SPARKS FOUNDATION

GRIP JUNE 2022

NAME: DHARANI GUNTUPALLI

DATA SCIENCE AND BUSINESS ANALYTICS INTERN

TASK 2: PREDICTION USING UNSUPERVISED ML

From the given 'Iris' dataset, predict the optimum number of clusters and represent it visually.

In [1]:
```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib as mpl
import seaborn as sns
from sklearn import datasets
```

In [2]:
```python
iris=pd.read_csv('Iris.csv')
```

In [3]:
```python
iris.head(10)
```

Out[3]:

|   | 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 |
| 5 | 6 | 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa |
| 6 | 7 | 4.6 | 3.4 | 1.4 | 0.3 | Iris-setosa |
| 7 | 8 | 5.0 | 3.4 | 1.5 | 0.2 | Iris-setosa |
| 8 | 9 | 4.4 | 2.9 | 1.4 | 0.2 | Iris-setosa |
| 9 | 10 | 4.9 | 3.1 | 1.5 | 0.1 | Iris-setosa |

In [4]:
```python
iris.tail(10)
```

Out[4]:

|   | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|-----|---------------|--------------|---------------|--------------|---------|
| 140 | 141 | 6.7 | 3.1 | 5.6 | 2.4 | Iris-virginica |
| 141 | 142 | 6.9 | 3.1 | 5.1 | 2.3 | Iris-virginica |
| 142 | 143 | 5.8 | 2.7 | 5.1 | 1.9 | Iris-virginica |
| 143 | 144 | 6.8 | 3.2 | 5.9 | 2.3 | Iris-virginica |
| 144 | 145 | 6.7 | 3.3 | 5.7 | 2.5 | Iris-virginica |
| 145 | 146 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 147 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 148 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 149 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 150 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

```
In [6]:  iris.isnull().sum()

Out[6]:  Id               0
         SepalLengthCm    0
         SepalWidthCm     0
         PetalLengthCm    0
         PetalWidthCm     0
         Species          0
         dtype: int64
```

```
In [7]:  iris.shape

Out[7]:  (150, 6)
```

```
In [8]:  iris.describe()
```

Out[8]:

|       | Id         | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-------|------------|---------------|--------------|---------------|--------------|
| count | 150.000000 | 150.000000    | 150.000000   | 150.000000    | 150.000000   |
| mean  | 75.500000  | 5.843333      | 3.054000     | 3.758667      | 1.198667     |
| std   | 43.445368  | 0.828066      | 0.433594     | 1.764420      | 0.763161     |
| min   | 1.000000   | 4.300000      | 2.000000     | 1.000000      | 0.100000     |
| 25%   | 38.250000  | 5.100000      | 2.800000     | 1.600000      | 0.300000     |
| 50%   | 75.500000  | 5.800000      | 3.000000     | 4.350000      | 1.300000     |
| 75%   | 112.750000 | 6.400000      | 3.300000     | 5.100000      | 1.800000     |
| max   | 150.000000 | 7.900000      | 4.400000     | 6.900000      | 2.500000     |

```
In [9]:  iris.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   Id             150 non-null    int64
 1   SepalLengthCm  150 non-null    float64
 2   SepalWidthCm   150 non-null    float64
 3   PetalLengthCm  150 non-null    float64
 4   PetalWidthCm   150 non-null    float64
 5   Species        150 non-null    object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```
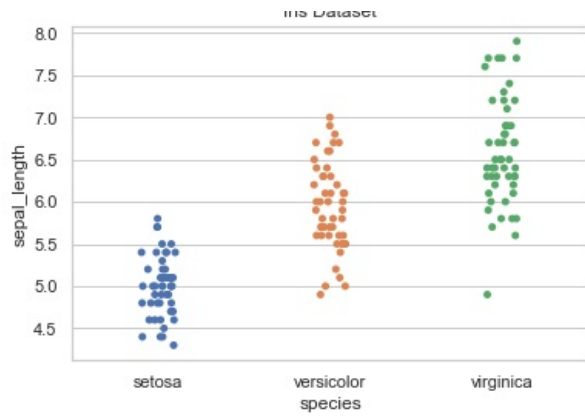
## CHECK FOR UNIQUE CLASS IN DATASET

