

# K-Means Clustering on IRIS Dataset

```
In [185]: %matplotlib inline

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [186]: df = pd.read_csv("IRIS.csv")
df.head()
```

Out[186]:

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

We have object data type so assignment and conversion integer type.

```
In [187]: df['species'].unique()
```

Out[187]: array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)

```
In [188]: df['species'] = df['species'].replace({'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica': 2})
df.head()
```

Out[188]:

	sepal_length	sepal_width	petal_length	petal_width	species
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

```
In [189]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
sepal_length    150 non-null float64
sepal_width     150 non-null float64
petal_length    150 non-null float64
petal_width     150 non-null float64
species         150 non-null int64
dtypes: float64(4), int64(1)
memory usage: 5.9 KB
```

```
In [190]: df.describe()
```

Out[190]:

	sepal_length	sepal_width	petal_length	petal_width	species
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667	1.000000
std	0.828066	0.433594	1.764420	0.763161	0.819232
min	4.300000	2.000000	1.000000	0.100000	0.000000
25%	5.100000	2.800000	1.600000	0.300000	0.000000
50%	5.800000	3.000000	4.350000	1.300000	1.000000
75%	6.400000	3.300000	5.100000	1.800000	2.000000
max	7.900000	4.400000	6.900000	2.500000	2.000000

```
In [191]: df.corr()
```

Out[191]:

	sepal_length	sepal_width	petal_length	petal_width	species
sepal_length	1.000000	-0.109369	0.871754	0.817954	0.782561
sepal_width	-0.109369	1.000000	-0.420516	-0.356544	-0.419446
petal_length	0.871754	-0.420516	1.000000	0.962757	0.949043
petal_width	0.817954	-0.356544	0.962757	1.000000	0.956464
species	0.782561	-0.419446	0.949043	0.956464	1.000000

```
In [192]: import seaborn as sns
sns.heatmap(df.corr(), annot=True)
```

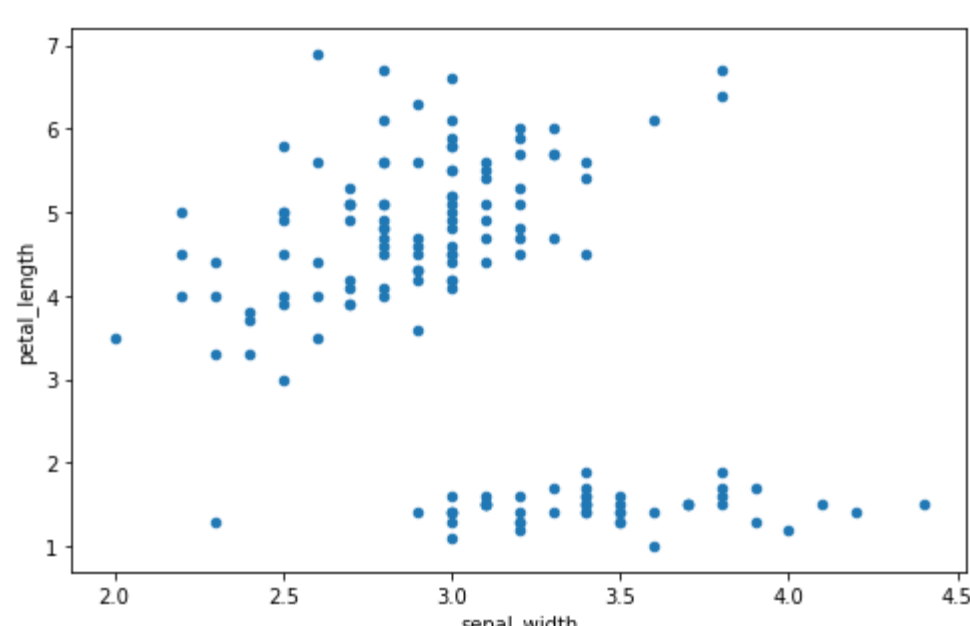
Out[192]: <matplotlib.axes.\_subplots.AxesSubplot at 0x211b9c8a780>



## Picking Dimensions which intuitively form clusters

```
In [193]: df.plot.scatter("sepal_width", "petal_length", figsize=(8,5))
```

Out[193]: <matplotlib.axes.\_subplots.AxesSubplot at 0x211b9c712e8>

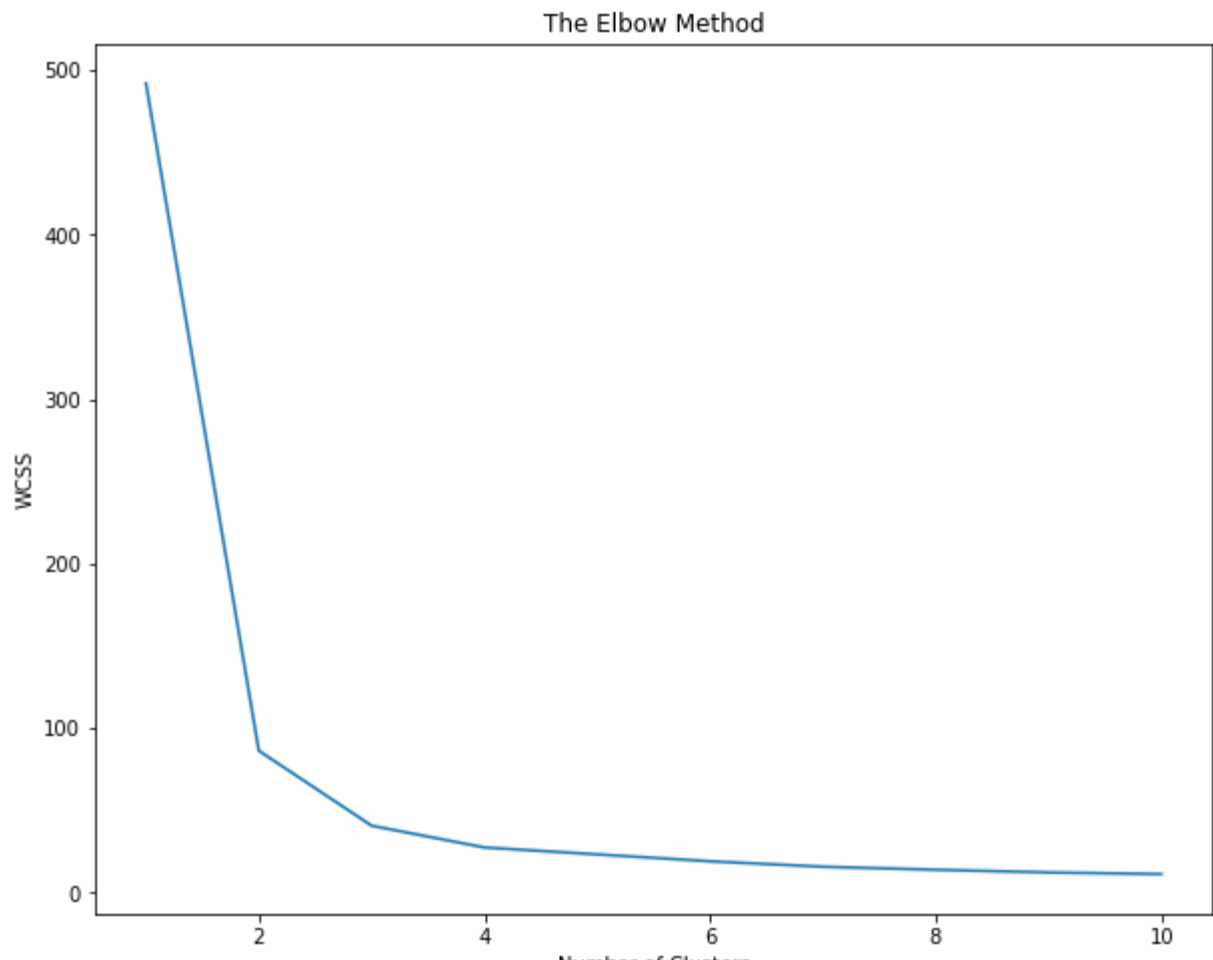


```
In [194]: X = df.iloc[:, [1, 2]].values
```

```
In [195]: from sklearn.cluster import KMeans

fig = plt.figure(figsize=(10, 8))
WCSS = []
for i in range(1, 11):
    clf = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    clf.fit(X)
    WCSS.append(clf.inertia_)

plt.plot(range(1, 11), WCSS)
plt.title('The Elbow Method')
plt.ylabel('WCSS')
plt.xlabel('Number of Clusters')
plt.show()
```



## Picking number of clusters on the basis of Elbow Method Graph

```
In [196]: clf = KMeans(n_clusters=2, init='k-means++', max_iter=300, n_init=10, random_state=0)
y_kmeans = clf.fit_predict(X)
```

```
In [197]: fig = plt.figure(figsize=(10, 8))
plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], color='red', s=60, label='Cluster 1', edgecolors='black')
plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], color='green', s=60, label='Cluster 2', edgecolor='black')
# cluster centres
plt.scatter(clf.cluster_centers_[0, 0], clf.cluster_centers_[0, 1], color='magenta', s=100, label='Centroid', edgecolors='black')
plt.legend()
plt.title('Clusters using KMeans')
plt.ylabel('Sepal Width')
plt.xlabel('Petal Length')
plt.show()
```

