Assignment-02                    Date: 08/08/2020

Title : Clustering.

Problem : Consider a suitable dataset for clustering
Statement   of data instances in different groups,
      apply different clustering techniques (min 2)
      Visualize clusters using suitable tools.

S/w & H/w : Rstudio / Jupyter Notebook.
Requirements   P IV , 2GB RAM, 500 GB HDD.

Learning Objectives: Use R functions / Scikit-learn functions
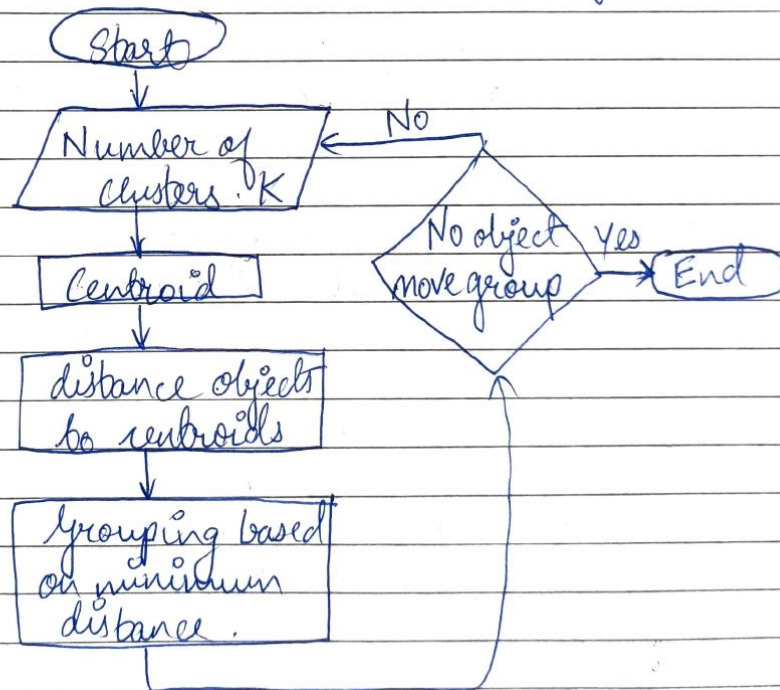      to create K-means clustering models & hehrachical
      clustering models.

Learning Outcomes: Visualize the effects of K means
      & hehrarchical clustering using graphic capabilities.

Theory :-
1] K-means Clustering
   • It is a type of unsupervised learning, which
is used when you have unlabelled data.
• The goal of this algorithm is to find groups
in the data, with number of groups represented
by the variable K.
• The algorithm works iteratively to assign each
data point to one of K group based on features
that are provided.
• Data points are clustered based on feature similarity.
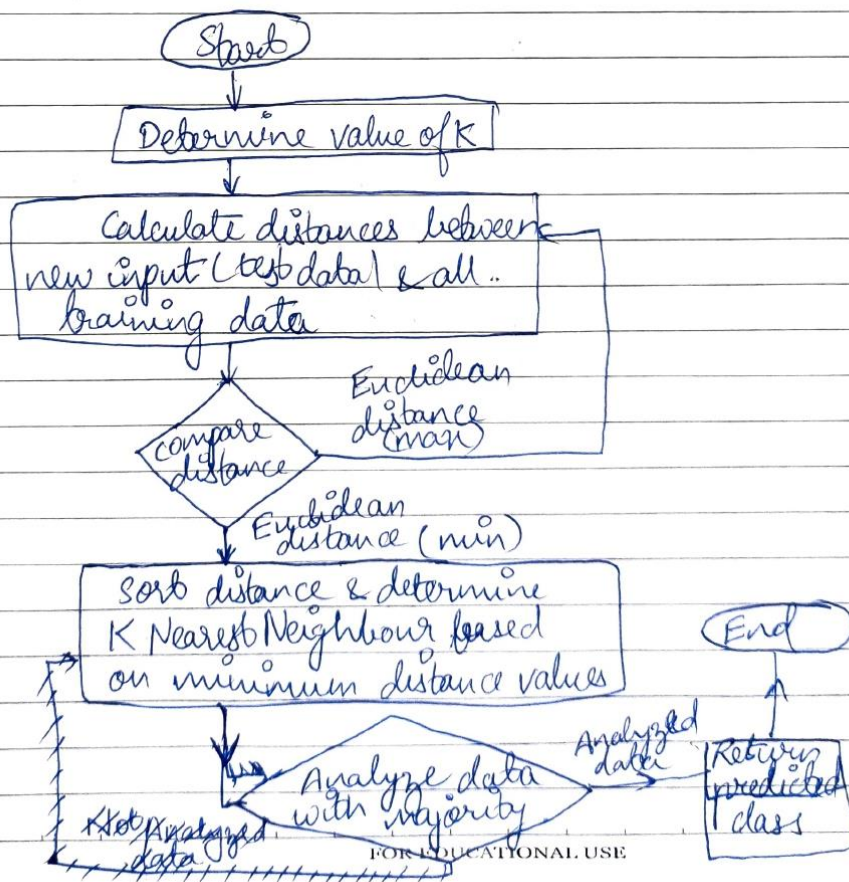• The results of K means clustering algorithm are:

1. The centroids of the K clusters, which can be used to label new data.

2. Labels for training data (each data point is assigned to a single cluster).

- Rather than defining groups before looking at the data, clustering allows you to find & analyze the groups that have been formed organically.

