

# Machine Learning Fundamentals Lab-9

Name: Gaurav Prasanna

Reg No: 19BEC1315

## Aim:

- To implement K Means Clustering Algorithm on given dataset and predict the outcomes without support of classes or labels.
- To implement, find the centroids of clusters and to plot them.

## Software Required:

- Jupyter Notebook
- Anaconda Navigator

**Libraries Required:** Numpy, Matplotlib, Seaborn, Sci-kit Learn, Pandas

## Code and Outputs:

- Train Ticket Dataset K Means

```
In [1]: # Dependencies
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import MinMaxScaler
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

In [2]: # Load the train and test datasets to create two DataFrames

train_url = "train.csv"
train = pd.read_csv(train_url)
test_url = "test.csv"
test = pd.read_csv(test_url)

In [3]: print("***** Train_Set *****")
print(train.head())
print("\n")
print("***** Test_Set *****")
print(test.head())
```

\*\*\*\* Train\_Set \*\*\*\*

|   | PassengerId | Survived | Pclass | \ |
|---|-------------|----------|--------|---|
| 0 | 1           | 0        | 3      |   |
| 1 | 2           | 1        | 1      |   |
| 2 | 3           | 1        | 3      |   |
| 3 | 4           | 1        | 1      |   |
| 4 | 5           | 0        | 3      |   |

|   | Name  | Sex    | Age  | SibSp | \ |
|---|---|--------|------|-------|---|
| 0 | Braund, Mr. Owen Harris                           | male   | 22.0 | 1     |   |
| 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1     |   |
| 2 | Heikkinen, Miss. Laina                            | female | 26.0 | 0     |   |
| 3 | Futrelle, Mrs. Jacques Heath (Lily May Peel)      | female | 35.0 | 1     |   |
| 4 | Allen, Mr. William Henry                          | male   | 35.0 | 0     |   |

|   | Parch | Ticket           | Fare    | Cabin | Embarked |
|---|-------|------------------|---------|-------|----------|
| 0 | 0     | A/5 21171        | 7.2500  | NaN   | S        |
| 1 | 0     | PC 17599         | 71.2833 | C85   | C        |
| 2 | 0     | STON/O2. 3101282 | 7.9250  | NaN   | S        |
| 3 | 0     | 113803           | 53.1000 | C123  | S        |
| 4 | 0     | 373450           | 8.0500  | NaN   | S        |

\*\*\*\* Test\_Set \*\*\*\*

|   | PassengerId | Pclass | Name   | Sex    | \ |
|---|-------------|--------|--|--------|---|
| 0 | 892         | 3      | Kelly, Mr. James                             | male   |   |
| 1 | 893         | 3      | Wilkes, Mrs. James (Ellen Needs)             | female |   |
| 2 | 894         | 2      | Myles, Mr. Thomas Francis                    | male   |   |
| 3 | 895         | 3      | Wirz, Mr. Albert                             | male   |   |
| 4 | 896         | 3      | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female |   |

|   | Age  | SibSp | Parch | Ticket  | Fare    | Cabin | Embarked |
|---|------|-------|-------|---------|---------|-------|----------|
| 0 | 34.5 | 0     | 0     | 330911  | 7.8292  | NaN   | Q        |
| 1 | 47.0 | 1     | 0     | 363272  | 7.0000  | NaN   | S        |
| 2 | 62.0 | 0     | 0     | 240276  | 9.6875  | NaN   | Q        |
| 3 | 27.0 | 0     | 0     | 315154  | 8.6625  | NaN   | S        |
| 4 | 22.0 | 1     | 1     | 3101298 | 12.2875 | NaN   | S        |

```
In [4]: print("**** Train_Set ****")
print(train.describe())
print("\n")
print("**** Test_Set ****")
print(test.describe())
```

\*\*\*\* Train\_Set \*\*\*\*

|       | PassengerId | Survived   | Pclass     | Age        | SibSp      | \ |
|-------|-------------|------------|------------|------------|------------|---|
| count | 891.000000  | 891.000000 | 891.000000 | 714.000000 | 891.000000 |   |
| mean  | 446.000000  | 0.383838   | 2.306642   | 29.699118  | 0.523008   |   |
| std   | 257.353842  | 0.486592   | 0.836071   | 14.526497  | 1.102743   |   |
| min   | 1.000000    | 0.000000   | 1.000000   | 0.420000   | 0.000000   |   |
| 25%   | 223.500000  | 0.000000   | 2.000000   | 20.125000  | 0.000000   |   |
| 50%   | 446.000000  | 0.000000   | 3.000000   | 28.000000  | 0.000000   |   |
| 75%   | 668.500000  | 1.000000   | 3.000000   | 38.000000  | 1.000000   |   |
| max   | 891.000000  | 1.000000   | 3.000000   | 80.000000  | 8.000000   |   |