```
In [10]:  print(iris.Species.nunique())
          print(iris.Species.value_counts())

3
Iris-versicolor    50
Iris-setosa        50
Iris-virginica     50
Name: Species, dtype: int64
```
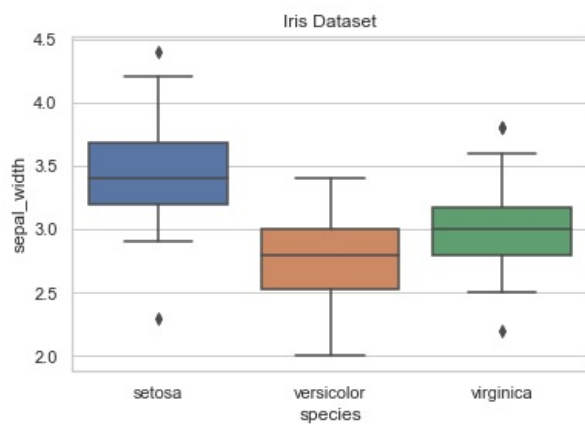
```
In [13]:  sns.set(style ='whitegrid')
          iris=sns.load_dataset('iris');
          ax=sns.stripplot(x='species',y='sepal_length',data= iris);
          plt.title('Iris Dataset')
          plt.show()

                      Iris Dataset
```
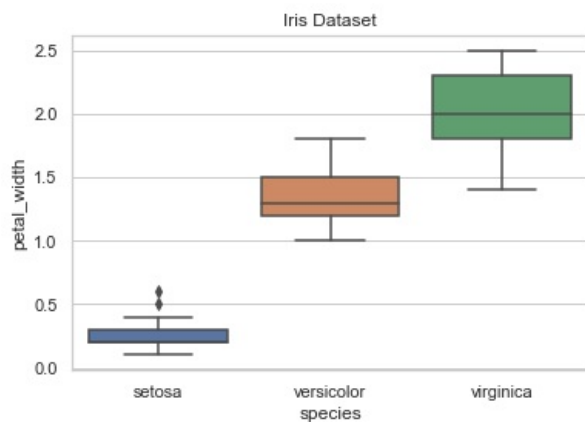
# BOX PLOT FOR SPECIES TO SEPAL_WIDTH

```
In [14]: sns.boxplot(x='species',y='sepal_width',data=iris)
         plt.title("Iris Dataset")
         plt.show()
```
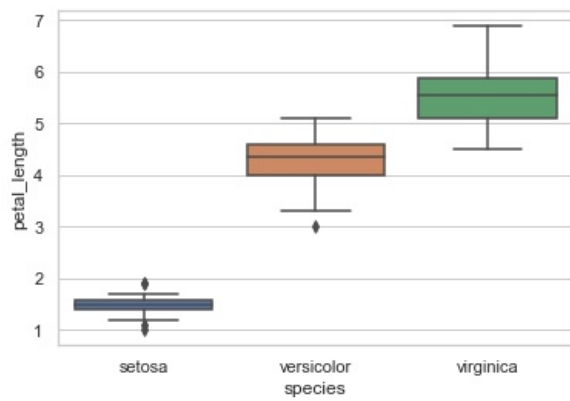


# BOX PLOT FOR SPECIES TO PETAL_WIDTH

```
In [15]: sns.boxplot(x='species',y='petal_width',data=iris)
         plt.title("Iris Dataset")
         plt.show()
```
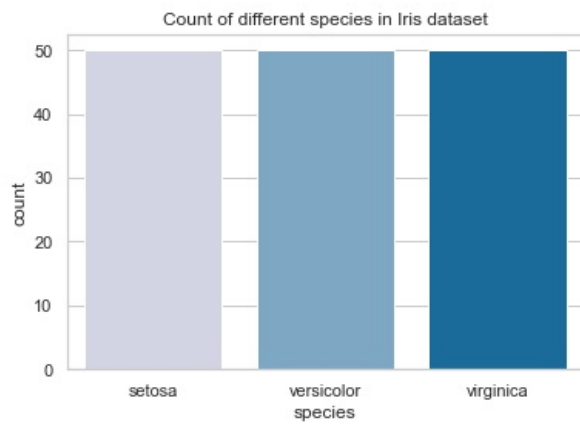


# BOX PLOT FOR SPECIES TO PETAL_LENGTH

