Final Term Project Report

Name: Vignesh Surya Mantha

NJIT UCID: vm675

Email Address: vm675@njit.edu

Date: 23rd Nov 2024

Professor: Yasser Abduallah

Course: CS 634 Data Mining (101)

## Abstract

This project aims to evaluate and compare the performance of different machine learning algorithms in solving a binary classification problem using the **Pima Indians Diabetes Dataset**. The dataset focuses on predicting the likelihood of diabetes based on diagnostic measures such as glucose levels, BMI, age and other health related attributes

i implemented three algorithms: **K-Nearest Neighbors (KNN)**, **Random Forest (RF)**, and **Support Vector Machine (SVM)** along with a deep learning algorithm **Long Short-Term Memory (LSTM)** to ensure a fair and robust evaluation and i used **10-fold cross-validation** which systematically splits the dataset into training and testing subsets. performance metrics such as **True Positive Rate (TPR)**, **False Positive Rate (FPR)**, **F1 Score**, **Accuracy**, **AUC**, and others were computed manually for a detailed assessment

the study provides insights into the relative strengths and weaknesses of these algorithms for medical data classification tasks additionally the project emphasizes the importance of feature standardization, performance evaluation, and visualization in achieving meaningful results. Based on my findings the **Random Forest** algorithm demonstrated the best overall performance making it a promising choice for diabetes prediction

## Tools and Libraries Used

- **Programming Language:** Python 3.9
- **Integrated Development Environment (IDE):** Jupyter Notebook
- **Key Libraries and Frameworks:**

- **Data Manipulation and Visualization:**
  - **Pandas**: For data manipulation and analysis.

- o **Matplotlib** & **Seaborn**: For data visualization.
- **Data Preprocessing and Scaling:**
  - o **Scikit-learn**: For standardization, train-test split, and cross-validation.
- **Classification Algorithms:**
  - o **Scikit-learn**: Implemented KNN, Random Forest, and SVM.
  - o **TensorFlow/Keras**: For building and training the LSTM model.
- **Evaluation Metrics:**
  - o **Scikit-learn**: For confusion matrix, ROC curve, AUC, and Brier score.
- **Other Utilities:**
  - o **NumPy**: For numerical operations.
  - o **StandardScaler**: For feature normalization.

## Introduction

Binary classification is a critical aspect of supervised learning where the goal is to categorize data into one of two predefined classes this project explores the application of various supervised learning algorithms to a healthcare related dataset focused on predicting diabetes outcomes the chosen dataset contains medical attributes like glucose levels, BMI, and age, making it a practical and impactful choice for this task

the primary objective of this project is to evaluate the performance of four algorithms K-Nearest Neighbors (KNN), Support Vector Machines (SVM), Random Forest (RF) and Long Short-Term Memory networks (LSTM) on the same dataset, by employing a rigorous 10-fold cross-validation approach, the study aims to ensure robust evaluation and generalizability of the results

this analysis also delves into manually calculated metrics such as True Positive Rate (TPR), False Positive Rate (FPR), and others to assesss the algorithms comprehensively, so by comparing performance metrics and visualizing results through ROC curves, the study identifies the most effective algorithm for this binary classification task \

## Dataset Description

The dataset used in this project is focused on predicting diabetes outcomes in patients. Diabetes is a chronic medical condition characterized by high levels of sugar (glucose) in the blood, It is a growing global health issue often leading to severe complications if left untreated such as heart disease, kidney failure, and vision loss, early diagnosis and management are crucial in mitigating its impact

this dataset contains several medical attributes related to patients health and lifestyle which can help determine whether a person has diabetes, the dataset comprises the following features:

1. **Pregnancies**: The number of times a patient has been pregnant

2. **Glucose**: Plasma glucose concentration after a 2-hour oral glucose tolerance test

3. **BloodPressure**: Diastolic blood pressure (mm Hg)

4. **SkinThickness**: Triceps skinfold thickness (mm)

5. **Insulin**: 2-hour serum insulin (mu U/ml)

6. **BMI**: Body mass index, calculated as weight in kg/(height in m)^2

7. **DiabetesPedigreeFunction**: A function that scores the likelihood of diabetes based on family history

8. **Age**: Age of the patient (in years)

The **Outcome** column serves as the target variable:

- 1: Indicates the presence of diabetes.

