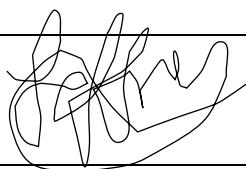




**C3879C - Capstone Project  
AY2021 Term 3  
Coursework Final Submission**

<b>Personal Details</b>	
<b>Name</b>	<b>WONG QI YUAN, JEFFREY</b>
<b>Student No.</b>	<b>20053371</b>

<b>Compliance Statement</b>	
<b>Plagiarism</b> I declare that this report is my original work. I understand that if I am suspected of plagiarism, my enrolment in the programme may be terminated.	<input checked="" type="checkbox"/>
<b>Retention of Backup Copy</b> I declare that I have a back-up electronic copy of this report for immediate submission.	<input checked="" type="checkbox"/>
<b>Signature</b>	
<b>IMPORTANT:</b> Non-compliance to these clauses will result in unconditional rejection of your submission	

DEVELOPMENT OF HEALTH WEB APP  
WITH INTEGRATION OF ENSEMBLE  
MACHINE LEARNING FOR EARLY  
DIAGNOSIS IN OBESITY-RELATED  
CHRONIC DISEASES

Date of Submission: 01-10-2021

Submitted By:

20053371

WONG QI YUAN, JEFFREY

# Table of Contents

<b>1.0 Abstract</b>	5
<b>2.0 Acknowledgement</b>	6
<b>3.0 Background</b>	7
3.0.1 Problem Statement	7
3.0.2 Overview of the Project Context	7
3.0.3 Basic Concepts of Obesity-related Chronic Diseases	8
<b>4.0 Methodology and Design</b>	9
<b>4.1 Datasets and SQL Database</b>	9
4.1.1 Diabetes	9
4.1.2 Hypertension (or High Blood Pressure)	12
4.1.3 Obesity	14
<b>4.2 Data Exploration on the Selected Attributes</b>	17
4.2.1 Diabetes	17
4.2.2 Hypertension	19
4.2.3 Obesity	21
<b>4.3 Data Cleansing</b>	24
4.3.1 Investigate Missing Data	24
4.3.2 Removal of Potential Outlier Detection using Isolation Forest	25
<b>4.4 Resampling Techniques for Imbalanced Training Sets</b>	29
<b>4.5 Baseline Model Selection and Evaluation via Shuffle-Split Cross-Validation (CV)</b>	34

<b>4.6</b>	<b>Hyperparameters Optimization for Proposed Baseline Models</b>	42
<b>4.7</b>	<b>Retrained Models with Best Hyperparameters and Evaluation</b>	48
<b>4.8</b>	<b>Evaluation of the Hyperparameters Classifier Models with Testing Sets</b>	61
<b>4.9</b>	<b>Development of Health Web Application using Streamlit</b>	69
<b>5.0</b>	<b>Conclusion</b>	83
<b>6.0</b>	<b>Recommendation</b>	84
<b>7.0</b>	<b>Project Implementation Scheduled</b>	85
<b>8.0</b>	<b>References</b>	87

## 1.0 Abstract

With the advancement of technologies in the 21<sup>st</sup> century, artificial intelligence and machine learning have been widely used in the healthcare fields by making use of the predictive models to predict the possible risks of developing obesity-related chronic diseases such as obesity, diabetes, and hypertension at the early stages. Thus, the intervention not only improves the overall healthcare quality but also improving the individual's overall health status by minimizing the detrimental effects caused by obesity-related chronic diseases. This capstone project developed some predictive ensemble machine learning algorithms such as random forest, gradient boosting, and/ or extra trees to provide an early prediction for obesity, diabetes, and hypertension based on individual risk factors data.

Prior to algorithms training, an isolation forest algorithm was used to detect and remove the potential outlier data, followed by synthetic minority oversampling technique Tomek link (SMOTE-Tomek) to balance the training data distribution such that minimized the biased and unreliable outcomes on the machine learning model performance, and the training of the predictive algorithms with the best hyperparameters to predict the respective diseases. Three datasets were utilized to implement the proposed predictive models and the attributes were chosen based on the relevance to the disease contexts. The outcome of the performance evaluation based on the testing sets showed that the proposed predictive models with hyperparameters for each respective disease were achievable with high accuracy of at least 90%, at least 95% for ROC-AUC, and less than 5% of false-positive rate and false-negative rate.

This project also developed a simple and user-friendly web application using Streamlit to simulate the practical application of the proposed predictive models. The developed web application collates all the risk factors from the user inputs and then sends them to the proposed prediction models to predict the respective diseases based on the class probabilities. The prediction outcome, as well as the constructive recommendation, are then displayed on the web application. Thus, prompt action can be taken swiftly to minimize the effects and prevent individual risks once undesirable health conditions such as obesity, diabetes, and hypertension are detected at the early stages.

## 2.0 Acknowledgement

The success and final outcome of this capstone project required a lot of guidance, assistance, and past-experience from many people, especially the lecturers from Specialist Diploma in Artificial Intelligence, who previously sharing their knowledge and skills in AI and Machine Learning. I am extremely privileged to have got this all along with the completion of this project and I would not forget to thank them.

## 3.0 Background

### 3.0.1 Problem Statement

In a recent report, the World Health Organization (WHO) mentioned that obesity-related chronic diseases have contributed to approximately about 41 million premature deaths annually, which is about 71% of all death worldwide. If the situation is unmitigated, we are expected that the overall number of obesity-related chronic diseases related-deaths to be increased up to 52 million yearly by the end of 2030. The most common obesity-related chronic diseases are the obesity, diabetes, and the hypertension.

Furthermore, as the world currently facing the crisis of the COVID-19 pandemic, hence it is expected that the number of people developing the obesity-related chronic diseases to rise as most people are working from home in order to curb the spread of the COVID-19 in the local community. Furthermore, based on a recent news article, about 1/3 of the Singaporeans have gained weight during the COVID-19 pandemic and many fingers point their fingers at how their lifestyles have changed to become more sedentary since the COVID-19 pandemic.

Thus, it is an alarming concern in the road of the covid-19 pandemic as the majority of the people are unaware of their overall health status and some of the people do not have a little bit of time to go for a simple health screening.

### 3.0.2 Overview of the Project Context

In this capstone project, the implementation of the isolation forest algorithm, the SMOTE-TOMEK, and the various classifier algorithms such as random forest, gradient boosting, or extra trees will be integrated to predict the respective possible risk of developing obesity-related chronic diseases such as obesity, diabetes, and hypertension.

Furthermore, three datasets were downloaded from various sources to be used to trained and evaluate the performance of the proposed classifier models for respective obesity-related chronic diseases.

In addition, a web application via Streamlit will be implemented to simulate the real-life practical application with the integration of the proposed classifier models such that offering individual's an effective and appropriate way to monitor their current health status.

### **How does this health web application benefit the user?**

This health web application is significant as it would provide an individual to know their likelihood of being developing the obesity-related chronic-diseases such as diabetes, hypertension, and obesity, and hence, reduce their time to visit the office clinic in person. Furthermore, being cost-effective, the users are able to know the diagnosis right on the spot and also provide the users the time to prevent and manage the chronic-diseases by making them aware of their present condition.

#### **3.0.3 Basic Concepts of Obesity-related Chronic Diseases**

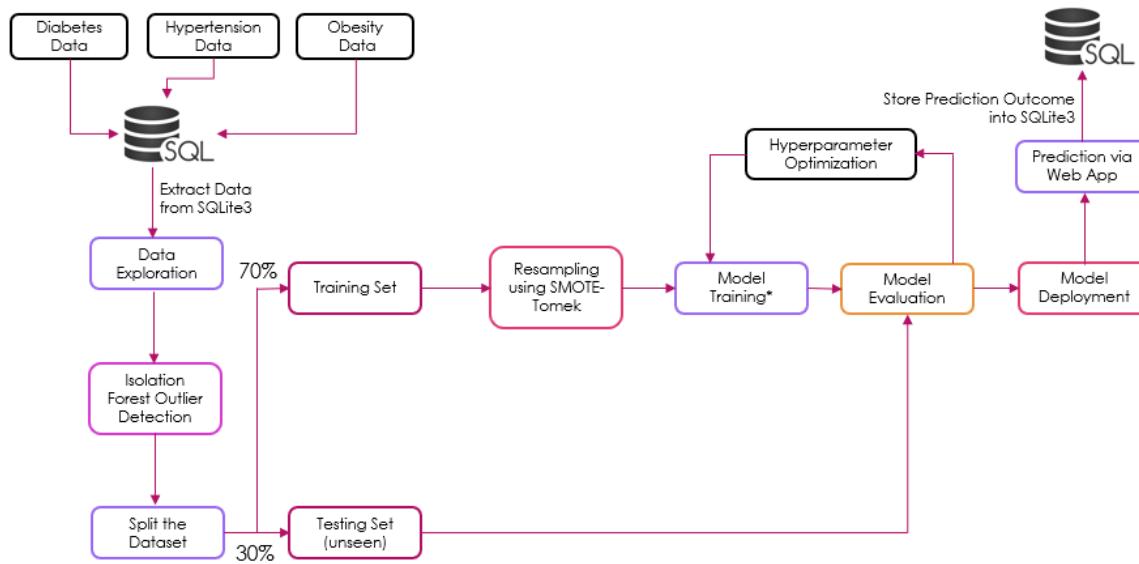
**Obesity** can be defined as a weight that is higher than what is considered healthy for a given weight. BMI is the main tool to screen for obesity. Obesity is a complex health issue resulting from a combination of causes and individuals' factors such as genetics and behaviour. Behaviours can include sedentary lifestyles, physical activity, dietary patterns, and other exposures. Obesity is serious because it is often associated with poor mental health and reduced quality of life. Furthermore, obesity is also associated with the leading causes of death, including diabetes, and high blood pressure.

**Diabetes** is a metabolic disorder that changes the blood sugar levels in the bloodstream. When there is insufficient insulin or cells stop responding to insulin, too much blood sugar will remain in your bloodstream. Over time, this usually leads to consequences of health issues such as vision loss, kidney failure, etc.

**Hypertension (High Blood Pressure)** is a blood pressure that is usually higher than the global health recommended blood pressure (< 140/100 mm Hg). However, having blood pressure that is constantly above the recommended range may result in the diagnosis of high blood pressure or hypertension. Hypertension usually develops over time and it could happen due to sedentary lifestyles or certain existing health conditions such as diabetes, or obesity.

## 4.0 Methodology and Design

Figure 4.0 illustrates a schematic architecture design of the ML predictive models for all the respective diseases (obesity, diabetes, and hypertension)



\* Model Training: Random Forest Classifier, Gradient Boosting Classifier, and Extra Trees Classifier

Figure 4.0: A schematic architecture design of the ML predictive models for all respective diseases

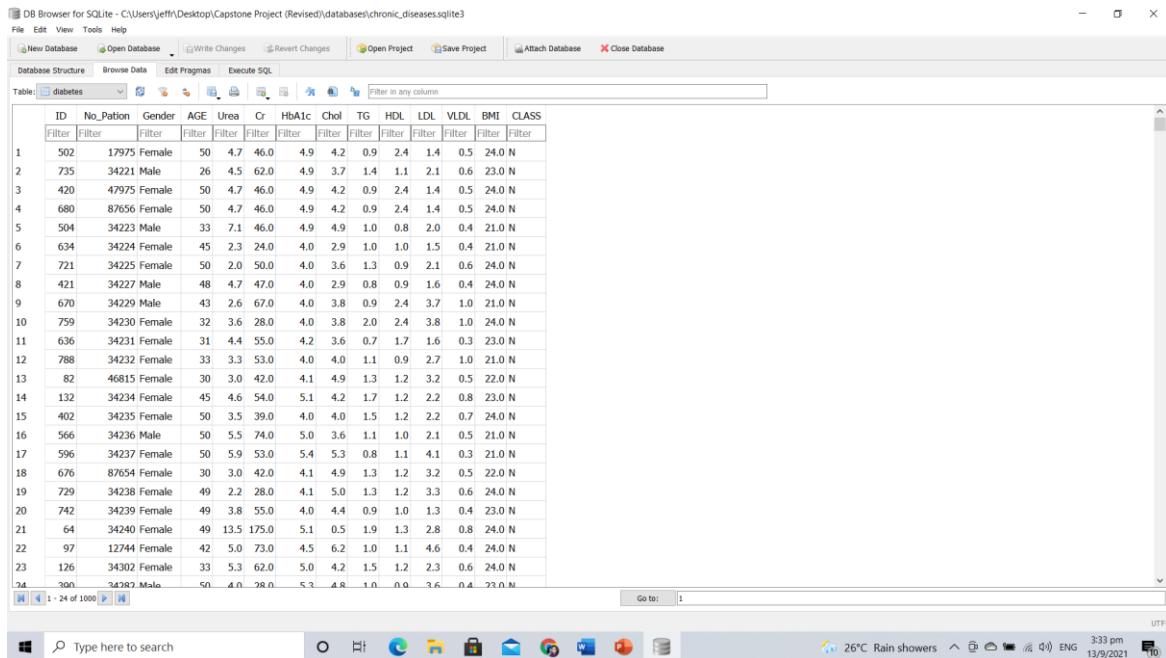
## 4.1 Datasets and SQL Database

### 4.1.1 Diabetes

The diabetes dataset was downloaded from **Mendeley Data** and the data were collected by the Iraqi society as were acquired from the laboratory of Medical City Hospital and the Specialization Center for Endocrinology and Diabetes AL Kindy Teaching Hospital. The original data consists of the patient ID, gender, age, urea, creatinine ratio (Cr), blood glucose level (HbA1c), cholesterol (chol), triglycerides (TG), high density lipid (HDL), low density lipoprotein (LDL), very-low density lipoprotein (VLDL), body-mass index (BMI), and the class outcomes (N: non-diabetic, P: pre-diabetic, or Y: diabetic).

## C3879C CAPSTONE PROJECT

Figure 4.1.1a shows the diabetes dataset (original) that was stored into SQLite3 database via Python.



The screenshot shows the DB Browser for SQLite interface with the 'diabetes' table selected. The table has 14 columns: ID, No\_Pation, Gender, AGE, Urea, Cr, HbA1c, Chol, TG, HDL, LDL, VLDL, BMI, and CLASS. The data consists of 760 rows, each representing a patient's medical profile. The columns are defined as follows:

- ID:** Integer, primary key.
- No\_Pation:** String.
- Gender:** String.
- AGE:** Integer.
- Urea:** Real.
- Cr:** Real.
- HbA1c:** Real.
- Chol:** Real.
- TG:** Real.
- HDL:** Real.
- LDL:** Real.
- VLDL:** Real.
- BMI:** Real.
- CLASS:** String.

Figure 4.1.1a: Diabetes Dataset (Original)

Table 4.1.1a shows the description of the attributes for the diabetes dataset (original)

Attribute	Description
ID	Patient identification number
No_Pation	N/A
Gender	Male or Female
AGE	Age (in years)
Urea	Measure the amount of urea in a blood sample
Cr	Measure the level of creatinine in a blood sample
HbA1C	Measure the amount of fasting blood sugar in a blood sample
Chol	Measure the amount of cholesterol and triglycerides in a blood sample
TG	Measure the amount of thyroglobulin in a blood sample
HDL	Measure the amount of cholesterol found inside high-density lipoprotein in a blood sample
LDL	Measure the amount of cholesterol found inside low-density lipoprotein in a blood sample

## C3879C CAPSTONE PROJECT

VLDL	Measures the amount of cholesterol found inside very-low-density lipoprotein in a blood sample
BMI	Measures the amount of body fat based on height and weight
CLASS	N (Class 0): No Diabetes, P (Class 1): Pre-Diabetes, Y (Class 2): Diabetes

Table 4.1.1a: Description of the Attributes for Diabetes Dataset (Original)

As such, the relevant attributes to be used for training phases were **Gender, HbA1c, BMI, and the Class Outcomes** as the rest of the attributes were irrelevant to the diabetes context.

Figure 4.1.1b illustrate a sample python code to extract the relevant column data from the SQL database for diabetes.

```
23 def load_datasets():
24     conn = sqlite3.connect(r"C:\Users\jeffr\Desktop\Capstone Project (Revised)\databases\chronic_diseases.sqlite3")
25     diabetes_df = pd.read_sql("SELECT Gender, HbA1C, BMI, CLASS FROM diabetes", conn)
26
27     # create a set of copy of diabetes dataset
28     copy_diabetes_df = diabetes_df.copy()
29
30     conn.close()
31     return copy_diabetes_df
```

Figure 4.1.1b: A sample code to extract the relevant column data from SQL database

Figure 4.1.1c illustrate the output of the relevant column data extracted from the SQL database via Python code.

diab_df - DataFrame				
Index	Gender	HbA1c	BMI	CLASS
0	Female	4.9	24	N
1	Male	4.9	23	N
2	Female	4.9	24	N
3	Female	4.9	24	N
4	Male	4.9	21	N
5	Female	4	21	N
6	Female	4	24	N
7	Male	4	24	N
8	Male	4	21	N
9	Female	4	24	N
10	Female	4.2	23	N
11	Female	4	21	N
12	Female	4.1	22	N

Figure 4.1.1c: An output of the relevant column data extracted from the SQL database via Python Code

### 4.1.2 Hypertension (or High Blood Pressure)

The hypertension (or high blood pressure) dataset was downloaded from **figshare** and the **original** dataset consists of the age, obese (Yes or No), body mass index (BMI), waist (WC), hip circumference (HC), waist-hip ratio (WHR), diastolic blood pressure (DBP), systolic blood pressure (SDB), and the class outcomes (Hypertension or No Hypertension).

Figure 4.1.2a shows the hypertension dataset (original) that was stored into SQLite3 database via Python.

The screenshot shows the DB Browser for SQLite interface with the 'hypertension' table selected. The table has columns: id, age, obese, bmi, wc, hc, whr, dbp, sbp, and outcome. The data consists of 399 rows, with the first few rows shown below:

	id	age	obese	bmi	wc	hc	whr	dbp	sbp	outcome
1	1	20	NO	27.0	94.0	95.0	112.0	85	120	REGULAR
2	1	31	NO	28.0	76.0	88.0	101.0	128	87	HYPERTENSION
3	2	19	NO	25.0	35.0	79.0	102.0	77	113	REGULAR
4	2	17	NO	18.0	8.0	78.0	95.0	100	82	REGULAR
5	3	20	NO	20.0	73.0	91.0	80.0	114	130	REGULAR
6	3	19	NO	28.0	48.0	87.0	111.0	140	78	HYPERTENSION
7	4	19	NO	24.0	54.0	83.0	98.0	85	130	REGULAR
8	4	19	NO	19.0	29.0	73.0	95.0	103	77	REGULAR
9	5	19	NO	24.0	66.0	81.0	99.0	82	136	REGULAR
10	5	21	NO	20.0	4.0	63.0	96.0	110	66	REGULAR
11	6	21	NO	20.0	61.0	79.0	96.0	82	130	REGULAR
12	6	27	NO	21.0	79.0	85.0	101.0	93	84	REGULAR
13	7	19	NO	18.0	65.0	69.0	85.0	81	100	REGULAR
14	7	21	NO	27.0	59.0	86.0	110.0	123	78	HYPERTENSION
15	8	31	NO	27.0	62.0	102.0	107.0	95	106	REGULAR
16	8	23	NO	22.0	45.0	72.0	104.0	90	69	REGULAR
17	9	22	NO	27.0	17.0	92.0	110.0	84	136	REGULAR
18	9	19	NO	23.0	16.0	80.0	100.0	113	80	REGULAR
19	10	27	NO	25.0	6.0	83.0	106.0	78	123	REGULAR
20	10	19	NO	21.0	9.0	74.0	94.0	110	79	REGULAR
21	11	21	YES	32.0	8.0	100.0	111.0	90	126	REGULAR
22	11	19	NO	21.0	95.0	70.0	94.0	105	74	REGULAR
23	12	22	NO	28.0	52.0	90.0	99.0	91	127	REGULAR
24	12	21	NO	23.0	58.0	78.0	100.0	140	78	HYPERTENSION

Figure 4.1.2a: Hypertension Dataset (Original)

Table 4.1.2a shows the description of the attributes for the hypertension dataset (original)

Attribute	Description
id	Patient identification numbers
age	Age (in years)
obese	Obesity Status
bmi	Measure the amount of body fat based on height and weight
wc	Measure the waist circumference (in cm)
hc	Measure the hip circumference (in cm)
whr	Ratio of the waist circumference (in cm) and the hip circumference (in cm)

## C3879C CAPSTONE PROJECT

dbp	Diastolic blood pressure is the bottom number found in the blood pressure measuring device. It measures the force your heart exerts on the wall of your arteries each time in between beats.
sbp	Systolic blood pressure is the top number found in the measuring blood pressure measuring device. It measures the force your heart exerts on the wall of your arteries each time it beats.
outcome	Class 0: Regular, Class 1: Hypertension

Table 4.1.2a: Description of the Attributes for Hypertension Dataset (Original)

As such, the relevant attributes to be used for training phases were **BMI, DBP, SBP, and the Class Outcomes** as the remaining attributes were irrelevant to the hypertension context.

Figure 4.1.2b illustrate a sample python code to extract the relevant column data from the SQL database for hypertension.

```

23 def load_datasets():
24     conn = sqlite3.connect(r"C:\Users\jeffr\Desktop\Capstone Project (Revised)\databases\chronic_diseases.sqlite3")
25     hypertension_df = pd.read_sql("SELECT bmi, dbp, sbp, outcome FROM hypertension", conn)
26
27     # create a set of copy of diabetes dataset
28     copy_hypertension_df = hypertension_df.copy()
29
30     conn.close()
31     return copy_hypertension_df

```

Figure 4.1.2b: A sample code to extract the relevant column data from SQL database

Figure 4.1.2c illustrate the output of the relevant column data extracted from the SQL database via Python code for hypertension.

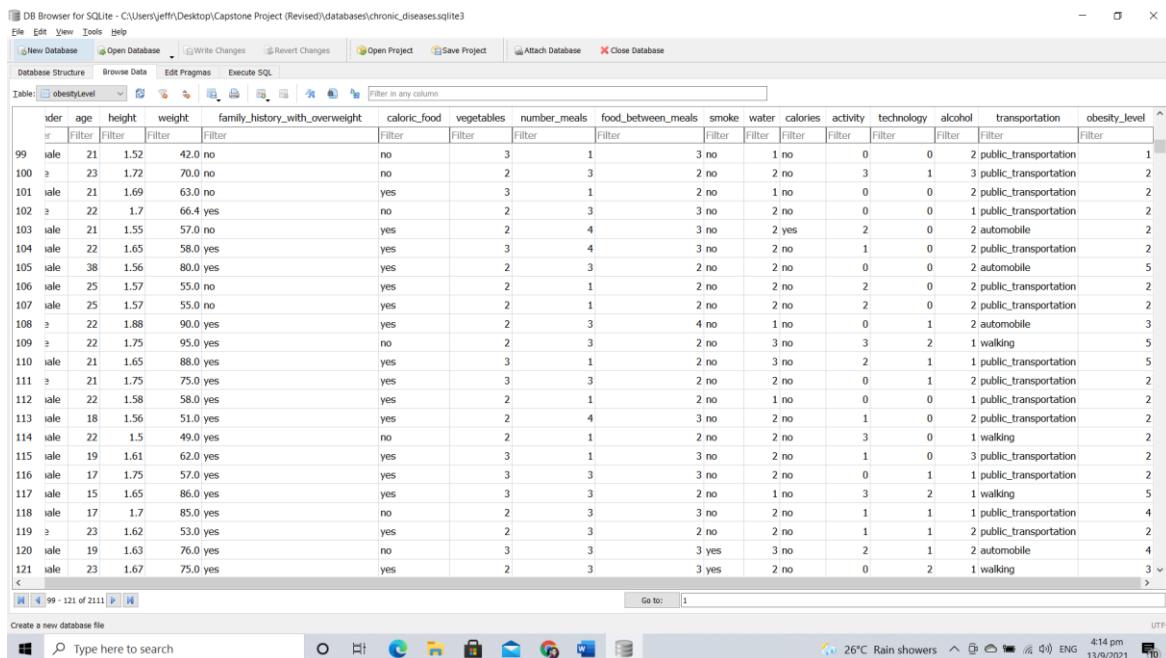
bp_df - DataFrame				
Index	bmi	dbp	sbp	outcome
0	27	85	120	REGULAR
1	18	81	100	REGULAR
2	27	95	106	REGULAR
3	24	80	120	REGULAR
4	25	83	126	REGULAR
5	22	81	170	HYPERTENSION
6	25	88	130	REGULAR
7	25	96	130	REGULAR
8	25	85	142	HYPERTENSION
9	25	81	166	REGULAR
10	25	81	146	REGULAR
11	27	80	136	REGULAR
12	19	84	110	REGULAR

Figure 4.1.2c: An output of the relevant column data extracted from the SQL database via Python Code

### 4.1.3 Obesity

The obesity (or high blood pressure) dataset was downloaded from **neuraldesigner** and the **original** dataset consists of the attributes of the gender, age, height, weight, family history with overweight, caloric food, vegetables, number meals, food between meals, smoke, water, calories, activity, technology, alcohol, transportation, and the class outcomes of the obesity level.

