

PREDICTION USING UNSUPERVISED MACHINE LEARNING

AUTHOR: Aman Agrawal

TASK

Iris dataset is given.We need to predict the optimum number of clusters and represent it visually.

```
In [1]: #filter warnings
import warnings
warnings.filterwarnings('ignore')

In [2]: from sklearn import datasets
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
import matplotlib.patches as mpatches
import sklearn.metrics as sm
from mpl_toolkits.mplot3d import Axes3D
from scipy.cluster.hierarchy import linkage, dendrogram
from sklearn.cluster import DBSCAN
from sklearn.decomposition import PCA

In [3]: #loading the csv data into a data frame
iris=pd.read_csv("Iris.csv")
print("Data is imported successfully")

Data is imported successfully
```

Exploratory Data Analysis

```
In [4]: #Head of the data
iris.head()

Out[4]:
```

	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 [5]: iris.tail()

Out[5]:
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
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]: # Checking the dimension of the data
iris.shape

Out[6]: (150, 6)

In [7]: # Checking the column information
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

In [8]: # Viewing the column names
iris.columns

Out[8]: Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm',
        'Species'],
        dtype='object')

In [9]: # Summary of all numerical data
iris.describe()

Out[9]:
```

	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 [10]: # Checkingfor missing values if any
iris.isnull().sum()

Out[10]: Id                0
SepalLengthCm            0
SepalWidthCm              0
PetalLengthCm            0
PetalWidthCm             0
Species                  0
dtype: int64

In [11]: iris=iris.drop(['Id'],axis=1)

In [12]: iris.head()

Out[12]:
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	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

```
In [13]: iris["Species"].value_counts()

Out[13]: Iris-setosa      50
Iris-versicolor    50
Iris-virginica     50
Name: Species, dtype: int64
```

Pre-processing standardization

```
In [14]: #Scaling the data before clustering

X=iris.iloc[:, :4]
X

Out[14]:
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
...
145	6.7	3.0	5.2	2.3
146	6.3	2.5	5.0	1.9
147	6.5	3.0	5.2	2.0
148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8

150 rows × 4 columns

```
In [15]: #scaling the features with standard scaler

from sklearn.preprocessing import StandardScaler

scale=StandardScaler()

#fitting
scale.fit(X)

Out[15]: StandardScaler
StandardScaler()

In [16]: x_scaled=scale.transform(X)

In [17]: x_scaled_df=pd.DataFrame(x_scaled,columns=X.columns)

In [18]: #data after scaling

x_scaled_df

Out[18]:
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
0	-0.900681	1.032057	-1.341272	-1.312977
1	-1.143017	-0.124958	-1.341272	-1.312977
2	-1.385353	0.337848	-1.398138	-1.312977
3	-1.506521	0.106445	-1.284407	-1.312977
4	-1.021849	1.263460	-1.341272	-1.312977
...
145	1.038005	-0.124958	0.819624	1.447956
146	0.553333	-1.281972	0.705893	0.922064
147	0.795669	-0.124958	0.819624	1.053537
148	0.432165	0.800654	0.933356	1.447956
149	0.068662	-0.124958	0.762759	0.790591

150 rows × 4 columns

```
In [19]: # Finding the optimum number of clusters for k-means classification

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_)

# 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()

The elbow method

700
600
500
400
300
200
100
0
0 2 4 6 8 10
Number of clusters
WCSS
```

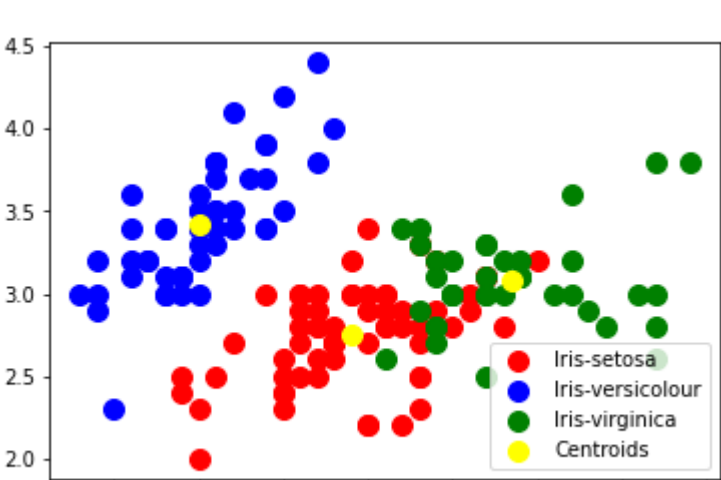
```
In [20]: # Applying kmeans clustering to the dataset
kmeans = KMeans(n_clusters = 3, init = 'k-means++',
                max_iter = 300, n_init = 10, random_state = 0)
y_kmeans = kmeans.fit_predict(x)

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(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:,1],
            s = 100, c = 'yellow', label = 'Centroids')

plt.legend()

Out[21]: <matplotlib.legend.Legend at 0x20abaafd180>
```



This shows the clusters present in the given dataset among species setosa, versicolour, virginica