|       | Parch      | Fare       |
|-------|------------|------------|
| count | 891.000000 | 891.000000 |
| mean  | 0.381594   | 32.204208  |
| std   | 0.806057   | 49.693429  |
| min   | 0.000000   | 0.000000   |
| 25%   | 0.000000   | 7.910400   |
| 50%   | 0.000000   | 14.454200  |
| 75%   | 0.000000   | 31.000000  |
| max   | 6.000000   | 512.329200 |

\*\*\*\* Test\_Set \*\*\*\*

|       | PassengerId | Pclass     | Age        | SibSp      | Parch      | Fare       |
|-------|-------------|------------|------------|------------|------------|------------|
| count | 418.000000  | 418.000000 | 332.000000 | 418.000000 | 418.000000 | 417.000000 |
| mean  | 1100.500000 | 2.265550   | 30.272590  | 0.447368   | 0.392344   | 35.627188  |
| std   | 120.810458  | 0.841838   | 14.181209  | 0.896760   | 0.981429   | 55.907576  |
| min   | 892.000000  | 1.000000   | 0.170000   | 0.000000   | 0.000000   | 0.000000   |
| 25%   | 996.250000  | 1.000000   | 21.000000  | 0.000000   | 0.000000   | 7.895800   |
| 50%   | 1100.500000 | 3.000000   | 27.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%   | 1204.750000 | 3.000000   | 39.000000  | 1.000000   | 0.000000   | 31.500000  |
| max   | 1309.000000 | 3.000000   | 76.000000  | 8.000000   | 9.000000   | 512.329200 |

```
In [5]: print(train.columns.values)

['PassengerId' 'Survived' 'Pclass' 'Name' 'Sex' 'Age' 'SibSp' 'Parch'
'Ticket' 'Fare' 'Cabin' 'Embarked']
```

```
In [6]: # For the train set
train.isna().head()
```

```
Out[6]:
```

|   | PassengerId | Survived | Pclass | Name  | Sex   | Age   | SibSp | Parch | Ticket | Fare  | Cabin | Embarked |
|---|-------------|----------|--------|-------|-------|-------|-------|-------|--------|-------|-------|----------|
| 0 | 0           | False    | False  | False | False | False | False | False | False  | False | True  | False    |
| 1 | 1           | False    | False  | False | False | False | False | False | False  | False | False | False    |
| 2 | 2           | False    | False  | False | False | False | False | False | False  | False | True  | False    |
| 3 | 3           | False    | False  | False | False | False | False | False | False  | False | False | False    |
| 4 | 4           | False    | False  | False | False | False | False | False | False  | False | True  | False    |

```
In [7]: # For the test set
test.isna().head()
```

```
Out[7]:
```

|   | PassengerId | Pclass | Name  | Sex   | Age   | SibSp | Parch | Ticket | Fare  | Cabin | Embarked |
|---|-------------|--------|-------|-------|-------|-------|-------|--------|-------|-------|----------|
| 0 | 0           | False  | False | False | False | False | False | False  | False | True  | False    |
| 1 | 1           | False  | False | False | False | False | False | False  | False | True  | False    |
| 2 | 2           | False  | False | False | False | False | False | False  | False | True  | False    |
| 3 | 3           | False  | False | False | False | False | False | False  | False | True  | False    |
| 4 | 4           | False  | False | False | False | False | False | False  | False | True  | False    |

```
In [8]: print("*****In the train set*****")
print(train.isna().sum())
print("\n")
print("*****In the test set*****")
print(test.isna().sum())
```

```
*****In the train set*****
```

```
PassengerId      0
Survived          0
Pclass            0
Name              0
Sex              0
Age              177
SibSp             0
Parch             0
Ticket            0
Fare              0
Cabin            687
Embarked          2
dtype: int64
```

```
*****In the test set*****
```

```
PassengerId      0
Pclass            0
Name              0
Sex              0
Age              86
SibSp             0
Parch             0
Ticket            0
Fare              1
Cabin            327
Embarked          0
dtype: int64
```