→ Steps to perform K-means Clustering.

```
                    ( Start )
                        |
                        v
    +-------------------+        No
--->| Number of         |<----------+
|   | clusters K        |           |
|   +-------------------+       /---------\
|            |                 /           \
|            v                / No object   \  Yes
|      +-----------+          \ move group  /------->( End )
|      | Centroid  |           \           /
|      +-----------+            \---------/
|            |                      ^
|            v                      |
|   +-----------------+             |
|   | distance object |             |
|   | to centroids    |             |
|   +-----------------+             |
|            |                      |
|            v                      |
|   +-----------------+             |
|   | Grouping based  |             |
+---| on minimum      |-------------+
    | distance.       |
    +-----------------+
```

B] K Nearest Neighbour (KNN) clustering
- It is a supervised classification algorithm.
- It takes bunch of labeled points and uses them to learn how to label other points. To label a

- To label a new point, it looks at the labeled point closest to that new point which are its nearest neighbors, and has those neighbors vote.
- So whichever label, the most of the neighbours have is the label, the most of the neighbors have is the label for the new point.
- Here 'K' in KNN is number of neighbors it checks.
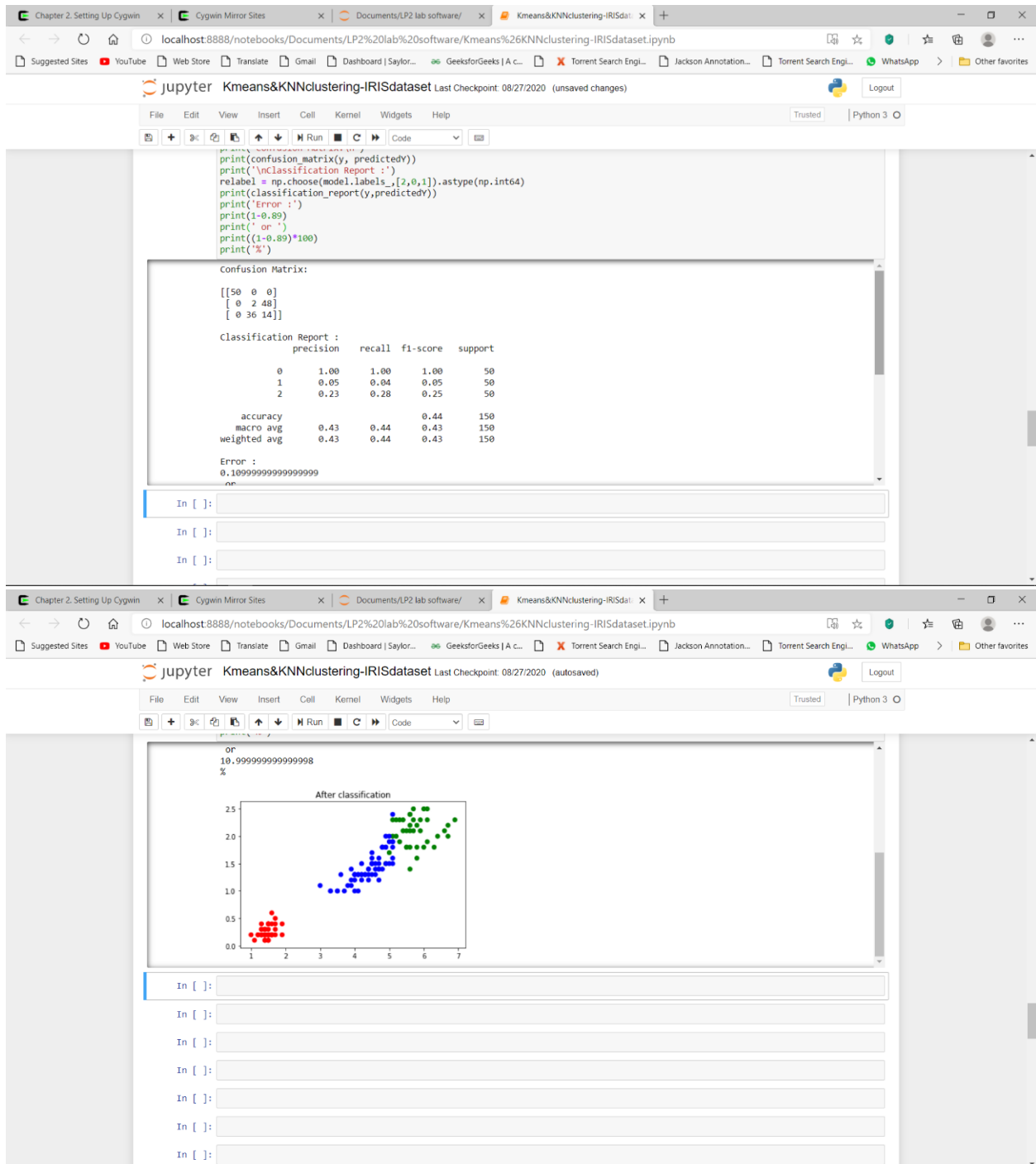- It is supervised because you are trying to classify a point based on the known classification of other points.

```
            ( Start )
                |
                v
    [ Determine value of K ]
                |
                v
    [ Calculate distances between
      new input (test data) & all
      training data ]                ---- Euclidean
                |                         distance
                v                         (max)
           < compare
             distance >
                |
                v  Euclidean
                   distance (min)
    [ Sort distance & determine
      K Nearest Neighbour based      ( End )
      on minimum distance values ]     ^
                |                       |
                v                    [ Return
         < Analyze data >  Analyzed   predicted
           with majority    data       class ]
                |
         Not analyzed
            data
```

