# K-means clustering

The notebook aims to study and implement a k-means clustering using "sklearn". The iris dataset will be used to identify clusters automatically using the K-means method.

## Acknowledgments

- Used dataset: https://archive.ics.uci.edu/ml/datasets/iris

- Inquiries: mauricio.antelis@tec.mx

## ⌄ Importing libraries

```
1 # Define where you are running the code: colab or local
2 RunInColab          = True     # (False: no  | True: yes)
3
4 # If running in colab:
5 if RunInColab:
6     # Mount your google drive in google colab
7     from google.colab import drive
8     drive.mount('/content/drive')
9
10    # Find location
11    #!pwd
12    #!ls
13    #!ls "/content/drive/My Drive/Colab Notebooks/MachineLearningWithPython/"
14
15    # Define path del proyecto
16    Ruta          = "/content/drive/MyDrive/ITC/5toSem/semanaTecAn/"
17
18 else:
19    # Define path del proyecto
20    Ruta          = ""
```

⇥ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_r

Suggested code may be subject to a license | | AnuragAnalog/Guided-Projects-Coursera | askorucuk/Randomforest_classification

```
1 # Import the packages that we will be using
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans
6 import seaborn as sns
```

## ⌄ Importing data

```
1 # Dataset url
2 url = "datasets/iris.csv"
3 # Load the dataset from HHDD
4 df = pd.read_csv(Ruta + url)
5 # change the name of the columns
6 df.columns = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width', 'flower']
7 # Visualize the dataset
8 df
```

|   | sepal_length | sepal_width | petal_length | petal_width | flower |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Setosa |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Virginica |

150 rows × 5 columns

Next steps:  [ Generate code with `df` ]   [ 🔘 View recommended plots ]   [ New interactive sheet ]

## Undertanding and preprocessing the data

1. Get a general 'feel' of the data

```
1 # get a general feel of the data
2 df
```

|   | sepal_length | sepal_width | petal_length | petal_width | flower |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Setosa |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Virginica |

150 rows × 5 columns

Next steps:  [ Generate code with `df` ]   [ 🔘 View recommended plots ]   [ New interactive sheet ]

2. Drop rows with any missing values

```
1 #drop rows with missing values
2 df.dropna(inplace=True)
3 df
```

|     | sepal_length | sepal_width | petal_length | petal_width | flower |
|-----|--------------|-------------|--------------|-------------|--------|
| 0   | 5.1          | 3.5         | 1.4          | 0.2         | Setosa |
| 1   | 4.9          | 3.0         | 1.4          | 0.2         | Setosa |
| 2   | 4.7          | 3.2         | 1.3          | 0.2         | Setosa |
| 3   | 4.6          | 3.1         | 1.5          | 0.2         | Setosa |
| 4   | 5.0          | 3.6         | 1.4          | 0.2         | Setosa |
| ... | ...          | ...         | ...          | ...         | ...    |
| 145 | 6.7          | 3.0         | 5.2          | 2.3         | Virginica |
| 146 | 6.3          | 2.5         | 5.0          | 1.9         | Virginica |
| 147 | 6.5          | 3.0         | 5.2          | 2.0         | Virginica |
| 148 | 6.2          | 3.4         | 5.4          | 2.3         | Virginica |
| 149 | 5.9          | 3.0         | 5.1          | 1.8         | Virginica |

150 rows × 5 columns

Next steps:  [ Generate code with `df` ]   [ 🔘 View recommended plots ]   [ New interactive sheet ]

### 3. Encoding the class label categorical column: from string to num

```
1 # Encoding the categorical column
2 df=df.replace({'flower': {'Setosa': 0, 'Versicolor': 1, 'Virginica': 2}})
3 #Visualize the dataset
4 df
```

|     | sepal_length | sepal_width | petal_length | petal_width | flower |
|-----|--------------|-------------|--------------|-------------|--------|
| 0   | 5.1          | 3.5         | 1.4          | 0.2         | 0      |
| 1   | 4.9          | 3.0         | 1.4          | 0.2         | 0      |
| 2   | 4.7          | 3.2         | 1.3          | 0.2         | 0      |
| 3   | 4.6          | 3.1         | 1.5          | 0.2         | 0      |
| 4   | 5.0          | 3.6         | 1.4          | 0.2         | 0      |
| ... | ...          | ...         | ...          | ...         | ...    |
| 145 | 6.7          | 3.0         | 5.2          | 2.3         | 2      |
| 146 | 6.3          | 2.5         | 5.0          | 1.9         | 2      |
| 147 | 6.5          | 3.0         | 5.2          | 2.0         | 2      |
| 148 | 6.2          | 3.4         | 5.4          | 2.3         | 2      |
| 149 | 5.9          | 3.0         | 5.1          | 1.8         | 2      |

150 rows × 5 columns

Next steps:  [ Generate code with `df` ]   [ 🔘 View recommended plots ]   [ New interactive sheet ]