---

```
In [9]: # Fill missing values with mean column values in the train set
train.fillna(train.mean(), inplace=True)
# Fill missing values with mean column values in the test set
test.fillna(test.mean(), inplace=True)
```

```
In [10]: print(train.isna().sum())
print(test.isna().sum())
```

```
PassengerId      0
Survived          0
Pclass           0
Name             0
Sex              0
Age             0
SibSp            0
Parch            0
Ticket           0
Fare            0
Cabin           687
Embarked         2
dtype: int64
PassengerId      0
Pclass           0
Name             0
Sex              0
Age             0
SibSp            0
Parch            0
Ticket           0
Fare            0
Cabin           327
dtype: int64
```

```
In [11]: train['Ticket'].head()
```

```
Out[11]: 0      A/5 21171
1      PC 17599
2      STON/O2. 3101282
3      113803
4      373450
Name: Ticket, dtype: object
```

```
In [12]: train['Cabin'].head()
```

```
Out[12]: 0      NaN
1      C85
2      NaN
3      C123
4      NaN
Name: Cabin, dtype: object
```

```
In [13]: train[['Pclass', 'Survived']].groupby(['Pclass'], as_index=False).mean().sort_values(by='Survived', ascending=False)
```

```
Out[13]:
```

|   | Pclass | Survived |
|---|--------|----------|
| 0 | 1      | 0.629630 |
| 1 | 2      | 0.472826 |
| 2 | 3      | 0.242363 |

```
In [14]: train[['Sex', 'Survived']].groupby(['Sex'], as_index=False).mean().sort_values(by='Survived', ascending=False)
```

```
Out[14]:
```

|   | Sex    | Survived |
|---|--------|----------|
| 0 | female | 0.742038 |
| 1 | male   | 0.188908 |

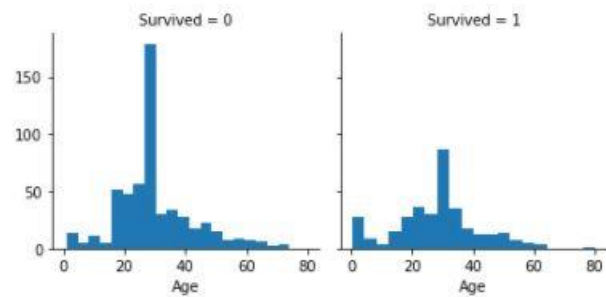
```
In [15]: train[['SibSp', 'Survived']].groupby(['SibSp'], as_index=False).mean().sort_values(by='Survived', ascending=False)
```

```
Out[15]:
```

