

Feature\_Selection-09



## **DEPARTMENT OF ELECTRICAL AND COMPUTER ENGINEERING**

### **Applied Machine Learning**

ELEC 8900-76-R-2021F

**Date: Nov 29<sup>th</sup>, 2021**

**“Feature Selection\_09”**

**(FS-09)**

**Instructor:**

Prof. Roozbeh Razavi Far

**Submitted by:**

**Kenil Doshi                  ID: 110031636**

**Meghal Kumar Darji        ID: 110041749**

**University of Windsor**

**ON, CANADA**

## **ABSTRACT:**

Artificial Intelligence is the most discussed topic in the industry at the moment. Compute scientists all over the world are working on research and development of Machine Learning algorithms for application in various domains. Machine learning is one of the subsets of Artificial Intelligence. Machine learning (ML) attempts to understand learn algorithms and mathematical models used in software systems from patterns observed in data. It tries to develop an algorithmic structure from training data to use it later on future inputs for a specific purpose. Machine learning applications can be seen in many fields such as Computer Vision, Financial Forecasting, Robotics, and many more. Machine learning is also being used in making predictions and data analysis. Machine learning has several learning approaches such as supervised, unsupervised, and reinforcement learning. Machine learning has a primary goal of optimization and to generalize the learning algorithm based on its training to get possible result on test data. The machine learning process starts with learning task or problem. Here, the dataset under consideration is used for training, where the algorithm tries to learn the unknown target function, to make predictions on the test dataset. This project focuses on 3 feature selection methods:

- i. Unsupervised Discriminative Feature Selection [UDFS]
- ii. Local Learning-based Clustering Feature Selection [LLCFS]
- iii. Correlation-based Feature Selection [CFS]

To perform classification the following algorithms are used:

- i. Naive Bayes Classifier [NBC]
- ii. Support Vector Machines [SVM]
- iii. Decision Trees [DT]
- iv. Multi-Layer Perceptron [MLP].

The main aim of this project is to try the different combinations of feature selection and classification methods to find the best results on the given data.

### **Keywords**

Classifier, Feature Selection, Parameter Tuning, Decision Tree, Support Vector machines, Hyper Parameter, Multi-layer Perceptron, Naive Bayes classifier, Machine Learning, F1 Score, confusion Matrix, MATLAB, python

### **Development Software**

MATLAB R2019b & MATLAB R2020a

Python, Jupyter Notebook

## **INTRODUCTION**

Machine learning algorithms try to learn patterns from the available data to make predictions on similar but new data. It generally used when dealing with large data sets that has some kind of pattern which cannot be explained by mathematical equations. Machine learning has two main approaches: supervised and unsupervised learning. Supervised learning involves learning from data for which the ground truth labels are known unlike unsupervised learning unsupervised learning where no information is available about the ground truth labels.. Supervised learning has 2 main categories of learning which are classification and regression. Classification simply classifies given data to give discrete responses. On the other hands, regression predicts continuous values.

Feature selection is also known as variable or attribute selection is in demand to the due formation of huge datasets and better machine learning techniques. Feature selection focuses on selecting a method automatic in nature for best possible attribute selection aligned of selecting relevant to data and predictive modeling. The benefits of feature selection are that it reduces over fitting and training time while providing an opportunity to improve model performance.

### **Feature Selection Methods:**

#### **1. Unsupervised Discriminative Feature Selection [UDFS]**

UDFS aims to select the most discriminative features for data representation, where manifold structure is considered making it different from the other feature selection algorithms. It is useful for feature selection in data without labels. UDFS algorithm selects the most discriminative feature subset from the whole feature set in batch mode. The UDFS algorithm aims to simultaneously exploiting discriminative information and feature correlations. Use one-step approach to select the most discriminative approach for data representation. UDFS analyzes features jointly and local structure of data distribution.

#### **2. Local Learning-based Clustering Feature Selection [LLCFS]**

The key idea here is to break an arbitrarily complex nonlinear problem into a set of locally linear sets through local learning method leading to learning feature relevance globally within the large margin framework. This Feature Selection model has a logarithmic space complexity with respect to a number of features.

#### **3. Correlation-based Feature Selection [CFS]**

The Correlation based Feature Selection (CFS) algorithm selects features for classification tasks based on the correlation between features and the class label. A feature is considered useful if it is correlated with or predictive of the output label and at the same time meekly correlated with other features in the dataset. CFS algorithm greedily adds feature to a feature subset and then evaluates the merit of this subset.