Test cases

| Sr no | Description | Expected O/P | Actual O/P |
|---|---|---|---|
| 1) | In KNN clustering method (un-Supervised algorithm) we created confusion matrix & classification report based on Euclidean distance K clusters are formed. | No. of clusters rendered = 5 | success |
| 2) | Visuals cluster using single, complete & average linkages. | Clusters displayed by means of scatter plot | success |
| 3) | While fitting K means to dataset, put random State = 42 | success | success |

Conclusion: Hence, we have successfully implemented hierarchical clustering and K means clustering algorithm in python using jupyter notebooks.

Output:

1)K-Means

# 2)KNN

Chapter 2. Setting Up Cygwin × | Cygwin Mirror Sites × | Documents/LP2 lab software/ × | Kmeans&KNNclustering-IRISdat: × | +

localhost:8888/notebooks/Documents/LP2%20lab%20software/Kmeans%26KNNclustering-IRISdataset.ipynb

Suggested Sites    YouTube    Web Store    Translate    Gmail    Dashboard | Saylor...    GeeksforGeeks | A c...    Torrent Search Engi...    Jackson Annotation...    Torrent Search Engi...    WhatsApp    >    Other favorites

jupyter    Kmeans&KNNclustering-IRISdataset Last Checkpoint: 08/27/2020  (unsaved changes)    Logout

File    Edit    View    Insert    Cell    Kernel    Widgets    Help                    Trusted    Python 3 O

```
Out[172]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                               metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                               weights='uniform')
```

```
In [173]: Y_pred = classifier.predict(X_test)
```

```
In [174]: from sklearn.metrics import classification_report, confusion_matrix
          print("Confusion Matrix:\n")
          print(confusion_matrix(Y_test, Y_pred))
          print("Classification Report:\n")
          print(classification_report(Y_test, Y_pred))
```

```
Confusion Matrix:

[[ 8  0  0]
 [ 0  8  0]
 [ 0  1 13]]
Classification Report:

                 precision    recall  f1-score   support

    Iris-setosa       1.00      1.00      1.00         8
Iris-versicolor       0.89      1.00      0.94         8
 Iris-virginica       1.00      0.93      0.96        14

       accuracy                           0.97        30
      macro avg       0.96      0.98      0.97        30
   weighted avg       0.97      0.97      0.97        30
```

```
In [175]: error = []

          # Calculating error for K values between 1 and 40
          for i in range(1, 40):
              knn = KNeighborsClassifier(n_neighbors=i)
              knn.fit(X_train, Y_train)
              pred_i = knn.predict(X_test)
```

Chapter 2. Setting Up Cygwin × | Cygwin Mirror Sites × | Documents/LP2 lab software/ × | Kmeans&KNNclustering-IRISdat: × | +

localhost:8888/notebooks/Documents/LP2%20lab%20software/Kmeans%26KNNclustering-IRISdataset.ipynb

Suggested Sites    YouTube    Web Store    Translate    Gmail    Dashboard | Saylor...    GeeksforGeeks | A c...    Torrent Search Engi...    Jackson Annotation...    Torrent Search Engi...    WhatsApp    >    Other favorites

jupyter    Kmeans&KNNclustering-IRISdataset Last Checkpoint: 08/27/2020  (unsaved changes)    Logout

File    Edit    View    Insert    Cell    Kernel    Widgets    Help                    Trusted    Python 3 O

```
          knn = KNeighborsClassifier(n_neighbors=i)
          knn.fit(X_train, Y_train)
          pred_i = knn.predict(X_test)
          error.append(np.mean(pred_i != Y_test))
```

```
In [176]: plt.figure(figsize=(12, 6))
          plt.plot(range(1, 40), error, color='red', linestyle='dashed', marker='o',
                   markerfacecolor='blue', markersize=10)
          plt.title('Error Rate K Value')
          plt.xlabel('K Value')
          plt.ylabel('Mean Error')
```

```
Out[176]: Text(0, 0.5, 'Mean Error')
```