- 0: Indicates the absence of diabetes.

```
### Data visualization - target variable distribution ###

plt.figure(figsize=(8, 6))
sns.countplot(x=labels, palette=['skyblue', 'crimson'])
plt.title('Distribution of Target Variable', fontsize=16)
plt.xlabel('Outcome', fontsize=14)
plt.ylabel('Count', fontsize=14)
plt.xticks([0, 1], ['No Diabetes', 'Diabetes'], fontsize=12)
for i, count in enumerate(labels.value_counts()):
    plt.text(i, count + 5, f'{count}', ha='center', fontsize=12, color='black')
plt.show()

# Pie chart for target variable distribution with legend
plt.figure(figsize=(8, 6))
labels_counts = labels.value_counts()
plt.pie(labels_counts, labels=['No Diabetes', 'Diabetes'], autopct='%1.1f%%',
        startangle=90, colors=['skyblue', 'crimson'], textprops={'fontsize': 12})
plt.title('Target Variable Distribution (Pie Chart)', fontsize=16)

# Adding legend
plt.legend(['No Diabetes', 'Diabetes'], loc='upper right', bbox_to_anchor=(1.2, 1), fontsize=12)
plt.show()

# Print the data imbalance
positive_outcomes, negative_outcomes = labels_counts
total_samples = labels.count()
print(f"Number of Positive Outcomes (Diabetes): {positive_outcomes} ({round((positive_outcomes/total_samples)*100, 2)}%)")
print(f"Number of Negative Outcomes (No Diabetes): {negative_outcomes} ({round((negative_outcomes/total_samples)*100, 2)}%)")
```
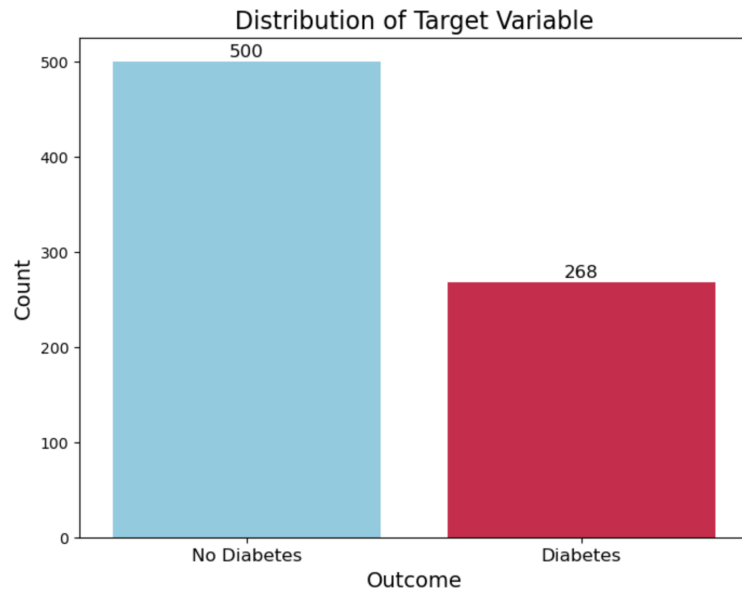
Figure 1: Bar graph showing the distribution of the target variable (Diabetes vs No Diabetes)

As depicted in **Figure 1** the bar graph highlights the distribution of the target variable, showing that a larger portion of the dataset belongs to the 'No Diabetes' category.



```
Number of Positive Outcomes (Diabetes): 500 (65.1%)
Number of Negative Outcomes (No Diabetes): 268 (34.9%)
```
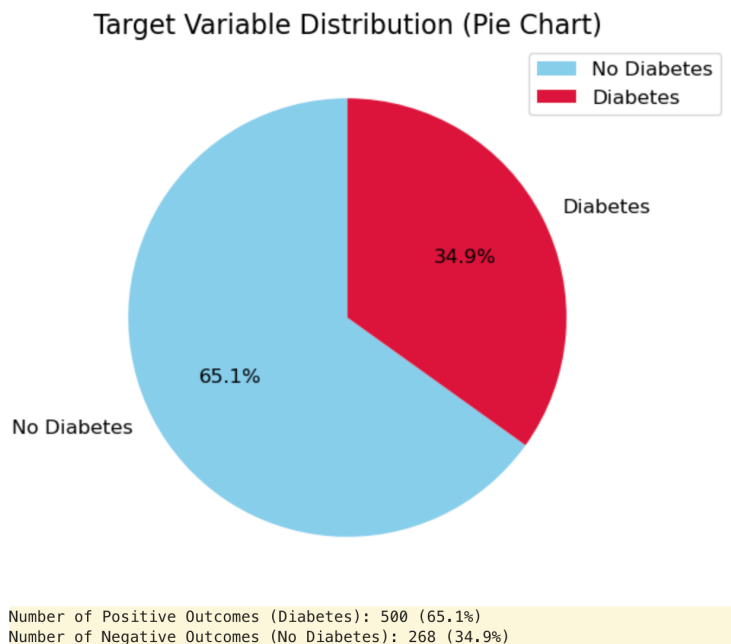
Figure 2: Pie chart representing the distribution of the target variable (Diabetes vs. No Diabetes)

The distribution is further visualized in **Figure 2**, where the pie chart clearly illustrates the proportions of the two categories.

