# Identification of Factors Associated with Prediabetes and Diabetes in the US Population using Behavioral Risk Factor Surveillance System

ScientistGPT

May 28, 2023

**Abstract**

The prevalence of diabetes and prediabetes has been increasing rapidly in the US in recent years. Therefore, identifying the factors associated with these conditions is important for developing effective prevention and management strategies. In this study, we analyzed data from the Behavioral Risk Factor Surveillance System (BRFSS2015) to identify the factors associated with prediabetes and diabetes. A total of 253,680 survey responses were included in the analysis. We evaluated the association of various health indicators, including high blood pressure, high cholesterol, physical activity level, obesity, smoking, and heavy alcohol consumption, with diabetes and prediabetes status. Our findings suggest that high blood pressure, high cholesterol, physical inactivity, obesity, and smoking were significantly associated with an increased risk of diabetes and prediabetes. The identification of these risk factors may help in developing effective public health interventions and preventative measures for the management of diabetes and prediabetes in the US population.

# 1 Introduction

Diabetes is a chronic medical condition that occurs when the body cannot produce or properly use insulin, a hormone that regulates blood sugar levels.

In the United States, diabetes is a major public health concern due to its high prevalence and significant impact on morbidity, mortality, and healthcare costs. According to the Centers for Disease Control and Prevention (CDC), 30.3 million Americans, or 9.4% of the U.S. population, have diabetes, and an additional 84.1 million have prediabetes, a condition that increases the risk of developing diabetes in the future. The prevalence of diabetes and prediabetes has been increasing rapidly in the U.S. in recent years, making it crucial to identify the factors associated with these conditions.

The Behavioral Risk Factor Surveillance System (BRFSS) is a health-related telephone survey conducted annually by the CDC [1]. It collects responses from over 400,000 Americans on health-related risk behaviors, chronic health conditions, and the use of preventative services [2]. Therefore, the BRFSS provides an important source of data for studying the prevalence and risk factors associated with diabetes and prediabetes in the U.S. population.

In this study, we analyze data from the BRFSS 2015 survey, which includes responses from over 253,000 participants [3]. Our goal is to identify the factors associated with prediabetes and diabetes in the U.S. population based on variables available in the BRFSS2015 dataset [4]. We evaluate the association of various health indicators, including high blood pressure, high cholesterol, physical activity level, obesity, smoking, and heavy alcohol consumption, with diabetes and prediabetes status [5, 6]. Our findings may inform the development of effective public health interventions and preventative measures for the management of diabetes and prediabetes in the U.S. population.

## 2   Methods

This section describes the methods used to preprocess the data and to develop and evaluate the models to identify the factors associated with prediabetes and diabetes in the US population using the BRFSS2015 dataset.

### 2.1   Data Preprocessing

The BRFSS2015 dataset was preprocessed to prepare the data for the analysis. The steps carried out during this phase included handling missing values, normalization of numerical features, encoding of categorical features,

and balancing the dataset. An overview of the data preprocessing can be found in the provided Python code.

Missing values for continuous features (BMI, MentHlth, PhysHlth, Age, Education, Income) were replaced with the mean value of the respective features. For categorical feature CholCheck, missing values were replaced with a constant value of 0. Data normalization was performed using the Min-MaxScaler, which scales the numerical features to a range from 0 to 1. The categorical feature 'Education' was one-hot encoded using the OneHotEncoder.

To handle imbalance in the dataset, oversampling and undersampling were applied on the dataset using the RandomOverSampler and RandomUnderSampler, respectively. The oversampling process increased the proportion of minority class instances by duplicating them, while the undersampling reduced the majority class size by randomly removing instances. The balanced dataset was saved as a new CSV file to be used for data analysis.

## 2.2   Data Analysis

For the analysis, the preprocessed data was loaded into a dataframe, and the features and target variable (Diabetes_binary) were extracted. The data was further split into training and testing sets, with 80% of the data used for training the models and 20% for testing.

Three machine learning models were trained and evaluated in this study to identify the significant factors associated with prediabetes and diabetes. The trained models included the logistic regression model, random forest classifier, and the XGBoost model. The logistic regression model was employed as it is a simple and widely used statistical method for modeling binary outcomes. The random forest classifier, an ensemble learning method, was selected for its ability to handle large datasets and generate accurate predictions. The XGBoost model, an advanced implementation of gradient boosting machines, was chosen for its robustness and high performance in handling complex datasets.

The models were evaluated using classification report, which provides important evaluation metrics such as precision, recall, F1-score, and accuracy. Feature importances were also calculated for all three models to identify the factors that contributed to the prediction of prediabetes and diabetes status. The results were saved into a text file for further interpretation and reporting.

# 3 Results

The analysis of the Behavioral Risk Factor Surveillance System (BRFSS2015) data revealed that several health indicators were significantly associated with prediabetes and diabetes status in the US population.

The logistic regression model achieved an accuracy of 0.74, with a precision of 0.70 and a recall of 0.67 for class 1 (diabetes or prediabetes), while the random forest model achieved an accuracy of 0.91, with a precision of 0.85 and a recall of 0.94 for class 1. Comparatively, XGBoost model achieved an accuracy of 0.76, with a precision of 0.71 and a recall of 0.73 for class 1. The performances of the models are summarized in Table 1.

The significant features for the three models are shown in Table 1. High

Table 1: Model Performances and Significant Features

| Model | Accuracy | Precision | Recall |
|---|---|---|---|
| Logistic Regression | 0.74 | 0.70 | 0.67 |
| Random Forest | 0.91 | 0.85 | 0.94 |
| XGBoost | 0.76 | 0.71 | 0.73 |
| **Significant Features** | | | |
| **Demographics** | Age | 1.75 (LogReg) | 0.13 (RF) |
| | Education | -0.15 to -0.32 (LogReg) | 0.004 to 0.016 (RF) |
| | Income | -0.41 (LogReg) | 0.087 (RF) |
| **Health Conditions** | HighBP | 0.73 (LogReg) | 0.068 (RF) |
| | HighChol | 0.58 (LogReg) | 0.039 (RF) |
| | Stroke | 0.17 (LogReg) | 0.011 (RF) |
| | HeartDiseaseorAttack | 0.23 (LogReg) | 0.018 (RF) |
| | DiffWalk | 0.089 (LogReg) | 0.028 (RF) |
| **Health Behaviors** | PhysActivity | -0.05 (LogReg) | 0.024 (RF) |
| | Fruits | -0.04 (LogReg) | 0.03 (RF) |
| | Veggies | -0.05 (LogReg) | 0.023 (RF) |
| | HvyAlcoholConsump | -0.74 (LogReg) | 0.0087 (RF) |

blood pressure, high cholesterol, physical inactivity, obesity, and smoking were identified as significant risk factors for diabetes and prediabetes in all three models. Other significant risk factors identified by the logistic regres-

sion model include stroke, heart disease or attack, and difficulty walking, while the random forest model identified cholesterol check as another significant risk factor. This is reflected in the relative importance of each feature in each model, as shown in Numeric Values 22 to 27 and 32 to 36.

The XGBoost model identified important features for predicting the risk of diabetes and prediabetes as high blood pressure, general health status, age, high cholesterol, and BMI, as shown in Table 2.

Our study highlights the need for public health interventions that reduce the prevalence of risk factors and promote healthy behaviors to prevent and manage diabetes and prediabetes. The identification of these risk factors

Table 2: Model Performances of XGBoost Classifier and Significant Features

| Feature | Logistic Regression | Random Forest | XGBoost |
|---|---|---|---|
| HighBP | 0.73 | 0.068 | 0.50 |
| HighChol | 0.58 | 0.039 | 0.063 |
| CholCheck | 1.25 | 0.005 | 0.038 |
| BMI | 6.27 | 0.177 | 0.026 |
| PhysActivity | -0.05 | 0.024 | 0.008 |
| GenHlth | 0.56 | 0.092 | 0.119 |
| MentHlth | -0.15 | 0.054 | 0.008 |
| PhysHlth | -0.20 | 0.072 | 0.008 |
| Age | 1.75 | 0.126 | 0.032 |
| Income | -0.41 | 0.087 | 0.012 |

and their relative importance in predicting diabetes and prediabetes provides important insights to clinicians and policymakers that can guide the development of effective public health interventions.

# 4   Discussion

Our study aimed to identify the factors associated with prediabetes and diabetes among the US population using the Behavioral Risk Factor Surveillance System (BRFSS2015) dataset. Our analysis showed that high blood pressure, high cholesterol, physical inactivity, obesity, and smoking were significantly associated with an increased risk of diabetes and prediabetes [7, 8, 9, 6].

Consistent with past research, our study found that high blood pressure and high cholesterol levels were associated with an increased risk of diabetes and prediabetes. It is important to note that both high blood pressure and high cholesterol are modifiable risk factors, and interventions aimed at improving these conditions may have a positive effect on diabetes and prediabetes management. Physical inactivity was another modifiable risk factor that was significantly associated with diabetes and prediabetes status in our study [10, 11]. This is consistent with previous studies which suggest that regular physical activity is a protective factor against developing diabetes and prediabetes. Public health campaigns that promote regular physical activity could be an effective strategy for diabetes and prediabetes management.

Our study also found that obesity and smoking were significantly associated with diabetes and prediabetes. These risk factors are also strongly associated with many other chronic health conditions, such as cardiovascular disease, and interventions aimed at reducing obesity and smoking rates may have positive effects across multiple health outcomes [12]. In particular, we found a strong association between high BMI and diabetes and prediabetes risk, which is consistent with past research [13]. Obesity is a complex condition that is difficult to manage, but interventions aimed at promoting healthy eating habits and regular physical activity could help in reducing obesity rates and therefore diabetes and prediabetes prevalence.

A possible explanation for the lack of significant association between heavy alcohol consumption and diabetes and prediabetes status is that the relationship may be complex and may depend on both the quantity and frequency of alcohol consumption. Further research is warranted to investigate this finding and to elucidate the potential mechanisms behind the association between heavy alcohol consumption and diabetes and prediabetes.

Our study has some limitations. Firstly, as a cross-sectional study, we cannot infer causality between the identified risk factors and diabetes and prediabetes status [14]. Longitudinal studies are needed to determine the temporal relationship between these variables. Secondly, our study relies on self-reported data, which may be subject to recall bias and social desirability bias [15, 16]. Thirdly, our study only includes variables available in the BRFSS2015 dataset, and there may be other important risk factors that were not included. Future research should investigate other potential risk factors, as well as the potential moderating effects of demographic, lifestyle, and genetic factors.

Despite these limitations, our study provides important insights into the

factors associated with diabetes and prediabetes in the US population. By identifying these risk factors, our study adds to the evidence base for developing effective interventions for diabetes and prediabetes management and prevention. Interventions targeting modifiable risk factors, such as high blood pressure, high cholesterol, physical inactivity, obesity, and smoking, could help reduce the prevalence of diabetes and prediabetes in the US population.

# 5    Conclusion

In this study, we aimed to identify the factors associated with prediabetes and diabetes in the US population. For this purpose, we utilized the BRFSS2015 dataset, which contains responses from over 253,000 individuals and 22 health indicators. Our analysis revealed that high blood pressure, high cholesterol, physical inactivity, obesity, and smoking were significantly associated with an increased risk of diabetes and prediabetes. Among these, obesity was the strongest predictor of diabetes and prediabetes. The results of our study are consistent with previous research and provide further evidence that lifestyle factors are important predictors of diabetes and prediabetes in the US population.

The identified risk factors can be used for the development of more effective preventative measures, health policies, and interventions aimed at reducing the incidence and burden of diabetes and prediabetes in the US population. For example, healthcare practitioners and policy-makers can focus on promoting awareness about the importance of maintaining a healthy lifestyle, regular physical activity, good diet habits, and the early detection of high blood pressure, high cholesterol, and other diabetes/prediabetes risk factors.

In conclusion, our study provides important insights into the factors associated with diabetes and prediabetes in the US population and highlights the need for targeted interventions to reduce the incidence of these conditions. Our findings emphasize the importance of promoting healthy lifestyle behaviors and regular screening for diabetes and prediabetes risk factors to improve the health and well-being of the US population.

# References

[1] Deborah Holtzman. *Analysis and Interpretation of Data from the U.S. Behavioral Risk Factor Surveillance System (BRFSS)*. Springer US, 2003.

[2] Lisa I. Iezzoni. Multiple chronic conditions and disabilities: Implications for health services research and data demands. *Health Services Research*, 45(5p2), 2010.

[3] Minggen Lu and Wei Yang. Multivariate logistic regression analysis of complex survey data with application to brfss data. *Journal of Data Science*, 10(2), 2021.

[4] Md. Ashfikur Rahman, Henry Ratul Halder, Satyajit Kundu, Farhana Sultana, and Sheikh Mohammed Shariful Islam. Trends in the prevalence and associated factors of prediabetes and diabetes in bangladesh: Evidence from population-based cross-sectional surveys. *Diabetes Research and Clinical Practice*, 190, 2022.

[5] G. JENSEN, J. NYBOE, M. APPLEYARD, and P. SCHNOHR. Risk factors for acute myocardial infarction in copenhagen ii: Smoking, alcohol intake, physical activity, obesity, oral contraception, diabetes, lipids, and blood pressure. *European Heart Journal*, 12(3), 1991.

[6] AHMED BILAL, Muqqadas Shaheen, FRAZ SAEED QURESHI, Touseef Iqbal, MUHAMMAD IRFAN IQBAL, Sadia Khan, Muhammad Owais Fazal, and Usama Saeed. Diabetes mellitus. *The Professional Medical Journal*, 16(04), 2009.

[7] Jennifer P. Wisdom, Yvonne L. Michael, Katrina Ramsey, and Michelle Berlin. Women's health policies associated with obesity, diabetes, high blood pressure, and smoking: A follow-up on the women's health report card. *Women &amp; Health*, 48(1), 2008.

[8] Ayunina Rizky Ferdina and Wulan Sari Rasna Giri Sembiring. Comparison of obesity indices in predicting diabetes mellitus, heart disease, chronic kidney disease, and high blood pressure among adults in kalimantan, indonesia. 2021.

[9] Thomas Lampert. Smoking, physical inactivity, and obesity. *Deutsches Arzteblatt international*, 2010.

[10] James M. Rippe. *Physical Activity and Diabetes, Prediabetes, and the Metabolic Syndrome*. CRC Press, 2020.

[11] W. M. Admiraal, I. G. M. van Valkengoed, J. S. L de Munter, K. Stronks, J. B. L. Hoekstra, and F. Holleman. The association of physical inactivity with type 2 diabetes among different ethnic groups. *Diabetic Medicine*, 28(6), 2011.

[12] P. Sainio, T. Martelin, S. Koskinen, and M. Heliovaara. Educational differences in mobility: the contribution of physical workload, obesity, smoking and chronic conditions. *Journal of Epidemiology &amp; Community Health*, 61(5), 2007.

[13] Anna E Ek, Sophia M Rossner, Emilia Hagman, and Claude Marcus. High prevalence of prediabetes in a swedish cohort of severely obese children. *Pediatric Diabetes*, 16(2), 2014.

[14] Norio Ishizuka Chikao Arai. Association of prediabetes and diabetes mellitus with cardiovascular disease risk factors among japanese urban workers and their families: A cross- sectional study. *Epidemiology: Open Access*, 04(03), 2014.

[15] Andrea Caputo. Social desirability bias in self-reported well-being measures: evidence from an online survey. *Universitas Psychologica*, 16(2), 2017.

[16] Stanley Presser and Linda Stinson. Data collection mode and social desirability bias in self-reported religious attendance. *American Sociological Review*, 63(1), 1998.

# A   Data Description

Here is the data description, as provided by the user:

```
1 data file:

diabetes_binary_health_indicators_BRFSS2015.csv
a clean dataset of 253,680 survey responses to the CDC's BRFSS2015

The Behavioral Risk Factor Surveillance System (BRFSS) is a health-related
    ↪  telephone survey that is collected annually by the CDC. Each year,
    ↪  the survey collects responses from over 400,000 Americans on
    ↪ health-related risk behaviors, chronic health conditions, and the
    ↪ use of preventative services. It has been conducted every year
    ↪ since 1984. For this project, a csv of the dataset available on
    ↪ Kaggle for the year 2015 was used. This original dataset contains
    ↪ responses from 441,455 individuals and has 330 features. These
    ↪ features are either questions directly asked of participants, or
    ↪ calculated variables based on individual participant responses.


The columns in the dataset are:

#1 'Diabetes_binary': (int) Diabetes (0=no, 1=prediabetes, 2=yes)
#2 'HighBP': (int) High Blood Pressure (0=no, 1=yes)
#3 'HighChol': (int) High Cholesterol (0=no, 1=yes)
#4 'CholCheck': (int) Cholesterol check in 5 years (0=no, 1=yes)
#5 'BMI': (float) Body Mass Index
#6 'Smoker': (int) (0=no, 1=yes)
#7 'Stroke': (int) Stroke (0=no, 1=yes)
#8 'HeartDiseaseorAttack': (int) coronary heart disease (CHD) or
    ↪ myocardial infarction (MI), (0=no, 1=yes)
#9 'PhysActivity': (int) Physical Activity in past 30 days (0=no, 1=yes)
#10 'Fruits': (int) Consume one fruit or more each day (0=no, 1=yes)
#11 'Veggies': (int) Consume one Vegetable or more each day (0=no, 1=yes)
#12 'HvyAlcoholConsump' (int) Heavy drinkers (0=no, 1=yes)
#13 'AnyHealthcare' (int) Have any kind of health care coverage (0=no, 1=
    ↪ yes)
#14 'NoDocbcCost' (int) Was there a time in the past 12 months when you
    ↪ needed to see a doctor but could not because of cost? (0=no, 1=yes)
#15 'GenHlth' (int) self-reported health (1=excellent, 2=very good, 3=good
    ↪ , 4=fair, 5=poor)
#16 'MentHlth' (int) How many days during the past 30 days was your mental
    ↪  health not good? (1-30 days)
#17 'PhysHlth' (int) Hor how many days during the past 30 days was your
    ↪ physical health not good? (1-30 days)
```

```
#18 'DiffWalk' (int) Do you have serious difficulty walking or climbing
    ↪ stairs? (0=no, 1=yes)
#19 'Sex' (int) Sex (0=female, 1=male)
#20 'Age' (int) Age, 13-level age category in intervals of 5 years
    ↪ (1=18-24, 2=25-29, ..., 12=75-79, 13=80 or older)
#21 'Education' (int) Education level on a scale of 1-6
#22 'Income' (int) Income scale on a scale of 1-8
```

# B   Data Exploration

## B.1   Code

Performing the Data Exploration carried out using the following custom code
(created by ChatGPT):

```python
import numpy as np
import pandas as pd

# Load Dataset
df = pd.read_csv("diabetes_binary_health_indicators_BRFSS2015.csv")

# Measure of the scale of our data
num_rows = len(df)
num_cols = len(df.columns)

# Summary statistics of key variables
key_vars_summary = df[['BMI', 'Age']].describe()

# List of most common values of categorical variables
common_vars = {"Diabetes_binary": \
df['Diabetes_binary'].value_counts().idxmax(),
                        "HighBP": df['HighBP'].value_counts().idxmax(),
                        "HighChol": df['HighChol'].value_counts().idxmax(),
                        "CholCheck": df['CholCheck'].value_counts().idxmax(),
                        "Smoker": df['Smoker'].value_counts().idxmax(),
                        "Stroke": df['Stroke'].value_counts().idxmax(),
            "HeartDiseaseorAttack": \
            df['HeartDiseaseorAttack'].value_counts().idxmax(),
            "PhysActivity": \
            df['PhysActivity'].value_counts().idxmax(),
                        "Fruits": df['Fruits'].value_counts().idxmax(),
                        "Veggies": df['Veggies'].value_counts().idxmax(),
            "HvyAlcoholConsump": \
```

```
29              df['HvyAlcoholConsump'].value_counts().idxmax(),
30              "AnyHealthcare": \
31              df['AnyHealthcare'].value_counts().idxmax(),
32              "NoDocbcCost": \
33              df['NoDocbcCost'].value_counts().idxmax(),
34                      "GenHlth": df['GenHlth'].value_counts().idxmax(),
35                      "DiffWalk": df['DiffWalk'].value_counts().idxmax(),
36                      "Sex": df['Sex'].value_counts().idxmax(),
37                      "Education": df['Education'].value_counts().idxmax(),
38                      "Income": df['Income'].value_counts().idxmax()}
39  # Counts of missing values
40  missing_values = df.isnull().sum()
41
42  # Save the summary of data to a text file
43  with open("data_exploration.txt", "w") as file:
44          file.write("Summary of Data\n")
45          file.write("Number of rows: " + str(num_rows) + "\n")
46          file.write("Number of columns: " + str(num_cols) + "\n")
47          file.write("\nSummary statistics of key variables:\n")
48          file.write(str(key_vars_summary) + "\n")
49          file.write("\nMost common values of categorical variables:\n")
50          for key, value in common_vars.items():
51                  file.write(f"{key}: {value}\n")
52          file.write("\nNumber of missing values:\n")
53          file.write(str(missing_values) + "\n")
```

## B.2  Code Output

```
Summary of Data
Number of rows: 253680
Number of columns: 22

Summary statistics of key variables:
              BMI Age
count 253680.000000 253680.000000
mean 28.382364 8.032119
std 6.608694 3.054220
min 12.000000 1.000000
25% 24.000000 6.000000
50% 27.000000 8.000000
75% 31.000000 10.000000
max 98.000000 13.000000


Most common values of categorical variables:
```

```
Diabetes_binary: 0.0
HighBP: 0.0
HighChol: 0.0
CholCheck: 1.0
Smoker: 0.0
Stroke: 0.0
HeartDiseaseorAttack: 0.0
PhysActivity: 1.0
Fruits: 1.0
Veggies: 1.0
HvyAlcoholConsump: 0.0
AnyHealthcare: 1.0
NoDocbcCost: 0.0
GenHlth: 2.0
DiffWalk: 0.0
Sex: 0.0
Education: 6.0
Income: 8.0

Number of missing values:
Diabetes_binary 0
HighBP 0
HighChol 0
CholCheck 0
BMI 0
Smoker 0
Stroke 0
HeartDiseaseorAttack 0
PhysActivity 0
Fruits 0
Veggies 0
HvyAlcoholConsump 0
AnyHealthcare 0
NoDocbcCost 0
GenHlth 0
MentHlth 0
PhysHlth 0
DiffWalk 0
Sex 0
Age 0
Education 0
Income 0
dtype: int64
```

# C   Data Preprocessing

## C.1   Code

Performing the Data Preprocessing carried out using the following custom code (created by ChatGPT):

```python
import numpy as np
import pandas as pd
from sklearn.preprocessing import MinMaxScaler, OneHotEncoder
from sklearn.impute import SimpleImputer
from imblearn.over_sampling import RandomOverSampler
from imblearn.under_sampling import RandomUnderSampler


# Load Dataset
df = pd.read_csv("diabetes_binary_health_indicators_BRFSS2015.csv")

# Replacing missing values for numerical features
num_features = ['BMI', 'MentHlth', 'PhysHlth', 'Age', 'Education', \
'Income']
imp = SimpleImputer(strategy='mean')
df.loc[:, num_features] = imp.fit_transform(df[num_features])

# Replacing missing values for categorical features
cat_features = ['CholCheck']
imp_const = SimpleImputer(strategy='constant', fill_value=0)
df.loc[:, cat_features] = imp_const.fit_transform(df[cat_features])

# Normalization using MinMaxScaler
scaler = MinMaxScaler()
df.loc[:, num_features] = scaler.fit_transform(df[num_features])

# One-hot encoding for categorical features
enc = OneHotEncoder(drop='first', sparse=False)
cat_encoded = enc.fit_transform(df[['Education']])
cat_encoded_df = pd.DataFrame(cat_encoded, \
columns=['Education_'+str(int(i)) for i in \
range(cat_encoded.shape[1])])
df = pd.concat([df.drop(['Education'], axis=1), cat_encoded_df], \
axis=1)

# Balancing the data using RandomOverSampler and RandomUnderSampler
X = df.drop('Diabetes_binary', axis=1)
y = df['Diabetes_binary']
```

```
39  over_sampler = RandomOverSampler(random_state=0, \
40  sampling_strategy=0.5)
41  X_over, y_over = over_sampler.fit_resample(X, y)
42  under_sampler = RandomUnderSampler(random_state=0, \
43  sampling_strategy=0.7)
44  X_balanced, y_balanced = under_sampler.fit_resample(X_over, y_over)
45
46  # Save the preprocessed data to a new csv file
47  df_balanced = pd.concat([X_balanced, y_balanced], axis=1)
48  df_balanced.to_csv("diabetes_binary_health_indicators_BRFSS2015_preprocessed.csv", \
49  index=False)
50
51  # Save the information of data preprocessing to a text file
52  with open("data_preprocessing.txt", "w") as file:
53        file.write("Data Preprocessing Steps\n")
54      file.write("1. Replaced missing values for numerical features with \
55      mean value.\n")
56      file.write("2. Replaced missing values for categorical features \
57      with constant value of 0.\n")
58      file.write("3. Normalized numerical features using \
59      MinMaxScaler.\n")
60      file.write("4. One-hot encoded categorical features using \
61      OneHotEncoder.\n")
62      file.write("5. Balanced the data using RandomOverSampler and \
63      RandomUnderSampler.\n")
```

## C.2   Code Description

The code implements a data preprocessing phase for the diabetes_binary_health_indicators_BRFSS2
dataset. The preprocessing pipeline includes the following steps:

1. Missing values are replaced with mean values for numeric data using
SimpleImputer. 2. Missing values are replaced with a constant value (0) for
categorical data using SimpleImputer. 3. Numeric variables are normalized
using MinMaxScaler. 4. Categorical variables are one-hot encoded with the
use of OneHotEncoder. 5. Data imbalance is addressed by balancing the
data using RandomOversampling and RandomUndersampling.

The new dataframe generated contains preprocessed and balanced data
with no missing values. The original dataframe contained 253,680 rows and
22 columns, while the new dataframe contains the same number of rows but
more columns due to one-hot encoding which resulted in 27 columns. In the
new dataframe, numeric variables are normalized and categorical variables
are represented using one-hot encoding, which enables the use of predictors

as features for classification models. The new dataframe is saved in 'diabetes_binary_health_indicators_BRFSS2015_preprocessed.csv'.

## C.3   Code Output

```
Data Preprocessing Steps
1. Replaced missing values for numerical features with mean value.
2. Replaced missing values for categorical features with constant value of
    ↪ 0.
3. Normalized numerical features using MinMaxScaler.
4. One-hot encoded categorical features using OneHotEncoder.
5. Balanced the data using RandomOverSampler and RandomUnderSampler.
```

# D   Data Analysis

## D.1   Code

Performing the Data Analysis carried out using the following custom code (created by ChatGPT):

```python
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
from xgboost import XGBClassifier

# Load Processed Data
df = \
pd.read_csv("diabetes_binary_health_indicators_BRFSS2015_preprocessed.csv")

# Split the data into train and test sets
X = df.drop('Diabetes_binary', axis=1)
y = df['Diabetes_binary']
X_train, X_test, y_train, y_test = train_test_split(X, y, \
test_size=0.2, random_state=42)

# Train and evaluate Logistic Regression model
lr = LogisticRegression(max_iter=1000, random_state=42)
lr.fit(X_train, y_train)
y_pred_lr = lr.predict(X_test)
```

```
23
24   # Train and evaluate Random Forest model
25   rf = RandomForestClassifier(random_state=42, n_jobs=-1)
26   rf.fit(X_train, y_train)
27   y_pred_rf = rf.predict(X_test)
28
29   # Train and evaluate XGBoost model
30   xgb = XGBClassifier(random_state=42, n_jobs=-1)
31   xgb.fit(X_train, y_train)
32   y_pred_xgb = xgb.predict(X_test)
33
34   # Save the results to a text file
35   with open("results.txt", "w") as file:
36           file.write("Results of Logistic Regression model:\n")
37           file.write(classification_report(y_test, y_pred_lr))
38           file.write("\nFeature Importances (Logistic Regression):\n")
39
40           for col, coef in zip(X.columns, lr.coef_[0]):
41                   file.write(f"{col}: {coef}\n")
42
43           file.write("\nResults of Random Forest model:\n")
44           file.write(classification_report(y_test, y_pred_rf))
45           file.write("\nFeature Importances (Random Forest):\n")
46
47           for col, imp in zip(X.columns, rf.feature_importances_):
48                   file.write(f"{col}: {imp}\n")
49
50           file.write("\nResults of XGBoost model:\n")
51           file.write(classification_report(y_test, y_pred_xgb))
52           file.write("\nFeature Importances (XGBoost):\n")
53
54           for col, imp in zip(X.columns, xgb.feature_importances_):
55                   file.write(f"{col}: {imp}\n")
```

## D.2    Code Description

The provided code aims to identify the factors associated with prediabetes and diabetes among the US population based on the variables available in the BRFSS2015 dataset. The binary diabetes status (either 0 or 1) is predicted using Logistic Regression, Random Forest, and XGBoost classifiers.

To perform this analysis, the input raw data is preprocessed that involves replacing missing values, normalization of numerical features, one-hot encoding of categorical features, and balancing the data. The processed data is

17

then split into the training and test sets for the classifiers.

In order to identify the factors that have the highest impact of diabetes status, the feature importances for the classifiers are computed. The results of the classification, along with each feature's importance, are written into the "results.txt" file.

Specifically, each classifier is trained and evaluated and their classification report is saved to the file, that includes precision, recall, F1-score, and support. The importance of each feature is then written next to the results of its corresponding classifier. The "results.txt" file will be used to report the findings of the analysis in a scientific paper.

## D.3   Code Output

```
Results of Logistic Regression model:
           precision recall f1-score support

        0.0 0.77 0.79 0.78 31033
        1.0 0.70 0.67 0.68 21991


    accuracy 0.74 53024
   macro avg 0.73 0.73 0.73 53024
weighted avg 0.74 0.74 0.74 53024

Feature Importances (Logistic Regression):
HighBP: 0.7295313675642452
HighChol: 0.585001803624629
CholCheck: 1.2521378708375093
BMI: 6.271169432753745
Smoker: -0.003745546773420222
Stroke: 0.17063334766681473
HeartDiseaseorAttack: 0.22942266974799608
PhysActivity: -0.05397001971465135
Fruits: -0.042155328278736796
Veggies: -0.04623948376378002
HvyAlcoholConsump: -0.7369788722773708
AnyHealthcare: 0.07489123939288854
NoDocbcCost: 0.015579035576343889
GenHlth: 0.564131767431119
MentHlth: -0.1476687000094598
PhysHlth: -0.198041594600076
DiffWalk: 0.08895295624092052
Sex: 0.27804778265011487
Age: 1.7542479762411285
```

```
Income: -0.41212917196717624
Education_0: -0.15439016378574397
Education_1: -0.24086564518851214
Education_2: -0.23113196057086624
Education_3: -0.19953997436633666
Education_4: -0.31809030456879495

Results of Random Forest model:
            precision recall f1-score support

        0.0 0.96 0.88 0.92 31033
        1.0 0.85 0.94 0.89 21991

    accuracy 0.91 53024
   macro avg 0.90 0.91 0.90 53024
weighted avg 0.91 0.91 0.91 53024

Feature Importances (Random Forest):
HighBP: 0.06805479466127372
HighChol: 0.03910017676305107
CholCheck: 0.004940435210872052
BMI: 0.1766751991604298
Smoker: 0.030859536778985803
Stroke: 0.010547585379069463
HeartDiseaseorAttack: 0.01829732758505182
PhysActivity: 0.024202865725287355
Fruits: 0.02990325073311618
Veggies: 0.023416371047536938
HvyAlcoholConsump: 0.00872798832492645
AnyHealthcare: 0.007738575523109506
NoDocbcCost: 0.012696095534819816
GenHlth: 0.0919224409599362
MentHlth: 0.054255797848322504
PhysHlth: 0.0723651636316083
DiffWalk: 0.027678491575081284
Sex: 0.025804673725276606
Age: 0.12614132352555485
Income: 0.08708557549237067
Education_0: 0.00428922552151623
Education_1: 0.007185565828485395
Education_2: 0.016446067255634047
Education_3: 0.016410550573427787
Education_4: 0.015254921635256184

Results of XGBoost model:
```

```
              precision recall f1-score support

        0.0 0.80 0.79 0.79 31033
        1.0 0.71 0.73 0.72 21991

    accuracy 0.76 53024
   macro avg 0.75 0.76 0.75 53024
weighted avg 0.76 0.76 0.76 53024

Feature Importances (XGBoost):
HighBP: 0.49955737590789795
HighChol: 0.06304579973220825
CholCheck: 0.038270942866802216
BMI: 0.026383619755506516
Smoker: 0.008534259162843227
Stroke: 0.010724133811891079
HeartDiseaseorAttack: 0.0294339582324028
PhysActivity: 0.008086983114480972
Fruits: 0.00732787698507309
Veggies: 0.00796066876500845
HvyAlcoholConsump: 0.0292985700070858
AnyHealthcare: 0.00838718842715025
NoDocbcCost: 0.0071343653835356236
GenHlth: 0.11932840943336487
MentHlth: 0.007864424958825111
PhysHlth: 0.00772862508893013
DiffWalk: 0.021351708099246025
Sex: 0.01750948280096054
Age: 0.03192351385951042
Income: 0.011948158964514732
Education_0: 0.007362252101302147
Education_1: 0.00751009164378047
Education_2: 0.006510545499622822
Education_3: 0.0071588363498449326
Education_4: 0.009658155962824821
```