## Feature\_Selection-09

Finally the feature subset with the highest merit score is selected when no improvements in the merit score are observed after 5 consecutive iterations. Merit of feature subset is calculated using the following formula

$$\text{Merit}_{S_k} = \frac{k\bar{r}_{cf}}{\sqrt{k + k(k-1)\bar{r}_{ff}}}.$$

CFS Optimization Problem:

$$\text{CFS} = \max_{S_k} \left[ \frac{r_{cf_1} + r_{cf_2} + \dots + r_{cf_k}}{\sqrt{k + 2(r_{f_1f_2} + \dots + r_{f_if_j} + \dots + r_{fkf_1})}} \right].$$

Where  $k$  is the number of features in the subset;

$r_{cf}$  is the average value of all feature-classification correlations

$r_{ff}$  is the average value of all feature-feature correlations

## Classification Methods:

### a. Naive Bayes Classifier [NBC]

The Naive Bayes Classifier is based on Bayes theorem with the assumption that features selected are independent. There are three main types of NBC's: Multinomial Naive Bayes, Bernoulli Naive Bayes, and Gaussian Naive Bayes. This project uses the Gaussian Naïve Bayes Classifier for all classification tasks.

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})} \text{ posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

$$\hat{y} = \underset{k \in \{1, \dots, K\}}{\text{argmax}} p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

### b. Support Vector Machines [SVM]

Support Vector Machines are non-probabilistic binary linear classifiers and are supervised learning models. It works by dividing well-defined wide gap points in space with separate categories.

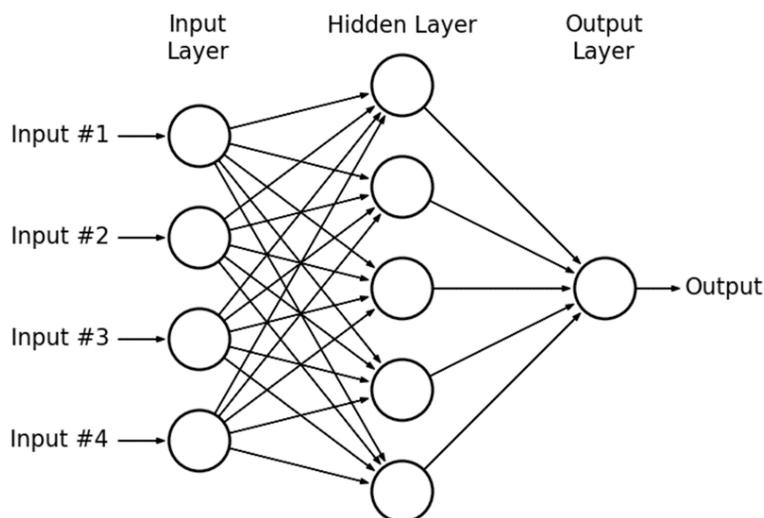
$$\left[ \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i - b)) \right] + \lambda \|w\|^2.$$

### c. Decision Trees [DT]

Decision Trees are used to display an algorithm with condition control statements and are decision support tools for best possible outcomes with a tree-like model of decisions. It has a structured algorithm with an easy understanding of data classification flow shows the possible consequences. Uses a predictive model on observations of elements to find target elements value

### d. Multi-Layer Perceptron [MLP]

Multi-Layer Perceptron (MLP) is a supplement of feed forward neural network. It consists of three types of layers—the input layer, output layer and hidden layer, as shown in Fig. 1. The input layer receives the input signal to be processed. The required task such as prediction and classification is performed by the output layer. An arbitrary number of hidden layers that are placed in between the input and output layer are the true computational engine of the MLP. Similar to a feed forward network in a MLP the data flows in the forward direction from input to output layer.



**Figure 1: Multi-Layer Perceptron**

## **Implementation Procedure:**

- For this project we primarily use Python 3 programming language in the Jupyter Notebook along with MATLAB for one of the feature selection methods. For certain computationally expensive tasks Google Colaboratory (Google Colab) was used for faster execution.
- Before using the given datasets for training and prediction, datasets were examined for missing values. In one the datasets missing value were found which required proper handling to avoid errors and exceptions.
- Next, the feature values and the class labels are read and separated from the provided data files. This training data is then used to obtain baseline classification accuracy and f1 scores with 10-fold stratified cross validation for all the classification algorithms mentioned above.
- After getting the baseline score, the next step in pipe line is to compute the feature rankings using the feature selection methods individually. Several features are then chosen based on the obtained rankings to form a new feature set. This reduced feature set is then used to again perform classification tasks using the same algorithms.
- All training and testing done in this project is done with 10-fold stratified cross-validation to ensure that class balance in the folds is maintained in-case of an imbalanced dataset.
- The accuracy and f1 scores for different classification algorithms are compared to analyze the improvements and identify the combination of feature selection and classification methods that provides improved results for the given datasets.

