



Faculty of Engineering and Technology
Electrical and Computer Engineering Department
First Semester 2024/2025
ENCS3340
Artificial Intelligence

**Predictive Medical Diagnosis of Diabetes Based on
Patient Symptoms Using Machine Learning**

Prepared by:

Fatima Dawabsheh 1210827

Shahed Shrateh 1210444

Raghd Isleem 1211326

Instructor: Dr. Adnan Yahya

Section: 2

Date: 2/1/ 2025

Abstract— This project develops a machine learning-based system for diabetes classification and diagnosis, utilizing Python in Virtual Studio and Weka applications. We employ diverse clinical data, perform preprocessing and feature extraction, and train classification models using algorithms like decision trees, neural networks, and Naive Bayes. By evaluating performance metrics (accuracy, precision, recall, F1-score), we demonstrate the system's effectiveness in predicting diabetes and classifying patients. This implementation provides healthcare professionals with an efficient tool for early detection and personalized treatment planning.

I. INTRODUCTION

• Machine Learning Overview

Machine learning also known as source, is a branch of artificial intelligence AI that deals with the possibility of the computer to learn. ML algorithms depends on the target output, and the input used for training.

• Model Development

This project also incorporates and discusses several machine learning models such as decision trees, random forests, Naive Bayes classifiers and even neural networks. The two development tools that are both used in development are Weka software and Python.

• Algorithm Description

- **Decision Tree Classifier:** Used to establish a decision tree which is preferable for classification tasks and represents choices and their results.
- **Random Forest Classifier:** The use of multiple decision trees in a way that optimizes a result function to obtain improved regression and classification predictions with lower variance.
- **Naive Bayes Classifier:** A derivative of Bayes' Theorem Probability Model that largely applies in classification challenges.
- **Neural Network Classifier:** A kind of a subcategory of the Deep Learning approach that attempts to mimic the functionality of the human brain. Due to the ability of learning complex patterns in the data, these models are appropriate for classification.

Evaluating the performance of models involves the use of accuracy, precision, recall and the F1-score.

II. PROCEDURE

A. Datasets

We obtained two different data sets from kaggle which are given below. The first dataset had 2768 entries as processed [1].and the second had 4304 entries [2]. each dataset is saved as .csv file. Initial data cleaning was performed, which included removing the ID column from the first dataset. we used some common classifiers for use in classification problems such as decision trees, random forests, naive Bayes, and neural networks. We had ensured that our data was clean and was fit for use in the machine learning models before training the machine. The finally selected algorithms were applied to build classification models to classify the existence of diabetes. We implemented this by using Weak and by using the scikit- learn Python library, which has optimized machine learning algorithms.

III. . PERFORMANCE METRICS:

We used the following metrics to assess the models performance:

- Precision:** As a measure of the model's ability to predict the positive cases correctly, defined as the total number of true positives divided by the total of all the positive predictions of the model.
- Recall:** Able to measure the totality of the True Positives discovered by a model as a fraction of the total count of actual true positives.
- Accuracy:** To summarize, the overall correctness of our predictions, defined as the ratio of the numbers of properly classified data points to the total number of data points within the test set.
- F1 Score:** Known as the harmonic mean of precision and recall, providing a balanced metric that is especially useful when classes are imbalanced.

