

Ejercicios resueltos de clase 2

In [2]:

```
import numpy as np
import matplotlib.pyplot as plt
```

Ejercicio 1

Obtener para cada fila en X, el índice de la fila en C con distancia euclídea más pequeña. Es decir, decir para cada fila en X a qué cluster pertenece en C.

In [3]:

```

# Ejercicio 1
def get_distances(X, C):
    """Obtiene para cada vector de X la distancia a cada vector de C.
    X -- Array. Cada fila es un vector
    C -- Array. Cada fila es un vector
    """

    """
    Ignorando el cuadrado y raíz por el momento, quiero calcular distances = X -
    C, tal que:

    ~~~
    D = [
        X[0] - C[0],
        X[0] - C[1],
        X[1] - C[0],
        X[1] - C[1],
        X[2] - C[0],
        X[2] - C[1]
    ]
    ~~~

    siendo:

    ~~~
    X[0] - C[0] = [1,2,3] - [1,0,0] = [0,-2,3],
    X[0] - C[1] = [1,2,3] - [1,0,0] = [0,2,3],
    ...
    ~~~

    Según https://scipy.github.io/old-wiki/pages/EricksBroadcastingDoc:

    > In order to broadcast, the size of the trailing axes for both arrays
    > in an operation must either be the same size or one of them must be one.

    Si intento restar directamente (3,3) - (2,3) Falla la regla de broadcast por
    que 3!=2
    y ninguno de los dos es 1.
    Pero si en cambio pruebo con (3,3)-(2,1,3) la regla de broadcasting permite
    calcular
    (1,3)-(1,3) sobre el eje 0 (filas):
        Operation Axis (0) Size      Trailing Axis Size
        3                          (3) automatically reshaped to (1,3)
        2                          (1,3)

    Nota: interpreto que es el (1,3) de dimensión 2 el que lleva a expandir el
    (3) a (1,3).
    """
    expanded_C = C[:, None]
    return np.sqrt(np.sum((expanded_C - X) ** 2, axis=-1))

```

In [4]:

```
# Ejercicio 1 - Test

X = np.array([
    [1,2,3],
    [4,5,6],
    [7,8,9]
])
C = np.array([
    [1,0,0],
    [0,1,1]
])
distances = get_distances(X,C)
```

Ejercicio 2

Obtener para cada fila en X, el índice de la fila en C con distancia euclídea más pequeña. Es decir, decir para cada fila en X a qué cluster pertenece en C. Por ejemplo, si el resultado anterior fue:

```
[[ 3.60555128  8.36660027 13.45362405]
 [ 2.44948974  7.54983444 12.72792206]]
```

El programa debería devolver [1, 1, 1] Hint: utilizar np.argmin

In [5]:

```
# Ejercicio 2
def get_nearest(distances):
    return np.argmin(distances,axis=0)
```

In [6]:

```
# Ejercicio 2 - Example
nearest = get_nearest(distances)
nearest
```

Out[6]:

```
array([1, 1, 1])
```

Nota: se puso el ejercicio 4 antes del 3 porque en el ejercicio 3 se usan funciones del ejercicio 4.

Ejercicio 4

Utilizar numpy para crear datos clusterizados A/B en 4 dimensiones.

Hint:

- Definir una matriz con centroides $[1,0,0,0]$ y $[0,1,0,0]$
- Utilizar una constante para separar o alejar los centroides entre si.
- Utilizar `np.repeat` para crear $n/2$ muestras de cada centroide.
- Sumar a cada centroide un vector aleatorio normal i.i.d. con media 0 y desvío (`np.random.normal`).
- Armar un arreglo que tenga n enteros indicado si la muestra pertenece a A o a B.

In [7]:

```
# Ejercicio 4

def build_cluster(n_samples, inv_overlap):
    """
    Genera muestras pertenecientes a dos clusters.
    n_samples -- Cantidad de muestras.
    inv_overlap -- Distancia de separación.
    """
    centroids = np.array([
        [1,0,0,0],
        [0,1,0,0],
    ], dtype=np.float32)
    centroids = centroids * inv_overlap
    data = np.repeat(centroids, n_samples / 2, axis=0)
    normal_noise = np.random.normal(loc=0, scale=1, size=(n_samples, 4))
    data = data + normal_noise
    cluster_ids = np.array([[0],[1],])
    cluster_ids = np.repeat(cluster_ids, n_samples / 2, axis=0)
    return data, cluster_ids
```

In [8]:

```
# Ejercicio 4 - Prueba de llamada (se usa en Ejercicio 12)
import matplotlib.pyplot as plt

data, cluster_ids = build_cluster(10, 1)
data
```

Out[8]:

```
array([[ 0.38713788,  0.40398896, -0.76608121,  1.27780392],
       [-0.66766558,  2.70097768, -0.7552887 , -1.0600832 ],
       [ 0.45757021,  0.91884714,  0.05339289, -0.19864588],
       [ 2.65521735, -0.75249413, -0.65006635,  1.38042781],
       [ 0.32082762, -1.30713403,  0.43342317, -0.70218285],
       [ 0.29218512,  0.57732234,  1.83623057, -1.54619309],
       [ 0.23419074,  2.18935715, -0.14258527, -1.12147595],
       [-0.8003339 ,  3.11161213, -0.02496614, -0.27303787],
       [-0.89941854, -0.35262111,  0.12247909, -1.04698691],
       [ 0.39131439,  2.00560484, -2.13241275,  0.98988417]])
```