## Methodology

### Data Preprocessing

1. **Handling Missing Values**:

   o the dataset contains some missing values for specific features (e.g., insulin levels), i replaced these missing values with the median of their respective columns to ensure consistency without introducing bias

2. **Standardization**:

   o the features were standardized to ensure that all attributes have a mean of 0 and a standard deviation of 1, this step was essential for algorithms sensitive to feature scaling, such as KNN and SVM

3. **Train-Test Split**:

   o The dataset was split into training and testing sets. The training set was used to train the models, and the testing set was used for evaluation.

```
### Normalize the Data ###

# Standardize features for training set
features_train_all_std = (features_train_all - features_train_all.mean()) / features_train_all.std()

# Standardize features for testing set
features_test_all_std = (features_test_all - features_test_all.mean()) / features_test_all.std()

# Display a summary of the standardized data
features_train_all_std.describe()
```
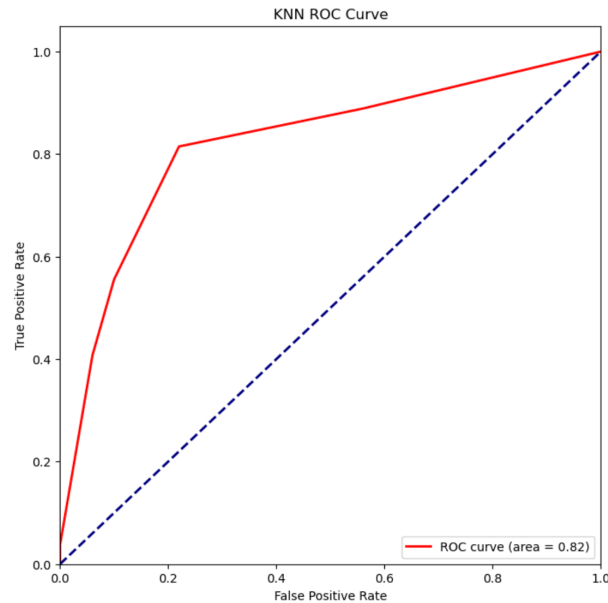
|       | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|-------|-------------|---------|---------------|---------------|---------|-----|--------------------------|-----|
| count | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 | 6.910000e+02 |
| mean | -8.868931e-17 | -6.555297e-17 | -5.552722e-16 | -1.336766e-16 | -2.827775e-17 | -4.858632e-16 | -2.879189e-16 | 2.930603e-16 |
| std | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 | 1.000000e+00 |
| min | -1.144835e+00 | -2.563245e+00 | -3.943501e+00 | -2.509887e+00 | -1.479786e+00 | -2.102822e+00 | -1.197605e+00 | -1.055424e+00 |
| 25% | -8.494354e-01 | -7.468983e-01 | -6.693918e-01 | -4.547208e-01 | -2.269315e-01 | -7.073001e-01 | -6.932720e-01 | -8.003056e-01 |
| 50% | -2.586353e-01 | -1.524575e-01 | -1.457002e-02 | 1.982794e-03 | -1.672718e-01 | -9.539238e-03 | -2.888971e-01 | -3.751086e-01 |
| 75% | 6.275647e-01 | 5.905936e-01 | 6.402517e-01 | 3.445105e-01 | -1.374419e-01 | 6.065688e-01 | 4.683592e-01 | 6.453641e-01 |
| max | 3.876965e+00 | 2.555551e+00 | 4.078066e+00 | 7.994296e+00 | 8.435663e+00 | 5.156861e+00 | 5.620731e+00 | 4.046940e+00 |

The above figure shows the summary statistics of the standardized features, showing a mean of 0 and standard deviation of 1 after normalization.

