

## Data Mining & Machine Learning

### **Group Members:**

1. Avik Das - MDS202112
2. Meghna Mondal – MDS202123

### **K-means clustering Models on Bag of Words dataset:**

In this assignment we have built K-means Clustering for 3 data sets:

- 1) KOS,
- 2) NIPS
- 3) ENRON.

Then we used suitable evaluation metric to compare the performance of the three classifiers.

### **Table for Comparing the classifier models:**

<b>Performance Measure</b>	<b>KOS</b>	<b>NIPS</b>	<b>ENRON (On 1% Sample)</b>
Time Taken(s)	<b>128.02</b>	<b>51.86</b>	<b>3962.09</b>
Space required(MiB)	<b>726.61</b>	<b>582.61</b>	<b>5115.92</b>

### **Procedure of fitting the Models:**

- At first, we read the dataset and skip the 1<sup>st</sup> 3 lines of each data. And then we take the matrix containing Doc\_ID, Word\_ID, Word\_Count and made the data frames.
- Then we create the sparse matrix with Doc\_IDs in the column and Word\_IDs in the rows. And it contains 1 at  $ij^{th}$  position if a  $i^{th}$  Word\_ID has appeared in a  $j^{th}$  Doc\_ID.
- Then we calculate the Jaccard Index for each of the document and made a symmetric matrix containing Jaccard index between  $i^{th}$  and  $j^{th}$  document at  $ij^{th}$  position.
- After that we pass that symmetric matrix to build Kmeans clustering model with different number of clusters and then by plotting the

graph of inertia vs # clusters, we find the optimum number of clusters by elbow method.

- Then we finally fit the Kmeans clustering model with the optimum number of clusters for KOS and NIPS dataset.
  - In KOS dataset the optimum number of clusters is 2.
  - In NIPS dataset the optimum number of clusters is 4.
- After that we perform dimension reduction for visualizing the clusters in 3-dimensional space.
- And at last, we observed the Doc\_IDs belonging to different clusters and give them as the final output in a dictionary format.
- In case of ENRON dataset, we at first did stratified sampling and take 1% of the data as sample grouping by Word\_counts(i.e., keeping the ratio of word frequencies same as the original data).
- Then we go through the same procedure, but here we didn't get a good elbow, so we can't able to fit the model in a good way.

### Visualizing the Clusters:





