

Type *Markdown* and LaTeX: α^2

Actividad Semanal -- 5 Repaso Transformación y reducción de dimensiones

CIENCIA Y ANALÍTICA DE DATOS

Profesor: Jobish Vallikavungal Devassia

ALBERTO NIEVES CISNEROS

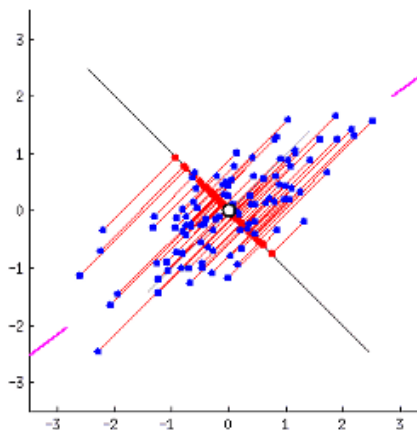
Bienvenido al notebook ¶

#Repaso de Reducción de dimensiones El objetivo es que entendamos de una manera visual, que es lo que pasa cuando nosotros seleccionamos cierto número de componentes principales o % de variabilidad de una base de datos.

Primero entenderemos, que pasa adentro de PCA que se basa en lo siguiente a grandes razgos:

Análisis de Componentes Principales

El análisis de datos multivariados involucra determinar transformaciones lineales que ayuden a entender las relaciones entre las características importantes de los datos. La idea central del Análisis de Componentes Principales (PCA) es reducir las dimensiones de un conjunto de datos que presenta variaciones correlacionadas, reteniendo una buena proporción de la variación presente en dicho conjunto. Esto se logra obteniendo la transformación a un nuevo conjunto de variables: los componentes principales (PC). Cada PC es una combinación lineal con máxima varianza en dirección ortogonal a los demás PC.



Para entender un poco más de PCA y SVD, visita el siguiente link: *Truco: Prueba entrar con tu cuenta del tec :)*

<https://towardsdatascience.com/pca-and-svd-explained-with-numpy-5d13b0d2a4d8>
(<https://towardsdatascience.com/pca-and-svd-explained-with-numpy-5d13b0d2a4d8>)

Basicamente, vamos a seguir los siguientes pasos:

