# MLBD Assignment-1

By:- Kwanit Gupta (B19EE046)
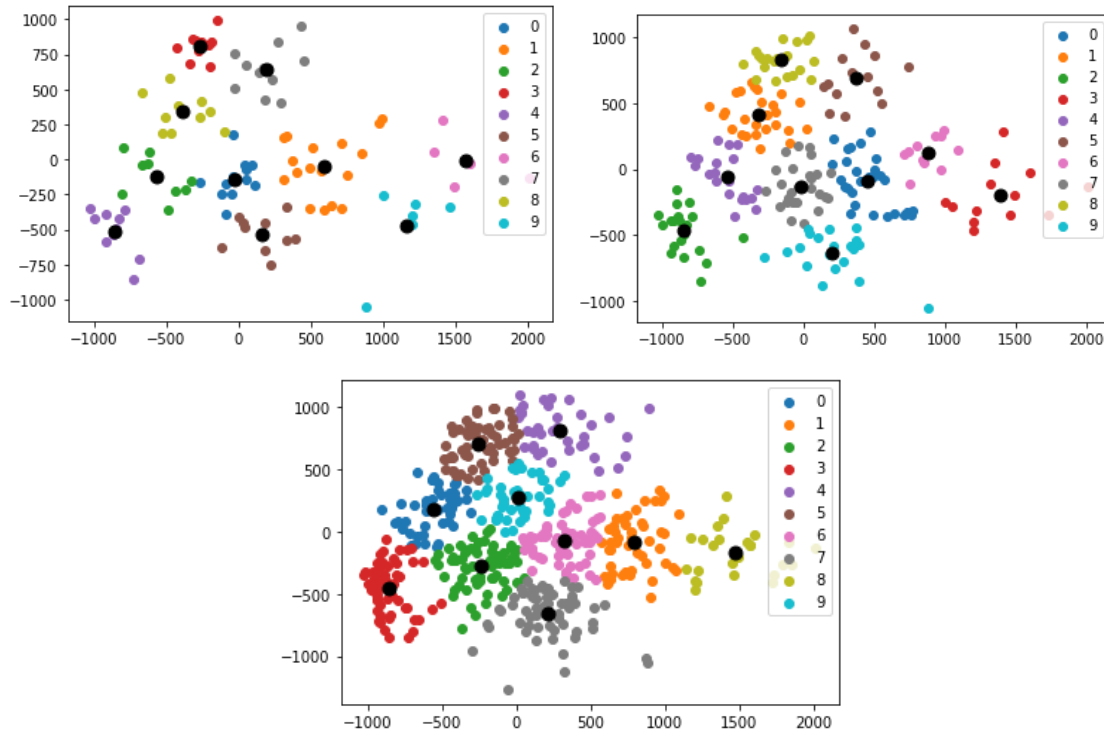
## Description:-

Download the MNIST dataset from here. Pick its training test set which contains 60K samples, each represented by a 784-dimensional vector. There are ten classes (the ten digits, `0' to `9'), and each sample is associated with one class. Preprocess the features using the appropriate technique(s), and store them in a file.

(1) Implement the BFR and CURE clustering algorithms on this data assuming that you can store `K1' samples in the main memory at a time. One may use any existing libraries for performing the initial k-means clustering in case of BFR, and agglomerative clustering in case of CURE. After this, every step needs to be implemented from scratch.

(2) For the BFR algorithm, ensure that the number of clusters obtained at the end of the training the process is 10.

(3) For the CURE algorithm, keep the number of clusters as 10.

(4) After clustering, calculate the percentage of samples from each class and convert it into probability values. Using these, calculate the entropy of each cluster. Also calculate the total entropy of all the clusters by summing the entropy of individual clusters.

(5) Re-run the two algorithms five times assuming K1 = {100,200,500}, and report the above result.

## Solution:-

For the pre-processing part, I utilized PCA (Principal Component Analysis), since variance has to be retained and escape from multiple 0s and 1s, since the visualization in the code wasn't evident, eventually. After bringing down the n_components to 2, for useful visualizations,

The Above 2 plots were the clustering results for K1=100 and K1=200. Whereas for the middle one, K1=500. Hereby, I observed the extensive intermixing of the PCA Components when K1 was increased.

