**Task 2: Prediction Using Unsupervised ML**

- Predict the optimal Number of Clusters and represent it visually.

**Done By: Harini G**

*Importing Required Libraries*

In [3]:
```python
1  #importing libraries
2
3  import numpy as np
4  import pandas as pd
5  from sklearn.preprocessing import StandardScaler, LabelEncoder
6
7  #importing visualization libraries
8
9  import matplotlib.pyplot as plt
10 import seaborn as sns
11 import plotly.express as px
12 import plotly.graph_objects as pgo
```

**Importing Dataset**

In [6]:
```python
1  df=pd.read_csv("D:/Harini(christ unniversity)/Internship/Iris.csv")
2  df.head(5)
```

Out[6]:

|   | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|----|---------------|--------------|---------------|--------------|---------|
| 0 | 1  | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2  | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3  | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4  | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5  | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

In [7]:
```python
1  ##dropping the column Id
2
3  df=df.drop('Id',axis=1)
```

In [8]:
```python
1  X=df.drop("Species",axis=1)
2  y=df['Species']
```

*Exploratory Data Analysis*

In [9]:
```python
1  ##Finding the shape of the dataset
2
3  X.shape
```

Out[9]: (150, 4)

The dataset contain 150 instances with 4 features.

In [10]:
```python
1  ##finding the information of the dataset
2
3  X.info()
```
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 4 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   SepalLengthCm  150 non-null    float64
 1   SepalWidthCm   150 non-null    float64
 2   PetalLengthCm  150 non-null    float64
 3   PetalWidthCm   150 non-null    float64
dtypes: float64(4)
memory usage: 4.8 KB
```