• TESTING AND RESULTS

1. Weka program

Initially, we evaluated dataset one using Weka, the results:-

A. Decision tree

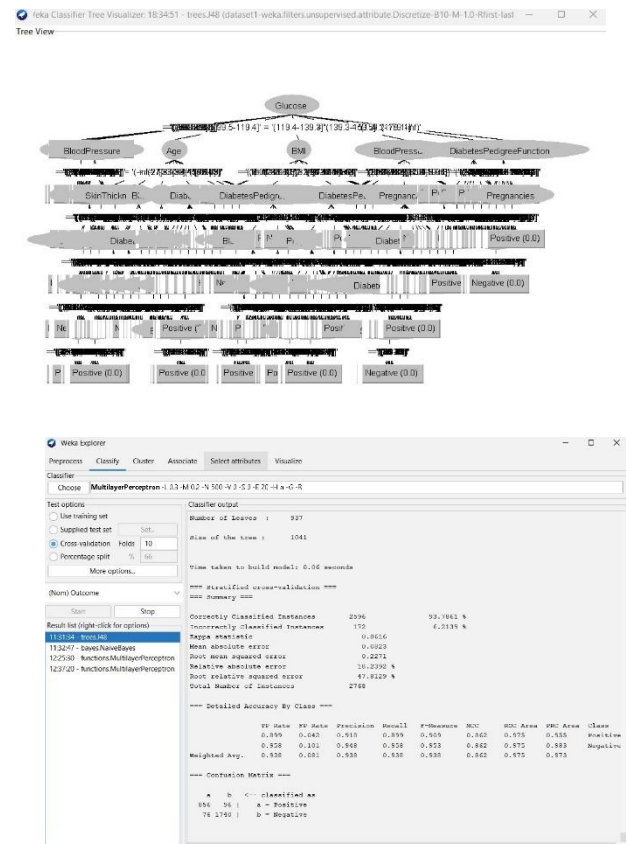


Figure 1:- Decision tree for dataset 1

The decision tree model's results showed a correctly classified instance rate of 93.78%, with a precision, recall, and F-measure equal to 0.938. The confusion matrix resulted in 856 true positives (TP), 96 false negatives (FN), 76 false positives (FP), and 1740 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9378.

B. Naïve Bayes

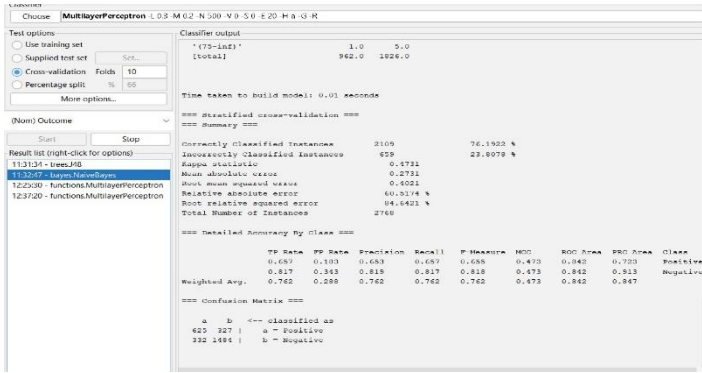


Figure 2:- Naïve Bayes for dataset 1

The Naïve Bayes model's results showed a correctly classified instance rate of 76.192%, with a precision, recall, and F-measure equal to 0.762. The confusion matrix resulted in 625 true positives (TP), 327 false negatives (FN), 332 false positives (FP), and 1484 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.7619.

C. Neural network

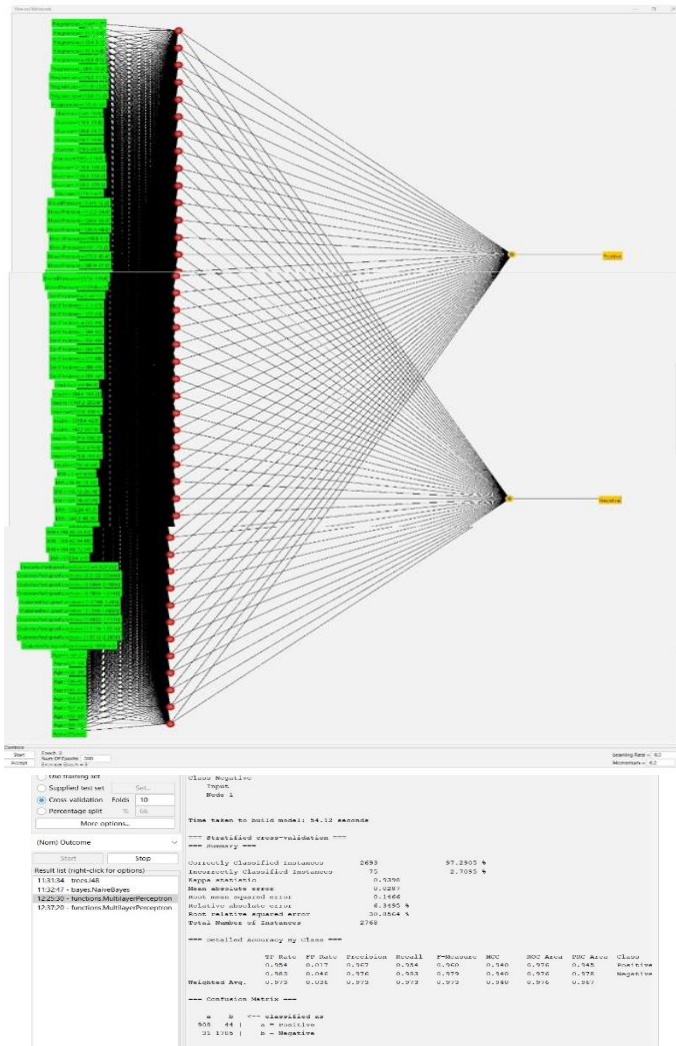


Figure 4:- Neural network for dataset 1

The Neural network model's results showed a correctly classified instance rate of 97.29%, with a precision, recall, and F-measure equal to 0.973. The confusion matrix resulted in 908 true positives (TP), 44 false negatives (FN), 31 false positives (FP), and 1785 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9729.

D. random forest

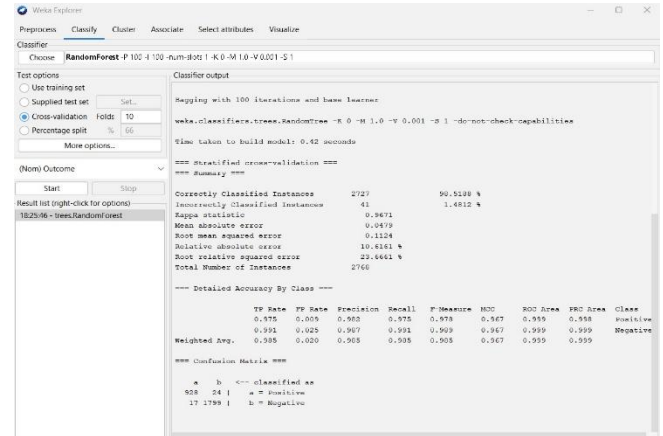


Figure 5:- random forest for dataset 1

The random forest model's results showed a correctly classified instance rate of 98.51% with a precision, recall, and F-measure equal to 0.985. The confusion matrix resulted in 928 true positives (TP), 24 false negatives (FN), 17 false positives (FP), and 1799 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9851.

Dataset 2:-

A. Decision tree

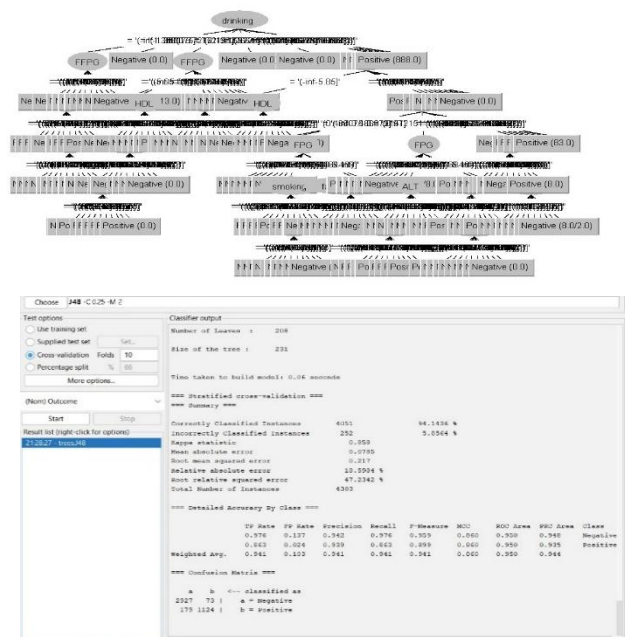


Figure 3:- decision tree for dataset 2

B. Naïve Bayes

Classifier: MultinomialPerceptron

Test options: ☒ Use training set

Classifier output	
0 (0.5 log%)	172 / 58
1 (0.5 log%)	172 / 58

[illegible]

The *Naïve Bayes* model's results showed a correctly classified instance rate of 94.39%, with a precision, recall, and F-measure equal to 0.944. The confusion matrix resulted in 2886 true positives (TP), 114 false negatives (FN), 127 false positives (FP), and 1176 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9439.

Channel

Test options

☐ Use training set
☐ Copied test set
☒ Cross-validation folds: 10
 Percentage split: 0.8

Get...

More options...

Classifier output

Class Positive

Input

Model 1

Time taken to build model: 466.00 seconds

--- Detailed cross-validation summary ---

Incorrectly classified instances

2404

54.815%

Incorrectly classified instances

240

5.762%

Range statistics

0.8005

Mean absolute error

0.399

Mean cross-entropy error

0.577

Relative absolute error

13.737%

Mean squared error

0.7041

Total number of instances

5922

--- Detailed Accuracy By Class ---

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.970	0.022	0.948	0.970	0.950	0.622	0.973	0.904	Negative
	0.078	0.020	0.327	0.078	0.102	0.072	0.073	0.070	Positive

--- Confusion Matrix ---

	A	B	Classified as
2510	80	0	Negative
3391	1	0	Positive

The Neural network model's results showed a correctly classified instance rate of 94.213%, with a precision, recall, and F-measure equal to 0.942. The confusion matrix resulted in 2910 true positives (TP), 90 false negatives (FN), 159 false positives (FP), and 1144 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9421.

[illegible]

The *random forest* model's results showed a correctly classified instance rate of 90.750%, with a precision, recall, and F-measure equal to 0.907. The confusion matrix resulted in 2829 true positives (TP), 171 false negatives (FN), 227 false positives (FP), and 1076 true negatives (TN). Based on these numbers, the overall accuracy, calculated as $(TP+TN)/(TP+FP+TN+FN)$, was found to be 0.9075.

2. A Web Interface for Diabetes Diagnosis

Confusion Matrix

	Actual Positive	Actual Negative
Classified Positive	100	7
Classified Negative	4	363

True Positives: 100 Accuracy: 0.9816 Precision: 0.9769 Recall: 0.9825 F1 Score: 0.9793

Do you want to do a test ?

Yes No

Diabetes Test Data Input

Pregnancies

Glucose

BloodPressure

SkinThickness

Insulin

BMI

DiabetesPedigreeFunction

Age

Sex

Submit

WORKING POINT
Introduction to Data Science (10/2020)

Figure 8:- Page 4

3. Python code

```
1 # all packages needed
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.neural_network import MLPClassifier
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix
8 import time
9 import pandas as pd
10
11 # load datasets into DataFrames
12 Dataset1 = pd.read_csv('dataset1.csv', delimiter=',')
13 Dataset2 = pd.read_csv('dataset2.csv', delimiter=',')
14
15 # target column
16 t_columns = {
17     '1': 'Outcome',
18     '2': 'Diabetes'
19 }
20
21 print("Dataset1 : 2768")
22 print("Dataset2 : 4303 , This dataset has more number of features")
23 choice = input("Enter the dataset you want to test 1 or 2 : ")
24 if choice == '1':
25     data = Dataset1
26 elif choice == '2':
27     data = Dataset2
28 else:
29     print("Invalid dataset")
30     exit()
```

Figure 13:- code 1

We import the necessary libraries, such as sklearn for machine learning models and metrics, and pandas for handling datasets. And to make the data ready for analysis we import the two files as dataset 1 and dataset 2 from local CSV files using the method `pd.read_csv()`.

```
32 # Verify the target column exists
33 target_column = t_columns.get(choice)
34 if target_column not in data.columns:
35     print(f"Error: The '{target_column}' column is missing in the selected dataset")
36     print("Available columns:", data.columns)
37     exit()
38
39 # Separate features (X) and labels (y)
40 X = data.drop(target_column, axis=1)
41 y = data[target_column]
42 # Split the dataset into 80% training and 20% testing sets
43 Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.2, random_state=42)
44
```

Figure 14:- code 2

The target column name is obtained by using the choice of the user as the key for accessing the value in the dictionary, `t_columns`. If the given target column is not in the selected dataset, we display which columns are available and then the program ceases operation.

```
45 while True:
46     print("\n")
47     print("Choose the algorithm :")
48     print("1. Decision Tree")
49     print("2. Naive Bayes")
50     print("3. Neural Network")
51     print("4. Random Forest")
52     print("\n> Exit")
53     choice = input("Enter your choice 1, 2, 3, 4, or 5: ")
54
55     if choice == '1':
56         model = DecisionTreeClassifier()
57     elif choice == '2':
58         model = GaussianNB()
59     elif choice == '3':
60         model = MLPClassifier(random_state=1, max_iter=300)
61     elif choice == '4':
62         model = RandomForestClassifier(n_estimators=100, random_state=42)
63     elif choice == '5':
64         break
65     else:
66         print("Invalid choice")
67         continue
```