```
In [16]: sns.boxplot(x='species',y='petal_length',data=iris)
         plt.title("Iris Dataset")
         plt.show()
```
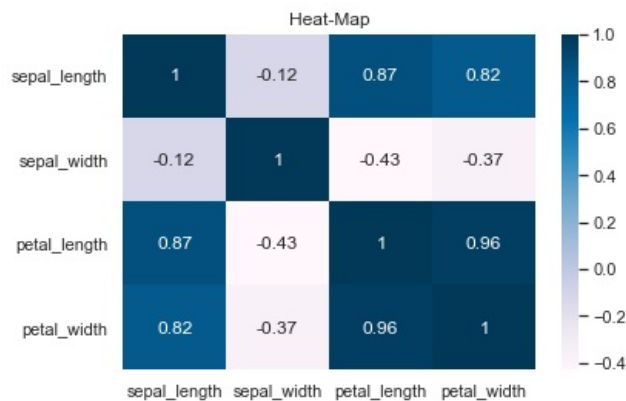
Iris Dataset

## COUNT PLOT

```python
sns.countplot(x='species',data=iris,palette="PuBu")
plt.title("Count of different species in Iris dataset")
plt.show()
```



## DETERMINING THE RELATIONSHIP BETWEEN TWO VARIABLE BY ANALYSING

```python
sns.heatmap(iris.corr(),annot=True,cmap='PuBu')
plt.title("Heat-Map")
plt.show()
```



## FINDING CLUSTERS USING K-MEANS

```python
x=iris.iloc[:,[0,1,2,3]].values
from sklearn.cluster import KMeans
wcss=[]
for i in range(1,11):
```
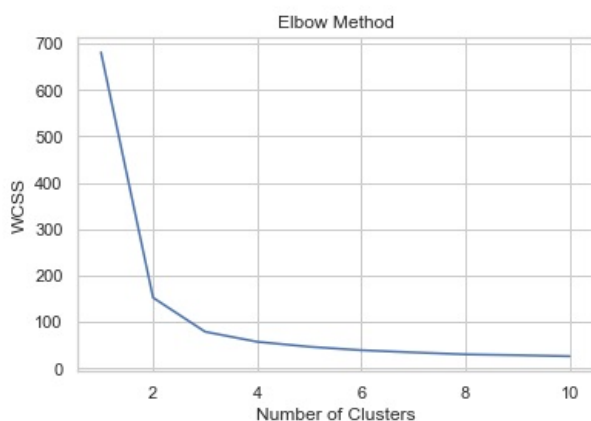
```
kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
kmeans.fit(x)
wcss.append(kmeans.inertia_)
print('k:',i ,"wcss:",kmeans.inertia_)
```

C:\Users\hp\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:881: UserWarning: KMeans is known to have a me
mory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the
environment variable OMP_NUM_THREADS=1.
  warnings.warn(