Now the label/category is numeric

## 4. Discard columns that won't be used

```
1 # # Drop out non necesary columns
2 # dataset.drop(['Sepal_Length', 'Sepal_Width'],axis='columns',inplace=True)
3 #
4 # #Visualize the dataset
5 # dataset
6
```
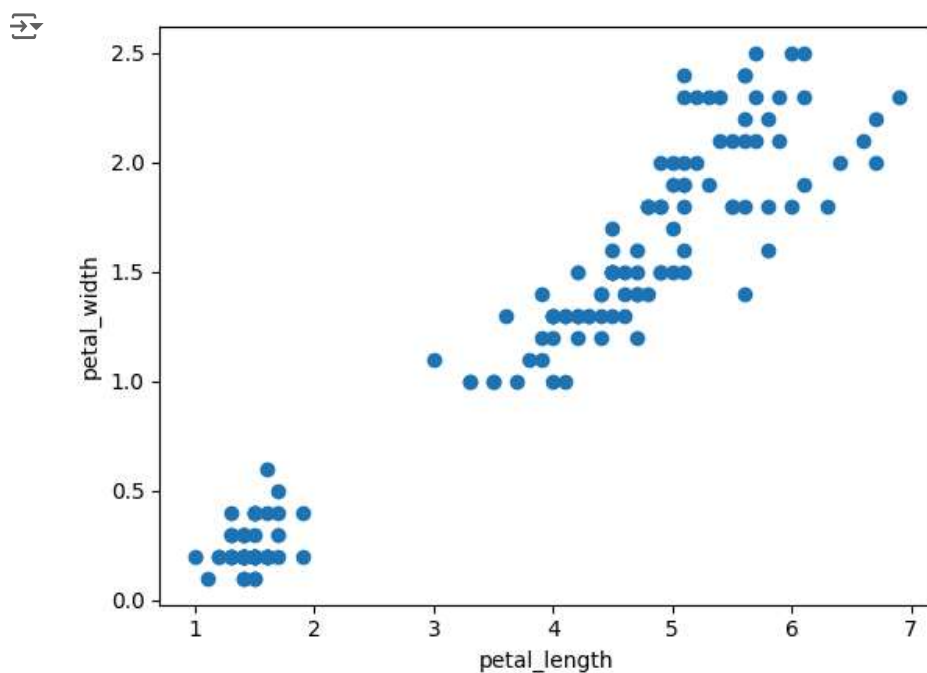
## 5. Scatter plot of the data

```
1 plt.scatter(df.petal_length, df.petal_width)
2 plt.xlabel('petal_length')
3 plt.ylabel('petal_width')
4 plt.show()
```
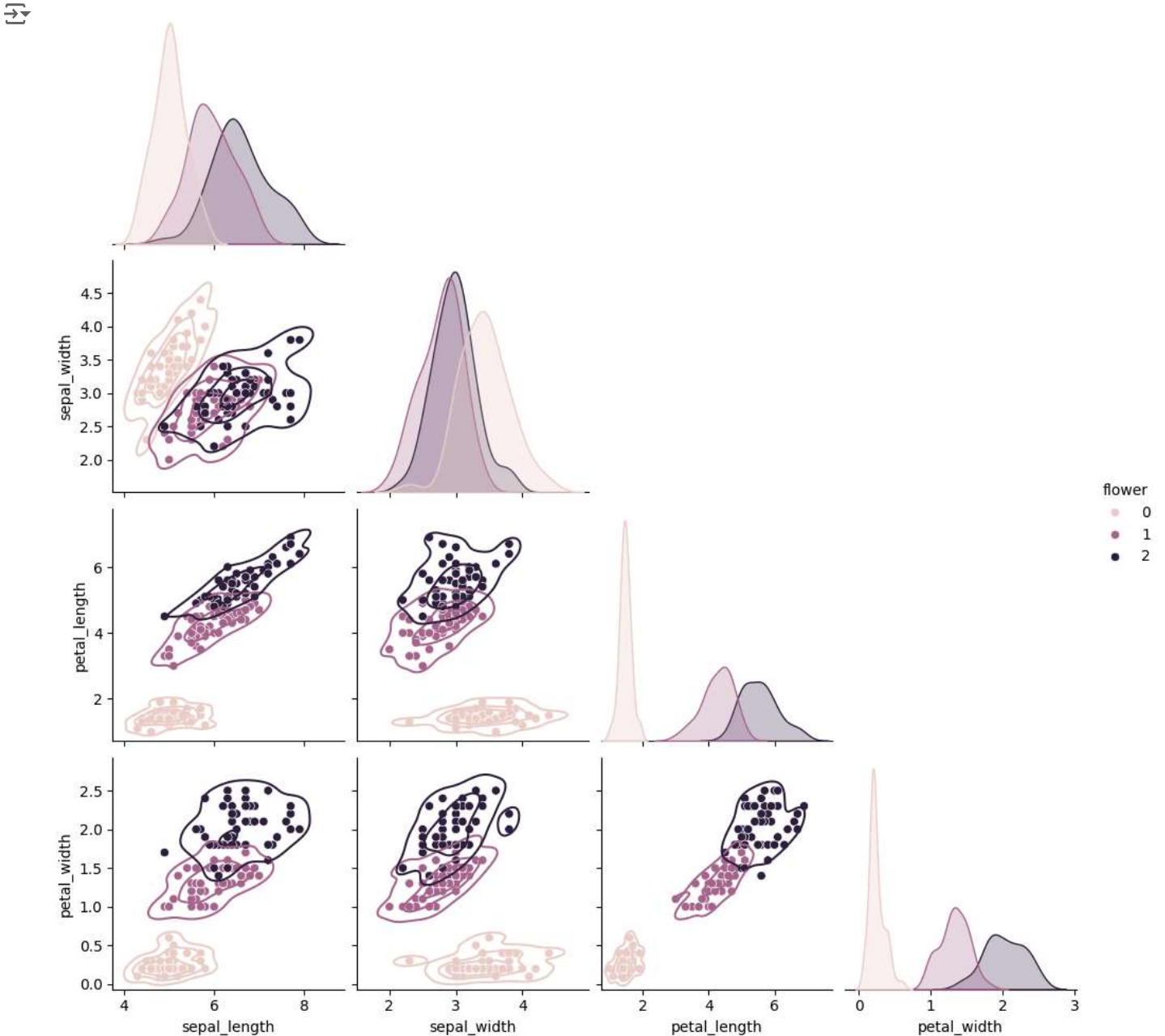


```
1 #PAIRPLOT: SCATTER PLOT OF ALL
2 g = sns.pairplot(df, corner = True, diag_kind = 'kde', hue = 'flower')
3 g.map_lower(sns.kdeplot, levels = 4, color = '.2')
4 plt.show()
```
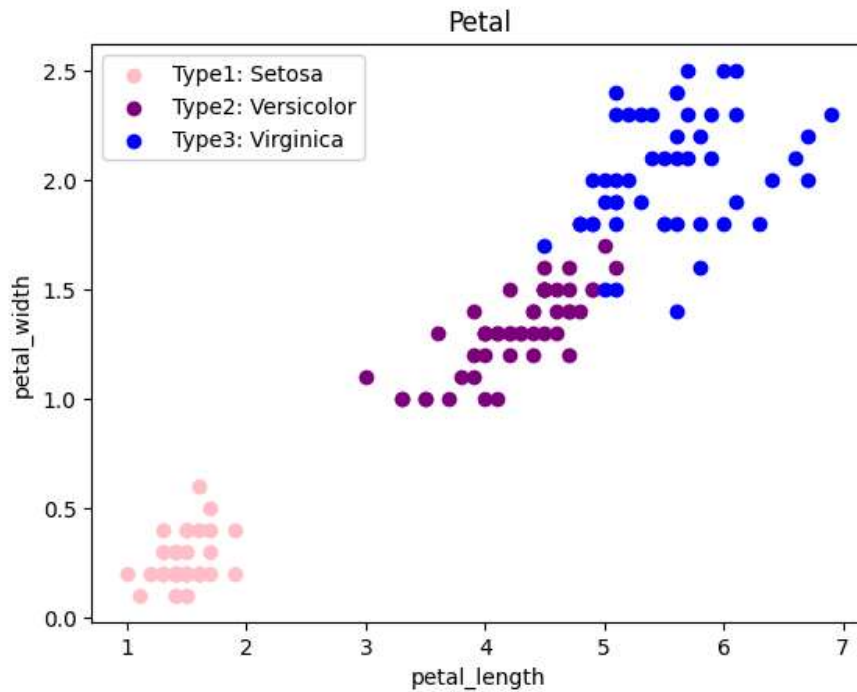
6. Scatter plot of the data asigning each point to the cluster it belongs to ¡¡

```
 1 # Get dataframes for each real cluster
 2 df0 = df[df.flower == 0]
 3 df1 = df[df.flower == 1]
 4 df2 = df[df.flower == 2]
 5 plt.scatter(df0.petal_length, df0.petal_width, color = 'pink')
 6 plt.scatter(df1.petal_length, df1.petal_width, color = 'purple')
 7 plt.scatter(df2.petal_length, df2.petal_width, color = 'blue')
 8 plt.legend(['Type1: Setosa', 'Type2: Versicolor', 'Type3: Virginica'])
 9 plt.xlabel('petal_length')
10 plt.ylabel('petal_width')
11 plt.title('Petal')
12 plt.show()
```

