

▼ Prediction using Unsupervised ML (Level - Beginner)

Er. Narayan Patil

TSF GRIP Data Science & Business Analytics

Tasks 2 :- From the given 'Iris' dataset, predict the optimum number of clusters and represent it visually.

Dataset Link :- <https://bit.ly/3kXTdox>

Importing all the libraries

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import datasets
```

Data available at the link - '<https://bit.ly/3kXTdox>'

Reading data from Github url

Loading and Reading the iris dataset

```
url = 'https://raw.githubusercontent.com/PatilNarayan/Data\_Science\_-\_Business\_Anal'
```

```
data = pd.read_csv(url)
print('Data import successfull')
```

Data import successfull

#loads the first five rows

```
print(data.head())
```

	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

Checking for Not a Number (NaN) values

```
print(data.isna().sum())
```

Id	0
SepalLengthCm	0
SepalWidthCm	0
PetalLengthCm	0
PetalWidthCm	0

```
Species      0
dtype: int64
```

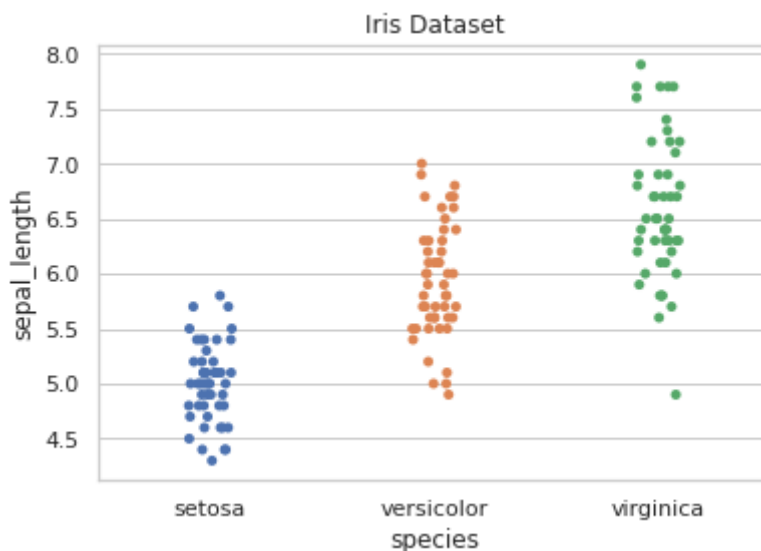
```
# Checking statistical description
print(data.describe())
```

	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

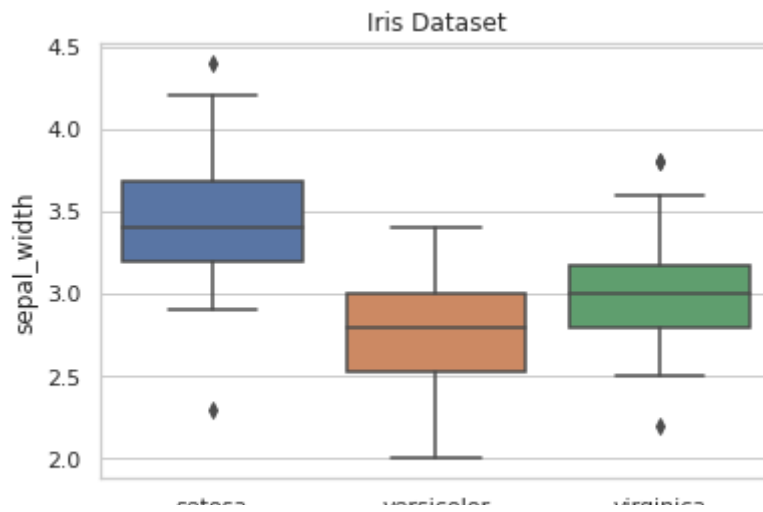
```
# Check for unique classes in the dataset.
print(data.Species.nunique())
print(data.Species.value_counts())
```