Ejercicio 3

K-means es uno de los algoritmos más básicos en Machine Learning no supervisado. Es un algoritmo de clusterización, que agrupa los datos que comparten características similares. Recordemos que entendemos datos como n realizaciones del vector aleatorio X . El algoritmo K-means funcione de la siguiente manera:

1. El usuario selecciona la cantidad de clusters a crear (n).
2. Se seleccionan n elementos aleatorios de X como posiciones iniciales de los centroides C .
3. Se calcula la distancia entre todos los puntos en X y todos los puntos en C .
4. Para cada punto en X se selecciona el centroide más cercano de C .
5. Se recalculan los centroides C a partir de usar las filas de X que pertenecen a cada centroide.
6. Se itera entre 3 y 5 una cantidad fija de veces o hasta que la posición de los centroides no cambie.

Implementar la función `def k_means(X, n)` de manera tal que al finalizar devuelva la posición de los centroides y a que cluster pertenece cada fila de X . Hint: para (2) utilizar funciones de `np.random`, para (3) y (4) usar los ejercicios anteriores, para (5) es válido utilizar un `for`. Iterar 10 veces entre (3) y (5).

In [9]:

```
# Ejercicio 3

MAX_ITERATIONS = 10

def k_means(X, n_clusters, max_iterations=MAX_ITERATIONS):
    centroids = np.eye(n_clusters, X.shape[1])
    #print(centroids)
    for i in range(max_iterations):
        #print("Iteration # {}".format(i))
        centroids, cluster_ids = k_means_loop(X, centroids)
        #print(centroids)
    return centroids, cluster_ids

def k_means_loop(X, centroids):
    # find labels for rows in X based in centroids values
    expanded_centroids = centroids[:, None]
    distances = np.sqrt(np.sum((expanded_centroids - X) ** 2, axis=2))
    arg_min = np.argmin(distances, axis=0)
    # recompute centroids
    for i in range(centroids.shape[0]):
        centroids[i] = np.mean(X[arg_min == i, :], axis=0)
    return centroids, arg_min
```

In [10]:

```
# Ejercicio 3 - Ejemplo de llamada (se prueba en Ejercicio 12)
data, cluster_ids = build_cluster(100, 0.1)
centroids, cluster_ids = k_means(data, 4)
```

Ejercicio 5

Utilizar numpy para simular una exponencial de parámetro lambda.

Hint:

- Hacer una función que genere n muestras de la variable aleatoria X
- Utilizar el resultado obtenido en la diapositiva anterior
- Utilizar `np.random.uniform`

In [11]:

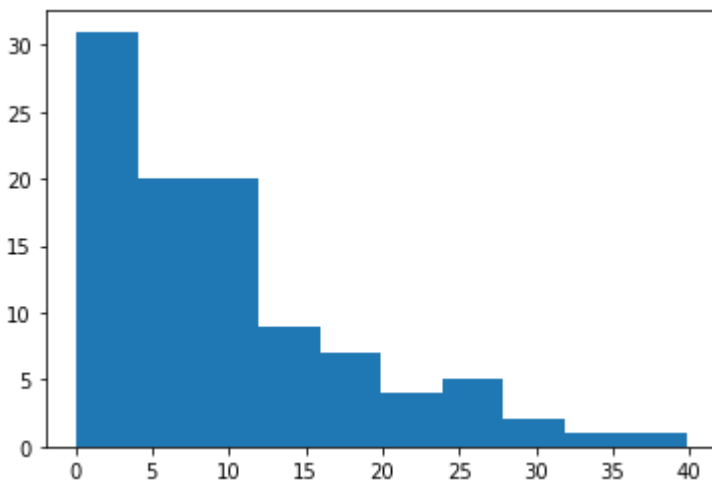
```
# Ejercicio 5
def exponential_random_variable(lambda_param, size):
    u = np.random.uniform(low=0.0, high=1.0, size=size)
    return (-1 / lambda_param) * np.log(1 - u)
```

In [12]:

```
# Ejercicio 5 - Prueba de llamada (se usa en ejercicio 12)
data = exponential_random_variable(0.1,100)
plt.hist(data)
```

Out[12]:

```
(array([31., 20., 20., 9., 7., 4., 5., 2., 1., 1.]),
 array([ 0.08647035,  4.0532191 ,  8.01996784, 11.98671659, 15.95346
534,
        19.92021408, 23.88696283, 27.85371157, 31.82046032, 35.78720
907,
        39.75395781])),
<a list of 10 Patch objects>)
```



Ejercicio 6

Calcular la inversa generalizada y simular. Sabiendo que la variable aleatoria es continua, la función de densidad de probabilidad se define según:

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f(t)dt$$

Para una variable aleatoria con función de densidad de probabilidad:

$$f_X(x) = 3x^2 \{0 < x < 1\}$$

Obtener la inversa generalizada y utilizar numpy para simular n muestras.

$$U = F_X(x) = x^3 \therefore x = U^{1/3}$$

In [13]:

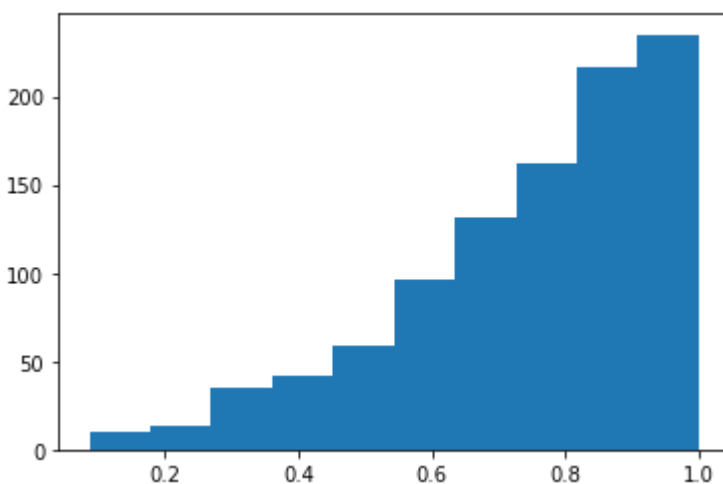
```
# Ejercicio 6
def cubic_random_variable(size):
    u = np.random.uniform(low=0.0, high=1.0, size=size)
    return u**(1/3)
```

In [14]:

```
# Ejercicio 6 - Test
dataset = cubic_random_variable(1000)
plt.hist(dataset)
```

Out[14]:

```
(array([ 10.,  13.,  35.,  42.,  59.,  96., 131., 162., 217., 23
5.]),
 array([0.08718682, 0.17841954, 0.26965225, 0.36088497, 0.45211768,
        0.5433504 , 0.63458312, 0.72581583, 0.81704855, 0.90828126,
        0.99951398]),
 <a list of 10 Patch objects>)
```



Ejercicio 7

Dado un dataset X de n muestras y m columnas, implementar un método en numpy para normalizar con z-score.

$$z = \frac{x - \mu}{\sigma}$$

Pueden utilizar np.mean() y np.std()

In [15]:

```
# Ejercicio 7
def zscore_norm(dataset):
    """ Normaliza un dataset c/ z-score
    -- dataset: mxn, m=muestras, n=features
    """
    feat_mean = np.mean(data,axis=0)
    feat_std = np.std(data,axis=0)
    return (data-feat_mean)/feat_std
```

In [16]:

```
# Ejercicio 7 - Example
data, _ = build_cluster(10, 1)
zscore_norm(data)
```

Out[16]:

```
array([[ 0.95389281,  0.38140895, -2.47113929, -1.05692773],
       [ 1.30845696, -1.88276473,  1.08697695,  0.77928785],
       [ 0.29639063, -0.51211016,  1.39064721,  0.717264  ],
       [ 0.11922888, -0.72601899,  0.04186609,  0.06776867],
       [-0.11163188,  1.20820222,  0.64533271, -0.48435756],
       [ 0.56138542, -0.5228208 , -0.443059 , -0.9166429 ],
       [ 0.91143724,  0.46309848, -0.13219506,  0.20892988],
       [-1.34291999,  1.79388043, -0.34534301, -1.83041657],
       [-0.76530638, -0.51023057,  0.1364949 ,  1.50733833],
       [-1.93093368,  0.30735517,  0.0904185 ,  1.00775602]])
```

Ejercicio 8

Siguiendo los pasos del slide anterior, se requiere utilizar numpy para calcular PCA del dataset de entrada X utilizando las 2 componentes más importantes.

```
x = np.array(
    [
        [0.4, 4800, 5.5], [0.7, 12104, 5.2],
        [1, 12500, 5.5], [1.5, 7002, 4.0]
    ])

```

Al finalizar la implementación en numpy corroborar obtener los mismos resultados que utilizando el código de la librería scikit-learn. Para comparar las matrices escribir un tests usando np.testing.assert_allclose

In [17]:

```
# Ejercicio 8
def my_pca(x):
    x2 = (x - x.mean(axis=0))
    cov_1 = np.cov(x2.T)
    w, v = np.linalg.eig(cov_1)
    idx = w.argsort()[::-1]
    w = w[idx]
    v = v[:,idx]
    return np.matmul(x2, v[:, :2])
```

In [18]:

```
# Ejercicio 8 - Test
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

x = np.array([ [0.4, 4800, 5.5], [0.7, 12104, 5.2], [1, 12500, 5.5], [1.5, 7002, 4.0] ])
x_pca_mine = my_pca(x)

pca = PCA(n_components=2)
x_std = StandardScaler(with_std=False).fit_transform(x)
x_pca_sk = pca.fit_transform(x_std)

x_pca_mine.shape, x_pca_sk.shape
```

Out[18]:

```
((4, 2), (4, 2))
```

In [19]:

```
np.testing.assert_allclose( x_pca_mine, x_pca_sk, rtol=2, atol=1 ) # FIXME, está bien tener que usar tanta tolerancia?
```

Ejercicio 9

Dado un dataset, hacer una funcion que utilizando numpy filtre las columnas y las filas que tienen NaNs/

