las observaciones estandarizadas a 'notas'
23 notas_sup = pd.concat([notas_estandarizadas, notas_S_estandarizadas])

```

Listing 20: Carga y tratamiento de observaciones complementarias.

Una vez tenemos las observaciones complementarias estandarizadas y agregadas a nuestra base de datos, procedemos a calcular sus CP.

```

1
2 # Calcular las componentes principales para el conjunto de datos combinado
3 componentes_principales_sup = pca.transform(notas_sup)
4
5 # Calcular las componentes principales para el conjunto de datos combinado
6 # y renombra las componentes
7 resultados_pca_sup = pd.DataFrame(fit.transform(notas_sup),
8                                   columns=['Componente {}'.format(i) for i in
9                                           range(1, fit.n_components_+1)],
                                   index=notas_sup.index)

```

Listing 21: Cálculo de las CP para las observaciones complementarias.

Podríamos representar a estos nuevos alumnos en los nuevos ejes, de la misma manera que lo hicimos sin ellos.

```

1
2 # Representacion observaciones + suplementarios
3 plot_pca_scatter(pca, notas_sup, fit.n_components)

```

Listing 22: Representacion observaciones + suplementarios.

A partir de lo observado en la Figura 20 podemos intuir, por ejemplo, que la observación 124 ha sacado notas (en general) por encima de la media, tanto en *letras* como en *ciencias*. Todo lo contrario esperaríamos de la observación o alumno 123.

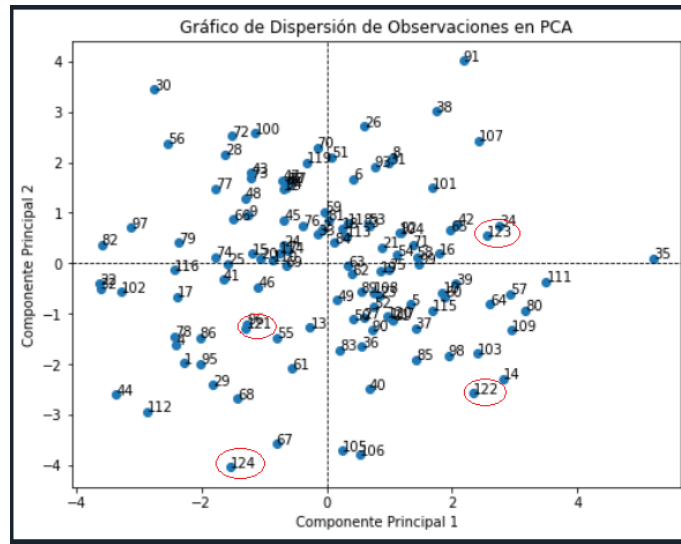


Figure 20: Dispersión observaciones y suplementarios.

Por último vamos a considerar la posibilidad de añadir una nueva variable categórica al análisis, como sabemos, las variables categóricas no pueden pertenecer al ACP, pero, ¿Cómo se relacionaría con dicho análisis?

Cuando incorporamos una variable categórica a un ACP previamente realizado, estamos enriqueciendo nuestra comprensión de los datos al considerar la influencia de categorías adicionales en la estructura de los datos. En este contexto, representar los centroides de estas categorías en los nuevos ejes de las componentes principales se convierte en una estrategia esencial. Los centroides actúan como puntos de referencia que resumen la tendencia central de cada categoría en el espacio de las CP. Este enfoque nos permite visualizar cómo las distintas categorías se relacionan entre sí y con las CP, lo que puede revelar patrones y relaciones más sutiles en los datos que antes no eran evidentes. En esencia, al considerar variables categóricas y representar sus centroides en el ACP, estamos agregando una dimensión adicional a nuestra comprensión de la estructura subyacente de los datos.

Volviendo a nuestro ejemplo, vamos a retomar la variable categórica *EXTRA\_ESC* para ver si hay algún patrón visible entre la actividad extraescolar de los alumnos y las CP.

```

1 # Incorporamos la variable categórica "EXTRA_ESC" en los datos
2 notas_componentes_sup= pd.concat([notas_sup, resultados_pca_sup], axis=1)
3
4
5 extra_sup = pd.concat([extra, extra_S], axis=0)
6 notas_componentes_sup_extra= pd.concat([notas_componentes_sup,
7                                         extra_sup], axis=1)

```

Listing 23: Representacion observaciones + suplementarios.

Para ver como se relacionan las categorías de la nueva variable, vamos a representar los centroides de dichas categorías calculados por cada par de ejes o CP (en nuestro ejemplo, solo dos).