Figure 4.1.3a shows the obesity dataset (original) that was stored into SQLite3 database via Python.



The screenshot shows the DB Browser for SQLite application interface. The title bar reads "DB Browser for SQLite - C:\Users\jeff\Desktop\Capstone Project (Revised)\databases\chronic\_diseases.sqlite". The main window displays a table named "obesityLevel" with 211 rows of data. The columns are: id, age, height, weight, family\_history\_with\_overweight, caloric\_food, vegetables, number\_meals, food\_between\_meals, smoke, water, calories, activity, technology, alcohol, transportation, and obesity\_level. The data includes various values such as age (e.g., 21, 23, 25), height (e.g., 1.52, 1.72, 1.69), weight (e.g., 42.0, 70.0, 63.0), and obesity levels (e.g., 1, 2, 3). The bottom status bar shows "99 - 121 of 2111" and "UTF-8".

	id	age	height	weight	family_history_with_overweight	caloric_food	vegetables	number_meals	food_between_meals	smoke	water	calories	activity	technology	alcohol	transportation	obesity_level
99	male	21	1.52	42.0	no	no	3	1	3 no	1 no	0	0	2	public_transportation	1		
100	♂	23	1.72	70.0	no	no	2	3	2 no	2 no	3	1	3	public_transportation	2		
101	male	21	1.69	63.0	no	yes	3	1	2 no	1 no	0	0	2	public_transportation	2		
102	♂	22	1.7	66.4	yes	no	2	3	3 no	2 no	0	0	1	public_transportation	2		
103	male	21	1.55	57.0	no	yes	2	4	3 no	2 yes	2	0	2	automobile	2		
104	male	22	1.65	58.0	yes	yes	3	4	3 no	2 no	1	0	2	public_transportation	2		
105	male	38	1.56	80.0	yes	yes	2	3	2 no	2 no	0	0	2	automobile	5		
106	male	25	1.57	55.0	no	yes	2	1	2 no	2 no	2	0	2	public_transportation	2		
107	male	25	1.57	55.0	no	yes	2	1	2 no	2 no	2	0	2	public_transportation	2		
108	♂	22	1.88	90.0	yes	yes	2	3	4 no	1 no	0	1	2	automobile	3		
109	♂	22	1.75	95.0	yes	no	2	3	2 no	3 no	3	2	1	walking	5		
110	male	21	1.65	88.0	yes	yes	3	1	2 no	3 no	2	1	1	public_transportation	5		
111	♂	21	1.75	75.0	yes	yes	3	3	2 no	2 no	0	1	2	public_transportation	2		
112	male	22	1.58	58.0	yes	yes	2	1	2 no	1 no	0	0	1	public_transportation	2		
113	male	18	1.56	51.0	yes	yes	2	4	3 no	2 no	1	0	2	public_transportation	2		
114	male	22	1.5	49.0	yes	no	2	1	2 no	2 no	3	0	1	walking	2		
115	male	19	1.61	62.0	yes	yes	3	1	3 no	2 no	1	0	3	public_transportation	2		
116	male	17	1.75	57.0	yes	yes	3	3	3 no	2 no	0	1	1	public_transportation	2		
117	male	15	1.65	86.0	yes	yes	3	3	2 no	1 no	3	2	1	walking	5		
118	male	17	1.7	85.0	yes	no	2	3	3 no	2 no	1	1	1	public_transportation	4		
119	♂	23	1.62	53.0	yes	yes	2	3	2 no	2 no	1	1	2	public_transportation	2		
120	male	19	1.63	76.0	yes	no	3	3	3 yes	3 no	2	1	2	automobile	4		
121	male	23	1.67	75.0	yes	yes	2	3	3 yes	2 no	0	2	1	walking	3	>	

Figure 4.1.3a: Obesity Dataset (Original)

Table 4.1.3a shows the description of the attributes for the obesity dataset (original)

Attribute	Description
gender	Male or Female
age	Age (in years)
height	Measure the height of an individual (in m)
weight	Measure the weight of an individual (in kg)
family_history_ow	To check whether if there is a family history with overweight. (0 = No and 1= Yes)

## C3879C CAPSTONE PROJECT

caloric_food	Does an individual consumed high caloric food? (0 = No and 1 = Yes)
vegetables	How frequent consumption of vegetables? (1 = Not always, 2 = Frequently, 3 = Often)
number_meals	What is the number of meals per day? (1, 2, 3, or 4)
food_between_meals	How often do you consume foods between meals? (1 = Not always, 2 = Sometimes, 3 = Frequently, 4 = Always)
smoke	Do you smoke? (0 = No, 1 = Yes)
water	How often do you drink water daily? (1 = Not always, 2 = Sometimes, 3 = Frequently)
calories	Do you often monitor your calories? (0 = No, 1 = Yes)
activity	How often do you exercise? (0 = Not always, 1 = Sometimes, 2 = Frequently, 3 = Always)
technology	How often do you use your electronic devices? (0 = Sometimes, 1 = Frequently, 2 = Always)
alcohol	Do you drink alcohol? (1 = Not always, 2 = Sometimes, 3 = Frequently, 4 = Always)
transportation	Types of transportation used: Automobile, Motorbike, Bike, Public Transportation, Walking
obesity_level	1 = Insufficient Weight, 2 = Normal Weight, 3 = Overweight Level I, 4 = Overweight Level II, 5 = Obesity Type I, 6 = Obesity Type II, 7 = Obesity Type III

Table 4.1.3a: Description of the Attributes for Obesity Dataset (Original)

As such, all the relevant attributes were used for training phases except for transportation.

## C3879C CAPSTONE PROJECT

Figure 4.1.3b illustrate a sample python code to extract the relevant column data from the SQL database for obesity.

```

24 def load_datasets():
25     conn = sqlite3.connect(r"C:\Users\jeffr\Desktop\Capstone Project (Revised)\databases\chronic_diseases.sqlite3")
26     obesity_df = pd.read_sql("SELECT gender, age, height, weight, family_history_with_overweight, \
27                             caloric_food, vegetables, number_meals, food_between_meals, smoke, water, \
28                             calories, activity, technology, alcohol, obesity_level FROM obesityLevel", conn)
29
30     # create a set of copy of diabetes dataset
31     copy_obesity_df = obesity_df.copy()
32
33     # rename the target column
34     copy_obesity_df.rename(columns = {'obesity_level': 'outcome', 'family_history_with_overweight': 'family_history_ow'}, inplace = True)
35
36     # manipulate some of the target values
37     copy_obesity_df['outcome'].replace([(3, 4), (5, 6, 7)], (3, 4), inplace = True)
38
39     conn.close()
40     return copy_obesity_df

```

Figure 4.1.3b: A sample code to extract the relevant column data from SQL database

Figure 4.1.3c illustrate the output of the relevant column data extracted from the SQL database via Python code for obesity.

Index	gender	age	height	weight	family_history	caloric_food	vegetable	number_meals	between_meals	smoke	water	calories	activity	technology	alcohol	outcome
0	Female	21	1.62	64	yes	no	2	3	2	no	2	no	0	1	1	2
1	Female	21	1.52	56	yes	no	3	3	2	yes	3	yes	3	0	2	2
2	Male	23	1.8	77	yes	no	2	3	2	no	2	no	2	1	3	2
3	Male	27	1.8	87	no	no	3	3	2	no	2	no	2	0	3	3
4	Male	22	1.78	89.8	no	no	2	1	2	no	2	no	0	0	2	3
5	Male	29	1.62	53	no	yes	2	3	2	no	2	no	0	0	2	2
6	Female	23	1.5	55	yes	yes	3	3	2	no	2	no	1	0	2	2
7	Male	22	1.64	53	no	no	2	3	2	no	2	no	3	0	2	2
8	Male	24	1.78	64	yes	yes	3	3	2	no	2	no	1	1	3	2
9	Male	22	1.72	68	yes	yes	2	3	2	no	2	no	1	1	1	2
10	Male	26	1.85	105	yes	yes	3	3	3	no	3	no	2	2	2	4
11	Female	21	1.72	80	yes	yes	2	3	3	no	2	yes	2	1	2	3
12	Male	22	1.66	55	no	no	2	3	2	no	2	no	0	0	2	2

Figure 4.1.3c: An output of the relevant column data extracted from the SQL database via Python Code

## 4.2 Data Exploration on the Selected Attributes

### 4.2.1 Diabetes

#### 4.2.1.1 Histogram

A **histogram** of each selected attributes in the **Diabetes DataFrame** was created by calling the **hist () function** to get an idea of the variable distribution. Figure 4.2.1.1 shows the histogram subplot and there are only **two numeric input** and **two class labels** in the diabetes DataFrame.

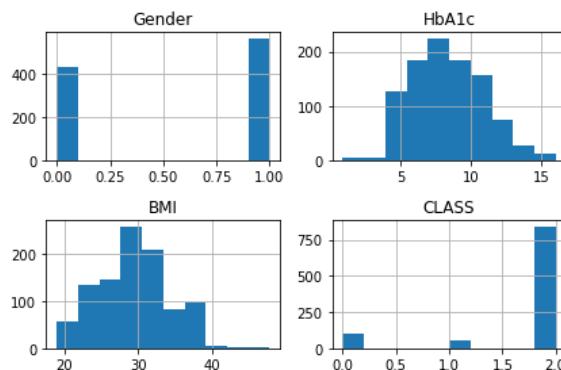


Figure 4.2.1.1: A histogram subplot of the diabetes

Based on the above histogram subplot, it would be expected to scale the distributions to the same range to be useful for modelling algorithms at the later stage, and perhaps the use of scaling techniques at the modelling stage.

#### 4.2.1.2 Heatmap and Pair-plot

Figure 4.2.1.2a shows the **heatmap** of the selected attributes in the Diabetes DataFrame and the objective is to visualize the correlation of each pair of attributes and the outcome variable.

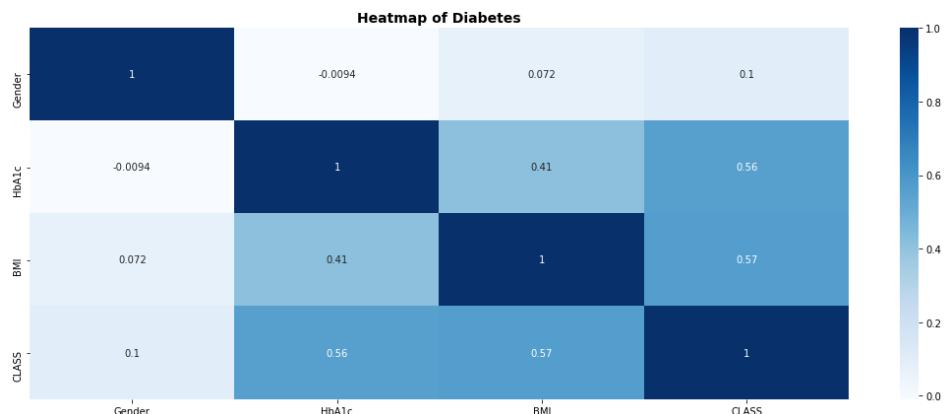


Figure 4.2.1.2a: A heatmap of the selected attributes in diabetes DataFrame

## C3879C CAPSTONE PROJECT

Based on the heatmap shown above, the darker colour indicates more correlation between variables. Even though, both HbA1c and BMI shows some moderate and positive correlation with the class outcome variables, but one must note that this is a classification context based. Therefore, there will be no significant impact on the overall performance during the algorithms analysis. In addition, there is also some weak positive correlation between the HbA1c and BMI variables, but the impact on the overall algorithm's performance is generally very small.

Figure 4.2.1.2b shows the pair-plot of the selected attributes in the Diabetes DataFrame and has the common objective as the Heatmap to visualize the correlation between variables.

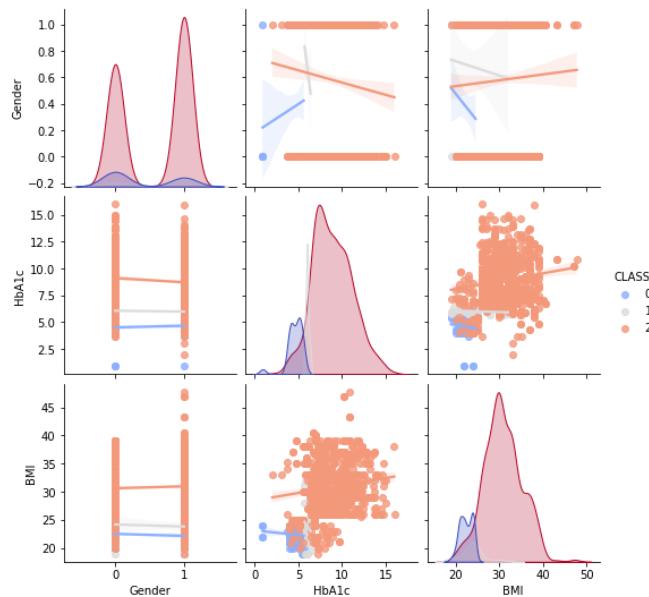


Figure 4.2.1.2b: A pair-plot of the selected attributes in the diabetes DataFrame

From the pair-plot, it can be observed that there is some weak correlation between the HbA1c and BMI attribute which can be proven from the above heatmap.

## 4.2.2 Hypertension

### 4.2.2.1 Histogram

Similar to the diabetes, a **histogram** of each selected attributes in the **Hypertension DataFrame** was created by calling the **hist () function** to get an idea of the variable distribution. Figure 4.2.2.1 shows the histogram subplot and there are only **three numeric input** and **one class label** in the hypertension DataFrame.

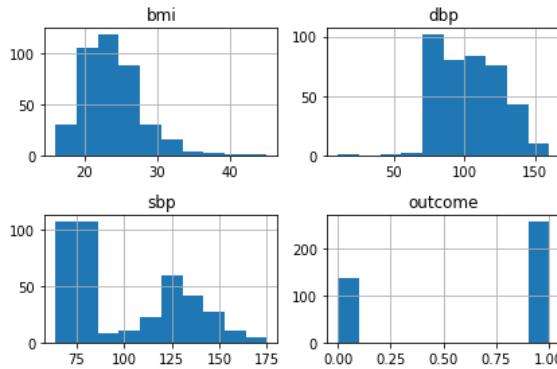


Figure 4.2.2.1: A histogram subplot of the Hypertension

Based on the above histogram subplot, one could be observed the BMI attribute is moderately skewed to the right, whereas the DBP attribute is also moderately skewed to the left. Furthermore, the SBP attribute consists of the bimodal distribution. As such, it would be expected to scale the distributions to the same range to be useful for algorithms analysis, and perhaps some power transforms involved at the modelling stage.

### 4.2.2.2 Heatmap and Pair-plot

Similar to the diabetes, Figure 4.2.2.2a shows the **heatmap** of the selected attributes in the Hypertension DataFrame and the objective is to visualize the correlation of each pair of attributes and the outcome variable.

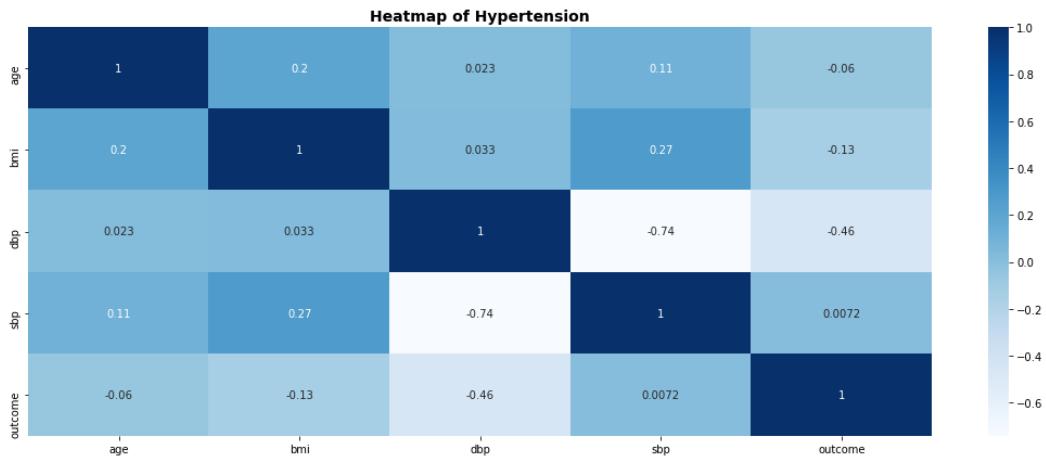


Figure 4.2.2.2a: A heatmap of the selected attributes in hypertension DataFrame

Based on the heatmap shown above, the darker colour indicates more correlation between variables. As such, no significant correlation was found between variables.

Figure 4.2.2.2b shows the pair-plot of the selected attributes in the Hypertension DataFrame and has the common objective as the Heatmap to visualize the correlation between variables.

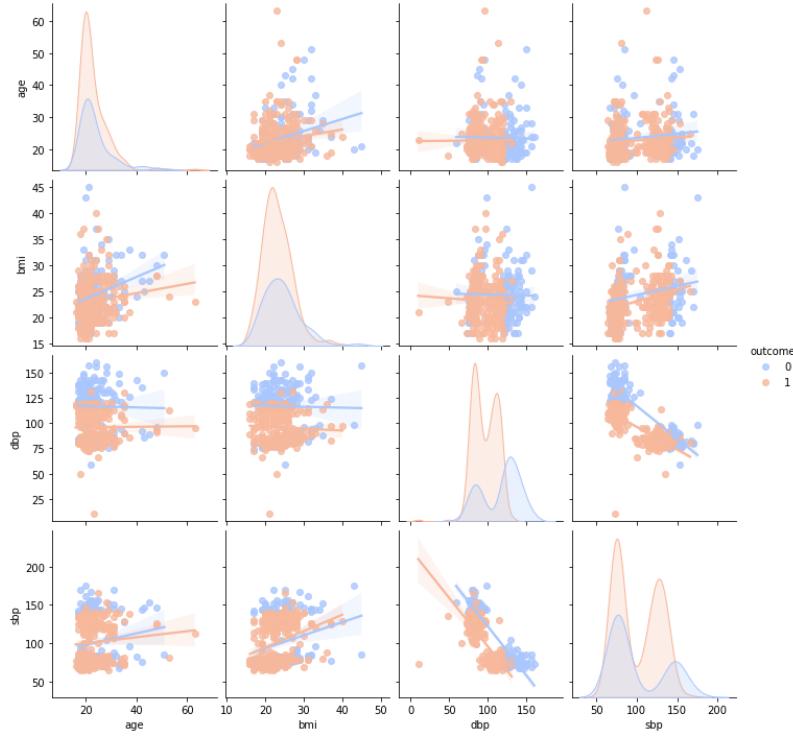


Figure 4.2.2.2b: A pair-plot of the selected attributes in the hypertension DataFrame

From the pair-plot, one can be observed that there is no significant correlation between the variables which can be proven from the above heatmap.

## 4.2.3 Obesity

### 4.2.3.1 Histogram

Similar to both diabetes and hypertension, a **histogram** of each selected attributes in the **Obesity DataFrame** was created by calling the **hist () function** to get an idea of the variable distribution. Figure 4.2.3.1 shows the histogram subplot and there are only **3 numeric input, 13 categorical variables and 1 class label** in the obesity DataFrame.

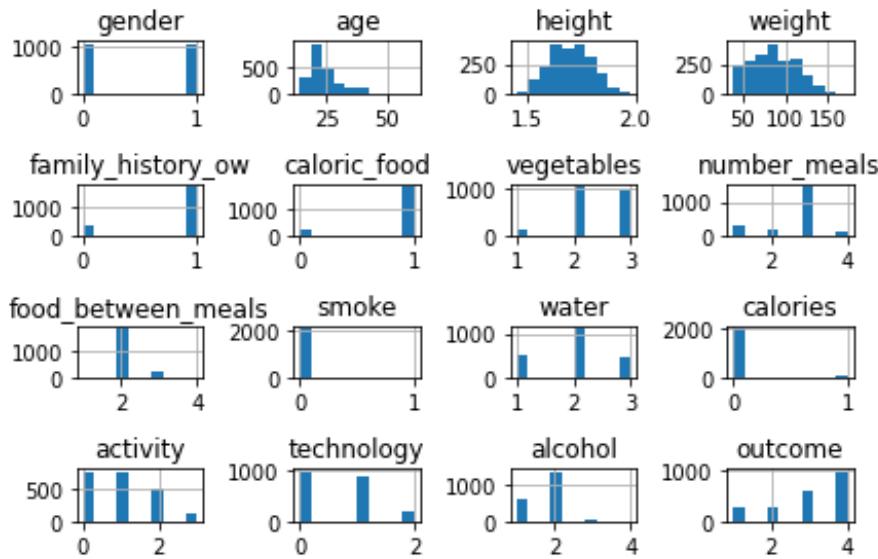


Figure 4.2.3.1: A histogram subplot of the Obesity

Based on the above histogram subplot, one could observe both the age and weight attributes are moderately skewed to the right, whereas the height attribute represents the Gaussian distribution. As such, it would be expected to scale the distributions to the same range to be useful for algorithms analysis, and perhaps some power transforms involved at the modelling stage.

### 4.2.3.2 Heatmap and Pair-plot

Similar to the both diabetes and hypertension, Figure 4.2.3.2a shows the **heatmap** of the selected attributes in the Obesity DataFrame and the objective is to visualize the correlation of each pair of attributes and the outcome variable.

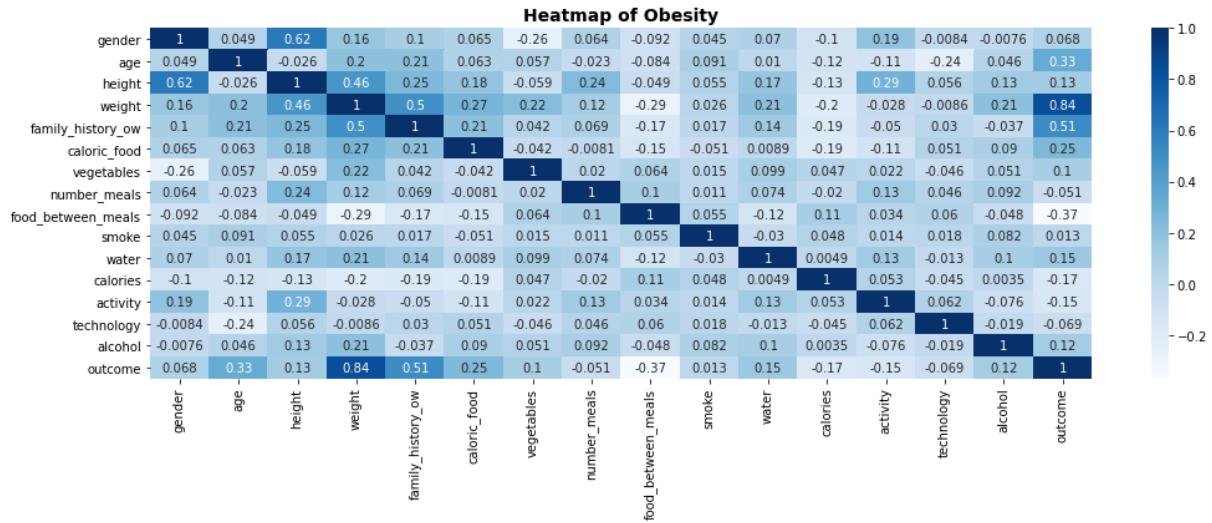


Figure 4.2.3.2a: A heatmap of the selected attributes in obesity DataFrame

Based on the heatmap shown above, the darker colour indicates more correlation between variables. Even though, the weight attribute shows moderately strong positive correlation between with the class outcomes, but one must note that this is part of the classification context-based. As such, there will be no significant impact on the overall performance of the algorithms analysis. Furthermore, there are some weak and positive correlation between the height, weight, and family history with overweight attributes. Generally, the impact on the algorithm's performance is very small.

Figure 4.2.3.2b shows the pair-plot of the selected attributes in the Obesity DataFrame and has the common objective as the Heatmap to visualize the correlation between variables.

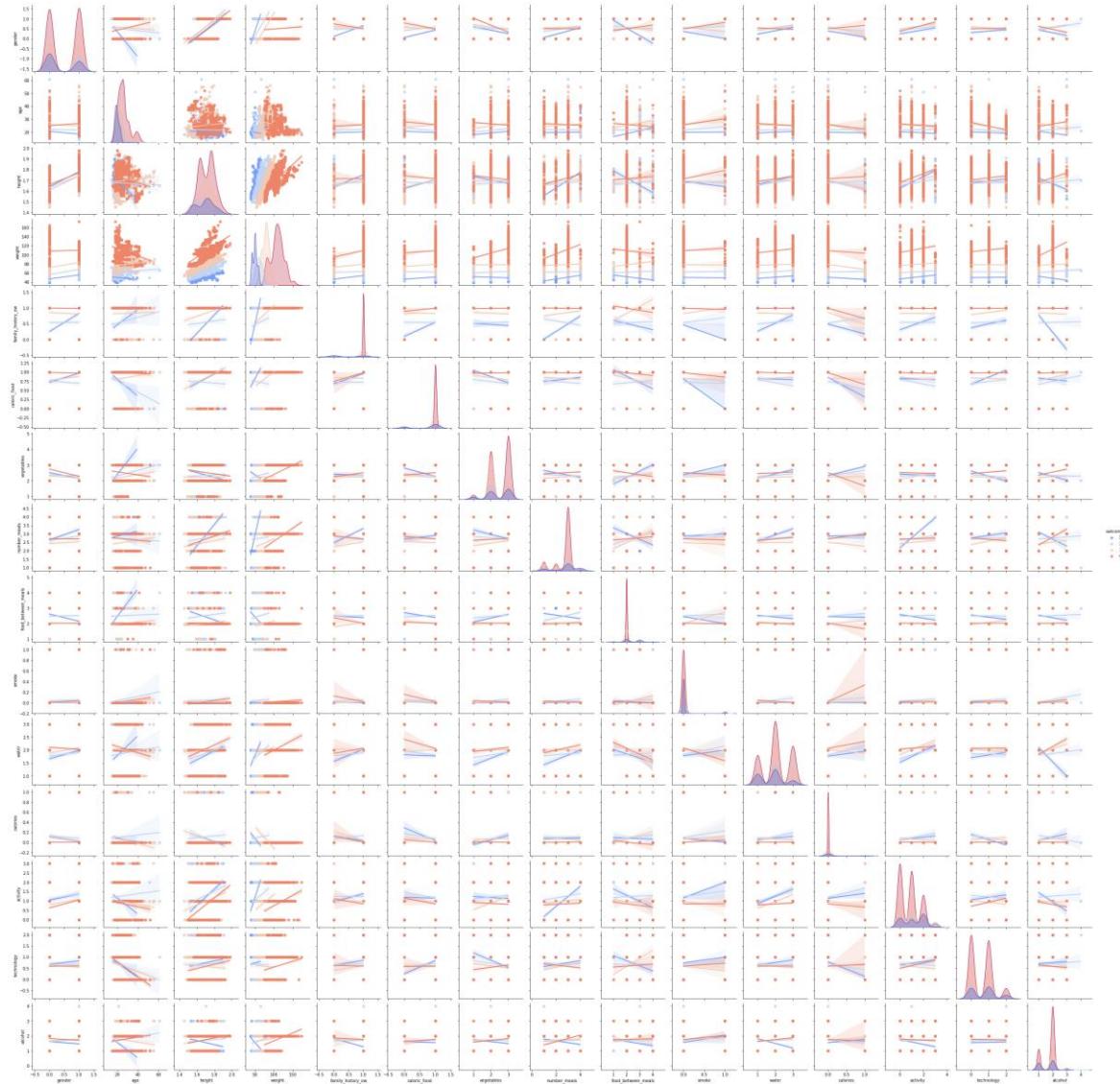


Figure 4.2.3.2b: A pair-plot of the selected attributes in the obesity DataFrame

From the pair-plot, it can observe that there is a weak and positive correlation between the height, weight, and family history with overweight variables which can be proven from the above heatmap.

## 4.3 Data Cleansing

### 4.3.1 Investigate Missing Data

Prior to the algorithms analysis, it is imperative to investigate for missing data during the pre-processing of the dataset as not all machine learning algorithms support missing values and might have a chance of leading to inaccurate and bias outcomes.

As such, the Seaborn heatmap was created for all three datasets to visualize the missing row data.

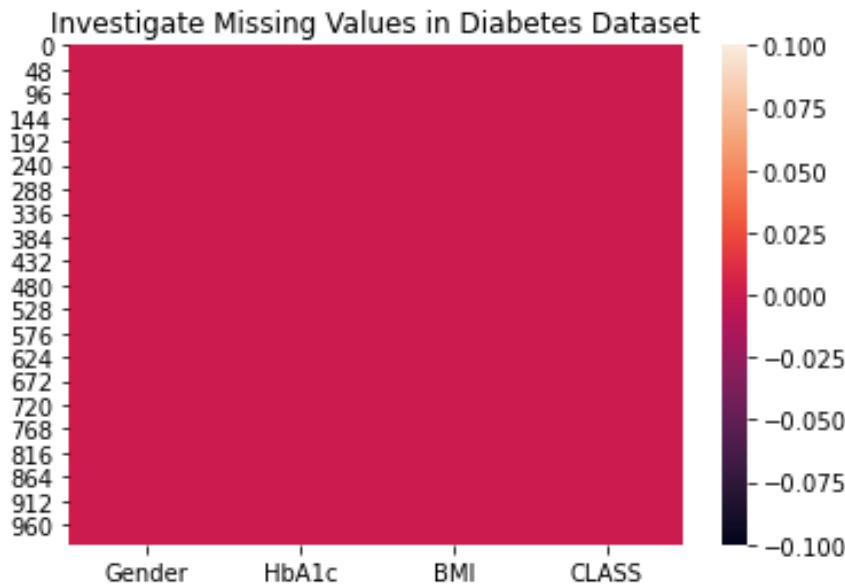


Figure 4.3.1a: Visualize the Missing Row Data in Diabetes Dataset

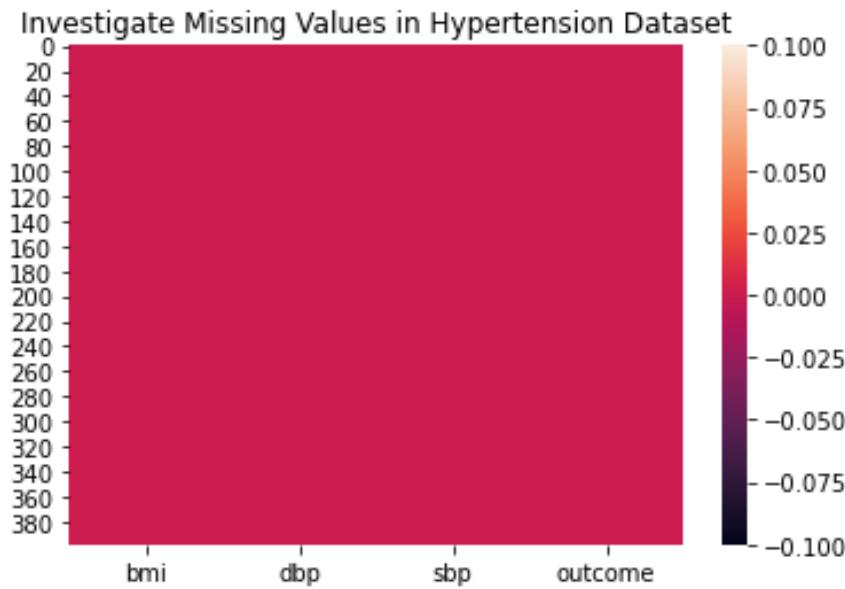


Figure 4.3.1b: Visualize the Missing Row Data in Hypertension Dataset



Figure 4.3.1c: Visualize the Missing Row Data in Obesity Dataset

From the above observations, there were no missing row data found in all three datasets.

### 4.3.2 Removal of Potential Outlier Detection using Isolation Forest

Isolation Forest (iForest) is the anomaly detection algorithm that identifies anomalies using isolation. As such, the iForest was used to identified the potential outlier data in the obesity, diabetes, and the hypertension datasets.

Basically, the isolation forest works by creating an ensemble of isolation forest trees for each dataset where outliers were as instances with short average length in the trees. The tree is then recursively created by dividing the dataset until all instances are isolated or specific tree height is achieved.

In this project, the isolation forest was implemented using Scikit-Learn Python library version 0.24.2.

## C3879C CAPSTONE PROJECT

After the potential outliers detected by the isolation forest, the outlier data from each dataset were eliminated and the remaining data were used for further modelling algorithms.

Figure 4.3.2a illustrates a sample code snippet for diabetes.

```

200 ##### Removing Outliers using Isolation Forest #####
201 # summarize the number of the dataset before removing outliers
202 X_features = diab_df2.drop("CLASS", axis = 1)
203 y_target = diab_df2['CLASS']
204
205 # define and identify outliers in the dataset
206 iso = IsolationForest(contamination = 0.06, random_state = 1)
207 iso_pred = iso.fit_predict(X_features)
208
209 # to obtain the anomaly score
210 iso_scores = iso.decision_function(X_features)
211
212 # display the anomaly score
213 sns.distplot(iso_scores, kde = False, color = 'green', bins = 50)
214 plt.title("Distribution of the Anomaly Scores using Isolation Forest for Diabetes", fontweight = 'bold', fontsize = 14)
215 plt.xlabel("Average Path Length", fontweight = 'bold', fontsize = 12)
216 plt.ylabel("Count", fontweight = 'bold', fontsize = 12)
217 plt.show()
218
219 # identify the outliers using the selected threshold value
220 outliers_index = np.where(iso_scores < -0.01)[0]
221
222 # filter out the outliers from the original dataset based on the outliers index
223 X_features_iso = X_features[~X_features.index.isin(outliers_index)]
224 y_target_iso = y_target[~y_target.index.isin(outliers_index)]
225
226 # convert the X_features_iso and y_target_iso into dataframe
227 X_features = pd.DataFrame(data = X_features_iso, columns = X_features.columns)
228 y_target = pd.Series(data = y_target_iso, name = 'CLASS')

```

Figure 4.3.2a: A sample code for outlier's detection using iForest for Diabetes

Figure 4.3.2b illustrates the distribution of the anomaly scores for diabetes.

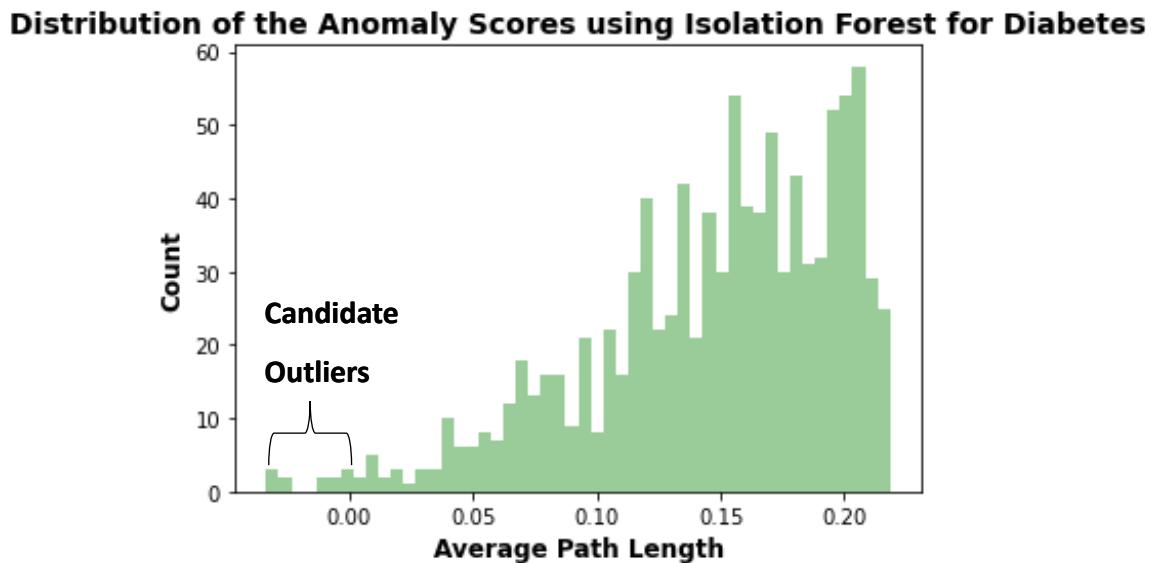


Figure 4.3.2b: The distribution of the anomaly scores using iForest for Diabetes

Figure 4.3.2c illustrates a sample code snippet for hypertension.

```

193 ##### Removing Outliers using Isolation Forest #####
194 # summarize the number of the dataset before removing outliers
195 X_features = bp_df.drop("outcome", axis = 1)
196 y_target = bp_df['outcome']
197
198 iso = IsolationForest(contamination = 0.01, random_state = 1)
199 iso_predict = iso.fit_predict(X_features)
200
201 # to obtain the anomaly score
202 iso_scores = iso.decision_function(X_features)
203
204 # display the anomaly score
205 sns.distplot(iso_scores, kde = False, bins = 50)
206 plt.title("Distribution of the Anomaly Scores using Isolation Forest for Hypertension", fontweight = 'bold', fontsize = 14)
207 plt.xlabel("Average Path Length", fontweight = 'bold', fontsize = 12)
208 plt.ylabel("Count", fontweight = 'bold', fontsize = 12)
209 plt.show()
210
211 # identify the outliers using the selected threshold value
212 outliers_index = np.where(iso_scores < -0.01)[0]
213
214 # filter out the outliers from the original dataset based on the outliers index
215 x_features_iso = X_features[~X_features.index.isin(outliers_index)]
216 y_target_iso = y_target[~y_target.index.isin(outliers_index)]
217
218 # convert the x_features_iso and y_target_iso into dataframe
219 X_features = pd.DataFrame(data = x_features_iso, columns = X_features.columns)
220 y_target = pd.Series(data = y_target_iso, name = 'outcome')
221

```

Figure 4.3.2c: A sample code for outlier's detection using iForest for Hypertension

Figure 4.3.2d illustrates the distribution of the anomaly scores for hypertension.

**Distribution of the Anomaly Scores using Isolation Forest for Hypertension**

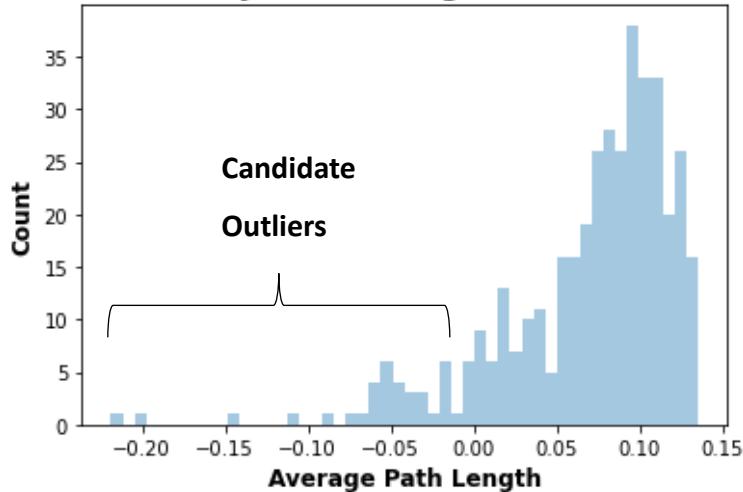


Figure 4.3.2d: The distribution of the anomaly scores using iForest for Hypertension

## C3879C CAPSTONE PROJECT

Figure 4.3.2e illustrates a sample code snippet for obesity.

```

216 ##### Removing Outliers using Isolation Forest #####
217 # summarize the number of the dataset before removing outliers
218 X_features = obs_df.drop("outcome", axis = 1)
219 y_target = obs_df['outcome']
220
221 # define and identify outliers in the dataset
222 iso = IsolationForest(contamination = 0.1, random_state = 1)
223 iso_pred = iso.fit_predict(X_features)
224
225 # to obtain the anomaly score
226 iso_scores = iso.decision_function(X_features)
227
228 # display the anomaly score
229 sns.distplot(iso_scores, kde= False, color = 'crimson', bins = 100)
230 plt.title("Distribution of the Anomaly Scores using Isolation Forest for Obesity", fontweight = 'bold', fontsize = 14)
231 plt.xlabel("Average Path Length", fontweight = 'bold', fontsize = 12)
232 plt.ylabel("Count", fontweight = 'bold', fontsize = 12)
233 plt.show()
234
235 # identify the outliers using the selected threshold value
236 outliers_index = np.where(iso_scores < -0.01)[0]
237 |
238 # filter out the outliers from the original dataset based on the outliers index
239 x_features_iso = X_features[~X_features.index.isin(outliers_index)]
240 y_target_iso = y_target[~y_target.index.isin(outliers_index)]
241
242 # convert the x_features_iso and y_target_iso into dataframe
243 X_features = pd.DataFrame(data = x_features_iso, columns = X_features.columns)
244 y_target = pd.Series(data = y_target_iso, name = 'outcome')
245
246 after_rOutlier = pd.concat([X_features, y_target], axis = 1)
247 boxplot(after_rOutlier)

```

Figure 4.3.2e: A sample code for outlier's detection using iForest for Obesity

Figure 4.3.2f illustrates the distribution of the anomaly scores for obesity.

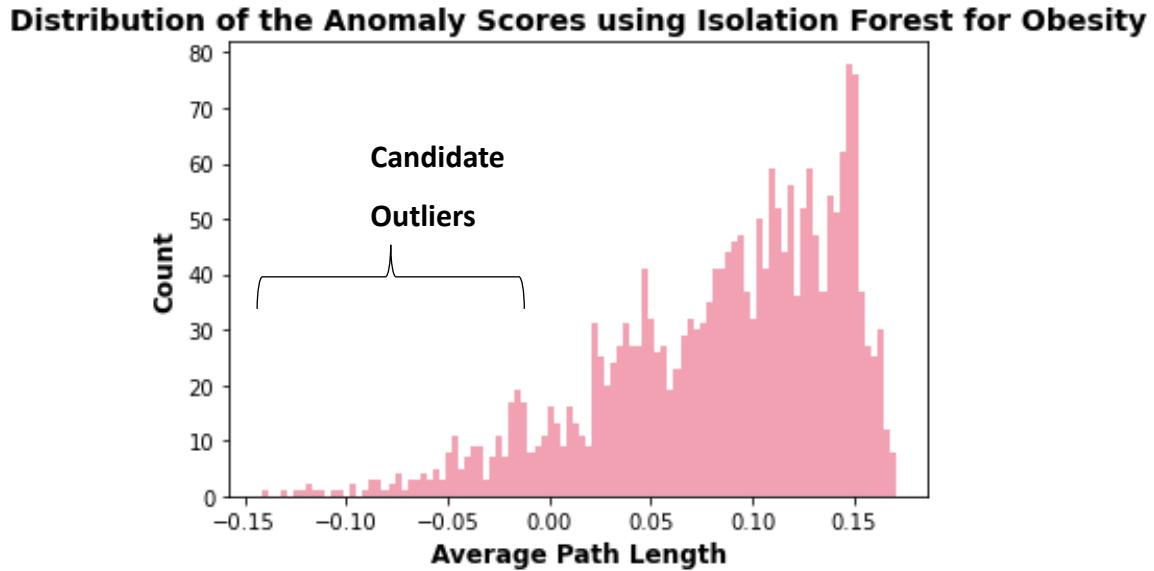


Figure 4.3.2f: The distribution of the anomaly scores using iForest for Obesity

Basically, the outliers with the lower anomaly scores are likely to be outliers. As such, the outlier threshold on the score will be set by the user preference to differentiate between the outlier and inlier. In this case, the outlier threshold will be set as below -0.01 for all data aspects.

## 4.4 Resampling Techniques for Imbalanced Training Sets

Prior to the resampling and algorithms analysis, all the respective pre-processed datasets were then divided into 70% training set and 30% testing set. Figure 4.4a illustrates a sample code snippet of splitting the pre-processed datasets and was applied to all the respective datasets.

```
228     X_train, X_test, y_train, y_test = train_test_split(X_features,
229                                     y_target,
230                                     test_size = 0.30,
231                                     random_state = 42,
232                                     stratify = y_target)
233
```

Figure 4.4a: A sample code of dividing the pre-processed data into 70% training set and 30% testing set

After splitting the datasets, it is also imperative to balance the datasets prior to the training phase. This is because if the imbalanced training data is not treated promptly prior to the training phase, this will degrade the performance of the classifier models. As such, most of the predictions will correspond to the majority class and treat the minority class features as noise in the data and resulting high bias in the models.

Therefore, resampling is one of the most commonly approach to achieve the aim and deal with the imbalanced training datasets.

In this case, the SMOTE-TOMEK resampling technique was chosen as a resampling technique on the imbalanced training datasets because the aims of this technique are to clean the overlapping data points for each of the classes distributed in the sample space. If oversampling is done by SMOTE alone, the class clusters may invade each other space and resulting in the classifier model being overfitting. Likewise, TOMEK links are the opposite paired samples that are closet neighbours to each other. Therefore, the majority of the class observations from these links are removed. Tomek links are applied to oversampled minority class samples performed by SMOTE. All in all, instead of removing the observations only from the majority class, only remove both class observations from TOMEK links.

One must note that only SMOTE-TOMEK technique was applied on the training set as testing set was not part during the training phase.

## C3879C CAPSTONE PROJECT

SMOTE-TOMEK technique was implemented using the Imbalanced-learn python library (V0.8.0). Figure 4.4b illustrates a sample code snippet to balance the training imbalanced datasets for all aspects.

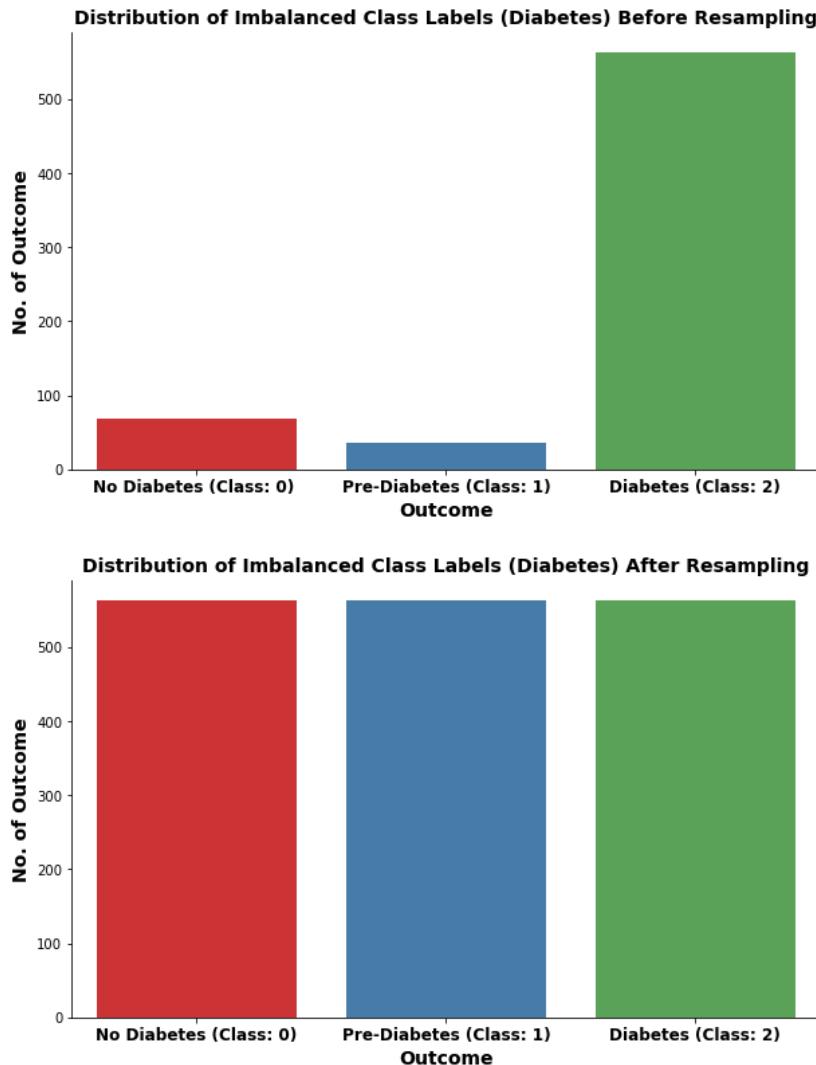
```

75  def balanced_datasets(X, y):
76      oversample_SMT = SMOTETomek(random_state = 1)
77      X_smt, y_smt = oversample_SMT.fit_resample(X, y)
78      return X_smt, y_smt
79
80 ###### call the function to balanced the datasets
81 X_smt, Y_smt = balanced_datasets(X_train, y_train)
82

```

Figure 4.4b: A sample code to balance the training imbalanced datasets using SMOTE-TOMEK technique

Figure 4.4c illustrates the initial minority class distributions **before** and **after** resampling for diabetes training datasets.

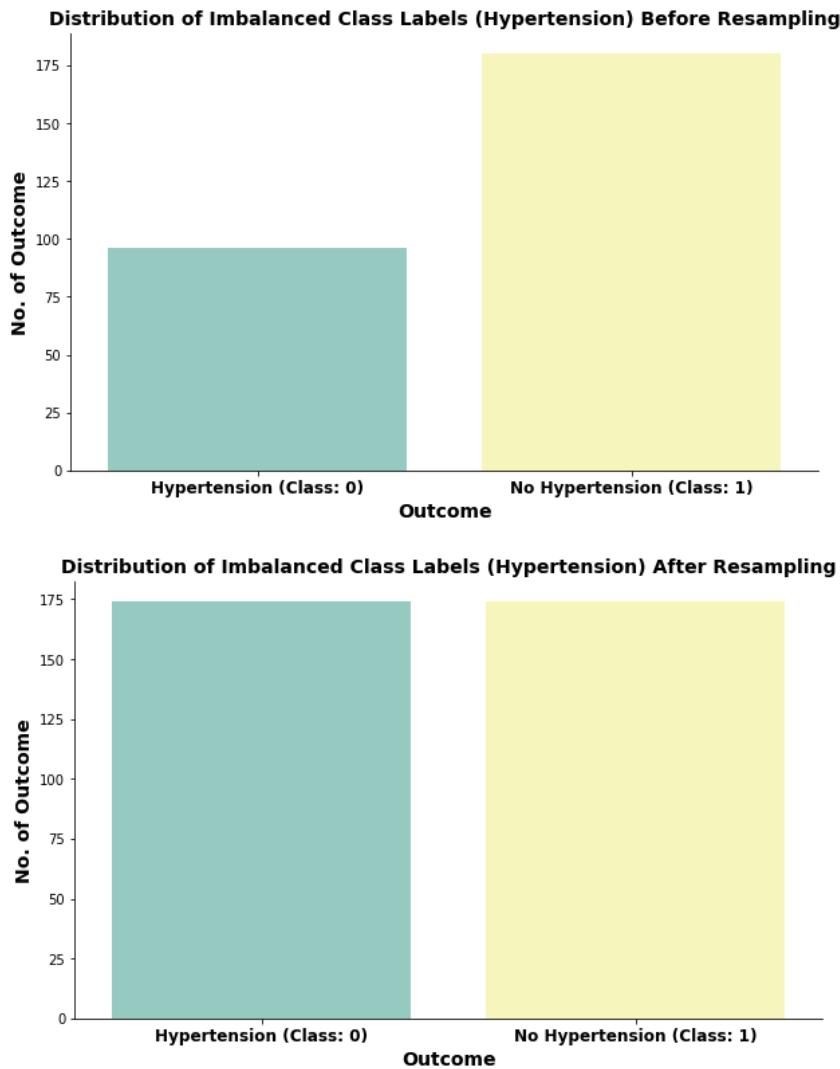


```
Before Resampling for Diabetes using SMOTE-Tomek:
0      68
1      36
2     563
Name: CLASS, dtype: int64

After Resampling for Diabetes using SMOTE-Tomek:
0     563
1     563
2     563
Name: CLASS, dtype: int64
```

Figure 4.4c: Before and After resampling for diabetes training datasets

Figure 4.4d illustrates the initial minority class distributions **before** and **after** resampling for hypertension training datasets.

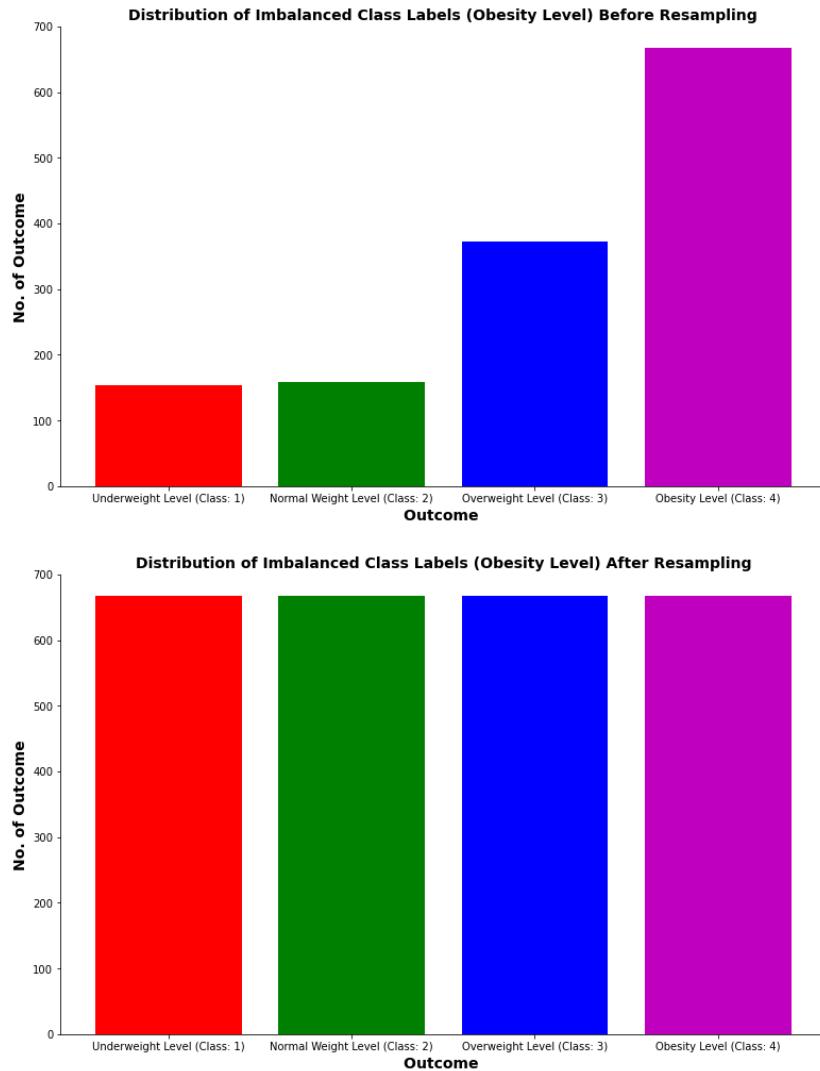


```
Before Resampling for Hypertension using SMOTE-Tomek:
0      96
1     180
Name: outcome, dtype: int64

After Resampling for Hypertension using SMOTE-Tomek:
0     174
1     174
Name: outcome, dtype: int64
```

Figure 4.4d: Before and After resampling for hypertension training datasets

Figure 4.4e illustrates the initial minority class distributions **before** and **after** resampling for obesity training datasets.



C3879C CAPSTONE PROJECT

---

```
Before Resampling for Obesity Level using SMOTE-Tomek:  
Class = 1, n = 158 (11.695%)  
Class = 2, n = 153 (11.325%)  
Class = 3, n = 373 (27.609%)  
Class = 4, n = 667 (49.371%)  
  
After Resampling for Obesity Level using SMOTE-Tomek:  
Class = 1, n = 667 (25.000%)  
Class = 2, n = 667 (25.000%)  
Class = 3, n = 667 (25.000%)  
Class = 4, n = 667 (25.000%)
```

Figure 4.4e: Before and After resampling for obesity training datasets

## 4.5 Baseline Model Selection and Evaluation via Shuffle-Split CV

Prior to the model building, the potential outlier data for the respective datasets (i.e., diabetes, hypertension, and obesity) were eliminated by the Isolation Forest algorithm as outlier detection, and the SMOTE-TOMEK was implemented to achieve the balance training datasets. Followed by, the classifier machine learning algorithms such as random forest, gradient boosting, and extra trees were implemented using Scikit-learn V0.24.2 and were used to learn from the respective training datasets.

Furthermore, the **baseline** of the classifier machine learning algorithms were built with the **default parameters** to investigate how well does each algorithm performed with the respective medical disease training datasets. Figure 4.5a illustrates a sample code of the baseline models.

```
311 # initialize the ensemble classifier models with their default parameters and add them into a model list
312 classifier_models = [('Random Forest', RandomForestClassifier(random_state = 42)),
313 ('Gradient Boosting', GradientBoostingClassifier(random_state = 42)),
314 ('Extra Trees', ExtraTreesClassifier(random_state = 42))]
```

Figure 4.5a: A sample code of the baseline models using the default parameters

Prior to the shuffle-split cross-validation, a pipeline object was built by providing it with a list of steps. Each step represents a tuple containing a string of name and an instance of an estimator.

For the hypertension training set, the first step was created with a ‘scaler’, which is an instance of MinMaxScaler (), followed by a ‘power’, is an instance of PowerTransformer (), and the final step was created with a ‘model’, is an instance of the respective classifier algorithms as mentioned above. In other words, each pipeline defines one (or more) data preparation techniques such as MinMaxScaler or StandardScaler, PowerTransformer, and ends with a defined model that takes the transformed training data as input.

The main objective of using the Power Transformer () was to change the distribution to be more alike of the Gaussian. Basically, the Power Transformer () performed the Yeo-Johnson transformation and automatically determined the best parameters that applied based on the dataset, e.g., how to best make each variable more Gaussian. Furthermore, the transformer will also standardize the dataset as part of the transform. However, the power transform may

C3879C CAPSTONE PROJECT

---

make use of the log () function, which does not really work for zero values. As such, a MinMaxScaler () function was included to scaled the dataset prior to the power transformer.

Figure 4.5b illustrates a sample code snippet of the model pipeline for **hypertension** training set.

```
330     # define pipeline
331     pipe = Pipeline(steps = [('scaler', MinMaxScaler()), ('power', PowerTransformer()), ('model', models)])
```

Figure 4.5b: A sample code of the model pipeline for hypertension training set

The model pipeline for obesity and diabetes can be created a similar process and achieve the same objective as hypertension, such that the first step was created with a ‘scaler’ of StandardScaler (), followed by created with a ‘model’, of the respective classifier algorithms. Figure 4.5c illustrates a sample code snippet of the model pipeline for both **obesity** and **diabetes** training set.

```
348     # build pipeline
349     pipe = Pipeline(steps = [('scaler', StandardScaler()), ('model', models)])
```

Figure 4.5c: A sample code of the model pipeline for both obesity and diabetes training set

After the baseline model pipelines were implemented, the **Shuffle-Split Cross Validation** was implemented to assess the performance and general errors of the entire classification algorithms. This cross-validation technique has been widely used and is one of the preferred validation techniques in machine learning as compared to the conventional split-sample approach (i.e., train-test split).

**Shuffle-Split Cross-Validation** allows for control over the number of iterations independently of the training and test sizes and also allows for using the only part of the data in each iteration, by providing the train size and test size setting that don’t add up to one. Therefore, shuffle-split cross-validation provides a more reliable result for classification tasks. For this case, the train size and test size has been set as 80% and 20% respectively. One must note that the test used in the cross-validation was not originated from the testing set. In other words, the test used in the cross-validation was actually the % split from the training set itself. Figure 4.5d illustrate a set of the sample code to evaluate each

## C3879C CAPSTONE PROJECT

of the classification machine learning algorithms via shuffle-split cross validation with model pipeline.

**A) For Diabetes:**

```

82 def cross_validation_eval(model, X_smt, Y_smt):
83
84     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
85     cv_result = cross_validate(model,
86                               X_smt,
87                               Y_smt,
88                               cv = shuffle_split,
89                               n_jobs = 1,
90                               scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro'),
91                               return_train_score = True)
92
93
94     # compute the ROC scores
95     myscore = make_scorer(roc_auc_score, multi_class='ovo', needs_proba=True)
96     roc_scores = cross_validate(model,
97                               X_smt,
98                               Y_smt,
99                               cv = shuffle_split,
100                              n_jobs = 1,
101                              scoring = myscore,
102                              return_train_score = True)
103
104
105     return cv_result, roc_scores

342 for name, models in classifier_models:
343
344     # build pipeline
345     pipe = Pipeline(steps = [('scaler', StandardScaler()), ('model', models)])
346
347     # call the function to perform cross-validation evaluation
348     cvScores, rocScores = cross_validation_eval(pipe, X_smt, Y_smt)

```

**B) For Hypertension:**

```

81 def cross_validation_eval(pipe, X_smt, Y_smt):
82
83     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
84     cv_result = cross_validate(pipe,
85                               X_smt,
86                               Y_smt,
87                               cv = shuffle_split,
88                               n_jobs = 1,
89                               scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro', 'roc_auc'),
90                               return_train_score = True,
91                               return_estimator = True)
92
93
94     return cv_result

328 for name, models in classifier_models:
329
330     # define pipeline
331     pipe = Pipeline(steps = [('scaler', MinMaxScaler()), ('power', PowerTransformer()), ('model', models)])
332
333
334     # call the function to perform cross validation evaluation
335     cvScores = cross_validation_eval(pipe, X_smt, Y_smt)
336

```

**C) For Obesity**

```

95 def cross_validation_eval(model, X_smt, Y_smt):
96
97     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
98
99     # compute the accuracy, precision, recall and f1-scores
100    cv_result = cross_validate(model,
101                                X_smt,
102                                Y_smt,
103                                cv = shuffle_split,
104                                n_jobs = 1,
105                                scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro'),
106                                return_train_score = True)
107
108    # compute the ROC scores
109    myscore = make_scorer(roc_auc_score, multi_class='ovo', needs_proba=True)
110    roc_scores = cross_validate(model,
111                                X_smt,
112                                Y_smt,
113                                cv = shuffle_split,
114                                n_jobs = 1,
115                                scoring = myscore,
116                                return_train_score = True)
117
118    return cv_result, roc_scores

341 for name, models in classifier_models:
342
343     # build pipeline
344     pipe = Pipeline(steps = [('scaler', StandardScaler()), ('model', models)])
345
346     # call function to compute cross validation evaluation
347     cvScores, rocScores = cross_validation_eval(pipe, X_smt, Y_smt)
348

```

Figure 4.5d: A sample code of Shuffle-Split Cross Validation with Model Pipeline

## C3879C CAPSTONE PROJECT

The following show the outcomes of the performance evaluation for the respective **classification machine learning algorithms (baseline)** for the respective medical diseases.

**A) For Diabetes**

Cross-Validation Scores for Diabetes:			
	Random Forest	Gradient Boosting	Extra Trees
Train Accuracy	0.928201	0.983124	0.725093
Validate Accuracy	0.918935	0.971598	0.701183
Train Precision	0.929950	0.983575	0.748202
Validate Precision	0.920668	0.972481	0.729633
Train Recall	0.928091	0.982962	0.725264
Validate Recall	0.919798	0.972397	0.700724
Train F1	0.927649	0.982974	0.729448
Validate F1	0.919090	0.971457	0.704465
Train ROC	0.968148	0.997788	0.909537
Validate ROC	0.960816	0.993452	0.896641

Figure 4.5e: Shuffle-Split Cross-Validation Scores for Diabetes

Classification Report for Random Forest			
	precision	recall	f1-score
No Diabetes	0.96	1.00	0.98
Pre-Diabetes	0.96	0.86	0.90
Diabetes	0.87	0.92	0.89
accuracy			0.92
macro avg	0.93	0.92	0.92
weighted avg	0.93	0.92	0.92

Figure 4.5f: Classification Report for Random Forest (Diabetes)

Classification Report for Gradient Boosting			
	precision	recall	f1-score
No Diabetes	0.96	1.00	0.98
Pre-Diabetes	0.98	0.99	0.99
Diabetes	0.99	0.94	0.97
accuracy			0.98
macro avg	0.98	0.98	0.98
weighted avg	0.98	0.98	0.98

Figure 4.5g: Classification Report for Gradient Boosting (Diabetes)

Classification Report for Extra Trees			
	precision	recall	f1-score
No Diabetes	0.59	0.78	0.67
Pre-Diabetes	0.66	0.55	0.60
Diabetes	0.99	0.86	0.92
accuracy			0.73
macro avg	0.75	0.73	0.73
weighted avg	0.75	0.73	0.73

Figure 4.5h: Classification Report for Extra Trees (Diabetes)

**B) For Hypertension**

	Cross-Validation Scores for Hypertension:		
	Random Forest	Gradient Boosting	Extra Trees
Train Accuracy	1.000000	1.000000	1.000000
Validate Accuracy	0.971429	0.974286	0.954286
Train Precision	1.000000	1.000000	1.000000
Validate Precision	0.971179	0.974958	0.955446
Train Recall	1.000000	1.000000	1.000000
Validate Recall	0.971604	0.974040	0.952646
Train F1	1.000000	1.000000	1.000000
Validate F1	0.970803	0.973696	0.952981
Train ROC	1.000000	1.000000	1.000000
Validate ROC	0.983561	0.979181	0.984008

Figure 4.5i: Shuffle-Split Cross-Validation Scores for Hypertension

Classification Report for Random Forest (Hypertension)				
	precision	recall	f1-score	support
No Hypertension	0.95	0.98	0.96	174
Hypertension	0.98	0.95	0.96	174
accuracy			0.96	348
macro avg	0.96	0.96	0.96	348
weighted avg	0.96	0.96	0.96	348

Figure 4.5j: Classification Report for Random Forest (Hypertension)

Classification Report for Gradient Boosting (Hypertension)				
	precision	recall	f1-score	support
No Hypertension	0.95	0.97	0.96	174
Hypertension	0.96	0.95	0.96	174
accuracy			0.96	348
macro avg	0.96	0.96	0.96	348
weighted avg	0.96	0.96	0.96	348

Figure 4.5k: Classification Report for Gradient Boosting (Hypertension)

Classification Report for Extra Trees (Hypertension)				
	precision	recall	f1-score	support
No Hypertension	0.94	0.97	0.96	174
Hypertension	0.97	0.94	0.96	174
accuracy			0.96	348
macro avg	0.96	0.96	0.96	348
weighted avg	0.96	0.96	0.96	348

Figure 4.5l: Classification Report for Extra Trees (Hypertension)

**C) For Obesity**

Cross-Validation Scores for Obesity:				
	Random Forest	Gradient Boosting	Extra Trees	
Train Accuracy	1.000000	1.000000	1.000000	
Validate Accuracy	0.989888	0.991011	0.976779	
Train Precision	1.000000	1.000000	1.000000	
Validate Precision	0.989712	0.990865	0.976325	
Train Recall	1.000000	1.000000	1.000000	
Validate Recall	0.989671	0.991119	0.976314	
Train F1	1.000000	1.000000	1.000000	
Validate F1	0.989662	0.990931	0.976187	
Train ROC	1.000000	1.000000	1.000000	
Validate ROC	0.999657	0.999787	0.998888	

Figure 4.5m: Shuffle-Split Cross-Validation Scores for Obesity

Classification Report for Random Forest				
	precision	recall	f1-score	support
Underweight	0.99	1.00	0.99	667
Normal Weight	0.97	0.98	0.98	667
Overweight	0.97	0.98	0.98	667
Obesity	1.00	0.99	0.99	667
accuracy			0.99	2668
macro avg	0.99	0.99	0.99	2668
weighted avg	0.99	0.99	0.99	2668

Figure 4.5n: Classification Report for Random Forest (Obesity)

Classification Report for Gradient Boosting				
	precision	recall	f1-score	support
Underweight	0.99	0.99	0.99	667
Normal Weight	0.99	0.98	0.98	667
Overweight	0.98	0.99	0.99	667
Obesity	1.00	0.99	1.00	667
accuracy			0.99	2668
macro avg	0.99	0.99	0.99	2668
weighted avg	0.99	0.99	0.99	2668

Figure 4.5o: Classification Report for Gradient Boosting (Obesity)

Classification Report for Extra Trees				
	precision	recall	f1-score	support
Underweight	0.99	1.00	0.99	667
Normal Weight	0.96	0.96	0.96	667
Overweight	0.96	0.96	0.96	667
Obesity	1.00	0.98	0.99	667
accuracy			0.97	2668
macro avg	0.97	0.97	0.97	2668
weighted avg	0.97	0.97	0.97	2668

Figure 4.5p: Classification Report for Extra Trees (Obesity)

---

C3879C CAPSTONE PROJECT

Based on the above observation, a **random forest classifier** was selected as the baseline model for both **obesity and hypertension** because of its model characteristics such that it is less biased and exhibits moderate variance in nature. Furthermore, the training and validation scores for all performance metrics showing not much difference with the rest of the classifiers, and the ROC scores also revealed that the random forest classifier is capable of distinguishing between classes. On the other hand, the **gradient boosting classifier** was selected as the baseline model for **diabetes** as it shows the only best performance metrics as compared to the rest of the other classifiers.

The above-selected classifier baseline models will proceed with the hyperparameter optimization process, which will be discussed in the next section, it is to achieve the best hyperparameters that maximize the predictive performance.

#### 4.6 Hyperparameters Optimization for Proposed Baseline Model

Similar as previous section, a pipeline object was built by providing it with a list of steps. Each steps represents a tuple containing a string of name and an instance of an estimator.

For example, in hypertension, the first step was created with a ‘scaler’, which is an instance of `MinMaxScaler()`, followed by a ‘power’, is an instance of `PowerTransformer()`, and the final step was created with a ‘rfc’, is an instance of the random forest classifier. As mentioned earlier, each pipeline defines one (or more) data preparation techniques such as `MinMaxScaler` or `StandardScaler`, `PowerTransformer`, and ends with a defined model that takes the transformed training data as input. Figure 4.6a illustrates a sample code snippet of the model pipeline for hypertension.

```
386 #build pipeline
387 pipe = Pipeline([('scaler', MinMaxScaler()), ('power', PowerTransformer()), ('rfc', RandomForestClassifier(random_state = 42,
388 max_leaf_nodes = None,
389 max_features = "sqrt"))])
390
```

Figure 4.6a: A sample code of the model pipeline for hypertension

For both **obesity** and **diabetes**, it can be created a similar process and achieve the same objective as hypertension, such that the first step was created with a ‘scaler’ of `StandardScaler ()`, followed by created with a ‘model’, of the respective classifier algorithms. Figure 4.6b and Figure 4.6c illustrates a sample code snippet of the model pipeline for both **obesity** and **diabetes** training set.

Figure 4.6b: A sample code of the model pipeline for diabetes

Figure 4.6c: A sample code of the model pipeline for obesity

The syntax to define a parameter grid for a pipeline is to specify for each parameter the step name, followed by `__`(a double underscore), followed by the parameter name. To search over the `n_estimators` parameter, for example, in the random forest classifier, the

C3879C CAPSTONE PROJECT

---

“rfc \_\_ n\_estimators” was used as key in the parameter grid dictionary, and similar for the rest of the other parameters. Figure 4.6d, Figure 4.6e, and Figure 4.6f illustrate the sample code of construction of the parameter grid for the respective medical disease.

```
401 # set parameter grid
402 param_grid = {'gbc_n_estimators': [50, 100, 150, 200, 250],
403             'gbc_max_depth': np.arange(1,5),
404             'gbc_learning_rate': [0.001, 0.01, 0.1]}
```

Figure 4.6d: A sample code of the parameter grid for Gradient Boosting Classifier  
(Diabetes)

```
391 # set parameter grid
392 param_grid = {'rfc_n_estimators': [50, 100, 150, 200, 250, 300, 350],
393                 'rfc_criterion': ['gini', 'entropy'],
394                 'rfc_max_depth': np.arange(5, 10),
395                 'rfc_min_samples_leaf': np.arange(2, 7)}
```

Figure 4.6e: A sample code of the parameter grid for Random Forest Classifier  
(Hypertension)

```
401 # set parameter grid
402 param_grid = {'rfc_n_estimators': [80, 100, 150, 200, 250, 300],
403                 'rfc_max_depth': np.arange(2, 7),
404                 'rfc_criterion': ['gini', 'entropy'],
405                 'rfc_min_samples_split': np.arange(2, 5)}
```

Figure 4.6f: A sample code of the parameter grid for Random Forest Classifier  
(Obesity)

In the hyperparameter optimization process, a GridSearchCV function was implemented to search the best or optimal hyperparameters based on the best scores. For each split in the cross-validation, the MinMaxScaler () or the StandardScaler () is refit only with the training splits and no information was leaked from the testing split into the parameter search. As mentioned before, the testing split was not based on the testing set. It is the defined amount of the % split from the training set. As such, all the respective test sets were not seen during the training phase as well as the hyperparameter optimization process. Furthermore, to get a finer control over the cross-validation, the CV parameter in the GridSearchCV process uses Shuffle-Split and the amount of split were defined as 80% training set and 20% testing (or validation) set.

## C3879C CAPSTONE PROJECT

The main important parameters of **gradient boosted tree models** to be optimized are the number of trees, **n\_estimators**, and the **learning\_rate**, which controls the degree to which each tree is allowed to correct the mistakes of the previous trees. These two parameters are highly interconnected such that a lower learning\_rate means that more trees are required to build a model in more complexity. Similar to the random forests, a higher n\_estimators are always better. But however, increasing n\_estimators in gradient boosting will leads to a model complexity, which may lead to overfitting. In addition, another important parameter is max\_depth, it is to minimizes the complexity of each tree. Usually, max\_depth is set to a very low split in gradient boosted models, often no deeper than five splits.

Figure 4.6g illustrates the sample code of the entire hyperparameter optimization process for Gradient Boosting Classifier for diabetes.

```

121 def hyperparameters_tuning(pipe, param_grid, X, y):
122
123     # split the dataset
124     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
125     # perform hyperparameter tuning using gridsearchcv
126     gs_res = GridSearchCV(pipe, param_grid = param_grid,
127                           cv = shuffle_split, scoring = 'accuracy', verbose = 3)
128     # fitted the model for grid search
129     gs_res.fit(X, y)
130     return gs_res
131
132
133 ###### Hyperparameter Optimization on Gradient Boosting #####
134
135 # #build pipeline
136 pipe = Pipeline([('scaler', StandardScaler()),("gbc", GradientBoostingClassifier(random_state = 42,
137                                         max_leaf_nodes = None))]
138
139 # set parameter grid
140 param_grid = {'gbc__n_estimators': [50, 100, 150, 200, 250],
141                 'gbc__max_depth': np.arange(1,5),
142                 'gbc__learning_rate': [0.001, 0.01, 0.1]}
143
144 # call the function to perform hyperparameter tuning
145 gridResult = hyperparameters_tuning(pipe, param_grid, X_smt, Y_smt)
146
147 # display the best parameters after tuning
148 print("\nHyperparameter Optimization Process for Gradient Boosting Classifier (Diabetes): ")
149 print("====")
150 print("The best score is {}".format(gridResult.best_score_))
151 print("The best params is {}".format(gridResult.best_params_))
152 print("The best estimator is {}".format(gridResult.best_estimator_))
153
154 # store the best params into Series
155 best_params_df = pd.Series(gridResult.best_params_)
156 print(best_params_df)

```

Figure 4.6g: A sample code of the hyperparameter optimization process for Gradient Boosting Classifier (Diabetes)

## C3879C CAPSTONE PROJECT

The important parameters to be optimized in **random forest tree models** are **n\_estimators**, **criterion**, **max\_depth**, and **min\_samples\_leaf**. For **n\_estimators**, the larger is always better, and averaging more trees will yield a more robust ensemble by minimizing overfitting. However, there are some drawbacks, more trees need more memory and more time to train. In general, a **smaller max\_features** reduce overfitting and it is a good rule of thumb to use the default values for max features: max features=sqrt(n features) for classification problems. Furthermore, adding **min\_samples\_leaf** might sometimes also improve performance.

Figure 4.6h and Figure 4.6i illustrates the sample code of the entire hyperparameter optimization process for Random Forest Classifier for hypertension and obesity respectively.

```

109 def hyperparameters_tuning(pipe, param_grid, X, y):
110
111     # split the dataset
112     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
113     # perform hyperparameter tuning using gridsearchcv
114     gs_res = GridSearchCV(pipe, param_grid = param_grid,
115                           cv = shuffle_split, scoring = 'accuracy', verbose = 3)
116     # fitted the model for grid search
117     gs_res.fit(X, y)
118     return gs_res
119
120
121 ###### Hyperparameter Optimization on Random Forest Classifier #####
122 #build pipeline
123 pipe = Pipeline([('scaler', MinMaxScaler()), ('power', PowerTransformer()), ('rfc', RandomForestClassifier(random_state = 42,
124
125                                         max_leaf_nodes = None,
126                                         max_features = "sqrt"))])
127
128 # set parameter grid
129 param_grid = {'rfc_n_estimators': [50, 100, 150, 200, 250, 300, 350],
130               'rfc_criterion': ['gini', 'entropy'],
131               'rfc_max_depth': np.arange(5, 10),
132               'rfc_min_samples_leaf': np.arange(2, 7)}
133
134
135 # call the function to perform hyperparameter tuning
136 gridResult = hyperparameters_tuning(pipe, param_grid, X_smt, Y_smt)
137
138 # display the best parameters after tuning
139 print("\nHyperparameter Optimization Process for Random Forest Classifier (Hypertension): ")
140 print("====")
141 print("The best score is {}".format(gridResult.best_score_))
142 print("The best params is {}".format(gridResult.best_params_))
143 print("The best estimator is {}".format(gridResult.best_estimator_))
144
145 # store the best params into Series
146 best_params_df = pd.Series(gridResult.best_params_)
147 print(best_params_df)

```

Figure 4.6h: A sample code of the hyperparameter optimization process for Random Forest Classifier (Hypertension)

## C3879C CAPSTONE PROJECT

```

138 def hyperparameters_tuning(pipe, param_grid, X, y):
139     # split the dataset
140     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
141     # perform hyperparameter tuning using gridsearchcv
142     gs_res = GridSearchCV(pipe, param_grid = param_grid,
143                           cv = shuffle_split, scoring = 'accuracy', verbose = 3)
144     # fitted the model for grid search
145     gs_res.fit(X, y)
146     return gs_res

147
148 ##### Hyperparameter Optimization on Random Forest #####
149 # build pipeline
150 pipe = Pipeline([('scaler', StandardScaler()),("rfc", RandomForestClassifier(random_state = 42,
151                                         max_leaf_nodes = None,
152                                         max_features = 'sqrt'))])

153 # set parameter grid
154 param_grid = {'rfc__n_estimators': [80, 100, 150, 200, 250, 300],
155               'rfc__max_depth': np.arange(2, 7),
156               'rfc__criterion': ['gini', 'entropy'],
157               'rfc__min_samples_split': np.arange(2, 5)}

158 # call the function to perform hyperparameter tuning
159 gridResult = hyperparameters_tuning(pipe, param_grid, X_smt, Y_smt)

160 # display the best parameters after tuning
161 print("\nHyperparameter Optimization Process for Random Forest Classifier (Obesity): ")
162 print("====")
163 print("The best score is {}".format(gridResult.best_score_))
164 print("The best params is {}".format(gridResult.best_params_))
165 print("The best estimator is {}".format(gridResult.best_estimator_))

166 # store the best params into Series
167 best_params_df = pd.Series(gridResult.best_params_)
168 print(best_params_df)

```

Figure 4.6i: A sample code of the hyperparameter optimization process for Random Forest Classifier (Obesity)

C3879C CAPSTONE PROJECT

---

As mentioned earlier, the hyperparameter optimization process derived the hyperparameters with the best scores. Figure 4.6j, Figure 4.6k, and Figure 4.6l illustrate the outcome of the hyperparameter optimization process for all the respective classifier models.

```
Hyperparameter Optimization Process for Gradient Boosting Classifier (Diabetes):
=====
The best score is 0.9863905325443787
The best params is {'gbc__learning_rate': 0.1, 'gbc__max_depth': 3, 'gbc__n_estimators': 100}
The best estimator is Pipeline(steps=[('scaler', StandardScaler()),
                                     ('gbc', GradientBoostingClassifier(random_state=42))])
```

Figure 4.6j: The outcome of the hyperparameter optimization process for Gradient Boosting Classifier (Diabetes)

```
Hyperparameter Optimization Process for Random Forest Classifier (Hypertension):
=====
The best score is 0.9742857142857144
The best params is {'rfc__criterion': 'gini', 'rfc__max_depth': 6, 'rfc__min_samples_leaf': 2, 'rfc__n_estimators': 150}
The best estimator is Pipeline(steps=[('scaler', MinMaxScaler()), ('power', PowerTransformer()),
                                     ('rfc',
                                       RandomForestClassifier(max_depth=6, max_features='sqrt',
                                                             min_samples_leaf=2, n_estimators=150,
                                                             random_state=42))])
```

Figure 4.6k: The outcome of the hyperparameter optimization process for Random Forest Classifier (Hypertension)

```
Hyperparameter Optimization Process for Random Forest Classifier (Obesity):
=====
The best score is 0.9573033707865168
The best params is {'rfc__criterion': 'entropy', 'rfc__max_depth': 6, 'rfc__min_samples_split': 2, 'rfc__n_estimators': 300}
The best estimator is Pipeline(steps=[('scaler', StandardScaler()),
                                     ('rfc',
                                       RandomForestClassifier(criterion='entropy', max_depth=6,
                                                             max_features='sqrt', n_estimators=300,
                                                             random_state=42))])
```

Figure 4.6l: The outcome of the hyperparameter optimization process for Random Forest Classifier (Obesity)

## 4.7 Retrained Model with Best Hyperparameter and Evaluation

The selected baseline classifier models (from the previous section) for each respective medical disease were re-trained with the best hyperparameters achieved during the hyperparameter optimization process.

Similar to the baseline model building and evaluation, the performance of the re-trained model with the best hyperparameters were evaluated via the shuffle-split cross-validation technique.

Figure 4.7a, Figure 4.7b, and Figure 4.8c illustrate the sample codes for the re-trained classifier models as well as the evaluation of the performance metrics.

```

82 def cross_validation_eval(model, X_smt, Y_smt):
83
84     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
85     cv_result = cross_validate(model,
86                               X_smt,
87                               Y_smt,
88                               cv = shuffle_split,
89                               n_jobs = 1,
90                               scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro'),
91                               return_train_score = True)
92
93
94     # compute the ROC scores
95     myscore = make_scorer(roc_auc_score, multi_class='ovo', needs_proba=True)
96     roc_scores = cross_validate(model,
97                               X_smt,
98                               Y_smt,
99                               cv = shuffle_split,
100                              n_jobs = 1,
101                              scoring = myscore,
102                              return_train_score = True)
103
104     return cv_result, roc_scores
105
106
107 ###### Rebuild Model (Gradient Boosting) using Best Parameters #####
108
109
110 # From Hyperparameter Tuning:
111 # learning rate = 0.001, max_depth = 2, n_estimators = 50
112
113 gbc_model = GradientBoostingClassifier(random_state = 42,
114                                         max_leaf_nodes = None,
115                                         n_estimators = int(best_params_df.loc['gbc_n_estimators']),
116                                         max_depth = int(best_params_df.loc['gbc_max_depth']),
117                                         learning_rate = best_params_df.loc['gbc_learning_rate'])
118
119
120 ##### Evaluate the Rebuild Model - Gradient Boosting #####
121 pipe_gbc = Pipeline([('scaler', StandardScaler()), ('gbc', gbc_model)])
122 gbc_cv_eval, gbc_rocScores = cross_validation_eval(pipe_gbc, X_smt, Y_smt)
123
124 gbc_cv_results = pd.DataFrame({
125     'Train Accuracy': gbc_cv_eval['train_accuracy'].mean(),
126     'Validate Accuracy': gbc_cv_eval['test_accuracy'].mean(),
127     'Train Precision': gbc_cv_eval['train_precision_macro'].mean(),
128     'Validate Precision': gbc_cv_eval['test_precision_macro'].mean(),
129     'Train Recall': gbc_cv_eval['train_recall_macro'].mean(),
130     'Validate Recall': gbc_cv_eval['test_recall_macro'].mean(),
131     'Train F1': gbc_cv_eval['train_f1_macro'].mean(),
132     'Validate F1': gbc_cv_eval['test_f1_macro'].mean(),
133     'Train ROC': gbc_rocScores['train_score'].mean(),
134     'Validate ROC': gbc_rocScores['test_score'].mean(), index = ['Gradient Boosting']).T
135
136
137 # display the cross-validation scores before and after hyperparameter optimization
138 print("\nCross-Validation Scores (Before Hyperparameter Optimization) for Gradient Boosting Classifier (Diabetes): ")
139 print(baseline_models_cv_eval['Gradient Boosting'])
140 print("\nCross-Validation Scores (After Hyperparameter Optimization) for Gradient Boosting Classifier (Diabetes): ")
141 print(gbc_cv_results)

```

Figure 4.7a: A sample code of the re-trained and re-evaluated of the Gradient Boosting Classifier with the Best Hyperparameters (Diabetes)

## C3879C CAPSTONE PROJECT

```

81  def cross_validation_eval(pipe, X_smt, Y_smt):
82
83      shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
84      cv_result = cross_validate(pipe,
85                                 X_smt,
86                                 Y_smt,
87                                 cv = shuffle_split,
88                                 n_jobs = 1,
89                                 scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro', 'roc_auc'),
90                                 return_train_score = True,
91                                 return_estimator = True)
92
93      return cv_result
94
95 ###### Rebuild Model (Random Forest) using Best Parameters #####
96
97 # From Hyperparameter Tuning:
98 # criterion = gini, max_depth = 5, n_estimators = 50, min_samples_leaf = 2
99
100 ▼ rfc_model = RandomForestClassifier(random_state = 42,
101                                       max_leaf_nodes = None,
102                                       max_features = "sqrt",
103                                       n_estimators = int(best_params_df.loc['rfc_n_estimators']),
104                                       criterion = best_params_df.loc['rfc_criterion'],
105                                       max_depth = int(best_params_df.loc['rfc_max_depth']),
106                                       min_samples_leaf = int(best_params_df.loc['rfc_min_samples_leaf']))
107
108 ###### Evaluate the Rebuild Model - Random Forest #####
109 pipe_rfc = Pipeline([('scaler', MinMaxScaler()), ('power', PowerTransformer()), ('rfc', rfc_model)])
110 rfc_cv_eval = pd.DataFrame(cross_validation_eval(pipe_rfc, X_smt, Y_smt))
111
112 ▼ rfc_cv_results = pd.DataFrame({
113     'Train Accuracy': rfc_cv_eval['train_accuracy'].mean(),
114     'Validate Accuracy': rfc_cv_eval['test_accuracy'].mean(),
115     'Train Precision': rfc_cv_eval['train_precision_macro'].mean(),
116     'Validate Precision': rfc_cv_eval['test_precision_macro'].mean(),
117     'Train Recall': rfc_cv_eval['train_recall_macro'].mean(),
118     'Validate Recall': rfc_cv_eval['test_recall_macro'].mean(),
119     'Train F1': rfc_cv_eval['train_f1_macro'].mean(),
120     'Validate F1': rfc_cv_eval['test_f1_macro'].mean(),
121     'Train ROC': rfc_cv_eval['train_roc_auc'].mean(),
122     'Validate ROC': rfc_cv_eval['test_roc_auc'].mean()}, index = ['Random Forest']).T
123
124
125
126 # display the cross-validation scores before and after hyperparameter optimization
127 print("\nCross-Validation Scores (Before Hyperparameter Optimization) for Random Forest Classifier (Hypertension): ")
128 print(baseline_models_CV_eval['Extra Trees'])
129 print("\nCross-Validation Scores (After Hyperparameter Optimization) for Random Forest Classifier (Hypertension): ")
130 print(rfc_cv_results)

```

Figure 4.7b: A sample code of the re-trained and re-evaluated of the Random Forest Classifier with the Best Hyperparameters (Hypertension)

```

95 def cross_validation_eval(model, X_smt, Y_smt):
96
97     shuffle_split = ShuffleSplit(test_size = 0.20, train_size = 0.80, n_splits = 5, random_state = 0)
98
99     # compute the accuracy, precision, recall and f1-scores
100    cv_result = cross_validate(model,
101                               X_smt,
102                               Y_smt,
103                               cv = shuffle_split,
104                               n_jobs = 1,
105                               scoring = ('accuracy', 'precision_macro', 'recall_macro', 'f1_macro'),
106                               return_train_score = True)
107
108    # compute the ROC scores
109    myscore = make_scoring(roc_auc_score, multi_class='ovo', needs_proba=True)
110    roc_scores = cross_validate(model,
111                               X_smt,
112                               Y_smt,
113                               cv = shuffle_split,
114                               n_jobs = 1,
115                               scoring = myscore,
116                               return_train_score = True)
117
118    return cv_result, roc_scores

```

## C3879C CAPSTONE PROJECT

```

423 ##### Rebuild Model (Random Forest) using Best Parameters #####
424
425 # From Hyperparameter Tuning:
426 # criterion = entropy, max_depth = 6, n_estimators = 80, min_samples_split = 3
427
428 rfc_model = RandomForestClassifier(random_state = 42,
429                                     max_features = 'sqrt',
430                                     max_leaf_nodes = None,
431                                     n_estimators = int(best_params_df.loc['rfc_n_estimators']),
432                                     max_depth = int(best_params_df.loc['rfc_max_depth']),
433                                     criterion = best_params_df.loc['rfc_criterion'],
434                                     min_samples_split = int(best_params_df.loc['rfc_min_samples_split']))
435
436
437 ##### Evaluate the Rebuild Model - Random Forest #####
438 pipe_rfc = Pipeline([('scaler', StandardScaler()), ('rfc', rfc_model)])
439 rfc_cv_eval, rfc_rocScores = cross_validation_eval(pipe_rfc, X_smt, Y_smt)
440
441
442 rfc_cv_results = pd.DataFrame({
443     'Train Accuracy': rfc_cv_eval['train_accuracy'].mean(),
444     'Validate Accuracy': rfc_cv_eval['test_accuracy'].mean(),
445     'Train Precision': rfc_cv_eval['train_precision_macro'].mean(),
446     'Validate Precision': rfc_cv_eval['test_precision_macro'].mean(),
447     'Train Recall': rfc_cv_eval['train_recall_macro'].mean(),
448     'Validate Recall': rfc_cv_eval['test_recall_macro'].mean(),
449     'Train F1': rfc_cv_eval['train_f1_macro'].mean(),
450     'Validate F1': rfc_cv_eval['test_f1_macro'].mean(),
451     'Train ROC': rfc_rocScores['train_score'].mean(),
452     'Validate ROC': rfc_rocScores['test_score'].mean(), index = ['Random Forest']).T
453
454
455 # display the cross-validation scores before and after hyperparameter optimization
456 print("\nCross-Validation Scores (Before Hyperparameter Optimization) for Random Forest Classifier (Obesity): ")
457 print(baseline_models_CV_eval['Random Forest'])
458 print("\nCross-Validation Scores (After Hyperparameter Optimization) for Random Forest Classifier (Obesity): ")
459 print(rfc_cv_results)

```

Figure 4.7c: A sample code of the re-trained and re-evaluated of the Random Forest Classifier with the Best Hyperparameters (Obesity)

## C3879C CAPSTONE PROJECT

Figure 4.7d, Figure 4.7e, Figure 4.7f and Figure 4.7g illustrate the outcomes of the **performance metrics before and after**, the **confusion matrix**, the **ROC-AUC** as well as the **classification report** of the Gradient Boosting Classifier with the best hyperparameters for diabetes.

Cross-Validation Scores (Before Hyperparameter Optimization) for Gradient Boosting Classifier (Diabetes):	
Train Accuracy	0.983124
Validate Accuracy	0.971598
Train Precision	0.983575
Validate Precision	0.972481
Train Recall	0.982962
Validate Recall	0.972397
Train F1	0.982974
Validate F1	0.971457
Train ROC	0.997788
Validate ROC	0.993452

Cross-Validation Scores (After Hyperparameter Optimization) for Gradient Boosting Classifier (Diabetes):	
Gradient Boosting	
Train Accuracy	0.994671
Validate Accuracy	0.986391
Train Precision	0.994773
Validate Precision	0.986391
Train Recall	0.994599
Validate Recall	0.986762
Train F1	0.994655
Validate F1	0.986402
Train ROC	0.999918
Validate ROC	0.997943

Figure 4.7d: The outcome of the performance metrics **before** and **after** for Gradient Boosting Classifier with best hyperparameters (Diabetes)

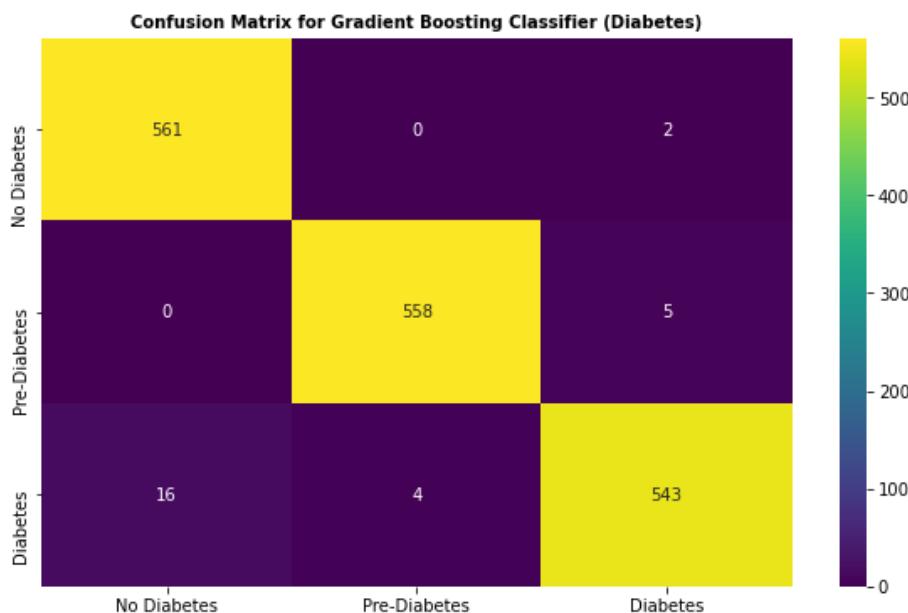


Figure 4.7e: The outcome of the confusion matrix for Gradient Boosting Classifier with best hyperparameters (Diabetes)

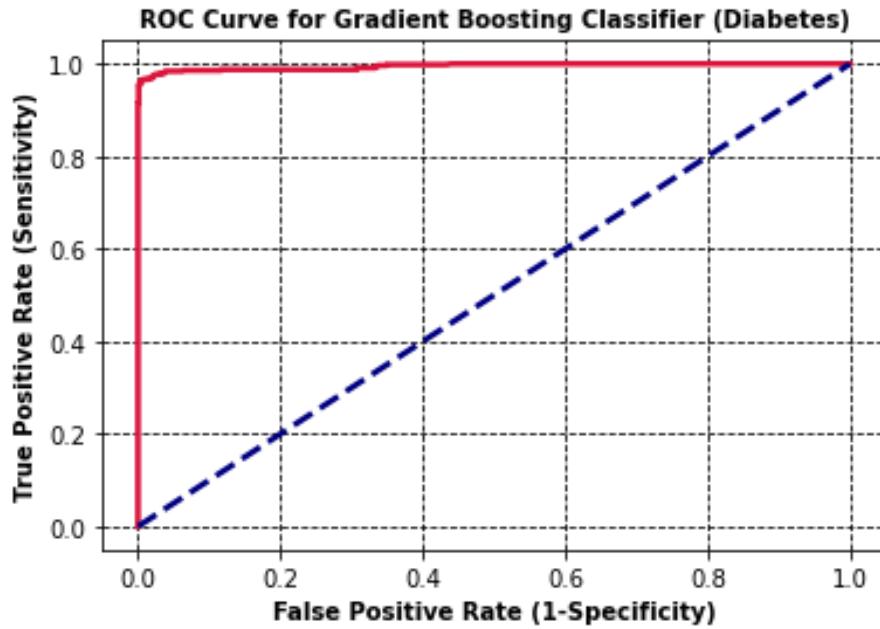


Figure 4.7f: The outcome of the ROC-AUC for Gradient Boosting Classifier with best hyperparameters (Diabetes)

Classification Report for Gradient Boosting with Hyperparameter Tuning (Diabetes)				
	precision	recall	f1-score	support
No Diabetes	0.97	1.00	0.98	563
Pre-Diabetes	0.99	0.99	0.99	563
Diabetes	0.99	0.96	0.98	563
accuracy			0.98	1689
macro avg	0.98	0.98	0.98	1689
weighted avg	0.98	0.98	0.98	1689

Figure 4.7g: The outcome of the classification report for Gradient Boosting Classifier with best hyperparameters (Diabetes)

Based on the illustration shown in figure 4.7d, the performance metrics for the gradient boosting classifier have improved after the best hyperparameters have been applied onto the gradient boosting classifier with the diabetes training set. Moreover, both training and validation ROC with 99% shows that the gradient boosting classifier model has the excellent capability to differentiate between the classes (no diabetes, pre-diabetes, or diabetes), which can be proven by the ROC-AUC shown in figure 4.7f. Furthermore, the gradient boosting classifier with hyperparameters has a very high accuracy rate of 99%, but it is also imperative to look at the other parameters to evaluate the performance of this classifier model. For example, high precision rate relates to the low false-positive rate, this is absolutely aggregable with the confusion matrix shown in figure 4.7e, which has about less

C3879C CAPSTONE PROJECT

---

than 5% of having false-positive rate. As such, this gradient boosting classifier with hyperparameters has a very high precision rate of nearly 99% based on the training and validation set. In addition, this gradient boosting classifier with hyperparameters also has an excellent recall rate of 99% based on the validation and training set, which is good for this hyperparameter model as it is above 50% and the high recall rate also relates to the low false-negative rate, which can be further proven by the confusion matrix as shown in figure 4.7e.

Based on the confusion matrix shown in figure 4.7e, both **false positive** and **false negative** are the major concern in the biomedical field. False-positive represented that these patients were predicted with diabetes, but in fact, these patients were having no diabetes. Conversely, false-negative represented that these patients were predicted with no diabetes, but in actual fact, these patients were having diabetes. As such, both false-positive and false-negative are the most important factors that we would like to minimize. One must also note that the diabetes training set consists of multi-class. As such, the computation for TP, FP, TN, and FN are slightly different from the binary classification, and below shows the computation for the TP, FP, TN, and FN for confusion matrix with multi-class.

**For Class 0: No Diabetes**

$$TP = 561$$

$$FP = 0 + 16 = 16$$

$$TN = 558 + 5 + 4 + 543 = 1110$$

$$FN = 0 + 2 = 2$$

**For Class 1: Pre-Diabetes**

$$TP = 558$$

$$FP = 0 + 4 = 4$$

$$TN = 561 + 2 + 16 + 543 = 1122$$

$$FN = 0 + 5 = 5$$

**For Class 2: Diabetes**

$$TP = 543$$

$$FP = 5 + 2 = 7$$

$$TN = 561 + 558 + 0 + 0 = 1119$$

$$FN = 16 + 4 = 20$$

## C3879C CAPSTONE PROJECT

Thus, both false-positive rate and false-negative rate for all classes are not too much of a concern for this case as it is likely to have less than 5% chances. In other words, the gradient boosting classifier with hyperparameters for diabetes works considerably well.

Figure 4.7h, Figure 4.7i, Figure 4.7j and Figure 4.7k illustrate the outcomes of the **performance metrics before and after**, the **confusion matrix**, the **ROC-AUC** as well as the **classification report** of the Random Forest Classifier with the best hyperparameters for hypertension.

Cross-Validation Scores (Before Hyperparameter Optimization) for Random Forest Classifier (Hypertension):	
Train Accuracy	1.000000
Validate Accuracy	0.954286
Train Precision	1.000000
Validate Precision	0.955446
Train Recall	1.000000
Validate Recall	0.952646
Train F1	1.000000
Validate F1	0.952981
Train ROC	1.000000
Validate ROC	0.984008

Cross-Validation Scores (After Hyperparameter Optimization) for Random Forest Classifier (Hypertension):	
Random Forest	
Train Accuracy	0.976259
Validate Accuracy	0.974286
Train Precision	0.976673
Validate Precision	0.974363
Train Recall	0.976329
Validate Recall	0.973351
Train F1	0.976230
Validate F1	0.973652
Train ROC	0.998363
Validate ROC	0.983529

Figure 4.7h: The outcome of the performance metrics **before** and **after** for Random Forest Classifier with best hyperparameters (Hypertension)

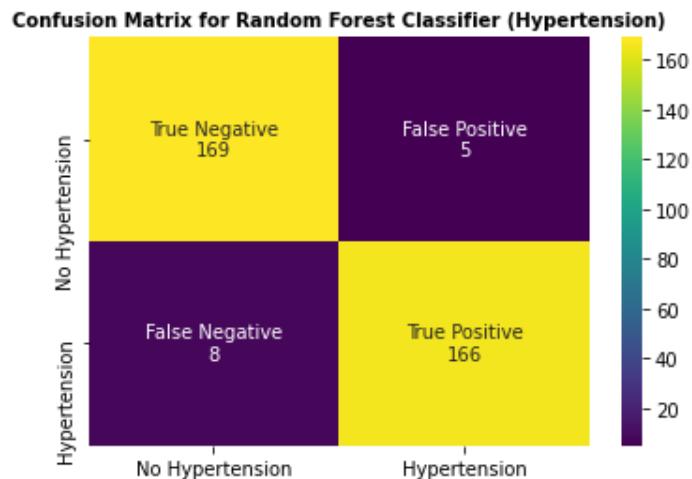


Figure 4.7i: The outcome of the confusion matrix for Random Forest Classifier with best hyperparameters (Hypertension)

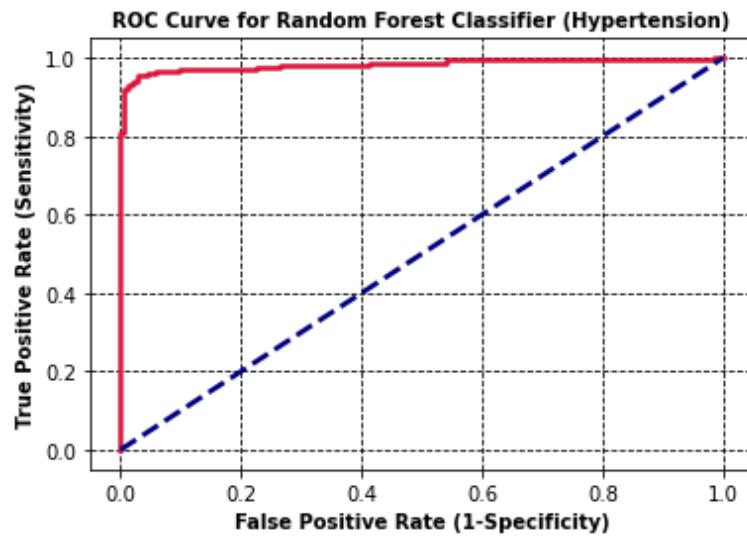


Figure 4.7j: The outcome of the ROC-AUC for Random Forest Classifier with best hyperparameters (Hypertension)

Classification Report for Random Forest with Hyperparameter Tuning (Hypertension)				
	precision	recall	f1-score	support
No Hypertension	0.95	0.97	0.96	174
Hypertension	0.97	0.95	0.96	174
accuracy			0.96	348
macro avg		0.96	0.96	348
weighted avg		0.96	0.96	348

Figure 4.7k: The outcome of the classification report for Random Forest Classifier with best hyperparameters (Hypertension)

Based on the illustration shown in figure 4.7h, prior to the hyperparameters optimization process, the random forest classifier accuracy score for the training set is slightly higher than the validation set. This phenomenon is known as **overfitting** and the cause of the overfitting could be mainly due to the **use of default parameters**. As such, after applying the best hyperparameters into the random forest classifier, the overfitting issues have been resolved. Furthermore, the rest of the performance metrics were also improved significantly. Moreover, both training and validation ROC with 98% shows that the random forest classifier model has also with the excellent capability to differentiate between the classes (no hypertension or hypertension), which can be proven by the ROC-AUC shown in figure 4.7j. As mentioned earlier, it is also imperative to look at the other parameters to evaluate the performance of this classifier model. For example, high precision rate relates to the low

---

C3879C CAPSTONE PROJECT

---

false-positive rate, this is absolutely aggregable with the confusion matrix shown in figure 4.7e, which has about less than 5% of having false-positive rate. As such, this random forest classifier with hyperparameters has a very high precision rate of nearly 97% based on the training and validation set. In addition, this random forest classifier with hyperparameters also performed exceptionally well with a recall rate of 97%, which is good for this hyperparameter model as it is above 50%, and high recall rate can also relate to the low false-negative rate, which can be proven by the confusion matrix as shown in Figure 4.7i.

Based on the confusion matrix shown in figure 4.7i, both **false-positive** and **false-negative** are the major concern in the biomedical field. False-positive represented that these patients were predicted with hypertension, but in fact, these patients were having no hypertension. Conversely, false-negative represented that these patients were predicted with no hypertension, but in actual fact, these patients were having hypertension. As such, both false positive and false negative are the most potent factors that we would like to minimize. For this case, both false positive and false negative for all classes are extremely low and most likely to have less than 5% chances. This means that the random forest classifier with hyperparameters for hypertension works generally well.

## C3879C CAPSTONE PROJECT

Figure 4.7l, Figure 4.7m, Figure 4.7n and Figure 4.7o illustrate the outcomes of the **performance metrics before** and **after**, the **confusion matrix**, the **ROC-AUC** as well as the **classification report** of the Random Forest Classifier with the best hyperparameters for obesity.

Cross-Validation Scores (Before Hyperparameter Optimization) for Random Forest Classifier (Obesity):	
Train Accuracy	1.000000
Validate Accuracy	0.989888
Train Precision	1.000000
Validate Precision	0.989712
Train Recall	1.000000
Validate Recall	0.989671
Train F1	1.000000
Validate F1	0.989662
Train ROC	1.000000
Validate ROC	0.999657

Cross-Validation Scores (After Hyperparameter Optimization) for Random Forest Classifier (Obesity):	
Random Forest	
Train Accuracy	0.968885
Validate Accuracy	0.957303
Train Precision	0.969301
Validate Precision	0.956819
Train Recall	0.969021
Validate Recall	0.956741
Train F1	0.969021
Validate F1	0.956539
Train ROC	0.996811
Validate ROC	0.995296

Figure 4.7l: The outcome of the performance metrics **before** and **after** for Random Forest Classifier with best hyperparameters (Obesity)

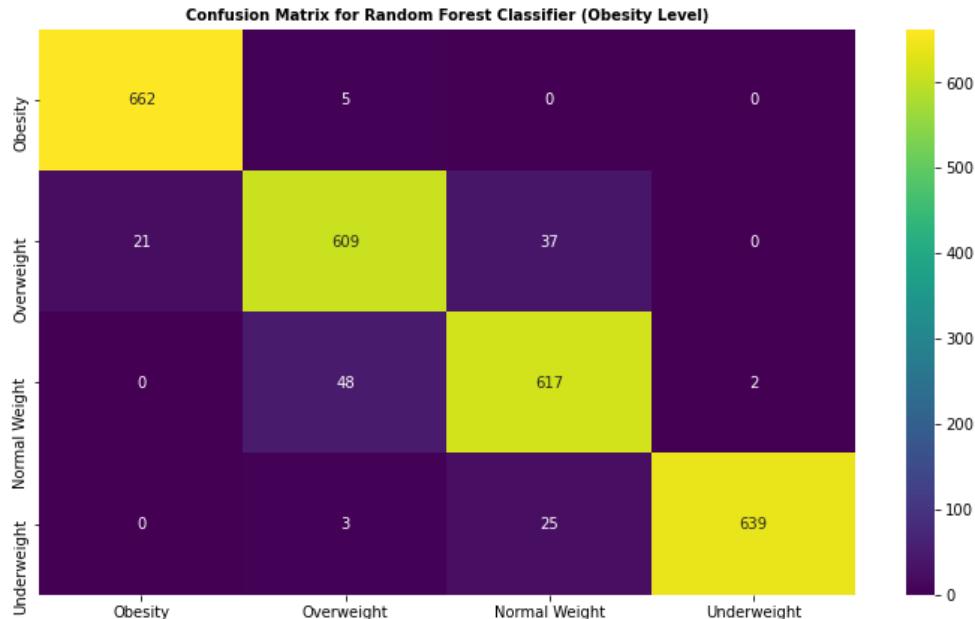


Figure 4.7m: The outcome of the confusion matrix for Random Forest Classifier with best hyperparameters (Obesity)

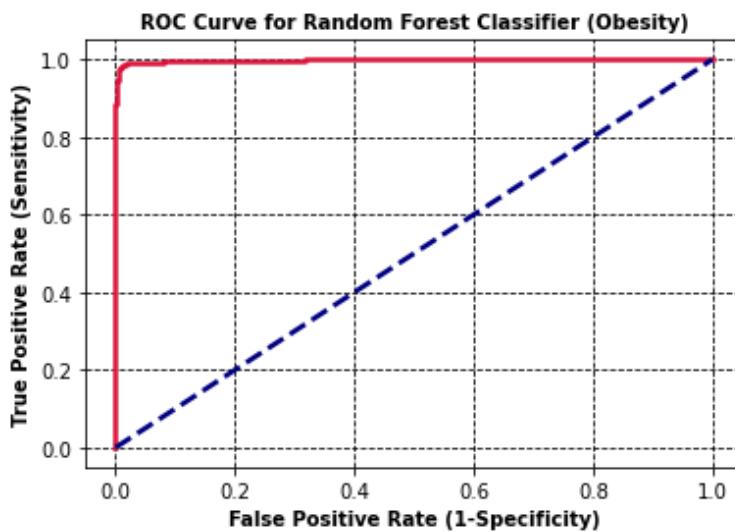


Figure 4.7n: The outcome of the confusion matrix for Random Forest Classifier with best hyperparameters (Obesity)

Classification Report for Random Forest with Hyperparameter Tuning (Obesity)				
	precision	recall	f1-score	support
Underweight	0.97	0.99	0.98	667
Normal Weight	0.92	0.91	0.91	667
Overweight	0.91	0.93	0.92	667
Obesity	1.00	0.96	0.98	667
accuracy			0.95	2668
macro avg	0.95	0.95	0.95	2668
weighted avg	0.95	0.95	0.95	2668

Figure 4.7o: The outcome of the classification report for Random Forest Classifier with best hyperparameters (Obesity)

Based on the illustration shown in figure 4.7l, prior to the hyperparameters optimization process, the random forest classifier accuracy score for the obesity training set is slightly higher than the obesity validation set. This phenomenon is known as **overfitting** and the cause of the overfitting could be mainly due to the **use of default parameters**. As such, after applying with the best hyperparameters into the random forest classifier, the overfitting issues has been minimized. Furthermore, the rest of the performance metrics were also improved significantly. Moreover, both training and validation ROC with 99% shows that the random forest classifier model has also with the excellent capability to differentiate between the classes (underweight, normal weight, overweight, or obesity), which can be

C3879C CAPSTONE PROJECT

---

proven by the ROC-AUC shown in figure 4.7n. As mentioned earlier, it is also imperative to look at the other parameters to evaluate the performance of this classifier model. For example, high precision relates to the low false positive rate, this is absolutely aggregable with the confusion matrix shown in figure 4.7m, which is about less than 5% of having false positive rate. As such, this random forest classifier with hyperparameters has a very high precision rate of nearly 96% based on the training and validation set. In addition, this random forest classifier with hyperparameters also performed exceptionally recall rate of 96% based on the validation and training set, which is good for this hyperparameter model as it is above 50% and high recall rate also relates to low false negative rate, which can be proven by the confusion matrix shown in Figure 4.7m.

Based on the confusion matrix shown in figure 4.7i, both **false positive** and **false negative** are the major concern in the biomedical field. False positive represented that these patients were predicted with obesity, but in fact, these patients were having no obesity. Conversely, false negative represented that these patients were predicted with no obesity, but in actual fact, these patients were having obesity. As such, both false positive and false negative are the most crucial factors that we would like to minimize. Similar as the diabetes, the obesity training set consists of multi-class. As such, the computation for TP, FP, TN, and FN are slightly different from the binary classification, and below shows the computation for the TP, FP, TN, and FN for confusion matrix with multi-class.

**For Class 1: Underweight**

$$TP = 662$$

$$FP = 21 + 0 + 0 = 21$$

$$TN = 609 + 37 + 0 + 48 + 617 + 2 + 3 + 25 + 639 = 1980$$

$$FN = 5 + 0 + 0 = 5$$

**For Class 2: Normal Weight**

$$TP = 609$$

$$FP = 48 + 3 + 5 = 56$$

$$TN = 662 + 0 + 0 + 0 + 0 + 617 + 2 + 25 + 639 = 1945$$

$$FN = 21 + 37 + 0 = 58$$

**For Class 3: Overweight**

TP = 617

FP =  $0 + 37 + 25 = 62$

TN =  $662 + 5 + 21 + 609 + 0 + 3 + 0 + 0 + 639 = 1939$

FN =  $48 + 2 = 50$

**For Class 4: Obesity**

TP = 639

FP =  $0 + 0 + 0 = 2$

TN =  $662 + 5 + 0 + 21 + 609 + 37 + 0 + 48 + 617 = 1999$

FN =  $3 + 25 + 0 = 28$

Generally, both false positive and false negative for all classes are not so much of the concern as they are likely to have less than 5% chances. In other words, this classifier model with hyperparameters for obesity works generally well.

## 4.8 Evaluation of the Hyperparameter Classifier Models using Testing Sets

Basically, a test set is a data set that is independent of the training set, but it follows the same probability distribution as the training data set. If a model fits into training set also fits into the test data as well. A test set is therefore a set of examples used to only assess the overall performance of the fully implemented classifiers. To achieve this, the final classifiers models were used to predict the classifications of examples in the test set. These predictions are then compared with the examples of true classifications in the test set to assess the model performances.

Figure 4.8a, Figure 4.8b, Figure 4.8c, and Figure 4.8d illustrate the outcome of the performance of the gradient boosting classifier with hyperparameter based on **diabetes testing set**.

```
Performance Metrics of the Gradient Boosting with Hyperparameter Tuning (Diabetes) using Testing Set:
=====
Accuracy Score: 0.9685314685314685
Precision Score: 0.8902255639097745
Recall Score: 0.9875518672199171
F1 Score: 0.9338360767139836
ROC Score: 0.996394632040826
```

Figure 4.8a: Performance metrics of Gradient Boosting Classifier with hyperparameters (Diabetes)

```
Classification Report for Gradient Boosting with Hyperparameter Tuning (Diabetes) using Testing Set:
      precision    recall    f1-score   support
No Diabetes      0.83     1.00      0.91       29
Pre-Diabetes     0.84     1.00      0.91       16
  Diabetes       1.00     0.96      0.98      241
          accuracy                           0.97      286
        macro avg       0.89     0.99      0.93      286
      weighted avg     0.97     0.97      0.97      286
```

Figure 4.8b: Classification report of Gradient Boosting Classifier with hyperparameters (Diabetes)

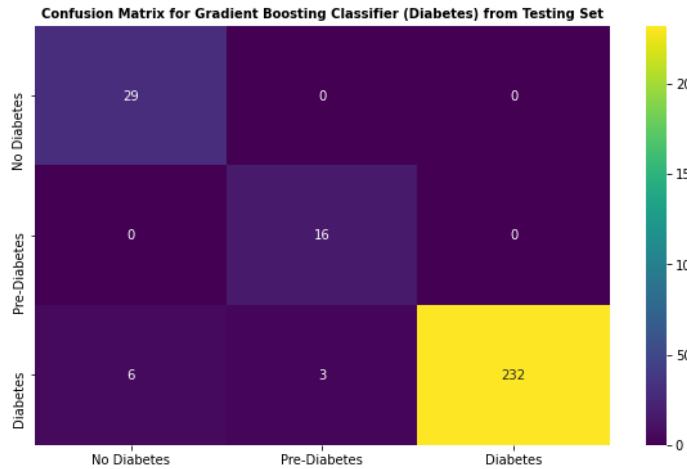


Figure 4.8c: Confusion matrix of Gradient Boosting Classifier with hyperparameters (Diabetes)

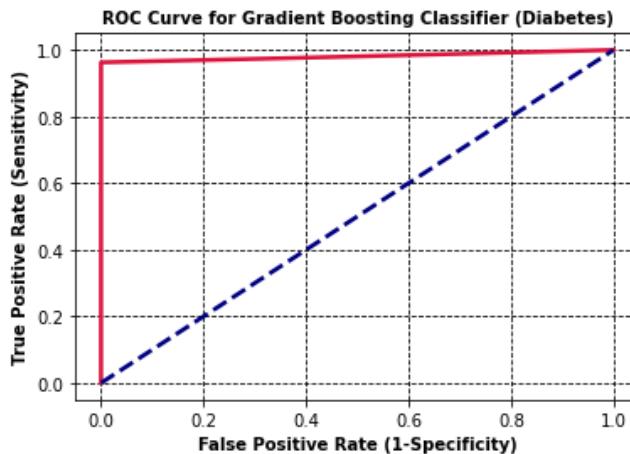


Figure 4.8d: ROC-AUC of Gradient Boosting Classifier with hyperparameters (Diabetes)

Based on the performance metrics and the classification report shown in Figure 4.8a and Figure 4.8b respectively, the gradient boosting classifier with hyperparameters has a very high accuracy rate of 97% as well as an excellent ROC-AUC score of about 99%, which means that this gradient boosting classifier with hyperparameters has the strong capability to distinguish between classes (i.e., no diabetes, pre-diabetes, or diabetes) and this can also be proven by the ROC-AUC shown in Figure 4.8d. Apart from the accuracy and ROC scores, it is also imperative to investigate other performance metrics as well. The precision score, recall score and f1 score shown in Figure 4.8a are actually based on macro-averaging. This means that the macro-averaging is performed by first computing the respective precision, recall, and f1-score of each class, and then taking the average of all the respective precision, recall and f1-scores. The reason of using macro-average is to ensure all classes

C3879C CAPSTONE PROJECT

---

are contributed equally regardless of how often they appear in the testing dataset. Based on the performance metrics as well as the classification report, it can see that the gradient boosting classifier with hyperparameters has done generally well with respect to the f1-score, precision and recall macro-averaging.

One can see that this is a multi-classification problem for diabetes. The computation of the TP, FP, TN, and FN in confusion matrix for multi-class is slightly indifferent from the binary classification. As such, the following below illustrate the calculation of the TP, FP, TN, and FN for multi-class.

**For Class 0: No Diabetes**

$$\text{TP} = 29$$

$$\text{FP} = 0 + 6 = 6$$

$$\text{TN} = 16 + 0 + 3 + 232 = 251$$

$$\text{FN} = 0 + 0 = 0$$

**For Class 1: Pre-Diabetes**

$$\text{TP} = 16$$

$$\text{FP} = 0 + 3 = 3$$

$$\text{TN} = 29 + 0 + 6 + 232 = 267$$

$$\text{FN} = 0 + 0 = 0$$

**For Class 2: Diabetes**

$$\text{TP} = 232$$

$$\text{FP} = 0 + 0 = 0$$

$$\text{TN} = 29 + 0 + 0 + 16 = 45$$

$$\text{FN} = 6 + 3 = 9$$

Based on the confusion matrix computation shown above, both **false positive (FP) rate** and **false-negative (FN) rate for all classes are generally low**, which is about less than 5%. Therefore, the low false positive rate also leading to a high precision rate and also a low false negative rate also leading to a high recall rate, which can be proven in the performance metrics shown in Figure 4.8a. Overall, this proposed classifier model is generally acceptable.

## C3879C CAPSTONE PROJECT

Figure 4.8e, Figure 4.8f, Figure 4.8g, and Figure 4.8h illustrate the outcome of the performance of the random forest classifier with hyperparameter based on **hypertension testing set**.

```
Performance Metrics of the Random Forest with Hyperparameter Tuning (Hypertension) using Testing Set:
-----
Accuracy Score: 0.9747899159663865
Precision Score: 0.9748417721518987
Recall Score: 0.9691994996873046
F1 Score: 0.971927341354093
ROC Score: 0.9691994996873046
```

Figure 4.8e: Performance metrics of Random Forest Classifier with hyperparameters (Hypertension)

Classification Report for Random Forest with Hyperparameter Tuning (Hypertension) using Testing Set:				
	precision	recall	f1-score	support
No Hypertension	0.97	0.95	0.96	41
Hypertension	0.97	0.99	0.98	78
accuracy			0.97	119
macro avg	0.97	0.97	0.97	119
weighted avg	0.97	0.97	0.97	119

Figure 4.8f: Classification report of Random Forest Classifier with hyperparameters (Hypertension)

**Confusion Matrix for Random Forest Classifier (Hypertension) from Testing Set**

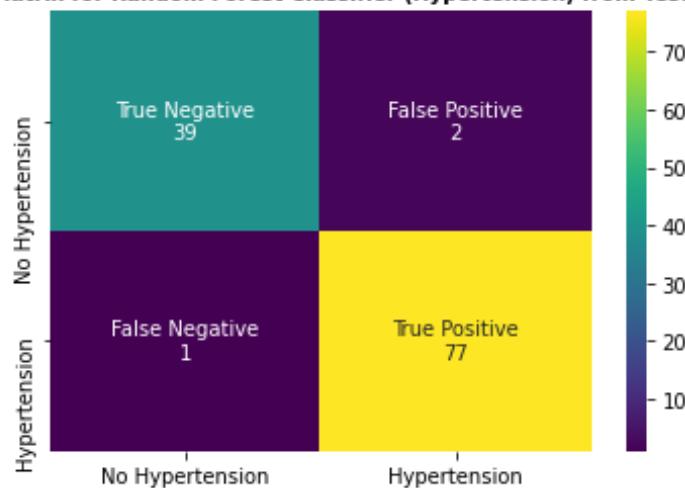


Figure 4.8g: Confusion matrix of Random Forest Classifier with hyperparameters (Hypertension)

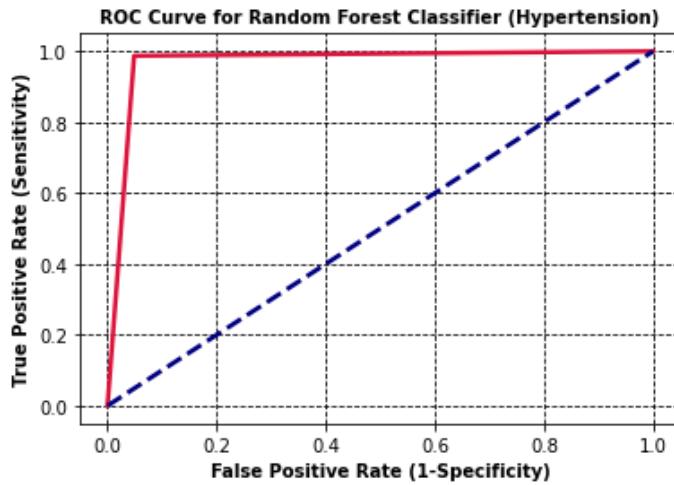


Figure 4.8h: ROC-AUC of Random Forest Classifier with hyperparameters  
(Hypertension)

Based on the performance metrics and the classification report shown in Figure 4.8e and Figure 4.8f respectively, the random forest classifier with hyperparameters has a very high accuracy rate of 97% as well as an excellent ROC-AUC score of about 97%, which means that this random forest classifier with hyperparameters has the strong capability to distinguish between classes (i.e., no hypertension or hypertension) and this can also be proven by the ROC-AUC shown in Figure 4.8h. As mentioned earlier, it is also imperative to investigate other performance metrics as well. The precision score, recall score and f1 score shown in Figure 4.8e are based on macro-averaging. Based on the performance metrics as well as the classification report, it can see that the random forest classifier with hyperparameters has done exceptionally well with respect to the scores of the f1-score, precision and recall macro-averaging.

Based on the confusion matrix shown in Figure 4.8g, both **false positive (FP) rate** and **false-negative (FN) rate** for **all classes are generally low**, which is about less than 5%. Therefore, the low false-positive rate leading to a high precision rate and also a low false negative-rate leading to a high recall rate which can be proven in the performance metrics shown in Figure 4.8e, and the proposed classifier model for hypertension is generally acceptable.

## C3879C CAPSTONE PROJECT

Figure 4.8i, Figure 4.8j, Figure 4.8k, and Figure 4.8l illustrate the outcome of the performance of the random forest classifier with hyperparameter based on **obesity testing set**.

```
Performance Metrics of the Random Forest with Hyperparameter Tuning (Obesity) using Testing Set:
```

```
=====
Accuracy Score: 0.9396551724137931
Precision Score: 0.9171024981257231
Recall Score: 0.9337673968188673
F1 Score: 0.9249661599872964
ROC Score: 0.9875320006419606
```

Figure 4.8i: Performance metrics of Random Forest Classifier with hyperparameters (Obesity)

```
Classification Report for Random Forest with Hyperparameter Tuning (Obesity) using Testing Set:
```

	precision	recall	f1-score	support
Underweight	0.96	0.97	0.96	68
Normal Weight	0.82	0.89	0.86	66
Overweight	0.91	0.91	0.91	160
Obesity	0.99	0.96	0.97	286
accuracy			0.94	580
macro avg	0.92	0.93	0.92	580
weighted avg	0.94	0.94	0.94	580

Figure 4.8j: Classification report of Random Forest Classifier with hyperparameters (Obesity)

Confusion Matrix for Random Forest Classifier (Obesity) from Testing Set

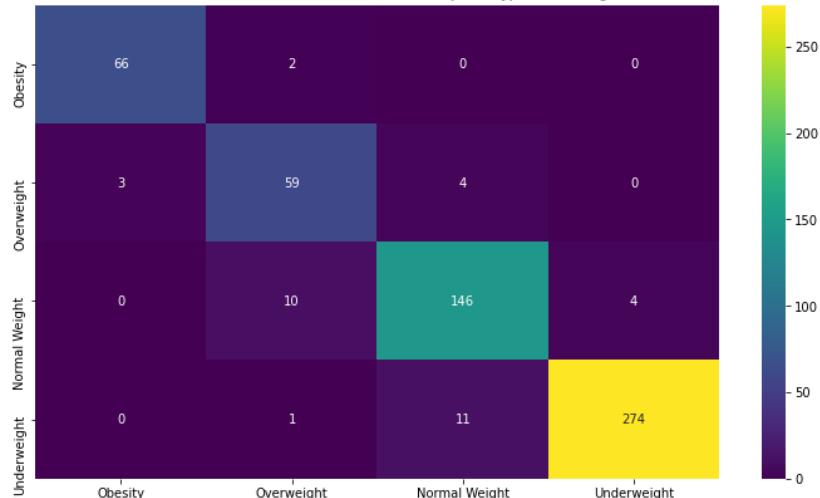


Figure 4.8k: Confusion matrix of Random Forest Classifier with hyperparameters (Obesity)

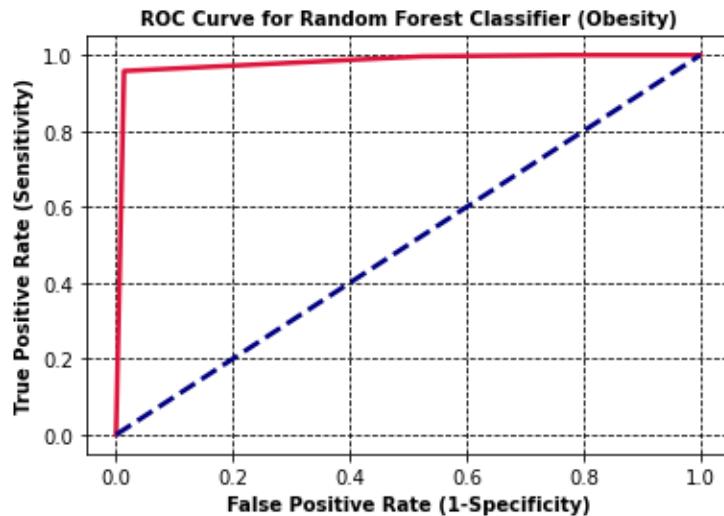


Figure 4.8l: ROC-AUC of Random Forest Classifier with hyperparameters  
(Obesity)

Based on the performance metrics and the classification report shown in Figure 4.8i and Figure 4.8j respectively, the random forest classifier with hyperparameters has a very high accuracy rate of 94% as well as an excellent ROC-AUC score of about 98%, which means that this random forest classifier with hyperparameters has the strong capability to distinguishes between classes (i.e., underweight, normal weight, overweight, or obesity) and this can also be proven by the ROC-AUC shown in Figure 4.8l. Furthermore, the precision score, recall score and f1 score shown in Figure 4.8i are actually based on macro-averaging. Based on the performance metrics as well as the classification report, this random forest classifier model has done generally well with respect to the scores of the f1-score, precision and recall macro-averaging.

Similar as diabetes, this is a multi-classification problem for obesity. As such, the following below illustrate the computation of the TP, FP, TN, and FN for multi-class.

### **For Class 1: Underweight**

$$TP = 66$$

$$FP = 0 + 0 + 3 = 3$$

$$TN = 59 + 4 + 0 + 10 + 146 + 4 + 1 + 11 + 274 = 509$$

$$FN = 2 + 0 + 0 = 2$$

**For Class 2: Normal Weight** $TP = 59$  $FP = 10 + 1 + 2 = 13$  $TN = 66 + 0 + 0 + 0 + 0 + 146 + 4 + 11 + 274 = 501$  $FN = 3 + 4 + 0 = 7$ **For Class 3: Overweight** $TP = 146$  $FP = 11 + 4 + 0 = 15$  $TN = 66 + 2 + 3 + 59 + 0 + 1 + 274 + 0 + 0 = 405$  $FN = 0 + 10 + 4 = 14$ **For Class 4: Obesity** $TP = 274$  $FP = 0 + 0 + 4 = 4$  $TN = 66 + 2 + 0 + 3 + 59 + 4 + 0 + 10 + 146 = 290$  $FN = 0 + 1 + 11 = 12$ 

Based on the confusion matrix computation shown above, both **false positive (FP) rate** and **false-negative (FN) rate** for **all classes are generally low**, which is about less than 5%. As such, the low false positive rate also leading to a high precision rate and also a low false negative rate also leading to a high recall rate, which can be proven in the performance metrics shown in Figure 4.8i. Thus, the proposed classifier model with hyperparameter for obesity is generally acceptable.

## 4.9 Development of Health Web Application using Streamlit

All the respective predictive algorithms that have been developed in Python using Spyder IDE were saved in the Pickle files and the web application that has been built for this project was using Streamlit. Thereafter, the proposed respective predictive algorithms have been integrated with the Streamlit app as its backbone. The development of this web app and disease predictive model could help users effectively minimize their hypertension, diabetes as well as obesity while all diseases are in their early stages. The outcome also included some recommendations such as try to achieve a healthy weight by balancing your caloric input or physical activity over the next 6 to 12 months and also consult your doctor as soon as possible, etc.

Figure 4.9a, Figure 4.9b, and Figure 4.9c illustrate the respective prediction outcome interface as well as recommendations for the **diabetes** via the Streamlit app.

### For Class 0: No Diabetes

The screenshot shows a Streamlit web application for "Diabetes Mellitus Prediction".

**Left Sidebar:**

- A dropdown menu titled "Please choose one to explore:" with "Diabetes Mellitus" selected.
- A section titled "What's so unique about this diabetes prediction app?" explaining the use of Gradient Boosting Algorithm (Supervised ML) with 90% accuracy, precision, recall, and AUC.
- An "Evaluation outcome from the testing set:" section showing a classification report for Gradient Boosting with Hyperparameter precision, recall, F1-score, and support. The report includes rows for No Diabetes, Prevalent Diabetes, and Diabetes.
- A note stating: "NOTE: The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase."
- A "What's Diabetes?" section with a colorful bar chart.

**Main Content Area:**

### Diabetes Mellitus Prediction

**Illustration:** An illustration of a doctor in a white coat interacting with a patient who is sitting on a yellow chair. A blood glucose meter displays a reading of 8.5. Medical icons like a shield with a cross, a stethoscope, and pills are scattered around.

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**Right Sidebar:**

- General guideline(s) to the user:**
  - You are **required** to fill up all the information in this form in less than 5 mins.
  - Some information may required your blood test results.
- A form for entering personal information: "Enter your name:" (Jefferson Wong) and "Enter your identification no.:".

## C3879C CAPSTONE PROJECT

**What's Diabetes?**

**Diabetes** is a chronic health condition that affects how your body generate food into energy.

Most of the food you eat is broken down into sugar and released into your bloodstream. When your blood sugar elevated, it signals your pancreas to produce insulin. Insulin acts like a key to let the blood sugar into your body cells to use it as energy.

If you have diabetes, your body either not producing enough insulin or cannot use the insulin it makes as well as it should. When there is not enough insulin or cell stop responding to insulin, too much blood sugar remains in your bloodstream. Overtime, that can cause serious health problems such as vision loss, or kidney disease, etc.

You may click the link below to know more about diabetes.  
<https://www.cdc.gov/diabetes/basics/diab...>

**How can we control or prevent**

Enter your name:  
Jefferson Wong

Enter your identification no.:  
TA101

Gender:  
 Male  
 Female

Enter your age:  
28

Enter your body mass index (BMI):  
19.50

Enter the amount of fasting blood sugar level (in mmol/L):  
4.5

**Predict**

**Outcome:**  
Hi Jefferson Wong! You have a [51.79532895)% chance at low risk of developing diabetes.

**Recommendation:**  
Please continue to achieve a healthy weight by balancing your caloric input and physical activity.

Figure 4.9a: A sample interface of the prediction outcome and recommendation in Streamlit app for **no diabetes**

**For Class 1: Pre-Diabetes**

Please choose one to explore:  
Diabetes Mellitus

**What's so unique about this diabetes prediction app?**

This application uses **Gradient Boosting Algorithm** (Supervised ML) with the best hyperparameters to predict the likelihood of developing diabetes in an individual. During the testing phase, the outcome shows that the gradient boosting algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC**.

**Evaluation outcome from the testing set:**

	precision	recall	f1-score	support
No Diabetes	0.83	1.00	0.91	29
Pre-Diabetes	0.83	0.99	0.90	29
Diabetes	1.00	0.99	0.99	243

accuracy      0.97      286  
precision      0.97      0.97      286  
recall      0.97      0.97      0.97      286  
f1-score      0.97      0.97      0.97      286

Based on the testing set, the ROC-Score for Gradient Boosting

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Diabetes?**

### Diabetes Mellitus Prediction

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

- 1) You are **required** to fill up all the information in this form in less than 5 mins.
- 2) Some information may required your blood test results.

Enter your name:  
May Phua

Enter your identification no.:

## C3879C CAPSTONE PROJECT

The screenshot shows a Streamlit application window. On the left, there's a sidebar with a section titled "How can we control or prevent the chances of developing diabetes?" featuring a list of tips like "Lose extra weight" and "Be more physically active". Below this is a "BMI Calculator" section. The main panel contains input fields for "Enter your name:" (May Phua), "Enter your identification no.:" (PA101), "Gender" (Female), "Enter your age:" (43), "Enter your body mass index (BMI):" (28.50), and "Enter the amount of fasting blood sugar level (in mmol/L):" (5.8). A "Predict" button is present. The "Outcome:" section displays a message: "Hi May Phua! You have a [99.7806879]% chance at moderate risk of developing pre-diabetes." The "Recommendation:" section suggests: "Please try to achieve a healthy weight by balancing your caloric input and physical activity over the next 6 to 12 months." At the bottom, the Windows taskbar shows the date as 16/9/2021.

Figure 4.9b: A sample interface of the prediction outcome and recommendation in Streamlit app for **pre-diabetes**

## For Class 2: Diabetes

This screenshot shows a Streamlit application for "Diabetes Mellitus Prediction". The sidebar includes a dropdown menu for "Please choose one to explore:" set to "Diabetes Mellitus" and a section titled "What's so unique about this diabetes prediction app?". It explains that the app uses Gradient Boosting Algorithm (Supervised ML) with 90% accuracy, precision, recall, and AUC. The main panel features a heading "Diabetes Mellitus Prediction" and an illustration of a doctor interacting with a patient. A "Classification Report for Gradient Boosting with Hyperparameters" table is shown:

	precision	recall	f1-score	support
No Diabetes	0.83	1.00	0.91	29
Pre-diabetes	0.83	0.93	0.77	23
Diabetes	1.00	0.90	0.90	241
<b>accuracy</b>	<b>0.99</b>	<b>0.99</b>	<b>0.97</b>	<b>286</b>
<b>macro avg</b>	<b>0.99</b>	<b>0.99</b>	<b>0.97</b>	<b>286</b>
<b>weighted avg</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>286</b>

A note states: "Based on the testing set, the ROC-Auc score for Gradient Boosting is 0.99". The "NOTE" section clarifies the testing set's purpose. The "What's Diabetes?" section includes a colorful icon of a person with a stethoscope.

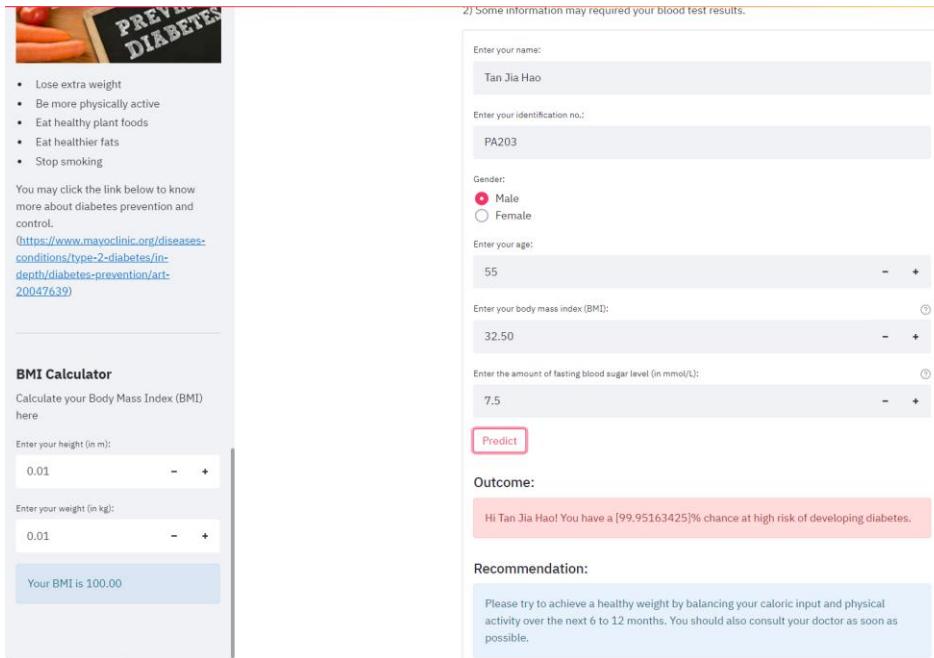
**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

- 1) You are **required** to fill up all the information in this form in less than 5 mins.
- 2) Some information may required your blood test results.

Input fields for "Enter your name:" (Tan Jia Hao) and "Enter your identification no.:" are also visible.

## C3879C CAPSTONE PROJECT



The screenshot shows a Streamlit application interface. On the left, there's a sidebar with a "PREVENT DIABETES" logo and a list of prevention tips. Below that is a "BMI Calculator" section where users can input height and weight to calculate their BMI. On the right, the main panel displays user information (name, ID, gender, age, BMI, and fasting blood sugar level), a "Predict" button, and outcome/recommendation sections.

**Prevention Tips:**

- Lose extra weight
- Be more physically active
- Eat healthy plant foods
- Eat healthier fats
- Stop smoking

You may click the link below to know more about diabetes prevention and control.  
<https://www.mayoclinic.org/diseases-conditions/type-2-diabetes/in-depth/diabetes-prevention/art-20047639>

**BMI Calculator**

Calculate your Body Mass Index (BMI) here

Enter your height (in m): 0.01

Enter your weight (in kg): 0.01

Your BMI is 100.00

2) Some information may required your blood test results.

Enter your name: Tan Jia Hao

Enter your identification no.: PA203

Gender: Male

Enter your age: 55

Enter your body mass index (BMI): 32.50

Enter the amount of fasting blood sugar level (in mmol/L): 7.5

Predict

**Outcome:**  
Hi Tan Jia Hao! You have a [99.95163425)% chance at high risk of developing diabetes.

**Recommendation:**  
Please try to achieve a healthy weight by balancing your caloric input and physical activity over the next 6 to 12 months. You should also consult your doctor as soon as possible.

Figure 4.9c: A sample interface of the prediction outcome and recommendation in Streamlit app for diabetes

**NOTE: Gender, BMI, and Fasting Blood Sugar Level are the attributes for the predictive process.**

## C3879C CAPSTONE PROJECT

Figure 4.9d, and Figure 4.9e illustrate the respective prediction outcome interface as well as recommendations for the **hypertension** via the Streamlit app.

**For Class 0: No hypertension**

Please choose one to explore:

Hypertension

**What's so unique about this hypertension prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the likelihood of developing hypertension in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC**.

**Evaluation outcome from the testing set:**

	precision	recall	f1-score	support
No Hypertension	0.99	0.99	0.99	28
Hypertension	0.97	0.99	0.98	28
<b>accuracy</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>135</b>
<b>macro avg</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>135</b>
<b>weighted avg</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>135</b>

Based on the testing set, the ROC-Score for Random Forest is 0.97.

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

### Hypertension Prediction



**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General instruction to the user:**

- 1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:  
Jefferson Wong

Enter your identification no.:  
TA103

**What's Hypertension?**



**Hypertension or high blood pressure** is a blood pressure that is higher than normal blood pressure. Your blood pressure changes throughout the day based on the level of activities. However, having blood pressure that is constantly above the normal level may result in the diagnosis of high blood pressure.

High blood pressure usually develops over time and it can happen because of the unhealthy lifestyle choices or certain health conditions such as diabetes and having obesity, can also increase the chance of developing high blood pressure.

You may click the link below to know more about hypertension.  
<https://www.cdc.gov/bloodpressure/about.html>

Enter your identification no.:  
TA103

Gender:  
 Male  
 Female

Enter your age:  
23

Enter your body mass index (BMI):  
22.50

Enter the systolic blood pressure (in mm Hg):  
135

Enter the diastolic blood pressure (in mm Hg):  
85

**Predict**

**Outcome:**  
Hi Jefferson Wong! You have a [92.85%] chance at low risk of developing hypertension

**Recommendation:**  
Please continue to achieve a healthy weight by balancing your caloric input and physical activity.

Figure 4.9d: A sample interface of the prediction outcome and recommendation in Streamlit app for **no hypertension**

## For Class 1: Hypertension

Please choose one to explore:

Hypertension

**What's so unique about this hypertension prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the likelihood of developing hypertension in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC.**

**Evaluation outcome from the testing set:**

	precision	recall	f1-score	support
No Hypertension	0.97	0.95	0.96	45
Hypertension	0.97	0.99	0.98	74
accuracy	0.97	0.97	0.97	119
macro avg	0.97	0.97	0.97	119
weighted avg	0.97	0.97	0.97	119

Based on the testing set, the ROC-Score for Random Forest

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Hypertension?**

### Hypertension Prediction



**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General instruction to the user:**

- 1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:  
May Phuu

Enter your identification no.:  
PA103

Enter your identification no.:  
PA103

Gender:  
 Male  
 Female

Enter your age:  
55

Enter your body mass index (BMI):  
32.50

Enter the systolic blood pressure (in mm Hg):  
155

Enter the diastolic blood pressure (in mm Hg):  
120

**Predict**

**Outcome:**  
Hi May Phuu! You have a [84.56097884]% chance at high risk of developing hypertension.

**Recommendation:**  
Please try to achieve a healthy weight by balancing your caloric input and physical activity over the next 6 to 12 months. You should also consult your doctor as soon as possible.

Figure 4.9e: A sample interface of the prediction outcome and recommendation in Streamlit app for hypertension

**NOTE: BMI, Systolic Blood Pressure (SBP), and Diastolic Blood Pressure (DBP) are the attributes for the predictive process.**

## C3879C CAPSTONE PROJECT

Figure 4.9f, Figure 4.9g, Figure 4.9h, and Figure 4.9i illustrate the respective prediction outcome interface as well as recommendations for the **obesity classification** via the Streamlit app.

### For Class 1: At low-risk for obesity-related diseases (underweight)

Please choose one to explore:

Obesity

**What's so unique about this obesity prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the obesity level in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC**.

Evaluation outcome from the testing set:

	precision	recall	F1-score	Support
Underweight	0.96	0.97	0.96	500
Normal weight	0.82	0.89	0.86	600
Overweight	0.97	0.91	0.93	2000
Obesity	0.99	0.98	0.97	2000
accuracy			0.94	5000
macro avg	0.92	0.93	0.92	5000
weighted avg	0.94	0.94	0.94	5000

Based on the testing set, the ROC-Score for Random Forest is 0.94.

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Obesity?**

Obesity can be defined as the weight that is higher than what is considered healthy for a given weight. BMI is a main screening tool for obesity.

Obesity is a complex health issue resulting from a combination of causes and individual factors such as genetics and behaviour. Behaviours can include physical activity, inactivity, dietary patterns, and other exposures.

Obesity is serious because it is often associated with poor mental health outcomes and reduced quality of life. Furthermore, obesity is also associated with the leading causes of death, including diabetes and high blood pressure.

You may click the link below to know more about obesity.  
<https://www.cdc.gov/obesity/adult/causes.html>

### Obesity Prediction

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:

Enter your identification no.:

Enter your age:

-
+

Enter your height (in m):

-
+

Enter your weight (in kg):

-
+

Does your family have a history of obesity?

Yes

No

Do you often consume high caloric foods?

Yes

No

How often do you consume fruits and vegetables?

I do not eat fruits and vegetables.

Frequently

Often

What is the number of main meals per day?

1

2

3

4

How often do you consume foods between meals?

I do not consume foods between meals.

Sometimes

Frequently

Always

Do you smoke?

Yes

No

## C3879C CAPSTONE PROJECT

<p>You may click the link below to know more about obesity. <a href="https://www.cdc.gov/obesity/adult/causes.html">https://www.cdc.gov/obesity/adult/causes.html</a></p> <p><b>How can we manage and prevent the chances of developing obesity?</b></p>  <ul style="list-style-type: none"> <li>• Have a balanced and calorie-controlled diet as recommended by GP or weight loss management health professional</li> <li>• Be more physically active</li> <li>• Eat healthy plant foods</li> <li>• Exercise regularly for 150 to 300 minutes</li> <li>• Less smoking</li> </ul> <p>You may click the link below to know more about management and prevention of obesity. <a href="https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity">https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity</a></p>	<p>Do you smoke?  <input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you drink water daily?  <input checked="" type="radio"/> I do not always drink water.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently</p> <p>Do you often monitor your own calories?  <input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you exercise?  <input checked="" type="radio"/> I do not exercise.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>How often do you use your electronic devices?  <input type="radio"/> Sometimes  <input checked="" type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you drink alcohol?  <input type="radio"/> I do not drink alcohol.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p><b>Predict</b></p> <p><b>Outcomes:</b></p> <div style="background-color: #ffffcc; padding: 5px;">           Hi Jefferson Wong! You have about [75.87822119]% chance at low risk for obesity-related diseases.         </div> <p><b>Recommendation:</b></p> <div style="background-color: #e0f2ff; padding: 5px;">           However, you are at risk of nutritional deficiency and osteoporosis. You are encouraged to eat a balanced meal and to seek medical advice if necessary.         </div>
---	---

Figure 4.9f: A sample interface of the prediction outcome and recommendation in Streamlit app for **at low-risk of having obesity-related diseases (underweight)**

## For Class 2: At low-risk for obesity-related diseases (normal weight)

Please choose one to explore:

Obesity

**What's so unique about this obesity prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the obesity level in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC**.

Evaluation outcome from the testing set:

	precision	recall	f1-score	support
Underweight	0.97	0.96	0.96	100
Normal Weight	0.89	0.93	0.91	100
Overweight	0.93	0.91	0.91	100
Obesity	0.99	0.97	0.97	200
accuracy	0.94	0.94	0.94	500
macro avg	0.92	0.91	0.91	500
weighted avg	0.94	0.94	0.94	500

Based on the testing set, the ROC-Score for Random Forest

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Obesity?**

### Obesity Prediction

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:

May Phua

Enter your identification no.:

TA103

Gender:

Male

Female

Enter your age:

29

Enter your height (in m):

1.60

Enter your weight (in kg):

56.00

Does your family have a history of obesity?

Yes

No

Do you often consume high caloric foods?

Yes

No

How often do you consume fruits and vegetables?

I do not eat fruits and vegetables.

Frequently

Often

What is the number of main meals per day?

1

2

3

4

How often do you consume foods between meals?

I do not consume foods between meals.

Sometimes

Frequently

...

## C3879C CAPSTONE PROJECT

<p><b>What's Obesity?</b></p>  <p>Obesity can be defined as the weight that is higher than what is considered healthy for a given weight. BMI is a main screening tool for obesity.</p> <p>Obesity is a complex health issue resulting from a combination of causes and individual factors such as genetics and behaviour. Behaviours can include physical activity, inactivity, dietary patterns, and other exposures.</p> <p>Obesity is serious because it is often associated with poor mental health outcomes and reduced quality of life. Furthermore, obesity is also associated with the leading causes of death, including diabetes and high blood pressure.</p> <p>You may click the link below to know more about obesity.  <a href="https://www.cdc.gov/obesity/adult/causes.html">https://www.cdc.gov/obesity/adult/causes.html</a></p>	<p>Do you smoke?  <input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you drink water daily?  <input checked="" type="radio"/> I do not always drink water.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently</p> <p>Do you often monitor your own calories?  <input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you exercise?  <input checked="" type="radio"/> I do not exercise.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>How often do you use your electronic devices?  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you drink alcohol?  <input type="radio"/> I do not drink alcohol.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p><b>Predict</b></p> <p><b>Outcomes:</b></p> <div style="background-color: #c8f7e4; padding: 5px;"> Hi May Phual You have about [62.13517619)% chance at low risk for obesity-related diseases. </div>
<p>You may click the link below to know more about obesity.  <a href="https://www.cdc.gov/obesity/adult/causes.html">https://www.cdc.gov/obesity/adult/causes.html</a></p> <p><b>How can we manage and prevent the chances of developing obesity?</b></p>  <ul style="list-style-type: none"> <li>• Have a balanced and calorie-controlled diet as recommended by GP or weight loss management health professional</li> <li>• Be more physically active</li> <li>• Eat healthy plant foods</li> <li>• Exercise regularly for 150 to 300 minutes</li> <li>• Less smoking</li> </ul> <p>You may click the link below to know more about management and prevention of obesity.  <a href="https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity">https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity</a></p>	<p><input checked="" type="radio"/> I do not exercise.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>How often do you use your electronic devices?  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you drink alcohol?  <input type="radio"/> I do not drink alcohol.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p><b>Predict</b></p> <p><b>Outcomes:</b></p> <div style="background-color: #c8f7e4; padding: 5px;"> Hi May Phual You have about [62.13517619)% chance at low risk for obesity-related diseases. </div> <p><b>Recommendation:</b></p> <div style="background-color: #d9eaf7; padding: 5px;"> Please continue to achieve a healthy weight by balancing your caloric input and physical activity. </div>

Made with Streamlit

Figure 4.9g: A sample interface of the prediction outcome and recommendation in Streamlit app for **at low-risk of having obesity-related diseases (normal weight)**

## C3879C CAPSTONE PROJECT

**For Class 3: At moderate-risk for obesity-related diseases (overweight)**

Please choose one to explore:

Obesity

**What's so unique about this obesity prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the obesity level in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC.**

Evaluation outcome from the testing set:

	precision	recall	f1-score	support
Underweight	0.97	0.97	0.96	500
Normal Weight	0.89	0.89	0.88	500
Overweight	0.91	0.91	0.91	1000
Obesity	0.99	0.96	0.97	2000
accuracy	0.92	0.94	0.93	5000
macro avg	0.92	0.91	0.91	5000
weighted avg	0.94	0.94	0.94	5000

Based on the testing set, the ROC-Score for Random Forest

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Obesity?**

Press F11 to exit full screen

## Obesity Prediction

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:  
Tan Jia Hao

Enter your identification no.:  
PA435

Gender:  
 Male  
 Female

Enter your age:  
20

Enter your height (in m):  
1.76

Enter your weight (in kg):  
87.90

Does your family have a history of obesity?  
 Yes  
 No

Do you often consume high caloric foods?  
 Yes  
 No

How often do you consume fruits and vegetables?  
 I do not eat fruits and vegetables.  
 Frequently  
 Often

What is the number of main meals per day?  
 1  
 2  
 3  
 4

How often do you consume foods between meals?  
 I do not consume foods between meals.  
 Sometimes  
 Frequently

## C3879C CAPSTONE PROJECT

**What's Obesity?**



Obesity can be defined as the weight that is higher than what is considered healthy for a given weight. BMI is a main screening tool for obesity.

Obesity is a complex health issue resulting from a combination of causes and individual factors such as genetics and behaviour. Behaviours can include physical activity, inactivity, dietary patterns, and other exposures.

Obesity is serious because it is often associated with poor mental health outcomes and reduced quality of life. Furthermore, obesity is also associated with the leading causes of death, including diabetes and high blood pressure.

You may click the link below to know more about obesity.  
<https://www.cdc.gov/obesity/adult/causes.html>

**How can we manage and prevent the chances of developing obesity?**



You may click the link below to know more about obesity.  
<https://www.cdc.gov/obesity/adult/causes.html>

**How often do you consume foods between meals?**

I do not consume foods between meals.  
 Sometimes  
 Frequently  
 Always

**Do you smoke?**

Yes  
 No

**How often do you drink water daily?**

I do not always drink water.  
 Sometimes  
 Frequently

**Do you often monitor your own calories?**

Yes  
 No

**How often do you exercise?**

I do not exercise.  
 Sometimes  
 Frequently  
 Always

**How often do you use your electronic devices?**

Sometimes  
 Frequently  
 Always

**Do you drink alcohol?**

I do not drink alcohol.  
 Sometimes  
 Frequently  
 Always

**Predict**

**Outcomes:**

I do not exercise.  
 Sometimes  
 Frequently  
 Always

**How often do you use your electronic devices?**

Sometimes  
 Frequently  
 Always

**Do you drink alcohol?**

I do not drink alcohol.  
 Sometimes  
 Frequently  
 Always

**Predict**

**Outcomes:**

Hi Tan Jia Hao! You have about [57.13606829)% chance at moderate risk for obesity-related diseases.

**Recommendation:**

Please try to aim to lose at least 5% to 10% of your body weight over 6 to 12 months by increasing your physical activity and reducing caloric intake. You are encouraged to use this app to check for possible risks of developing diabetes and hypertension.

Made with Streamlit

Figure 4.9h: A sample interface of the prediction outcome and recommendation in Streamlit app for **at moderate-risk of having obesity-related diseases (overweight)**

## For Class 4: At high-risk for obesity-related diseases (obesity)

Please choose one to explore:

Obesity

**What's so unique about this obesity prediction app?**

This application uses **Random Forest Algorithm** (Supervised ML) with the best hyperparameters to predict the obesity level in an individual. During the testing phase, the outcome shows that the random forest algorithm can be achievable **with at least 90% accuracy, precision, recall, and AUC**.

Evaluation outcome from the testing set:

	precision	recall	f1-score	support
Underweight	0.97	0.96	0.96	100
Normal weight	0.89	0.91	0.90	100
Overweight	0.93	0.91	0.92	100
Obesity	0.99	0.96	0.97	200
accuracy	0.92	0.94	0.93	500
macro avg	0.92	0.93	0.92	500
weighted avg	0.94	0.94	0.94	500

Based on the testing set, the ROC-Score for Random Forest

**NOTE:** The objective of the testing set is to evaluate the performance of the trained model and was unseen during the training phase.

**What's Obesity?**

**Obesity** can be defined as the weight that is higher than what is considered healthy for a given weight. BMI is a main screening tool for obesity.

Obesity is a complex health issue resulting from a combination of causes and individual factors such as genetics and behaviour. Behaviours can include physical activity, inactivity, dietary patterns, and other exposures.

Obesity is serious because it is often associated with poor mental health outcomes and reduced quality of life. Furthermore, obesity is also associated with the leading causes of death, including diabetes and high blood pressure.

You may click the link below to know more about obesity.  
<https://www.cdc.gov/obesity/adult/causes.html>

### Obesity Prediction

**Medical Disclaimer:** This platform is not serve as an alternative to medical advice from medical professional healthcare provider. If you have any specific questions about any medical matter, you should consult your doctor or other medical professional healthcare provider.

**General guideline(s) to the user:**

1) You are **required** to fill up all the information in this form in less than 5 mins.

Enter your name:  
Gina Tan

Enter your identification no.:  
TA234

Gender:  
 Male  
 Female

Enter your age:  
38

Enter your height (in m):  
1.62

Enter your weight (in kg):  
80.00

Does your family have a history of obesity?  
 Yes  
 No

Do you often consume high caloric foods?  
 Yes  
 No

How often do you consume fruits and vegetables?  
 I do not eat fruits and vegetables.  
 Frequently  
 Often

What is the number of main meals per day?  
 1  
 2  
 3  
 4

How often do you consume foods between meals?  
 I do not consume foods between meals.  
 Sometimes  
 Frequently  
 ...

## C3879C CAPSTONE PROJECT

<p>You may click the link below to know more about obesity. <a href="https://www.cdc.gov/obesity/adult/causes.html">https://www.cdc.gov/obesity/adult/causes.html</a></p> <p><b>How can we manage and prevent the chances of developing obesity?</b></p>  <ul style="list-style-type: none"> <li>• Have a balanced and calorie-controlled diet as recommended by GP or weight loss management health professional</li> <li>• Be more physically active</li> <li>• Eat healthy plant foods</li> <li>• Exercise regularly for 150 to 300 minutes</li> <li>• Less smoking</li> </ul> <p>You may click the link below to know more about management and prevention of obesity. <a href="https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity">https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity</a></p>	<p>How often do you consume foods between meals?</p> <p><input type="radio"/> I do not consume foods between meals.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you smoke?</p> <p><input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you drink water daily?</p> <p><input type="radio"/> I do not always drink water.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently</p> <p>Do you often monitor your own calories?</p> <p><input type="radio"/> Yes  <input checked="" type="radio"/> No</p> <p>How often do you exercise?</p> <p><input checked="" type="radio"/> I do not exercise.  <input type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>How often do you use your electronic devices?</p> <p><input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you drink alcohol?</p> <p><input type="radio"/> I do not drink alcohol.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p><b>Predict</b></p> <p><b>Outcomes:</b></p> <p><input type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>How often do you use your electronic devices?</p> <p><input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p>Do you drink alcohol?</p> <p><input type="radio"/> I do not drink alcohol.  <input checked="" type="radio"/> Sometimes  <input type="radio"/> Frequently  <input type="radio"/> Always</p> <p><b>Predict</b></p> <p><b>Outcomes:</b></p> <div style="background-color: #f8d7da; padding: 5px;"> <p>Hi Gina Tan! You have about [55.89949949]% chance at high risk for obesity-related diseases.</p> </div> <p><b>Recommendation:</b></p> <div style="background-color: #e1f5fe; padding: 5px;"> <p>Please try to aim to loss at least 5% to 10% of your body weight over 6 to 12 months by increasing your physical activity and reducing caloric intake. You should go for regular health screening to keep co-morbid conditions in check or use this app to check for possible risks of developing diabetes and hypertension.</p> </div>
---	---

Figure 4.9i: A sample interface of the prediction outcome and recommendation in Streamlit app for **at high-risk of having obesity-related diseases (obesity)**

## 5.0 Conclusion

The development of the respective predictive models for respective obesity-related chronic diseases involved in the integration of the isolation forest to detect the potential outliers present in the original datasets, the SMOTE-TOMEK to balanced the training sets to minimize the bias between classes, and the investigation of the baseline models and the hyperparameters optimization process for the proposed machine learning models. The performance evaluation for the respective predictive models for respective obesity-related chronic diseases was done using the unseen testing sets. It was observed that all of the proposed predictive models for respective chronic diseases such as diabetes, hypertension, and obesity were found to have at least 90% for accuracy, at least 95% for ROC-AUC, and less than 5% of false-positive rate and false-negative rate based on the confusion matrix computation.

To simulate a real-life practical application, it has also developed a simple and friendly web application via Streamlit with the integration of the proposed predictive models with optimal hyperparameters for the respective obesity-related chronic diseases such as obesity, diabetes, and hypertension to diagnosed the user at their convenience. The web application comprises of respective obesity-related chronic diseases platform whereby the respective platform collates all the user risk factors input and then transmitted to the respective predictive model with optimal hyperparameters to performed the prediction. Hence, the class outcome with the highest probability represents the final outcome diagnosis of the respective obesity-related chronic diseases. The final outcome of the respective diagnosis appeared on the user's web app inclusive of the general health recommendation, which would allow them to take prompt and appropriate action to minimize and/ or prevent the risks further.

## 6.0 Recommendation

Future studies or expansion will include a more adequate large amount of data, more relevant attributes, and features, the implementation of the architecture design of the neural network to predict obesity-related chronic diseases, and enhancement of the web application based on a user perspective.

Firstly, the collection of a huge amount of data not only enhances the knowledge base of the predictive model but also improves the accurateness of the system to a more reliable extent.

Secondly, apart from the medical factors like BMI, blood sugar level, blood pressure level, etc., it will also put into consideration of other high-risk factors that have been found in both diabetes and hypertension in recent research studies such as the duration of exercise, dietary habits, smoking habits, and drinking of alcohol habits for the predictive models. Furthermore, different age groups could also be further investigated.

Thirdly, it will also be implementing and architecting the prediction systems using some kinds of neural networks such as ANN, etc., to predict the respective obesity-related chronic diseases as the neural networks are generally considerably better than the machine learning classifiers.

Lastly, an extensive evaluation of the usability and user experience of the web application should also be carried out in the near future. Furthermore, some additional feature(s) will be can also be included in the existing web application such as the transmission of the user's personal info with disease predictive outcomes to their respective professional healthcare providers. This will benefit to those users with unexpected prediction outcomes or possibly under threat of both diabetes and hypertension.

## 7.0 Project Implementation Scheduled

s/n	Task(s)	Scheduled Date								
		27/6/2021 to 4/7/2021	5/7/2021 to 11/7/2021	12/7/2021 to 18/7/2021	19/7/2021 to 25/7/2021	26/7/2021 to 1/8/2021	2/8/2021 to 8/8/2021	9/8/2021 to 15/8/2021	16/8/2021 to 22/8/2021	23/8/2021 to 29/8/2021
1	Collect of 1-dimensional datasets from various sources	Planned								
2	Store the 1-dimensional datasets into SQLite3 via Python	Planned	Actual							
3	Data Exploration on selected features for all respective obesity-related chronic diseases			Actual	Actual					
4	Data Cleansing (Investigate Missing Data) for all Datasets (Obesity, Diabetes, and Hypertension)				Actual	Actual				
5	Data Cleansing (Removal of Potential Outlier Detection using Isolation Forest) for all Datasets (Obesity, Diabetes, and Hypertension)			Planned	Actual	Actual				
6	Resampling techniques for imbalanced training datasets using SMOTE-Tomek		Planned			Actual	Actual			
7	Data Splitting, Model Selection (Baseline) and Evaluation via Shuffle-Split Cross-Validation for all respective Balanced Training Sets (Obesity, Diabetes, and Hypertension)				Actual	Actual	Actual			
8	Hyperparameters Optimization for Proposed Baseline Classifier Models				Actual	Actual	Actual	Actual		
9	Re-trained and Re-evaluation of the Proposed Classifier Models with the Best Hyperparameters via Shuffle-Split Cross-Validation						Actual	Actual	Actual	
10	Performance Evaluation of the Proposed Classifier Models with Hyperparameters using Testing Sets						Actual	Actual	Actual	
11	Development of Web App using Streamlit for Respective Obesity-related Chronic Diseases (Obesity, Diabetes, and Hypertension) with Integration of Predictive Modelling									
12	White-Box Testing on Web App									
13	Capstone Report Writing and Poster					Actual	Actual	Actual	Actual	
						Planned	Actual	Actual	Actual	

C3879C CAPSTONE PROJECT

## 8.0 References

- [1] Rashid, A. (2020, Jul 18). Diabetes Dataset (Mendeley Data, version 1). Retrieved from <https://data.mendeley.com/datasets/wj9rwkp9c2/1>
- [2] Machine Learning Examples (Neural Designer). Retrieved from <https://www.neuraldesigner.com/learning/examples/Obesity-level>
- [3] Golino, H. (2013, Nov 11). Men's dataset from the "Predicting increased blood pressure using Machine Learning" paper. Retrieved from [https://figshare.com/articles/dataset/Men\\_s\\_dataset\\_from\\_the\\_Predicting\\_increased\\_blood\\_pressure\\_using\\_Machine\\_Learning\\_paper/845665/1](https://figshare.com/articles/dataset/Men_s_dataset_from_the_Predicting_increased_blood_pressure_using_Machine_Learning_paper/845665/1)
- [4] Golino, H. (2013, Nov 11). Women's dataset from the "Predicting increased blood pressure using Machine Learning" paper. Retrieved from [https://figshare.com/articles/dataset/Women\\_s\\_dataset\\_from\\_the\\_Predicting\\_increased\\_blood\\_pressure\\_using\\_Machine\\_Learning\\_paper/845664/1](https://figshare.com/articles/dataset/Women_s_dataset_from_the_Predicting_increased_blood_pressure_using_Machine_Learning_paper/845664/1)
- [5] Fitriyano, N.L. Syafdrudin, M. Alfian, G. Rhee, J. (2019, Oct 17). Development of Disease Prediction Model Based on Ensemble Learning Approach for Diabetes and Hypertension. Retrieved from [https://www.researchgate.net/publication/336223791\\_Development\\_of\\_Disease\\_Prediction\\_Model\\_Based\\_on\\_Ensemble\\_Learning\\_Approach\\_for\\_Diabetes\\_and\\_Hypertension](https://www.researchgate.net/publication/336223791_Development_of_Disease_Prediction_Model_Based_on_Ensemble_Learning_Approach_for_Diabetes_and_Hypertension)
- [6] Streamlit Documentation. Retrieved from [https://docs.streamlit.io/en/stable/getting\\_started.html](https://docs.streamlit.io/en/stable/getting_started.html)
- [7] [Scikit-learn: Machine Learning in Python](#), Pedregosa *et al.*, JMLR 12, pp. 2825-2830, 2011.
- [8] Dhiraj, K. (2019). Anomaly Detection Using Isolation Forest in Python. Retrieved from <https://blog.paperspace.com/anomaly-detection-isolation-forest/>
- [9] Brownlee, J. (2020, Jan 22). How to Combine Oversampling and Undersampling for Imbalanced Classification. Retrieved from <https://machinelearningmastery.com/combine-oversampling-andundersampling-for-imbalanced-classification/>
- [10] Satpathy, S. (2020, Oct 6). Overcoming Class Imbalance using SMOTE Techniques. Retrieved from <https://www.analyticsvidhya.com/blog/2020/10/overcoming-class-imbalance-using-smote-techniques/>

C3879C CAPSTONE PROJECT

---

- [11] API Reference. Retrieved from <https://imbalanced-learn.org/dev/references/generated/imblearn.combine.SMOTETomek.html>
- [12] Scikit Learn on Gradient Boosting Classifier. Retrieved from <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html>
- [13] Scikit Learn on Random Forest Classifier. Retrieved from <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- [14] Scikit Learn on Extra Trees Classifier. Retrieved from <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html>
- [15] Khan, N.S. Muaz, M.H. Kabir, A. Islam, M.N. (December 2017). Diabetes Predicting mHealth Application Using Machine Learning. Retrieved from [https://www.researchgate.net/publication/327817574\\_Diabetes\\_Predicting\\_mHealth\\_Application\\_Using\\_Machine\\_Learning](https://www.researchgate.net/publication/327817574_Diabetes_Predicting_mHealth_Application_Using_Machine_Learning)
- [16] Qing, A. (2021, Jan 23). Nearly one-third of Singaporeans gained weight during COVID-19 pandemic: Survey. Retrieved from <https://www.straitstimes.com/singapore/nearly-one-third-of-singaporeans-gain-weight-during-covid-19-pandemic-survey>
- [17] Division of Nutrition, Physical Activity, and Obesity, National Center for Chronic Diseases Prevention and Health Promotion. (2021, Mar 22). Adult Obesity Causes & Consequences. Retrieved from <https://www.cdc.gov/obesity/adult/causes.html>
- [18] NHS Inform. (2021, May 17). Obesity. Retrieved from <https://www.nhsinform.scot/illnesses-and-conditions/nutritional/obesity>
- [19] Centers for Disease Control and Prevention. (2020, Jun 11). What is Diabetes? Retrieved from <https://www.cdc.gov/diabetes/basics/diabetes.html>
- [20] Mayo Clinic Staff. (2021, Jun 25). Diabetes prevention: 5 tips for taking control. Retrieved from <https://www.mayoclinic.org/diseases-conditions/type-2-diabetes/in-depth/diabetes-prevention/art-20047639>
- [21] National Center for Chronic Diseases Prevention and Health Promotion, Division for Heart Disease and Stroke Prevention. (2021, May 18). High Blood Pressure Symptoms and Causes. Retrieved from <https://www.cdc.gov/bloodpressure/about.htm>

C3879C CAPSTONE PROJECT

---

- [22] National Center for Chronic Diseases Prevention and Health Promotion, Division for Heart Disease and Stroke Prevention. (2020, Feb 24). Prevent High Blood Pressure. Retrieved from <https://www.cdc.gov/bloodpressure/prevent.htm>
- [23] National Center for Chronic Diseases Prevention and Health Promotion, Division for Heart Disease and Stroke Prevention. (2020, Feb 24). Prevent High Blood Pressure. Retrieved from <https://www.cdc.gov/bloodpressure/manage.htm>
- [24] Andreas C.Muller, Sarah Guido. Introduction to Machine Learning with Python Books (A Guide for Data Scientists).