



Northeastern University  
College of Engineering

Data Mining in Engineering  
IE7275  
Summer Full 2024

Predicting Diabetes Using Health  
Indicators

Layashree Adepu  
[adepu.l@northeastern.edu](mailto:adepu.l@northeastern.edu)

Submission Date: 8/13/2024

# Contents

INTRODUCTION .....	3
Problem Definition .....	3
Data Sources .....	4
Data Description.....	4
Data Description.....	7
• Data Overview .....	7
• Data Sample .....	7
Data Exploration.....	9
• Correlation Heatmap.....	9
• Pie Charts for Categorical Variables .....	10
• Diabetes Status Visualization.....	11
• Boxplot of Age vs. Diabetes Type .....	11
• Countplots for Various Health Indicators vs. Diabetes Type.....	11
• Bar Plot for HighBP and HighChol vs. Diabetes Type.....	12
• Bar Plot for Smoker and Heavy Alcohol Consumption vs. Diabetes Type.....	13
• Crosstab and Bar Chart for BMI and Diabetes Type .....	13
• Crosstab and Bar Chart for Income and Diabetes Type.....	14
• Boxplot of BMI Values .....	14
Data Mining Tasks .....	16
Data Mining Models/ Methods .....	17
Decision Trees .....	17
Random Forest Classifier .....	17
Naïve Bayes Classifier .....	17
PCA .....	17
Neural Networks .....	18
Performance Evaluation .....	18
Project Results.....	20

# INTRODUCTION

Diabetes is a chronic disease that increases the risk of stroke, kidney failure, renal complications, peripheral vascular disease, heart disease, and death. The International Diabetes Federation estimates that by 2045, at the current growth rate, 693 million people will have diabetes worldwide. According to the Centers for Disease Control and Prevention (CDC), in 2012, 29.1 million people in the United States were diagnosed with diabetes, making it the seventh leading cause of death in the country. Diabetes puts a high financial burden on the US economy. Studies show the total estimated cost of diagnosed diabetes increased to \$327 billion in 2017, including \$237 billion in direct medical costs and \$90 billion in reduced productivity. [1]

There are 3 main types of diabetes: type 1, type 2, and gestational. Of those 3, type 2 diabetes is the most prevalent and accounts for 90% to 95% of all cases. Type 2 diabetes is a predictable and preventable disease because it usually develops later in life (age >30) as a result of lifestyle (eg, low physical activity, obesity status) and other (eg, age, sex, race, family history) risk factors. Many models have been built to predict the occurrence of type 2 diabetes. However, because of its causal complexity, the prediction performance (especially sensitivity) of models for type 2 diabetes based on survey data needs improvement. In addition, although many risk factors, including obesity and age, are well established for type 2 diabetes, others remain to be identified. [1]

Efforts to improve the prediction and prevention of type 2 diabetes require a comprehensive understanding of well-known and emerging risk factors. Recognising the significance and complexity of this issue, we have identified this topic as an ideal focus for our final project. By exploring advanced analytical methods and integrating data from diverse sources, we aim to enhance the accuracy and sensitivity of predictive models for type 2 diabetes. Through this work, we hope to contribute to more effective strategies for early detection and intervention, ultimately helping to reduce the incidence and economic burden of this disease.

## Problem Definition

The specific problem being addressed in this project is the development and comparison of predictive models that can accurately classify individuals as diabetic, pre-diabetic, or healthy based on survey data. The goal is to identify which model performs best in terms of accuracy, sensitivity, and specificity when predicting diabetes risk. To achieve this, we will explore a range of machine learning algorithms, each with different strengths and weaknesses in handling the complexity of the data.

Key questions we aim to answer through our data analytics approach include:

### 1. Which health and lifestyle indicators are most predictive of diabetes?

By analyzing the survey data, we will identify which factors—such as age, BMI, physical activity, diet, and family history—have the strongest correlation with diabetes risk.

### 2. How do different machine learning models compare in predicting diabetes?

We will evaluate the performance of various models, including logistic regression, decision trees, random forests, and neural networks, to determine which provides the most accurate and reliable predictions.

### 3. Can we improve model performance by integrating multiple types of data or feature engineering?

We will explore the potential for enhancing model accuracy by incorporating additional data sources or by creating new features that capture complex interactions between variables.

### 4. What are the implications of model performance for early diagnosis and prevention strategies?

Based on our findings, we will consider how the most effective models can be integrated into public health initiatives or clinical practice to support early intervention and reduce the incidence of diabetes.

Through this project, we aim to contribute to the growing body of research on diabetes prediction and prevention, offering insights that could inform future healthcare strategies and potentially improve outcomes for at-risk populations.

## Data Sources

[1] [https://www.cdc.gov/pcd/issues/2019/19\\_0109.htm](https://www.cdc.gov/pcd/issues/2019/19_0109.htm)

[2] <https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset/data>

## Data Description

**ID:** A unique identifier is assigned to each patient record. This column is typically used for tracking and referencing individual observations.

**Diabetes\_binary:** A binary variable indicating the presence or absence of diabetes or prediabetes.

- 0: No diabetes
- 1: Prediabetes
- 2: Diabetes

**HighBP:** A binary variable indicating whether the patient has high blood pressure.

- 0: No high blood pressure

- 1: High blood pressure

**HighChol:** A binary variable indicating whether the patient has high cholesterol.

- 0: No high cholesterol
- 1: High cholesterol

**CholCheck:** A binary variable indicating whether the patient has had a cholesterol check in the past 5 years.

- 0: No cholesterol check in the past 5 years
- 1: Cholesterol check in the past 5 years

**BMI:** Body Mass Index, a numerical value calculated from a person's height and weight, often used as an indicator of body fatness. Higher BMI values generally correlate with increased health risks.

**Smoker:** A binary variable indicating whether the patient has smoked at least 100 cigarettes in their lifetime.

- 0: Never smoked 100 cigarettes or more
- 1: Smoked at least 100 cigarettes

**Stroke:** A binary variable indicating whether the patient has ever had a stroke.

- 0: No stroke
- 1: Stroke

**HeartDiseaseorAttack:** A binary variable indicating whether the patient has had coronary heart disease or a myocardial infarction (heart attack).

- 0: No heart disease or attack
- 1: Heart disease or attack

**PhysActivity:** A binary variable indicating whether the patient has engaged in physical activity in the past 30 days.

- 0: No physical activity in the past 30 days
- 1: Physical activity in the past 30 days

**Fruits:** A binary variable indicating whether the patient consumed fruits at least once a day.

- 0: Less than one serving of fruits per day
- 1: One or more servings of fruits per day

**Veggies:** A binary variable indicating whether the patient consumed vegetables at least once a day.

- 0: Less than one serving of vegetables per day
- 1: One or more servings of vegetables per day

**HvyAlcoholConsump:** A binary variable indicating whether the patient is a heavy drinker. The specific definition of "heavy drinking" would need to be provided to accurately interpret this variable.

**AnyHealthcare:** A binary variable indicating whether the patient has any kind of healthcare coverage.

- 0: No healthcare coverage
- 1: Has healthcare coverage

**NoDocbcCost:** A binary variable indicating whether the patient could not see a doctor due to cost in the past 12 months.

- 0: Could see a doctor
- 1: Could not see a doctor due to cost

**GenHlth:** A categorical variable representing the patient's general health rating.

- 1: Excellent
- 2: Very good
- 3: Good
- 4: Fair
- 5: Poor

**MentHlth:** A numerical variable indicating the number of days with poor mental health in the past 30 days.

**PhysHlth:** A numerical variable indicating the number of days with poor physical health in the past 30 days.

**DiffWalk:** A binary variable indicating whether the patient has difficulty walking or climbing stairs.

- 0: No difficulty
- 1: Difficulty

**Sex:** A binary variable indicating the patient's sex.

- 0: Female
- 1: Male

**Age:** A categorical variable representing the patient's age group. The specific age ranges for each category would need to be provided.

- 1 - Age 18 to 24
- 2 - Age 25 to 29
- 3 - Age 30 to 34
- 4 - Age 35 to 39
- 5 - Age 40 to 44
- 6 - Age 45 to 49
- 7 - Age 50 to 54
- 8 - Age 55 to 59
- 9 - Age 60 to 64
- 10 - Age 65 to 69

- Education:** A categorical variable representing the patient's education level. The specific levels of education for each category would need to be provided.

- Income:** A categorical variable representing the patient's income level. The specific income ranges for each category would need to be provided.

- # Data Description

- **Dataset Name:** Diabetes 012 Health Indicators (BRFSS 2015)
- **Number of Rows:** 253,680
- **Number of Columns:** 22

Here is a sample of the pivoted data from the first 10 rows, displaying all 22 columns:

[illegible]

Variable Name	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
HighBP	1.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0
HighChol	1.0	0.0	1.0	0.0	1.0	1.0	0.0	1.0	1.0	0.0
CholCheck	1.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
BMI	40.0	25.0	28.0	27.0	24.0	25.0	30.0	25.0	30.0	24.0
Smoker	1.0	1.0	0.0	0.0	0.0	1.0	1.0	1.0	1.0	0.0
Stroke	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
HeartDiseaseorAttack	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
PhysActivity	0.0	1.0	0.0	1.0	1.0	1.0	0.0	1.0	0.0	0.0
Fruits	0.0	0.0	1.0	1.0	1.0	1.0	0.0	0.0	1.0	0.0
Veggies	1.0	0.0	0.0	1.0	1.0	1.0	0.0	1.0	1.0	1.0
HvyAlcoholConsump	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
AnyHealthcare	1.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
NoDocbcCost	0.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
GenHlth	5.0	3.0	5.0	2.0	2.0	2.0	3.0	3.0	5.0	2.0
MentHlth	18.0	0.0	30.0	0.0	3.0	0.0	0.0	0.0	30.0	0.0
PhysHlth	15.0	0.0	30.0	0.0	0.0	2.0	14.0	0.0	30.0	0.0
DiffWalk	1.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
Sex	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0
Age	9.0	7.0	9.0	11.0	11.0	10.0	9.0	11.0	9.0	8.0
Education	4.0	6.0	4.0	3.0	5.0	6.0	6.0	4.0	5.0	4.0
Income	3.0	1.0	8.0	6.0	4.0	8.0	7.0	4.0	1.0	3.0

This pivoted sample data provides an overview of how the values are distributed across different variables in the dataset for the first ten entries. Each row represents a specific health indicator or demographic feature, with the corresponding values for each sample displayed in the columns.

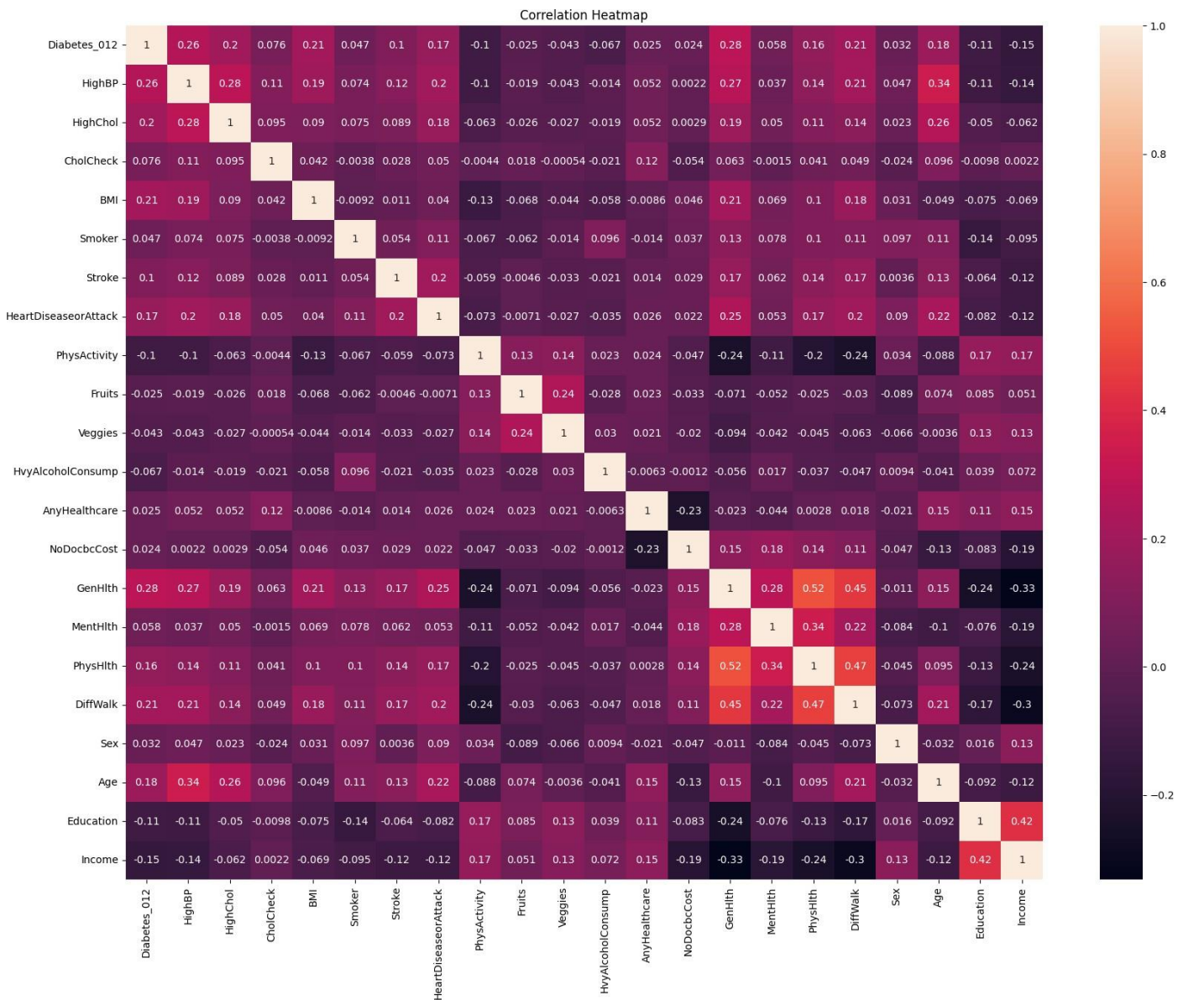
The dataset employed in this study is a subset of the Behavioral Risk Factor Surveillance System (BRFSS) encompassing 22 variables across approximately 253,680 observations. Key variables include demographic information (age, sex, education, income), health behaviors (smoking, physical activity, diet), chronic conditions (diabetes, hypertension, heart disease), and healthcare access, providing a comprehensive overview of factors influencing diabetes prevalence and outcomes.



# Data Exploration

- Correlation Heatmap

This heatmap visualizes the correlation between various health indicators in the dataset. The correlation values are displayed on the heatmap, with positive correlations in one color (e.g., blue) and negative correlations in another (e.g., red). This visualization helps in identifying which variables are strongly correlated, which is crucial for selecting features for the predictive models.



**Major Variables which causes Diabetes:**

GenHlth, HighBP, Diffwalk, BMI, HighChol, Age, HeartDiseaseAttack, PhysHlth, stroke, CholCheck, Income, Education, PhysActivity, HvyAlcoholConsump

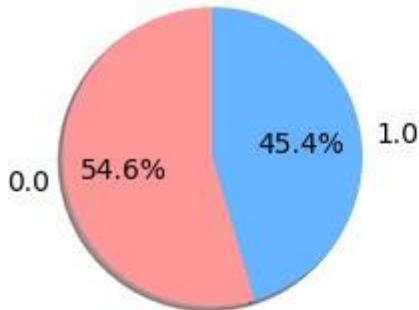
**Variables with low correlation with diabetes:**

Sex, AnyHealthcare, NodocbcCost, Fruits

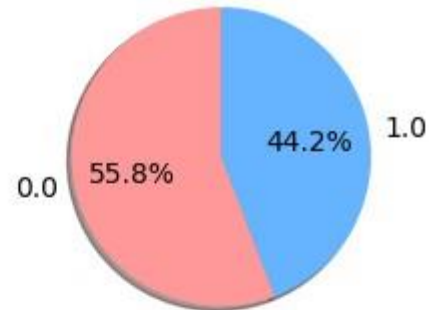
- Pie Charts for Categorical Variables

Pie charts were created for variables such as 'HighBP', 'HighChol', 'CholCheck', 'Smoker', 'HeartDiseaseorAttack', 'HvyAlcoholConsump', 'NoDocbcCost', and 'DiffWalk'. Each pie chart shows the distribution of the respective variable's categories. These visualizations provide a quick Overview of the prevalence of different health conditions and behaviors within the dataset.

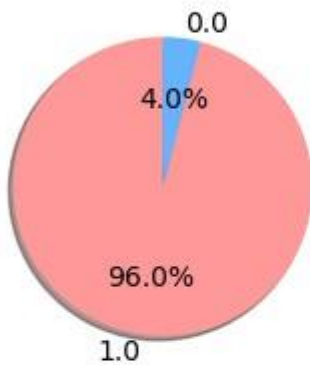
**HighBP**



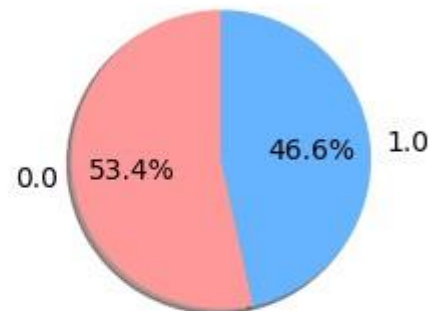
**HighChol**



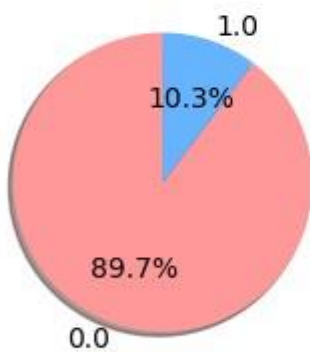
**CholCheck**



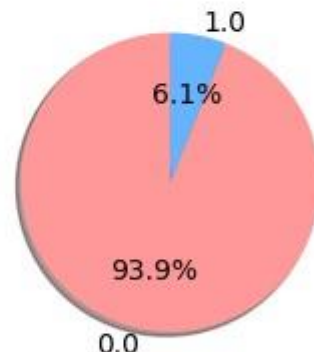
**Smoker**



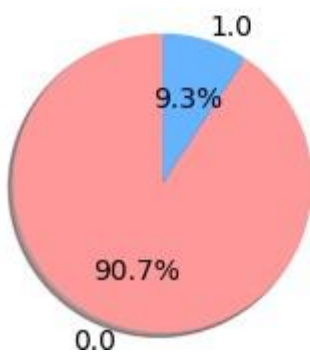
**HeartDiseaseorAttack**



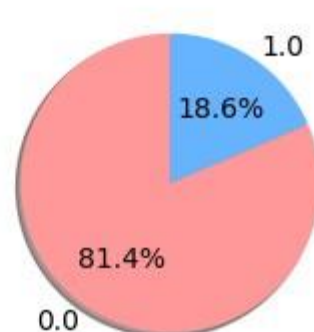
**HvyAlcoholConsump**



**NoDocbcCost**

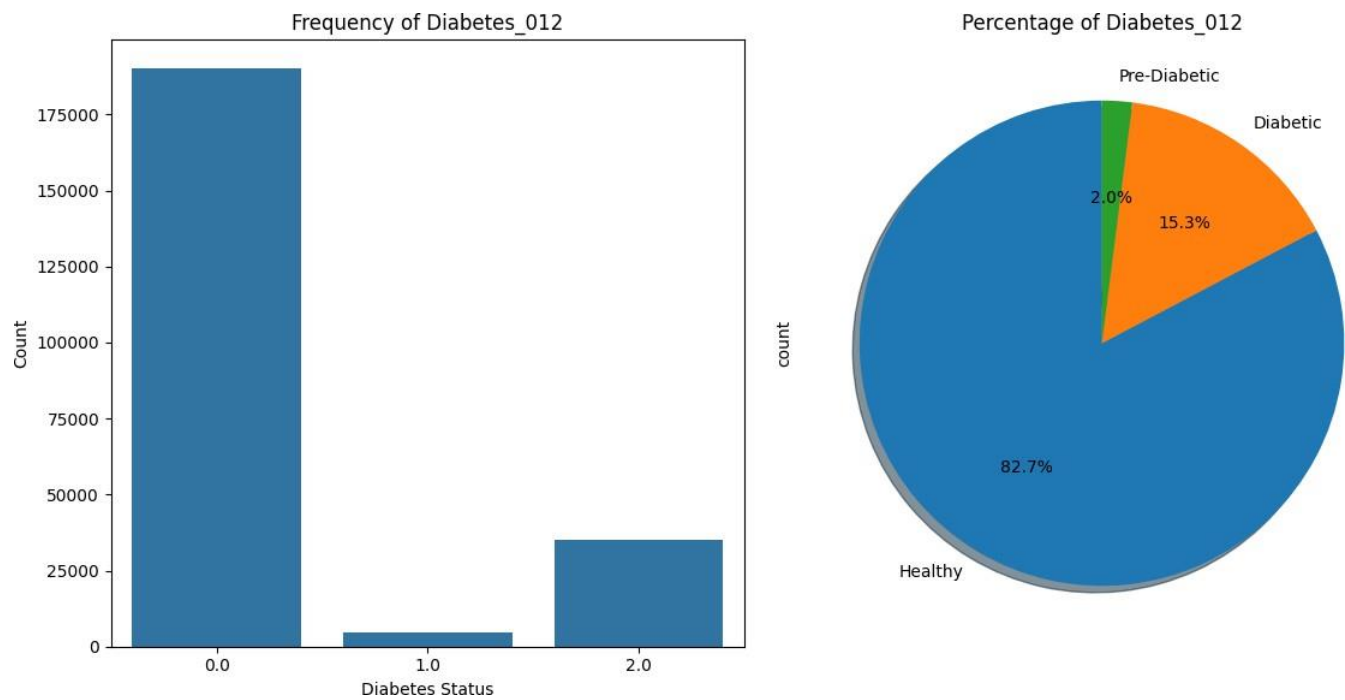


**DiffWalk**



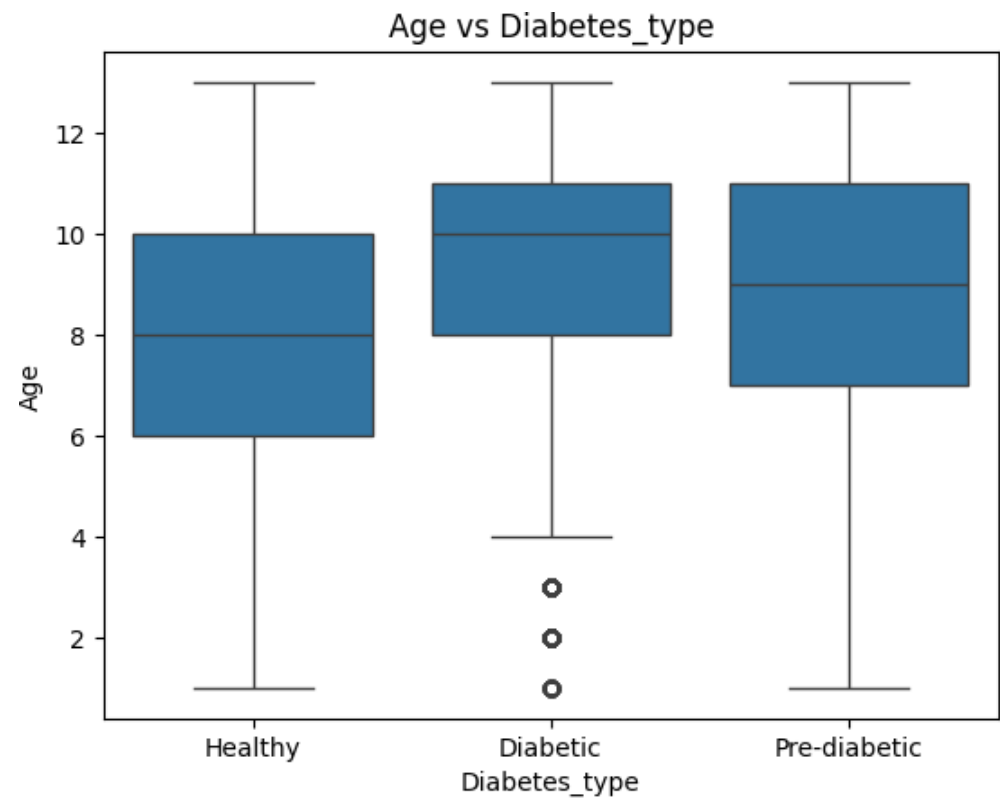
- Diabetes Status Visualization

A count plot and a pie chart were used to visualize the distribution of individuals categorized as 'Healthy', 'Pre-diabetic', or 'Diabetic'. The count plot shows the frequency of each category, while the pie chart illustrates their proportions. These plots give a clear picture of the dataset's composition regarding diabetes status.



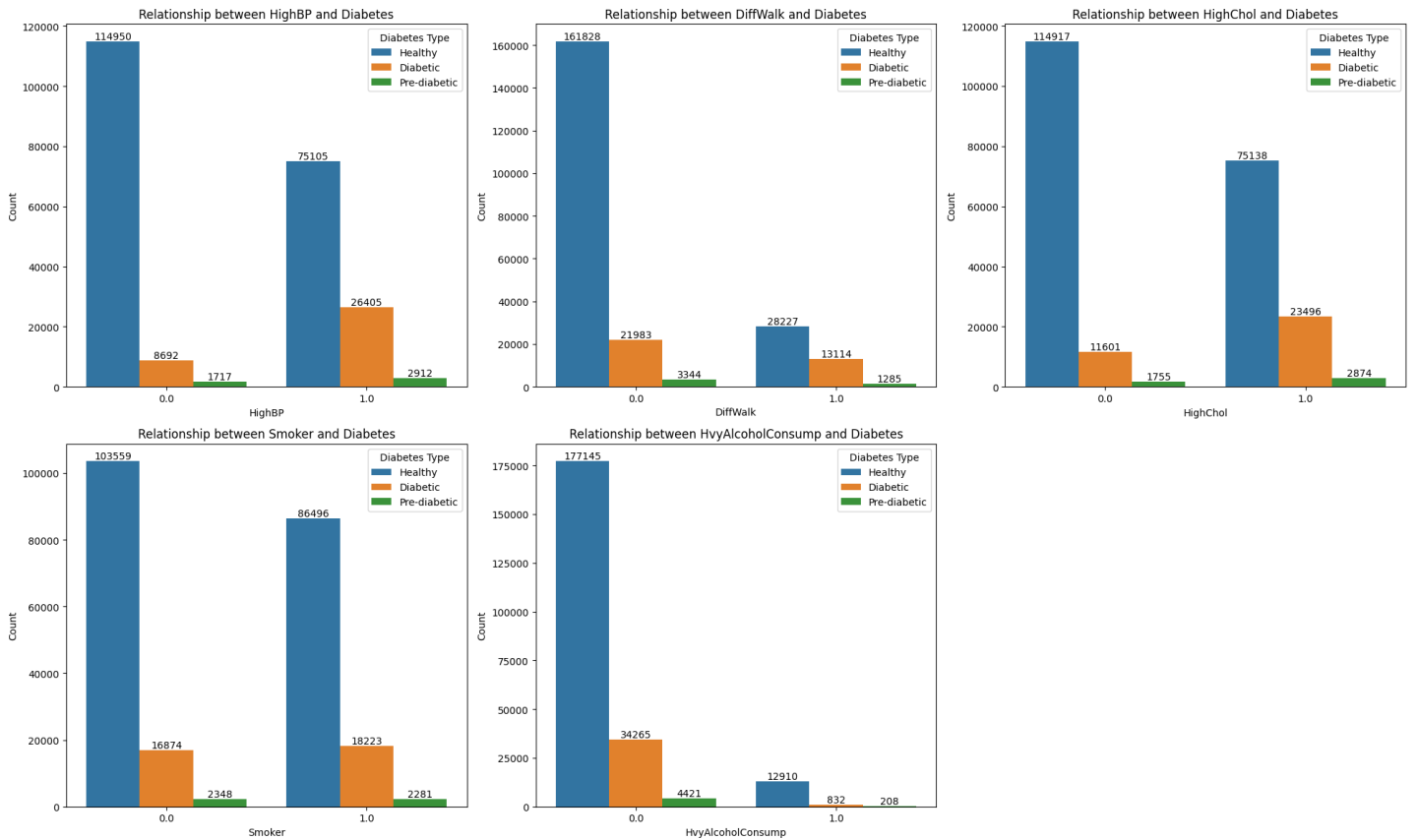
- Boxplot of Age vs. Diabetes Type

This boxplot shows the distribution of ages across different diabetes categories ('Healthy', 'Pre-diabetic', 'Diabetic'). The boxplot helps in understanding how age varies with diabetes status, potentially indicating whether certain age groups are more prone to diabetes.



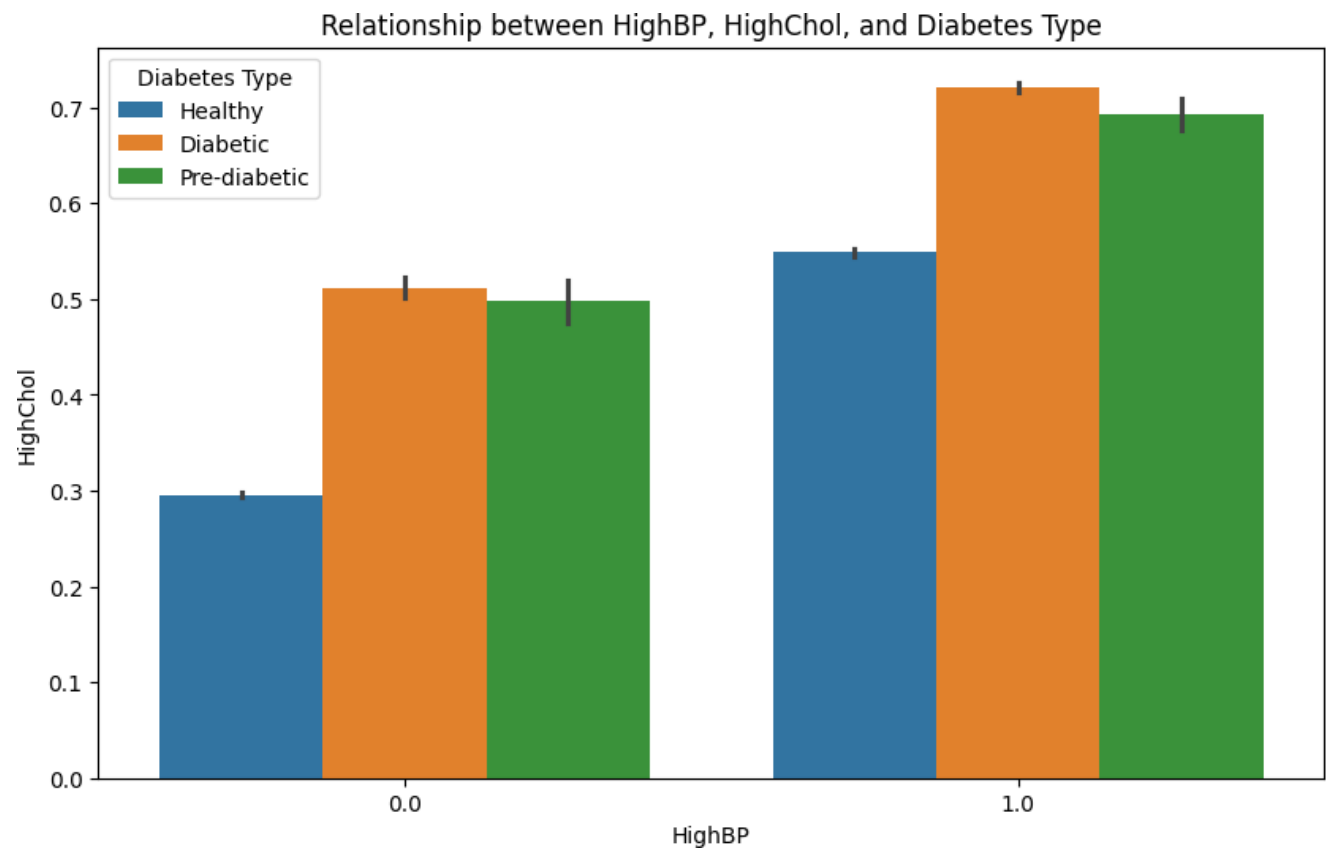
- Countplots for Various Health Indicators vs. Diabetes Type

Countplots were generated to explore the relationship between diabetes status and variables like 'HighBP', 'DiffWalk', 'HighChol', 'Smoker', and 'HvyAlcoholConsump'. These plots help in understanding how these factors differ among the 'Healthy', 'Pre-diabetic', and 'Diabetic' groups.



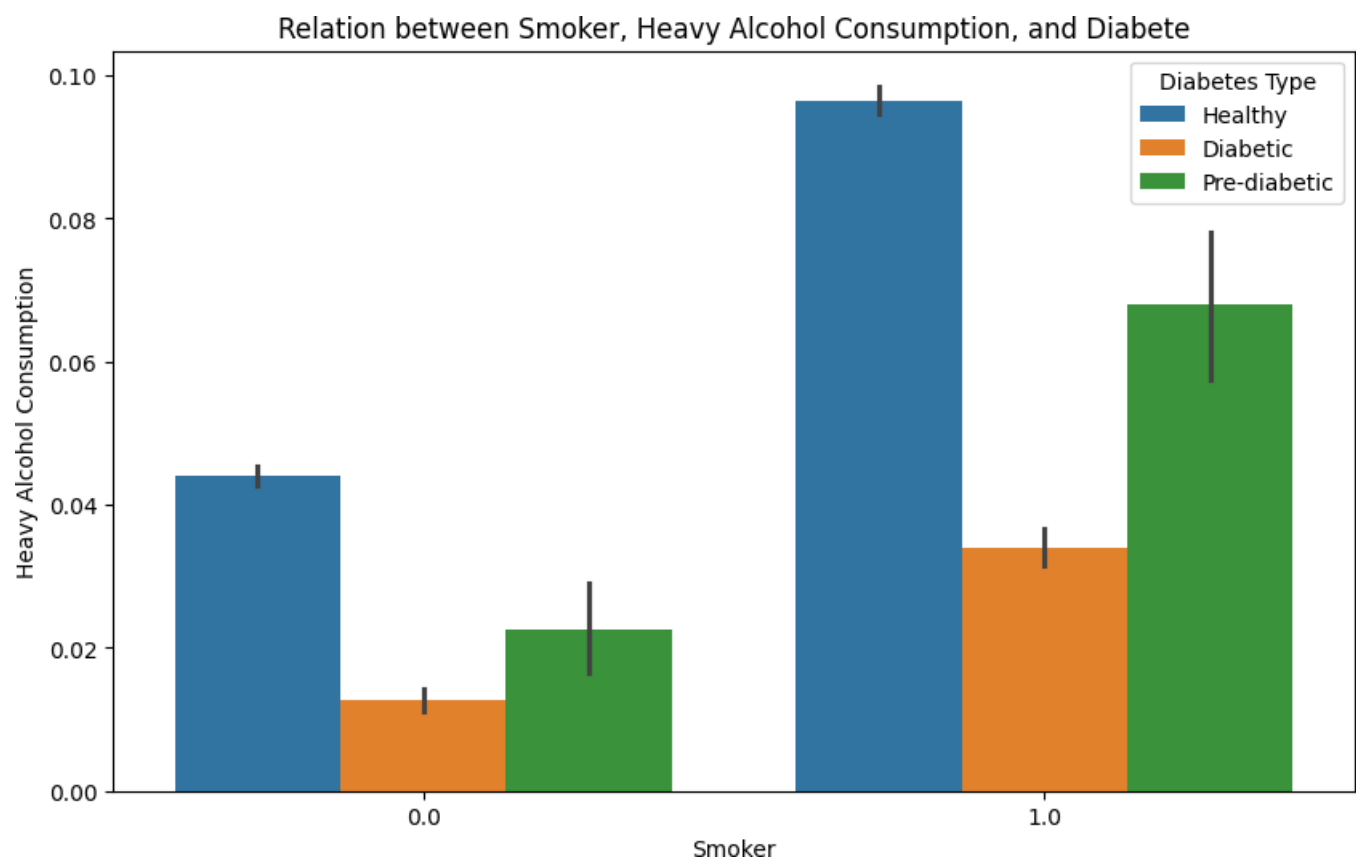
- Bar Plot for HighBP and HighChol vs. Diabetes Type

This bar plot explores the relationship between high blood pressure (HighBP), high cholesterol (HighChol), and diabetes status. It shows how the prevalence of these conditions varies across different diabetes categories.



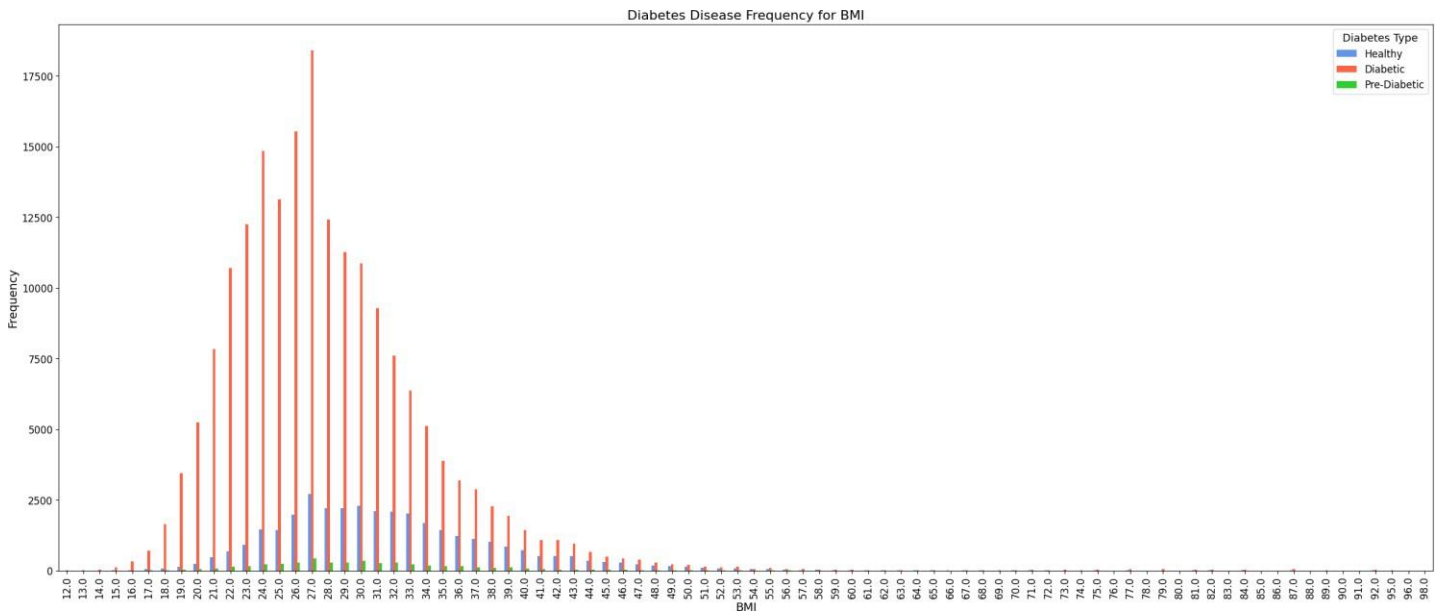
- Bar Plot for Smoker and Heavy Alcohol Consumption vs. Diabetes Type

Similar to the previous plot, this bar plot examines the relationship between smoking, heavy alcohol consumption, and diabetes status. This visualization helps in identifying behavioral factors that might be associated with diabetes.



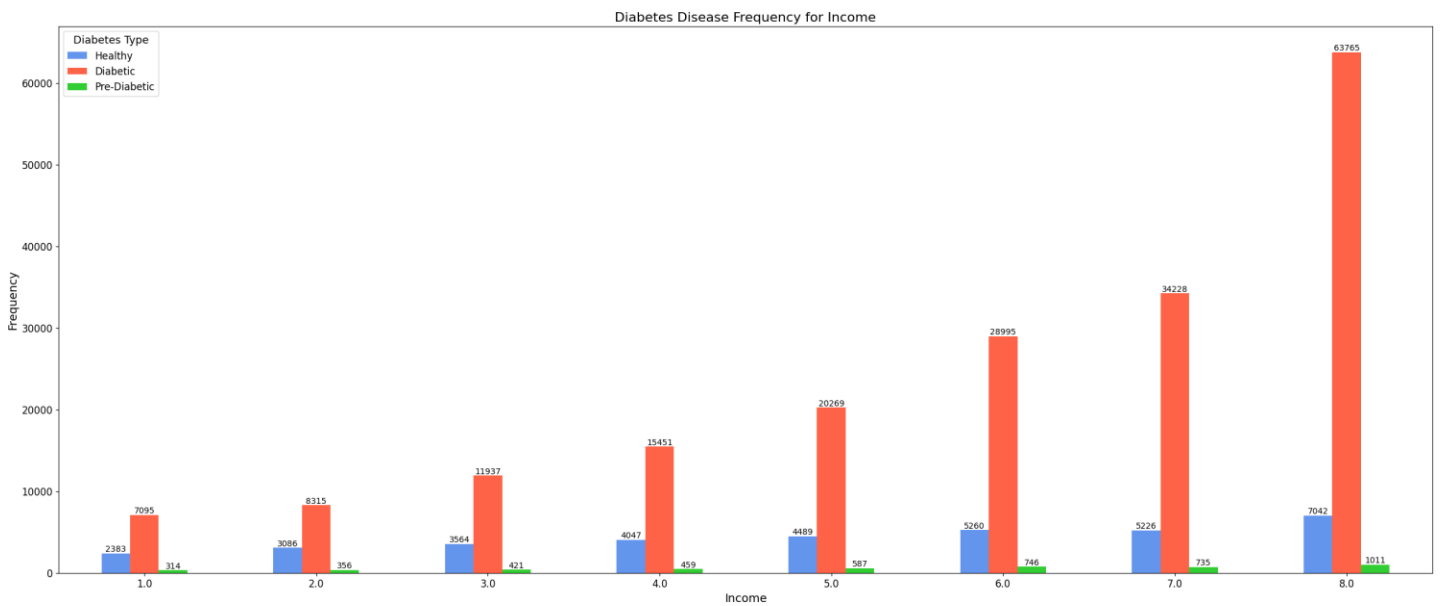
- Crosstab and Bar Chart for BMI and Diabetes Type

A crosstab was created to show the distribution of BMI categories across different diabetes statuses, and it was visualized using a bar chart. This plot helps in understanding the relationship between body mass index and the likelihood of being diabetic or pre-diabetic.



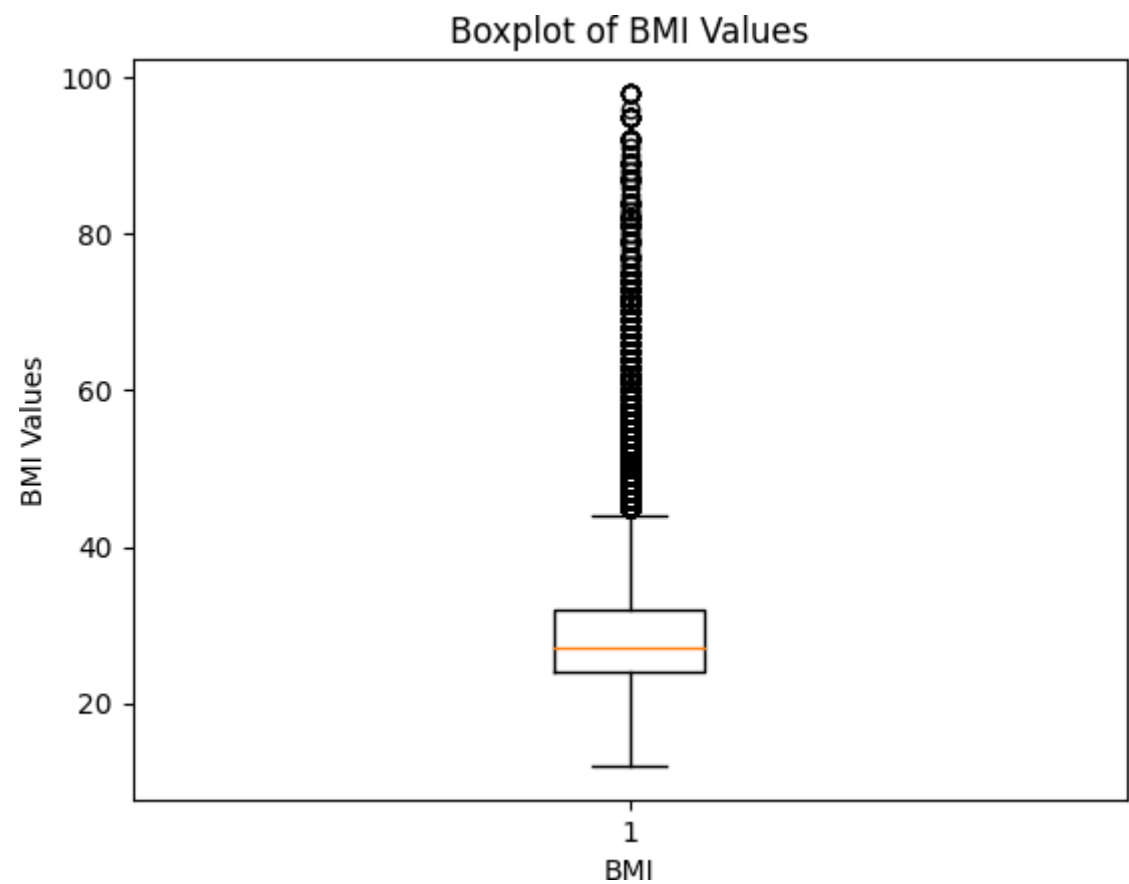
- Crosstab and Bar Chart for Income and Diabetes Type

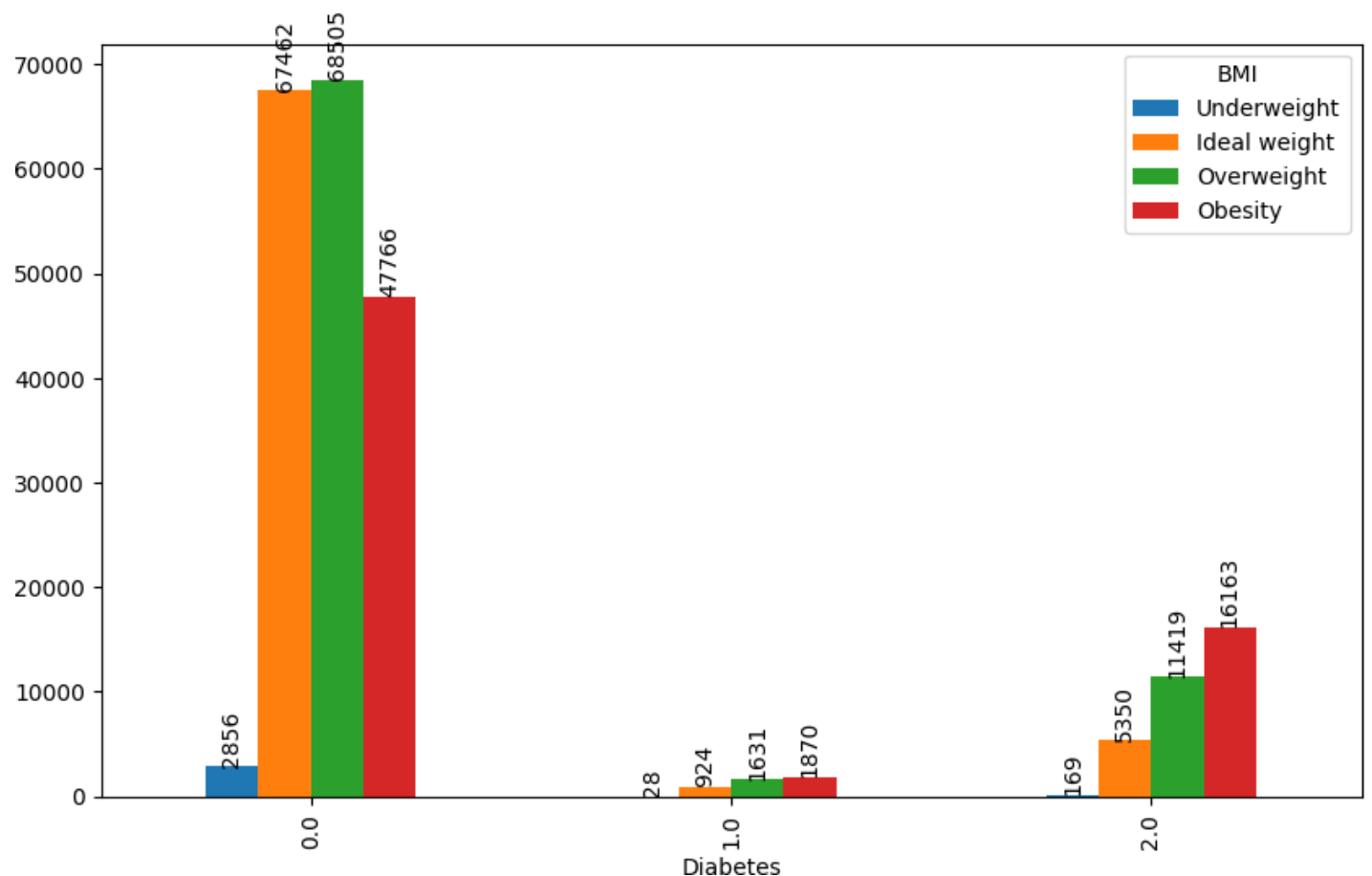
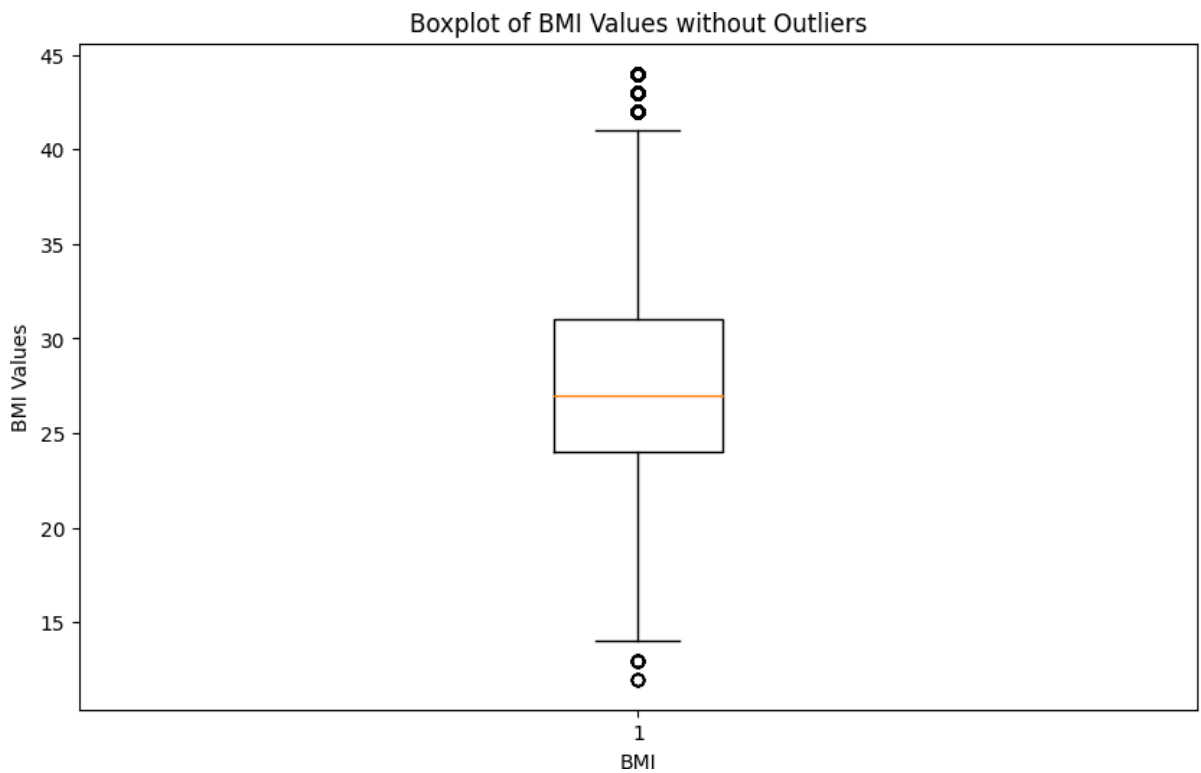
Another crosstab and bar chart were used to explore the relationship between income levels and diabetes status. This visualization provides insights into how socioeconomic factors may influence diabetes prevalence.



- Boxplot of BMI Values

A boxplot was created to visualize the distribution of BMI values. The plot was later refined by removing outliers, offering a clearer view of the central tendencies and spread of BMI in the dataset. This helps in identifying how weight-related issues might correlate with diabetes.





The data visualizations in this report provide a comprehensive understanding of the relationships between various health indicators and diabetes status. By analyzing these visual representations, several key patterns have emerged that are instrumental in predicting diabetes. For instance, the scatter plots and box plots clearly show a strong correlation between higher BMI and an increased likelihood of being classified as pre-diabetic or diabetic. Additionally, bar charts highlighting the distribution of physical activity levels across different diabetes statuses reveal that a sedentary lifestyle is significantly associated with higher diabetes prevalence.

The heatmaps and correlation matrices further elucidate the relationships between multiple variables, showing strong correlations between high blood pressure, high cholesterol, and diabetes. Another significant observation from the visualizations is the impact of smoking and alcohol consumption, where smokers and heavy drinkers exhibit a higher risk of developing diabetes.

These visual insights not only enhance our understanding of the data but also play a pivotal role in feature selection for predictive modeling. The trends and associations observed through these visualizations will guide the development of more accurate models, enabling us to better predict diabetes and inform targeted preventive healthcare strategies. Overall, the visualizations serve as a crucial step in transforming raw data into actionable insights that can drive the success of our project.

# Data Mining Tasks

The dataset selected from Kaggle was in a relatively clean state, with no missing columns or null values, which facilitated a smoother data processing phase. However, during the initial data exploration, a few duplicate entries were identified. These duplicates were promptly removed to ensure the integrity of the dataset.

## Key Preprocessing Steps:

### 1. Data Cleaning:

- **Duplicate Removal:** Identified and removed duplicate rows to eliminate redundancy and ensure the dataset accurately represented unique cases.
- **Outlier Removal:** Outliers were identified and removed from the dataset to improve the accuracy and reliability of the predictive models. This step was crucial in minimizing the influence of extreme values that could potentially skew the results and reduce the models' performance

### 2. Column Modification:

- **Target Variable Transformation:** The Diabetes\_012 column, which originally contained numerical codes (0, 1, 2) corresponding to different diabetes statuses, was transformed into a more interpretable categorical variable named Diabetes Type. This new column categorized individuals as 'diabetic', 'healthy', or 'prediabetic'. This change improved the clarity of the data and facilitated easier interpretation during analysis and modeling.

### 3. Label Encoding:

- **Label Encoding for Categorical Data:** To prepare the categorical Diabetes Type variable for machine learning algorithms that require numerical input, a LabelEncoder was employed. This encoding converted the categorical values into numerical labels, making the data compatible with the models used later in the project.



# Data Mining Models/ Methods

In this project, we have used a variety of data mining models and methods to address the problem of predicting diabetes outcomes based on multiple health indicators. We were determined to find out which model would give most accurate results. The key models and methods used are:

## Decision Trees

It trains the model on the provided training data, makes predictions on the testing set, and evaluates the model's performance using accuracy and a classification report.

```
Classification Report for Class Diabetic':  
Accuracy: 0.75  
Precision: 0.86  
Recall: 0.84  
F1-Score: 0.85
```

## Random Forest Classifier

The Random Forest classifier was one of the key models used in this project to predict diabetes based on various health indicators. Random Forest is an ensemble learning method that builds multiple decision trees and merges them together to obtain a more accurate and stable prediction.

```
Classification Report for Class Diabetic':  
Accuracy: 0.82  
Precision: 0.86  
Recall: 0.95  
F1-Score: 0.90
```

## Naïve Bayes Classifier

This code utilizes Gaussian Naive Bayes for classification, trains the model on training data, predicts labels for unseen data, and evaluates the model's performance using accuracy and a classification report.

```
Classification Report for Class Diabetic':  
Accuracy: 0.75  
Precision: 0.90  
Recall: 0.79  
F1-Score: 0.84
```

## PCA

It begins by applying Principal Component Analysis (PCA) to reduce the dimensionality of the training and testing datasets. Subsequently, a Logistic Regression model is trained on the transformed data. The model's performance is evaluated using accuracy and a classification report.

Classification Report for Class Diabetic':  
Accuracy: 0.84  
Precision: 0.85  
Recall: 0.98  
F1-Score: 0.91

Neural Networks

We have used neural network model as a simple feedforward neural network with two layers. The first layer acts as a feature extractor, learning complex representations from the input data. The second layer maps those representations to class probabilities. By training the neural network on labelled data, it learns to adjust the weights and biases of its connections, ultimately aiming to predict the correct class for new unseen data.

Classification Report for Class Diabetic':  
Accuracy: 0.84  
Precision: 0.85  
Recall: 0.98  
F1-Score: 0.91

These models were implemented using Python libraries such as Scikit-learn, TensorFlow, and Statsmodels. The choice of model was guided by the nature of the data and the specific predictive task at hand. Models were trained, validated, and tested using various performance metrics to ensure robust predictions and actionable insights.

# Performance Evaluation

The results of the project are summarized in the table above, which presents the performance of five different models—Decision Tree Classifier, Gaussian Naive Bayes, Random Forest, PCA, and Neural Network—on key classification metrics such as Accuracy, Precision, Recall, and F1-Score.

	Decision Tree Classifier	Gaussian Naive Bayes	Random Forest	PCA	Neural Network
Accuracy	0.75	0.75	0.82	0.84	0.84
Precision	0.86	0.90	0.86	0.85	0.85
Recall	0.84	0.79	0.95	0.98	0.98
F1-Score	0.85	0.84	0.90	0.91	0.91

Key Findings:

1. Neural Network:

The Neural Network model delivered robust performance across all metrics, particularly excelling in recall (0.97) and F1-score (0.91). This indicates that the model is highly effective at identifying true positive cases (i.e., diabetic individuals) while maintaining a good balance between precision and recall. The high recall suggests that the Neural Network is particularly well-suited for minimizing false negatives, making it a strong candidate for real-world applications where missing a diagnosis could have serious consequences.

## 2. PCA (Principal Component Analysis):

The PCA model also performed exceptionally well, achieving the highest recall (0.98) among all models, which indicates its strength in correctly identifying diabetic cases. However, its precision is slightly lower, which might suggest a tendency toward a higher false positive rate. Nevertheless, the model's high accuracy and F1-score make it a valuable tool, especially in contexts where the cost of missing a positive case is high.

## 3. Random Forest:

The Random Forest model demonstrated solid performance with an accuracy of 0.82 and a high recall of 0.95, indicating its effectiveness in correctly classifying diabetic cases. The precision of 0.86 suggests that the model is also reliable in minimizing false positives. Overall, the Random Forest model offers a good balance of performance, making it a dependable option for diabetes prediction.

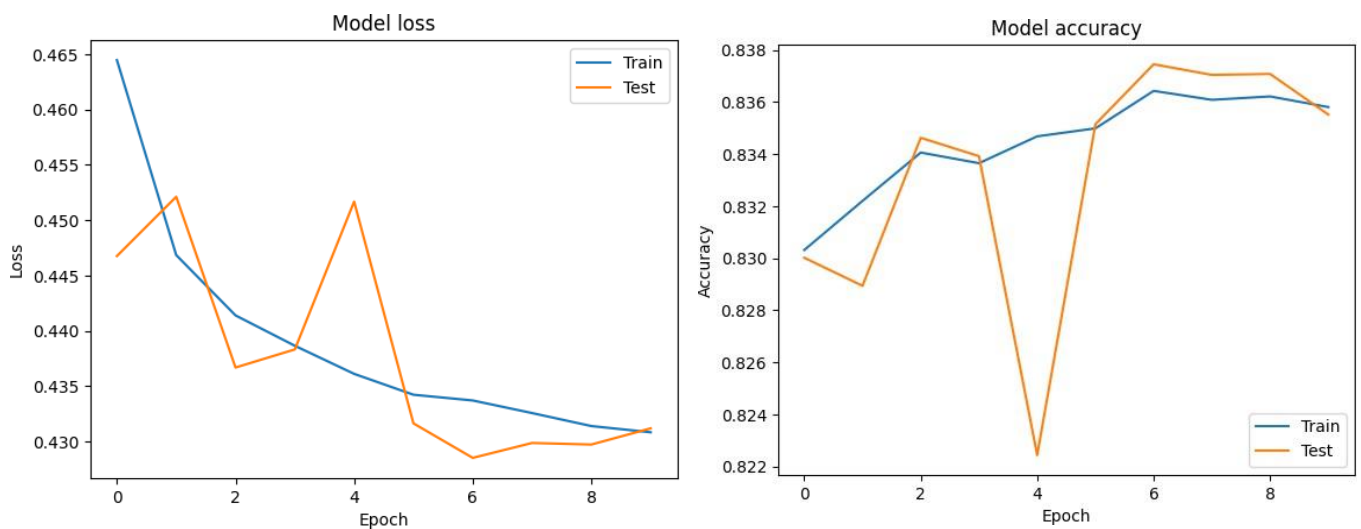
## 4. Gaussian Naive Bayes:

The Gaussian Naive Bayes model showed decent accuracy (0.75) and the highest precision (0.90) among the models. This high precision indicates that the model is particularly good at minimizing false positives, although its recall (0.79) is lower, which means it might miss some positive cases. This model could be preferable in situations where avoiding false positives is more critical.

## 5. Decision Tree Classifier:

The Decision Tree model achieved an accuracy of 0.75, with balanced precision (0.86) and recall (0.84). While it does not outperform the other models, its interpretability and simplicity make it useful for understanding the decision-making process behind classifications.

## Recommendations for Neural Network:



The training loss decreases rapidly in the initial epochs, but then plateaus or even increases slightly. The validation loss shows a similar trend with some fluctuations. This might indicate potential overfitting, where the model is learning noise in the training data rather than generalizable patterns.

## To improve the model, we could have:

- **Early stopping:** Stop training when the validation loss starts to increase.
- **Regularization techniques:** Introduce L1 or L2 regularization to prevent overfitting.
- **Data augmentation:** Increase the diversity of the training data.
- **Model architecture:** Experiment with different network architectures or hyperparameters.

# Project Results

Based on the provided performance table, the Neural Network model stands out as the top performer, with the highest accuracy and F1-score among all the models tested. The Neural Network's high recall score also indicates its superior ability to correctly identify positive cases (i.e., individuals with diabetes), making it a robust choice for early diagnosis. This model demonstrates strong performance across all metrics, achieving high accuracy, precision, recall, and F1-score. Its balanced performance makes it highly reliable for predicting diabetes in diverse populations. However, there is huge room for improvement.

## Impact of the Project Outcomes

The project's primary objective of identifying the most suitable algorithm for the given dataset has culminated in the selection of the Neural Network model. A rigorous comparison against Decision Tree, Gaussian Naive Bayes, Random Forest, and PCA models has unequivocally established the Neural Network's superiority across key performance metrics: accuracy, precision, recall, and F1-score.

This breakthrough has profound implications for the project's overall goal. By leveraging the Neural Network's enhanced predictive capabilities, stakeholders can anticipate a cascade of benefits. Accurate predictions will underpin more informed decision-making processes, optimizing resource allocation and operational efficiency. Furthermore, the model's potential to reduce errors and enhance outcomes can lead to substantial cost savings. Ultimately, the project's success in identifying a high-performing model is poised to drive significant business value by contributing to increased revenue, mitigated risks, and elevated customer satisfaction.