1. Obtener la covarianza. OJO: X tiene sus datos centrados :)

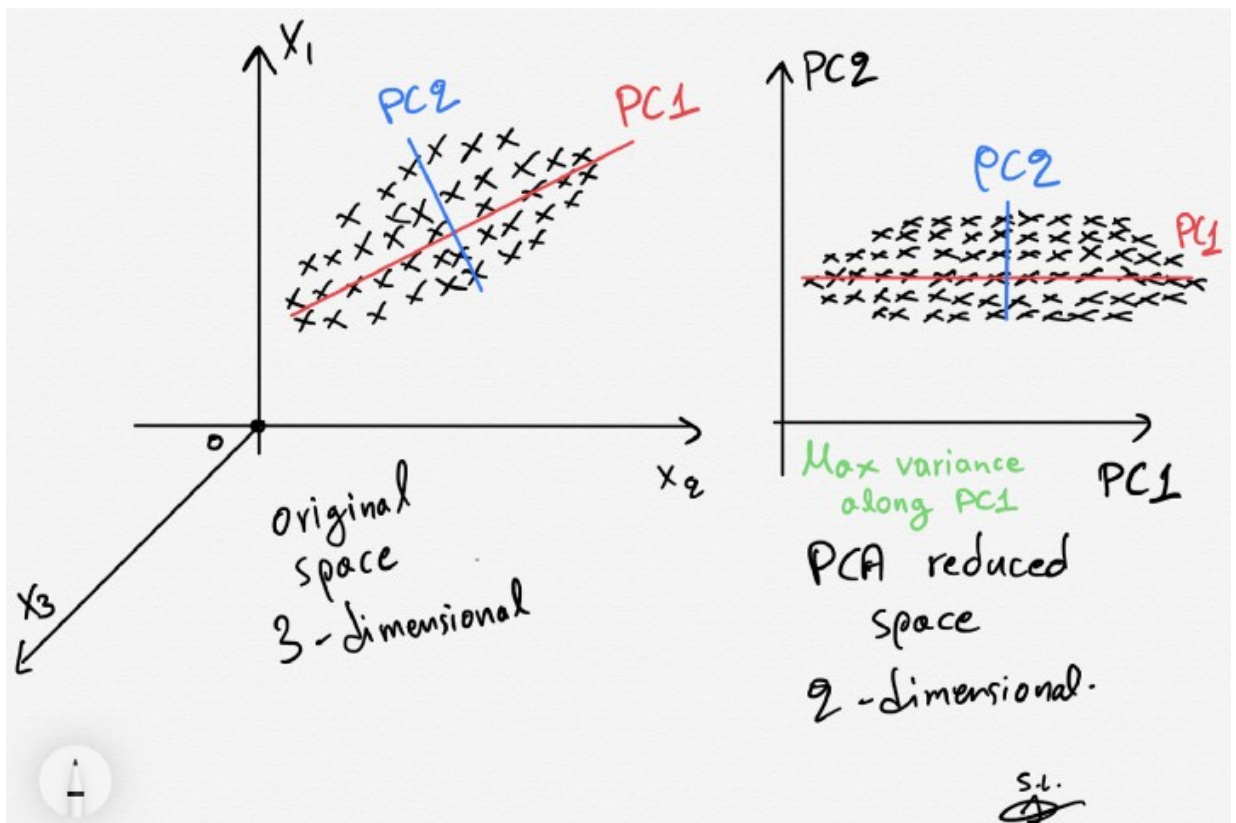
$$C = \frac{X^T X}{n - 1}$$

2. Los componentes principales se van a obtener de la eigen descomposicion de la matriz de covarianza.

$$C = W \Lambda W^{-1}$$

3. Para la reducci3n de dimensiones vamos a seleccionar k vectores de W y proyectaremos nuestros datos.

$$X_k = XW_k$$



Ejercicio 1, Descomposici3n y composici3n

Descomposici3n

Encuentra los eigenvalores y eigenvectores de las siguientes matrices

$$A = \begin{pmatrix} 3, 0, 2 \\ 3, 0, -2 \\ 0, 1, 1 \end{pmatrix} \quad A2 = \begin{pmatrix} 1, 3, 8 \\ 2, 0, 0 \\ 0, 0, 1 \end{pmatrix} \quad A3 = \begin{pmatrix} 5, 4, 0 \\ 1, 0, 1 \\ 10, 7, 1 \end{pmatrix}$$

y reconstruye la matriz original a través de las matrices $W D W^{-1}$ (OJO. Esto es lo mismo de la ecuación del paso 2 solo le cambiamos la variable a la matriz diagonal)

Eigenvalores y eigenvectores

```
In [ ]: 1  ###-----EJEMPLO DE EIGENVALORES
2  import numpy as np
3  from numpy import array
4  from numpy.linalg import eig
5  # define la matriz
6  A = array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
7  print("-----Matriz original-----")
8  print(A)
9  print("-----")
10 # calcula la eigendescomposición
11 values, vectors = eig(A)
12 print(values) #D
13 print(vectors) #W
14
15 #Ejemplo de reconstrucción
16
17
18 values, vectors = np.linalg.eig(A)
19
20 W = vectors
21 Winv = np.linalg.inv(W)
22 D = np.diag(values)
23 #La matriz B tiene que dar igual a A
24 #reconstruye la matriz
25 print("-----Matriz reconstruida-----")
26 # Realiza la reconstruccion de B=W*D*Winv, te da lo mismo de A?
27 #ojo, estas multiplicando matrices, no escalares ;)
28 #TU CODIGO AQUI-----
29 #B=
30 B = W@D@Winv
31 print(B)
32 print("-----")
```

-----Matriz original-----

```
[[1 2 3]
 [4 5 6]
 [7 8 9]]
```

```
[ 1.61168440e+01 -1.11684397e+00 -1.30367773e-15]
[[-0.23197069 -0.78583024  0.40824829]
 [-0.52532209 -0.08675134 -0.81649658]
 [-0.8186735   0.61232756  0.40824829]]
```

-----Matriz reconstruida-----

```
[[1. 2. 3.]
 [4. 5. 6.]
 [7. 8. 9.]]
```

```
In [ ]: 1 A = array([[3, 0, 2], [3, 0, -2], [0, 1, 1]])
        2 values, vectors = eig(A)
        3 print(values) #D
        4 print(vectors) #W
```

```
[3.54451153+0.j          0.22774424+1.82582815j 0.22774424-1.82582815j]
[[-0.80217543+0.j          -0.04746658+0.2575443j  -0.04746658-0.2575443j ]
 [-0.55571339+0.j          0.86167879+0.j          0.86167879-0.j          ]
 [-0.21839689+0.j          -0.16932106-0.40032224j -0.16932106+0.40032224j]]
```

```
In [ ]: 1 #Matriz 1
        2 A1 = array([[3,0,2],[3,0,-2],[0,1,1]])
        3 print("-----Matriz original-----")
        4 print(A1)
        5 print("-----")
        6 # calcula la eigendescomposición
        7 values1, vectors1 = eig(A1)
        8 print(values1) #D
        9 print(vectors1) #W
       10
       11 #Reconstrucción
       12
       13
       14 values1,vectors1 = np.linalg.eig(A1)
       15
       16 W1 = vectors1
       17 Winv1 = np.linalg.inv(W1)
       18 D1 = np.diag(values1)
       19 #La matriz B tiene que dar igual a A
       20 #reconstruye la matriz
       21 print("-----Matriz reconstruida-----")
       22 # Realiza la reconstruccion de B=W*D*Winv, te da lo mismo de A?
       23 #ojo, estas multiplicando matrices, no escalares ;)
       24 #TU CODIGO AQUI-----
       25 #B=
       26 B1 = np.round(np.real(W1@D1@Winv1),10)
       27 print(B1)
       28 print("-----")
       29
```

```
-----Matriz original-----
```

```
[[ 3  0  2]
 [ 3  0 -2]
 [ 0  1  1]]
```

```
-----
```

```
[3.54451153+0.j          0.22774424+1.82582815j 0.22774424-1.82582815j]
[[-0.80217543+0.j          -0.04746658+0.2575443j  -0.04746658-0.2575443j ]
 [-0.55571339+0.j          0.86167879+0.j          0.86167879-0.j          ]
 [-0.21839689+0.j          -0.16932106-0.40032224j -0.16932106+0.40032224j]]
```

```
-----Matriz reconstruida-----
```

```
[[ 3.  0.  2.]
 [ 3.  0. -2.]
 [ 0.  1.  1.]]
```

```
-----
```