```
k: 1 wcss: 681.3705999999996
k: 2 wcss: 152.34795176035797
k: 3 wcss: 78.851441426146
k: 4 wcss: 57.22847321428572
k: 5 wcss: 46.47223015873018
k: 6 wcss: 39.03998724608725
k: 7 wcss: 34.299712121212146
k: 8 wcss: 30.063110617452732
k: 9 wcss: 28.27172172856384
k: 10 wcss: 26.094324740540422
```

## PLOTTING RESULTS ON LINE GRAPH

In [22]:
```
plt.plot(range(1,11),wcss)
plt.title('Elbow Method')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS')
plt.show()
```



## FITTING K-MEANS TO THE DATA SET

In [23]:
```
kmeans = KMeans(n_clusters = 3, init = 'k-means++',max_iter = 300, n_init = 10, random_state = 0)
y_kmeans = kmeans.fit_predict(x)
```
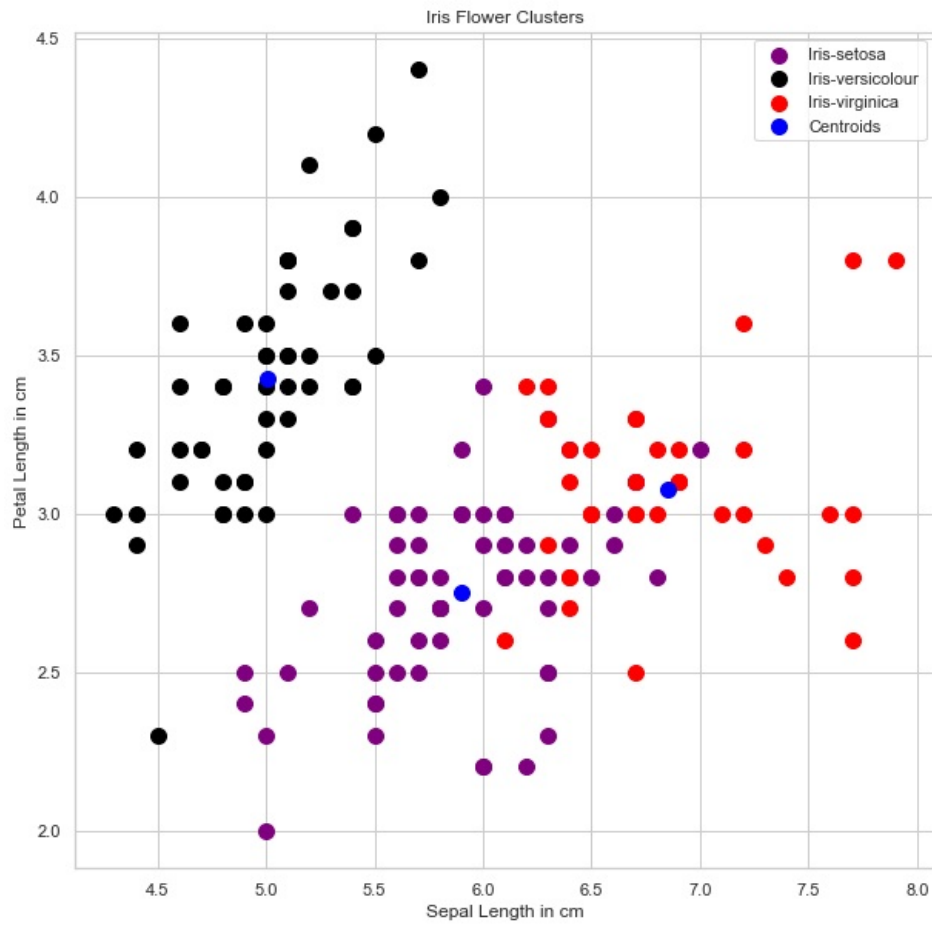
In [26]:
```
y_kmeans
```

Out[26]:
```
array([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, 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 2, 2, 2, 2, 0, 2, 2, 2,
       2, 2, 2, 0, 0, 2, 2, 2, 2, 0, 2, 0, 2, 0, 2, 2, 0, 0, 2, 2, 2, 2,
       2, 0, 2, 2, 2, 2, 0, 2, 2, 2, 0, 2, 2, 2, 0, 2, 2, 0])
```

## VISUALISING AND PLOTTING THE CLUSTERS

In [24]:
```
plt.figure(figsize=(10,10))
plt.scatter(x[y_kmeans==0,0],x[y_kmeans==0,1],s=100,c='purple',label='Iris-setosa')
plt.scatter(x[y_kmeans==1,0],x[y_kmeans==1,1],s=100,c='black',label='Iris-versicolour')
```

```
plt.scatter(x[y_kmeans==2,0],x[y_kmeans==2,1],s=100,c='red',label='Iris-virginica')
plt.scatter(kmeans.cluster_centers_[:,0],kmeans.cluster_centers_[:,1],s=100,c='blue',label='Centroids')
plt.title('Iris Flower Clusters')
plt.xlabel('Sepal Length in cm')
plt.ylabel('Petal Length in cm')
plt.legend()
plt.show()
```



In [ ]: