

Predictive Modeling of Diabetes Using Decision Tree, K-Nearest Neighbour and Logistic Regression

Importing Dataset and Libraries

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 !pip install ydata-profiling
7 from ydata_profiling import ProfileReport
8
9 from sklearn.preprocessing import LabelEncoder
10 from sklearn.preprocessing import MinMaxScaler
11 from sklearn.preprocessing import StandardScaler
12
13 from sklearn.tree import plot_tree
14 from sklearn.model_selection import train_test_split
15 from sklearn.model_selection import StratifiedKFold
16 from sklearn.model_selection import cross_val_predict
17 from sklearn.tree import DecisionTreeClassifier
18 from sklearn.linear_model import LogisticRegression
19 from sklearn.neighbors import KNeighborsClassifier
20
21 from sklearn.metrics import confusion_matrix
22 from sklearn.metrics import classification_report
23 from sklearn.metrics import precision_score, accuracy_score, recall_score, f1_score

```

Collecting ydata-profiling

Downloading ydata_profiling-4.12.0-py2.py3-none-any.whl.metadata (20 kB)

Requirement already satisfied: scipy<1.14,>=1.4.1 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (1.13.1)

Requirement already satisfied: pandas<1.4.0,<3,>1.1 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (2.2.2)

Requirement already satisfied: matplotlib<3.10,>=3.5 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (3.8.0)

Requirement already satisfied: pydantic<=2 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (2.9.2)

Requirement already satisfied: PyYAML<6.1,>=5.0.0 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (6.0.2)

Requirement already satisfied: Jinja2<3.2,>=2.11.1 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (3.1.4)

Collecting visions<0.7.7,>=0.7.5 (from visions[type_image_path]<0.7.7,>=0.7.5->ydata-profiling)

Downloading visions-0.7.6-py3-none-any.whl.metadata (11 kB)

Requirement already satisfied: numpy<2.2,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (1.26.4)

Collecting htmlmin==0.1.12 (from ydata-profiling)

Downloading htmlmin-0.1.12.tar.gz (19 kB)

Preparing metadata (setup.py) ... done

Collecting phik<0.13,>=0.11.1 (from ydata-profiling)

Downloading phik-0.12.4-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (5.6 kB)

Requirement already satisfied: requests<3,>=2.24.0 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (2.32.3)

Requirement already satisfied: tqdm<5,>=4.48.2 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (4.66.6)

Requirement already satisfied: seaborn<0.14,>=0.10.1 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (0.13.2)

Collecting multimethod<2,>=1.4 (from ydata-profiling)

Downloading multimethod-1.12-py3-none-any.whl.metadata (9.6 kB)

Requirement already satisfied: statsmodels<1,>=0.13.2 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (0.14.4)

Requirement already satisfied: typeguard<5,>=3 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (4.4.1)

Collecting imagehash==4.3.1 (from ydata-profiling)

Downloading ImageHash-4.3.1-py2.py3-none-any.whl.metadata (8.0 kB)

Requirement already satisfied: wordcloud==1.9.3 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (1.9.4)

Collecting dacite==1.8 (from ydata-profiling)

Downloading dacite-1.8.1-py3-none-any.whl.metadata (15 kB)

Requirement already satisfied: numba<1,>=0.56.0 in /usr/local/lib/python3.10/dist-packages (from ydata-profiling) (0.60.0)

Collecting PyWavelets (from imagehash==4.3.1->ydata-profiling)

Downloading pywavelets-1.7.0-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (9.0 kB)

Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-packages (from imagehash==4.3.1->ydata-profiling) (11.0.0)

Requirement already satisfied: MarkupSafe==2.0 in /usr/local/lib/python3.10/dist-packages (from Jinja2<3.2,>=2.11.1->ydata-profiling) (3.0.2)

Requirement already satisfied: contourpy==1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (1.3.1)

Requirement already satisfied: cycler==0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (0.12.1)

Requirement already satisfied: fonttools==4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (4.54.1)

Requirement already satisfied: kiwisolver==1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (1.4.7)

Requirement already satisfied: packaging==20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (24.2)

Requirement already satisfied: pyparsing==2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (3.2.0)

Requirement already satisfied: python-dateutil==2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata-profiling) (2.8.2)

Requirement already satisfied: llvmlite<0.44,>=0.43.0dev0 in /usr/local/lib/python3.10/dist-packages (from numba<1,>=0.56.0->ydata-profiling) (0.43.0)

Requirement already satisfied: pytz==2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas<1.4.0,<3,>1.1->ydata-profiling) (2024.2)

Requirement already satisfied: tzdata==2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas<1.4.0,<3,>1.1->ydata-profiling) (2024.2)

Requirement already satisfied: joblib==0.14.1 in /usr/local/lib/python3.10/dist-packages (from phik<0.13,>=0.11.1->ydata-profiling) (1.4.2)

Requirement already satisfied: annotated-types==0.6.0 in /usr/local/lib/python3.10/dist-packages (from pydantic<=2->ydata-profiling) (0.7.0)

Requirement already satisfied: pydantic-core==2.23.4 in /usr/local/lib/python3.10/dist-packages (from pydantic<=2->ydata-profiling) (2.23.4)

Requirement already satisfied: typing-extensions==4.6.1 in /usr/local/lib/python3.10/dist-packages (from pydantic<=2->ydata-profiling) (4.12.2)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata-profiling) (3.4.0)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata-profiling) (3.10)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata-profiling) (2.2.3)

Requirement already satisfied: certifi==2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata-profiling) (2024.8.30)

Requirement already satisfied: patsy==0.5.6 in /usr/local/lib/python3.10/dist-packages (from statsmodels<1,>=0.13.2->ydata-profiling) (1.0.1)

Requirement already satisfied: attrs==19.3.0 in /usr/local/lib/python3.10/dist-packages (from visions<0.7.7,>=0.7.5->visions[type_image_path]<0.7.7,>=visions[type_image_path]<0.7.7,>=0.7.5->ydata-profiling) (24.2)

Requirement already satisfied: networkx==2.4 in /usr/local/lib/python3.10/dist-packages (from visions<0.7.7,>=0.7.5->visions[type_image_path]<0.7.7,>=visions[type_image_path]<0.7.7,>=0.7.5->ydata-profiling) (3.0)

Requirement already satisfied: six==1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil==2.7->matplotlib<3.10,>=3.5->ydata-profiling) (1.16.0)

Downloading ydata_profiling-4.12.0-py2.py3-none-any.whl (390 kB)

390.6/390.6 KB 7.2 MB/s eta 0:00:00