## **Key MATLAB and ‘python’ files used for Program:**

The project is completed using the following code files:

1. **baseline.ipynb**: Jupyter notebook code file to obtain baseline classification accuracy and f1 scores on all datasets using the mentioned classification algorithms.
2. **CFS.ipynb**: Jupyter notebook code file to perform Correlation Based Classification.
3. **Classification\_after\_udfs.ipynb**: Jupyter notebook code file to perform classification using the feature rankings obtained from UDFS method.
4. **UDFS.ipynb**: Jupyter notebook code file to perform UDFS on the given datasets.
5. **dataset\_sampling.ipynb**: Jupyter notebook code file to sample data from the given datasets to reduce dataset size. (To successfully perform computationally heavy feature selection methods).

## Feature\_Selection-09

6. **MLP.ipynb:** Jupyter notebook code file to perform classification and obtain baseline measures using MLP on all data.
7. **MLP\_UDFS.ipynb:** Jupyter notebook code file to perform classification after feature selection by UDFS method.
8. **FS09\_demo.ipynb:** Jupyter notebook code file for class demonstration and to obtain classification results after feature selection for all methods and classifier algorithms.
9. **K\_demo.m:** MATLAB code file to obtain feature rankings using LLCFS.

### **Key Parameters for tuning the Classifiers:**

All the classifier methods in this project are implemented using Scikit-learn library in python 3.

#### **1) Support Vector Machines: `SVC(kernel='rbf', gamma=0.5, C=10000)`**

SVM's **kernel** shows which type of function we use and according to that our model decide it to linear or nonlinear. **Degree** gives information for the nonlinear degree of the function, and **max iteration** gives info about spoke for higher accurate model.

#### **2) Decision Tree: `DecisionTreeClassifier()`**

Here, no tuning parameters are used for Decision Tree Classifier.

#### **3) Naïve Bayes Classifier: `GaussianNB()`**

In our case, we do not use any tuning parameter. But the Gaussian Naive bayes performs poorly when compared to other methods in this project.

#### **4) Multilayer Perceptron:**

`MLPClassifier(hidden_layer_sizes=(10,10,10), max_iter=1000)`

The MLP used here consists of three hidden layers with 10 neurons in each hidden layer. Furthermore, the MLP Iterated for a maximum of 1000 times during training before moving in to the next epoch.

## **Feature Selection Algorithms:**

### **1) UDFS (Unsupervised Discriminative Feature Selection.)**

for i = 1 to n do

- 1 Bi = ( ~XTi~Xi + λI)−1
- 2 Mi = SiHk+1BiHk+1STi ;
- 3 M = X\_ni=1MiXT ;
- 4 Set t = 0 and initialize D0 ∈ Rd×d as an identity matrix;
- 5 Repeat
- 6 Pt = M + γDt;
- 7 Wt = [p1, ..., pc] where p1, ..., pc are the eigenvectors of Pt corresponding to the first c smallest eigenvalues;
- 8 Update the diagonal matrix Dt+1 as Dt+1 = |||12\_w1t\_2...12\_wdt\_2|||;
- 9 t = t + 1;
- 10 until Convergence;
- 11 Sort each feature fi |di=1 according to \_ wit \_ 2indescending order and select the top ranked ones.

### **2) LLCFS ( Feature selection for local learning-based clustering algorithm.)**

**input:** X ¼fxign i¼1, size of the neighborhood k, trade-off parameter (Beta)

**output:** Y;

- 1 Initialize \_1 ¼ 1d , for 1 ¼ 1; . . . ; d;
- 2 while not converge do
- 3 Find k-mutual neighbors for fxign i¼1, using the metric defined);
4. Construct the matrix M in \_i given ,
5. then solve the problem to obtain Y;
6. Compute wc \_ i ; 8i; c and update it for next epoch);
7. End

### **3) CFS (Correlation Feature Selection)**

A correlation-based approach towards hypothesis focused on good feature sets containing features that are high correlated with the class, yet being uncorrelated with each other

CFS is an algorithm using a feature evaluation formula based on ideas from test theory that provides an operational definition of this hypothesis with an appropriate correlation measure and a heuristic search strategy CFS executes faster than other wrapper algorithms Merit of feature subset

## Feature\_Selection-09

On the base of following merits it find the correlation between the set and returns the feature which highly correlate with the dataset.

$$\text{Merit}_{S_k} = \frac{k\bar{r}_{cf}}{\sqrt{k + k(k-1)\bar{r}_{ff}}}.$$

### ❖ Common Algorithm of Classification:

#### Steps:

1. Discretize a copy of training set using Fayyas and Irani's method.
2. Search the feature space using either 1 of the 3 searching scheme: forward, backward, or best first
3. Evaluate the merit of a feature using Merit(S<sub>k</sub>) equation .
4. Use the features subset with the highest merit to reduce the dimensionality of the training and test set.
5. Perform training and testing for the machine learning task using the reduced data sets.

## Feature\_Selection-09

### Sample Results:

#### Accuracy and F1 Score with SVM classifier:

The screenshot shows a Jupyter Notebook interface with the title "jupyter Baseline Results". The code cell contains Python code for an SVM classifier across 10 folds. The output cell displays accuracy and F1 scores for each fold.

```
In [ ]: # SVM with C=10000 (Good results)
for i in range(1,14):
    Data = pd.read_csv('Datasets/D'+str(i)+'.csv', header=None)
    print('\nDataset D'+str(i))
    if i==1:
        Data.drop([92], axis=1, inplace=True)
    Xi = Data[Data.columns[:-1]]
    Yi = Data[Data.columns[-1]]
    kf = StratifiedKFold(n_splits=10, shuffle=True, random_state=1)
    svm = SVC(kernel='rbf', gamma=0.5, C=10000)
    j=1
    for train_idx, test_idx in kf.split(Xi,Yi):
        x_train, x_test = Xi.iloc[train_idx],Xi.iloc[test_idx]
        y_train, y_test = Yi.iloc[train_idx],Yi.iloc[test_idx]
        svm.fit(x_train, y_train)
        y_pred = svm.predict(x_test)
        acc = svm.score(x_test, y_test)
        f1 = f1_score(y_test, y_pred, average='weighted')
        print('\t\tacc fold {} : {}'.format(j), acc)
        print('\t\tf1 fold{} : {}'.format(j), f1)
        print()
        j+=1

Dataset D1
acc fold 1: 0.8497142857142858
f1 fold1: 0.8512810266952605

acc fold 2: 0.8594285714285714
f1 fold2: 0.85969633365640317

acc fold 3: 0.8582857142857143
f1 fold3: 0.8585677637853654

acc fold 4: 0.8571428571428571
f1 fold4: 0.8569298276994294

acc fold 5: 0.8565714285714285
f1 fold5: 0.8558166938992802

acc fold 6: 0.8771428571428571
f1 fold6: 0.8765651903170625
```

#### Accuracy and F1 Score with SVM classifier:

The screenshot shows a Jupyter Notebook interface with the title "jupyter Project MLP". The code cell contains Python code for an SVM classifier across 10 folds. The output cell displays accuracy and F1 scores for each fold.

```
In [ ]: #train, y_test = y.iloc[train_idx], y.iloc[test_idx]
scaler.fit(x_train)
x_train = scaler.transform(x_train)
x_test = scaler.transform(x_test)
mlp.fit(x_train, y_train)
y_pred = mlp.predict(x_test)
acc = accuracy_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred, average='weighted')
print('\t\tacc:', acc)
print('\t\tf1:', f1)
j+=1

Data D1
fold1
    ACC: 0.8011428571428572
    F1: 0.79838902436652
fold2
    ACC: 0.8154285714285714
    F1: 0.8137804043540515
fold3
    ACC: 0.8262857142857143
    F1: 0.826971002018799
fold4
    ACC: 0.8097142857142857
    F1: 0.8063211947517024
fold5
    ACC: 0.8525714285714285
    F1: 0.8505494497920253
fold6
    ACC: 0.8194285714285714
    F1: 0.8158276603371556
```

## Feature\_Selection-09

### Classification accuracy and f1 scores after feature selection using UDFS:

In [81]: # Decision Trees  
n\_features\_t = 25  
Dt = DecisionTreeClassifier()  
  
d1 = pd.read\_csv('Datasets/D1.csv', header=None)  
x1 = d1[d1.columns[:-1]]  
y1 = d1[d1.columns[-1]]  
with open('udfs\_results/D1.txt') as file:  
 data = file.read()  
L = data.replace('\n','\r\n').split('\r\n')  
sel = [int(l) for l in L if l != '']  
sel = [int(l) for l in sel if l != '']  
  
X\_selected = x1.iloc[:,sel[n\_features\_dt:]]  
kf = StratifiedKFold(n\_splits=10, random\_state=1, shuffle=True)  
j = 1  
for train\_idx, test\_idx in kf.split(X\_selected,y1):  
 print('fold {}.'.format(j))  
 x\_train, x\_test = X\_selected.iloc[train\_idx], X\_selected.iloc[test\_idx]  
 y\_train, y\_test = y1.iloc[train\_idx], y1.iloc[test\_idx]  
 Dt.fit(x\_train, y\_train)  
 y\_pred = Dt.predict(x\_test)  
 f1 = f1\_score(y\_test, y\_pred, average='weighted')  
 acc = accuracy\_score(y\_test, y\_pred)  
 print('({}acc:{} acc{})'.format(j, acc))  
 print('({}f1:{} f1{})'.format(j, f1))  
 j+=1  
  
fold 1  
Acc: 0.96  
F1: 0.9602978013216178  
fold 2  
Acc: 0.9537142857142857  
F1: 0.9539376327791682  
fold 3  
Acc: 0.9531428571428572  
F1: 0.9531866059542204  
fold 4  
Acc: 0.9554285714285714  
F1: 0.9550824659551795  
fold 5  
Acc: 0.9622857142857143  
F1: 0.9620849634589707  
fold 6  
Acc: 0.956  
F1: 0.956782532096203  
fold 7  
Acc: 0.957142857142857  
F1: 0.9573604089130801  
fold 8  
Acc: 0.9497142857142857  
F1: 0.949787009813855  
fold 9  
Acc: 0.9638571428571429  
F1: 0.962949062359083  
fold 10  
Acc: 0.9594285714285714  
F1: 0.9594270727261591

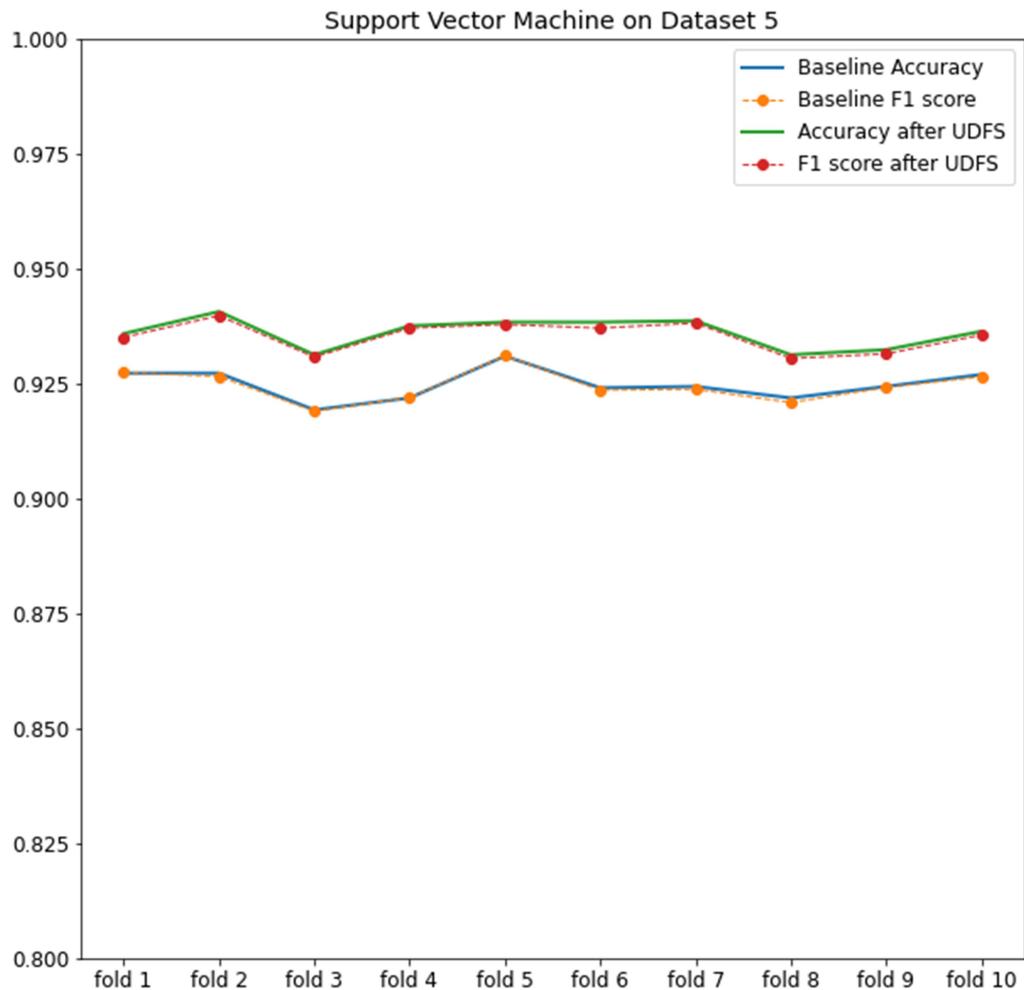
### Baseline Accuracy and F1 Score for NB Classifier:

In [\*]: # separate classifier for all datasets  
for i in range(1,14):  
 filename = 'Datasets/D' + str(i) + '.csv'  
 Data = pd.read\_csv(filename, header=None)  
 if i==1:  
 Data.drop([92], axis=1, inplace=True)  
 X1 = Data[Data.columns[:-1]]  
 Y1 = Data[Data.columns[-1]]  
  
 c1f = GaussianNB()  
  
 kf = StratifiedKFold(n\_splits=10, shuffle=True)  
 print("Dataset {} split into folds. Training now...".format(i))  
 j = 1  
 for train\_idx, test\_idx in kf.split(X1, Y1):  
 x\_train, x\_test = X1.iloc[train\_idx], X1.iloc[test\_idx]  
 y\_train, y\_test = Y1.iloc[train\_idx], Y1.iloc[test\_idx]  
  
 c1f.fit(x\_train, y\_train)  
 y\_pred = c1f.predict(x\_test)  
 acc = accuracy\_score(y\_test, y\_pred)  
 f1 = f1\_score(y\_test, y\_pred, average='weighted')  
 print('({}Accuracy:{} acc{})'.format(j, acc))  
 print('({}f1:{} f1{})'.format(j, f1))  
 j+=1  
 print('\n')  
  
Dataset D1 split into folds. Training now...  
Accuracy(fold 1): 0.5365714285714286  
f1 (fold1) 0.5181082317881291  
Accuracy(fold 2): 0.5194285714285715  
f1 (fold2) 0.508405827739860  
Accuracy(fold 3): 0.545059171428572  
f1 (fold3) 0.5328135158905793  
Accuracy(fold 4): 0.5234285714285715  
f1 (fold4) 0.50399170427362  
Accuracy(fold 5): 0.552  
f1 (fold5) 0.517348571428572322307  
Accuracy(fold 6): 0.5405714285714286  
f1 (fold6) 0.5249557844058634  
Accuracy(fold 7): 0.5217142857142857  
f1 (fold7) 0.5173485714285728  
Accuracy(fold 8): 0.521734857142857  
f1 (fold8) 0.505350146795752  
Accuracy(fold 9): 0.5194285714285715  
f1 (fold9) 0.4979122703724286  
Accuracy(fold 10): 0.5068571428571429  
f1 (fold10) 0.48909581616129846

## Feature\_Selection-09

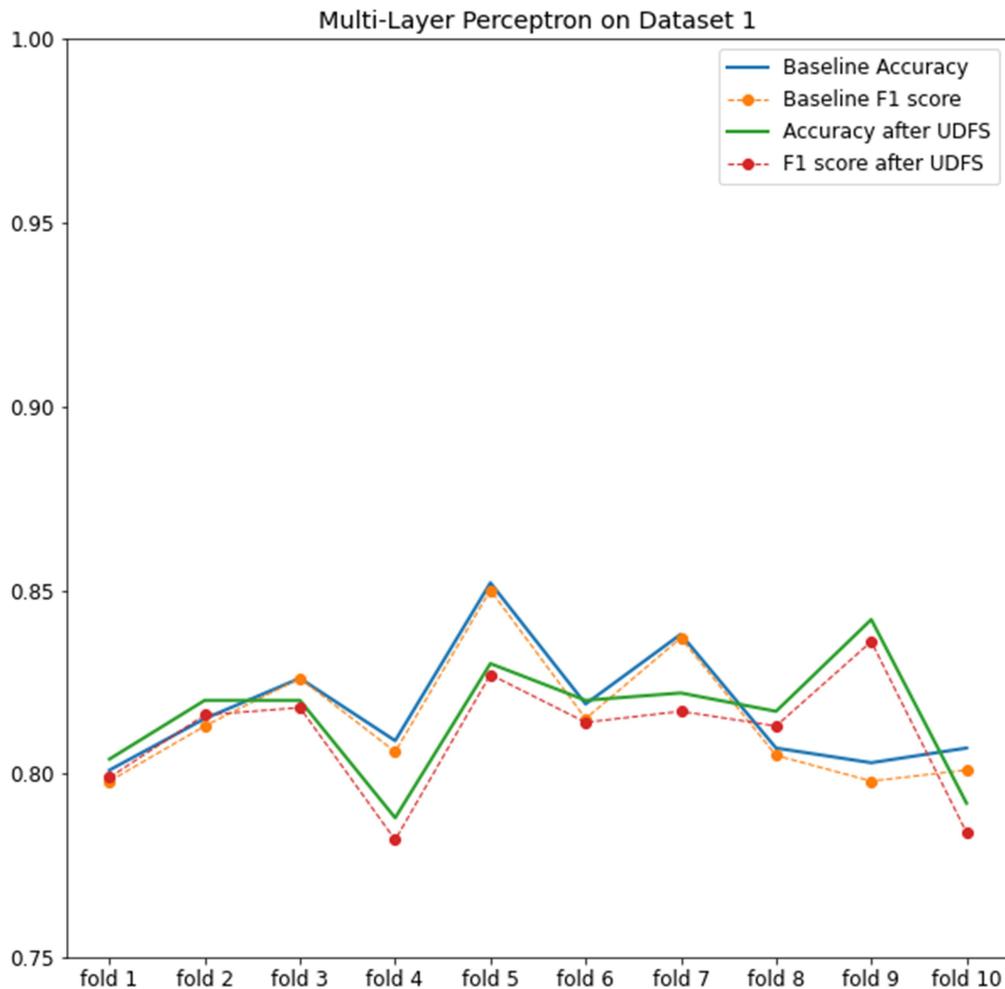
### Visualized accuracy and f1 score measures:

The graph below shows the accuracy and f1 score measures on dataset D5 before and after feature selection. It can be seen that performance improves after feature selection.



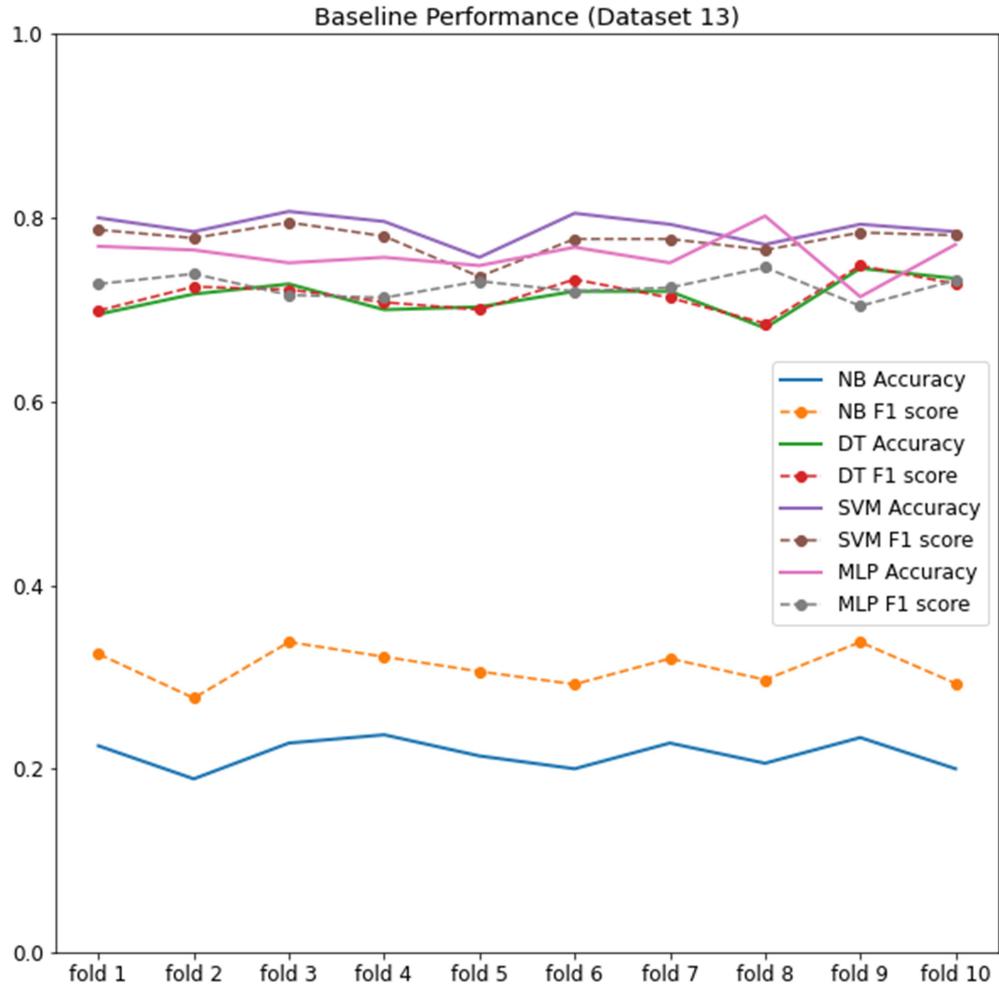
## Feature\_Selection-09

The graph below shows the accuracy and f1 score measures for dataset D1 and MLP classifier using UDFS.