```

1
2
3 def plot_pca_scatter_with_categories(datos_componentes_sup_var,
4   componentes_principales_sup, n_components):
5     """
6     Genera gráficos de dispersión de observaciones en pares de componentes
7     principales seleccionados con categorías.
8
9     Args:
10        datos_componentes_sup_var (pd.DataFrame): DataFrame que contiene
11        las categorías.
12        componentes_principales_sup (np.ndarray): Matriz de componentes
13        principales.
14        n_components (int): Número de componentes principales seleccionadas
15
16    """
17    # Obtener las categorías únicas
18    categorias = datos_componentes_sup_var["EXTRA_ESC"].unique() #Modificar
19    por el nombre de la variable categórica
20
21    for i in range(n_components):
22        for j in range(i + 1, n_components): # Evitar pares duplicados
23            # Crear un gráfico de dispersión de las observaciones en las
24            dos primeras componentes principales
25            plt.figure(figsize=(8, 6)) # Ajustar la dimensió de la figura
26            si es necesario
27            plt.scatter(componentes_principales_sup[:, i],
28                componentes_principales_sup[:, j])
29
30            for categoria in categorias:
31                # Filtrar las observaciones por categoría
32                observaciones_categoria = componentes_principales_sup[
33                    datos_componentes_sup_var["EXTRA_ESC"] == categoria]
34                # Calcular el centroide de la categoría
35                centroide = np.mean(observaciones_categoria, axis=0)
36                plt.scatter(centroide[i], centroide[j], label=categoria, s
37                    =100, marker='o')
38
39            # Mostrar etiquetas a las observaciones
40            etiquetas_de_observaciones = list(datos_componentes_sup_var.
41                index)
42
43            for k, label in enumerate(etiquetas_de_observaciones):
44                plt.annotate(label, (componentes_principales_sup[k, i],
45                    componentes_principales_sup[k, j]))
46                # Dibujar líneas discontinuas que representen los ejes
47
48            # Dibujar líneas discontinuas que representen los ejes
49            plt.axhline(0, color='black', linestyle='--', linewidth=0.8)

```

```

37 plt.axvline(0, color='black', linestyle='--', linewidth=0.8)
38
39 # Etiquetar los ejes
40 plt.xlabel(f'Componente Principal {i + 1}')
41 plt.ylabel(f'Componente Principal {j + 1}')
42
43 # Establecer el título del gráfico
44 plt.title('Gráfico de Dispersión de Observaciones en PCA')
45
46 # Mostrar la leyenda para las categorías
47 plt.legend()
48 plt.show()
49
50 plot_pca_scatter_with_categories(notas_componentes_sup_extra,
    componentes_principales_sup, fit.n_components)

```

Listing 24: Representacion observaciones + centroides.

Podemos intuir a partir de la Figura 21 que los estudiantes que practican tanto música como deporte en las actividades extraescolares tienen en media valores más altos tanto en las asignaturas de *letras* como en *ciencias*, por otro lado, los estudiantes que solo practican música tienen, en media, valores próximos al promedio en las variables de *ciencias* pero algo superiores en las de *letras*. Misma situación para los que solo realizan deporte, pero con notas en promedio inferiores a la media en las asignaturas de *letras*.

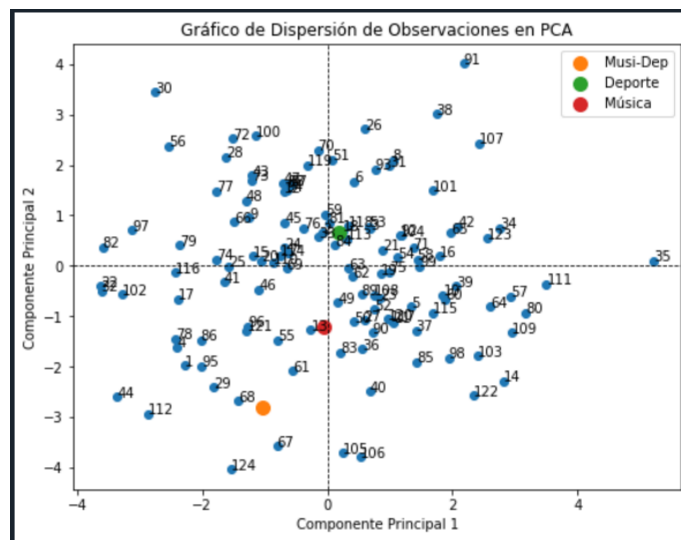


Figure 21: Dispersión centroides.