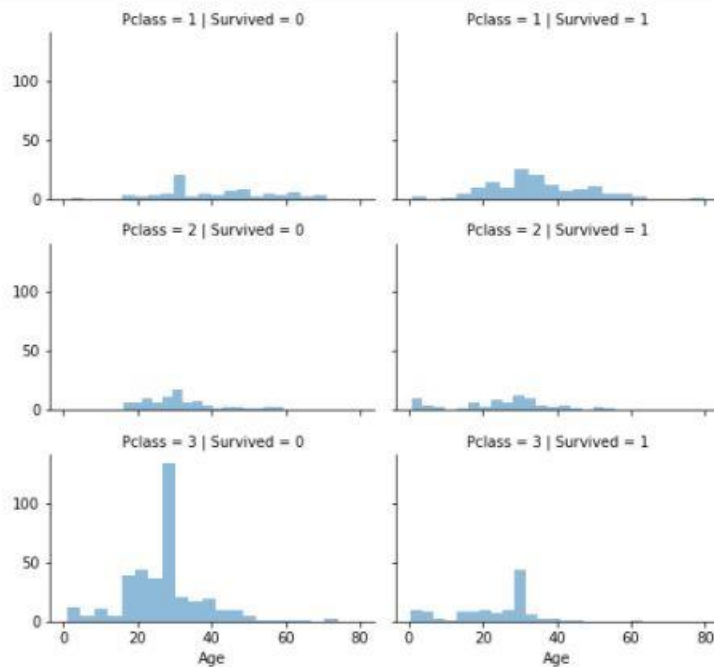
|   | SibSp | Survived |
|---|-------|----------|
| 1 | 1     | 0.535885 |
| 2 | 2     | 0.464286 |
| 0 | 0     | 0.345395 |
| 3 | 3     | 0.250000 |
| 4 | 4     | 0.166667 |
| 5 | 5     | 0.000000 |
| 6 | 8     | 0.000000 |

```
In [16]: g = sns.FacetGrid(train, col='Survived')
g.map(plt.hist, 'Age', bins=20)
```

```
Out[16]: <seaborn.axisgrid.FacetGrid at 0x187eef29160>
```



```
In [17]: grid = sns.FacetGrid(train, col='Survived', row='Pclass', height=2.2, aspect=1.6)
grid.map(plt.hist, 'Age', alpha=.5, bins=20)
grid.add_legend();
```





```
In [18]: train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column             Non-Null Count  Dtype
---  ---
0   PassengerId         891 non-null    int64
1   Survived            891 non-null    int64
2   Pclass              891 non-null    int64
3   Name                891 non-null    object
4   Sex                 891 non-null    object
5   Age                 891 non-null    float64
6   SibSp               891 non-null    int64
7   Parch              891 non-null    int64
8   Ticket              891 non-null    object
9   Fare                891 non-null    float64
10  Cabin               204 non-null    object
11  Embarked            889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
In [19]: train = train.drop(['Name', 'Ticket', 'Cabin', 'Embarked'], axis=1)
test = test.drop(['Name', 'Ticket', 'Cabin', 'Embarked'], axis=1)
```

```
In [20]: labelEncoder = LabelEncoder()
labelEncoder.fit(train['Sex'])
labelEncoder.fit(test['Sex'])
train['Sex'] = labelEncoder.transform(train['Sex'])
test['Sex'] = labelEncoder.transform(test['Sex'])
```

```
In [21]: train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 8 columns):
#   Column             Non-Null Count  Dtype
---  ---
0   PassengerId         891 non-null    int64
1   Survived            891 non-null    int64
2   Pclass              891 non-null    int64
3   Sex                 891 non-null    int32
4   Age                 891 non-null    float64
5   SibSp               891 non-null    int64
6   Parch              891 non-null    int64
7   Fare                891 non-null    float64
dtypes: float64(2), int32(1), int64(5)
memory usage: 52.3 KB
```

```
In [22]: test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 7 columns):
#   Column             Non-Null Count  Dtype
---  ---
0   PassengerId         418 non-null    int64
1   Pclass              418 non-null    int64
2   Sex                 418 non-null    int32
3   Age                 418 non-null    float64
4   SibSp               418 non-null    int64
5   Parch              418 non-null    int64
6   Fare                418 non-null    float64
```

```
In [23]: X = np.array(train.drop(['Survived'], 1).astype(float))
        y = np.array(train['Survived'])
```

```
In [24]: train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   PassengerId      891 non-null    int64
1   Survived         891 non-null    int64
2   Pclass           891 non-null    int64
3   Sex              891 non-null    int32
4   Age              891 non-null    float64
5   SibSp            891 non-null    int64
6   Parch           891 non-null    int64
7   Fare             891 non-null    float64
dtypes: float64(2), int32(1), int64(5)
memory usage: 52.3 KB
```

```
In [25]: kmeans = KMeans(n_clusters=2) # You want cluster the passenger records into 2: Survived or Not survived
        kmeans.fit(X)
```

```
Out[25]: KMeans(n_clusters=2)
```

```
In [26]: correct = 0
        for i in range(len(X)):
            predict_me = np.array(X[i].astype(float))
            predict_me = predict_me.reshape(-1, len(predict_me))
            prediction = kmeans.predict(predict_me)
            if prediction[0] == y[i]:
                correct += 1

        print(correct/len(X))

0.5084175084175084
```

```
In [27]: kmeans = KMeans(n_clusters=2, max_iter=600, algorithm = 'auto')
        kmeans.fit(X)
```

```
Out[27]: KMeans(max_iter=600, n_clusters=2)
```

```
In [28]: correct = 0
        for i in range(len(X)):
            predict_me = np.array(X[i].astype(float))
            predict_me = predict_me.reshape(-1, len(predict_me))
            prediction = kmeans.predict(predict_me)
            if prediction[0] == y[i]:
                correct += 1

        print(correct/len(X))

0.49158249158249157
```

```
In [29]: scaler = MinMaxScaler()
        X_scaled = scaler.fit_transform(X)
        kmeans.fit(X_scaled)
```

```
Out[29]: KMeans(max_iter=600, n_clusters=2)
```

```
In [30]: correct = 0
        for i in range(len(X)):
            predict_me = np.array(X[i].astype(float))
            predict_me = predict_me.reshape(-1, len(predict_me))
            prediction = kmeans.predict(predict_me)
            if prediction[0] == y[i]:
                correct += 1

        print(correct/len(X))

0.37373737373737376
```

