

Weekly Report 2 - K-Means

Ganji Varshitha
AI20BTECH11009

Introduction

K-Means is an unsupervised learning algorithm which performs partitioned clustering. Clustering helps us to understand the structure of the data by grouping it into distinct sub-groups.

Algorithm

It is a parametric method where we need to specify the number of clusters K , which we want to divide the data into.

It is a simple and iterative algorithm described as follows:

```
Initialize K random points from the dataset as centroids
repeat
    Form K clusters by assigning all points to the closest centroid
    Recompute the centroid of each cluster
until The centroids don't change
```

It can be seen as Expectation-Maximisation problem where the E-step is assigning the data points to the closest cluster and the M-step is computing the centroid of each cluster. It minimises the intra-cluster sum of squared distance from its centroid and keeps the clusters distant from each other.

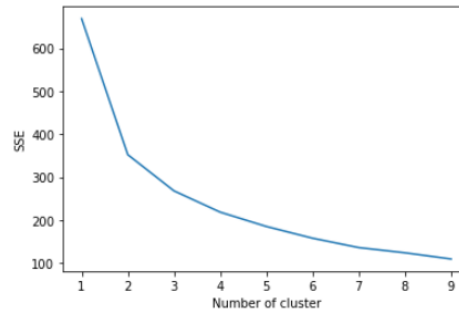
Evaluation metric

SSE - Elbow method

Since it is an unsupervised learning model, we don't have ground truth values to evaluate its performance. Hence, we select the model with a K value which minimises sum of squared distances between data points and its centroids.

When SSE is plotted against K , select the K where the graph flatten out and forms an elbow. This method is referred to as elbow method.

Figure 1: Elbow method plot example



Code Example

Generated data using gaussian distribution and took user defined values for K and input dimensions.

Data generation

```
In [1]: 1 NUM_PTS=3
2 #Accepting input dimensions(p) and number of clusters(k)
3 Y=np.empty((p,k*NUM_PTS),int)
4 p=int(input("Enter input dimensions: "))
5 k=int(input("Enter number of clusters:"))
6 #contains the data points of all the clusters
7 #declaring positive semi definite covariance matrix
8 cov=[]
9 y=[]
10
11
12 #Generating clusters using Gaussian distribution
13 for i in range(k):
14     #declaring mean matrix
15     mean=np.random.randint(i+1,p+i+1,p,int)
16     cov=np.diag([i for j in range(p)])
17     y=np.random.multivariate_normal(mean,cov,NUM_PTS).T
18     Y=np.concatenate((Y,y),axis=1)
```

Declaring stopping condition and implementing the algorithm

```
In [2]: 1 #Setting stopping condition
2 #epsilon = 0.001
3 # Initialize error to a large value
4 #error = 10000.0
5 # Initialize centroids - assume 4 of them from a 2D Gaussian
6     distribution (zero mean, unit variance)
7 D = np.random.multivariate_normal(np.zeros(p), np.identity(p), k).T
```

```

8 # Initialize iteration count to 0
9 count = 0
10 # Initialize cluster size to 0
11 num=np.zeros(k)
12 # Initialize centroid update to 0
13 centr= np.zeros((p, k))
14
15 #Initialise distance_vector
16 distance_vector=[]
17
18
19 while (error > epsilon):
20
21     # Update clusters based on distance
22     for idx in range(k*NUM_PTS):
23         cur_y = Y[:,idx]
24         print(cur_y)
25         for j in range(k):
26             cur_d=D[:,j]
27             d1=np.sqrt(np.sum((cur_y-cur_d)**2))
28             distance_vector.append(d1)
29         # Find closest centroid
30         min_idx = distance_vector.index(np.min(distance_vector))
31         for i in range(k):
32             if (min_idx==i):
33                 num[i]+=1
34                 for j in range(p):
35                     centr[j,i]+=cur_y[j]
36                 if (num[i]!=0):
37                     centr[j,i]=centr[j,i]/num[i]
38
39     error=np.sqrt(np.sum((centr-D)**2))
40     for i in range(k):
41         for j in range(p):
42             D[j,i]=centr[j,i]
43     count+=1

```

Key points

- The algorithm clusters the data into distinct sub groups which will not work for overlapping clusters.
- Normalization is required as it deals with distances.
- It assumes spherical shapes of clusters with center as centroid and fails in case of complex designs or even elliptical shape.
- The algorithm works well with large datasets and easy to implement.
- The centroids are sensitive to outliers.
- The algorithm is stochastic as the clusters formed depend on centroid initialisation.

Questions

1. Does the cost function of the algorithm has an optimized solution? **Solution:**
Objective function is given by

$$\arg \min \sum_{i=1}^K \sum_{x \in S_i} ||x - C_i||^2$$

This function is a NP hard problem. Therefore, we only use Lloyd algorithm to stop iterating if the problem converges.

Since most of the convergence takes place in first few iterations, stopping condition is changed to 'Until relatively few points change clusters'.

2. Illustrate an example where K means fails when clusters have complex design.
Solution:

Figure 2: K means for non spherical clusters