In []:

```
1  #Matriz 2
2
3  A2 = array([[1, 3, 8], [2, 0, 0], [0,0, 1]])
4  print("-----Matriz original-----")
5  print(A2)
6  print("-----")
7  # calcula la eigendescomposición
8  values2, vectors2 = eig(A2)
9  print(values2) #D
10 print(vectors2) #W
11
12 #Reconstrucción
13
14
15 values2, vectors2 = np.linalg.eig(A2)
16
17 W2 = vectors2
18 Winv2 = np.linalg.inv(W2)
19 D2 = np.diag(values2)
20 #La matriz B tiene que dar igual a A
21 #reconstruye la matriz
22 print("-----Matriz reconstruida-----")
23 # Realiza la reconstruccion de B=W*D*Winv, te da lo mismo de A?
24 #ojo, estas multiplicando matrices, no escalares ;)
25 #TU CODIGO AQUI-----
26 #B=
27 B2 = W2@D2@Winv2
28 print(np.round(B2,10))
29 print("-----")
30
```

-----Matriz original-----

```
[[1 3 8]
 [2 0 0]
 [0 0 1]]
```

```
[ 3. -2.  1.]
[[ 0.83205029 -0.70710678 -0.42399915]
 [ 0.5547002   0.70710678 -0.8479983 ]
 [ 0.          0.          0.31799936]]
```

-----Matriz reconstruida-----

```
[[1. 3. 8.]
 [2. 0. 0.]
 [0. 0. 1.]]
```

In []:

```
1
2 #Matriz 3
3 A3 = array([[5,4,0], [1, 0, 1], [10, 7,1]])
4 print("-----Matriz original-----")
5 print(A3)
6 print("-----")
7 # calcula la eigendescomposición
8 values3, vectors3 = eig(A3)
9 print(values3) #D
10 print(vectors3) #W
11
12 #Reconstrucción
13
14
15 values3, vectors3 = np.linalg.eig(A3)
16
17 W3 = vectors3
18 Winv3 = np.linalg.inv(W3)
19 D3 = np.diag(values3)
20 #La matriz B tiene que dar igual a A
21 #reconstruye la matriz
22 print("-----Matriz reconstruida-----")
23 # Realiza la reconstruccion de B=W*D*Winv, te da lo mismo de A?
24 #ojo, estas multiplicando matrices, no escalares ;)
25 #TU CODIGO AQUI-----
26 #B=
27 B3 = W3@D3@Winv3
28
29 print(np.round(B3,10))
30 print("-----")
31
```

```
-----Matriz original-----
[[ 5  4  0]
 [ 1  0  1]
 [10  7  1]]
-----
[[ 6.89167094 -0.214175  -0.67749594]
 [ 0.3975395  0.55738222  0.57580768]
 [ 0.18800348 -0.72657211 -0.81728644]
 [ 0.89811861 -0.40176864 -0.02209943]]
-----Matriz reconstruida-----
[[ 5.  4. -0.]
 [ 1. -0.  1.]
 [10.  7.  1.]]
-----
```

¿Qué significa reducir dimensiones?

Esto será cuando proyectemos a ese espacio de los componentes principales pero no los seleccionemos todos, solo los más importantes y viajemos de regreso a nuestras unidades a través de una proyección.

Es decir: Unidades-PC PC-Unidades

Veámoslo gráficamente, ¿qué pasa con esa selección de los PCs y su efecto?.

Para ello usaremos Singular Value Descomposition (SVD).

Singular Value Descomposition(SVD)

Es otra descomposición que tambien nos ayudara a reducir dimensiones.