In [20]:

```
# Función helper para poblar dataset con NaNs
def corrupt_dataset(dataset, percent, feat_range):
    """Corrompe un dataset con NaNs
        -- percent: porcentaje de 0=0%, 1=100% a corromper
        -- feat_range: tupla con intervalo de columnas a corromper
    """
    dataset_sz = dataset.shape[0]
    dataset_feat = dataset.shape[1]
    alt_row_indexes = np.random.randint(low=0, high=dataset_sz, size=int(dataset_sz*percent), dtype='l')
    alt_col_indexes = np.random.randint(
        low=feat_range[0],
        high=feat_range[1],
        size=int(dataset_sz*percent), dtype='l')
    dataset[alt_row_indexes, alt_col_indexes] = None
    return dataset
```

In [21]:

```
def drop_rows_with_nan(corrupted_dataset):
    """Elimina filas que tengan NaNs
    """
    return corrupted_dataset[np.all(np.isfinite(corrupted_dataset), axis=1)]

def drop_cols_with_nan(corrupted_dataset):
    """Elimina columnas que tengan NaNs
    """
    return corrupted_dataset[:, np.all(np.isfinite(corrupted_dataset), axis=0)]
```

In [22]:

```
# Ejercicio 9 - Example
dataset, _ = build_cluster(10, 1)
corrupted_dataset = corrupt_dataset(dataset, 0.4, (0,3))
corrupted_dataset
```

Out[22]:

```
array([[ 1.03551555,  0.68113311, -1.51051018, -0.23626792],
       [ 0.86242874, -0.5907845 ,          nan,  0.60171878],
       [ 2.24686909,  0.41950051,  0.50789615,  0.71105728],
       [          nan,  0.38145443,          nan,  0.27856089],
       [ 0.3822718 ,  2.09631142, -1.55116273,  0.2896872 ],
       [-2.81659897,  0.36331974, -1.07954851, -0.00774502],
       [-0.53674946, -1.14710732,  0.39007052, -1.22235032],
       [ 0.28400259,  2.86234799, -0.40558007,  0.30633819],
       [-1.47202197,  1.87482797,  1.52372873,  1.9360993 ],
       [          nan,  2.08841174,  0.99584428,  1.1426333 ]])
```

In [23]:

```
drop_rows_with_nan(corrupted_dataset)
```

Out[23]:

```
array([[ 1.03551555,  0.68113311, -1.51051018, -0.23626792],
       [ 2.24686909,  0.41950051,  0.50789615,  0.71105728],
       [ 0.3822718 ,  2.09631142, -1.55116273,  0.2896872 ],
       [-2.81659897,  0.36331974, -1.07954851, -0.00774502],
       [-0.53674946, -1.14710732,  0.39007052, -1.22235032],
       [ 0.28400259,  2.86234799, -0.40558007,  0.30633819],
       [-1.47202197,  1.87482797,  1.52372873,  1.9360993 ]])
```

In [24]:

```
dataset, _ = build_cluster(10, 1)
corrupted_dataset = corrupt_dataset(dataset, 0.4, (0,2))
corrupted_dataset
```

Out[24]:

```
array([[ nan,  1.18579247, -2.03976202,  1.72677752],
       [-0.26823567,  0.51586667,  2.37299614,  1.27457127],
       [ 1.13710743,  nan, -0.22919357, -1.35943673],
       [ 0.78500502, -1.15012096, -0.1959834 ,  1.08142597],
       [ 0.17296403,  1.41866399,  0.93422206, -0.33762249],
       [ 0.57318583,  0.35498939,  1.64424717, -0.29477832],
       [-0.39673013, -0.03822542, -1.58473381, -0.82797422],
       [-0.60074557,  0.13408455,  0.16226646,  0.30287211],
       [ nan, -0.59791761, -0.03308032, -1.3093717 ],
       [-0.56755086,  nan,  1.31810199, -0.68754121]])
```

In [25]:

```
drop_cols_with_nan(corrupted_dataset)
```

Out[25]:

```
array([[ -2.03976202,  1.72677752],
       [ 2.37299614,  1.27457127],
       [-0.22919357, -1.35943673],
       [-0.1959834 ,  1.08142597],
       [ 0.93422206, -0.33762249],
       [ 1.64424717, -0.29477832],
       [-1.58473381, -0.82797422],
       [ 0.16226646,  0.30287211],
       [-0.03308032, -1.3093717 ],
       [ 1.31810199, -0.68754121]])
```

Ejercicio 10

Reemplazar NaNs por la media de la columna. Dato un dataset, hacer una función que utilizando numpy reemplace los NaNs por la media de la columna.