Recall that for this dataset we know in advance the class to which each point belongs to

## ⌄ Kmeans clustering

Kmeans clustering

```
1 # Import sklearn KMeans
2 from sklearn.cluster import KMeans
3 from sklearn.metrics import silhouette_score
4 # Define number of clusters
5 k = 3
6 # Do k-means clustering (assign each point in the dataset to a cluster)
7 km = KMeans(n_clusters=k, n_init="auto") #modelo ML here
8 FlowerPredicted = km.fit_predict(df[['sepal_length', 'sepal_width','petal_length', 'petal_width']]) #entrena y predi
9 # Print estimated cluster of each point in the dataset
10 FlowerPredicted
```

```
array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1,
       1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1,
       1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0], dtype=int32)
```

NOTE: the lables of the estimated clusters do not agree with the lables in the real labels, therefore, it will be important to pair the labels of the real and estimated clusters

```
1 # Manual pairing the labels of the real and estimated clusters
```

```
1 # Add cluster information to dataset
2 df['flower_predictedK3'] = FlowerPredicted
3 # Show the first few rows to check the cluster assignment
4 df.head()
```

| | sepal_length | sepal_width | petal_length | petal_width | flower | flower_predictedK3 |
|---|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | 0 | 2 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | 0 | 2 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | 0 | 2 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | 0 | 2 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | 0 | 2 |