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|} \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline \text{orange} & 0 & 0 \\ \hline 0 & \text{orange} & 0 \\ \hline 0 & 0 & \text{yellow} \\ \hline 0 & 0 & 0 \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline \text{light blue} & \text{light blue} & \text{light blue} \\ \hline \text{purple} & \text{purple} & \text{purple} \\ \hline \text{pink} & \text{pink} & \text{pink} \\ \hline \end{array} \\
 \mathbf{X} & = & \mathbf{U} & \mathbf{\Sigma} & \mathbf{V}^* \\
 n \times m & & n \times n & n \times m & m \times m
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|} \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \text{teal} & \text{green} & \text{blue} & \text{green} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{teal} & \text{teal} & \text{teal} & \text{teal} \\ \hline \text{green} & \text{green} & \text{green} & \text{green} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{green} & \text{green} & \text{green} & \text{green} \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|} \hline 1 & 0 & 0 & 0 \\ \hline 0 & 1 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} \\
 \mathbf{U} & \mathbf{U}^* & = & \mathbf{I}_n
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{|c|c|c|} \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \text{light blue} & \text{purple} & \text{pink} \\ \hline \end{array} & = & \begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 1 \\ \hline \end{array} \\
 \mathbf{V} & \mathbf{V}^* & = & \mathbf{I}_m
 \end{array}$$

#Ejercicio 2 Juega con Lucy, una cisne, ayudala a encontrar cuantos valores singulares necesita para no perder calidad a través de SVD. Posteriormente usa 3 imágenes de tu preferencia y realiza la misma acción :D

A esto se le llama **compresión de imagenes** :o

Type *Markdown* and LaTeX: α^2

```

In [ ]: 1 from six.moves import urllib
        2 from PIL import Image
        3 import matplotlib.pyplot as plt
        4 import numpy as np
        5
        6 plt.style.use('classic')
        7 img = Image.open(urllib.request.urlopen('https://biblioteca.acropolis.org/wp
        8 #img = Image.open('Lucy.jpg')
        9 imggray = img.convert('LA')
       10 imgmat = np.array(list(imggray.getdata(band=0)),float)
       11
       12 print(imgmat)
       13
       14 imgmat.shape = (imggray.size[1],imggray.size[0])
       15
       16 plt.figure(figsize=(9,6))
       17 plt.imshow(imgmat,cmap='gray')
       18 plt.show()
       19 print(img)

```

[72. 73. 74. ... 48. 47. 47.]



<PIL.Image.Image image mode=LA size=1024x660 at 0x7F72A186E450>

```
In [ ]: 1 U,D,V = np.linalg.svd(imgmat)
        2 imgmat.shape
        3
```

Out[7]: (660, 1024)

```
In [ ]: 1 U.shape
```

Out[8]: (660, 660)

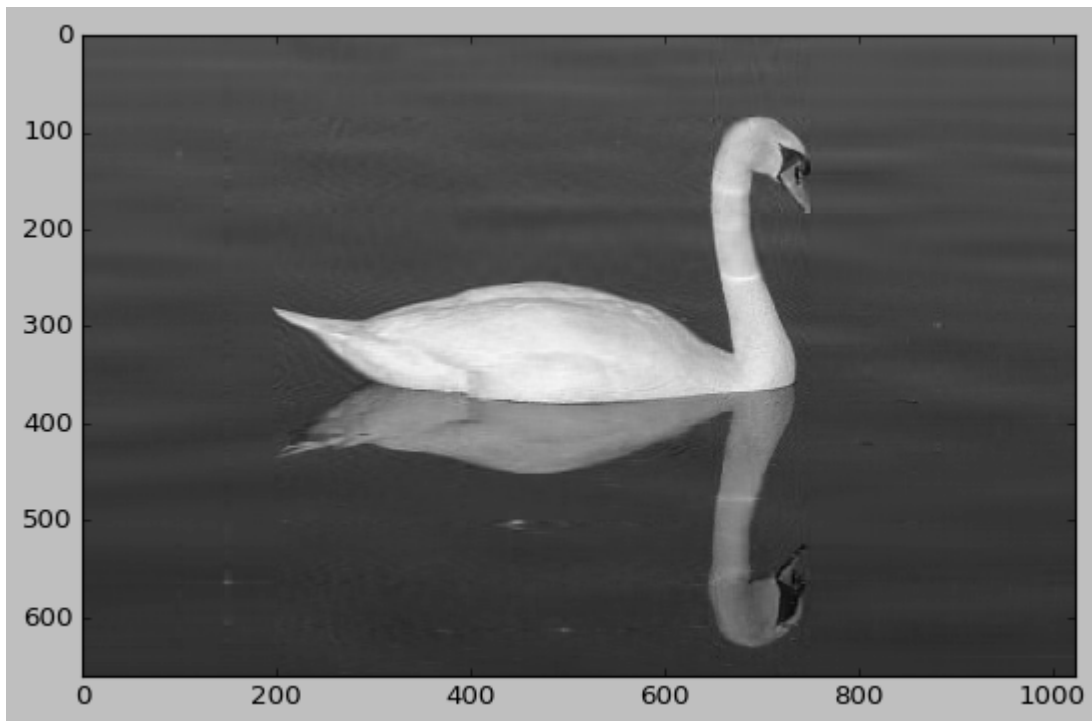
```
In [ ]: 1 V.shape
```

Out[9]: (1024, 1024)

```

In [ ]: 1 #Cuantos valores crees que son necesarios?
        2 #A=U*D*V
        3 #aqui los elegiremos-----
        4 # por las dimensiones de este caso en particular
        5 #iremos de 0-660, siendo 660 como normalmente están los datos
        6 #con 50 podemos observar que Lucy se ve casi igual, es decir conservamos aquí
        7 # realidad estaba aportando a la imagen en este caso :D por medio de la vari
        8 #juega con el valor nvalue y ve que pasa con otros valores
        9 nvalue = 50
       10 #-----
       11 reconstimg = np.matrix(U[:, :nvalue])*np.diag(D[:nvalue])*np.matrix(V[:nvalue]
       12 #ve las dimensiones de la imagen y su descomposición
       13 #660x1024= U(660X660)D(660X1024)V(1024x1024)
       14         #=U(660Xnvalues)D(nvaluesXnvalue)V(nvaluesx1024)
       15
       16         #=U(660X50)(50X50)(50X1024)
       17 plt.imshow(reconstimg, cmap='gray')
       18 plt.show()
       19 print("Felicidades la imagen está comprimida")

```



Felicidades la imagen está comprimida

¡Ahora es tu turno!, comprime 3 imágenes

```
In [ ]: 1 #imagen 1
2
3 plt.style.use('classic')
4 img = Image.open(urllib.request.urlopen('https://biblioteca.acropolis.org/wp
5 imggray = img.convert('LA')
6 imgmat = np.array(list(imggray.getdata(band=0)),float)
7
8 print(imgmat)
9
10 imgmat.shape = (imggray.size[1],imggray.size[0])
11
12 plt.figure(figsize=(9,6))
13 plt.imshow(imgmat,cmap='gray')
14 plt.show()
15 print(img)
```

[16. 17. 17. ... 151. 148. 147.]



<PIL.Image.Image image mode=LA size=1200x800 at 0x7F729DC24CD0>

```
In [ ]: 1 U,D,V = np.linalg.svd(imgmat)
        2 imgmat.shape
```

Out[12]: (800, 1200)

```
In [ ]: 1 U.shape
```

Out[13]: (800, 800)

```
In [ ]: 1 V.shape
```

Out[14]: (1200, 1200)

```
In [ ]: 1 #A=U*D*V
        2 nvalue = 70
        3 #-----
        4 reconstimg = np.matrix(U[:, :nvalue])*np.diag(D[:nvalue])*np.matrix(V[:nvalue
        5 #ve las dimensiones de la imagen y su descomposicion
        6 #800x1200= U(800,800)D(800X1200)V(1200,1200)
        7         #=U(800Xnvalues)D(nvaluesXnvalue)V(nvaluesx1200)
        8
        9         #=U(800X70)(70X70)(70X1200)
       10 plt.imshow(reconstimg,cmap='gray')
       11 plt.show()
       12 print("Felicidades la imagen está comprimida")
```



Felicidades la imagen está comprimida

```
In [ ]: 1 #imagen 2
2
3 plt.style.use('classic')
4 img = Image.open(urllib.request.urlopen('https://biblioteca.acropolis.org/wp
5 imggray = img.convert('LA')
6 imgmat = np.array(list(imggray.getdata(band=0)),float)
7
8 print(imgmat)
9
10 imgmat.shape = (imggray.size[1],imggray.size[0])
11
12 plt.figure(figsize=(9,6))
13 plt.imshow(imgmat,cmap='gray')
14 plt.show()
15 print(img)
```

[34. 34. 34. ... 88. 92. 95.]



<PIL.Image.Image image mode=LA size=1200x800 at 0x7F729DC24B50>

```
In [ ]: 1 U,D,V = np.linalg.svd(imgmat)
2 imgmat.shape
```

Out[17]: (800, 1200)

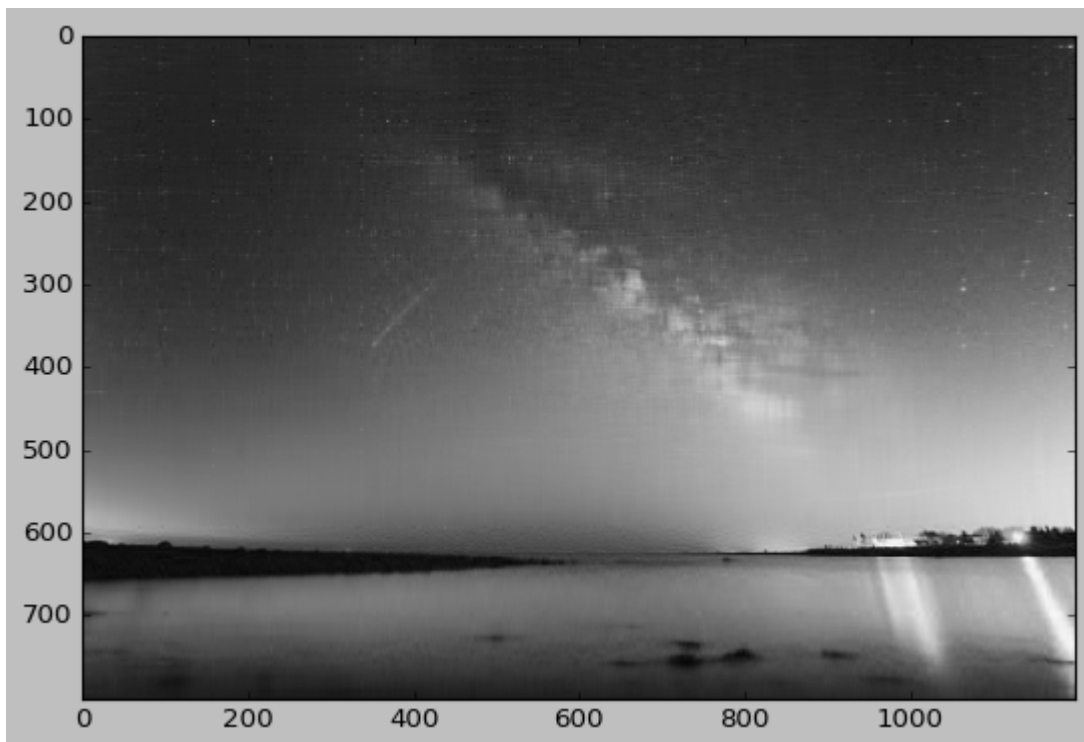
```
In [ ]: 1 U.shape
```

```
Out[18]: (800, 800)
```

```
In [ ]: 1 V.shape
```

```
Out[19]: (1200, 1200)
```

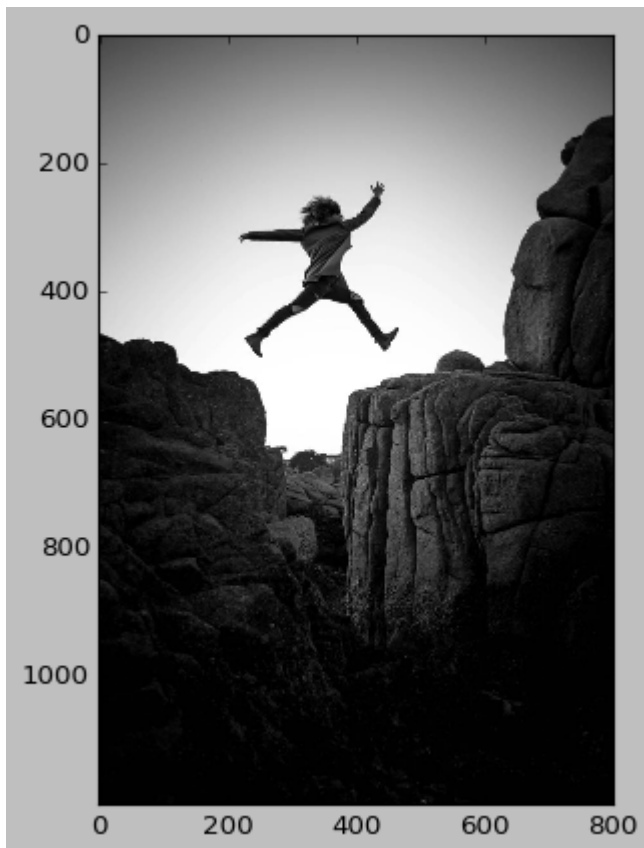
```
In [ ]: 1 nvalue = 40
2 #-----
3 reconstimg = np.matrix(U[:, :nvalue])*np.diag(D[:nvalue])*np.matrix(V[:nvalue
4 #ve las dimensiones de la imagen y su descomposicion
5 #800x1200= U(800,800)D(800X1200)V(1200,1200)
6         #=U(660Xnvalues)D(nvaluesXnvalue)V(nvaluesx1024)
7
8         #=U(800X40)(40X40)(40X1200)
9 plt.imshow(reconstimg,cmap='gray')
10 plt.show()
11 print("Felicidades la imagen está comprimida")
```



Felicidades la imagen está comprimida

```
In [ ]: 1 #imagen 3
2 plt.style.use('classic')
3 img = Image.open(urllib.request.urlopen('https://biblioteca.acropolis.org/wp
4 #img = Image.open('Lucy.jpg')
5 imggray = img.convert('LA')
6 imgmat = np.array(list(imggray.getdata(band=0)),float)
7
8 print(imgmat)
9
10 imgmat.shape = (imggray.size[1],imggray.size[0])
11
12 plt.figure(figsize=(9,6))
13 plt.imshow(imgmat,cmap='gray')
14 plt.show()
15 print(img)
16
17
```

[71. 71. 71. ... 1. 1. 1.]



<PIL.Image.Image image mode=LA size=801x1200 at 0x7F729F160950>

```
In [ ]: 1 U,D,V = np.linalg.svd(imgmat)
        2 imgmat.shape
```

Out[27]: (1200, 801)

```
In [ ]: 1 U.shape
```

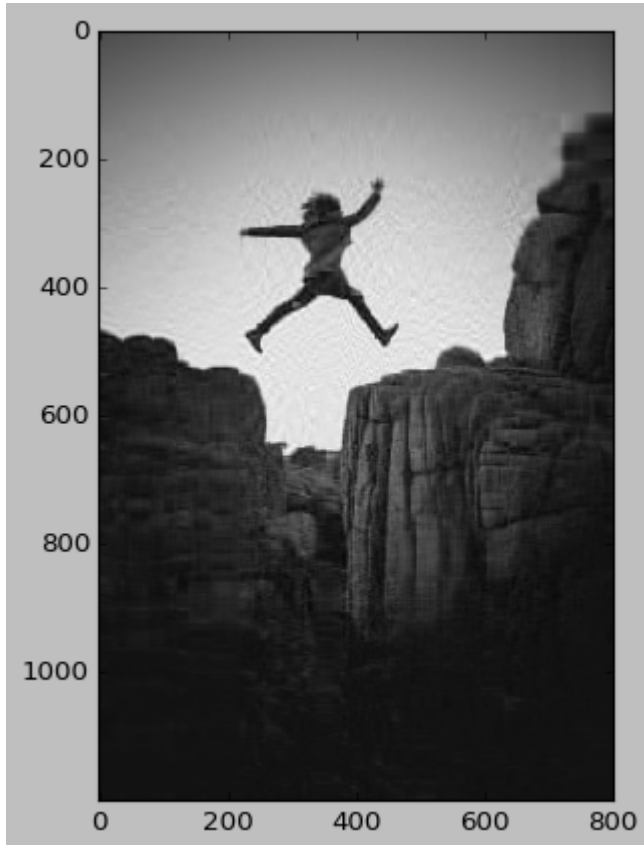
Out[28]: (1200, 1200)

```
In [ ]: 1 V.shape
```

Out[29]: (801, 801)

In []:

```
1
2 nvalue = 40
3 #-----
4 reconstimg = np.matrix(U[:, :nvalue])*np.diag(D[:nvalue])*np.matrix(V[:nvalue
5 #ve las dimensiones de la imagen y su descomposicion
6 #800x1200= U(800, 800)D(800X1200)V(1200, 1200)
7         #=U(800Xnvalues)D(nvaluesXnvalue)V(nvaluesx1200)
8
9         #=U(800X40)(40X40)(40X1200)
10 plt.imshow(reconstimg, cmap='gray')
11 plt.show()
12 print("Felicidades la imagen está comprimida")
```



Felicidades la imagen está comprimida

Ejercicio 3

Feature importances

Para este ejercicio, te pediremos que sigas el tutorial de la siguiente pagina:

<https://towardsdatascience.com/pca-clearly-explained-how-when-why-to-use-it-and-feature-importance-a-guide-in-python-7c274582c37e> (<https://towardsdatascience.com/pca-clearly-explained-how-when-why-to-use-it-and-feature-importance-a-guide-in-python-7c274582c37e>)