Talking about the algorithms BFR with K-Means and CURE with Agglomerative Clustering, the dedicated code bases will be attached within the zip. Following are some of the screenshots for the working (Final conclusions didn't came since bugs weren't resolved, but some of the verbose can extensively show the promise of the implementations):-

1. BFR with K-Means Clustering:-

```
        rs_idx, cs_idx = process_cluster(group_cluster(km.labels_, rs_points))
        if len(cs_idx) > 0:
            if len(cs_clusters) > 0:
                new_cs_clusters = summarise_info_clusters(km.labels_[tuple(cs_idx)], rs_points[tuple(cs_idx)])
            else:
                cs_clusters = summarise_info_clusters(km.labels_[tuple(cs_idx)], rs_points[tuple(cs_idx)])
            if len(new_cs_clusters):
                cs_clusters = merge_clusters(cs_clusters, new_cs_clusters, threshold, return_two=False)
        if len(rs_idx) > 0:
            rs_points = rs_points[rs_idx]
        if chunk < (K - 1):
            stats.append(statistics(ds_clusters, cs_clusters, rs_points))
```

```
{'DS Clusters': 10, 'DS Points': 990, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 10, 'SUM': 1000}
{'DS Clusters': 10, 'DS Points': 1484, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 16, 'SUM': 1500}
{'DS Clusters': 10, 'DS Points': 1975, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 25, 'SUM': 2000}
{'DS Clusters': 10, 'DS Points': 2471, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 29, 'SUM': 2500}
{'DS Clusters': 10, 'DS Points': 2956, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 44, 'SUM': 3000}
{'DS Clusters': 10, 'DS Points': 3447, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 53, 'SUM': 3500}
{'DS Clusters': 10, 'DS Points': 3941, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 59, 'SUM': 4000}
{'DS Clusters': 10, 'DS Points': 4434, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 66, 'SUM': 4500}
{'DS Clusters': 10, 'DS Points': 4930, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 70, 'SUM': 5000}
{'DS Clusters': 10, 'DS Points': 5425, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 75, 'SUM': 5500}
{'DS Clusters': 10, 'DS Points': 5920, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 80, 'SUM': 6000}
{'DS Clusters': 10, 'DS Points': 6415, 'CS Clusters': 1, 'CS Points': 2, 'RS Points': 83, 'SUM': 6500}
-----------------------------------------------------------------------------
IndexError                          Traceback (most recent call last)
<ipython-input-118-4451ed035fb6> in <module>
     16     if len(cs_idx) > 0:
     17         if len(cs_clusters) > 0:
---> 18             new_cs_clusters = summarise_info_clusters(km.labels_[tuple(cs_idx)], rs_points[tuple(cs_idx)])
     19         else:
     20             cs_clusters = summarise_info_clusters(km.labels_[tuple(cs_idx)], rs_points[tuple(cs_idx)])

IndexError: too many indices for array: array is 1-dimensional, but 2 were indexed
```

SEARCH STACK OVERFLOW

```
ds_clusters, cs_clusters = merge_clusters(ds_clusters, cs_clusters, threshold, return_two=True)
stats.append(statistics(ds_clusters, cs_clusters, rs_points))
predictions = predict(K, ds_clusters, tr_splits, threshold)
```

```
{'DS Clusters': 10, 'DS Points': 6914, 'CS Clusters': 0, 'CS Points': 0, 'RS Points': 86, 'SUM': 7000}
```

The above 2 was for K1=100, K1=200, while the middle one is for K1=500.

2.  CURE with Agglomerative Clustering:-

The above 2 was for K1=100, K1=200, while the middle one is for K1=500. The middle one was on run during the submission.

Due to recurring bugs, I wasn't able to solve them timely, but following the logic of the course slides, I was able to extensively understand the core components of the program chunks. That's why I won't be able to present the observations and maybe the last part too, since I also tried to run these algorithms 5 times, the initiations of the clusters were different each time, but bugs being the same.