



FACULTY OF ENGINEERING AND TECHNOLOGY
DEPARTMENT OF COMPUTER ENGINEERING
Artificial Intelligence - ENCS434

Prepared by:

Leyan Burait #1211439

Haneen Odeh #1210716

Instructor: Dr. Adnan Yahya

Date: 26/1/2024

Section: 1

Abstract:

Diabetes is a complex chronic disease that requires accurate classification and prompt diagnosis for effective management. This project aims to utilize machine learning techniques and the Weka program and python language to develop a robust classification and diagnostic system for diabetes. By employing a diverse dataset containing clinical and laboratory features, we preprocess and transform the data to extract meaningful information. We then utilize various machine learning algorithms available in the Weka program, including decision trees, and neural networks, to train and optimize our classification model. And decision tree, naïve bias, neural network and random forest using Python language. By evaluating the model's performance metrics, including accuracy, precision, recall, and F1-score, we demonstrate the effectiveness of our system in predicting diabetes onset and categorizing patients into relevant diabetes types. The implementation of this system using the Weka program (Also python language) provides a user-friendly interface and efficient tool for healthcare professionals to aid in early diabetes detection and personalized treatment strategies.

Table of Contents

1. INTRODUCTION	4
1.1 Machine learning.....	4
1.2 Models Training	4
1.3 Data set.....	5
2. Methodology	5
3.TESTING AND RESULTS	6
3.1 Weka program	6
3.2 Python Language	10

1. INTRODUCTION

1.1 Machine learning

Machine learning is a branch of artificial intelligence that concerns how the computer can learn and adapt to new circumstances. There are many learning techniques based on the desired outcome from the technique and the input available at the training process.

1.2 Models Training

In this project, various machine learning techniques, including decision trees, random forest, and neural networks, are utilized to train and optimize a classification model by both the Weka program and Python language.

1. A Decision Tree: This algorithm is employed as a powerful tool for constructing a tree-like model of decisions and their possible consequences. It allows for the representation of rules and relationships between variables in a structured manner. This algorithm is particularly useful for classification tasks where the data can be split into distinct categories.

2. Random Forest: The Random Forest algorithm is a versatile and widely used machine learning technique that operates by constructing multiple decision trees during training and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees.

3. A Naive Bayes classifier: It is a probabilistic machine learning model that's used for classification tasks. It is based on the Bayes theorem.

4. A Neural Network: A type of deep learning algorithm, are also employed in this project. Neural networks are inspired by the structure and function of the human brain and consist of interconnected nodes or "neurons." They are capable of learning complex patterns and relationships in the data, making them suitable for classification tasks.

The classification model is trained and optimized using these techniques, taking into account the specific requirements and characteristics of the dataset. The performance of the model is evaluated using different metrics to evaluate its accuracy and effectiveness in classifying new data points.

By training and testing datasets using decision trees, random forest, neural networks, and naïve Bayes we will calculate accuracy, precision, recall & F1 score.

1.3 Data set

This project was implemented using 1 dataset for Weka program the link below to obtain data set for patients, some of whom have diabetes and some of whom do not have it.

Link:

- 1) Pore, N. (2023, August 23). *Healthcare diabetes dataset*. Kaggle.
<https://www.kaggle.com/datasets/nanditapore/healthcare-diabetes>

For Python language we used 3 different datasets

- 1) Pore, N. (2023, August 23). *Healthcare diabetes dataset*. Kaggle.
<https://www.kaggle.com/datasets/nanditapore/healthcare-diabetes>
- 2) ACUR, S. (2020a, March 21). *Diabetes*. Kaggle.
<https://www.kaggle.com/datasets/salihacur/diabetes>
- 3) Darabi, P. K. (2023b, July 28). *Diabetes_dataset_with_18_features*. Kaggle.
<https://www.kaggle.com/datasets/pkdarabi/diabetes-dataset-with-18-features>

2. Methodology

In this section, we will briefly summarize all the steps we took to complete our project to detect diabetes We followed the following steps:

1. We first read the raw data and performed all the preprocessing
2. All pre-processed data was rewritten with the extracted data into a CSV file.
3. Next, we decided which model we would practice on, and made sure we had a balanced approach data and then trained the model.
4. We stored all trained models for future use and new test samples.

3.TESTING AND RESULTS

3.1 Weka program

The dataset we used contains 2768 entries, we divided it 80% for training data, and 20% for testing data.

Before going in depth with the different outputs, we need to consider the following metrics:

1. Precision: is the number of correctly classified positive examples divided by the total number of examples that are classified as positive.
2. Recall: is the number of correctly classified positive examples divided by the total number of actual positive examples in the test set.
3. Accuracy: is the number of correct classifications divided by the total number of test cases.
4. F1 score: is the harmonic mean of precision and recall, and it tends to closer the smaller of the two

The diagram illustrates a decision tree for predicting gestational diabetes. The root node is 'Glucose' with a split at 127. The left branch leads to 'BMI' (split at 26.4), which then splits into 'Pregnancies' (0, 2) and 'age' (split at 28). The right branch leads to 'BMI' (split at 29.9), which splits into 'age' (split at 29.9) and 'Pressure' (split at 157). The tree continues with numerous internal nodes and leaf nodes containing numerical values and variable names, representing a highly complex model.

TP=181 , FP= 10 , FN=2 , TN = 361

Precision= $TP/TP+FP = 181/181+10 = 0.9476$

Recall = $TP/TP+FN=181/181+2 = 0.9891$

F1 score= $2P*R/P+R = 0.96791$

3.2 Naive bais

1. traning

Total Number of Instances: 2214

Correctly classified instances: 1691

In correctly classified instances 523

2. test

Total Number of Instances: 554

Correctly classified instances: 422

In correctly classified instances 132

Accuracy= correctly classified Instances/ Total number of instances in data set

Accuracy= $422/554 = 0.7617$

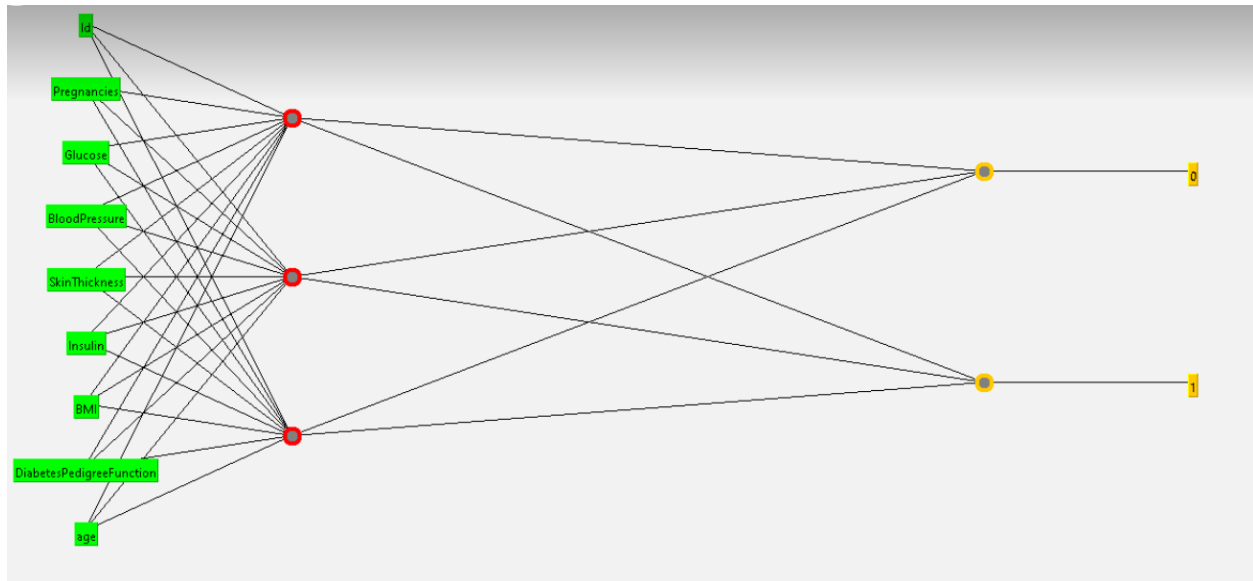
TP=117 , FP= 66 , FN=66 , TN = 305

Precision= $TP/TP+FP = 117/554= 0.6393$

Recall = $TP/TP+FN=171/171+66 = 0.6393$

F1 score= $2P*R/P+R = 0.6395$

3.3 A Neural Network:



1.training

Total Number of Instances: 2214

Correctly classified instances: 1825

In correctly classified instances 389

2.testing

Total Number of Instances: 554

Correctly classified instances: 455

In correctly classified instances 99

Accuracy= correctly classified Instances/ Total number of instances in data set

Accuracy=455/554 =0.8213

TP=128 , FP= 44 , FN=55 , TN = 327

Precision= TP/TP+FP = 128/128+44 = 0.7442

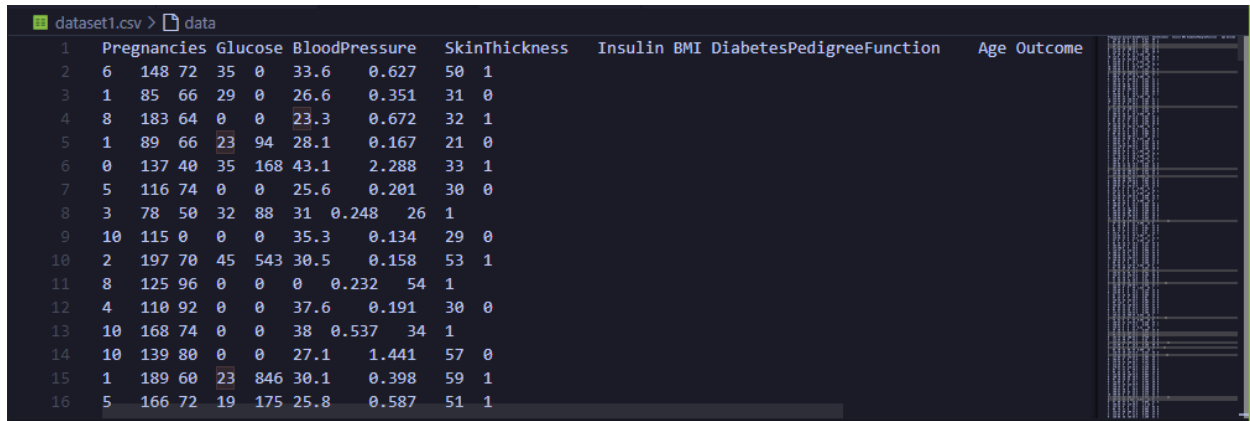
Recall = TP/TP+FN=128/128+55 = 0.6995

F1 score= 2P*R/P+R = 0.7211

3.2 Python Language

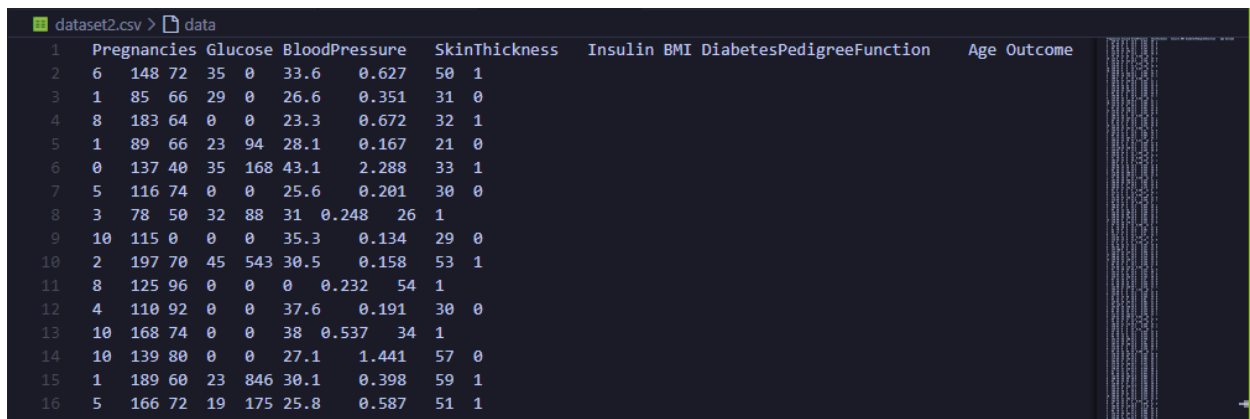
We used 3 different datasets with different number of entries to see the difference,

A short picture for Dataset1 as .csv



	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
1	6	148	72	35	0	33.6	0.627	50	1
2	1	85	66	29	0	26.6	0.351	31	0
3	8	183	64	0	0	23.3	0.672	32	1
4	1	89	66	23	94	28.1	0.167	21	0
5	0	137	40	35	168	43.1	2.288	33	1
6	5	116	74	0	0	25.6	0.201	30	0
7	3	78	50	32	88	31	0.248	26	1
8	10	115	0	0	0	35.3	0.134	29	0
9	2	197	70	45	543	30.5	0.158	53	1
10	8	125	96	0	0	0	0.232	54	1
11	4	110	92	0	0	37.6	0.191	30	0
12	10	168	74	0	0	38	0.537	34	1
13	10	139	80	0	0	27.1	1.441	57	0
14	1	189	60	23	846	30.1	0.398	59	1
15	5	166	72	19	175	25.8	0.587	51	1

A short picture for Dataset2 as .csv



	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
1	6	148	72	35	0	33.6	0.627	50	1
2	1	85	66	29	0	26.6	0.351	31	0
3	8	183	64	0	0	23.3	0.672	32	1
4	1	89	66	23	94	28.1	0.167	21	0
5	0	137	40	35	168	43.1	2.288	33	1
6	5	116	74	0	0	25.6	0.201	30	0
7	3	78	50	32	88	31	0.248	26	1
8	10	115	0	0	0	35.3	0.134	29	0
9	2	197	70	45	543	30.5	0.158	53	1
10	8	125	96	0	0	0	0.232	54	1
11	4	110	92	0	0	37.6	0.191	30	0
12	10	168	74	0	0	38	0.537	34	1
13	10	139	80	0	0	27.1	1.441	57	0
14	1	189	60	23	846	30.1	0.398	59	1
15	5	166	72	19	175	25.8	0.587	51	1

A short picture for Dataset3 as .csv

1	Age	Gender	BMI	SBP	DBP	FPG	Chol	Tri HDL	LDL	ALT	BUN	CCR	FFPG	smoking	drinking	family_histroy		
2	26	1	20.1	119	81	5.8	4.36	0.86	0.9	2.43	12	5.4	63.8	5.4	3	3	0	0
3	40	1	17.7	97	54	4.6	3.7	1.02	1.5	2.04	9.2	3.7	70.3	4.1	1	1	0	0
4	40	2	19.7	85	53	5.3	5.87	1.29	1.75	3.37	10.1	4.1	61.1	4.85	3	3	0	0
5	43	1	23.1	111	71	4.5	4.05	0.74	1.27	2.6	36.5	4.38	73.4	5.3	2	3	0	0
6	36	1	26.5	130	82	5.54	6.69	3.49	0.91	3.64	69.3	3.86	67.5	5.53	3	3	0	0
7	46	2	20.5	88	63	5.76	4.6	1	1.32	2.78	15	4.19	59	4.8	3	3	0	0
8	52	1	31.7	129	84	5.9	6.14	2.18	1.15	3.43	26	4.7	79	5.48	1	3	0	0
9	33	1	22.9	129	92	5.17	6.02	3.9	1.09	3.12	39.6	4.48	68.3	5.84	3	3	0	0
10	42	1	27.1	109	56	5.06	4.73	1.02	1.15	2.82	11.5	2.98	80.2	5.2	3	3	0	0
11	37	2	20.08	128	75	4.67	6.75	0.61	2.4	4.25	8.3	5.03	62.7	5.19	3	3	0	0
12	51	1	22	109	84	4.61	4.95	2.4	1.07	2.95	15.5	3.67	79.4	4.68	3	3	0	0
13	38	2	20.3	95	58	4.96	3.95	1.24	0.96	2.4	10.4	3	50.8	4.99	3	3	0	0
14	50	1	21.4	103	77	4.8	4.58	1.36	1.53	2.31	21.9	5.2	80	5.35	1	2	0	0
15	53	1	21.7	130	86	5.38	4.4	3.47	1.17	2.25	25.8	5.7	71	5.52	2	3	0	0
16	47	1	20.76	100	79	4.69	4.49	0.51	1.52	2.66	22.1	5.93	95.6	4.92	3	3	0	0

We had to install many packages to perform all the algorithms we used, such as

```
main2.py > ...
1  #!/ Installing 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
8  from sklearn.metrics import confusion_matrix
9  import pandas as pd
10 import time
11
```

At first, we had to make the program read all 3 datasets and let the user decide which dataset to use

```
11
12 #!/ Load our datasets into DataFrames
13 data1 = pd.read_csv('dataset1.csv', delimiter='\t')
14 data2 = pd.read_csv('dataset2.csv', delimiter='\t')
15 data3 = pd.read_csv('dataset3.csv', delimiter='\t')
16
17 #!/ Letting the user decide which dataset he/she wants
18 print("Last ID for dataset1:2768")
19 print("Last ID for dataset2:768")
20 print("Last ID for dataset3:4303, Note:(This dataset has more number of features)")
21
```

Based on the user's dataset, we divided it 80% training set and 20% testing set

```
34 #!/ Separate features (X) and Labels (y), X for all features except Outcome, y for Outcome feature
35 X = data.drop('Outcome', axis=1)
36 y = data['Outcome']
37
38 #!/ Split the dataset into training and testing sets (80% training and 20% testing)
39 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
40
```

We let the user decide which algorithm he wants to perform by showing this menu

```
40
41 while True:
42     #! Ask the user to choose the algorithm
43     print("\n\t")
44     print("Choose the algorithm you want to test:")
45     print("1. Decision Tree")
46     print("2. Naive Bayes")
47     print("3. Neural Network")
48     print("4. Random Forest")
49     print("5. Exit")
50     choice = input("Enter your choice (1, 2, 3, 4, or 5): ")
51
```

We had to initialize each algorithm so that we can perform and run it later

```
51
52     if choice == '1':
53         #! Initialize the decision tree model
54         model = DecisionTreeClassifier()
55     elif choice == '2':
56         #! Initialize the Naive Bayes model
57         model = GaussianNB()
58     elif choice == '3':
59         #! Initialize the Neural Network model
60         model = MLPClassifier(random_state=1, max_iter=300)
61     elif choice == '4':
62         #! Initialize the Random Forest model
63         model = RandomForestClassifier(n_estimators=100, random_state=42)
64     elif choice == '5':
65         break
```

After that we trained our model, made predictions on the test set, and calculated the confusion matrix

```
70     #! Train the model
71     model.fit(X_train, y_train)
72
73     #! Start time
74     start_time = time.time()
75
76     #! Make predictions on the test set
77     y_pred = model.predict(X_test)
78
79     #! Calculate confusion matrix
80     conf_matrix = confusion_matrix(y_test, y_pred)
81
```

After that, we calculated the evaluation metrics (Precision, Accuracy, Recall and F1-score), and rounded each one to 5 decimal places

```
82      #! Calculate evaluation metrics
83      accuracy = accuracy_score(y_test, y_pred)
84      precision = precision_score(y_test, y_pred)
85      recall = recall_score(y_test, y_pred)
86      f1 = f1_score(y_test, y_pred)
87
88      #! Round the evaluation metrics to 5 decimal places
89      accuracy = round(accuracy, 5)
90      precision = round(precision, 5)
91      recall = round(recall, 5)
92      f1 = round(f1, 5)
93
```

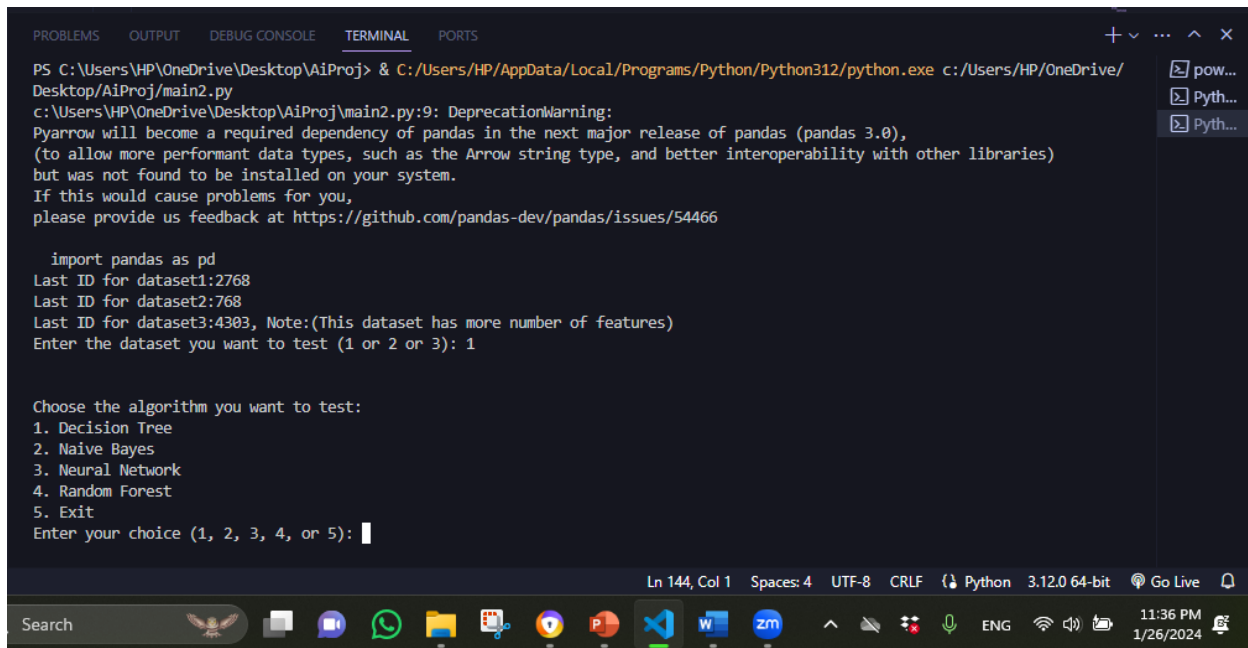
And finally, printing all the results to the screen

```
107      #! Print confusion matrix
108      print("Confusion Matrix:")
109      print("\t\t\tActual Positive\t\t\tActual Negative")
110      print("Classified Positive\t\t", conf_matrix[1, 1], "\t\t\t", conf_matrix[1, 0])
111      print("Classified Negative\t\t", conf_matrix[0, 1], "\t\t\t", conf_matrix[0, 0])
112
113      #! Print evaluation metrics
114      print("\n")
115      print("Time elapsed:", time_elapsed, "seconds")
116      print("Accuracy:", accuracy)
117      print("Precision:", precision)
118      print("Recall:", recall)
119      print("F1-score:", f1)
120      print("\n")
121
```

The final step was to let the user enter his own values to see what will the output be for the diabetes prediction

```
122     #! Ask the user to enter values for each feature
123     value1 = input(f"Do you want to test values from your own of this algorithm (y/n)")
124     if value1.lower() == 'y':
125         print("\nEnter values for each feature:")
126         features = {}
127         for column in X.columns:
128             value = input(f"Enter value for {column}: ")
129             features[column] = [value]
130
131     #! Create a DataFrame with the user-provided values
132     user_data = pd.DataFrame(features)
133
134     #! Make predictions using the trained model
135     user_pred = model.predict(user_data)
136
```

Some pictures for the results:-



```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
PS C:\Users\HP\OneDrive\Desktop\AiProj> & C:/Users/HP/AppData/Local/Programs/Python/Python312/python.exe c:/Users/HP/OneDrive/
Desktop/AiProj/main2.py
c:\Users\HP\OneDrive\Desktop\AiProj\main2.py:9: DeprecationWarning:
Pyarrow will become a required dependency of pandas in the next major release of pandas (pandas 3.0),
(to allow more performant data types, such as the Arrow string type, and better interoperability with other libraries)
but was not found to be installed on your system.
If this would cause problems for you,
please provide us feedback at https://github.com/pandas-dev/pandas/issues/54466

import pandas as pd
Last ID for dataset1:2768
Last ID for dataset2:768
Last ID for dataset3:4303, Note:(This dataset has more number of features)
Enter the dataset you want to test (1 or 2 or 3): 1

Choose the algorithm you want to test:
1. Decision Tree
2. Naive Bayes
3. Neural Network
4. Random Forest
5. Exit
Enter your choice (1, 2, 3, 4, or 5):
```

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
Choose the algorithm you want to test:
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:
Classified Positive    Actual Positive    Actual Negative
Classified Negative    180              7
                      4              363

Time elapsed: 0.0175 seconds
Accuracy: 0.98014
Precision: 0.97826
Recall: 0.96257
F1-score: 0.97035

Do you want to test values from your own of this algorithm (y/n)
```

Ln 144, Col 1 Spaces: 4 UTF-8 CRLF Python 3.12.0 64-bit Go Live

Search [Taskbar icons: File Explorer, Edge, WhatsApp, etc.] 11:37 PM 1/26/2024

```
Do you want to test values from your own of this algorithm (y/n)

Enter values for each feature:
Enter value for Pregnancies: 3
Enter value for Glucose: 30
Enter value for BloodPressure: 20
Enter value for SkinThickness: 34.5
Enter value for Insulin: 10.4
Enter value for BMI: 32
Enter value for DiabetesPedigreeFunction: 52
Enter value for Age: 30
The expected result for values entered is to (NOT) have diabetes :)
```

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
Choose the algorithm you want to test:
1. Decision Tree
2. Naive Bayes
3. Neural Network
4. Random Forest
5. Exit
Enter your choice (1, 2, 3, 4, or 5): 2

Naive Bayes Algorithm:-

Confusion Matrix:
Classified Positive    Actual Positive    Actual Negative
Classified Negative    111              76
                      52              315

Time elapsed: 0.00852 seconds
Accuracy: 0.76895
Precision: 0.68098
Recall: 0.59358
F1-score: 0.63429

Do you want to test values from your own of this algorithm (y/n)
```

Ln 144, Col 1 Spaces: 4 UTF-8 CRLF Python 3.12.0 64-bit Go Live

Search [Taskbar icons: File Explorer, Edge, WhatsApp, etc.] 11:38 PM 1/26/2024

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
Choose the algorithm you want to test:
1. Decision Tree
2. Naive Bayes
3. Neural Network
4. Random Forest
5. Exit
Enter your choice (1, 2, 3, 4, or 5): 3

Neural Network Algorithm:-

Confusion Matrix:
Classified Positive    Actual Positive    Actual Negative
Classified Negative    76              111
                      25              342

Time elapsed: 0.0172 seconds
Accuracy: 0.75451
Precision: 0.75248
Recall: 0.40642
F1-score: 0.52778

Do you want to test values from your own of this algorithm (y/n)

Ln 144, Col 1 Spaces: 4 UTF-8 CRLF Python 3.12.0 64-bit Go Live
Search [Icons] 11:39 PM 1/26/2024
```

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
Choose the algorithm you want to test:
1. Decision Tree
2. Naive Bayes
3. Neural Network
4. Random Forest
5. Exit
Enter your choice (1, 2, 3, 4, or 5): 4

Random Forest Algorithm:-

Confusion Matrix:
Classified Positive    Actual Positive    Actual Negative
Classified Negative    180              7
                      3              364

Time elapsed: 0.02767 seconds
Accuracy: 0.98195
Precision: 0.98361
Recall: 0.96257
F1-score: 0.97297

Do you want to test values from your own of this algorithm (y/n)

Ln 144, Col 1 Spaces: 4 UTF-8 CRLF Python 3.12.0 64-bit Go Live
Search [Icons] 11:39 PM 1/26/2024
```