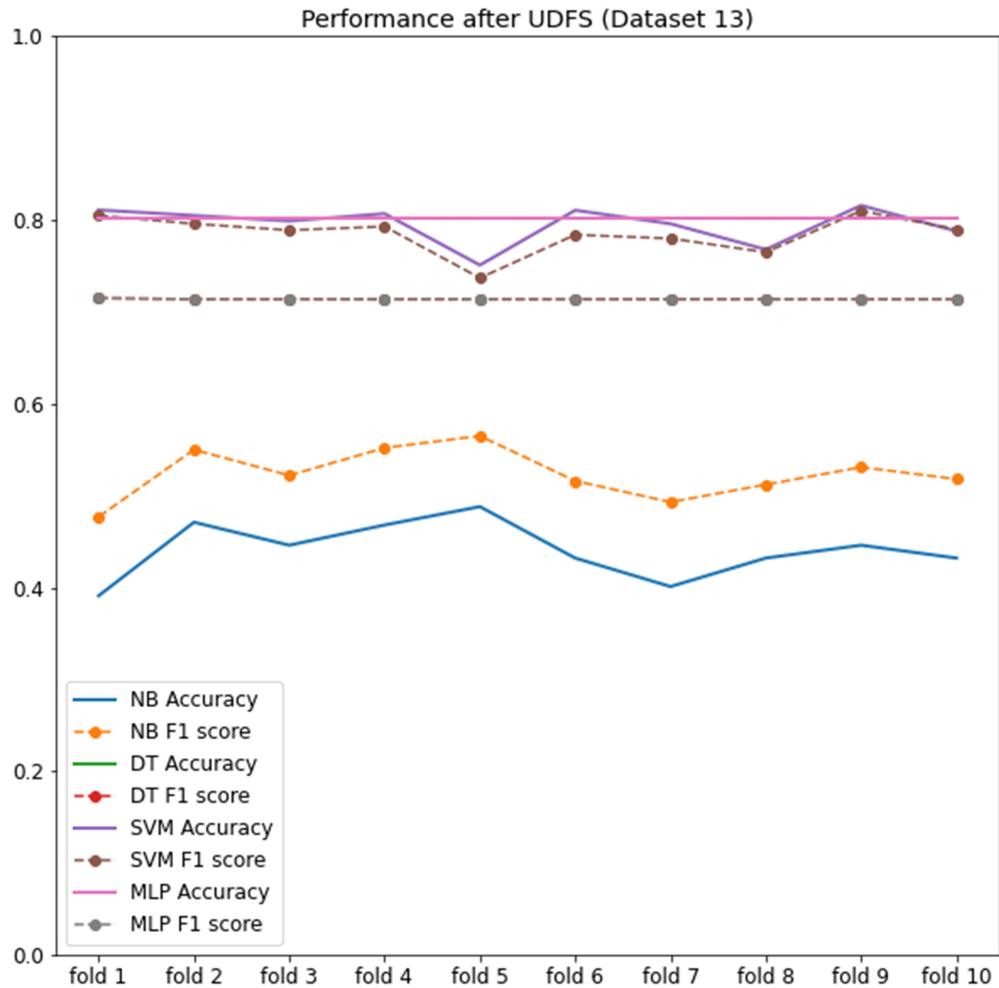


The graphs below compare the accuracy and f1 score measures for all classification methods on Dataset D13 before and after performing feature selection using UDFS.

## Feature\_Selection-09



## Feature\_Selection-09



**References:**

- [1] [ Yi Yang1, Heng Tao Shen1, Zhigang Ma2, Zi Huang1, Xiaofang Zhou1 (2019). *Norm Regularized Discriminative Feature Selection for Unsupervised Learning*. [online] Ijcai.org. Available at: <https://www.ijcai.org/Proceedings/11/Papers/267.pdf> [Accessed 14 Jul. 2019].
- [2] YM, Z. (2019). *Feature Selection and Kernel Learning for Local Learning-Based Clustering*. - PubMed - NCBI. [online] Ncbi.nlm.nih.gov. Available at: <https://www.ncbi.nlm.nih.gov/pubmed/21135434> [Accessed 16 Jul. 2019].
- [3] Hall, M. (2019). *Correlation-based Feature Selection for Machine Learning*. [online] Cs.waikato.ac.nz. Available at: <https://www.cs.waikato.ac.nz/~mhall/thesis.pdf> [Accessed 14 Jul. 2019].