```
1  ##finding the NaN values of the dataset
2
3  X.isna().sum()
```

Out[11]: 
```
SepalLengthCm    0
SepalWidthCm     0
PetalLengthCm    0
PetalWidthCm     0
dtype: int64
```
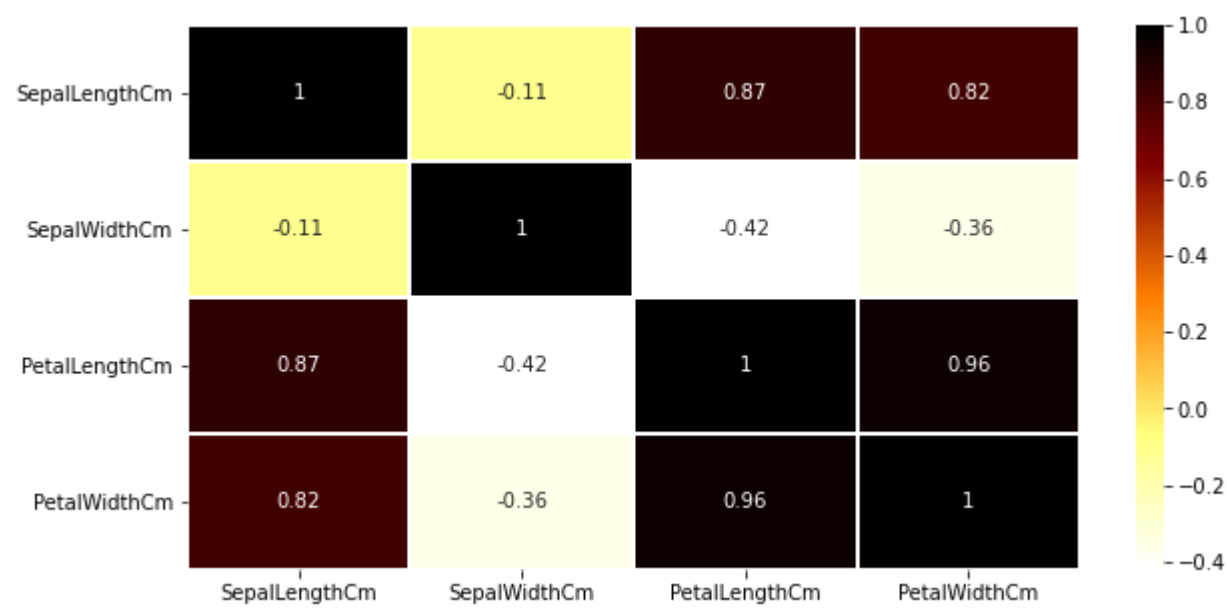
There is no Null values or NaN values in the dataset.

In [12]:
```
1  ##Taking the description of the dataset
2
3  X.describe()
```

Out[12]:

|       | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-------|---------------|--------------|---------------|--------------|
| count | 150.000000    | 150.000000   | 150.000000    | 150.000000   |
| mean  | 5.843333      | 3.054000     | 3.758667      | 1.198667     |
| std   | 0.828066      | 0.433594     | 1.764420      | 0.763161     |
| min   | 4.300000      | 2.000000     | 1.000000      | 0.100000     |
| 25%   | 5.100000      | 2.800000     | 1.600000      | 0.300000     |
| 50%   | 5.800000      | 3.000000     | 4.350000      | 1.300000     |
| 75%   | 6.400000      | 3.300000     | 5.100000      | 1.800000     |
| max   | 7.900000      | 4.400000     | 6.900000      | 2.500000     |

*Correlation*

In [13]:
```
1  plt.figure(figsize = (10, 5))
2  sns.heatmap(X.corr(), linecolor = 'white', linewidths = 1, cmap ='afmhot_r', annot = True)
3  plt.show()
```



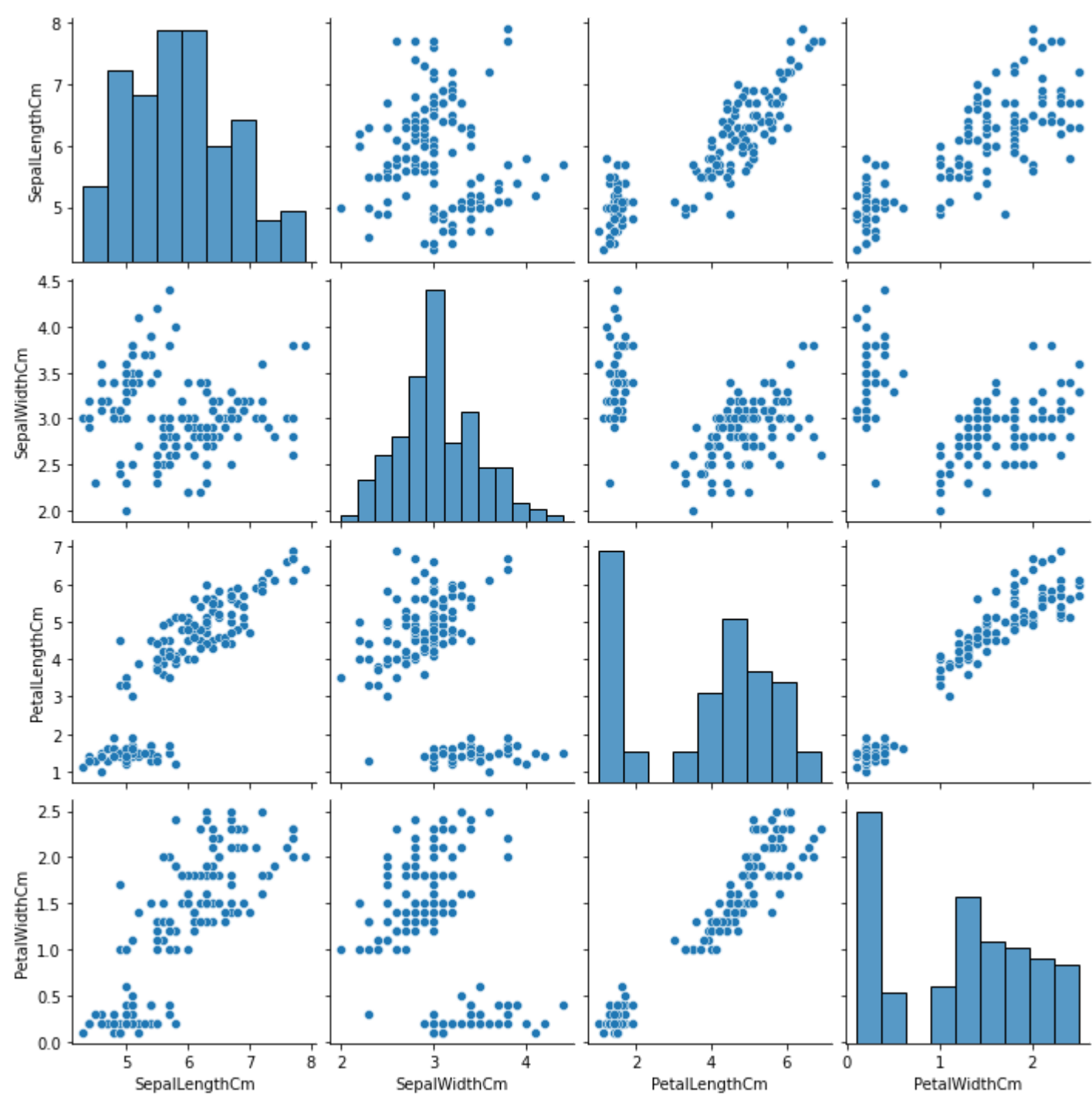Here, Petal Length and Petal Width are highly correlated.

*Pairplot*

```
In [14]:    1  sns.pairplot(X)
```

Out[14]:   `<seaborn.axisgrid.PairGrid at 0x20233583f98>`


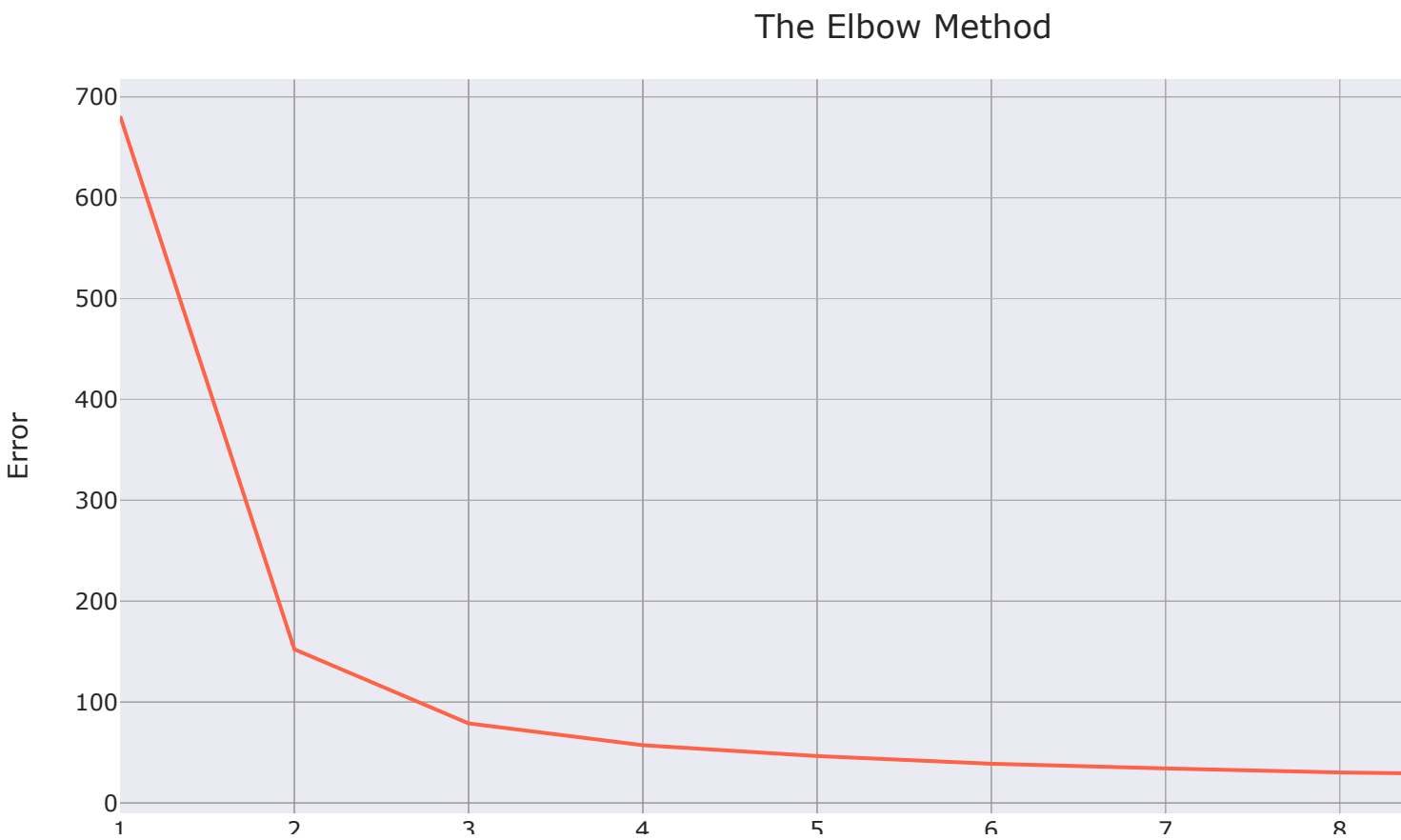
```
In [15]:    1  x = X.iloc[:, [0, 1, 2, 3]].values
```

```
In [16]:    1  import warnings
            2  warnings.filterwarnings("ignore")
```

**Finding the optimum number of clusters for k-means clustering**

```
In [17]:    1  #Finding the optimum number of clusters for k-means clustering
            2  from sklearn.cluster import KMeans
            3  Error =[]
            4  for i in range(1, 11):
            5      kmeans = KMeans(n_clusters = i).fit(x)
            6      kmeans.fit(X)
            7      Error.append(kmeans.inertia_)
```

```
In [18]:   1  uu=pd.DataFrame(Error)
           2  uu.columns=['Error']
```

```
In [19]:   1  fig =px.line(x=range(1, 11),y=uu['Error'],color_discrete_sequence=[ "tomato"])
           2
           3  fig.update_layout(title = {'text':'The Elbow Method',
           4                                 'y':0.95,
           5                                 'x':0.5},
           6                      xaxis_title='No.of Clusters',
           7                      yaxis_title='Error',
           8                      template='seaborn'
           9                      )
          10
          11  fig.show()
```
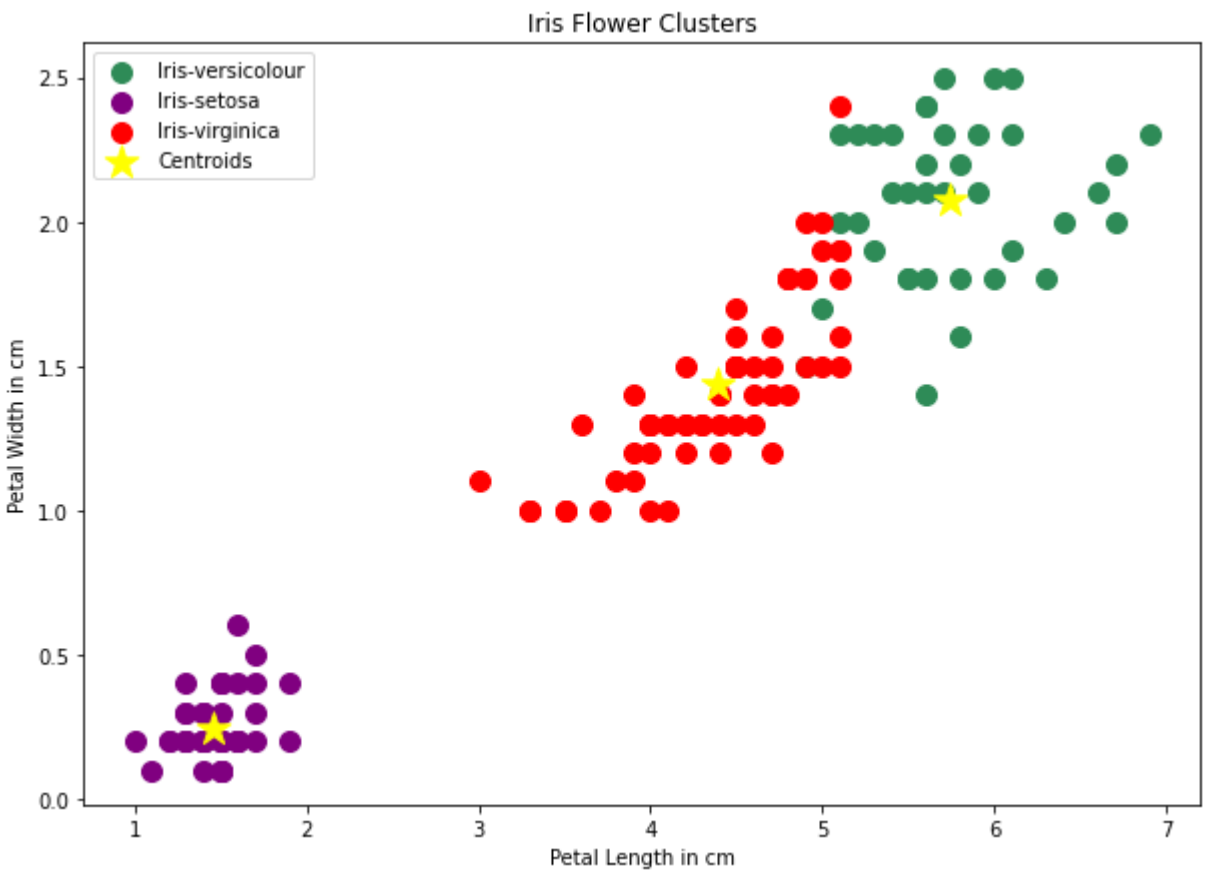


The Optimal Number of Cluter is 3.

*K-Means Clustering*

```
In [20]:   1  kmeans = KMeans(n_clusters=3)
           2  y_kmeans = kmeans.fit_predict(x)
           3  print(y_kmeans)
           4
           5  kmeans.cluster_centers_
```

```
[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 0 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2 2 0 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 0 2 0 0 0 0 2 0 0 0 0
 0 0 2 2 0 0 0 0 2 0 2 0 2 0 0 2 2 0 0 0 0 0 2 0 0 0 0 2 0 0 0 2 0 0 0 2 0
 0 2]
```

```
Out[20]: array([[6.85      , 3.07368421, 5.74210526, 2.07105263],
                [5.006     , 3.418     , 1.464     , 0.244     ],
                [5.9016129 , 2.7483871 , 4.39354839, 1.43387097]])
```

```
1  fig = plt.figure(figsize=(10, 7))
2  plt.title('Clusters with Centroids',fontweight ='bold', fontsize=20)
3  plt.scatter(x[y_kmeans == 0, 2], x[y_kmeans == 0, 3], s = 100, c = 'seagreen', label = 'Iris-v
4  plt.scatter(x[y_kmeans == 1, 2], x[y_kmeans == 1, 3], s = 100, c = 'purple', label = 'Iris-set
5  plt.scatter(x[y_kmeans == 2, 2], x[y_kmeans == 2, 3],s = 100, c = 'red', label = 'Iris-virgini
6  plt.scatter(kmeans.cluster_centers_[:, 2], kmeans.cluster_centers_[:,3], s = 300, c = 'yellow'
7             label = 'Centroids')
8  plt.title('Iris Flower Clusters')
9  plt.ylabel('Petal Width in cm')
10 plt.xlabel('Petal Length in cm')
11 plt.legend()
```

Out[21]:  <matplotlib.legend.Legend at 0x20236c7ef60>



In [ ]:

```
1
```