

```
In [35]: # Importing the necessary Libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
```

```
In [36]: df = pd.read_csv("Iris.csv")
```

```
In [37]: # Reading the head of our data
df.head()
```

Out[37]:

	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 [38]: df.Species.value_counts()
```

Out[38]: Iris-setosa 50
Iris-versicolor 50
Iris-virginica 50
Name: Species, dtype: int64

```
In [5]: df.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
```

```
In [39]: # Checking for Null Values
df.isnull().sum()
```

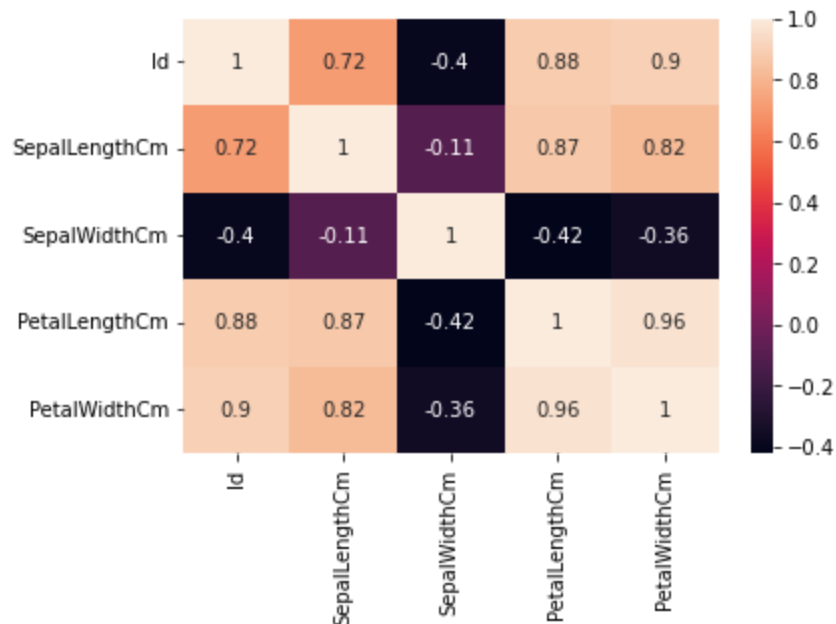
Out[39]: Id 0
SepalLengthCm 0
SepalWidthCm 0
PetalLengthCm 0
PetalWidthCm 0
Species 0
dtype: int64

```
In [40]: # Dropping the duplicate rows if any
```

```
df.drop_duplicates(inplace=True)
```

```
In [42]: # Checking for correlation
corr = df.corr()
sns.heatmap(corr, annot=True)
```

Out[42]: <AxesSubplot:>



```
In [43]: # Selecting the values for predicting the clusters
X = df.iloc[:, :-1].values
```

```
In [11]: # Importing the cluster
from sklearn.cluster import KMeans
```

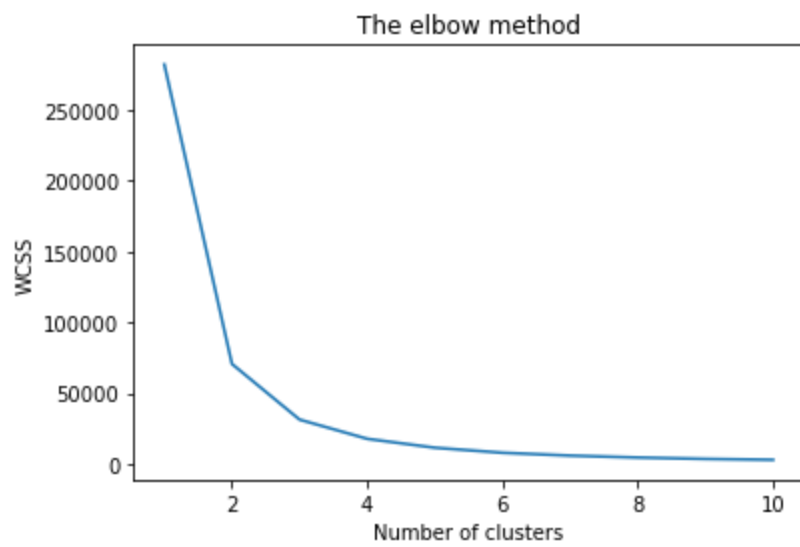
```
In [12]: # Finding the appropriate number of clusters using Elbow Method

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_)

# Plotting the results onto a line graph,
# `allowing us to observe 'The elbow'
plt.plot(range(1, 11), wcss)
plt.title('The elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS') # Within cluster sum of squares
plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036: UserWarning: K Means is known to have a memory 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(



After the third point we can see there is a smooth decline in the number of clusters so its safe to say that appropriate number of clusters for this data is 3

We will be applying the number of clusters as 3 to the data

```
In [13]: kmeans = KMeans(n_clusters = 3, init = 'k-means++',
                        max_iter = 300, n_init = 10, random_state = 0)
kmeans.fit(X)
```

```
Out[13]: KMeans(n_clusters=3, random_state=0)
```

```
In [14]: y_kmeans = kmeans.predict(X)
```

```
In [15]: y_kmeans
```

```
Out[15]: array([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, 0, 0, 0, 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, 2, 2, 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, 1])
```

```
In [16]: labels = kmeans.labels
```

```
In [17]: labels
```

```
Out[17]: array([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, 0, 0, 0, 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, 2, 2, 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, 1])
```

```
In [18]: kmeans.inertia
```

```
Out[18]: 31326.886799999997
```

In [19]:

```
centroids = kmeans.cluster_centers_
```

```
In [20]: centroids
```

```
Out[20]: array([[ 25.5 ,  5.006,  3.418,  1.464,  0.244],  
        [125.5 ,  6.588,  2.974,  5.552,  2.026],  
        [ 75.5 ,  5.936,  2.77 ,  4.26 ,  1.326]])
```

```
In [21]: # Visualising the clusters - On the first two columns  
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(centroids[:, 0], centroids[:,1],  
            s = 100, c = 'yellow', label = 'Centroids')  
  
plt.legend()
```

```
Out[21]: <matplotlib.legend.Legend at 0x18292bcaa30>
```