```
1 from google.colab import files
2 uploaded = files.upload()
```

Choose Files

Dataset of Diabetes .csv

Dataset of Diabetes .csv(text/csv) - 49511 bytes, last modified: 10/4/2024 - 100% done

Saving Dataset of Diabetes .csv to Dataset of Diabetes .csv

```
1 db = pd.read_csv("Dataset of Diabetes .csv")
2 db
```

	ID	No_Patient	Gender	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS		
	0	502	17975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N	<div><div></div><div></div></div>
	1	735	34221	M	26	4.5	62	4.9	3.7	1.4	1.1	2.1	0.6	23.0	N	<div><div></div><div></div></div>
	2	420	47975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N	
	3	680	87656	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N	
	4	504	34223	M	33	7.1	46	4.9	4.9	1.0	0.8	2.0	0.4	21.0	N	
	
	995	200	454317	M	71	11.0	97	7.0	7.5	1.7	1.2	1.8	0.6	30.0	Y	
	996	671	876534	M	31	3.0	60	12.3	4.1	2.2	0.7	2.4	15.4	37.2	Y	
	997	669	87654	M	30	7.1	81	6.7	4.1	1.1	1.2	2.4	8.1	27.4	Y	
	998	99	24004	M	38	5.8	59	6.7	5.3	2.0	1.6	2.9	14.0	40.5	Y	
	999	248	24054	M	54	5.0	67	6.9	3.8	1.7	1.1	3.0	0.7	33.0	Y	

1000 rows × 14 columns

Next steps:

Generate code with db

View recommended plots

New interactive sheet

Stage 1: Data Preparation

Data Formatting

```
1 #renamed columns to remove capitalization, correct spelling and have overall consistency
2
3 db.columns = ['id', 'no_patient', 'gender', 'age', 'urea', 'cr',
4               'hba1c', 'chol', 'tg', 'hdl',
5               'ldl', 'vldl', 'bmi', 'db_class']
6
7 db.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   id          1000 non-null   int64
1   no_patient  1000 non-null   int64
2   gender      1000 non-null   object
3   age         1000 non-null   int64
4   urea        1000 non-null   float64
5   cr          1000 non-null   int64
6   hba1c       1000 non-null   float64
7   chol        1000 non-null   float64
8   tg          1000 non-null   float64
9   hdl         1000 non-null   float64
10  ldl         1000 non-null   float64
11  vldl        1000 non-null   float64
12  bmi         1000 non-null   float64
13  db_class    1000 non-null   object
dtypes: float64(8), int64(4), object(2)
memory usage: 109.5+ KB
```

Data Cleaning

```
1 #checking for missing values
2
3 print(db.isnull().sum())
```

```
id          0
no_patient  0
gender      0
age         0
```

```

urea      0
cr        0
hba1c     0
chol      0
tg         0
hdl        0
ldl        0
vldl      0
bmi        0
db_class  0
dtype: int64

```

```

1 #checking for duplicate values in id column
2
3 duplicate_id = db['id'].duplicated().sum()
4 duplicate_id

```

```
200
```

```

1 #identifying the duplicate id vlaues
2
3 duplicate_num_id = db['id'].value_counts()
4 duplicate_num_id = duplicate_num_id[duplicate_num_id > 1]
5 print(duplicate_num_id)

```

```

id
76      2
108     2
57      2
26      2
69      2
..
150     2
49      2
144     2
145     2
147     2
Name: count, Length: 200, dtype: int64

```

```

1 #Random sample #1 - Checking if the ids are duplicate values only or duplicate patient record.
2 #Since the other columns have different values, this might be duplicate id values only.
3
4 display(db[db['id'] == 76])

```

```

id no_patient gender age urea cr hba1c chol tg hdl ldl vldl bmi db_class
249 76      9903    M  73   4.3 79   6.0  5.3 1.4 1.5 3.2 0.6 27.0    Y
910 76      8978    M  60   5.4 64  10.4  3.8 1.5 0.8 2.3 0.6 31.0    Y

```

```

1 #Random sample #2
2 #Helps to confirm that these are duplicate id values rather than duplicate patient records.
3
4 display(db[db['id'] == 150])

```

```

id no_patient gender age urea cr hba1c chol tg hdl ldl vldl bmi db_class
88  150      45382   F  40   6.3 79   4.9  4.3 0.8 0.8 1.8 1.2 22.0    N
184 150      34287   F  42   2.8 39   4.6  4.7 2.5 1.3 2.4 1.1 25.0    Y

```

```

1 #checking for duplicate values in patient number column
2
3 duplicate_no_patient = db['no_patient'].duplicated().sum()
4 duplicate_no_patient

```

```
39
```

```

1 #identifying the duplicate patient numbers
2
3 duplicate_num_patient = db['no_patient'].value_counts()
4 duplicate_num_patient = duplicate_num_patient[duplicate_num_patient > 1]
5 print(duplicate_num_patient)

```

```

no_patient
454316    19
856        2
87654      2
71741      2
34290      2
14389      2
34517      2
48362      2
45646      2
44835      2

```

```

24033    2
34514    2
23972    2
34325    2
234      2
45370    2
34516    2
34518    2
2345     2
34515    2
34545    2
345      2
Name: count, dtype: int64

```

```

1 #Random sample #1 - Checking if the patient number are duplicate values only or duplicate patient record.
2 #Since the other columns have different values, this might be duplicate values only.
3
4 display(db[db['no_patient'] == 45646])

```

```

↗

```

	id	no_patient	gender	age	urea	cr	hba1c	chol	tg	hdl	ldl	vldl	bmi	db_class
936	111	45646	F	63	4.7	55	12.60	7.9	5.0	1.3	2.1	0.7	35.0	Y
938	113	45646	F	55	2.1	23	9.96	4.1	4.2	1.2	1.4	1.3	29.0	Y

```

1 #Random sample #2
2 #Helps to confirm that these are duplicate patient number rather than duplicate patient records.
3
4 display(db[db['no_patient'] == 34290])

```

```

↗

```

	id	no_patient	gender	age	urea	cr	hba1c	chol	tg	hdl	ldl	vldl	bmi	db_class
34	699	34290	F	47	5.6	67	5.1	6.5	1.5	0.9	4.9	0.7	23.0	N
260	510	34290	M	73	4.3	79	6.9	5.3	1.4	1.5	3.2	0.6	28.0	Y

```

1 #checking the gender column
2
3 gender_frequency = db['gender'].value_counts()
4 print(gender_frequency)

```

```

↗
gender
M    565
F    434
f      1
Name: count, dtype: int64

```

```

1 #Converting the lowercase f to uppercase and updating the gender_frequency variable
2
3 db['gender'] = db['gender'].replace({'f': 'F'})
4 gender_frequency = db['gender'].value_counts()
5 print(gender_frequency)

```

```

↗
gender
M    565
F    435
Name: count, dtype: int64

```

```

1 #checking the class column
2
3 class_frequency = db['db_class'].value_counts()
4 print(class_frequency)

```

```

↗
db_class
Y    840
N   102
P    53
Y     4
N     1
Name: count, dtype: int64

```

```

1 #There is a space trailing after some of the Y and N values (eg: 'Y ' and 'N '). This is causing the data to be read as different values.
2 #The below converts these values to Y and N respectively
3
4 db['db_class'] = db['db_class'].replace({'N ': 'N', 'Y ': 'Y'})
5 class_frequency = db['db_class'].value_counts()
6 print(class_frequency)

```

```

↗
db_class
Y    844
N   103
P    53
Name: count, dtype: int64

```

Stage 2: Exploratory Data Analysis

Statistical Summary

```
1 #Statistical summary of the attributes
2
3 db.describe()
```

	id	no_patient	age	urea	cr	hba1c	cho1	tg	hdl	ldl	vldl	
count	1000.000000	1.000000e+03	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000
mean	340.500000	2.705514e+05	53.528000	5.124743	68.943000	8.281160	4.862820	2.349610	1.204750	2.609790	1.854700	29.510000
std	240.397673	3.380758e+06	8.799241	2.935165	59.984747	2.534003	1.301738	1.401176	0.660414	1.115102	3.663599	4.960000
min	1.000000	1.230000e+02	20.000000	0.500000	6.000000	0.900000	0.000000	0.300000	0.200000	0.300000	0.100000	19.000000
25%	125.750000	2.406375e+04	51.000000	3.700000	48.000000	6.500000	4.000000	1.500000	0.900000	1.800000	0.700000	26.000000
50%	300.500000	3.439550e+04	55.000000	4.600000	60.000000	8.000000	4.800000	2.000000	1.100000	2.500000	0.900000	30.000000
75%	550.250000	4.538425e+04	59.000000	5.700000	73.000000	10.200000	5.600000	2.900000	1.300000	3.300000	1.500000	33.000000
max	800.000000	7.543566e+07	79.000000	38.900000	800.000000	16.000000	10.300000	13.800000	9.900000	9.900000	35.000000	47.700000

The average age of the patients in the dataset is 53.53, suggesting the population is mostly middle age. The creatinine variable has a large standard deviation (59.98) and a large range (from 6 to 800), suggesting the presence of extreme values.

Correlation Analysis

```
1 #Using Panda Profile Report for further analysis
2
3 report = ProfileReport(db)
4 report
```

Show hidden output

Report confirms there are no missing values or duplicate rows within the dataset.

Both HbA1c and Cholesterol exhibit a normal bell shaped distribution, indicating a relatively even spread of data around the mean. The other features appear skewed, suggesting that their data distributions are more concentrated around certain values, and the presence of outliers.

Cholesterol, triglycerides, HDL and LDL have variances below 2, with the lowest being 0.44 for the HDL variable. This suggests that these variables exhibit small spread with most of the values being close to the mean.

There are several correlations observed among the variable. Triglycerides(tg) has a high positive correlation with vldl (0.599), and creatinine has a high positive correlation with urea (0.568). BMI exhibits moderately positive correlation with a coefficient of 0.417, suggesting elevated HbA1c levels are associated with BMI.

The target variable, db_class, has a strong positive correlation with bmi (0.549) and hba1c (0.662). This suggests that both bmi and hba1c are influential and important in determining the likelihood of diabetes.

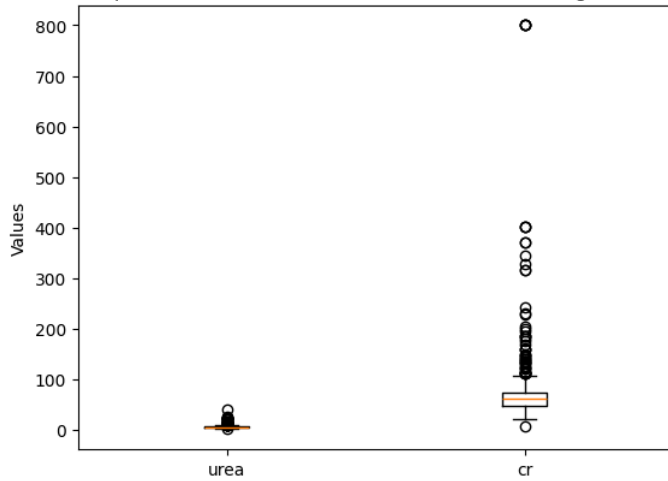
There is a significant class imbalance with 844 individuals classified as diabetic, 103 as non-diabetic, and only 53 as pre-diabetic. This 51.5% imbalance may affect the performance of the predictive algorithms leading to bias towards the majority class.

Data Visualization

```
1 #Observing distribution of values in Urea and Creatinine
2
3 plt.boxplot(db[['urea','cr']])
4 plt.title('Boxplot of Urea Levels and Creatinine Ratio Among Patients')
5
6 plt.xticks([1, 2], ['urea','cr'])
7 plt.ylabel("Values")
8
9 plt.show()
```



Boxplot of Urea Levels and Creatinine Ratio Among Patients



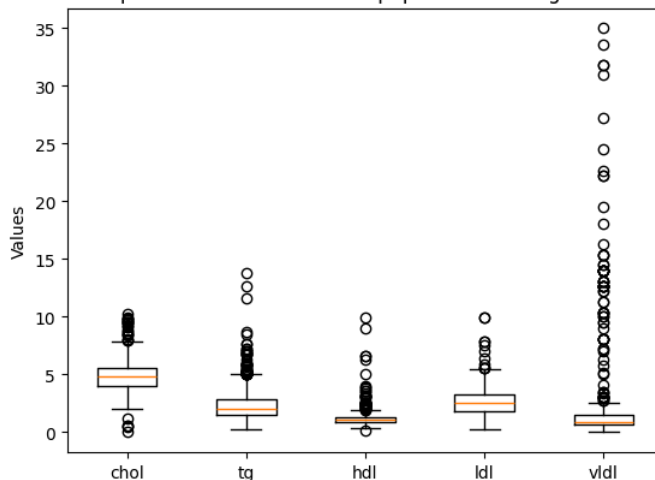
```

1 #Observing distribution of values in Urea and Creatinine
2
3 plt.boxplot(db[['chol','tg','hdl','ldl','vldl']])
4 plt.title('Boxplot for Cholesterol and Lipoproteins Among Patients')
5
6 plt.xticks([1, 2, 3, 4, 5], ['chol','tg','hdl','ldl','vldl'])
7 plt.ylabel("Values")
8
9 plt.show()

```



Boxplot for Cholesterol and Lipoproteins Among Patients



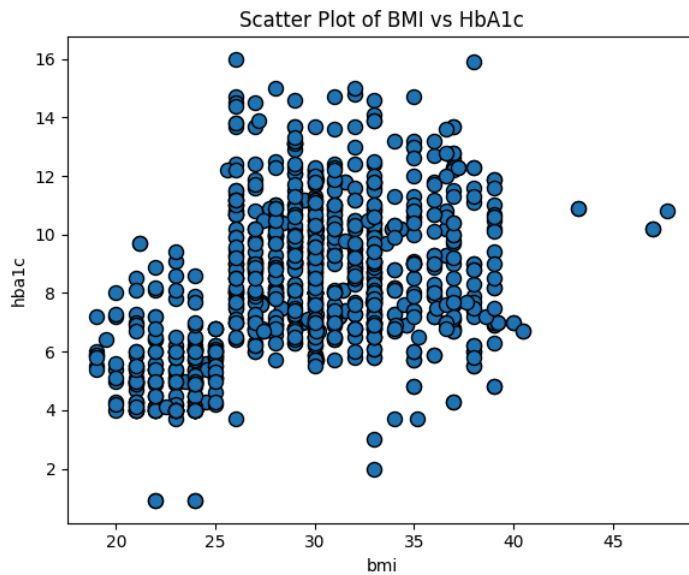
There are noticeable outliers present in the cholesterol, lipoprotein, urea, and creatinine attributes. This suggests that several patients may have extreme values for these attributes.

Urea and creatinine ratio are clinical measurements used to assess kidney impairment. While cholesterol, the lipoproteins and triglycerides levels are used to determine cardiovascular damage. These outliers therefore suggest there may be patients with kidney and cardiovascular damage.

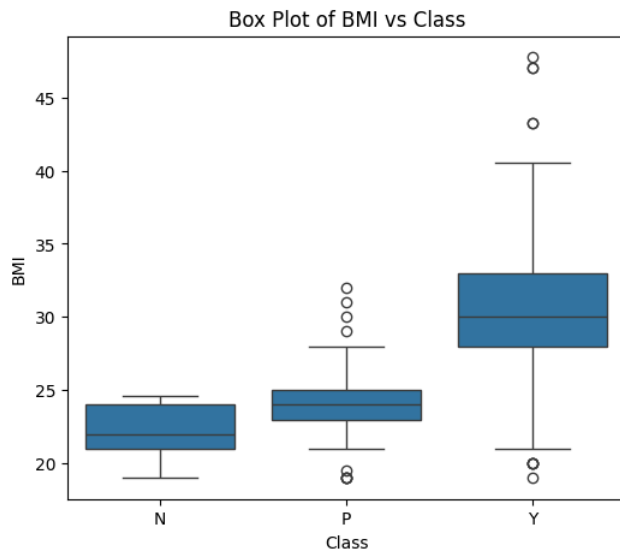
```

1 #Graph showing relationship between BMI and Glucose
2
3 plt.figure(figsize=(6, 5))
4 plt.scatter(db['bmi'], db['hba1c'], s=75, edgecolor='k')
5
6 plt.title('Scatter Plot of BMI vs HbA1c')
7 plt.xlabel('bmi')
8 plt.ylabel('hba1c')
9 plt.tight_layout()
10 plt.show()

```



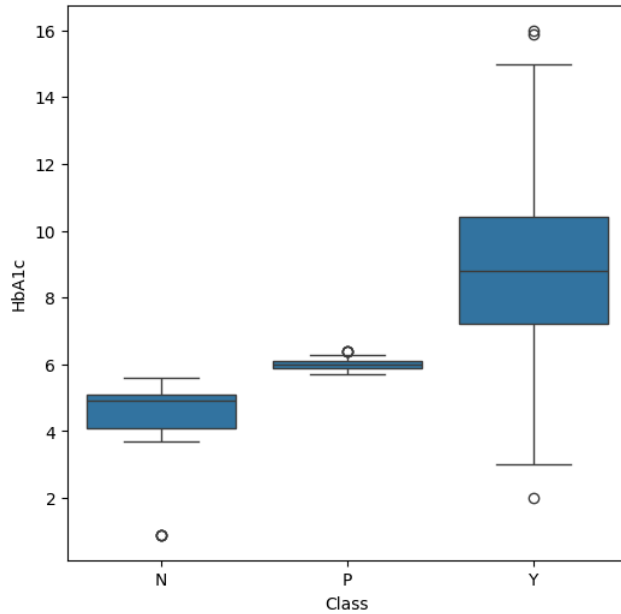
```
1 #Graph showing relationship between BMI and Class
2
3 plt.figure(figsize=(6, 5))
4 sns.boxplot(x='db_class', y='bmi', data=db)
5 plt.title('Box Plot of BMI vs Class')
6 plt.xlabel('Class')
7 plt.ylabel('BMI')
8 plt.show()
```



```
1 #Graph showing relationship between BMI and Class
2
3 plt.figure(figsize=(6, 6))
4 sns.boxplot(x='db_class', y='hba1c', data=db)
5 plt.title('Box Plot of HbA1c vs Class')
6 plt.xlabel('Class')
7 plt.ylabel('HbA1c')
8 plt.show()
```



Box Plot of HbA1c vs Class



Stage 3: Data Preprocessing

Label Encoding

```

1 #Tried to apply DecisionTreeClassifier to dataset, however got an error - ValueError: could not convert string to float: 'F'. This indicates that the Deci
2 #Applied LabelEncoder which allows each string in the categorical variables (gender and db_class columns) to be replaced with a numerical value
3
4 db_encode=db.copy()
5
6 label_encoder = LabelEncoder()
7
8 db_encode['gender'] = label_encoder.fit_transform(db_encode['gender'])
9 db_encode['db_class'] = label_encoder.fit_transform(db_encode['db_class'])
10
11 db_encode.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   id           1000 non-null   int64
1   no_patient   1000 non-null   int64
2   gender       1000 non-null   int64
3   age         1000 non-null   int64
4   urea        1000 non-null   float64
5   cr          1000 non-null   int64
6   hba1c       1000 non-null   float64
7   chol        1000 non-null   float64
8   tg          1000 non-null   float64
9   hdl         1000 non-null   float64
10  ldl         1000 non-null   float64
11  vldl        1000 non-null   float64
12  bmi         1000 non-null   float64
13  db_class    1000 non-null   int64
dtypes: float64(8), int64(6)
memory usage: 109.5 KB

```

```

1 #Compared the count with db_class count performed during cleaning to ensure consistency and understand the new format.
2 #Class Attributes are now: 0 = Not diabetic (previously N) ; 1 = Prediabetic (previously P) ; 2 = Diabetic (previously Y)
3 db_encode['db_class'].value_counts()

```




count	
db_class	
2	844
0	103
1	53

dtypes: int64

▼ Data Normalization

```
1 #Normalizing to scale the dataset to a standard range, and reduce influence of outliers:
2
3 def normalize(x):
4     return ((x - min(x)) / (max(x) - min(x)))
```

```
1 X = list(set(list(db_encode)) - set(['db_class']))
2 db_norm = db_encode.iloc[:,].copy()
3 db_norm[X] = db_norm[X].apply(normalize)
4
5 db_norm.head(5)
```



	id	no_patient	gender	age	urea	cr	hba1c	chol	tg	hdl	ldl	vldl	bmi	db_class
0	0.627034	0.000237	0.0	0.508475	0.109375	0.050378	0.264901	0.407767	0.044444	0.226804	0.114583	0.011461	0.173913	0
1	0.918648	0.000452	1.0	0.101695	0.104167	0.070529	0.264901	0.359223	0.081481	0.092784	0.187500	0.014327	0.139130	0
2	0.524406	0.000634	0.0	0.508475	0.109375	0.050378	0.264901	0.407767	0.044444	0.226804	0.114583	0.011461	0.173913	0
3	0.849812	0.001160	0.0	0.508475	0.109375	0.050378	0.264901	0.407767	0.044444	0.226804	0.114583	0.011461	0.173913	0
4	0.629537	0.000452	1.0	0.220339	0.171875	0.050378	0.264901	0.475728	0.051852	0.061856	0.177083	0.008596	0.069565	0

Next steps:

[Generate code with db_norm](#)[View recommended plots](#)[New interactive sheet](#)

▼ Data Splitting

```
1 train_norm_set, test_norm_set = train_test_split(db_norm, train_size = 0.7, random_state = 89)
2
3 # Separating feature variables from the target variable
4 X_train = train_norm_set.drop('db_class', axis=1)
5 y_train = train_norm_set['db_class']
6
7 X_test = test_norm_set.drop('db_class', axis=1)
8 y_test = test_norm_set['db_class']
```

▼ Stage 4: Data Modelling

▼ Cross Validation of Base Models

```
1 #Decided to use StratifiedKFold for cross validation given the class imbalance within the dataset. StratifiedKFold ensures that each fold maintains the sa
2 #This ensure that the cross-validation process is more reliable in representing the class distribution in the original dataset, and prevents misleading r
3
4 strat_fold = StratifiedKFold(n_splits=10, random_state=10, shuffle=True)
5
6 dc_model = DecisionTreeClassifier(random_state=10)
7 pred_dc_val = cross_val_predict(dc_model, X_train, y_train, cv=strat_fold)
8
9 lgr_model = LogisticRegression(random_state=10, multi_class='multinomial')
10 pred_lgr_val = cross_val_predict(lgr_model, X_train, y_train, cv=strat_fold)
11
12 knn_model = KNeighborsClassifier(n_neighbors=5)
13 pred_knn_val = cross_val_predict(knn_model, X_train, y_train, cv=strat_fold)
```



```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
```

```

/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(

```

▼ Training Second Level Model

```

1 # Stacking the validation predictions from each base model to form the training set for the second level model
2
3 X_seclvl_train = np.column_stack((pred_dc_val, pred_lgr_val, pred_knn_val))

1 # Training the second level model - Logistic Regression
2
3 seclvl_model = LogisticRegression(random_state=20)
4 seclvl_model.fit(X_seclvl_train, y_train)

```

LogisticRegression

LogisticRegression(random_state=20)

▼ Final Prediction

```

1 # Training the base models on all the training dataset to adequately make predictions on the test set
2
3 dc_model.fit(X_train, y_train)
4 pred_dc_test = dc_model.predict(X_test)
5
6 lgr_model.fit(X_train, y_train)
7 pred_lgr_test = lgr_model.predict(X_test)
8
9 knn_model.fit(X_train, y_train)
10 pred_knn_test = knn_model.predict(X_test)

```

```

/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be r
warnings.warn(

```

```

1 # Stacking the predictions made by the base models on the test set
2 X_seclvl_test = np.column_stack((pred_dc_test, pred_lgr_test, pred_knn_test))
3
4 # Final prediction using the second level model
5 pred_final = seclvl_model.predict(X_seclvl_test)

```

Double-click (or enter) to edit

▼ Stage 5: Model Evaluation

▼ Base Model Evaluation

```

1 cf_dc = confusion_matrix(y_test, pred_dc_test)
2 print("Confusion Matrix for Decision Tree Base Model:\n")
3 print(cf_dc)

```

Confusion Matrix for Decision Tree Base Model:

```

[[ 27  0  1]
 [  0 16  2]
 [  0  4 250]]

```

```

1 print("Confusion Matrix for Decision Tree Base Model:\n")
2 print(classification_report(y_test, pred_dc_test))
3 dc_accuracy = accuracy_score(y_test, pred_dc_test)
4 print("Decision Tree Accuracy:", dc_accuracy)

```

Confusion Matrix for Decision Tree Base Model:

	precision	recall	f1-score	support
0	1.00	0.96	0.98	28
1	0.80	0.89	0.84	18
2	0.99	0.98	0.99	254
accuracy			0.98	300
macro avg	0.93	0.95	0.94	300
weighted avg	0.98	0.98	0.98	300

Decision Tree Accuracy: 0.9766666666666667

```
1 cf_lgr =confusion_matrix(y_test, pred_lgr_test)
2 print("Confusion Matrix for Logisitic Regression Base Model:\n")
3 print(cf_lgr)
```

Confusion Matrix for Logisitic Regression Base Model:

```
[[ 27  0  1]
 [  4  0 14]
 [  8  0 246]]
```

```
1 print("Confusion Matrix for Logisitic Regression Base Model:\n")
2 print(classification_report(y_test, pred_lgr_test))
3 lgr_accuracy = accuracy_score(y_test, pred_lgr_test)
4 print("Logisitic Regression Accuracy:", lgr_accuracy)
```

Confusion Matrix for Logisitic Regression Base Model:

	precision	recall	f1-score	support
0	0.69	0.96	0.81	28
1	0.00	0.00	0.00	18
2	0.94	0.97	0.96	254
accuracy			0.91	300
macro avg	0.54	0.64	0.59	300
weighted avg	0.86	0.91	0.88	300

Logisitic Regression Accuracy: 0.91
 /usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 i
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
 /usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 i
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
 /usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 i
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

```
1 cf_knn =confusion_matrix(y_test, pred_knn_test)
2 print("Confusion Matrix for kNN Base Model:\n")
3 print(cf_knn)
```

Confusion Matrix for kNN Base Model:

```
[[ 24  2  2]
 [  7  6  5]
 [  6  4 244]]
```

```
1 print("Confusion Matrix for kNN Base Model:\n")
2 print(classification_report(y_test, pred_knn_test))
3 knn_accuracy = accuracy_score(y_test, pred_knn_test)
4 print("kNN Accuracy:", knn_accuracy)
```

Confusion Matrix for kNN Base Model:

	precision	recall	f1-score	support
0	0.65	0.86	0.74	28
1	0.50	0.33	0.40	18
2	0.97	0.96	0.97	254
accuracy			0.91	300
macro avg	0.71	0.72	0.70	300
weighted avg	0.91	0.91	0.91	300

kNN Accuracy: 0.9133333333333333

Second Level Model Evaluation

```
1 cf_final =confusion_matrix(y_test, pred_final)
2 print("Confusion Matrix for Second Level Model:\n")
3 print(cf_final)
```

 Confusion Matrix for Second Level Model: