



University  
of Windsor  
Faculty of Engineering

## **ELEC8900-76-R-2021F: Applied Machine Learning**

**Fall 2021**

### **Project DR03**

Ishita Joshi (110028901) and Sanjay Kumar Gengaiyan (110038152)

E-mails: {joshi14; gengaiy} @uwindsor.ca

**Submitted to:**

Professor: Roozbeh Razavi Far

Date: 29-11-2021

## I. INTRODUCTION

This report summarizes the results of implementing dimensionality reduction techniques: Linear Discriminant Analysis (LDA), Sammon Mapping (SM) and Isomap Embedding (IM) on 13 given datasets with 17.5K, or 35K, or 3541 samples, 354 or 220 features, and 70 or 10 classes. The reduced datasets obtained using these algorithms are then passed on to four classifiers: Naive Bayes Classifier (NBC), Support Vector Machines (SVM), Decision Trees (DT), and Multi-layer Perceptron<sup>\*</sup> (MLP).

## II. DIMENSIONALITY REDUCTION ALGORITHMS

Dimensionality reduction techniques are applied during data preprocessing and the goal is to project the dataset onto a lower dimension for better visualization and easier classification. The Dimensionality reductions algorithms used in the project are briefly explained in the following sections.

### A. Linear Discriminant Analysis (LDA)

LDA is a supervised classification algorithm that is commonly used as a dimensionality reduction technique. It projects a d-dimensional feature space to a k-dimensional feature subspace. LDA is a linear transformation technique just like Principal Component Analysis (PCA), but it also maximizes separability between classes. Fisher [1] first introduced LDA for binary classification in 1936 and Rao [2] extended it for multiclass classification in 1948.

Conceptually, LDA creates a new axis and projects the data onto it in a way that the separation between means of each class is maximized and the scatter within each class is minimized. LDA extracts eigenvectors and eigenvalues from the scatter matrices or from covariance matrices. It then arranges the eigenvectors in decreasing order and chooses the top-k among them. The concept behind this is that features with eigenvalues nearer to zero are less informative and hence, can be discarded. Algorithm 1 represents the pseudo-code for LDA [3].

---

#### Algorithm 1 LDA

---

**Input:** Dataset  $X_i$ , lower dimension  $d$

**Output:**  $k$  dimensional feature subspace  $Y_i$  such that  $(k < d)$

1. Compute  $d$ -dimensional mean vectors  $\mathbf{m}_i$

$$\mathbf{m}_i = \frac{1}{N_i} \sum_{x \in D_i}^N x_k$$

2. Compute the within-class and between-class scatter matrices  $\mathbf{S}_W$  and  $\mathbf{S}_B$ :

$$\mathbf{S}_W = \sum_{i=1}^c \mathbf{S}_i$$

Where  $\mathbf{S}_i$  is the scatter matrix for every class

$$\mathbf{S}_i = \sum_{x \in D_i}^N (x - \mathbf{m}_i)(x - \mathbf{m}_i)^T$$

And between-class scatter matrix

$$\mathbf{S}_B = \sum_{i=1}^c N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$$

Where  $\mathbf{m}$  is the overall mean

3. Solve the matrix  $\mathbf{S}_W^{-1} \mathbf{S}_B$  to generate eigenvectors and eigenvalues.
  4. Sort the eigenvectors in decreasing order of eigenvalues, select the top-k among them and create a  $d \times k$ -dimensional transformation matrix  $\mathbf{W}$
  5. Transform the data to new  $k$ -dimensional subspace  $Y = X \times \mathbf{W}$
- 

### B. Sammon Mapping

Sammon's mapping algorithm is a nonlinear mapping algorithm proposed by John W. Sammon in 1969 [4]. It maps the given dataset to a lower dimension in a way that the structure is preserved, and it does

<sup>\*</sup>Initially assigned Recurrent Neural Network but changed to MLP.

so by conserving the Euclidean distance between the data points.

If we have L-dimensional data containing N data points  $X_i$  to be mapped on a d-dimensional subspace ( $d < L$ )  $Y_i$ , then  $d_{ij}^*$  and  $d_{ij}$  are the distances between vectors  $X_i$  and  $X_j$  in the L-subspace and d-subspace respectively. Distance between all these data points of  $X_i$  and  $Y_i$  is calculated, and an error function E is defined that gives an idea about structure lost in the transformation.

$$E = \frac{1}{\sum_{i < j} [d_{ij}^*]} \sum_{i < j}^N \frac{[d_{ij}^* - d_{ij}]^2}{d_{ij}^*}$$

The goal here is to minimize this error using a steepest descent technique. Algorithm 2 describes the pseudo-code for Sammon Mapping.

---

#### Algorithm 2 SM

---

**Input:** Given data  $X_i$ , k-nearest neighbors, dimension of output d,

**Output:**  $Y_i$  in the lower dimension d.

1. Compute all pair-wise distances in the d-dimensional space  $d_y(i,j)$ .
  2. Initialize  $Y_i$  either randomly or using a linear transformation like PCA.
  3. Using the steepest descent technique, minimize error E.
  4. Update the co-ordinates of all data points of  $Y_i$
  5. **Repeat** until E is optimized
  6. **Return**  $Y_i$
- 

#### C. Isomap

Isomap is the isomap mapping algorithm proposed by Tenenbaum *et al.* in 2000 <sup>[5]</sup>. It is built upon the multidimensional scaling (MDS) which computes the low-dimensional representation based on pairwise geodesic distances. Geodesic is defined as the shortest path on the given surface. MDS computes low-dimensional representation which is usually measured using straight-line Euclidean distance whereas isomap mapping computes shortest path between points in the constructed neighbourhood graph.

The goal is to estimate the distance between the faraway point with the given input-space distances. The input-space distances provide a good approximation for nearby points and distance between the faraway points can be approximated by short-hops between the neighbouring points. After estimating the geodesic distances between all the pairs of points on the manifold classic MDS is applied to the matrix to construct a d-dimensional embedding that best preserves the intrinsic geometry. The coordinate vectors are chosen to minimize the cost function  $E = ||\tau(D_G) - \tau(D_Y)||_{L^2}$  where  $D_Y$  is the matrix of Euclidean distances  $D_Y(i,j) = ||y_i - y_j||$  and  $||A||_{L^2}$  is the matrix form of  $\sqrt{\sum_{i,j} A_{ij}^2}$ . The  $\tau$  operator converts distances to inner products which characterize the geometry of data that supports efficient optimization. The global minimum is achieved by setting the coordinates  $y_i$  to the top  $d$  eigenvectors of the matrix  $\tau(D_G)$ . <sup>[5]</sup>

---

#### Algorithm 3 IM

---

**Input:** Dataset X with N data points, k nearest neighbours

**Output:**  $Y = \{y_1, y_2, \dots, y_n\}$

1. Construct a neighborhood graph. This step determines the neighbor points on the manifold M based on the distance  $d_x(i,j)$  between the pairs of inputs  $i,j$  in the input space X. The set of neighbors are determined by k nearest neighbors or e-radius.
2. In the second step, the algorithm estimates the geodesic distances  $d_M(i,j)$  between all pairs on the manifold M by computing the shortest path lengths  $d_G(i,j)$  in graph G <sup>[6]</sup>.

$$\begin{aligned} d_G(i,j) &= d_x(i,j) \text{ neighboring } i,j \\ d_G(i,j) &= \infty \text{ otherwise} \end{aligned}$$

For each value of  $k=1,2,\dots,N$ , replace  $d_G(i,j)$  by  $\min \{d_G(i,j), d_G(i,k) + d_G(k,j)\}$ . This gives the final values of  $D_G = d_G(i,j)$

3. Apply classical MDS to construct a lower dimensional embedding. In the given matrix  $D$ , MDS attempts to find  $n$  data points  $y_1, y_2, \dots, y_n$  in  $d$  dimension such that, if  $\hat{d}_{ij}$  is the Euclidean distance between  $y_i$  and  $y_j$ , then  $\hat{D}$  is similar to  $D$  [7].

MDS minimizes the following function:

$$\min_y \sum_{i=1}^n \sum_{j=1}^n (d_{ij}^{(X)} - d_{ij}^{(Y)})^2 \text{ where}$$

$$d_{ij}^{(X)} = ||x_i - x_j||^2 \text{ and } d_{ij}^{(Y)} = ||y_i - y_j||^2$$

4. Obtain Gram Matrix by double centering [7]  
 $X^T X = -\frac{1}{2} H D^X H$  where  $H = I - \frac{1}{n} e e^T$  and  
 $e$  is a column vector of all 1s.

The equation is reduced to

$$\min_y \sum_{i=1}^n \sum_{j=1}^n (x_i^T x_j - y_i^T y_j)^2$$

5. Compute eigenvalue decomposition of Gram matrix  $S = U \Lambda U^T$ .

The output of the MDS is  $Y = \Lambda^{1/2} V^T$  where  $V$  is the eigen vectors of  $X^T X$  corresponding to  $d$  eigenvalues and  $\Lambda$  is the top  $d$  eigenvalues of  $X^T X$ . [7]

6. A  $p$ -dimensional representation of  $Y$  is obtained by  $\hat{X} = I_{p \times N} \Lambda^{1/2} V^T$  [7]
- 

### III. IMPLEMENTATION

This section will cover the languages and tools used, the implementation procedure and some challenges faced during implementation.

#### A. Language and tools used

- Python 3.8.8
- Scikit-learn 0.24.1
- Pandas 1.2.4
- NumPy 1.20.1
- Sammon mapping function by Tom Pollard [8]
- Editor used: Jupyter Notebook

#### B. Implementation Procedure

1. Import libraries and dataset.
2. Split dataset into input  $X$  and output  $y$ .
3. Perform baseline classification.
4. Apply DR technique (LDA, SM, IM) to extract 2 features ( $d=2$ ) and compute variance of each feature.
5. Apply NBC classifier, check accuracy and F1-score.
6. Manipulate the value of  $d$  and find a value that gives highest accuracy and F1-score.
7. Compute variance of  $d$  extracted features.
8. Apply remaining classifiers SVM, DT and MLP to the reduced dataset.
9. Repeat for all datasets.

#### C. Challenges Faced

- Recurrent Neural Network models (even the lightest one) was too bulky and took hours to run. It was later replaced with MLP which was slow as well but better than RNN.
- Isomap and Sammon mapping algorithms are computationally heavy and had to be implemented on a sampled dataset.
- Isomap and Sammon mapping being distance-based algorithms, their explained variance does not

say much about the ‘goodness of fit’. Sklearn does not provide a function to calculate the explained variance of each feature.

- MLP needs a lot of fine-tuning for layer size, number of hidden layers, learning rate, optimizer, etc.

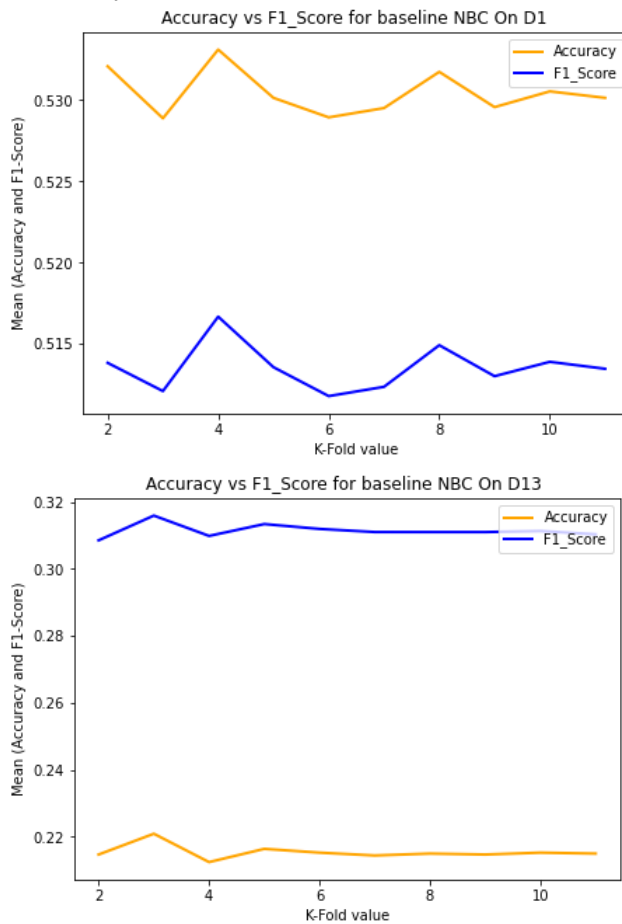
## IV. EXPERIMENTAL RESULTS

In this section, the graphical results generated during our experiments to find the best combination of hyperparameters will be discussed. Firstly, baseline classification was performed on our classifiers NBC, SVM, DT and MLP and generated graphs for one of the 17k, 35k datasets each and D13. Once the baseline classification results were ready, LDA was implemented followed by IM and SM. For dimensionality reduction, the most important parameter is the number of extracted features. We tried different values of  $d$  for each algorithm on D1 and implemented it on the remaining datasets.

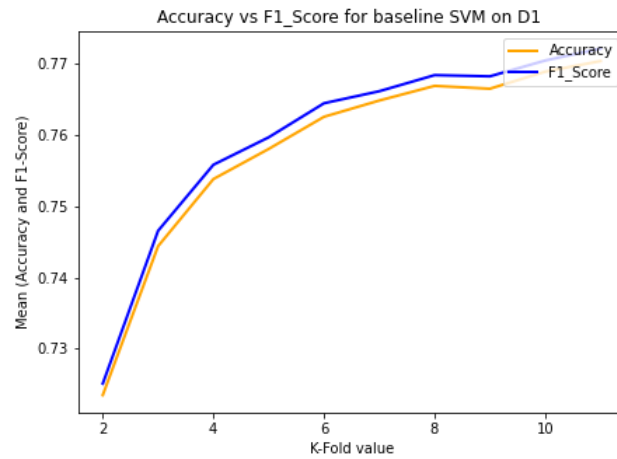
A general pattern observed throughout our implementation process is that datasets with same number of samples have closely similar accuracies and require the same number of extracted features. The classifier results after applying dimensionality reduction algorithms give a reduced accuracy than the baseline except for dataset 13.

### A. Baseline classification

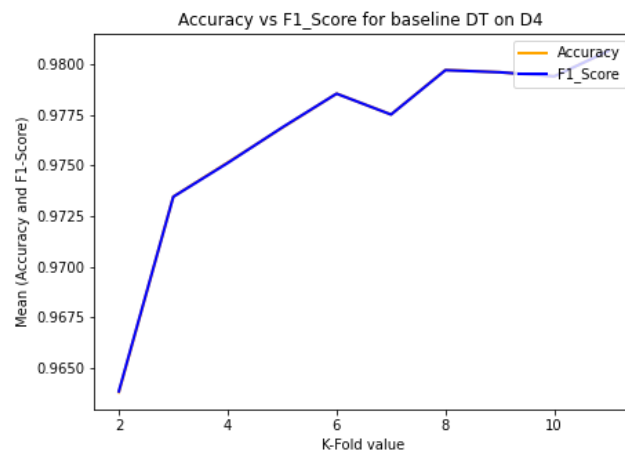
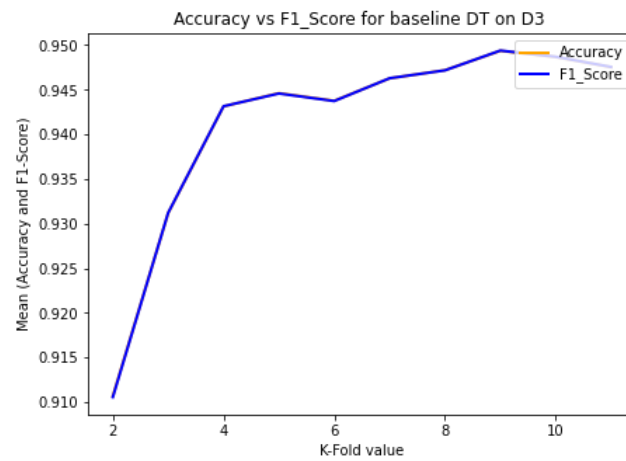
#### i. Naive Bayes Classifier

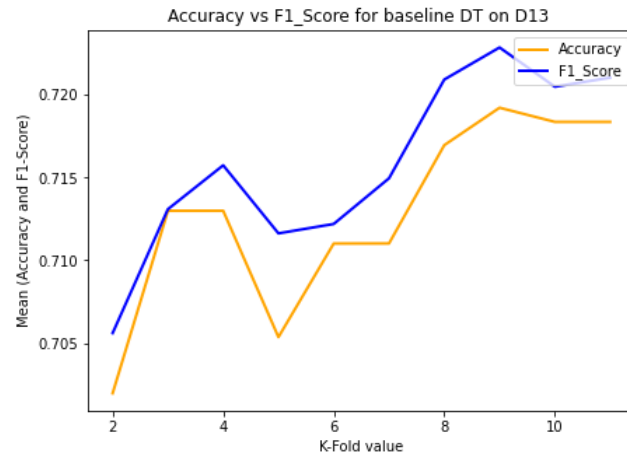


## ii. Support Vector Machine



## iii. Decision Tree Classifier

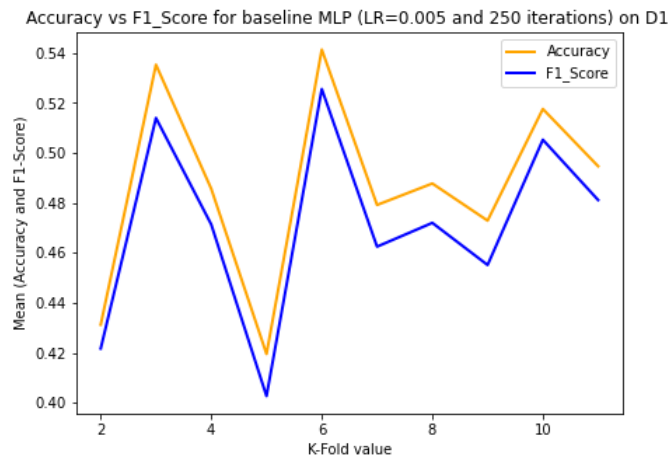


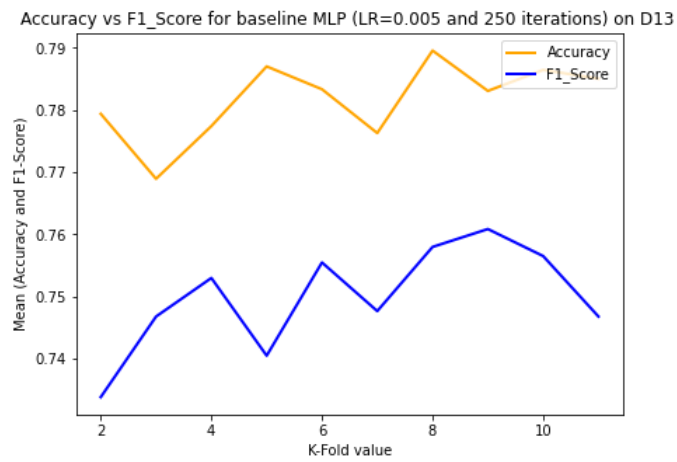
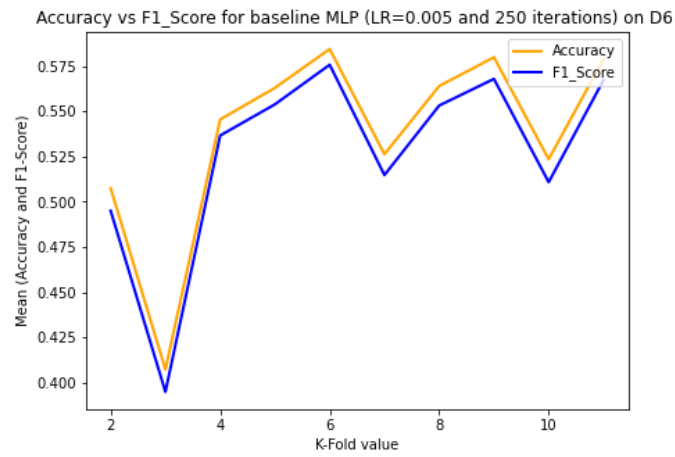


#### iv. Multi-layer Perceptron

For MLP, different combinations of hyperparameters were experimented on the data. The parameters of the final model are as mentioned:

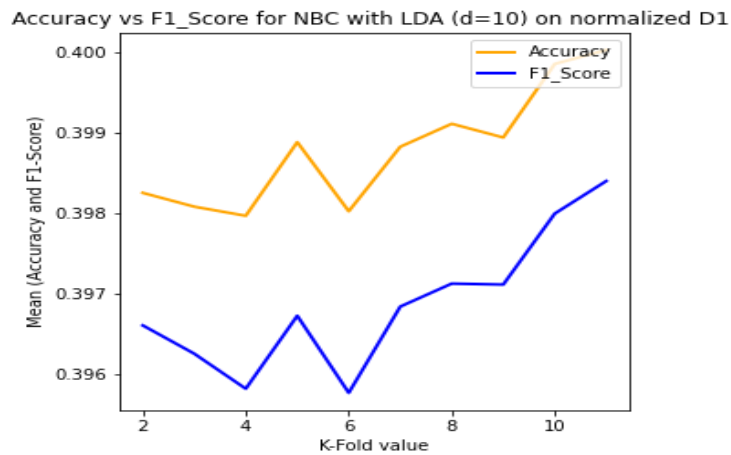
- 3 hidden layers of sizes (256,128,128)
- Learning rate = 0.005  
Values 0.001 and 0.01 were tried during experimentation and these values were either too slow or too fast hence failed to converge within the maximum number of iterations.
- Optimizer: Adam
- Maximum iterations: 250



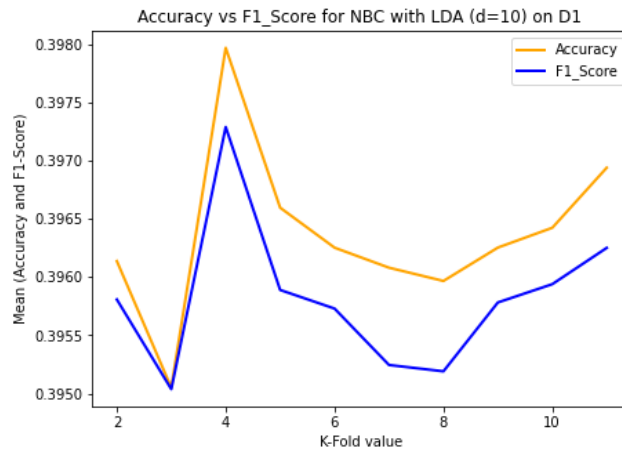


## B. Linear Discriminant Analysis (LDA)

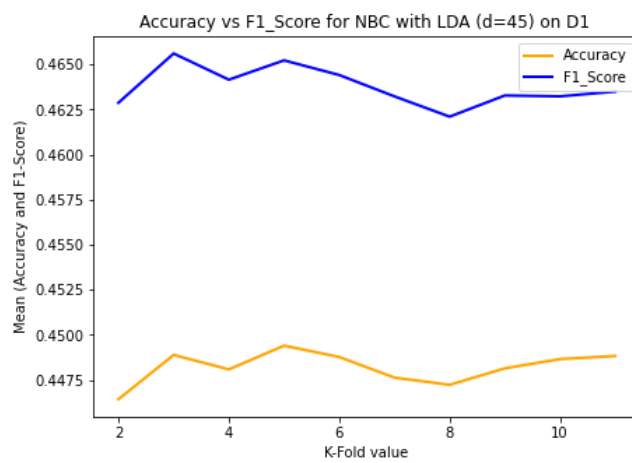
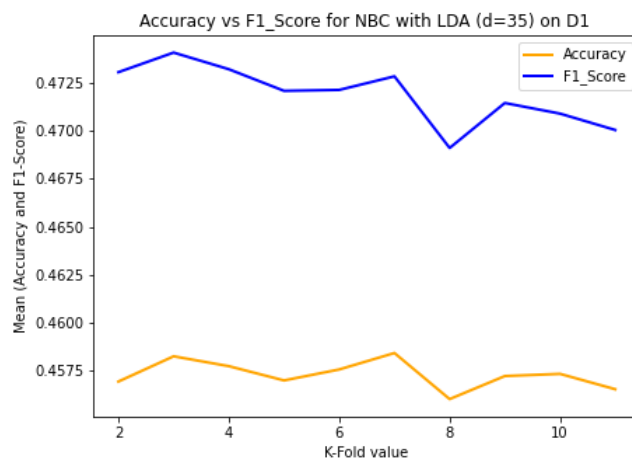
- i. LDA with and without data normalization

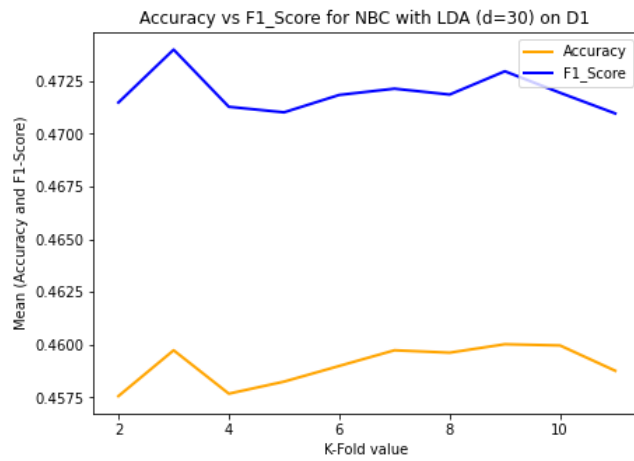
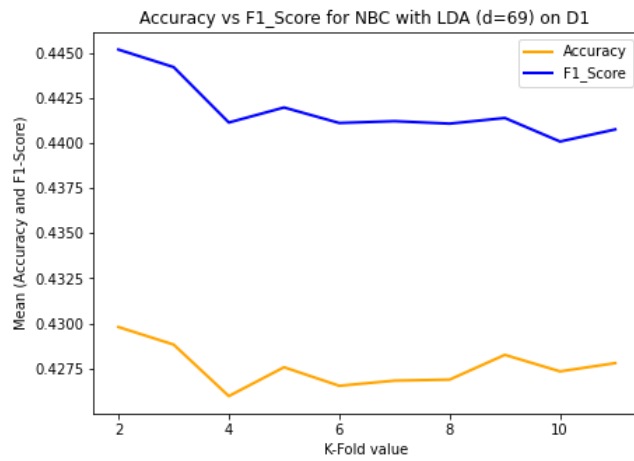
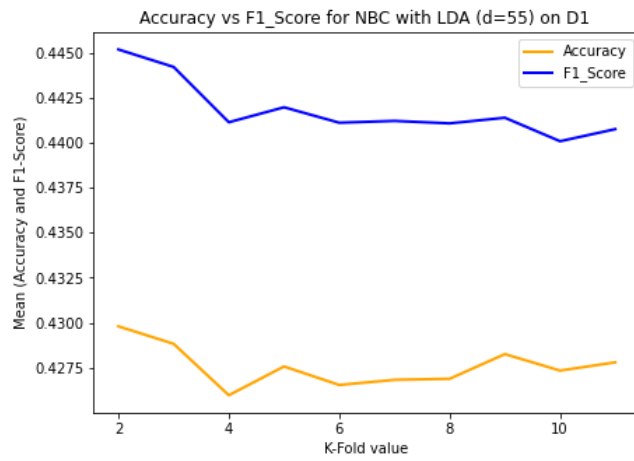


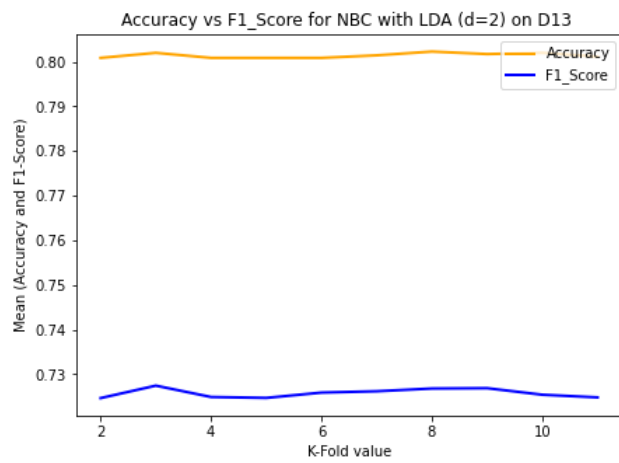
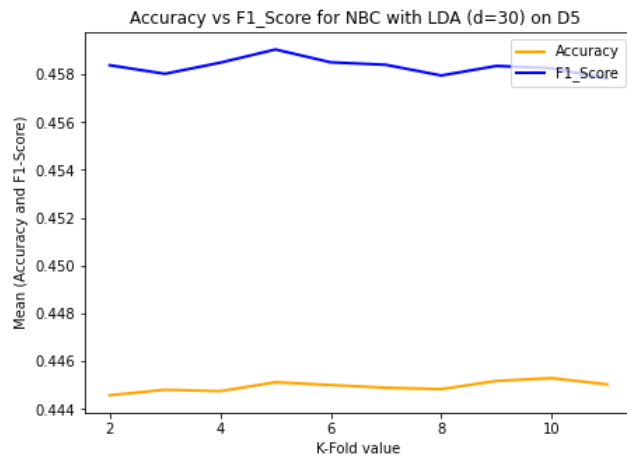




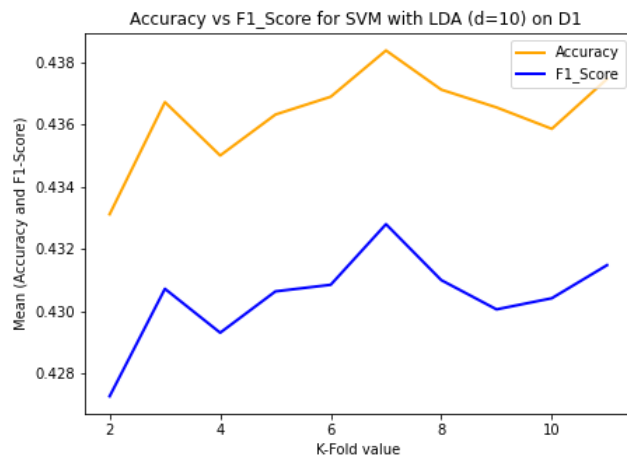
ii. LDA for different values of  $d$  with NBC on D1

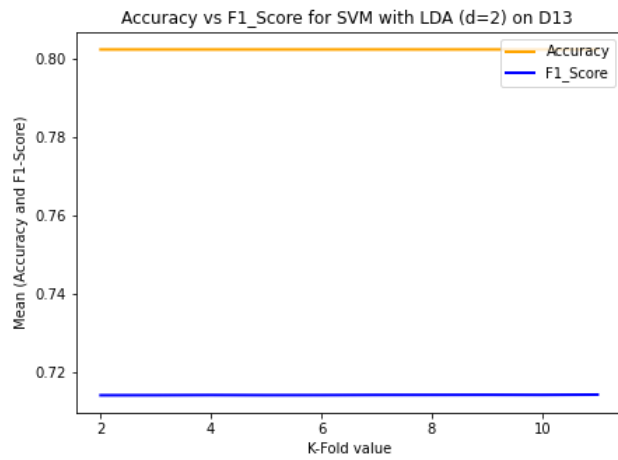
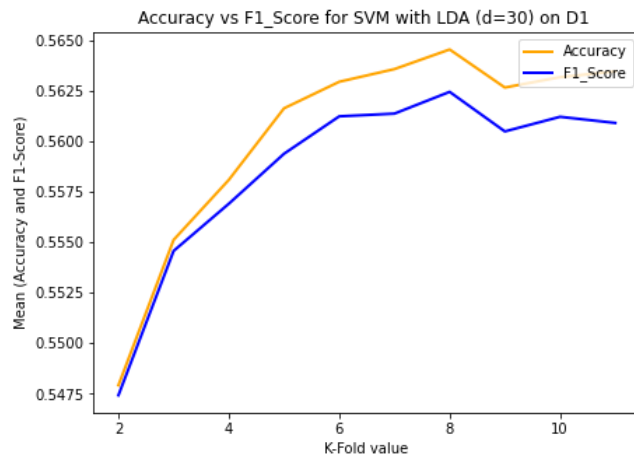




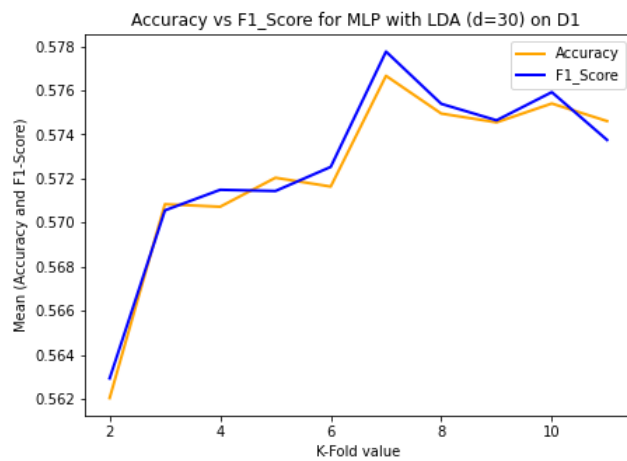


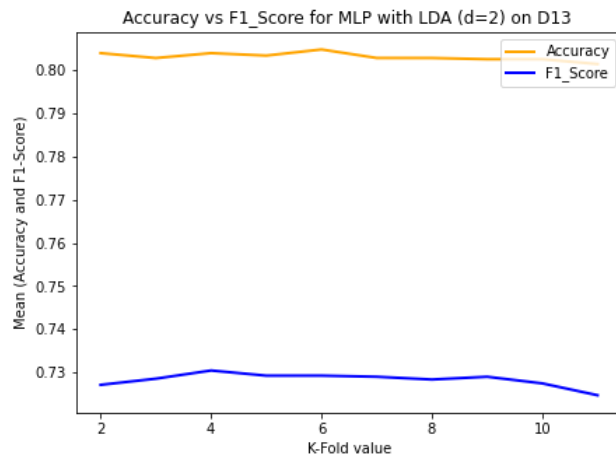
iii. LDA with SVM classifier



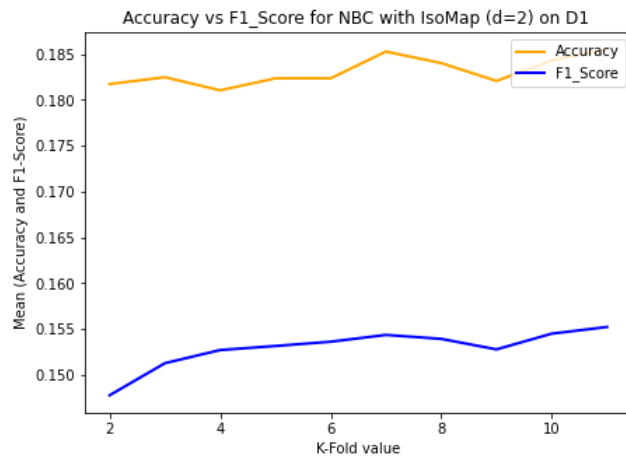
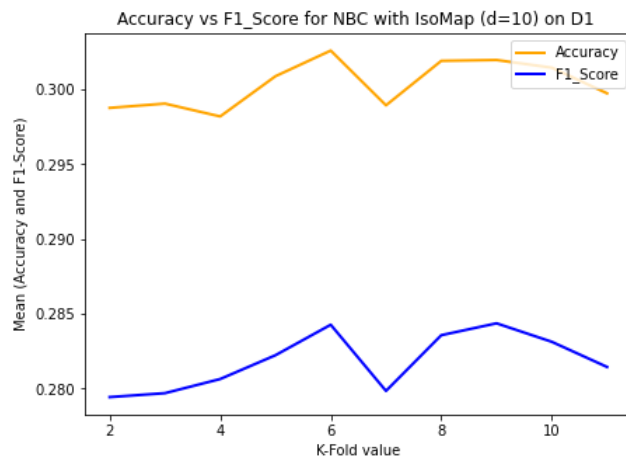


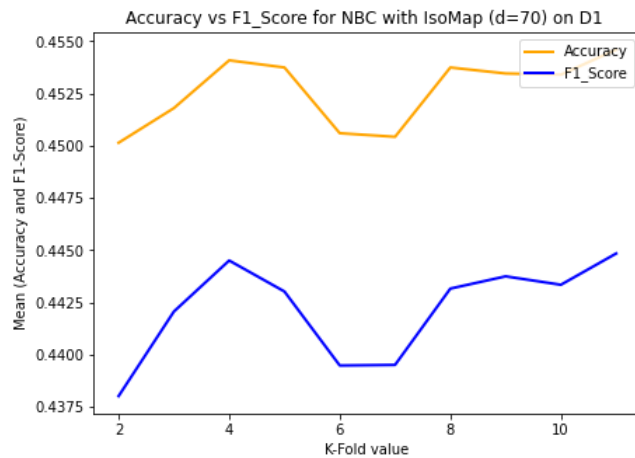
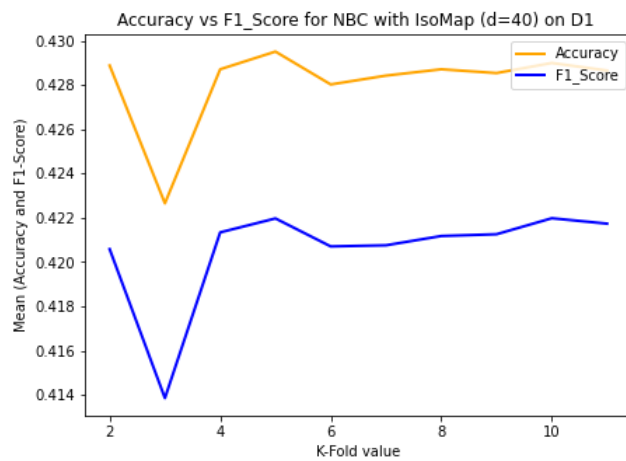
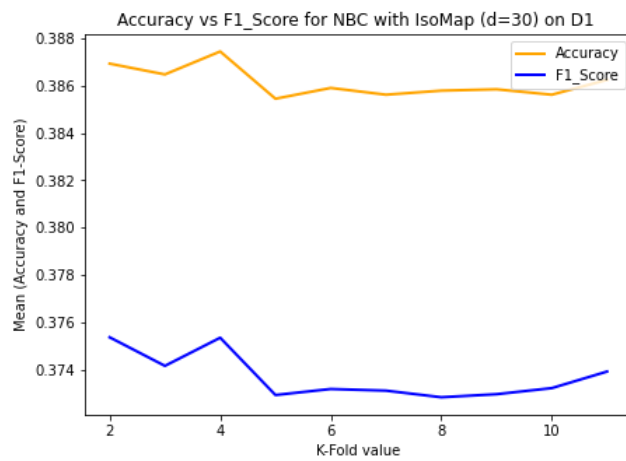
iv. LDA with MLP classifier

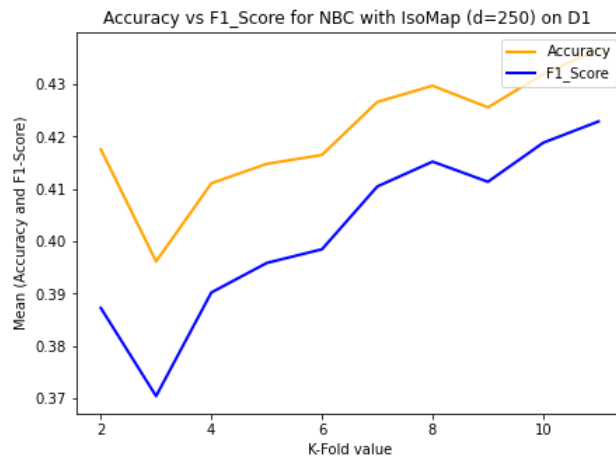
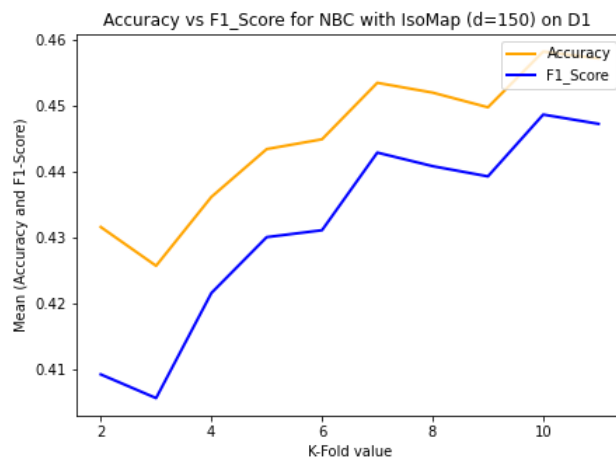
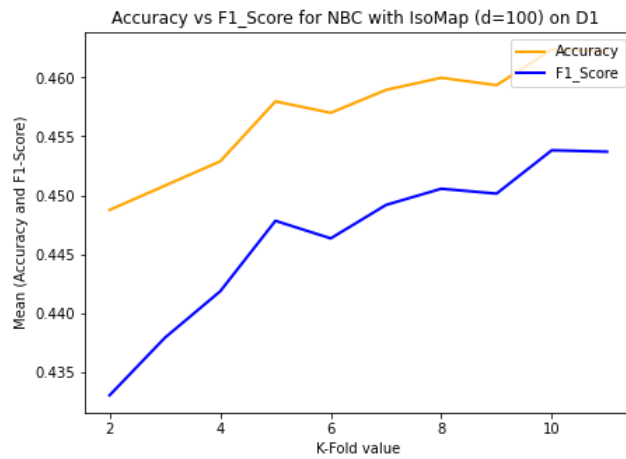


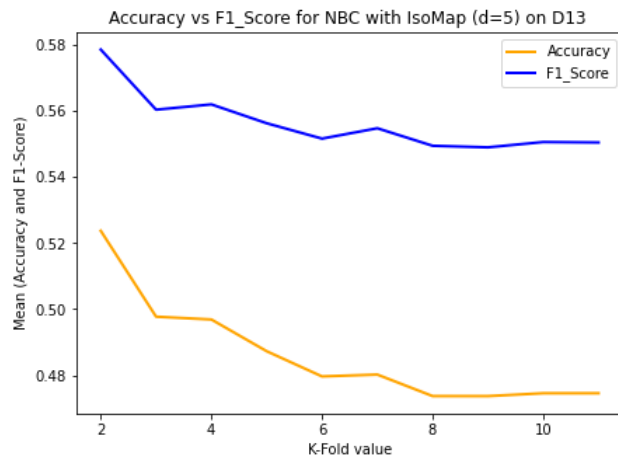
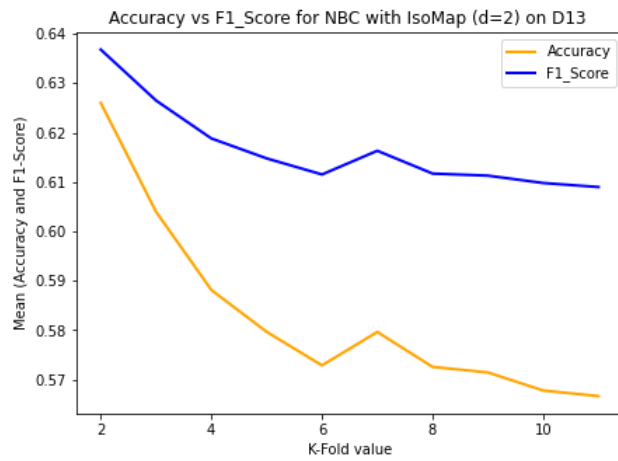
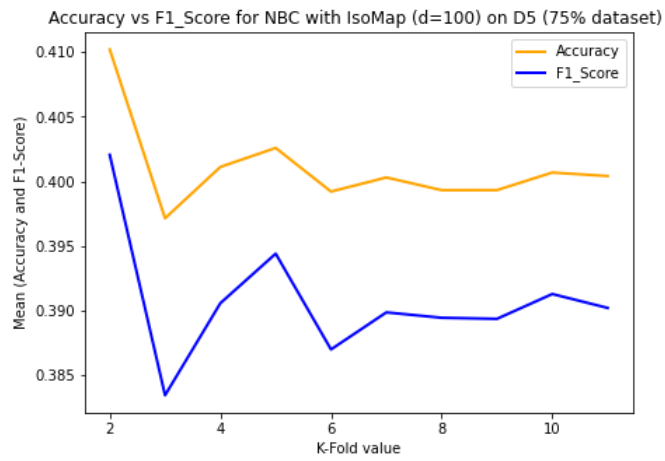


### C. Isomap Embedding

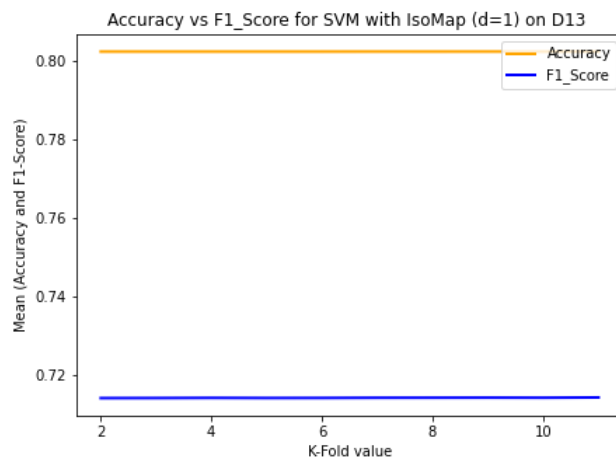
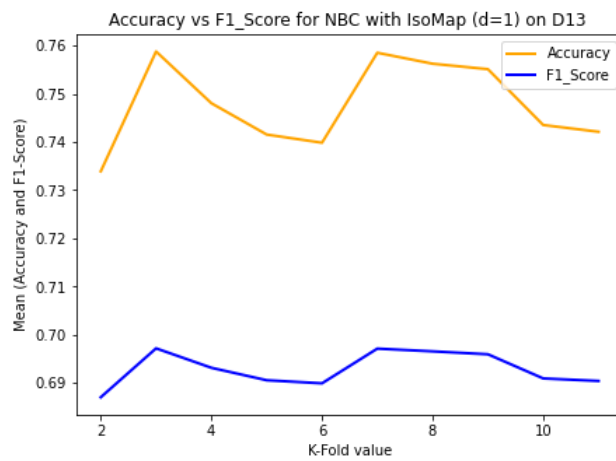




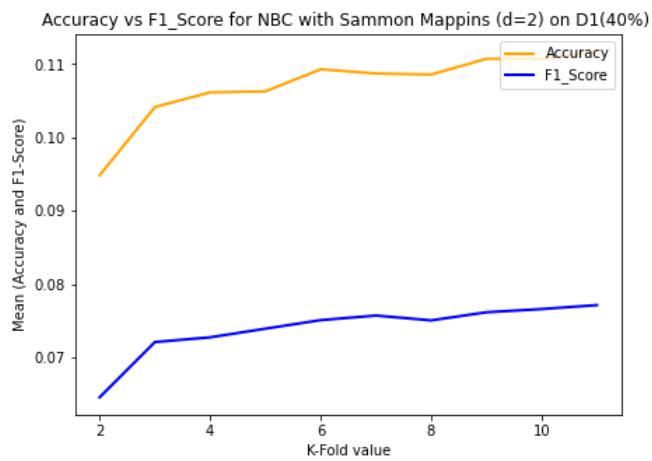


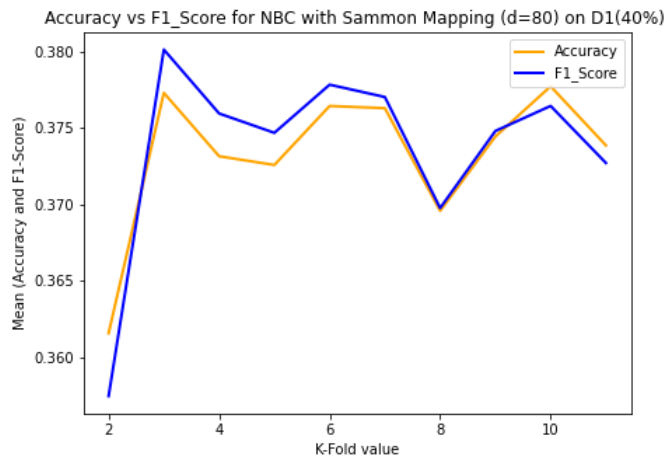
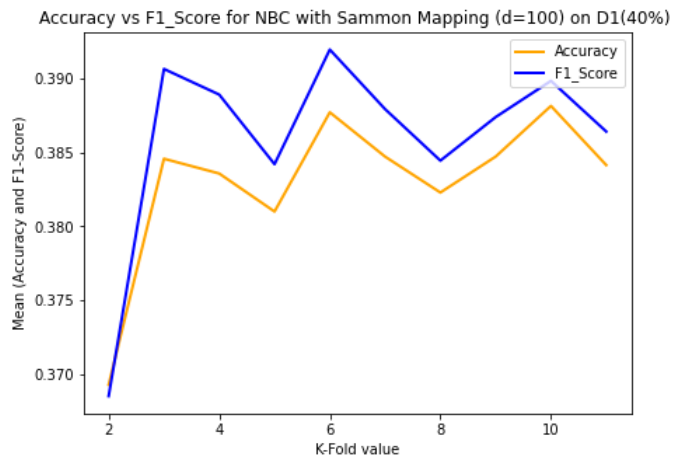
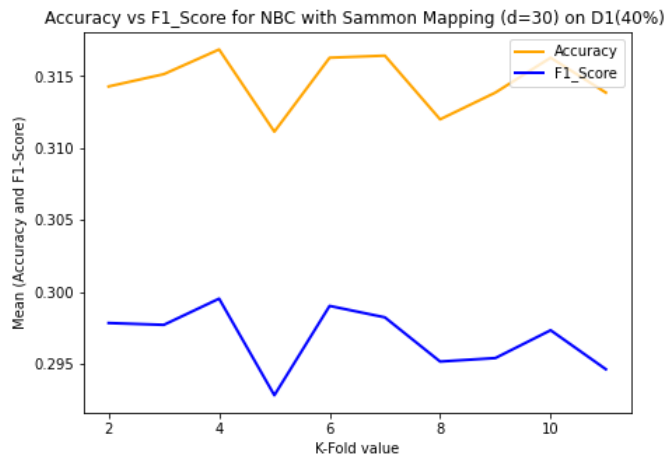


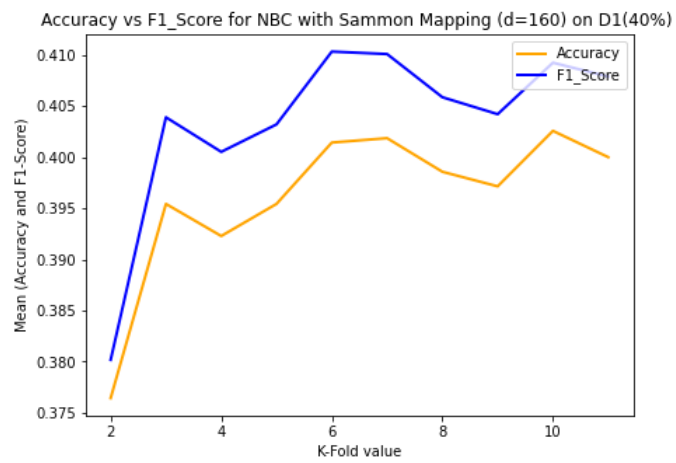
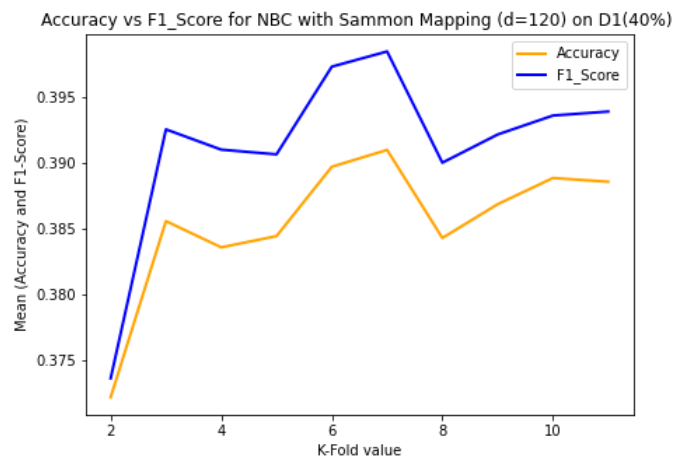




#### D. Sammon Mapping







## REFERENCES

- [1] R.A. Fisher, "The Statistical Utilization of Multiple Measurements," *Ann. Eugenics*, vol. 8 pp. 376-386, 1938
- [2] C.R. Rao, "The Utilization of Multiple Measurements in Problems of Biological Classification," *J. Royal Statistical Soc., B*, vol. 10, pp. 159-203, 1948.
- [3] S. Raschka, "Linear Discriminant Analysis," *sebastianraschka.com*, 03-Aug-2014.
- [4] J. W. Sammon, "A Nonlinear Mapping for Data Structure Analysis," in *IEEE Transactions on Computers*, vol. C-18, no. 5, pp. 401-409, May 1969, doi: 10.1109/T-C.1969.222678.
- [5] J. B. Tenenbaum, V. de Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [6] O. Samko, A.D. Marshall, and P.L. Rosin, "Selection of the optimal parameter value for the Isomap algorithm," *Pattern Recognition Letters*. 27. 968-979. 10.1016/j.patrec.2005.11.017.
- [7] L. F. Terán Tamayo, *Smartparticipation: A fuzzy-based recommender system for political community-building*. Not Avail, 2014.
- [8] Pollard, T. (2014, April 18). Sammon mapping in Python. computer software. Retrieved from <https://github.com/tompollard/sammon>.