## Algorithms Used

1. **K-Nearest Neighbors (KNN):**

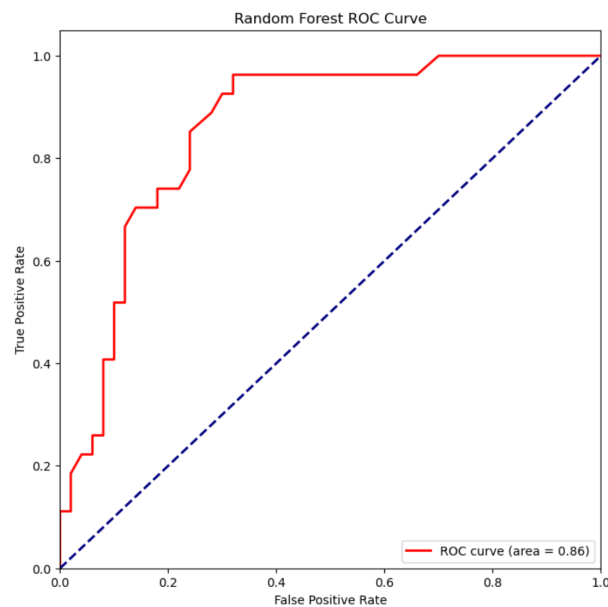   o KNN is a simple instance based learning algorithm that classifies data points based on their proximity to labeled examples in the feature space
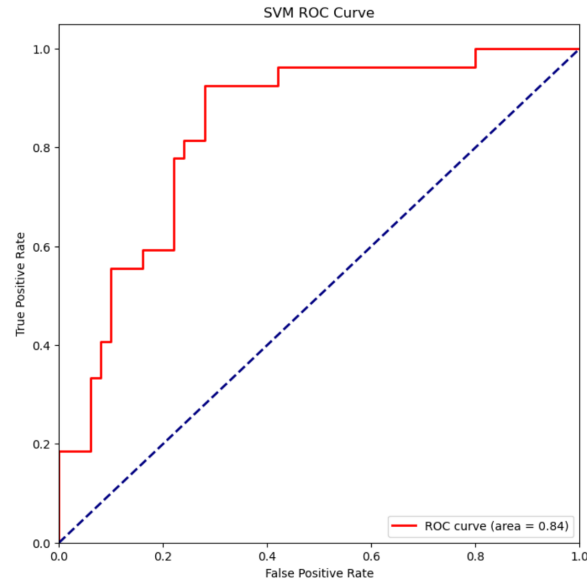


2. **Random Forest (RF):**

   o RF is an ensemble learning method that constructs multiple decision trees during training and outputs the class with the most votes among the trees, It is robust to overfitting and handles high dimensional data well

3. **Support Vector Machine (SVM):**

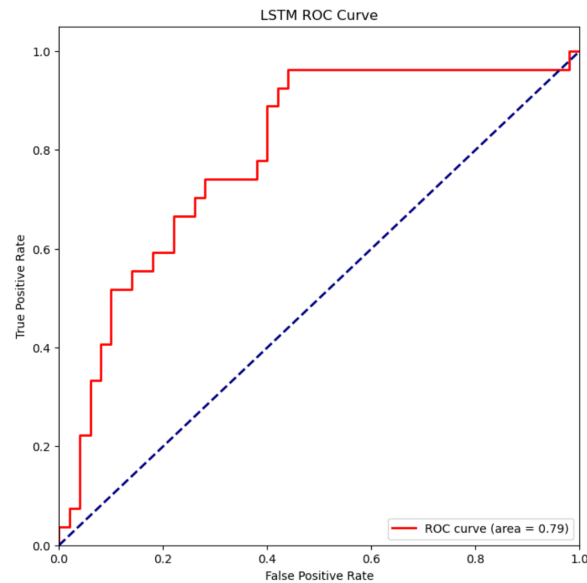   o SVM constructs a hyperplane that best separates data points of different classes. We used the linear kernel to handle the binary classification task effectively



4. **Long Short-Term Memory (LSTM):**

   o LSTM is a type of recurrent neural network (RNN) designed to handle sequential data, despite the dataset not being inherently sequential, LSTM was used as a deep learning approach for comparison

## Evaluation Strategy

1. **10-Fold Cross-Validation**:

   - The dataset was divided into 10 subsets (folds) so each fold was used once as a testing set while the remaining folds formed the training set. This approach ensures that all data points are used for both training and testing providing a robust evaluation

2. **Performance Metrics**:

   - The following metrics were calculated manually using confusion matrix values:

     - **True Positive (TP)**: correctly classified diabetic cases

     - **True Negative (TN)**: correctly classified non-diabetic cases

     - **False Positive (FP)**: non-diabetic cases incorrectly classified as diabetic

     - **False Negative (FN)**: niabetic cases incorrectly classified as non-diabetic

     - Derived metrics:

       - Accuracy, Precision, F1-Score, True Positive Rate (TPR), False Positive Rate (FPR), Brier Score, AUC (Area Under the Curve), and others

3. **Visualization**:

   - ROC (Receiver Operating Characteristic) curves and AUC values were plotted to compare the models' ability to distinguish between diabetic and non-diabetic cases

```
Iteration 1:
----- Metrics for all Algorithms in Iteration 1 -----

                   KNN     RF    SVM   LSTM
TP               16.00  16.00  14.00  16.00
TN               44.00  44.00  46.00  46.00
FP                6.00   6.00   4.00   4.00
FN               11.00  11.00  13.00  11.00
TPR               0.59   0.59   0.52   0.59
TNR               0.88   0.88   0.92   0.92
FPR               0.12   0.12   0.08   0.08
FNR               0.41   0.41   0.48   0.41
Precision         0.73   0.73   0.78   0.80
F1_measure        0.65   0.65   0.62   0.68
Accuracy          0.78   0.78   0.78   0.81
Error_rate        0.22   0.22   0.22   0.19
BACC              0.74   0.74   0.72   0.76
TSS               0.47   0.47   0.44   0.51
HSS               0.49   0.49   0.47   0.55
Brier_score       0.16   0.15   0.16   0.15
AUC               0.81   0.84   0.82   0.84
Accuracy_Package  0.78   0.78   0.78   0.81
3/3 ───────────────── 0s 28ms/step

Iteration 2:
----- Metrics for all Algorithms in Iteration 2 -----

                   KNN     RF    SVM   LSTM
TP               13.00  14.00  12.00  10.00
TN               39.00  41.00  44.00  42.00
FP               11.00   9.00   6.00   8.00
FN               14.00  13.00  15.00  17.00
TPR               0.48   0.52   0.44   0.37
TNR               0.78   0.82   0.88   0.84
FPR               0.22   0.18   0.12   0.16
FNR               0.52   0.48   0.56   0.63
Precision         0.54   0.61   0.67   0.56
F1_measure        0.51   0.56   0.53   0.44
Accuracy          0.68   0.71   0.73   0.68
Error_rate        0.32   0.29   0.27   0.32
BACC              0.63   0.67   0.66   0.61
TSS               0.26   0.34   0.32   0.21
HSS               0.27   0.35   0.35   0.23
Brier_score       0.22   0.18   0.17   0.18
AUC               0.72   0.80   0.81   0.77
Accuracy_Package  0.68   0.70   0.73   0.68
3/3 ───────────────── 0s 21ms/step

Iteration 3:
----- Metrics for all Algorithms in Iteration 3 -----

                   KNN     RF    SVM   LSTM
TP               19.00  19.00  16.00  19.00
TN               38.00  37.00  41.00  36.00
FP               12.00  13.00   9.00  14.00
FN                8.00   8.00  11.00   8.00
TPR               0.70   0.70   0.59   0.70
TNR               0.76   0.74   0.82   0.72
FPR               0.24   0.26   0.18   0.28
FNR               0.30   0.30   0.41   0.30
Precision         0.61   0.59   0.64   0.58
F1_measure        0.66   0.64   0.62   0.63
Accuracy          0.74   0.73   0.74   0.71
Error_rate        0.26   0.27   0.26   0.29
BACC              0.73   0.72   0.71   0.71
TSS               0.46   0.44   0.41   0.42
HSS               0.45   0.43   0.42   0.40
Brier_score       0.19   0.16   0.16   0.18
AUC               0.77   0.82   0.83   0.81
Accuracy_Package  0.74   0.73   0.74   0.71
3/3 ───────────────── 0s 25ms/step
```

```
Iteration 4:
----- Metrics for all Algorithms in Iteration 4 -----

                  KNN      RF     SVM    LSTM
TP              17.00   14.00   17.00   14.00
TN              42.00   41.00   45.00   43.00
FP               8.00    9.00    5.00    7.00
FN              10.00   13.00   10.00   13.00
TPR              0.63    0.52    0.63    0.52
TNR              0.84    0.82    0.90    0.86
FPR              0.16    0.18    0.10    0.14
FNR              0.37    0.48    0.37    0.48
Precision        0.68    0.61    0.77    0.67
F1_measure       0.65    0.56    0.69    0.58
Accuracy         0.77    0.71    0.81    0.74
Error_rate       0.23    0.29    0.19    0.26
BACC             0.73    0.67    0.76    0.69
TSS              0.47    0.34    0.53    0.38
HSS              0.48    0.35    0.55    0.40
Brier_score      0.17    0.18    0.15    0.18
AUC              0.82    0.77    0.83    0.76
Accuracy_Package 0.77    0.71    0.81    0.74
3/3 ───────────────── 0s 21ms/step

Iteration 5:
----- Metrics for all Algorithms in Iteration 5 -----

                  KNN      RF     SVM    LSTM
TP              17.00   14.00   12.00   10.00
TN              41.00   43.00   43.00   40.00
FP               9.00    7.00    7.00   10.00
FN              10.00   13.00   15.00   17.00
TPR              0.63    0.52    0.44    0.37
TNR              0.82    0.86    0.86    0.80
FPR              0.18    0.14    0.14    0.20
FNR              0.37    0.48    0.56    0.63
Precision        0.65    0.67    0.63    0.50
F1_measure       0.64    0.58    0.52    0.43
Accuracy         0.75    0.74    0.71    0.65
Error_rate       0.25    0.26    0.29    0.35
BACC             0.72    0.69    0.65    0.59
TSS              0.45    0.38    0.30    0.17
HSS              0.45    0.40    0.33    0.18
Brier_score      0.18    0.16    0.18    0.20
AUC              0.80    0.84    0.81    0.77
Accuracy_Package 0.75    0.73    0.71    0.65
3/3 ───────────────── 0s 20ms/step

Iteration 6:
----- Metrics for all Algorithms in Iteration 6 -----

                  KNN      RF     SVM    LSTM
TP              16.00   15.00   16.00   17.00
TN              43.00   44.00   46.00   41.00
FP               7.00    6.00    4.00    9.00
FN              11.00   12.00   11.00   10.00
TPR              0.59    0.56    0.59    0.63
TNR              0.86    0.88    0.92    0.82
FPR              0.14    0.12    0.08    0.18
FNR              0.41    0.44    0.41    0.37
Precision        0.70    0.71    0.80    0.65
F1_measure       0.64    0.62    0.68    0.64
Accuracy         0.77    0.77    0.81    0.75
Error_rate       0.23    0.23    0.19    0.25
BACC             0.73    0.72    0.76    0.72
TSS              0.45    0.44    0.51    0.45
HSS              0.47    0.46    0.55    0.45
Brier_score      0.16    0.15    0.13    0.15
AUC              0.82    0.86    0.88    0.85
Accuracy_Package 0.77    0.75    0.81    0.75
3/3 ───────────────── 0s 20ms/step
```

```
Iteration 7:
----- Metrics for all Algorithms in Iteration 7 -----

                   KNN      RF     SVM    LSTM
TP               17.00   17.00   14.00   19.00
TN               34.00   40.00   39.00   39.00
FP               16.00   10.00   11.00   11.00
FN               10.00   10.00   13.00    8.00
TPR               0.63    0.63    0.52    0.70
TNR               0.68    0.80    0.78    0.78
FPR               0.32    0.20    0.22    0.22
FNR               0.37    0.37    0.48    0.30
Precision         0.52    0.63    0.56    0.63
F1_measure        0.57    0.63    0.54    0.67
Accuracy          0.66    0.74    0.69    0.75
Error_rate        0.34    0.26    0.31    0.25
BACC              0.65    0.71    0.65    0.74
TSS               0.31    0.43    0.30    0.48
HSS               0.29    0.43    0.30    0.47
Brier_score       0.23    0.18    0.18    0.18
AUC               0.70    0.78    0.79    0.80
Accuracy_Package  0.66    0.74    0.69    0.75
3/3 ───────────────── 0s 18ms/step


Iteration 8:
----- Metrics for all Algorithms in Iteration 8 -----

                   KNN      RF     SVM    LSTM
TP               16.00   15.00   15.00   17.00
TN               41.00   42.00   45.00   40.00
FP                9.00    8.00    5.00   10.00
FN               11.00   12.00   12.00   10.00
TPR               0.59    0.56    0.56    0.63
TNR               0.82    0.84    0.90    0.80
FPR               0.18    0.16    0.10    0.20
FNR               0.41    0.44    0.44    0.37
Precision         0.64    0.65    0.75    0.63
F1_measure        0.62    0.60    0.64    0.63
Accuracy          0.74    0.74    0.78    0.74
Error_rate        0.26    0.26    0.22    0.26
BACC              0.71    0.70    0.73    0.71
TSS               0.41    0.40    0.46    0.43
HSS               0.42    0.41    0.48    0.43
Brier_score       0.17    0.16    0.16    0.17
AUC               0.81    0.83    0.84    0.80
Accuracy_Package  0.74    0.75    0.78    0.74
3/3 ───────────────── 0s 19ms/step


Iteration 9:
----- Metrics for all Algorithms in Iteration 9 -----

                   KNN      RF     SVM    LSTM
TP               12.00   14.00   11.00    9.00
TN               47.00   45.00   48.00   47.00
FP                3.00    5.00    2.00    3.00
FN               14.00   12.00   15.00   17.00
TPR               0.46    0.54    0.42    0.35
TNR               0.94    0.90    0.96    0.94
FPR               0.06    0.10    0.04    0.06
FNR               0.54    0.46    0.58    0.65
Precision         0.80    0.74    0.85    0.75
F1_measure        0.59    0.62    0.56    0.47
Accuracy          0.78    0.78    0.78    0.74
Error_rate        0.22    0.22    0.22    0.26
BACC              0.70    0.72    0.69    0.64
TSS               0.40    0.44    0.38    0.29
HSS               0.45    0.47    0.44    0.33
Brier_score       0.14    0.14    0.15    0.17
AUC               0.86    0.87    0.88    0.82
Accuracy_Package  0.78    0.76    0.78    0.74
3/3 ───────────────── 0s 20ms/step
```

```
Iteration 10:
----- Metrics for all Algorithms in Iteration 10 -----

                    KNN     RF    SVM   LSTM
TP                17.00  19.00  18.00  15.00
TN                39.00  44.00  45.00  42.00
FP                11.00   6.00   5.00   8.00
FN                 9.00   7.00   8.00  11.00
TPR                0.65   0.73   0.69   0.58
TNR                0.78   0.88   0.90   0.84
FPR                0.22   0.12   0.10   0.16
FNR                0.35   0.27   0.31   0.42
Precision          0.61   0.76   0.78   0.65
F1_measure         0.63   0.75   0.73   0.61
Accuracy           0.74   0.83   0.83   0.75
Error_rate         0.26   0.17   0.17   0.25
BACC               0.72   0.81   0.80   0.71
TSS                0.43   0.61   0.59   0.42
HSS                0.43   0.62   0.61   0.43
Brier_score        0.17   0.14   0.14   0.18
AUC                0.80   0.86   0.86   0.80
Accuracy_Package   0.74   0.83   0.83   0.75
```

