K-means clustering algorithm performed on Iris data.

First data preprocessing, of reading the data from iris.txt and then extracting the class label
column from the whole data and X_train has the all the features excluding class label column.

```
3
4      # Fetching Iris Data
5      data = pand.read_csv('iris.txt')
6
7      # Extracting Class label from the data
8      col = data.columns[data.columns.str.startswith('I')]
9
0      # Seperating classlabel and data
1      classlabel = data[col]
2      features = data.drop(data.columns[4], axis=1)
3
4      X_train = features.as_matrix()
5
```

I have used k= 3 because we can evaluate, as there are 3 labels (0 Iris-setosa, 1 Iris-versicolor, 2
Iris-virginica) in the dataset. This can be done intuitively. Also, the data is uniformly distributed
over all the clusters. A good clustering with smaller k such has 3 can have a lower Sum of Squared
Error (SSE) than a poor clustering with higher k.

The below function will calculate the initial centroids to start our algorithm.

```
16     # Fuction for Calculating Initial Centriods
17     def initialCentriod(k=3):
18         np.random.seed(0)
19         cen = np.random.random((k,4))
20         return cen
21
22     initialcen =initialCentriod()*5
23
```

Function to calculate the Euclidean distance from centroid to each data point so as to form a
neighborhood which easily identifies which data point belongs to which cluster.

```
24     # function to calculate the distance between centriod and data points
25     from math import sqrt
26     def pairwisedistance(datapoint,cent):
27         return sqrt(np.sum((datapoint-cent)*(datapoint-cent)))
28
```

KmeansClustering is the function which performs the kmeans clustering algorithm (mentioned
below) using 2 parameters the training data and number of iterations.

1: Select $K$ points as the initial centroids.

2: **repeat**

3:     Form $K$ clusters by assigning all points to the closest centroid.

4:     Recompute the centroid of each cluster.
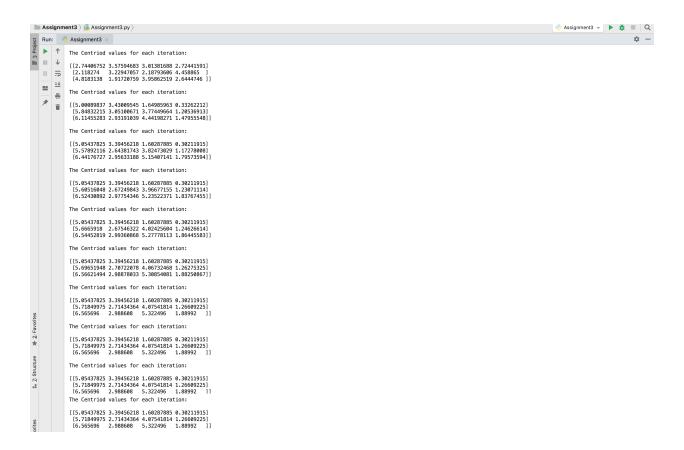
5: **until** The centroids don't change

The Complexity of the above algorithm is O( n * K * I * d ) where n=number of points(), K= number of clusters, I = number of iterations, d = number of attributes.

```python
29      # Function for calculating Kmeans along with training data and iterations
30      def kmeansClustering(X_train,iterations):
31          size = X_train.shape[0]
32          centriod = initialCentriod(3)*5
33          kvalue = centriod.shape[0]
34          distance = np.zeros([size, kvalue])
35          classAssign = np.zeros([size, ])
36          cenn = centriod
37          temp = np.zeros([1, 4])
38          count = 0
39          for t in range(iterations):
40              # print(centriod)
41              for r in range(0, size):
42                  for c in range(0, kvalue):
43                      distance[r][c] = pairwisedistance(X_train[r], centriod[c])
44              classAssign = (np.argmin(distance, axis=1)).reshape((-1,))
45              cenn = np.concatenate((cenn, centriod))
46              print("\nThe Centriod values for each iteration:\n")
47              print(centriod)
48              for c in range(0, kvalue):
49                  temp = np.zeros([1, 4])
50                  count = 0
51                  for r in range(0, size):
52                      temp = temp + (0.98) * (classAssign[r] == c) * X_train[r] + (0.02) * (classAssign[r] != c) * X_train[r]
53                      count = count + (0.98) * (classAssign[r] == c) + (0.02) * (classAssign[r] != c)
54                  centriod[c] = (temp.reshape((-1,))) / (count)
55          return centriod, classAssign, cenn
```

The kmeansClustering results in centroids computed in each iterations along with the assignments of the data points to each cluster.

```python
57      centroid ,classAssign,cenn = kmeansClustering(X_train,iterations=10)
58      print("\nKmeans Result Assignments:\n")
59      print(classAssign)
60
61
```

10 Iterations centroids are shown below:

```
Run:      Assignment3 ×                                                                          ⚙ ─

      The Centriod values for each iteration:

      [[2.74406752 3.57594683 3.01381688 2.72441591]
       [2.118274   3.22947057 2.18793606 4.458865  ]
       [4.8183138  1.91720759 3.95862519 2.6444746 ]]

      The Centriod values for each iteration:

      [[5.00089837 3.43009545 1.64985963 0.33262212]
       [5.84832215 3.05100671 3.77449664 1.20536913]
       [6.11455283 2.93191039 4.44198271 1.47955548]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.57892116 2.64381743 3.82473029 1.17278008]
       [6.44176727 2.95633188 5.15407141 1.79573594]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.60516048 2.67249843 3.96677155 1.23071114]
       [6.52430892 2.97754346 5.23522371 1.83767455]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.6665918  2.67546322 4.02425604 1.24626614]
       [6.54452819 2.99360868 5.27778113 1.86445583]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.69651948 2.70722078 4.06732468 1.26275325]
       [6.56621494 2.98878033 5.30854081 1.88250867]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.71849975 2.71434364 4.07541814 1.26609225]
       [6.565696   2.988608   5.322496   1.88992   ]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.71849975 2.71434364 4.07541814 1.26609225]
       [6.565696   2.988608   5.322496   1.88992   ]]

      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.71849975 2.71434364 4.07541814 1.26609225]
       [6.565696   2.988608   5.322496   1.88992   ]]
      The Centriod values for each iteration:

      [[5.05437825 3.39456218 1.60287885 0.30211915]
       [5.71849975 2.71434364 4.07541814 1.26609225]
       [6.565696   2.988608   5.322496   1.88992   ]]
```

## The result of our algorithm is:

```
Kmeans Result Assignments:

[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 1 2 1 2 1 2 1 1 1 1 1 1 2 1 1 1 1 2 1 2 1 1
 1 2 2 1 1 1 1 2 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 1 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]

Process finished with exit code 0
```

## References:

http://madhugnadig.com/articles/machine-learning/2017/03/04/implementing-k-means-clustering-from-scratch-in-python.html

http://benalexkeen.com/k-means-clustering-in-python/

https://www.kaggle.com/andyxie/k-means-clustering-implementation-in-python

GitHub & Wikipedia