



Faculty of Engineering and Technology

Department of Electrical and Computer Engineering

ENCS5141—Intelligent Systems Laboratory

**Assignment #1—Data Cleaning and Feature Engineering for the Bike Sharing
Dataset and Comparative Analysis of Classification Techniques**

Prepared by: Rahaf Naser—1201319

Instructor: Dr. Aziz Qaroush

Assistant teacher: Eng.Ahmad Abbas

Section#: 2

Date of Submission: November 30, 2024

BIRZEIT

November– 2024

Abstract

This study evaluates the impact of data preprocessing on model performance and compares the effectiveness of three classification techniques (Random Forest (RF), Support Vector Machine (SVM), and Multilayer Perceptron (MLP)) using the Bike Sharing dataset. The analysis involves comprehensive data cleaning and feature engineering, including handling missing values, encoding categorical variables, scaling numerical features, and applying dimensionality reduction, to assess its influence on model performance compared to raw data. RF, SVM, and MLP models are trained on the preprocessed data and evaluated based on accuracy, precision, recall, F1-score, computational time, and memory usage, with parameter tuning conducted for optimization. Results show that data preprocessing significantly improves model accuracy, consistency, and efficiency. While SVM demonstrates high accuracy and efficiency, RF handles complex patterns effectively but is computationally intensive, and MLP shows robust performance with sensitivity to parameter tuning. Overall, SVM strikes the best balance between predictive performance and computational efficiency for this dataset.

Table of Contents

Abstract.....	I
1.Introduction.....	1
1.1.Motivation.....	1
1.2.Background.....	1
1.3.Objective.....	1
2.Procedure and Discussion	2
2.1.Part1.....	2
2.1.1.Load the Bike Sharing Dataset	2
2.1.2.Data Exploration:.....	2
2.1.3.Handling missing values.....	2
2.1.4.Handling outliers.....	6
2.1.5.Data visualization.....	8
2.1.6.Correlation matrix.....	10
2.1.7.Variance Threshold.....	17
2.1.8.Information gain.....	18
2.1.9.Encode categorical features and normalize numerical features.....	18
2.1.10.Split the dataset into training and testing subsets.....	19
2.1.11.Train a Random Forest model on raw data.....	19
2.1.12.Dimensionality reduction PCA.....	20
2.1.13.Train a Random Forest model on the preprocessed data.....	21
2.1.14. Comparing the performance of the model trained on preprocessed data vs. raw data.....	22
2.2.Part2.....	23
2.2.1.Training SVM Model on preprocessed data.....	23
2.2.2.Parameter tuning of SVM model	24
2.2.3. Training MLP Model on preprocessed data.....	24
2.2.4. Parameter tuning of MLP model.....	25
2.2.5. Training Random Forest Model on preprocessed data.....	25
2.2.6. Parameter tuning of Random Forest model.....	26
2.2.7.Compare the performance of RF, SVM, and MLP on the same classification problem.....	27
3.Conclusion	28

Table of Figures

Figure 2.1.2.1: shape of dataset.....	2
Figure 2.1.2.2: Dataset information.....	2
Figure 2.1.2.3: Dataset description.....	3
Figure 2.1.2.4: Number of missing value.....	4
Figure 2.1.2.5: Names of columns in the dataset.....	5
Figure 2.1.2.6: Number of unique values in each column.....	6
Figure 2.1.3.1: Number of missing values after handling them.....	7
Figure 2.1.4.1: Boxplot of Waist Circumference before outliers removal.....	8
Figure 2.1.4.2: Boxplot of Pulmonary Function before outlier removal.....	8
Figure 2.1.4.3: Computing percentiles, interquartile range, lower and upper bound.....	9
Figure 2.1.4.4: Number of records before outlier removal.....	9
Figure 2.1.4.5: Boxplot of Pulmonary Function after outlier removal.....	9
Figure 2.1.4.6: Boxplot of Waist Circumference after outlier removal.....	10
Figure 2.1.5.1: Histogram of Distribution of Insulin Levels.....	10
Figure 2.1.5.2: Histogram of Distribution of Blood Glucose Levels.....	11
Figure 2.1.5.3: Histogram of Distribution of BMI.....	11
Figure 2.1.5.4: Histogram of Distribution of Age.....	12
Figure 2.1.5.5: Histogram of Distribution of Target.....	12
Figure 2.1.5.6: Histogram of Distribution of Physical Activity.....	13
Figure 2.1.5.7: Box plot of Blood Glucose Levels by Target.....	13
Figure 2.1.5.8: Scatter plot of Age vs Blood Pressure.....	14
Figure 2.1.5.9: Stacked histogram of BMI Distribution by Target.....	14
Figure 2.1.5.10: Scatter plot of Age vs Blood Pressure by Target.....	15
Figure 2.1.5.11: pie chart about Percentage Distribution of Target.....	15
Figure 2.1.5.12: Stacked Histogram of Blood Glucose Levels by Target.....	16
Figure 2.1.5.13: Stacked Histogram Blood Glucose Levels by Target and Smoking States.....	16
Figure 2.1.5.14: Scatter plot of Avg vs Blood Pressure.....	17
Figure 2.1.6.1: Correlation Matrix.....	17
Figure 2.1.7.1: Variance Threshold result.....	18
Figure 2.1.8.1: Information gain result.....	18
Figure 2.1.9.1: Encoding and Normalization result.....	19
Figure 2.1.11.1: Training Random Forest model on raw data Result.....	20

Figure 2.1.11.2: Top 10 Features by Importance.....	21
Figure 2.1.11.3: Confusion Matrix Random Forest on raw data.....	21
Figure 2.1.13.1: Train a Random Forest model on the preprocessed data result.....	22
<i>Figure 2.1.13.2: Confusion Matrix Random Forest after PCA</i>	<i>22</i>
<i>Figure 2.2.1.1: Result of Training SVM Model on preprocessed data.....</i>	<i>23</i>
<i>Figure 2.2.1.2: Confusion matrix of SVM Model.....</i>	<i>24</i>
Figure 2.2.2.1: Training Result of Best Parameter of SVM Model.....	24
<i>Figure 2.2.3.1: Result of Training MLP Model on preprocessed data.....</i>	<i>25</i>
<i>Figure 2.2.3.2: Confusion matrix of MLP Model after PCA</i>	<i>25</i>
Figure 2.2.4.1: Training Result of Best Parameter of MLP Model.....	25
<i>Figure 2.2.5.1: Result of Training RF Model on preprocessed data.....</i>	<i>26</i>
<i>Figure 2.2.5.2: Confusion Matrix of Random Forest after PCA</i>	<i>26</i>
Figure 2.2.6.1: Training Result of Best Parameter of Random Forest Model.....	27

Table of Tables

Table 2.1.11.1: Training Random Forest model on raw data Result.....	20
Table 2.1.13.1: Train a Random Forest model on the preprocessed data result.....	22
Table 2.1.14.1: Comparing the performance of the model trained on preprocessed data vs. raw data	23
Table 2.2.1.1: Result of Training SVM Model on preprocessed data.....	24
Table 2.2.3.1: Result of Training MLP Model on preprocessed data.....	24
Table 2.2.7.1: Comparing the performance of RF, SVM, and MLP.....	27

1.Introduction

1.1.Motivation

The motivation of this study is to gain insights into the comparative performance of three classification algorithms (Random Forest (RF), Support Vector Machine (SVM), and Multilayer Perceptron (MLP)) and to understand the impact of data preprocessing techniques such as handling missing values, encoding categorical variables, scaling, and dimensionality reduction on their performance. By doing so, the study aims to provide empirical evidence to guide the selection of models and preprocessing strategies based on the specific requirements of bike-sharing demand prediction tasks.

1.2.Background

Machine learning has become a cornerstone of predictive analytics, with various models like Random Forest (RF), Support Vector Machine (SVM), and Multilayer Perceptron (MLP) are three widely used machine learning models for classification tasks, each with unique strengths. Random Forest is an ensemble method that combines multiple decision trees to improve performance and reduce overfitting, making it robust and simple to use, especially with high-dimensional data. It excels at handling both large datasets and complex feature interactions, though it can be computationally expensive and harder to interpret. SVM, on the other hand, is effective for high-dimensional data and non-linear classification, using a hyperplane to separate classes. It performs well in situations where data is not linearly separable, and kernel functions allow it to work in higher-dimensional spaces. However, SVM can be slow with large datasets and sensitive to noise. MLP, a type of artificial neural network, is powerful for modeling complex, non-linear patterns, and is particularly effective when dealing with large datasets and deep learning tasks. However, MLP models require careful tuning of parameters such as the number of layers and neurons, and they are computationally expensive, with a risk of overfitting if not regularized properly. Each of these models has its advantages and weaknesses, making them suitable for different types of problems in predictive analytics.

Preprocessing techniques, including handling missing values, encoding categorical variables, and scaling, are essential to prepare raw datasets for effective model training. The Bike Sharing dataset provides an ideal platform to study the impact of these methods due to its complexity and real-world relevance.

1.3.Objective

The study aims to investigate the impact of data preprocessing on the performance of machine learning models in predicting bike-sharing usage. It also seeks to determine which classification model (Random Forest (RF), Support Vector Machine (SVM), or Multilayer Perceptron (MLP)) achieves the best results in terms of accuracy, precision, recall, and computational efficiency on the dataset. Additionally, the study aims to explore how parameter tuning affects the performance of these models.

2.Procedure and Discussion

2.1.Part1

2.1.1.Load the Bike Sharing Dataset

The "diabetes_dataset " csv file was loaded using the read_csv function from the pandas library.

2.1.2.Data Exploration:

I used the code df.shape to retrieve the dimensions of the dataset, which returns a tuple representing the number of rows and columns in the DataFrame.

(70000, 34)

Figure 2.1.2.1: shape of dataset

The result in Figure 2.1.2.1 shows that the dataset contains 70,000 rows, corresponding to the number of data entries or samples, and 34 columns, representing the number of features associated with each data entry. This gives an overview of the dataset's size, which is important for understanding the data structure and determining appropriate preprocessing steps.

I used the df.info() method to display a summary of the DataFrame, which includes key details about the dataset.

```
RangeIndex: 70000 entries, 0 to 69999
Data columns (total 34 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Target                                70000 non-null  object
1   Genetic Markers                       70000 non-null  object
2   Autoantibodies                       70000 non-null  object
3   Family History                       70000 non-null  object
4   Environmental Factors                 70000 non-null  object
5   Insulin Levels                       70000 non-null  int64
6   Age                                  70000 non-null  int64
7   BMI                                  70000 non-null  int64
8   Physical Activity                     70000 non-null  object
9   Dietary Habits                       70000 non-null  object
10  Blood Pressure                       70000 non-null  int64
11  Cholesterol Levels                   70000 non-null  int64
12  Waist Circumference                 70000 non-null  int64
13  Blood Glucose Levels                 70000 non-null  int64
14  Ethnicity                           70000 non-null  object
15  Socioeconomic Factors               70000 non-null  object
16  Smoking Status                      70000 non-null  object
17  Alcohol Consumption                 70000 non-null  object
18  Glucose Tolerance Test               70000 non-null  object
19  History of PCOS                     70000 non-null  object
20  Previous Gestational Diabetes        70000 non-null  object
21  Pregnancy History                   70000 non-null  object
22  Weight Gain During Pregnancy         70000 non-null  int64
23  Pancreatic Health                   70000 non-null  int64
24  Pulmonary Function                  70000 non-null  int64
25  Cystic Fibrosis Diagnosis            70000 non-null  object
26  Steroid Use History                 70000 non-null  object
27  Genetic Testing                     70000 non-null  object
28  Neurological Assessments            70000 non-null  int64
29  Liver Function Tests                70000 non-null  object
30  Digestive Enzyme Levels              70000 non-null  int64
31  Urine Test                          70000 non-null  object
32  Birth Weight                        70000 non-null  int64
33  Early Onset Symptoms                70000 non-null  object
```

Figure 2.1.2.2: Dataset information

Figure 2.1.2.2 shows that All columns have 70,000 non-null values, meaning there are no missing values in any of the columns. The dataset contains 13 columns with numerical data types (int64) and 21 columns with categorical data types (object). The dataset uses approximately 18.2 MB of memory.

I used the `df.describe()` method to generate summary statistics for the numerical columns in the dataset. This method provides key statistical measures that help understand the distribution and range of values in the dataset.

	Insulin Levels	Age	BMI	Blood Pressure	Cholesterol Levels	Circu	BMI	Blood Pressure	Cholesterol Levels	Waist Circumference	Blood Glucose Levels	Weight Gain During Pregnancy	Pancreatic Health	Pulmonary Function	Neurological Assessments	Digestive Enzyme Levels	Birth Weight
count	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000
mean	21.607443	32.020700	24.782943	111.339543	194.867200	34.782943	111.339543	194.867200	35.051657	160.701657	15.496414	47.564243	70.264671	1.804157	46.420529	3097.061071	
std	10.785852	21.043173	6.014236	19.945000	44.532466	6.014236	19.945000	44.532466	6.803461	48.165547	9.633096	19.984683	11.965600	0.680154	19.391089	713.837300	
min	5.000000	0.000000	12.000000	60.000000	100.000000	22.000000	60.000000	100.000000	20.000000	80.000000	0.000000	10.000000	30.000000	1.000000	10.000000	1500.000000	
25%	13.000000	14.000000	20.000000	99.000000	163.000000	30.000000	99.000000	163.000000	30.000000	121.000000	7.000000	32.000000	63.000000	1.000000	31.000000	2629.000000	
50%	19.000000	31.000000	25.000000	113.000000	191.000000	35.000000	113.000000	191.000000	34.000000	152.000000	18.000000	46.000000	72.000000	2.000000	48.000000	3103.000000	
75%	28.000000	49.000000	29.000000	125.000000	225.000000	39.000000	125.000000	225.000000	39.000000	194.000000	22.000000	64.000000	79.000000	2.000000	61.000000	3656.250000	
max	49.000000	79.000000	39.000000	149.000000	299.000000	59.000000	149.000000	299.000000	54.000000	299.000000	39.000000	99.000000	89.000000	3.000000	99.000000	4499.000000	

Figure 2.1.2.3: Dataset description

Figure 2.1.2.3 shows the total number of non-null entries for each column (70,000 for all numerical columns in this case), The average value for each numerical column, The measure of how spread out the values are for each column, The minimum value in each numerical column, The value below which 25% of the data falls, The middle value, also known as the 50th percentile, The value below which 75% of the data falls, The maximum value in each numerical column.

By running `df.isnull().sum()`, I am able to get the total count of missing (null) values for each column in the dataset.

	0
Target	0
Genetic Markers	0
Autoantibodies	0
Family History	0
Environmental Factors	0
Insulin Levels	0
Age	0
BMI	0
Physical Activity	0
Dietary Habits	0
Blood Pressure	0
Cholesterol Levels	0
Waist Circumference	0
Blood Glucose Levels	0
Ethnicity	0
Socioeconomic Factors	0
Smoking Status	0
Alcohol Consumption	0
Glucose Tolerance Test	0
History of PCOS	0
Previous Gestational Diabetes	0
Pregnancy History	0
Weight Gain During Pregnancy	0
Pancreatic Health	0
Pulmonary Function	0
Cystic Fibrosis Diagnosis	0
Steroid Use History	0
Genetic Testing	0
Neurological Assessments	0
Liver Function Tests	0
Digestive Enzyme Levels	0
Urine Test	0
Birth Weight	0
Early Onset Symptoms	0

dtype: int64

Figure 2.1.2.4: Number of missing value

Figure 2.1.2.4 shows a count of zero for all columns, it indicates that there are no missing values in the dataset.

```
Numeric Columns:
Index(['Insulin Levels', 'Age', 'BMI', 'Blood Pressure', 'Cholesterol Levels',
      'Waist Circumference', 'Blood Glucose Levels',
      'Weight Gain During Pregnancy', 'Pancreatic Health',
      'Pulmonary Function', 'Neurological Assessments',
      'Digestive Enzyme Levels', 'Birth Weight'],
      dtype='object')

Categorical Columns:
Index(['Target', 'Genetic Markers', 'Autoantibodies', 'Family History',
      'Environmental Factors', 'Physical Activity', 'Dietary Habits',
      'Ethnicity', 'Socioeconomic Factors', 'Smoking Status',
      'Alcohol Consumption', 'Glucose Tolerance Test', 'History of PCOS',
      'Previous Gestational Diabetes', 'Pregnancy History',
      'Cystic Fibrosis Diagnosis', 'Steroid Use History', 'Genetic Testing',
      'Liver Function Tests', 'Urine Test', 'Early Onset Symptoms'],
      dtype='object')
```

Figure 2.1.2.5: Names of columns in the dataset

Figure 2.1.2.5 shows the Numeric Columns and Categorical Columns in the dataset.

	0
Target	13
Genetic Markers	2
Autoantibodies	2
Family History	2
Environmental Factors	2
Insulin Levels	45
Age	80
BMI	28
Physical Activity	3
Dietary Habits	2
Blood Pressure	90
Cholesterol Levels	200
Waist Circumference	35
Blood Glucose Levels	220
Ethnicity	2
Socioeconomic Factors	3
Smoking Status	2
Alcohol Consumption	3
Glucose Tolerance Test	2
History of PCOS	2
Previous Gestational Diabetes	2
Pregnancy History	2
Weight Gain During Pregnancy	40
Pancreatic Health	90
Pulmonary Function	60
Cystic Fibrosis Diagnosis	2
Steroid Use History	2
Genetic Testing	2
Neurological Assessments	3
Liver Function Tests	2
Digestive Enzyme Levels	90
Urine Test	4
Birth Weight	3000
Early Onset Symptoms	2

Figure 2.1.2.6: Number of unique values in each column

Figure 2.1.2.6 shows how many unique values exist in each column of the dataset.

2.1.3. Handling missing values

I used two methods in the code to handle missing values in both numeric and non-numeric columns of the dataset. I fill any missing (NaN) values in each numeric column with the median value of that column. The median is often preferred for numeric imputation because it is less sensitive to outliers than the mean.

I fill any missing (NaN) values in each categorical column with the mode (most frequent value) of that column. The mode is used for categorical imputation because it represents the most common value in the column, which is a reasonable assumption for missing data in categorical features.

	0
Target	0
Genetic Markers	0
Autoantibodies	0
Family History	0
Environmental Factors	0
Insulin Levels	0
Age	0
BMI	0
Physical Activity	0
Dietary Habits	0
Blood Pressure	0
Cholesterol Levels	0
Waist Circumference	0
Blood Glucose Levels	0
Ethnicity	0
Socioeconomic Factors	0
Smoking Status	0
Alcohol Consumption	0
Glucose Tolerance Test	0
History of PCOS	0
Previous Gestational Diabetes	0
Pregnancy History	0
Weight Gain During Pregnancy	0
Pancreatic Health	0
Pulmonary Function	0
Cystic Fibrosis Diagnosis	0
Steroid Use History	0
Genetic Testing	0
Neurological Assessments	0
Liver Function Tests	0
Digestive Enzyme Levels	0
Urine Test	0
Birth Weight	0
Early Onset Symptoms	0

dtype: int64

Figure 2.1.3.1: Number of missing values after handling them

Figure 2.1.3.1 shows that dataset cleaned from missing values.

2.1.4.Handling outliers

I wrote code to detect outliers in the specified columns of the dataset by visualizing them using boxplots.

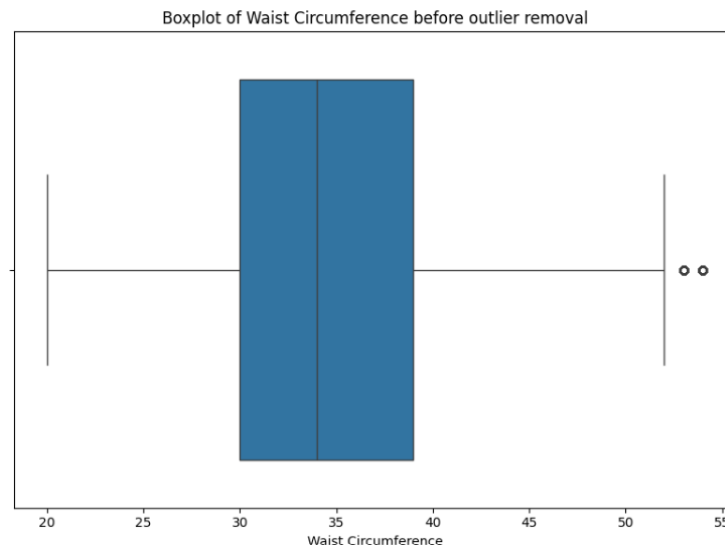


Figure 2.1.4.1: Boxplot of Waist Circumference before outliers removal

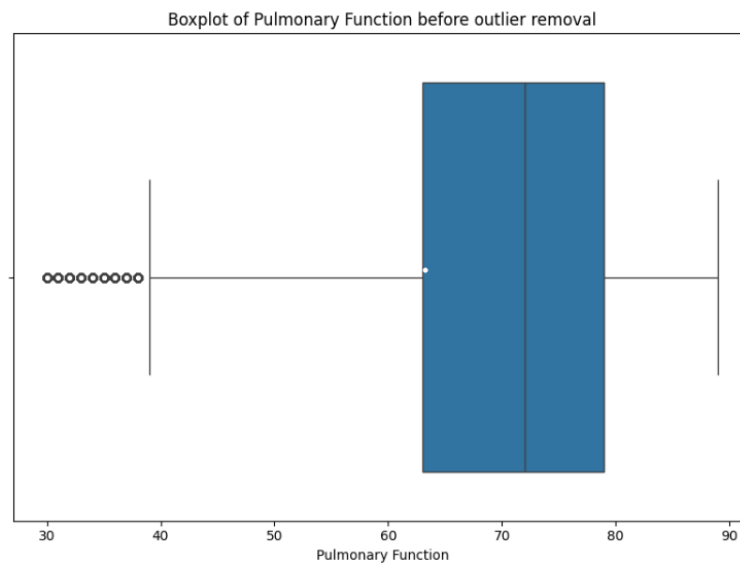


Figure 2.1.4.2: Boxplot of Pulmonary Function before outlier removal

Figures above detect outliers in the 'Waist Circumference' 'Pulmonary Function' columns.

```

25th Percentile: 30.0
50th Percentile: 34.0
75th Percentile: 39.0
Interquartile Range (IQR): 9.0
Lower Bound for Waist Circumference = 16.5, and Upper Bound for fare = 52.5
25th Percentile for Pulmonary Function: 63.0
50th Percentile for Pulmonary Function: 72.0
75th Percentile for Pulmonary Function: 79.0
Interquartile Range (IQR) for Pulmonary Function: 16.0
Lower Bound for Pulmonary Function = 39.0, and Upper Bound for Pulmonary Function = 103.0

```

Figure 2.1.4.3: Computing percentiles, interquartile range, lower and upper bound

Figure 2.1.4.3 shows the outlier detection results for **Waist Circumference** and **Pulmonary Function** were calculated using the Interquartile Range (IQR) method. For **Waist Circumference**, the IQR is 9.0, with lower and upper bounds of 16.5 and 52.5, respectively, meaning values outside this range are potential outliers. For **Pulmonary Function**, the IQR is 16.0, with lower and upper bounds of 39.0 and 103.0, respectively. Any data points falling below or above these bounds are considered outliers for each respective feature.

```

Number of outliers based on the Interquartile Range for 'Pulmonary_Function' column: 1206 (1.72%)
Number of outliers based on the Interquartile Range for 'Waist_Circumference' column: 522 (0.75%)

```

Figure 2.1.4.4: Number of records before outlier removal

The results in Figure 2.1.4.4 show that there are **1,206 outliers** (1.72%) in the **Pulmonary Function** column and **522 outliers** (0.75%) in the **Waist Circumference** column based on the Interquartile Range (IQR) method. This indicates that a small percentage of the data points in both columns fall outside the calculated bounds, suggesting the presence of extreme values that may need to be addressed in further analysis.

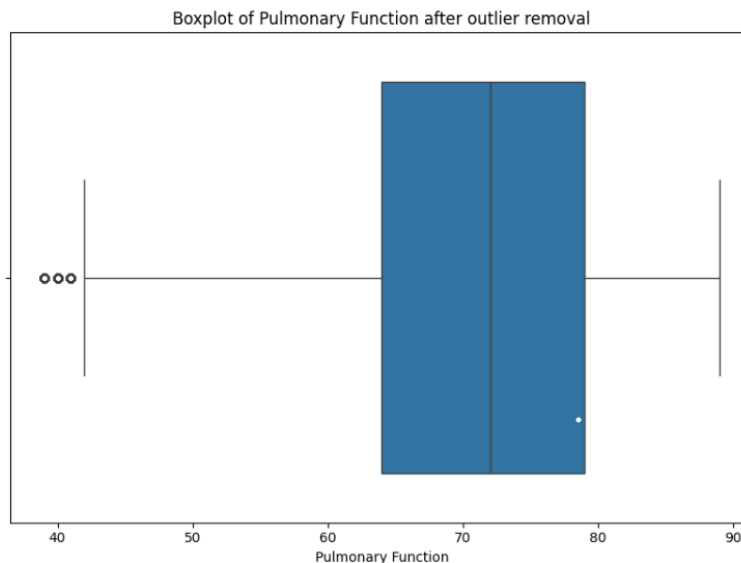


Figure 2.1.4.5: Boxplot of Pulmonary Function after outlier removal

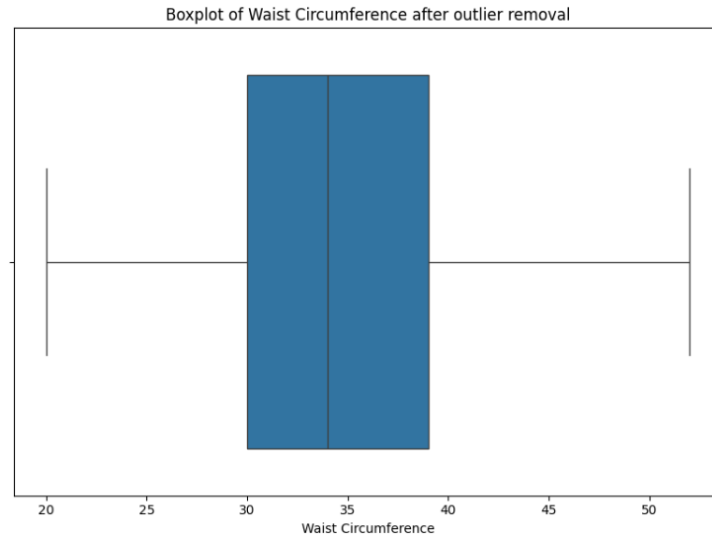


Figure 2.1.4.6: Boxplot of Waist Circumference after outlier removal

The above figures show a Boxplot after outlier removal in the two columns.

2.1.5.Data visualization

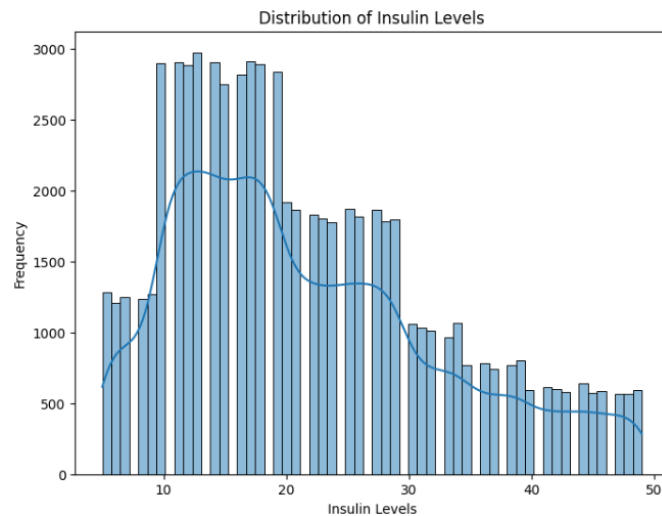


Figure 2.1.5.1: Histogram of Distribution of Insulin Levels

The above histogram and its fitted distribution curve show a unimodal distribution of insulin levels, peaking around 15-20 units. The frequency declines steadily as levels increase beyond the peak, indicating that most individuals fall within a specific range. This suggests that insulin levels tend to cluster around a central value and gradually become less common towards the extremes.

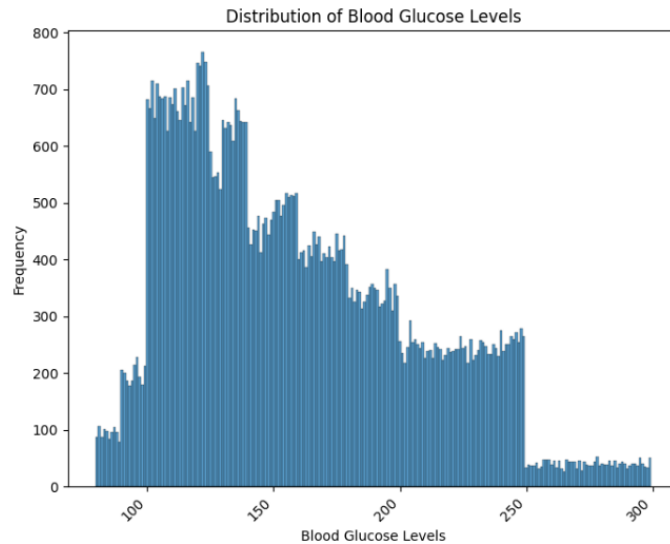


Figure 2.1.5.2: Histogram of Distribution of Blood Glucose Levels

Figure 2.1.5.2 shows the distribution of blood glucose levels. The x-axis represents the blood glucose levels, and the y-axis represents the frequency. The histogram shows that the distribution of blood glucose levels is unimodal, with a peak around 150 mg/dL. The distribution is skewed to the right, indicating that there are more individuals with lower blood glucose levels than those with higher levels.

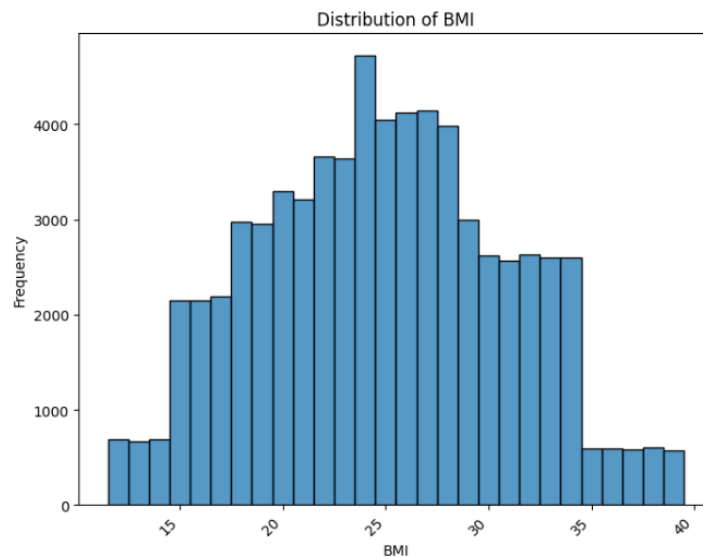


Figure 2.1.5.3: Histogram of Distribution of BMI

Figure 2.1.5.3 shows histogram with a unimodal distribution, with a peak around a BMI of 25. This indicates that the majority of individuals have a BMI in this range. The distribution appears to be slightly skewed to the right, suggesting that there are more individuals with higher BMI values than lower BMI values.

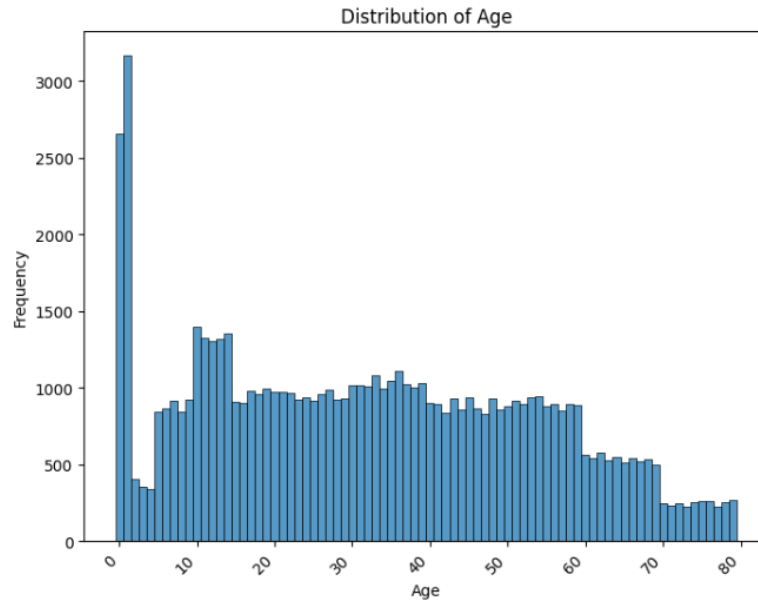


Figure 2.1.5.4: Histogram of Distribution of Age

The histogram in Figure 2.1.5.4 shows a right-skewed distribution, indicating that there are more younger individuals in the dataset than older individuals. The peak of the distribution occurs at around 2-3 years old, suggesting a high frequency of individuals in that age group.

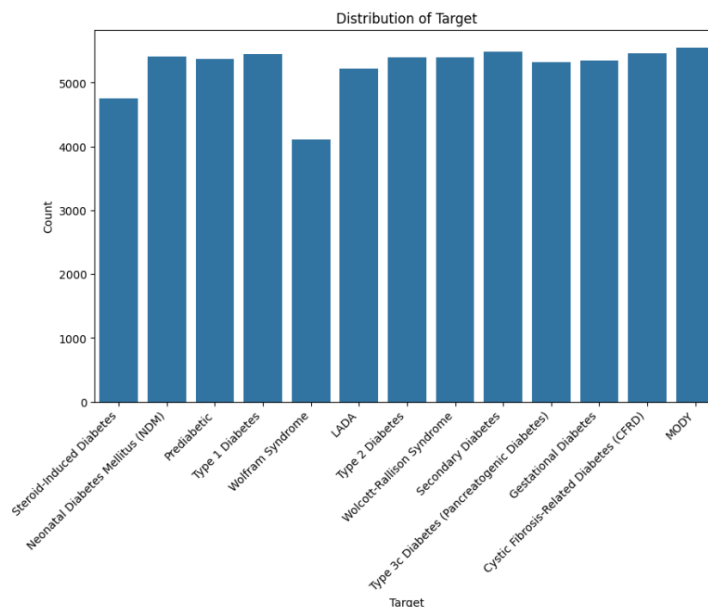


Figure 2.1.5.5: Histogram of Distribution of Target

Figure 2.1.5.5 shows the distribution of different types of diabetes in a dataset. The most common type of diabetes is **Prediabetic**, followed by **Type 1 Diabetes**, **Steroid-Induced Diabetes**, **Neonatal Diabetes Mellitus (NDM)**, and **Type 2 Diabetes**. The rarest type of diabetes is **MODY**.

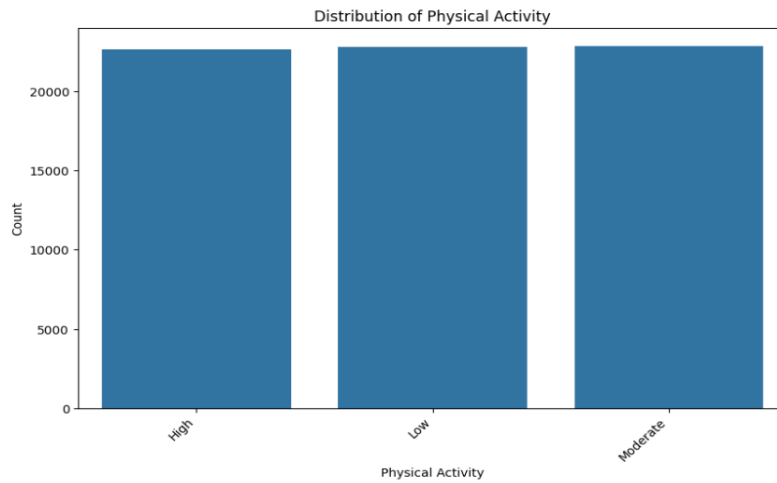


Figure 2.1.5.6: Histogram of Distribution of Physical Activity

This histogram in Figure 2.1.5.8 shows the distribution of physical activity levels in a dataset. The x-axis represents the levels of physical activity, which are categorized as High, Low, and Moderate. The y-axis shows the count of individuals within each category. The graph reveals that there are similar numbers of individuals in each physical activity category.

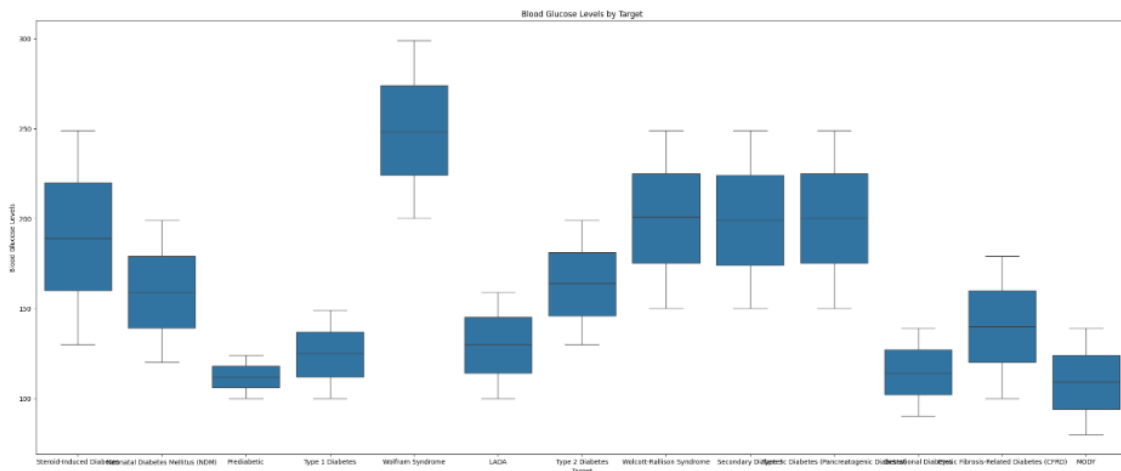


Figure 2.1.5.7: Box plot of Blood Glucose Levels by Target

The boxplot in Figure 2.1.5.9 shows the distribution of blood glucose levels for individuals with different types of diabetes. The boxplots show that individuals with Wolfram Syndrome, LADA, Type 1 Diabetes, and Type 2 Diabetes have the highest blood glucose levels. The boxplots show that individuals with Secondary Diabetic (Pancreatogenic) and Fibrosis-Related Diabetes (CFRD) have moderate blood glucose levels. The boxplots show that individuals with MODY, Neonatal Diabetes Mellitus (MDM), and Prediabetic have the lowest blood glucose levels. There are some outliers in the data, particularly for Wolfram Syndrome, LADA, Type 1 Diabetes, and Type 2 Diabetes. This indicates that some individuals with these conditions have much higher blood glucose levels than others.

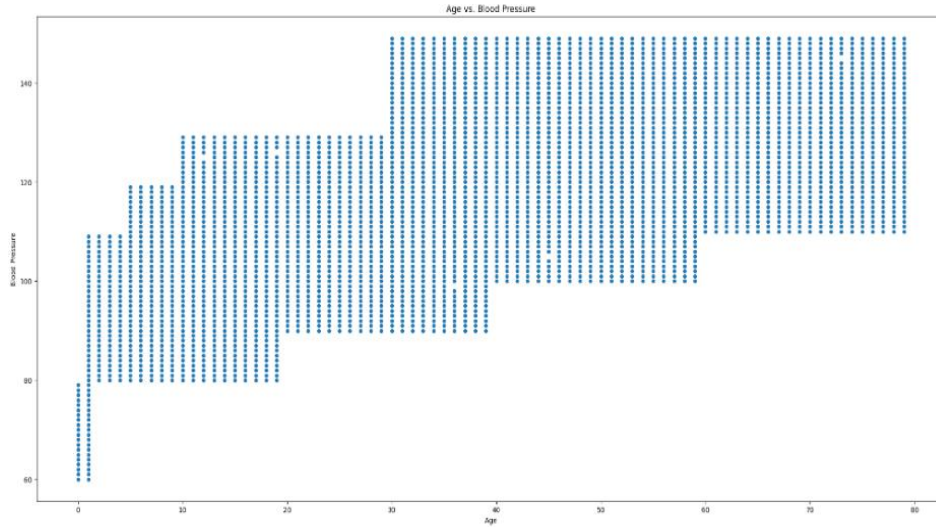


Figure 2.1.5.8: Scatter plot of Age vs Blood Pressure

The Scatter plot in Figure 2.1.5.10 shows clusters of points for each age group, implying that several individuals share similar blood pressure readings for a specific age.

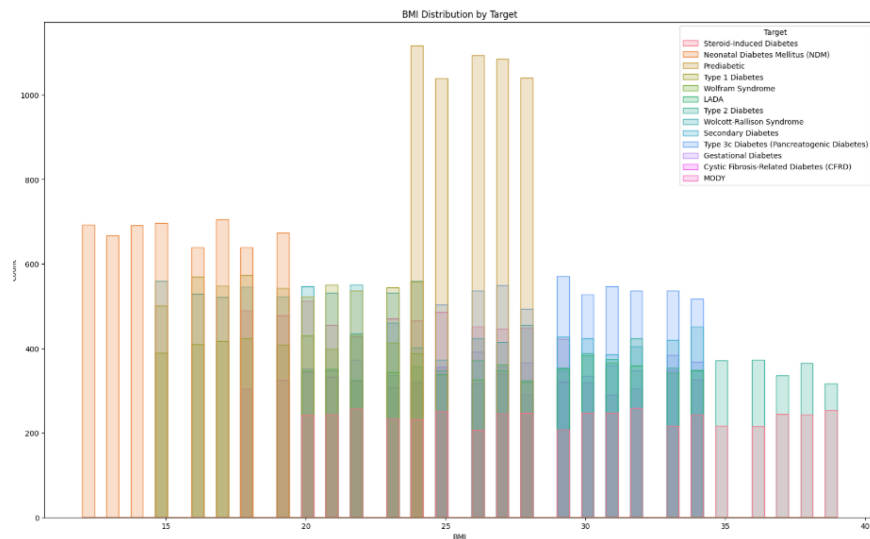


Figure 2.1.5.9: Stacked histogram of BMI Distribution by Target

Figure 2.1.5.11 shows the distribution of Body Mass Index (BMI) across various types of diabetes. It demonstrates that **Type 1 Diabetes** and **Type 2 Diabetes** are the most prevalent types, with the highest number of individuals. The figure also highlights that **Prediabetic** individuals are present in a substantial number, while **Gestational Diabetes** and **MODY** are the least prevalent. The distribution suggests that different types of diabetes may have different BMI profiles.

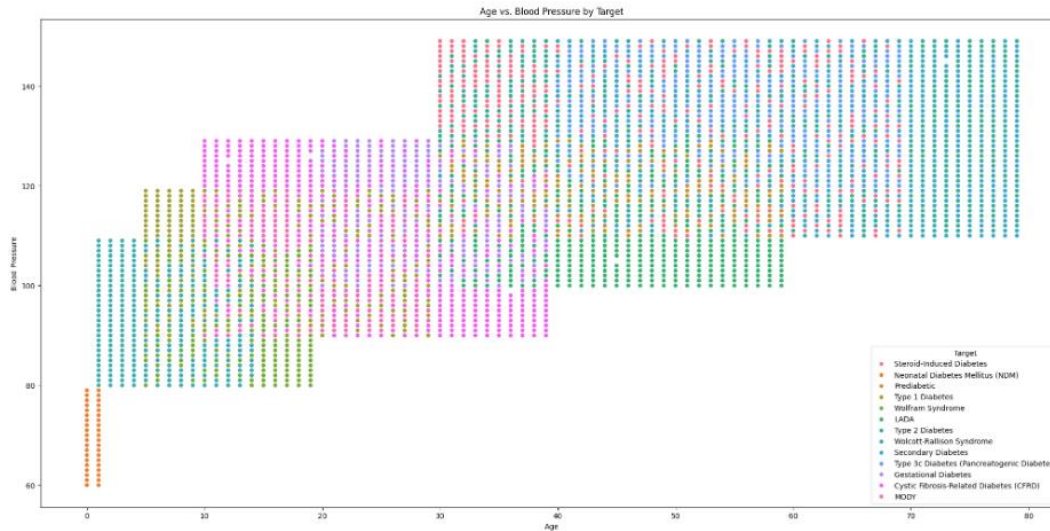


Figure 2.1.5.10: Scatter plot of Age vs Blood Pressure by Target

The Scatter plot in Figure 2.1.5.12 shows that blood pressure generally increases with age. This trend is similar across different types of diabetes, although the specific relationship varies. For example, MODY patients tend to have a higher blood pressure than prediabetic patients.

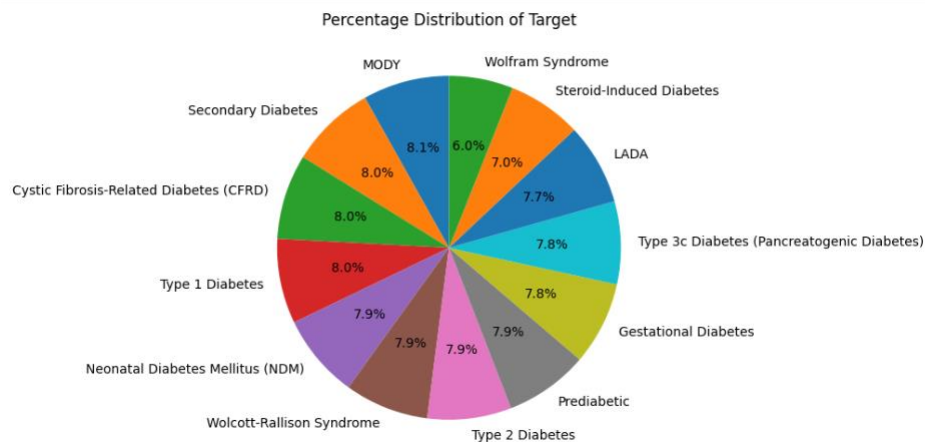


Figure 2.1.5.11: pie chart about Percentage Distribution of Target

The pie chart in Figure 2.1.5.13 shows the percentage distribution of different types of diabetes and related conditions. The most common type is Type 2 Diabetes (7.9%), followed by Type 3c Diabetes (Pancreatogenic Diabetes) (7.8%), and Gestational Diabetes (7.8%). Other conditions, such as MODY (8.1%) and Cystic Fibrosis-Related Diabetes (CFRD) (8.0%), are also represented.

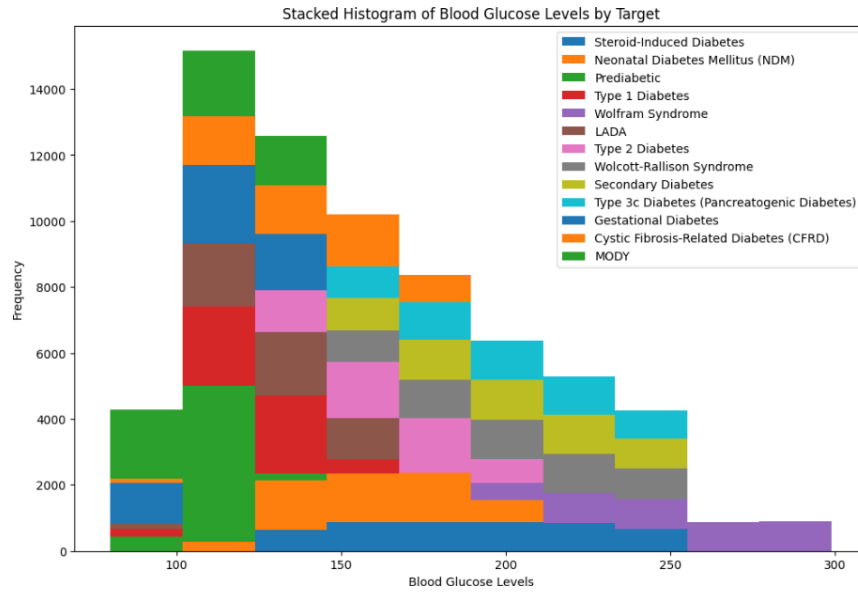


Figure 2.1.5.12: Stacked Histogram of Blood Glucose Levels by Target

Figure 2.1.5.14 shows a stacked histogram of blood glucose levels by target. The different colors represent different types of diabetes. The histogram shows that the distribution of blood glucose levels varies depending on the type of diabetes. For example, people with type 1 diabetes have a higher frequency of blood glucose levels above 200 mg/dL compared to people with type 2 diabetes.