Next steps:   [ Generate code with `df` ]   [ ● View recommended plots ]   [ New interactive sheet ]

Double-click (or enter) to edit

```
1 # Display cluster labels
2 FlowerPredicted = df['flower_predictedK3']
3 print("Labels of the estimated clusters:", FlowerPredicted.unique())
```

Labels of the estimated clusters: [2 1 0]

```
1 # Get the centroids of the clusters
2 print("Cluster centroids:\n", km.cluster_centers_)
```

Cluster centroids:
 [[5.88360656 2.74098361 4.38852459 1.43442623]
 [6.85384615 3.07692308 5.71538462 2.05384615]
 [5.006      3.428      1.462      0.246     ]]

```
1 # Sum of squared error (sse) of the final model
2 sse = km.inertia_
3 print("Sum of squared error (SSE):", sse)
```

Sum of squared error (SSE): 78.85566582597727

```
1 # The number of iterations required to converge
2 print("Number of iterations:", km.n_iter_)
```

Number of iterations: 12

**Important remarks**

- The number of each cluster is randomly assigned
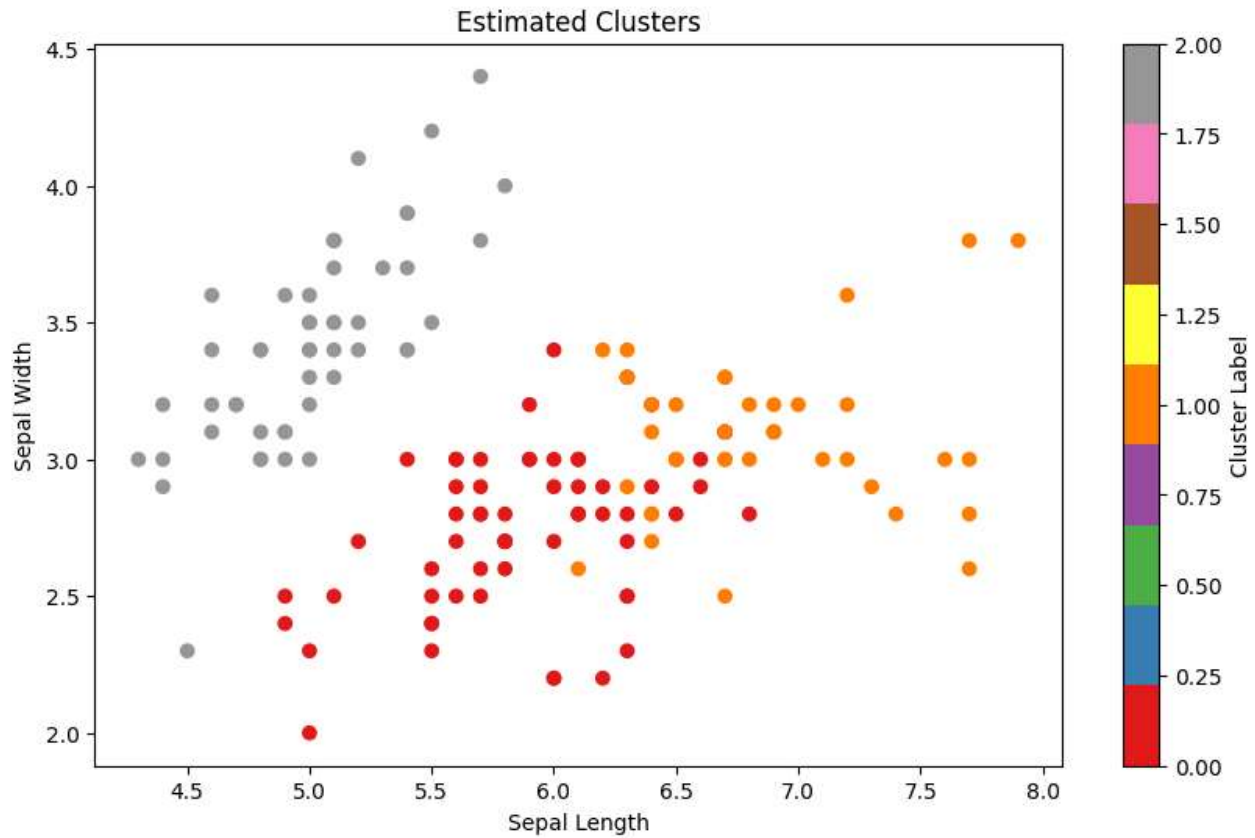- The order of the numer in each cluster is random

## ⌄ Plot estimated clusters

Plot estimated clusters

```
1 plt.figure(figsize=(10, 6))
2 plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower_predictedK3'], cmap='Set1', marker='o')
3 plt.title('Estimated Clusters')
4 plt.xlabel('Sepal Length')
5 plt.ylabel('Sepal Width')
6 plt.colorbar(label='Cluster Label')
7 plt.show()
```
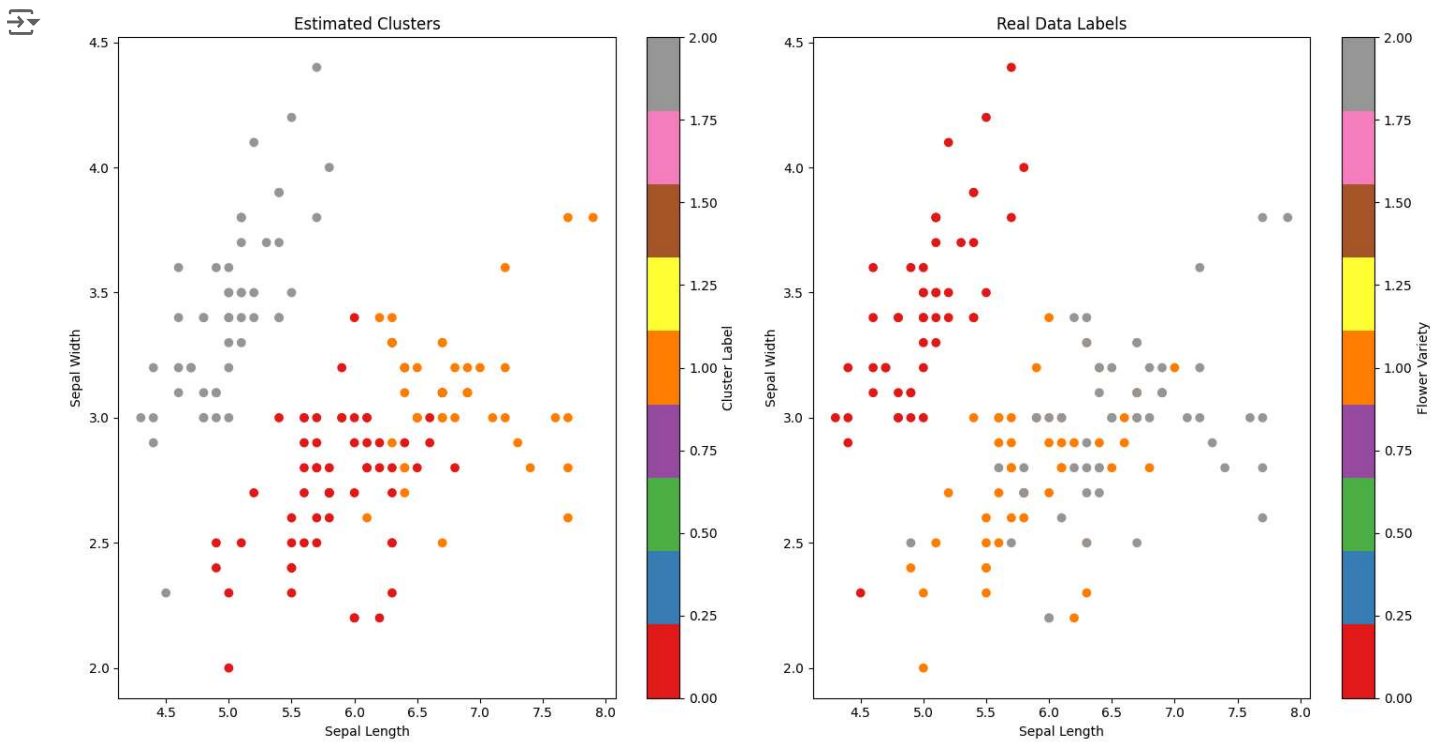
Estimated Clusters

## Plot both real and estimated clusters to check for errors

```
1  #GRAPH
2  plt.figure(figsize=(15, 8))
3
4  # Plot estimated clusters
5  plt.subplot(1, 2, 1)
6  scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower_predictedK3'], cmap='Set1', marker='o')
7  plt.title('Estimated Clusters')
8  plt.xlabel('Sepal Length')
9  plt.ylabel('Sepal Width')
10 plt.colorbar(label='Cluster Label')
11
12 # Plot real labels
13 plt.subplot(1, 2, 2)
14 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower'].astype('category').cat.codes, cmap='Set1
15 plt.title('Real Data Labels')
16 plt.xlabel('Sepal Length')
17 plt.ylabel('Sepal Width')
18 plt.colorbar(label='Flower Variety')
19
20 plt.tight_layout()
21 plt.show()
```

Choose the k after which the sse is minimally reduced

**Important remarks**

- Note that for K=2 ...
- Note that for K=4 ...
- Note that for K=5 ...

**EXAMPLE WITH K=4**

```
 1 print("NUMBER OF CLUSTERS = 4")
 2 # Define number of clusters
 3 k = 4
 4 # Do k-means clustering (assign each point in the dataset to a cluster)
 5 km = KMeans(n_clusters=k, n_init="auto") #modelo ML here
 6 FlowerPredicted = km.fit_predict(df[['sepal_length', 'sepal_width','petal_length', 'petal_width']]) #entrena y predice
 7 # Print estimated cluster of each point in the dataset
 8 FlowerPredicted
 9 # Add cluster information to dataset
10 df['flower_predictedK4'] = FlowerPredicted
11 # Show the first few rows to check the cluster assignment
12 print(df.head())
13 # Display cluster labels
14 FlowerPredicted = df['flower_predictedK4']
15 print("Labels of the estimated clusters:", FlowerPredicted.unique())
16 # Get the centroids of the clusters
17 print("Cluster centroids:\n", km.cluster_centers_)
18 # Sum of squared error (sse) of the final model
19 sse = km.inertia_
20 print("Sum of squared error (SSE):", sse)
21 # The number of iterations required to converge
22 print("Number of iterations:", km.n_iter_)
23
24 #GRAPH
```

```
25 plt.figure(figsize=(15, 8))
26
27 # Plot estimated clusters
28 plt.subplot(1, 2, 1)
29 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower_predictedK4'], cmap='Set1', marker='o')
30 plt.title('Estimated Clusters')
31 plt.xlabel('Sepal Length')
32 plt.ylabel('Sepal Width')
33 plt.colorbar(label='Cluster Label')
34
35 # Plot real labels
36 plt.subplot(1, 2, 2)
37 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower'].astype('category').cat.codes, cmap='Set1'
38 plt.title('Real Data Labels')
39 plt.xlabel('Sepal Length')
40 plt.ylabel('Sepal Width')
41 plt.colorbar(label='Flower Variety')
42
43 plt.tight_layout()
44 plt.show()
```

```
NUMBER OF CLUSTERS = 4
   sepal_length  sepal_width  petal_length  petal_width  flower  \
0           5.1          3.5           1.4          0.2       0
1           4.9          3.0           1.4          0.2       0
2           4.7          3.2           1.3          0.2       0
3           4.6          3.1           1.5          0.2       0
4           5.0          3.6           1.4          0.2       0


   flower_predictedK3  flower_predictedK4
0                   2                   1
1                   2                   1
2                   2                   1
3                   2                   1
4                   2                   1
Labels of the estimated clusters: [1 0 2 3]
Cluster centroids:
 [[6.23658537 2.85853659 4.80731707 1.62195122]
 [5.006      3.428      1.462      0.246     ]
 [5.52962963 2.62222222 3.94074074 1.21851852]
 [6.9125     3.1        5.846875   2.13125   ]]
Sum of squared error (SSE): 57.25600931571815
Number of iterations: 3
```
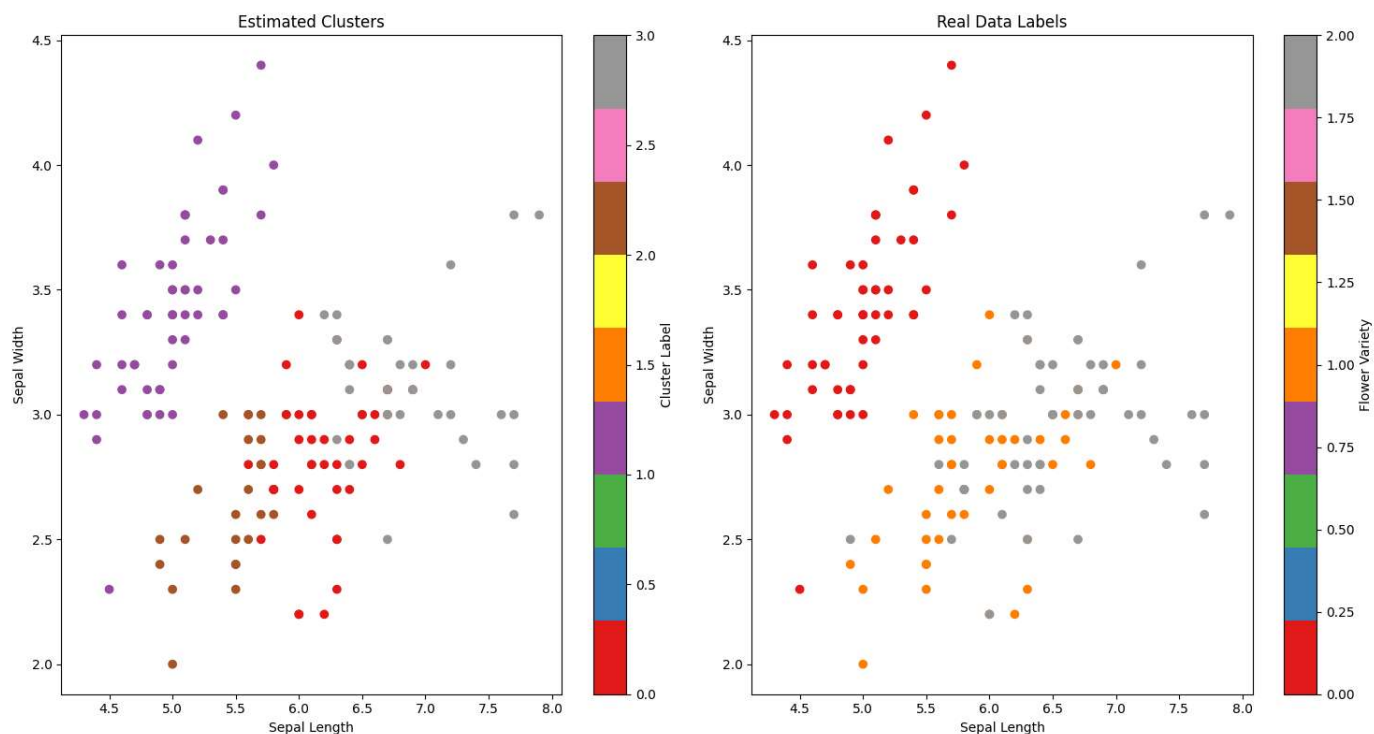


## EXAMPLE WITH K=5

```
 1 print("NUMBER OF CLUSTERS = 5")
 2 # Define number of clusters
 3 k = 5
 4 # Do k-means clustering (assign each point in the dataset to a cluster)
 5 km = KMeans(n_clusters=k, n_init="auto") #modelo ML here
 6 FlowerPredicted = km.fit_predict(df[['sepal_length', 'sepal_width','petal_length', 'petal_width']]) #entrena y predi
 7 # Print estimated cluster of each point in the dataset
 8 FlowerPredicted
 9 # Add cluster information to dataset
10 df['flower_predictedK5'] = FlowerPredicted
11 # Show the first few rows to check the cluster assignment
12 print(df.head())
13 # Display cluster labels
14 FlowerPredicted = df['flower_predictedK5']
15 print("Labels of the estimated clusters:", FlowerPredicted.unique())
16 # Get the centroids of the clusters
17 print("Cluster centroids:\n", km.cluster_centers_)
18 # Sum of squared error (sse) of the final model
19 sse = km.inertia_
20 print("Sum of squared error (SSE):", sse)
21 # The number of iterations required to converge
22 print("Number of iterations:", km.n_iter_)
23
24 #GRAPH
25 plt.figure(figsize=(15, 8))
26
27 # Plot estimated clusters
28 plt.subplot(1, 2, 1)
29 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower_predictedK5'], cmap='Set1', marker='o')
30 plt.title('Estimated Clusters')
31 plt.xlabel('Sepal Length')
32 plt.ylabel('Sepal Width')
33 plt.colorbar(label='Cluster Label')
34
35 # Plot real labels
36 plt.subplot(1, 2, 2)
37 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower'].astype('category').cat.codes, cmap='Set1'
38 plt.title('Real Data Labels')
39 plt.xlabel('Sepal Length')
40 plt.ylabel('Sepal Width')
41 plt.colorbar(label='Flower Variety')
42
43 plt.tight_layout()
44 plt.show()
```

```
NUMBER OF CLUSTERS = 5
   sepal_length  sepal_width  petal_length  petal_width  flower  \
0           5.1          3.5           1.4          0.2       0
1           4.9          3.0           1.4          0.2       0
2           4.7          3.2           1.3          0.2       0
3           4.6          3.1           1.5          0.2       0
4           5.0          3.6           1.4          0.2       0

   flower_predictedK3  flower_predictedK4  flower_predictedK5
0                   2                   1                   0
1                   2                   1                   0
2                   2                   1                   0
3                   2                   1                   0
4                   2                   1                   0
Labels of the estimated clusters: [0 2 3 1 4]
Cluster centroids:
 [[5.006      3.428      1.462      0.246     ]
 [6.52916667 3.05833333 5.50833333 2.1625    ]
 [6.20769231 2.85384615 4.74615385 1.56410256]
 [5.508      2.6        3.908      1.204     ]
 [7.475      3.125      6.3        2.05      ]]
Sum of squared error (SSE): 46.44618205128204
Number of iterations: 5
```
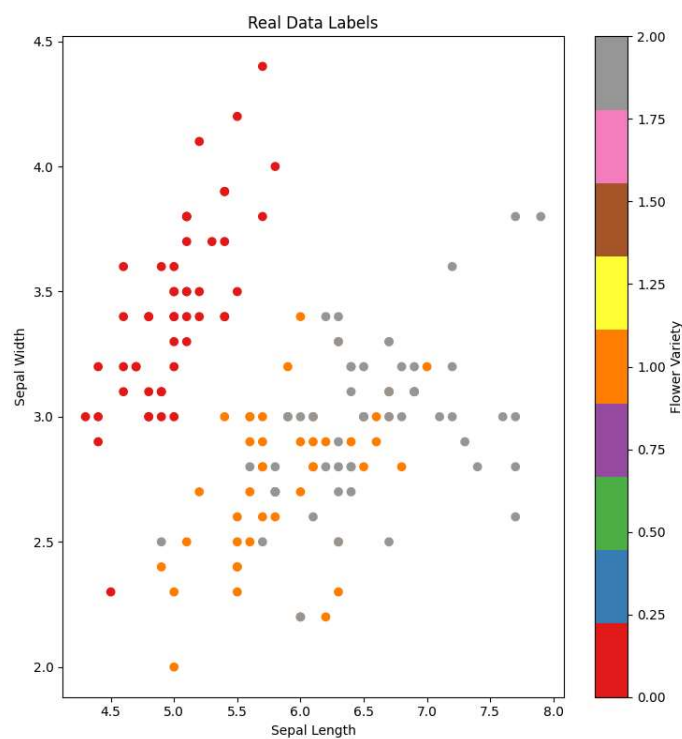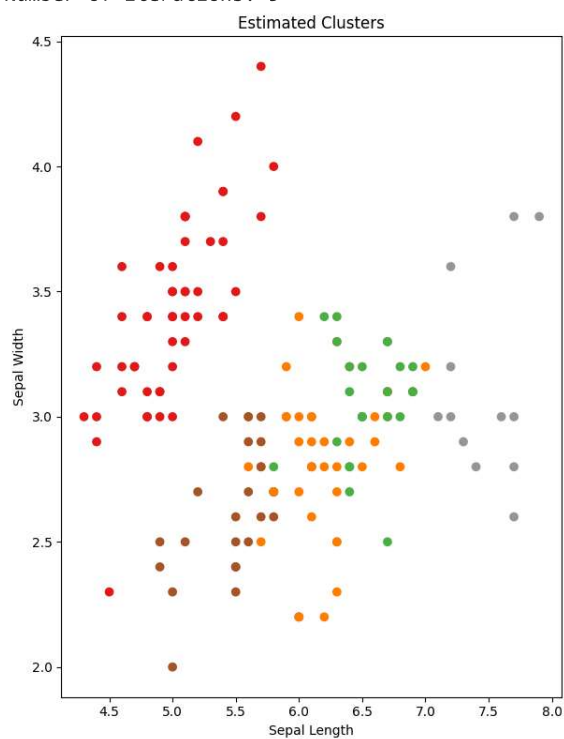


## EXAMPLE WITH K=2

```
1  print("NUMBER OF CLUSTERS = 2")
2  # Define number of clusters
3  k = 2
4  # Do k-means clustering (assign each point in the dataset to a cluster)
5  km = KMeans(n_clusters=k, n_init="auto") #modelo ML here
6  FlowerPredicted = km.fit_predict(df[['sepal_length', 'sepal_width','petal_length', 'petal_width']]) #entrena y predi
7  # Print estimated cluster of each point in the dataset
8  FlowerPredicted
9  # Add cluster information to dataset
10 df['flower_predictedK2'] = FlowerPredicted
11 # Show the first few rows to check the cluster assignment
12 print(df.head())
13 # Display cluster labels
14 FlowerPredicted = df['flower_predictedK2']
15 print("Labels of the estimated clusters:", FlowerPredicted.unique())
16 # Get the centroids of the clusters
17 print("Cluster centroids:\n", km.cluster_centers_)
18 # Sum of squared error (sse) of the final model
19 sse = km.inertia_
20 print("Sum of squared error (SSE):", sse)
21 # The number of iterations required to converge
22 print("Number of iterations:", km.n_iter_)
23
24 #GRAPH
25 plt.figure(figsize=(15, 8))
26
27 # Plot estimated clusters
28 plt.subplot(1, 2, 1)
29 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower_predictedK2'], cmap='Set1', marker='o')
30 plt.title('Estimated Clusters')
31 plt.xlabel('Sepal Length')
32 plt.ylabel('Sepal Width')
33 plt.colorbar(label='Cluster Label')
34
35 # Plot real labels
36 plt.subplot(1, 2, 2)
37 scatter = plt.scatter(df['sepal_length'], df['sepal_width'], c=df['flower'].astype('category').cat.codes, cmap='Set1
38 plt.title('Real Data Labels')
39 plt.xlabel('Sepal Length')
40 plt.ylabel('Sepal Width')
41 plt.colorbar(label='Flower Variety')
42
43 plt.tight_layout()
44 plt.show()
```

```
NUMBER OF CLUSTERS = 2
   sepal_length  sepal_width  petal_length  petal_width  flower  \
0           5.1          3.5           1.4          0.2       0
1           4.9          3.0           1.4          0.2       0
2           4.7          3.2           1.3          0.2       0
3           4.6          3.1           1.5          0.2       0
4           5.0          3.6           1.4          0.2       0

   flower_predictedK3  flower_predictedK4  flower_predictedK5  \
0                   2                   1                   0
1                   2                   1                   0
2                   2                   1                   0
3                   2                   1                   0
4                   2                   1                   0

   flower_predictedK2
0                   0
1                   0
2                   0
```

## ⌄ Selecting K: elbow plot

Check the acurracy of the model using k-fold cross-validation

```
Sum of squared error (SSE): 152.34795176035797

1  # Intialize a list to hold sum of squared error (sse)
2  # Initialize a list to hold SSE values
3  sse_list = []
4
5  # Define values of k
6  k_values = range(1, 11)  # Testing k from 1 to 10
7
8  # For each k
9  # For each k, fit KMeans and calculate SSE
10 for k in k_values:
11     km = KMeans(n_clusters=k, n_init="auto")
12     km.fit(df[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']])
13     sse_list.append(km.inertia_)
14
15 # Plot SSE vs. k
16 plt.figure(figsize=(8, 5))
17 plt.plot(k_values, sse_list, marker='o')
18 plt.title('Elbow Plot for K-Means Clustering')
19 plt.xlabel('Number of clusters (k)')
20 plt.ylabel('Sum of Squared Errors (SSE)')
21 plt.show()
```

### Elbow Plot for K-Means Clustering