Figure 15:- code 3

A while True loop ensures the user can repeatedly choose an algorithm until they decide to exit and Based on the user's input, the appropriate machine learning model is instantiated

```
68 # Train the model
69 model.fit(Xtrain, ytrain)
70
71 # Start time
72 start_time = time.time()
73
74 # Make predictions on the test set
75 ypred = model.predict(Xtest)
76
77 # calculate confusion matrix
78 Conf_Matrix = confusion_matrix(ytest, ypred)
79
80 # calculate evaluation metrics
81 Accuracy = accuracy_score(ytest, ypred)
82 precision = precision_score(ytest, ypred)
83 Recall = recall_score(ytest, ypred)
84 f1 = f1_score(ytest, ypred)
85
86 # Round to 5 decimal places
87 Accuracy = round(Accuracy,5)
88 precision = round(precision,5)
89 Recall = round(Recall,5)
90 f1 = round(f1,5)
91 Time_elapsed = round(time.time()-start_time,5)
92
93 # print algorithm name
94 if choice == '1':
95     print("Decision Tree Algorithm")
96 elif choice == '2':
97     print("Naive Bayes Algorithm")
98 elif choice == '3':
99     print("Neural Network Algorithm")
100 elif choice == '4':
101     print("Random Forest Algorithm")
102
```

Figure 16:- code 4

We trained the previously selected machine learning model using the training data (features in `Xtrain` and labels in `ytrain`). This was where the model learned the patterns in the data. And we recorded the current time before the prediction process started, used for measuring the time it takes to make a prediction. And we used the trained model to predict the labels (`ypred`) for the test dataset's features (`Xtest`).

```
104 # Print confusion matrix
105 print("Confusion Matrix:")
106 print("\n\t\tActual Positive\t\tActual Negative")
107 print("Classified Positive\t\t", Conf_Matrix[1, 1], "\t\t", Conf_Matrix[1, 0])
108 print("Classified Negative\t\t", Conf_Matrix[0, 1], "\t\t", Conf_Matrix[0, 0])
109
110 # Print evaluation metrics
111 print("\n")
112 print("Time elapsed :", Time_elapsed, "seconds")
113 print("Accuracy :", Accuracy)
114 print("Precision :", precision)
115 print("Recall :", Recall)
116 print("F1-score :", f1)
117 print("\n")
118
119 # Test for user
120 value = input("Do you want to test values from your own of this algorithm (y/n)? ")
121 if value.lower() == 'y':
122     Features = {}
123     for column in X.columns:
124         value = input("Enter value for [column]: ")
125         Features[column]=value
126
127 # Create a DataFrame
128 user_data = pd.DataFrame(Features)
129
130 # Make predictions
131 user_pred = model.predict(user_data)
132
133 # prediction
134 if user_pred[0] == 1:
135     print("The expected result is to have diabetes")
136 else:
137     print("The expected result is not to have diabetes")
138 else:
139     continue
```

Figure 17:- code 5

We ask the user if they want to input new values. And if the user choice yes, we prompt them to enter values for each column in the data. Then we create a new DataFrame using the values provided by the user. Then we use the trained model to predict the outcome based on the user input. And print whether the expected result indicates diabetes or not.

❖ Here is an example of what the output looks like when we run this code

```
Dataset1 : 2768
Dataset2 : 4303 , This dataset has more number of features
Enter the dataset you want to test 1 or 2 : 1

Choose the algorithm :
1.Decision Tree
2.Naive Bayes
3.Neural Network
4.Random Forest
5.Exit
Enter your choice 1, 2, 3, 4, or 5: 1

Decision Tree Algorithm
Confusion Matrix:
               Actual Positive      Actual Negative
Classified positive              188              7
Classified Negative              4              363

Time elapsed : 0.0 seconds
Accuracy : 0.98814
Precision : 0.97626
Recall : 0.98257
F1-score : 0.97935

Do you want to test values from your own of this algorithm (y/n)? y
Enter value for Pregnancies: 6
Enter value for Glucose: 148
Enter value for BloodPressure: 72
Enter value for SkinThickness: 35
Enter value for Insulin: 0
Enter value for BMI: 31.6
Enter value for DiabeticPedigreeFunction: 0.627
Enter value for Age: 50
The expected result is to have diabetes
```

Figure 18:- result 1

REFERENCES

- [1] <https://www.kaggle.com/datasets/salihacur/diabetes>
- [2] <https://www.kaggle.com/datasets/pkdarabi/diabetes-dataset-with-18-features>