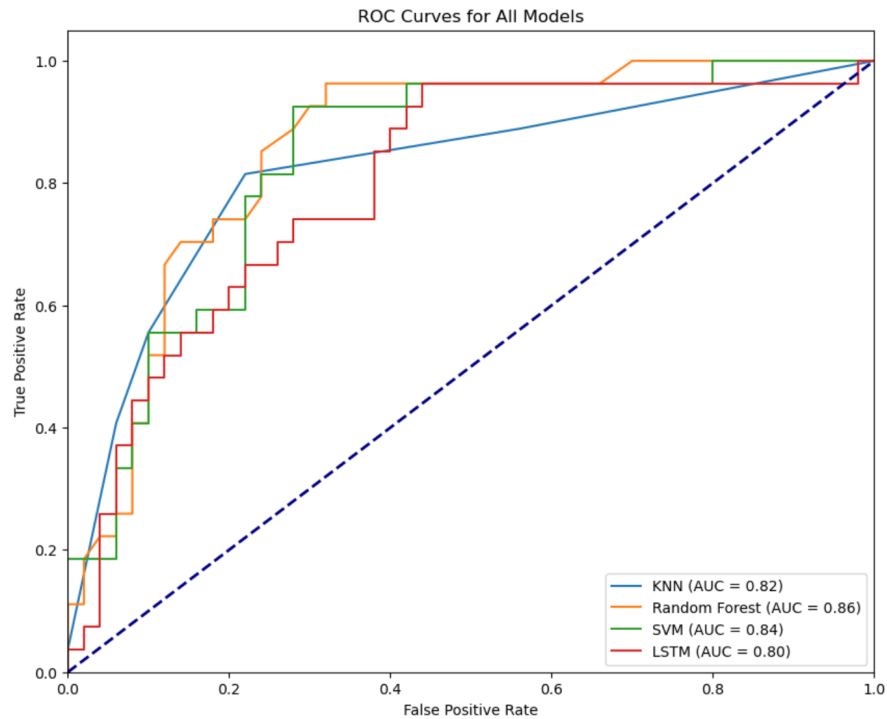The figure above is the 10-Fold Cross-Validation Metrics Table.