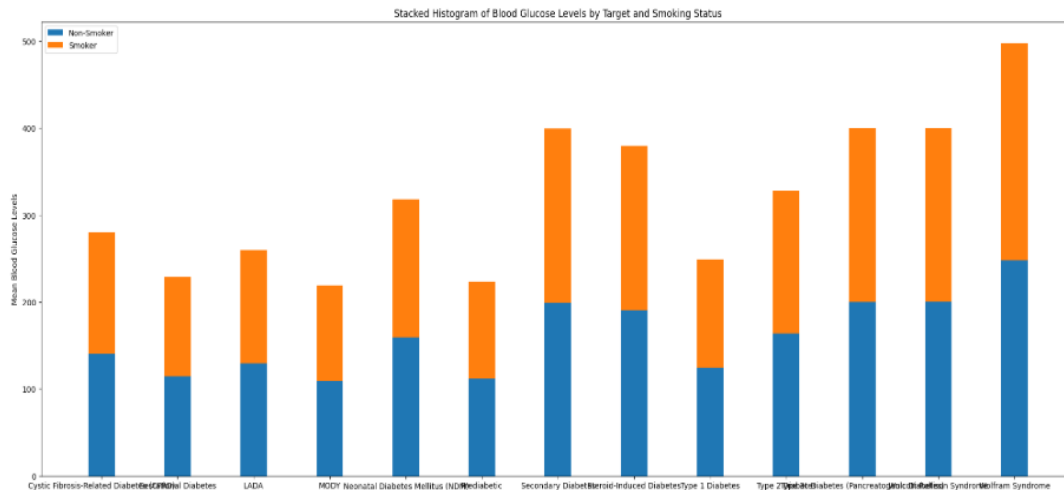


Figure 2.1.5.13: Stacked Histogram Blood Glucose Levels by Target and Smoking States

Figure 2.1.5.15 shows the mean blood glucose levels in different types of diabetes, categorized by smoking status (non-smoker vs. smoker). **Smokers generally have higher blood glucose levels across all diabetes types compared to non-smokers. Wolfram syndrome and Type 2 diabetes show the highest blood glucose levels among both smokers and non-smokers. LADA (latent autoimmune diabetes in adults) appears to have lower blood glucose levels compared to other types.**

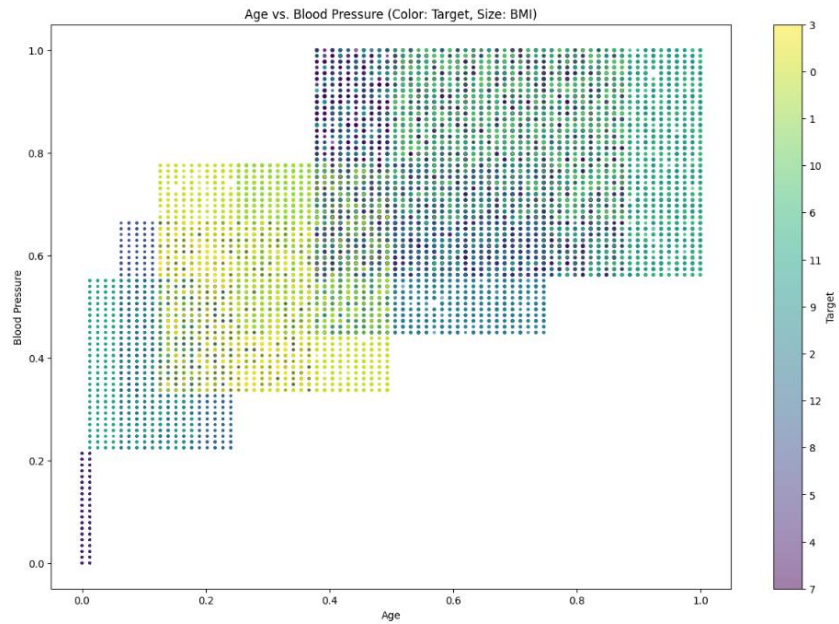


Figure 2.1.5.14: Scatter plot of Avg vs Blood Pressure

The figure displays the relationship between age and blood pressure, with the color representing the target variable and the size representing BMI. The data points are clustered in distinct regions, indicating a potential correlation between the variables. The clustering suggests that younger individuals with lower blood pressure have a higher target value (yellow). Conversely, older individuals with higher blood pressure have a lower target value (purple).

2.1.6. Correlation matrix

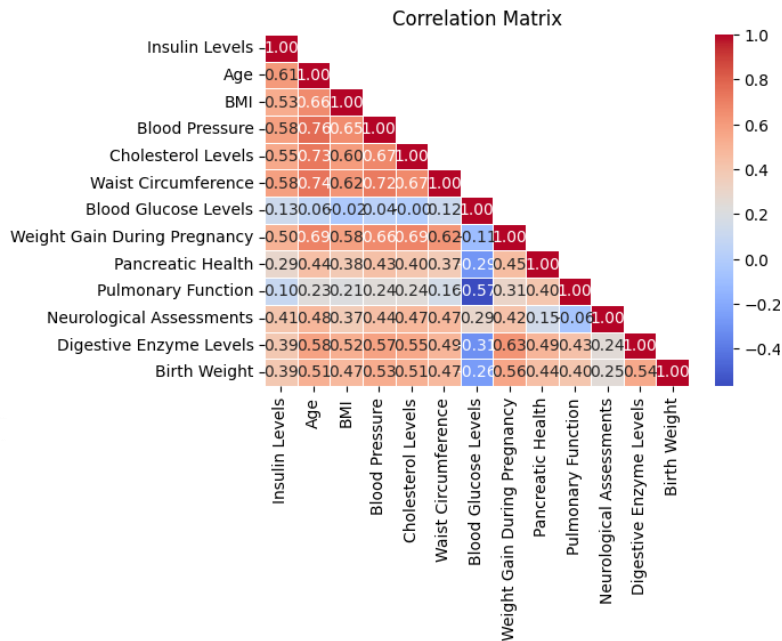


Figure 2.1.6.1: Correlation Matrix

the correlation matrix in Figure 2.1.6.1 provides insights into the relationships between various factors, highlighting the strong associations and areas where further investigation may be warranted. Age and BMI have a strong positive correlation (0.66). This suggests that as Age increases, BMI also tends to increase. Similarly, Age and Blood Pressure show a strong positive correlation (0.76). Age and Insulin Levels have a strong negative correlation (-0.61). This indicates that as Age increases, Insulin Levels tend to decrease. Blood Glucose Levels and Birth Weight have a weak correlation (0.11). This suggests there is little relationship between these two factors.

2.1.7.Variance Threshold

The variance of a feature is a measure of its spread or dispersion. A low variance indicates that the data points are clustered closely together, while a high variance indicates that the data points are spread out over a wider range.

```
Array of variances:
[4.40000000e+01 7.90000000e+01 2.70000000e+01 8.90000000e+01
 1.99000000e+02 3.20000000e+01 2.19000000e+02 3.90000000e+01
 8.90000000e+01 5.00000000e+01 4.62522695e-01 8.90000000e+01
 2.99900000e+03]
Lowest variance: 0.4625226946797927
Column number: 10
Column name: Blood Pressure
```

Figure 2.1.7.1: Variance Threshold result

Figure 2.1.7.1 shows that the feature "Blood Pressure" has the lowest variance among all features. This means that the data points for blood pressure are clustered very closely together, suggesting that blood pressure values are relatively similar across the dataset.

2.1.8.Information gain

Information Gain is a method used in machine learning for feature selection. It measures the amount of information a particular feature provides about the target variable. Higher information gain suggests that the feature is more relevant for predicting the target variable.

```
Feature selection using Information Gain:
Feature Information Gain
12 Birth Weight 7.884586
11 Digestive Enzyme Levels 0.393194
7 Weight Gain During Pregnancy 0.366187
1 Age 0.307593
4 Cholesterol Levels 0.282703
3 Blood Pressure 0.280317
8 Pancreatic Health 0.241361
5 Waist Circumference 0.238151
2 BMI 0.204372
9 Pulmonary Function 0.200439
0 Insulin Levels 0.154767
6 Blood Glucose Levels 0.143196
10 Neurological Assessments 0.051726
```

Figure 2.1.8.1: Information gain result

Figure 2.1.8.1 shows the information gain for each feature. The higher the information gain, the more useful the feature is in predicting the target variable.

Birth Weight has the highest information gain (7.884586), indicating it is the most important feature for predicting birth weight. This is expected as birth weight is the target variable, so it strongly influences itself. Other features with significant information gain include: **Digestive Enzyme Levels** (0.393194), **Weight Gain During Pregnancy** (0.366187), **Age** (0.307593)

Features with low information gain, such as **Neurological Assessments** (0.051726), are less useful for predicting birth weight.

2.1.9. Encode categorical features and normalize numerical features

I make MinMaxScaler (normalization) For numerical columns and label encoding for categorical columns the result was shown below in Figure 2.1.9.1

Transformed DataFrame:

	Target	Genetic Markers	Autoantibodies	Family History	Environmental Factors	Insulin Levels	Age	BMI	Physical Activity	Dietary Habits	...	Pulmonary Function	Cystic Fibrosis Diagnosis	Steroid Use History	Genetic Testing	Neurological Assessments	Liver Function Tests	Digestive Enzyme Levels	Urine Test	Birth Weight	Early Onset Symptoms
0	7	1	0	0	1	0.795455	0.558982	0.982983	0	0	...	0.779881	0	0	1	1.0	1	0.518854	1	0.378459	0
1	4	1	0	0	1	0.181818	0.012658	0.185185	0	0	...	0.508475	1	0	0	0.0	1	0.202247	0	0.127042	1
2	5	1	1	1	1	0.500000	0.455898	0.444444	0	1	...	0.847458	1	0	0	0.0	0	0.505818	1	0.707589	1
3	8	0	1	0	1	0.088182	0.088808	0.148148	1	1	...	1.000000	1	0	1	0.5	0	0.581798	1	0.880894	0
4	12	0	0	1	1	0.272727	0.126582	0.185185	0	0	...	0.188441	0	0	1	0.0	1	0.157303	3	0.090030	0

5 rows x 34 columns

Figure 2.1.9.1: Encoding and Normalization result

2.1.10. Split the dataset into training and testing subsets

I first import the necessary libraries for data handling, model training, and evaluation. Then, I define the target column ('Target') which the model will predict. Next, I separate the dataset into features (X) and the target variable (y). Finally, I split the data into training and testing sets using `train_test_split`, with 80% of the data for training and 20% for testing. This prepares the data for training a machine learning model, ensuring it learns on one set and is evaluated on another to check its performance.

2.1.11. Train a Random Forest model on raw data

```

Training Time (in seconds): 15.175531387329102

Accuracy of the Random Forest model: 0.8996704503844744

Classification Report:
              precision    recall  f1-score   support

     0           1.00       0.89       0.94       1057
     1           0.89       0.92       0.91       1103
     2           0.95       0.94       0.94       1025
     3           0.96       0.86       0.91       1073
     4           1.00       1.00       1.00       1057
     5           0.97       1.00       0.98       1079
     6           0.79       0.78       0.79       1101
     7           0.81       0.78       0.80       1013
     8           0.87       1.00       0.93       1110
     9           0.90       0.72       0.80       1093
    10           0.80       1.00       0.89       1122
    11           0.96       0.86       0.91       1053
    12           0.83       0.95       0.89       769

 accuracy          0.90          0.90          0.90       13655
 macro avg         0.90          0.90          0.90       13655
 weighted avg      0.90          0.90          0.90       13655

```

Figure 2.1.11.1: Training Random Forest model on raw data Result

Figure 2.1.11.1 shows that the Random Forest model shows good performance with an accuracy of 0.89, meaning that it correctly classified 89% of the data points. The model also has high precision and recall scores across different classes, indicating that it is able to accurately identify relevant data points and minimize false positives and false negatives.

Table 2.1.11.1: Training Random Forest model on raw data Result

Model	Accuracy	recall	F1-score	precision	Training time(s)
Random Forest model on raw data	89%	89%	94%	100%	15.17

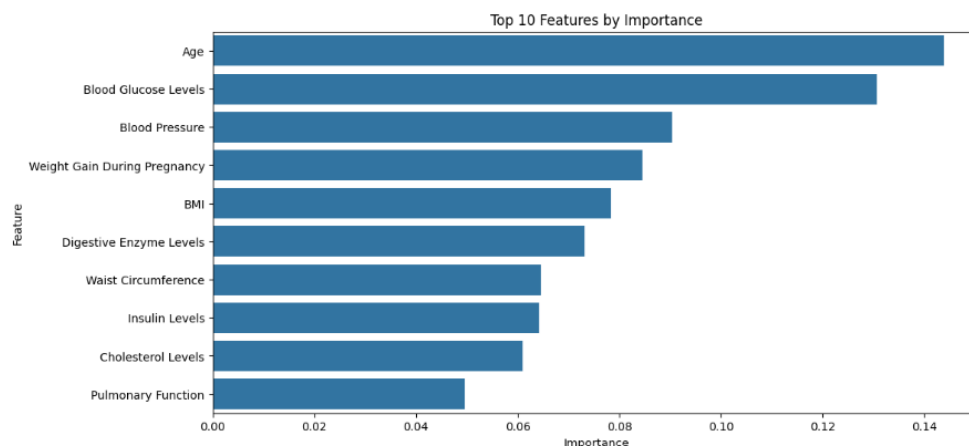


Figure 2.1.11.2: Top 10 Features by Importance

		Confusion Matrix Random Forest												
True		Predicted												
		Cystic Fibrosis-Related Diabetes (CFRD)	Gestational Diabetes	LADA	MODY	Neonatal Diabetes Mellitus (NDM)	Prediabetic	Secondary Diabetes	Steroid-induced Diabetes	Type 1 Diabetes	Type 2 Diabetes	Type 3c Diabetes (Pancreatogenic Diabetes)	Wolcott-Rallison Syndrome	Wolfram Syndrome
Cystic Fibrosis-Related Diabetes (CFRD)	-	946	72	17	13	0	0	0	0	9	0	0	0	0
Gestational Diabetes	-	3	1015	35	24	0	22	0	0	4	0	0	0	0
LADA	-	0	52	959	0	0	12	0	1	0	0	1	0	0
MODY	-	0	0	0	925	0	0	0	0	148	0	0	0	0
Neonatal Diabetes Mellitus (NDM)	-	0	0	0	0	1057	0	0	0	0	0	0	0	0
Prediabetic	-	0	0	0	0	0	1079	0	0	0	0	0	0	0
Secondary Diabetes	-	0	0	0	0	0	0	863	78	0	50	110	0	0
Steroid-induced Diabetes	-	0	0	2	0	0	0	45	789	0	40	137	0	0
Type 1 Diabetes	-	0	0	0	4	0	0	0	0	1106	0	0	0	0
Type 2 Diabetes	-	1	0	0	0	0	0	180	102	0	784	26	0	0
Type 3c Diabetes (Pancreatogenic Diabetes)	-	0	0	0	0	0	0	0	0	0	0	1122	0	0
Wolcott-Rallison Syndrome	-	0	0	0	0	0	0	0	0	0	0	0	906	147
Wolfram Syndrome	-	0	0	0	0	0	0	0	0	0	0	0	35	734

Figure 2.1.11.3: Confusion Matrix Random Forest on raw data

2.1.12.Dimensionality reduction PCA

PCA is a technique used for **dimensionality reduction**. It transforms the original features of a dataset into new, uncorrelated variables called **principal components**. The goal is to reduce the number of dimensions (features) while retaining as much of the original variance (information) as possible. In my code, PCA is used to reduce the feature set to a smaller number of components (with 70% of variance retained) to simplify the dataset for further analysis or model training, which can help improve performance and reduce computational complexity.

2.1.13.Train a Random Forest model on the preprocessed data

Training Time (in seconds): 13.040186882019043

Accuracy of the Random Forest model after PCA: 0.9989747345294764

Classification Report:				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	1057
1	1.00	1.00	1.00	1103
2	1.00	1.00	1.00	1025
3	1.00	1.00	1.00	1073
4	1.00	1.00	1.00	1057
5	1.00	1.00	1.00	1079
6	1.00	1.00	1.00	1101
7	1.00	1.00	1.00	1013
8	0.99	1.00	1.00	1110
9	1.00	1.00	1.00	1093
10	1.00	1.00	1.00	1122
11	1.00	1.00	1.00	1053
12	1.00	1.00	1.00	769
accuracy			1.00	13655
macro avg	1.00	1.00	1.00	13655
weighted avg	1.00	1.00	1.00	13655

Figure 2.1.13.1: Train a Random Forest model on the preprocessed data result

Figure 2.1.13.1 shows that **the Random Forest model performs extremely well on the preprocessed data with a high accuracy and perfect scores on most classification metrics.** The model's effectiveness suggests that the preprocessing steps have successfully improved the data quality, leading to high accuracy and a well-performing model.

Table 2.1.13.1: Train a Random Forest model on the preprocessed data result

Model	Accuracy	recall	F1-score	precision	Training time(s)
Random Forest model on preprocessed data	99%	100%	100%	100%	13.04

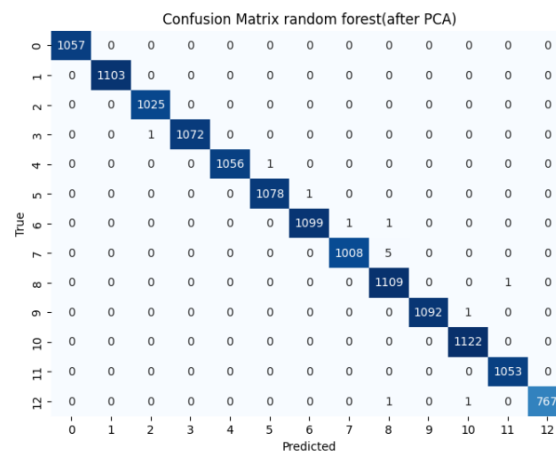


Figure 2.1.13.2: Confusion Matrix Random Forest after PCA

2.1.14. Comparing the performance of the model trained on preprocessed data vs. raw data

Random Forest on raw data: Represents the model trained using the original, unprocessed data.

Random Forest on preprocessed data: Represents the model trained using data that has been preprocessed for better performance.

Table 2.1.14.1 below shows that preprocessing the data before training the random forest model resulted in significantly better performance, with higher accuracy, recall, precision, and F1-score, as well as a slightly faster training time (less training time). This suggests that preprocessing can be a valuable step in improving model performance.

Table 2.1.14.1: Comparing the performance of the model trained on preprocessed data vs. raw data

Model	Accuracy	Recall	precision	F1-score	Training time (in seconds)
Random Forest on raw data	89%	89%	100%	94%	15.17
Random Forest on preprocessed data	99%	100%	100%	100%	13.04

2.2.Part2

2.2.1.Training SVM Model on preprocessed data

```
Support Vector Machine (SVM):
Training Time: 2.0099 seconds
Peak Memory Usage: 2.3114 MB
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1-Score: 1.0000
```

Figure 2.2.1.1: Result of Training SVM Model on preprocessed data

Figure 2.2.1.1 shows that the Support Vector Machine (SVM) model exhibits excellent performance based on the provided metrics:

Training Time: 2.0099 seconds

Indicates a relatively fast training process.

Peak Memory Usage: 2.3114 MB

Suggests a modest memory footprint during training.

Accuracy: 1.0000

Perfect accuracy, indicating the model correctly classified all instances.

Precision: 1.0000

Perfect precision, implying that all predicted positive instances were truly positive.

Recall: 1.0000

Perfect recall, meaning the model identified all actual positive instances.

F1-Score: 1.0000

The F1-Score, being the harmonic mean of precision and recall, also reaches 1.0000, confirming the model's exceptional performance in both precision and recall.

Table 2.2.1.1: Result of Training SVM Model on preprocessed data

Model	Accuracy	Recall	Precision	F1-score	Training time (in seconds)	Peak Memory usage (in MB)
SVM	100%	100%	100%	100%	2.0099	2.3114

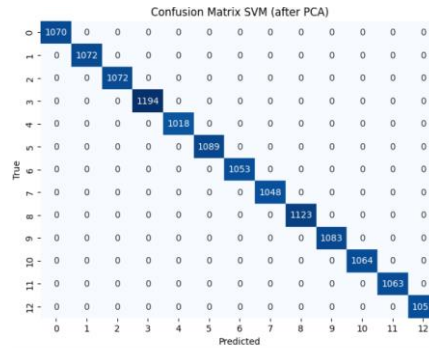


Figure 2.2.1.2: Confusion matrix of SVM Model

This confusion matrix suggests the SVM model (after PCA) performs exceptionally well, achieving high accuracy and classifying instances correctly across all classes.

2.2.2. Parameter tuning of SVM model

Fitting 2 folds for each of 10 candidates, totalling 20 fits

```
Support Vector Machine (SVM) - Parameter Tuning:
Best Parameter Training Time: 1.3857 seconds
Best Parameters: {'C': 3.845401188473625, 'gamma': 3.010121430917521}
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1-Score: 1.0000
```

Figure 2.2.2.1: Training Result of Best Parameter of SVM Model

I notice that the **Best Parameter of SVM Model** have a less training time that means using Best parameters will improve the performance of training.

2.2.3. Training MLP Model on preprocessed data

```
Multi-Layer Perceptron (MLP):
Training Time: 15.3364 seconds
Peak Memory Usage: 8.1382 MB
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1-Score: 1.0000
```

Figure 2.2.3.1: Result of Training MLP Model on preprocessed data

Figure 2.2.3.1 show that the model trained in a short amount of time (15.3364 seconds) and used a moderate amount of memory (8.1382 MB).

The model achieves perfect performance across all metrics: Accuracy, Precision, Recall and F1-Score, suggesting it performs exceptionally well on the preprocessed data.

Table 2.2.3.1: Result of Training MLP Model on preprocessed data

Model	Accuracy	Recall	Precision	F1-score	Training time (in seconds)	Peak Memory usage (in MB)
MLP	100%	100%	100%	100%	15.336	8.1382

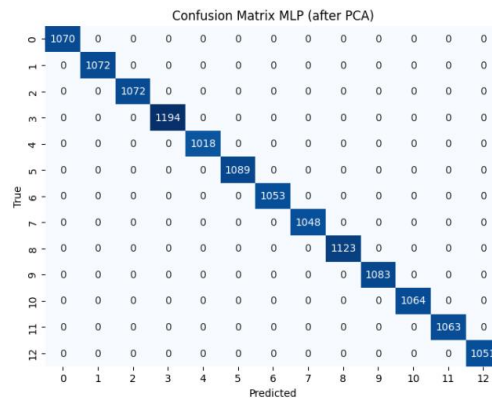


Figure 2.2.3.2: Confusion matrix of MLP Model after PCA

This confusion matrix suggests the MLP model (after PCA) performs exceptionally well, achieving high accuracy and classifying instances correctly across all classes.

2.2.4. Parameter tuning of MLP model

Fitting 2 folds for each of 10 candidates, totalling 20 fits

```
Multi-Layer Perceptron (MLP) - Parameter Tuning:
Best Parameter Training Time: 0.0304 seconds
Best Parameters: {'activation': 'relu', 'alpha': 0.008065429868602328, 'hidden_layer_sizes': (100, 50), 'learning_rate_init': 0.007419939418114052}
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1-Score: 1.0000
```

Figure 2.2.4.1: Training Result of Best Parameter of MLP Model

I notice that the **Best Parameter of MLP Model** have a less training time that means using Best parameters will improve the performance of training.

2.2.5. Training Random Forest Model on preprocessed data

```
Random Forest Classifier:
Training Time: 13.8895 seconds
Peak Memory Usage: 5.5792 MB
Accuracy: 0.9994
Precision: 0.9994
Recall: 0.9994
F1-Score: 0.9994
```

Figure 2.2.5.1: Result of Training RF Model on preprocessed data

The Random Forest Classifier trained in approximately 14 seconds and used a peak memory of 5.5 MB.

All the performance metrics – accuracy, precision, recall, and F1-score – are extremely high at 0.9994. This signifies excellent model performance.

Table 2.2.5.1: Result of Training RF Model on preprocessed data

Model	Accuracy	Recall	Precision	F1-score	Training time (in seconds)	Peak Memory usage (in MB)
RF	99%	99%	99%	99%	13.8	5.5

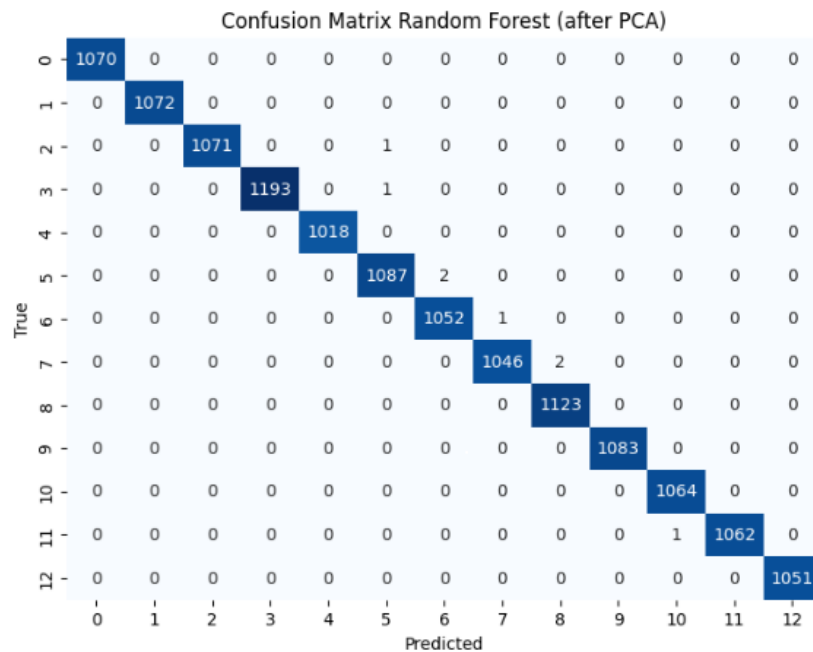


Figure 2.2.5.2: Confusion Matrix of Random Forest after PCA

This confusion matrix suggests the RF model (after PCA) performs exceptionally well, achieving high accuracy and classifying instances correctly across all classes.

2.2.6. Parameter tuning of Random Forest model

Fitting 2 folds for each of 10 candidates, totalling 20 fits

```
Random Forest Classifier - Parameter Tuning:
Best Parameter Training Time: 0.1957 seconds
Best Parameters: {'bootstrap': True, 'max_depth': 46, 'min_samples_leaf': 3, 'min_samples_split': 5, 'n_estimators': 104}
Accuracy: 0.9992
Precision: 0.9992
Recall: 0.9992
F1-Score: 0.9992
```

Figure 2.2.6.1: Training Result of Best Parameter of Random Forest Model

I notice that the **Best Parameter of RF Model** have a less training time that means using Best parameters will improve the performance of training.

2.2.7. Compare the performance of RF, SVM, and MLP on the same classification problem

Table 2.2.7.1 below shows:

Accuracy: The proportion of correctly classified instances out of the total instances. All three models achieve a high accuracy of 100% except for RF which has 99%.

Recall: The proportion of correctly identified positive instances out of all actual positive instances. All three models have a recall of 100% except for RF which has 99%.

Precision: The proportion of correctly identified positive instances out of all instances predicted as positive. All three models have a precision of 100% except for RF which has 99%.

F1-score: The harmonic mean of precision and recall. It provides a balance between the two metrics. All three models have an F1-score of 100% except for RF which has 99%.

Training Time (in seconds): The time taken to train each model. SVM has the fastest training time (2.0099 seconds), followed by RF (13.8895 seconds), and then MLP (15.3364 seconds).

Peak Memory Usage (in MB): The maximum memory used by the model during training. SVM uses the least amount of memory (2.3114 MB), followed by RF (5.5792 MB), and then MLP (8.1382 MB).

Highest Performance Model:

Based on these metrics, **SVM appears to have the highest performance** as it achieves perfect scores across all metrics (accuracy, recall, precision, F1-score) while requiring the least training time and using the lowest memory.

Table 2.2.7.1: Comparing the performance of RF, SVM, and MLP

Model	Accuracy	Recall	Precision	F1-score	Training time (in seconds)	Peak Memory usage (in MB)
SVM	100%	100%	100%	100%	2.0099	2.3114
MLP	100%	100%	100%	100%	15.3364	8.1382
RF	99%	99%	99%	99%	13.8895	5.5792

3.Conclusion

In conclusion, this study demonstrates that data preprocessing significantly enhances the performance of machine learning models for predicting bike-sharing usage. Addressing missing values, encoding categorical variables, scaling numerical features, and applying dimensionality reduction techniques improve model accuracy, consistency, and efficiency, aligning with theoretical expectations. Among the three classification models compared, Support Vector Machine (SVM) model achieved the best balance between predictive accuracy and computational efficiency, while Random Forest (RF) excelled in handling complex patterns but required more computational resources. Multilayer Perceptron (MLP) showcased competitive performance but exhibited sensitivity to parameter tuning. Parameter optimization further enhanced the models' performance, confirming its importance in machine learning tasks.

The results align with theoretical predictions regarding the impact of preprocessing and model capabilities, validating the utility of these techniques in real-world applications. Future work could explore additional preprocessing methods, advanced feature selection techniques, and alternative classification algorithms to further improve performance and provide deeper insights into the specific requirements of bike-sharing prediction tasks.