In [26]:

```
# Ejercicio 10
def replace_nans_with_mean(dataset):
    filtered_rows = drop_rows_with_nan(dataset)
    col_avg = filtered_rows.mean(axis=0)
    nan_idx = np.where(np.isnan(dataset))
    dataset[nan_idx] = np.take(col_avg, nan_idx[1])
    return dataset
```

In [27]:

```
# Ejercicio 10 - Test
dataset, _ = build_cluster(10, 1)
corrupted_dataset = corrupt_dataset(dataset, 0.4, (0,2))
corrupted_dataset
```

Out[27]:

```
array([[ 2.17241388,  0.6139108 ,  0.96638639, -0.11291219],
       [-0.43861291, -1.01444412,  0.33203728, -1.55047366],
       [          nan,          nan, -1.40891645, -0.80665885],
       [ 1.21827696,  0.17410948, -0.1542649 , -0.69962163],
       [          nan,  0.42908096, -1.13490616, -0.6379495 ],
       [ 1.34832952,  2.59910225, -0.37849503, -0.45890284],
       [-2.14646426,  0.42101137,  2.0373931 ,  1.38939284],
       [-0.40132761,  1.89586126, -0.40137734, -0.03362552],
       [-0.72598669,  1.72806003,  0.53010083, -0.54071209],
       [          nan,  1.50038239,  0.73480828,  0.39209539]])
```

In [28]:

```
replace_nans_with_mean(corrupted_dataset)
```

Out[28]:

```
array([[ 2.17241388,  0.6139108 ,  0.96638639, -0.11291219],
       [-0.43861291, -1.01444412,  0.33203728, -1.55047366],
       [ 0.14666127,  0.91680158, -1.40891645, -0.80665885],
       [ 1.21827696,  0.17410948, -0.1542649 , -0.69962163],
       [ 0.14666127,  0.42908096, -1.13490616, -0.6379495 ],
       [ 1.34832952,  2.59910225, -0.37849503, -0.45890284],
       [-2.14646426,  0.42101137,  2.0373931 ,  1.38939284],
       [-0.40132761,  1.89586126, -0.40137734, -0.03362552],
       [-0.72598669,  1.72806003,  0.53010083, -0.54071209],
       [ 0.14666127,  1.50038239,  0.73480828,  0.39209539]])
```

Ejercicio 11

Dado un dataset X separarlo en 70 / 20 / 10.

Hint: A partir de utilizar `np.random.permutation` hacer un método que dado un dataset, devuelva los 3 datasets como nuevos numpy arrays.

In [29]:

```
# Ejercicio 11
def split_dataset(dataset, training=0.7, validation=0.2, testing=0.1):
    """ Dado un dataset, lo particiona en training, validation y testing y devuelve estos tres arrays.
    """
    dataset_sz = dataset.shape[0]
    i0 = int(dataset_sz*training)
    i1 = int(dataset_sz*(training+validation))
    indices = np.random.permutation(dataset_sz)
    training_idx, validation_idx, testing_idx = indices[0:i0], indices[i0:i1], indices[i1:]
    return dataset[training_idx], dataset[validation_idx], dataset[testing_idx]
```

In [30]:

```
# Ejercicio 11 - Test
dataset, _ = build_cluster(10, 1)
split_dataset(dataset)
```

Out[30]:

```
(array([[ -0.75359605,  1.97825486,  2.03993678, -1.7233503 ],
        [ 2.18928342, -0.96271816, -0.1458563 , -0.10078731],
        [ 0.83588296,  0.80378508, -1.50966575, -0.12042366],
        [ 1.11514604, -1.34605049,  1.02340521,  0.6475141 ],
        [ 1.04728188, -0.48129955, -0.35753219, -0.96871213],
        [-0.54316648,  0.7457398 , -1.88577976, -0.8922986 ],
        [-0.24294314, -0.54136055,  0.53761065,  0.84812656]]),
 array([[ 1.07312957,  0.50290546, -0.19125833, -0.84881579],
        [ 0.21506568,  1.35439905, -2.14566108,  0.90164097]]),
 array([[ -0.71091229,  1.29990521, -1.22699549,  1.63428364]]))
```

Ejercicio 12 (integrador)

1. Generar un dataset sintético que clusterice data en 4 clusters utilizando números random.
 - a. Utilizar 4 dimensiones.
 - b. Generar un dataset con 100K de muestras.

In [31]:

```
N_SAMPLES = 100000
CLUSTER_INV_OVERLAP = 10
dataset, clusters = build_cluster(N_SAMPLES, CLUSTER_INV_OVERLAP)
dataset.shape
```

Out[31]:

(100000, 4)

1. Cambiar algunos puntos de manera aleatoria y agregar NaN (0.1% del dataset).

2.1 Agregar ruido a índices elegidos de manera aleatoria.

In [32]:

```
NOISE_TEMP = 5.0
alt_indexes = np.random.randint(low=0, high=N_SAMPLES-1, size=int(N_SAMPLES*0.1), dtype='l')
noise = NOISE_TEMP*np.random.normal(size=(alt_indexes.shape[0],4))
dataset[alt_indexes] += noise
```

2.1 Corromper con NaNs.

In [33]:

```
dataset = corrupt_dataset(dataset, 0.1,(0,3))
```

1. Guardar el dataset en un .pkl

In [34]:

```
import pickle

PKL_FILENAME = "dataset.pkl"
with open(PKL_FILENAME, 'wb') as file:
    pickle.dump(dataset, file, protocol=pickle.HIGHEST_PROTOCOL)
    print(f"Dumped to {PKL_FILENAME}")
```

Dumped to dataset.pkl

1. Cargar el dataset con Numpy desde el .pkl

In [35]:

```
dataset = None
dataset
```

In [36]:

```
print(f"Loading from PKL {PKL_FILENAME}")
with open(PKL_FILENAME, 'rb') as file:
    dataset = pickle.load(file)
```