```
3
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
Name: Species, dtype: int64
```

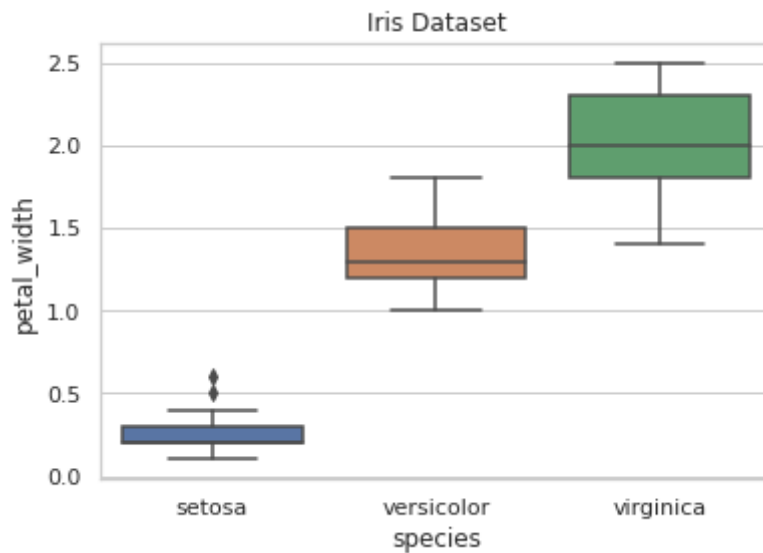
```
# Data Visualization
sns.set(style = 'whitegrid')
iris = sns.load_dataset('iris');
ax = sns.stripplot(x = 'species', y = 'sepal_length', data = iris);
plt.title('Iris Dataset')
plt.show()
```



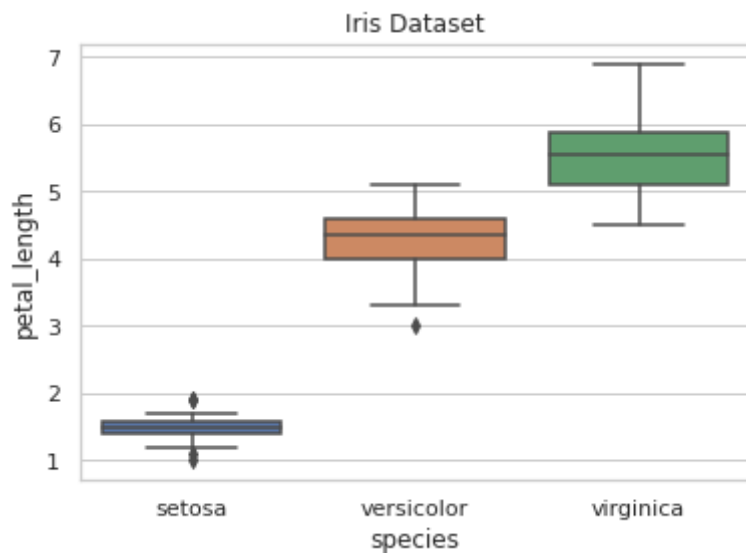
```
sns.boxplot(x='species',y='sepal_width',data=iris)
plt.title("Iris Dataset")
plt.show()
```



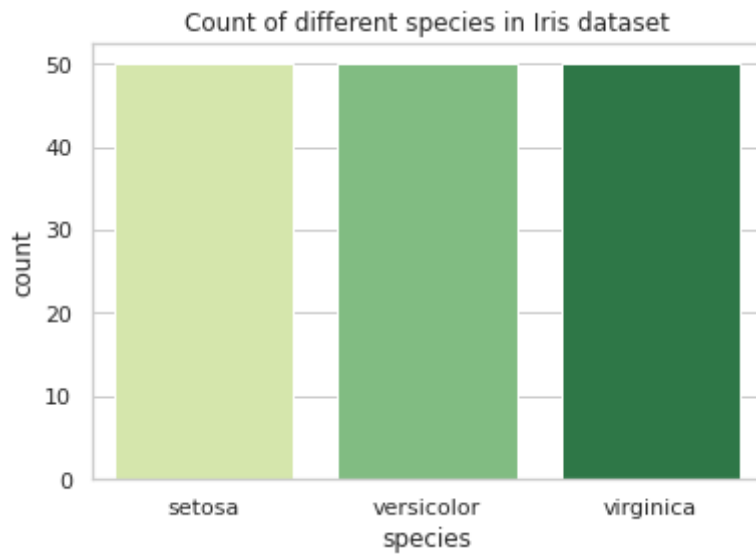
```
sns.boxplot(x='species',y='petal_width',data=iris)
plt.title("Iris Dataset")
plt.show()
```



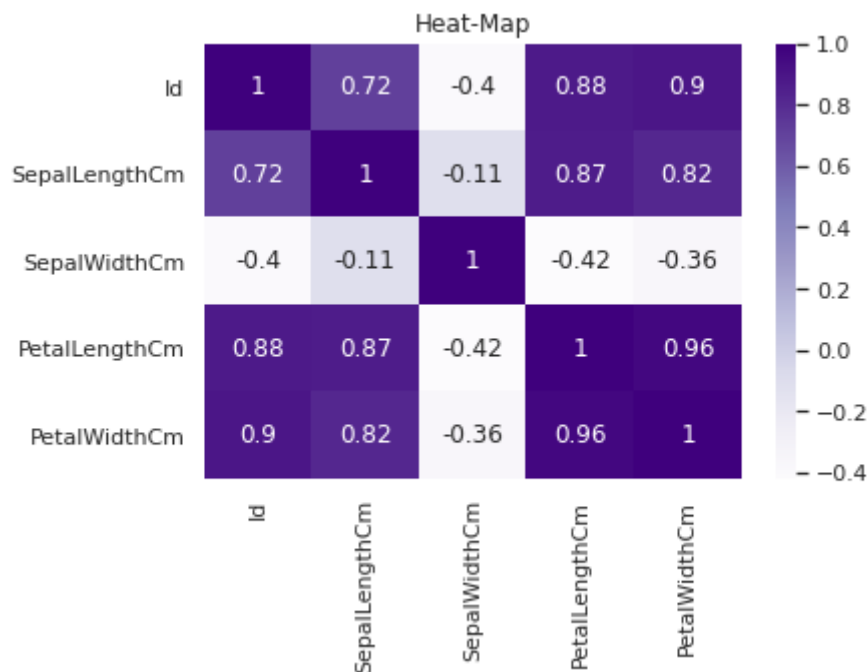
```
sns.boxplot(x='species',y='petal_length',data=iris)
plt.title("Iris Dataset")
plt.show()
```



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



```
# Heat Map
sns.heatmap(data.corr(), annot=True, cmap='Purples')
plt.title("Heat-Map")
plt.show()
```



```
# Finding the optimum number of clusters using k-means
```

```
x = data.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)
```

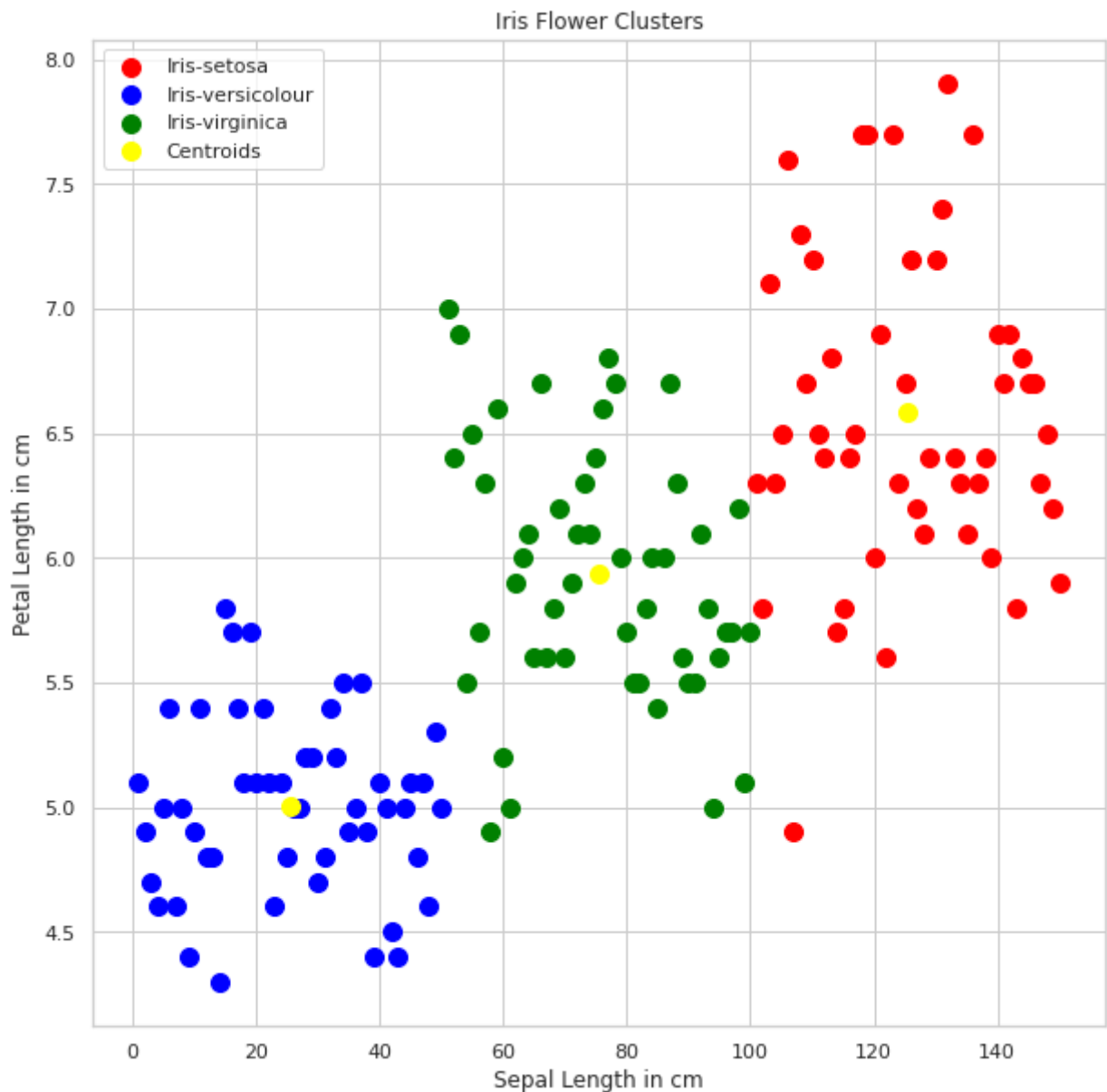
```
# Plotting the results onto a line graph, allowing us to observe 'The elbow'
```

Number of Clusters	WCSS (Approximate)
1	280,000
2	70,000
3	30,000
4	15,000
5	10,000
6	8,000
7	6,000
8	5,000
9	4,000
10	3,000

```
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, 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, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], dtype=int32)
```

```
# Visualising the clusters
plt.figure(figsize=(10,10))
plt.scatter(x[y_kmeans==0,0],x[y_kmeans==0,1],s=100,c='red',label='Iris-setosa')
plt.scatter(x[y_kmeans==1,0],x[y_kmeans==1,1],s=100,c='blue',label='Iris-versicolour')
plt.scatter(x[y_kmeans==2,0],x[y_kmeans==2,1],s=100,c='green',label='Iris-virginica')

# Plotting the centroids of the clusters
plt.scatter(kmeans.cluster_centers_[0,0],kmeans.cluster_centers_[0,1],s=100,c='yellow',label='Centroids')
plt.title('Iris Flower Clusters')
plt.xlabel('Sepal Length in cm')
plt.ylabel('Petal Length in cm')
plt.legend()
plt.show()
```



✓

0s

completed at 15:39

×