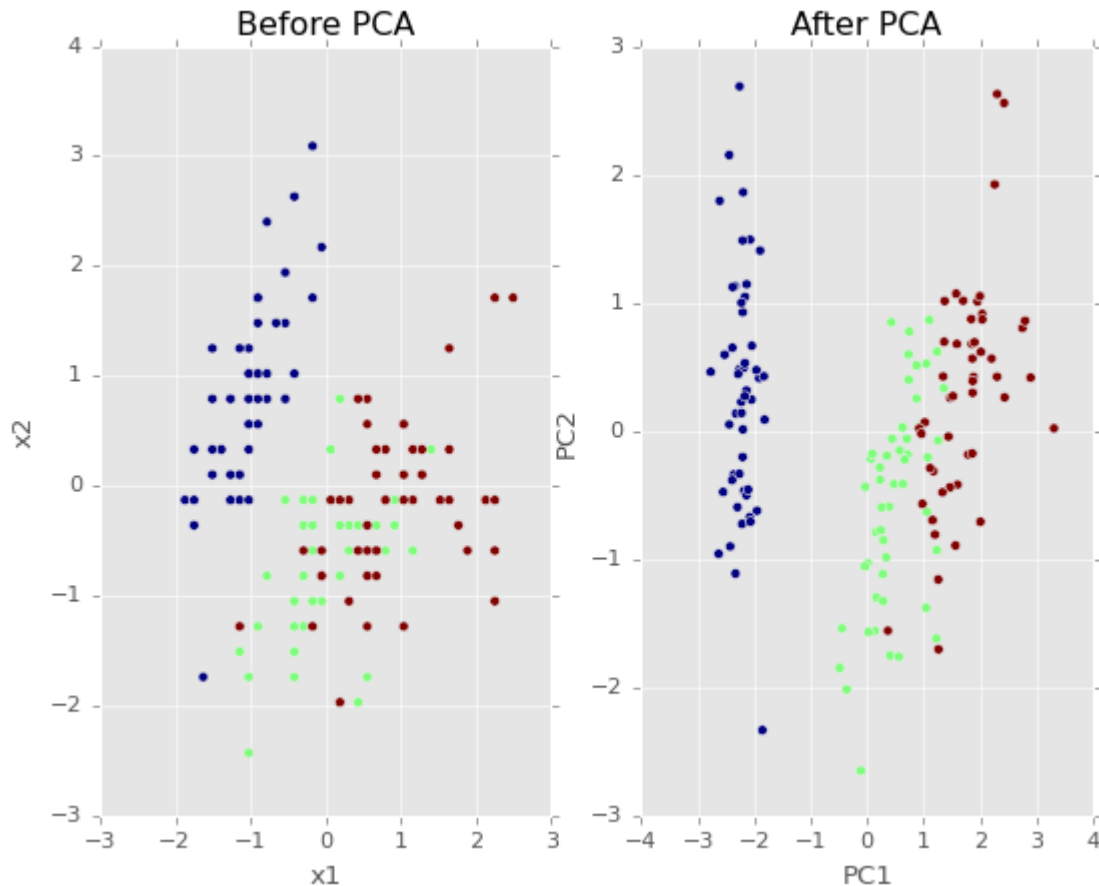
In []:

```
1  #tu codigo aqui
2  #Se vió el siguiente código en el ejercicio
3  import numpy as np
4  import matplotlib.pyplot as plt
5  from sklearn import datasets
6  from sklearn.decomposition import PCA
7  import pandas as pd
8  from sklearn.preprocessing import StandardScaler
9  plt.style.use('ggplot')
10 # Load the data
11 iris = datasets.load_iris()
12 X = iris.data
13 y = iris.target
14 # Z-score the features
15 scaler = StandardScaler()
16 scaler.fit(X)
17 X = scaler.transform(X)
18 # The PCA model
19 pca = PCA(n_components=2) # estimate only 2 PCs
20 X_new = pca.fit_transform(X) # project the original data into the PCA space
```

```

In [ ]: 1 fig, axes = plt.subplots(1,2)
        2 axes[0].scatter(X[:,0], X[:,1], c=y)
        3 axes[0].set_xlabel('x1')
        4 axes[0].set_ylabel('x2')
        5 axes[0].set_title('Before PCA')
        6 axes[1].scatter(X_new[:,0], X_new[:,1], c=y)
        7 axes[1].set_xlabel('PC1')
        8 axes[1].set_ylabel('PC2')
        9 axes[1].set_title('After PCA')
       10 plt.show()

```



```

In [ ]: 1 print(pca.explained_variance_ratio_)
        2

```

[0.72962445 0.22850762]

```

In [ ]: 1 np.cov(X_new.T)
        2

```

Out[39]: array([[2.93808505e+00, 5.33928780e-16],
[5.33928780e-16, 9.20164904e-01]])

```

In [ ]: 1 pca.explained_variance_
        2

```

Out[38]: array([2.93808505, 0.9201649])

In []:

```
1 print(abs( pca.components_ ))  
2
```

```
[[0.52106591 0.26934744 0.5804131 0.56485654]  
 [0.37741762 0.92329566 0.02449161 0.06694199]]
```

¿Qué es feature importance y para que nos sirve? la importancia de cada 'feature' está reflejada por la magnitud de los valores correspondientes en los eigenvectores (una más alta magnitud implica una mayor importancia) En resumen, se buscan los valores absolutos de los eigenvectores de los componentes correspondientes a los k mayores eigenvalores. A medida que estos valores absolutos sean mayores, mayor es la contribución de un feature a un componente principal

¿Qué hallazgos fueron los más relevantes durante el análisis del ejercicio? Se realizó una explicación de qué son y que función tienen el análisis de componentes principales y el análisis de SVD

¿Dónde lo aplicarías o te sería de utilidad este conocimiento? Este análisis se puede emplear de manera no supervisada, para realizar una reducción de dimensionalidad, en caso de que se cuenten con muchas variables 'features'. Estos se puede realizar para acelerar el entrenamiento de un modelo. también para visuzalización de datos

Describe lo relevante del ejercicio y que descubriste de las variables analizadas.

Se muestra paso a paso, el procedimiento para realizar el análisis de componentes principales (PCA), de detallan sus aplicaciones para machine learning, también se explora su aplicación en python a través de la librería sklearn.