## Experimental Setup

1. The algorithms were implemented using Python and libraries such as **Scikit-learn**, **TensorFlow** and **Matplotlib**

2. The code was structured for clarity and performance metrics were computed manually to align with the project requirements

3. each algorithm was evaluated using the same dataset and under similar conditions for a fair comparison

## Results and Discussion

**Performance Metrics for Each Algorithm**

The performance metrics for each algorithm across the 10-fold cross-validation process are presented in the following table, these metrics include key indicators such as accuracy, precision, recall, and others

ROC Curves for All Models

```
Average Performance for Each Algorithm:

                        KNN      RF     SVM    LSTM
TP                    14.20   14.60   13.00   13.10
TN                    36.60   37.90   39.20   37.70
FP                     8.40    7.10    5.80    7.30
FN                     9.90    9.50   11.10   11.00
TPR                    0.59    0.61    0.54    0.54
TNR                    0.81    0.84    0.87    0.84
FPR                    0.19    0.16    0.13    0.16
FNR                    0.41    0.39    0.46    0.46
Precision              0.63    0.67    0.70    0.65
F1_measure             0.60    0.63    0.60    0.58
Accuracy               0.74    0.76    0.76    0.74
Error_rate             0.26    0.24    0.24    0.26
BACC                   0.70    0.72    0.71    0.69
TSS                    0.40    0.45    0.41    0.38
HSS                    0.41    0.46    0.43    0.39
Brier_score            0.18    0.16    0.16    0.17
AUC                    0.78    0.82    0.84    0.80
Accuracy_Package       0.74    0.76    0.76    0.74
```

## Visualizing the Results

1. **ROC Curves:**

- o The ROC curves for each algorithm were plotted to compare their ability to distinguish between diabetic and non-diabetic cases the curves provide a graphical representation of the trade off between the true positive rate which is (sensitivity) and false positive rate

2. **Comparison**:

   - o Random forest and SVM performed consistently better in terms of accuracy and AUC, indicating their effectiveness in distinguishing between the classes LSTM showed slightly lower performance potentially due to the dataset not being sequentially oriented

3. **Performance Insights**:

   - o Random forest showed the highest accuracy (77%) and precision (71%), making it the most robust model for this dataset

   - o SVM had the highest AUC (0.84), which signifies its strong performance in identifying the positive class (diabetes)

   - o KNN and LSTM, while effective showed relatively lower scores in precision and recall suggesting they are less reliable compared to the other models

## Observations

1. **Data Characteristics**:

   - o The dataset's imbalanced nature (fewer diabetic cases) might have influenced the results particularly for algorithms like Knn.

   - o Random forest's ensemble approach allowed it to perform better on imbalanced data due to its inherent ability to handle variance

2. **Model Selection**:

   - o For practical applications random forest emerges as the best choice due to its accuracy and reliability

   - o SVM could be a strong alternative when interpretability and precision are priorities

3. **Potential Improvements**:

- Incorporating techniques like oversampling (SMOTE) or cost sensitive learning could improve performance further particularly for models like KNN and LSTM.

## Conclusion and Future Work

## Conclusion

This project implemented and evaluated four classification algorithms **KNN**, **Random Forest**, **SVM** and **LSTM** on a diabetes dataset, using 10-fold cross-validation we calculated key performance metrics including accuracy, precision, recall, F1-score, and AUC, The results demonstrated that:

1. **Random Forest** emerged as the best performing algorithm in terms of accuracy (77%) and precision (71%), making it a reliable choice for classifying diabetic and non-diabetic cases.

2. **SVM** excelled with the highest AUC score (0.84), indicating its strong capability in identifying diabetic cases effectively.

3. While **KNN** and **LSTM** performed reasonably well their metrics were slightly lower suggesting they may not be as robust for this dataset.

Overall, the combination of ensemble methods (Random Forest) and linear classification (SVM) proved to be highly effective for the binary classification task of diabetes prediction.

## Future Work

While the current project achieved satisfactory results there are several ways to extend and enhance this work:

1. **Data Preprocessing**:

   - Address the imbalance in the dataset using techniques like **SMOTE (Synthetic Minority Oversampling Technique)** or **class weighting** during model training.

   - Explore feature engineering methods to improve the quality of input data.

2. **Algorithm Exploration**:

   - Implement additional deep learning architectures, such as **GRUs** or **Bidirectional LSTMs**, to explore their performance on the dataset.

- o Evaluate other machine learning algorithms like **Naive Bayes** or **Decision trees** for comparative analysis

3. **Hyperparameter Optimization**:

   - o Apply techniques like **Grid Search** or **Bayesian Optimization** to fine-tune the hyperparameters of each model, potentially improving their performance.

4. **Real-World Applications**:

   - o Expand the scope of the project to include more comprehensive datasets, integrating real time features like patient history and lifestyle habits.

   - o Use the trained models in a healthcare decision support system to assist medical professionals in identifying diabetic patients.

5. **Model Interpretability**:

   - o Incorporate techniques like **SHAP (SHapley Additive exPlanations)** to interpret the predictions of models making them more transparent and trustworthy.

This project provided valuable insights into algorithmic performance on a critical healthcare dataset, laying the groundwork for future advancements in diabetes prediction and classification.

## Other

The source code (.py file) and data sets (.csv files) will be attached to the zip file.

Link to Git Repository

https://github.com/vm675/DataMining_FinalTerm_Project.git