## b) Cluster Centroids

```
In [7]: from sklearn.cluster import KMeans
from sklearn import metrics
import numpy as np
import matplotlib.pyplot as plt
```

```
In [16]: x1 = np.array([3, 1, 1, 2, 1, 6, 6, 6, 5, 6, 7, 8, 9, 8, 9, 9, 8])
x2 = np.array([5, 4, 6, 6, 5, 8, 6, 7, 6, 7, 1, 2, 1, 2, 3, 2, 3])

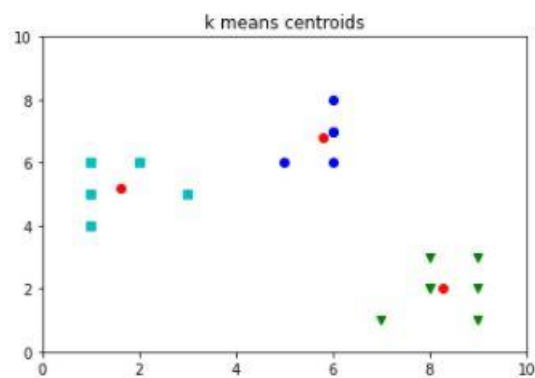
for K in range(2,6,1):
    plt.plot()
    X = np.array(list(zip(x1, x2))).reshape(len(x1), 2)
    colors = ['b', 'g', 'c', 'r', 'y']
    markers = ['o', 'v', 's', '*', '+']

    K = 3
    kmeans_model = KMeans(n_clusters=K).fit(X)
    print(kmeans_model.cluster_centers_)
    centers = np.array(kmeans_model.cluster_centers_)
    plt.plot()
    plt.title('k means centroids')

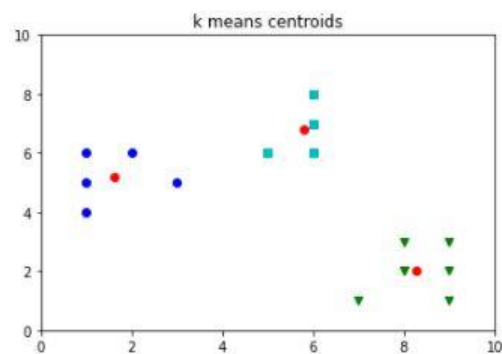
    for i,l in enumerate(kmeans_model.labels_):
        plt.plot(x1[i], x2[i], color=colors[l], marker=markers[l], ls='None')
        plt.xlim([0, 10])
        plt.ylim([0, 10])

    plt.scatter(centers[:, 0], centers[:, 1], color='r')
    plt.show()
```

```
[[5.8      6.8      ]
 [8.28571429 2.      ]
 [1.6      5.2      ]]
```

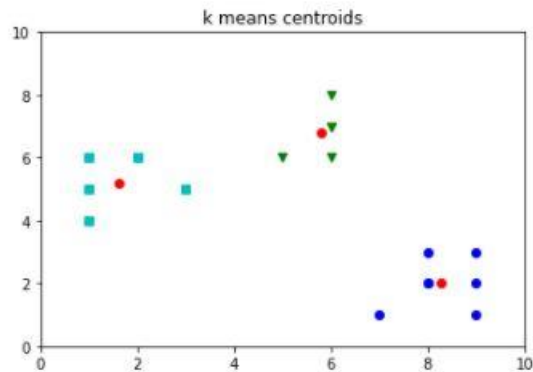


```
[[1.6      5.2      ]
 [8.28571429 2.      ]
 [5.8      6.8      ]]
```

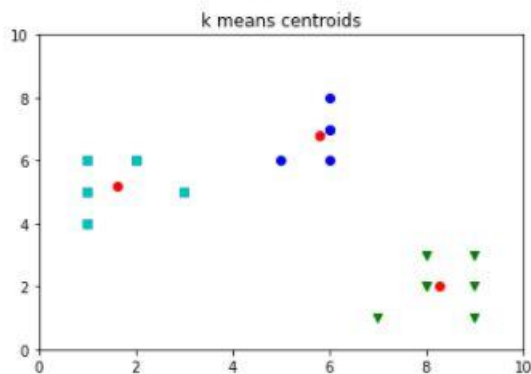




```
[[8.28571429 2. ]
 [5.8       6.8 ]
 [1.6       5.2 ]]
```



```
[[5.8       6.8 ]
 [8.28571429 2. ]
 [1.6       5.2 ]]
```



### Inference:

- From the first one we infer that K means clustering is an unsupervised learning algorithm and it can be used in classifying the classes without labels by looking into the features, and even the prediction along with accuracy of predictions is given.
- From second one we see the plot of centroids has the value of K changes from 2 to 6.

**Result:** K Means on train tickets dataset is implemented and the required plots are visualized using Jupyter notebook.