Loading from PKL dataset.pkl

In [37]:

```
dataset.shape
```

Out[37]:

```
(100000, 4)
```

1. Completar NaN con la media de cada feature.

In [38]:

dataset

Out[38]:

```
array([[11.56189071, -0.41420008,  0.02096634,  1.47585016],
       [10.43573303, -0.81087638, -0.47967898,  0.51947519],
       [10.51852598, -0.70749926, -0.44358223,  0.62866123],
       ...,
       [-0.23293592,          nan, -2.59374293,  1.40502653],
       [ 0.54274824,          nan,  1.94628679, -1.01505868],
       [ 0.74418936,  9.67459149,  0.82896338, -0.61364884]])
```

In [39]:

```
dataset = replace_nans_with_mean(dataset)
dataset
```

Out[39]:

```
array([[11.56189071, -0.41420008,  0.02096634,  1.47585016],
       [10.43573303, -0.81087638, -0.47967898,  0.51947519],
       [10.51852598, -0.70749926, -0.44358223,  0.62866123],
       ...,
       [-0.23293592,  4.99831595, -2.59374293,  1.40502653],
       [ 0.54274824,  4.99831595,  1.94628679, -1.01505868],
       [ 0.74418936,  9.67459149,  0.82896338, -0.61364884]])
```

1. Calcular la norma l2, la media y el desvío de cada feature con funciones numpy vectorizadas.

In [40]:

```
import sys
sys.path.append("../")
from clase_1.ex1_norm import vector_norm_l2
```

Nota: el transpuesto es porque la implementación de `vector_norm_l2()` de la clase anterior se calculaba sobre las filas.

In [41]:

```
norms = vector_norm_l2(dataset.T)
norms
```

Out[41]:

```
array([2290.47122668, 2290.69189506,  574.72288469,  583.08470669])
```

In [42]:

```
feat_mean = np.mean(data,axis=0)
feat_mean
```

Out[42]:

```
array([ 1.26806786,  0.89029501, -0.1842862 , -0.69102905])
```

In [43]:

```
feat_std = np.std(data,axis=0)
feat_std
```

Out[43]:

```
array([1.08475562, 1.19342098, 0.87889799, 0.8599748 ])
```

1. Agregar una columna a partir de generar una variable aleatoria exponencial a todos los puntos.

In [44]:

```
dataset= np.hstack((dataset, np.zeros((dataset.shape[0], 1), dtype=dataset.dtype
)))
dataset.shape
```

Out[44]:

```
(100000, 5)
```

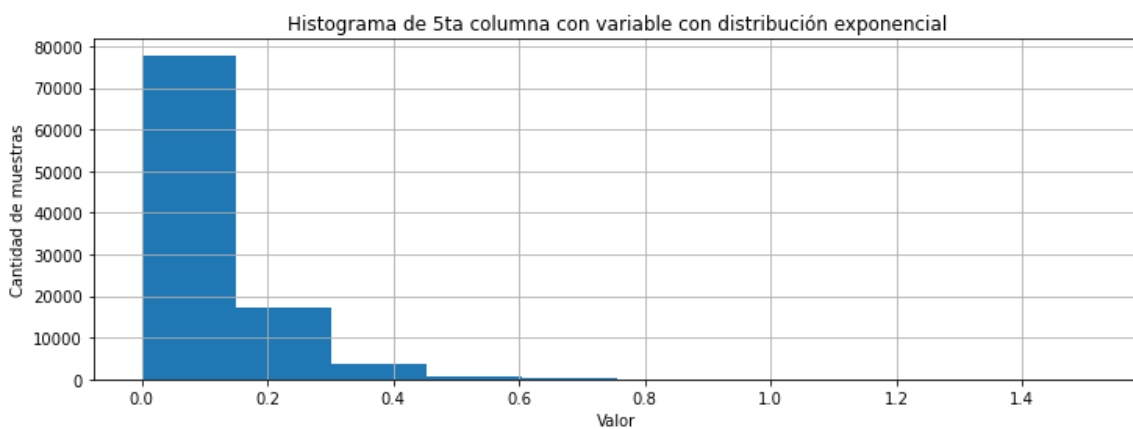
In [45]:

```
lambda_param = 10
dataset[:,4] = exponential_random_variable(lambda_param, dataset.shape[0])
```

1. Hacer el histograma de la distribución exponencial.

In [46]:

```
plt.figure(figsize=(12,4))
plt.hist(dataset[:,4])
plt.title("Histograma de 5ta columna con variable con distribución exponencial")
plt.xlabel("Valor")
plt.ylabel("Cantidad de muestras")
plt.grid(which="Both")
plt.show()
```



1. Aplicar PCA al dataset reduciendo a 2 dimensiones y graficar el cluster.

In [47]:

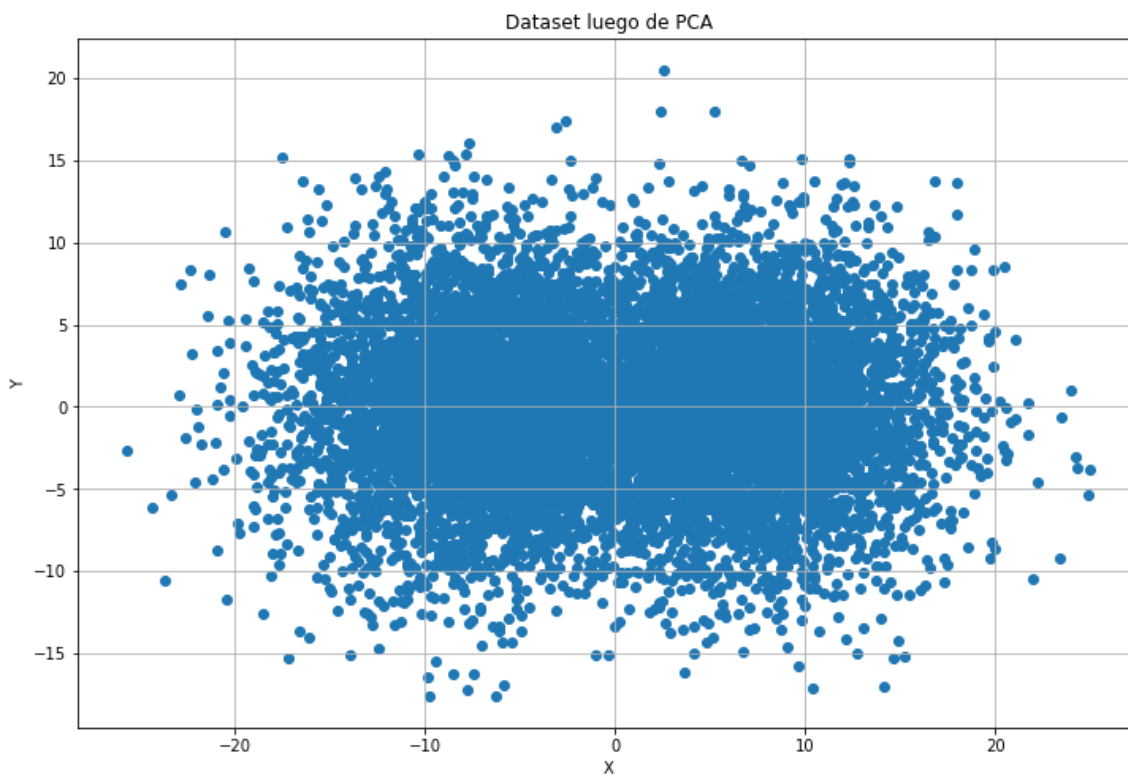
```
compressed_dataset = my_pca(dataset)
compressed_dataset.shape
```

Out[47]:

(100000, 2)

In [48]:

```
plt.figure(figsize=(12,8))
plt.scatter(compressed_dataset[:,0],compressed_dataset[:,1])
plt.title("Dataset luego de PCA")
plt.xlabel("X")
plt.ylabel("Y")
plt.grid(which="Both")
plt.show()
```



1. Hacer la clusterización con el k-means desarrollado en clase.

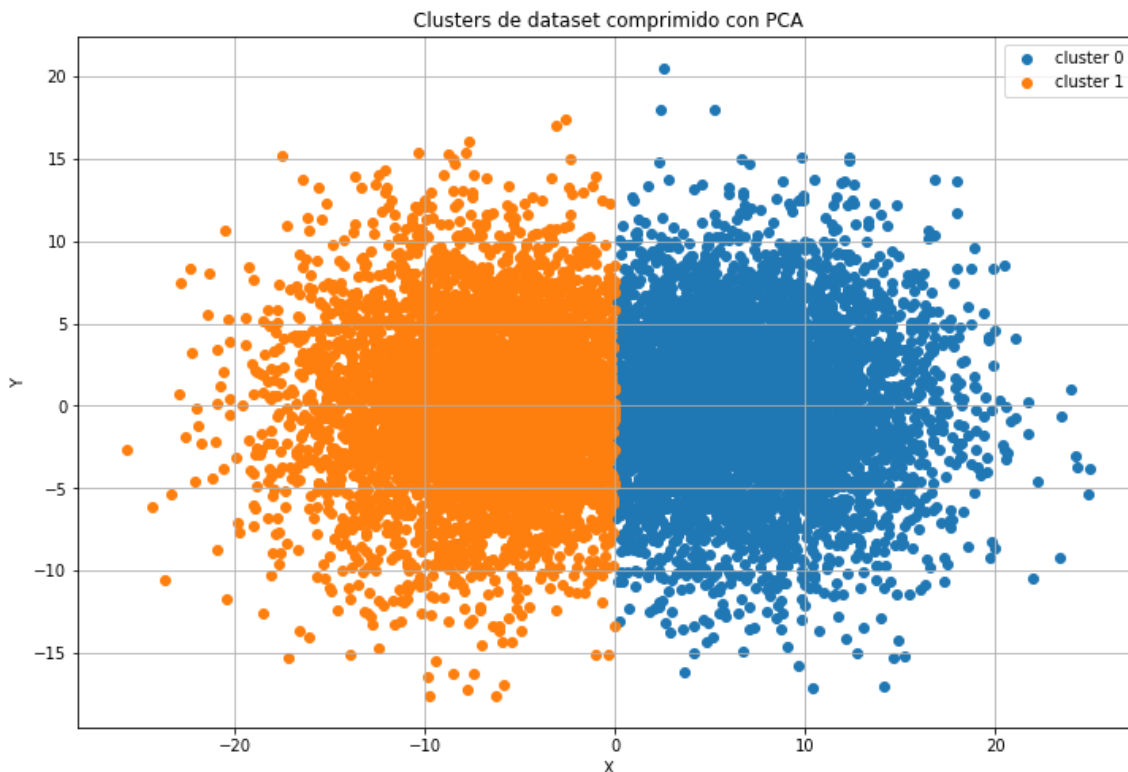
In [49]:

```
centroids, cluster_ids = k_means(compressed_dataset,2)
```

1. Volver a graficar el cluster con lo obtenido en (10) y comparar resultados con (9).

In [50]:

```
def plot_clusters(compressed_dataset, cluster_ids):  
    plt.figure(figsize=(12,8))  
    plt.scatter(compressed_dataset[cluster_ids==0,0],compressed_dataset[cluster_  
ids==0,1])  
    plt.scatter(compressed_dataset[cluster_ids==1,0],compressed_dataset[cluster_  
ids==1,1])  
    plt.legend(["cluster 0", "cluster 1"])  
    plt.title("Clusters de dataset comprimido con PCA")  
    plt.xlabel("X")  
    plt.ylabel("Y")  
    plt.grid(which="Both")  
    plt.show()  
  
plot_clusters(compressed_dataset, cluster_ids)
```

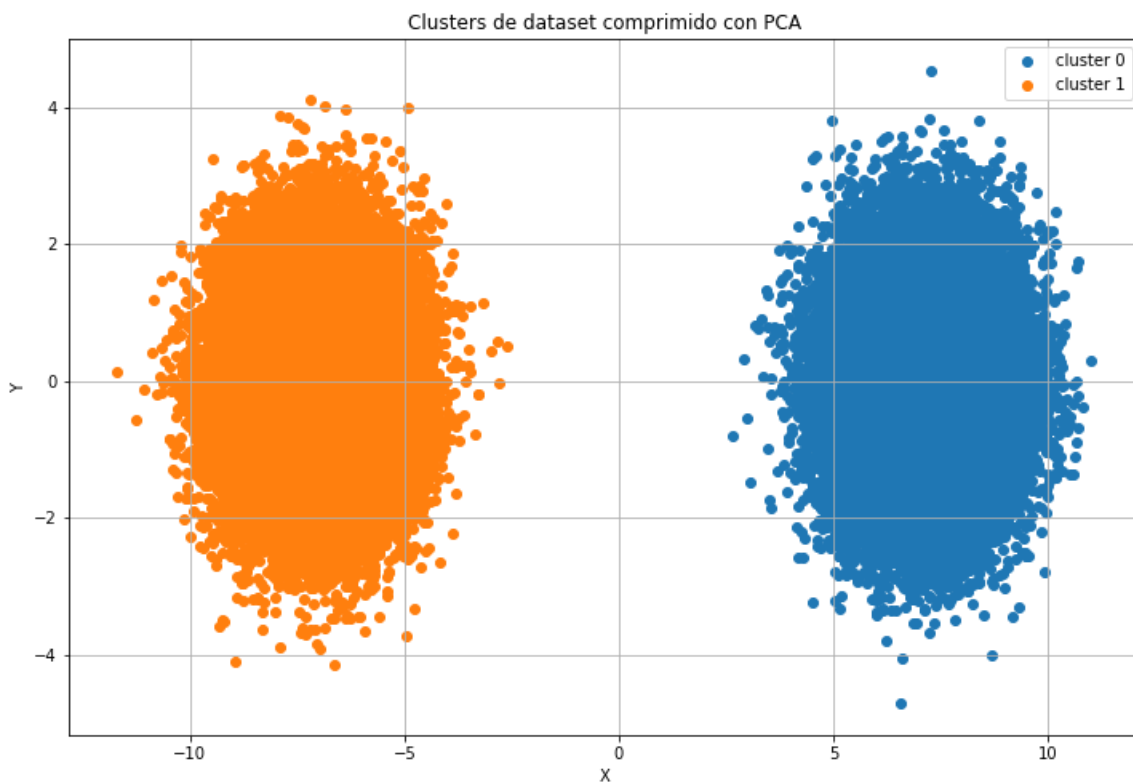


1. Analizar qué pasa si los clusters comienzan a tener overlapping.

Cuando los clusters están alejados es clara la separación de los centroides A/B y la asignación de cada punto a uno de ellos.

In [51]:

```
inv_overlap = 10  
dataset, clusters = build_cluster(N_SAMPLES, inv_overlap)  
compressed_dataset = my_pca(dataset)  
centroids, cluster_ids = k_means(compressed_dataset, 2)  
plot_clusters(compressed_dataset, cluster_ids)
```



Cuando los clusters están próximos igual se produce una asignación de cada punto a un cluster, pero esta categorización podría no ser representativa para el problema que se quiere resolver.

In [52]:

```
inv_overlap = 0.00001  
dataset, clusters = build_cluster(N_SAMPLES, inv_overlap)  
compressed_dataset = my_pca(dataset)  
centroids, cluster_ids = k_means(compressed_dataset, 2)  
plot_clusters(compressed_dataset, cluster_ids)
```

