



Predicting Diabetes Probability using Machine Learning

Using machine learning techniques to forecast the likelihood of individuals developing diabetes based on various factors.

Introduction



Presentation Overview

Brief introduction to the presentation, highlighting the key topics to be covered.



Diabetes Prediction

Using machine learning models to predict the probability of an individual developing diabetes.



Machine Learning Approach

Outlining the specific machine learning algorithms and techniques employed in the analysis.

Diabetes: The Challenge

Diabetes is a chronic condition characterized by high blood sugar levels, which can have serious health consequences if left unmanaged. Early detection is crucial, as it allows for timely intervention and the implementation of effective management strategies. Accurate prediction models can play a vital role in identifying individuals at risk, enabling proactive measures to prevent or delay the onset of the disease.



Data Collection and Preprocessing

- BRFSS Data

The Behavioral Risk Factor Surveillance System (BRFSS) is a health-related telephone survey that is collected annually by the CDC. For this project, a csv of the dataset available on Kaggle for the year 2015 was used.

- Feature Engineering

Explanation of the process of creating new features from the raw data, such as calculating derived metrics, encoding categorical variables, and handling date and time-related features.

- Preprocessing

Description of the steps taken to ensure data quality, including handling missing values, removing outliers, and addressing inconsistencies in data formats and units.

Machine Learning Algorithms

Logistic Regression

A widely used algorithm for binary classification tasks, such as predicting whether a patient has diabetes or not.

Decision Tree

A tree-based algorithm that creates a series of rules to make predictions, useful for both classification and regression problems.

Random Forest

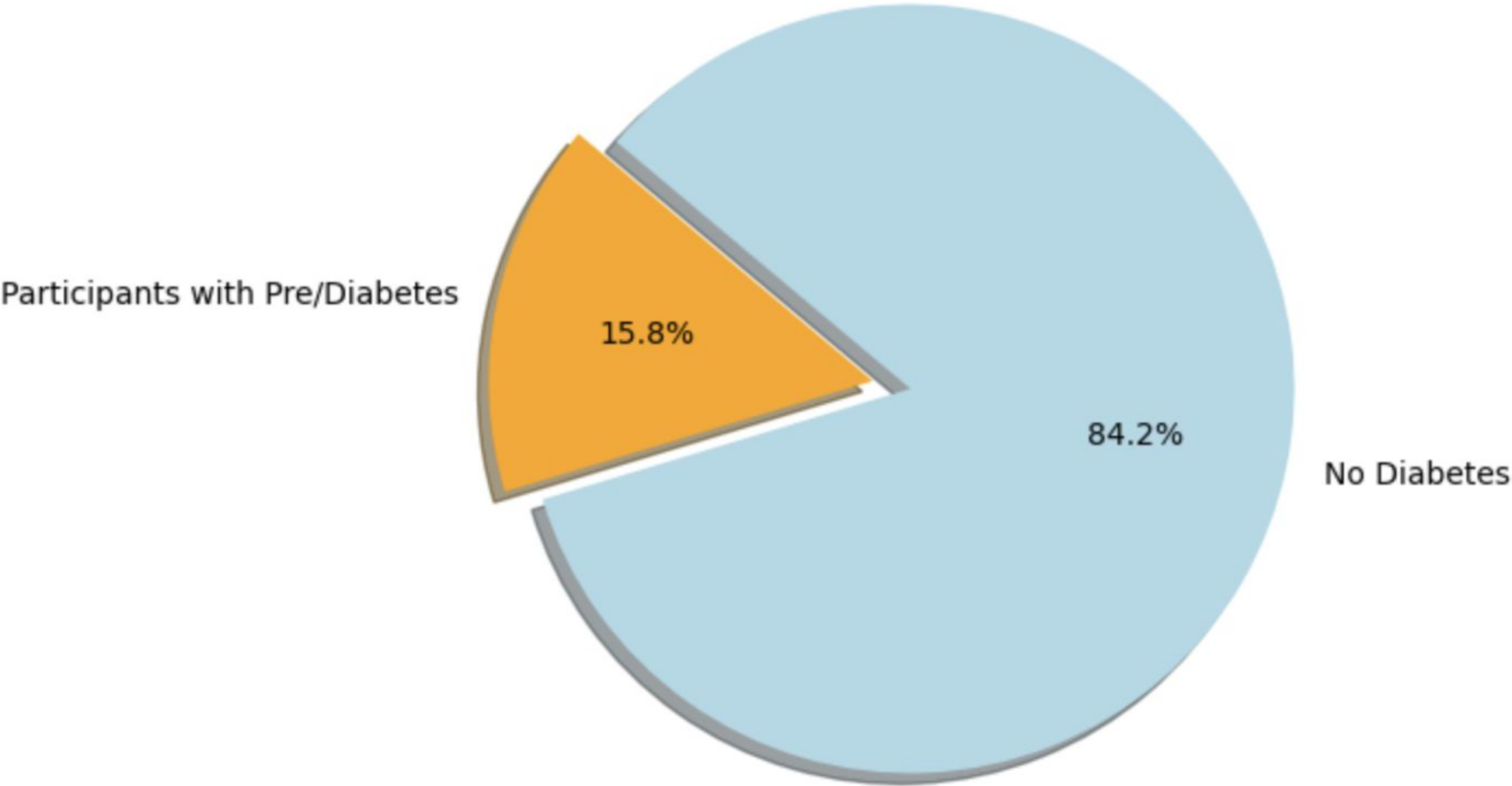
An ensemble learning method that combines multiple decision trees to improve the accuracy and robustness of the predictions.

K-Means

Unsupervised machine learning algorithm used for clustering data. Clustering is the task of grouping data points into clusters or groups based on their similarity.

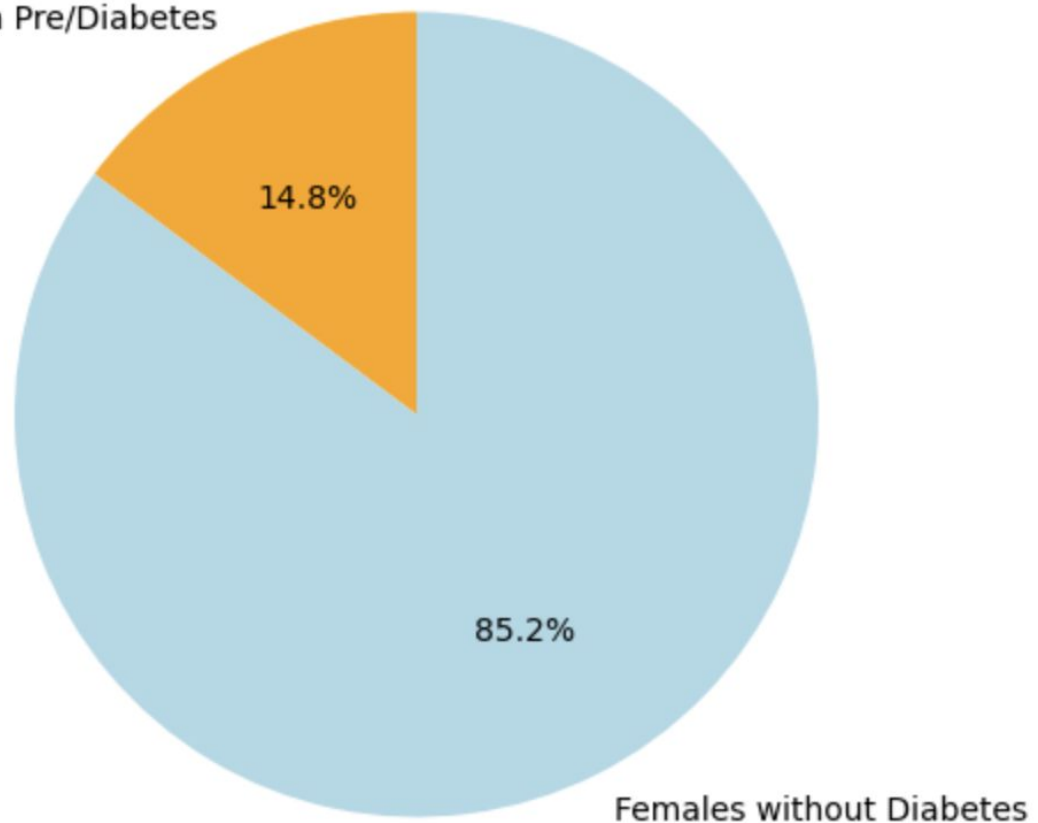
How well does the data set provide accurate predictions of whether an individual has diabetes?

Overall Percentage of Participants with Pre/Diabetes



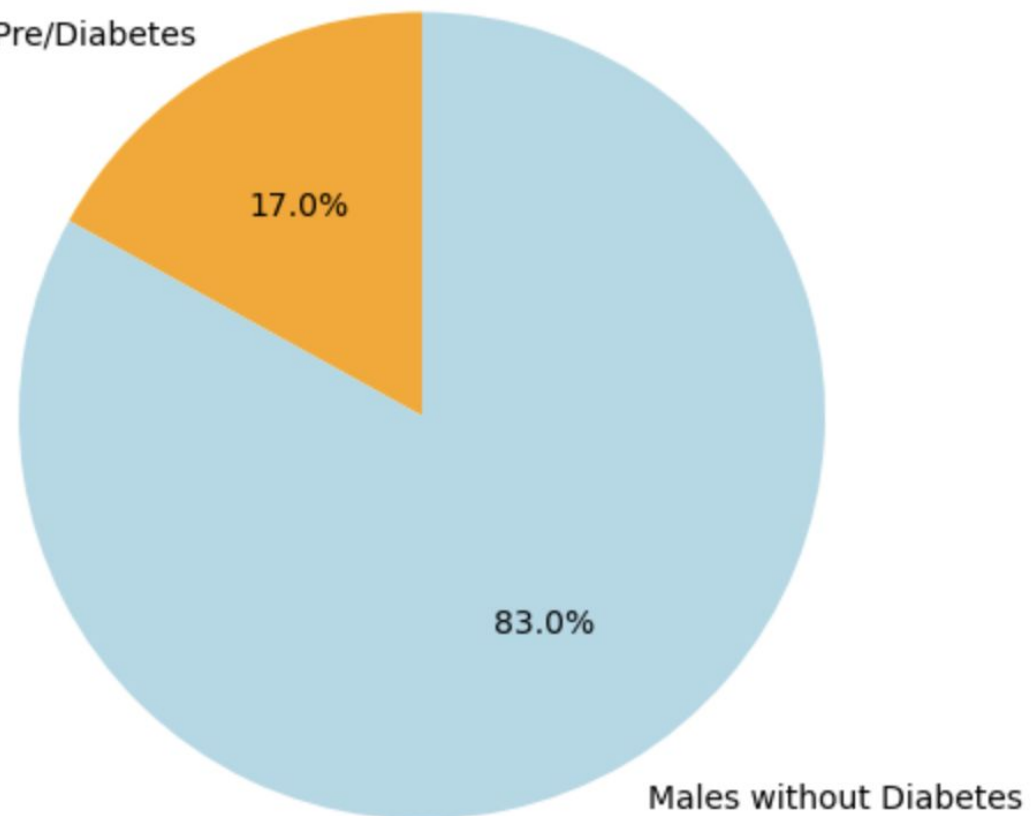
Pre/Diabetes Distribution Among Female Participants

Females with Pre/Diabetes

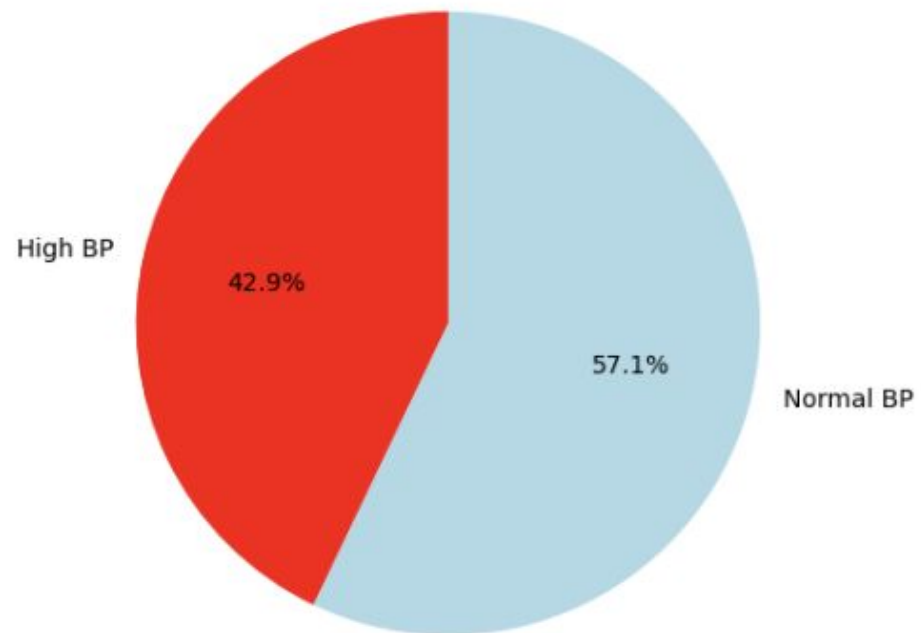


Pre/Diabetes Distribution Among Male Participants

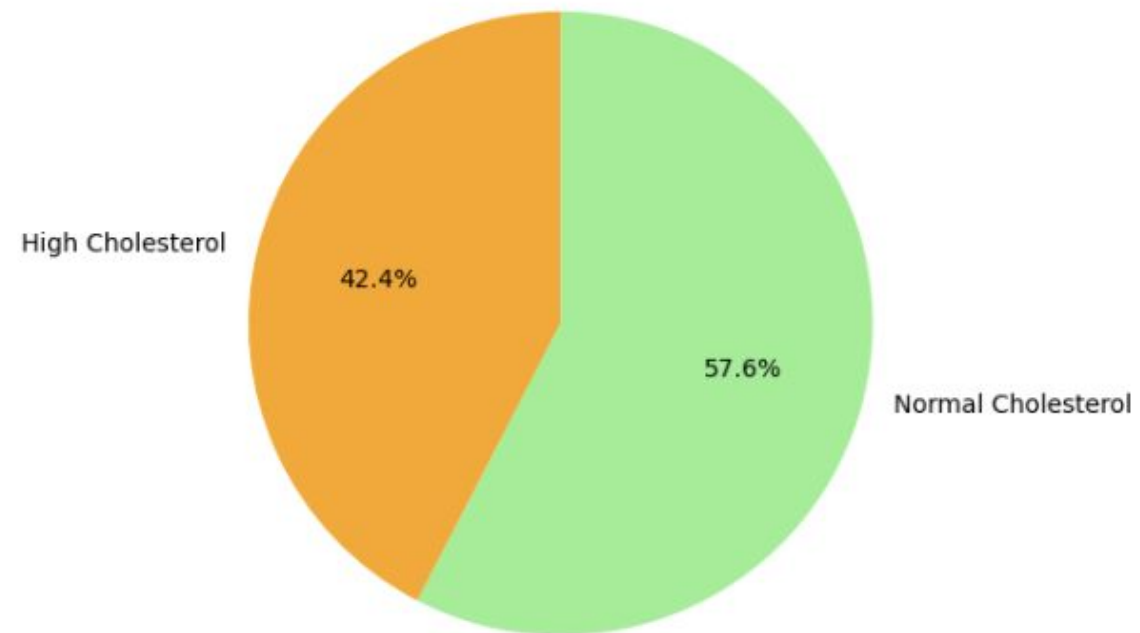
Males with Pre/Diabetes



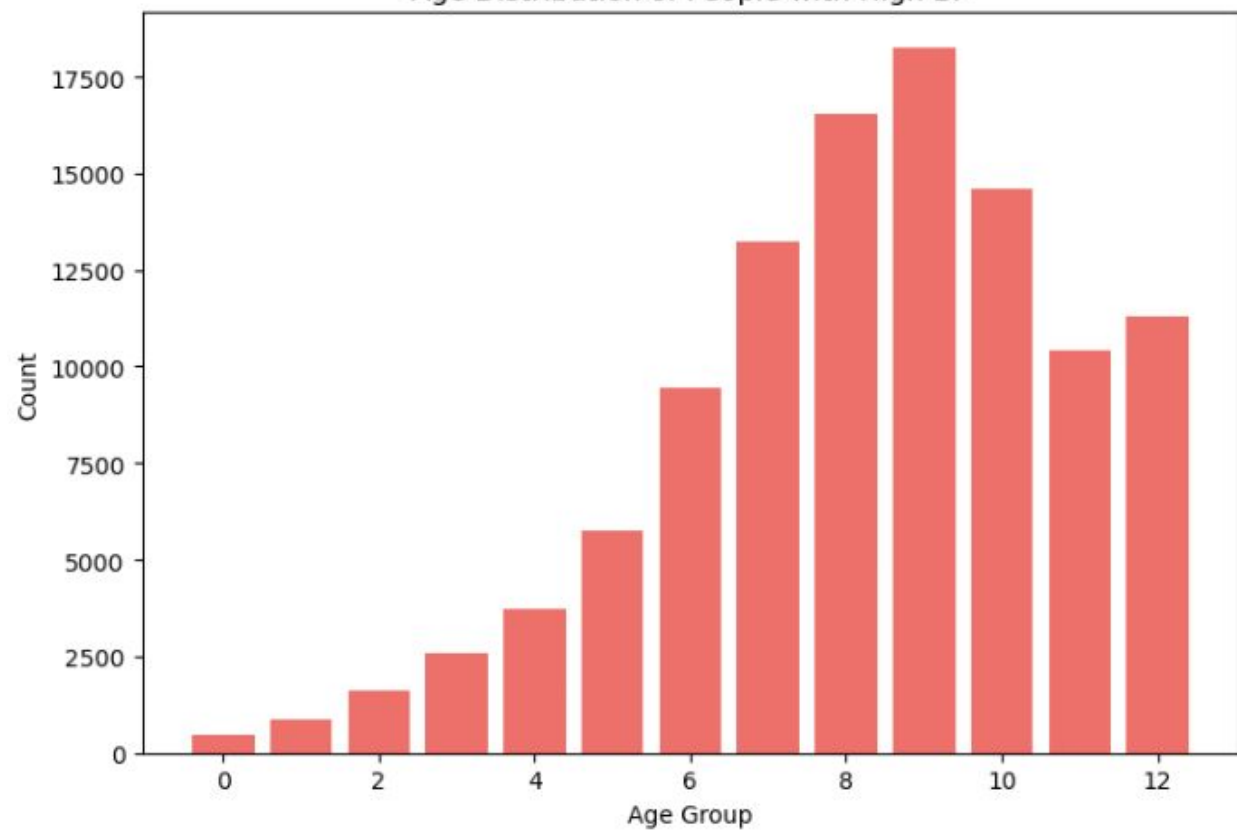
Distribution of High Blood Pressure in Population



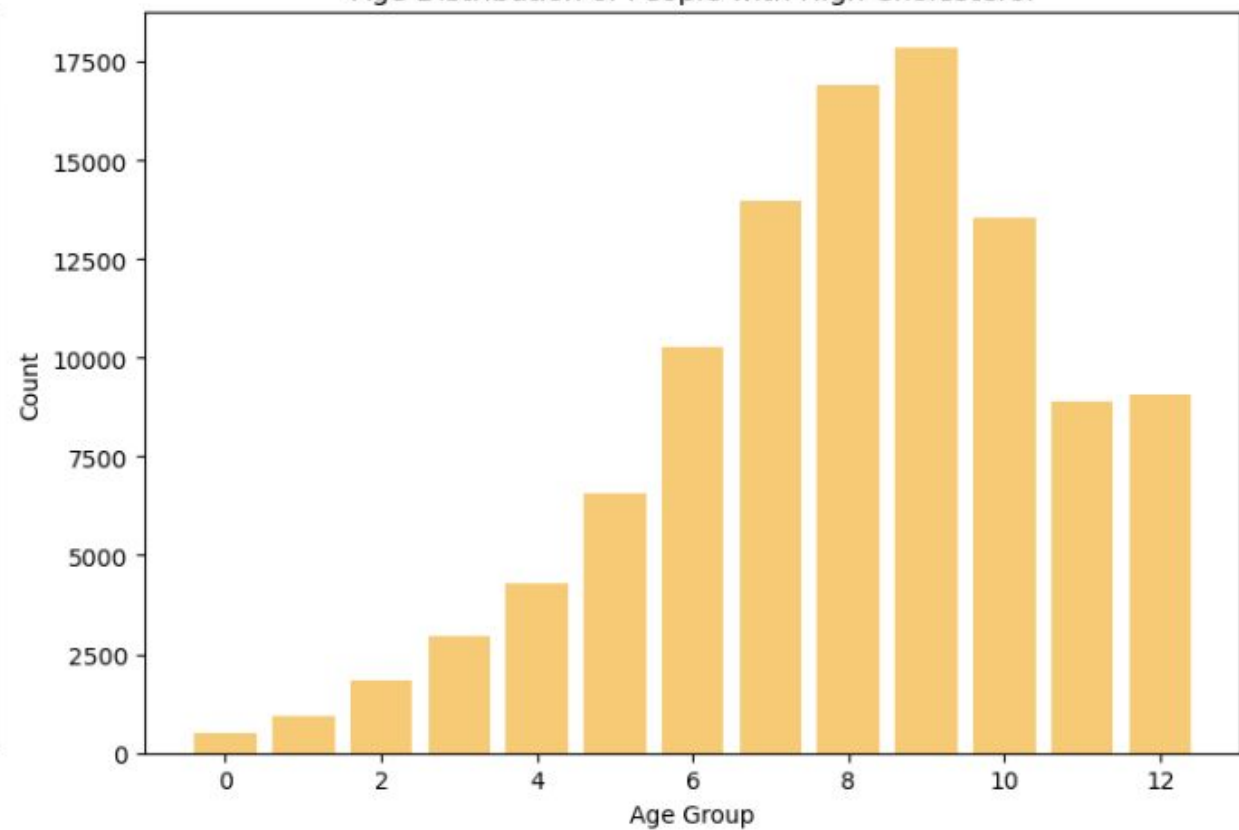
Distribution of High Cholesterol in Population



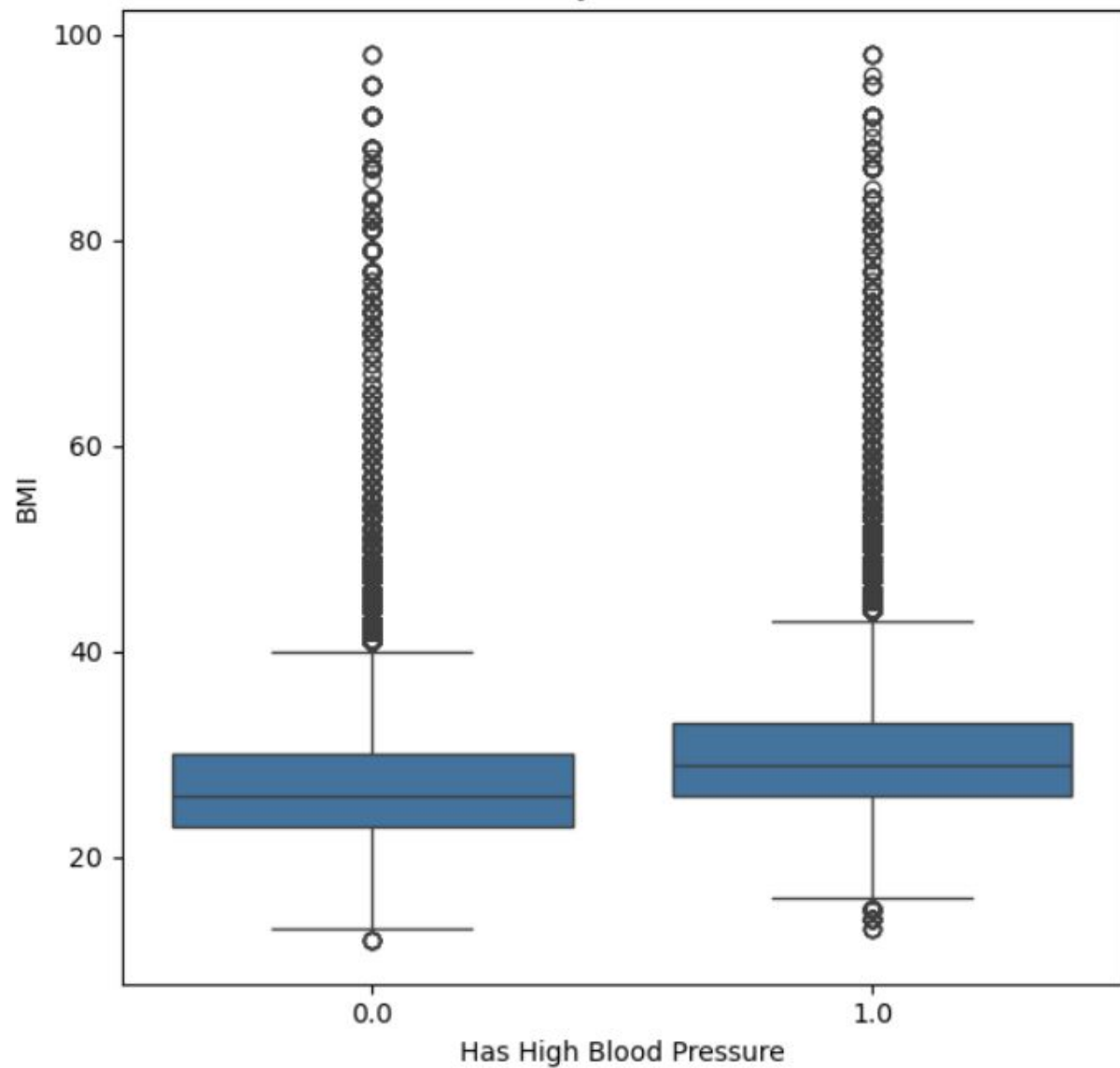
Age Distribution of People with High BP



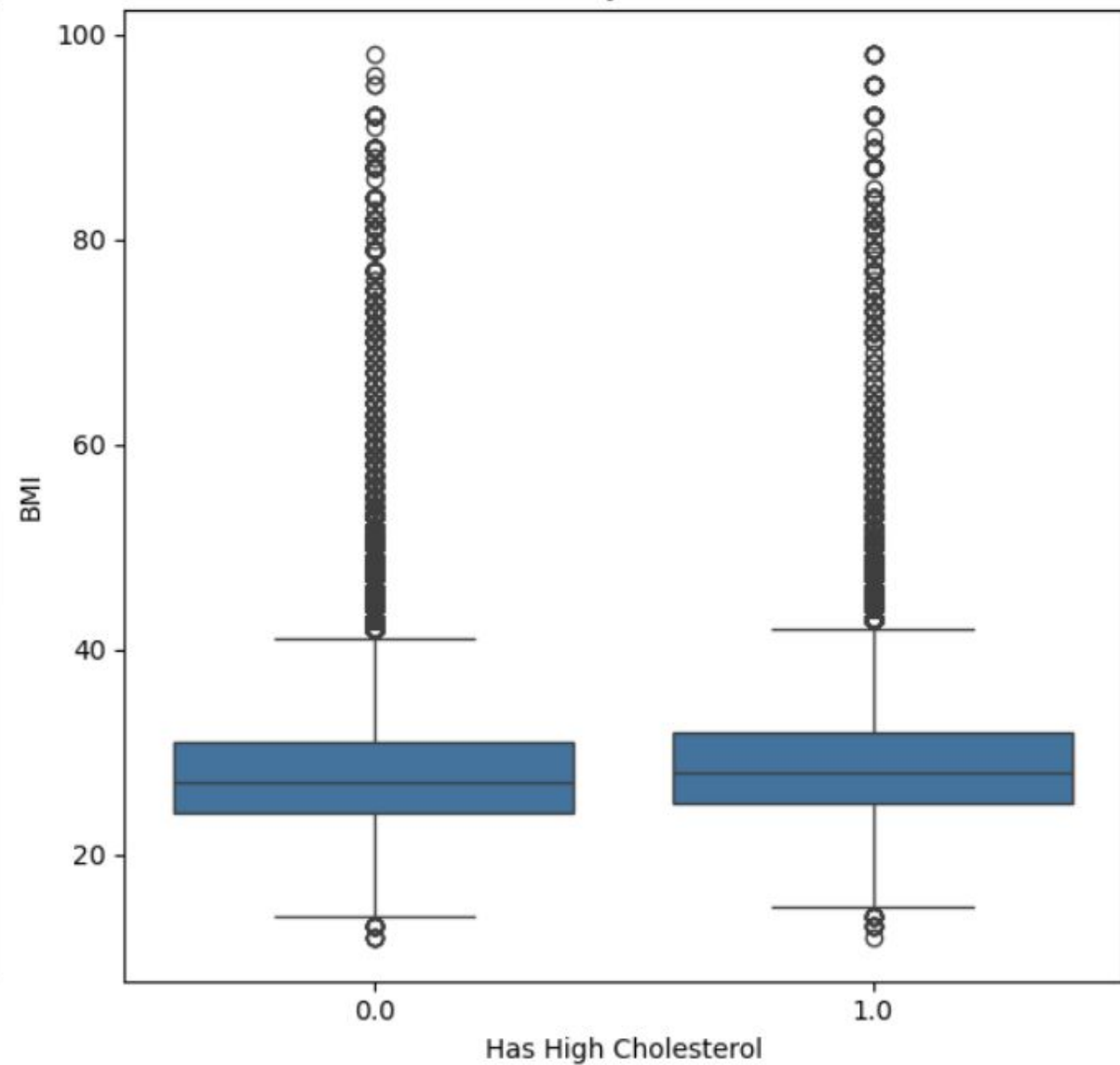
Age Distribution of People with High Cholesterol



BMI Distribution by Blood Pressure Status

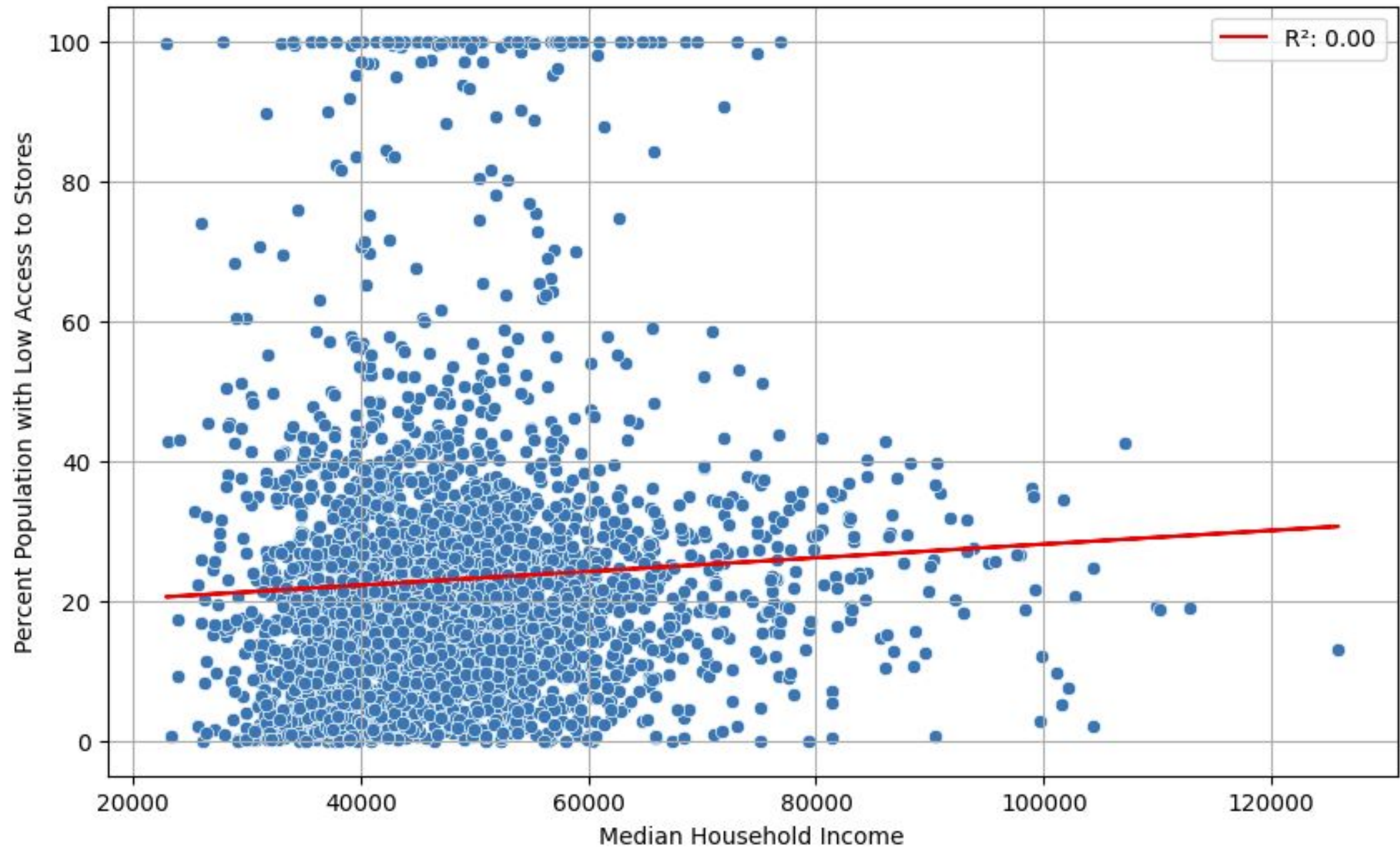


BMI Distribution by Cholesterol Status

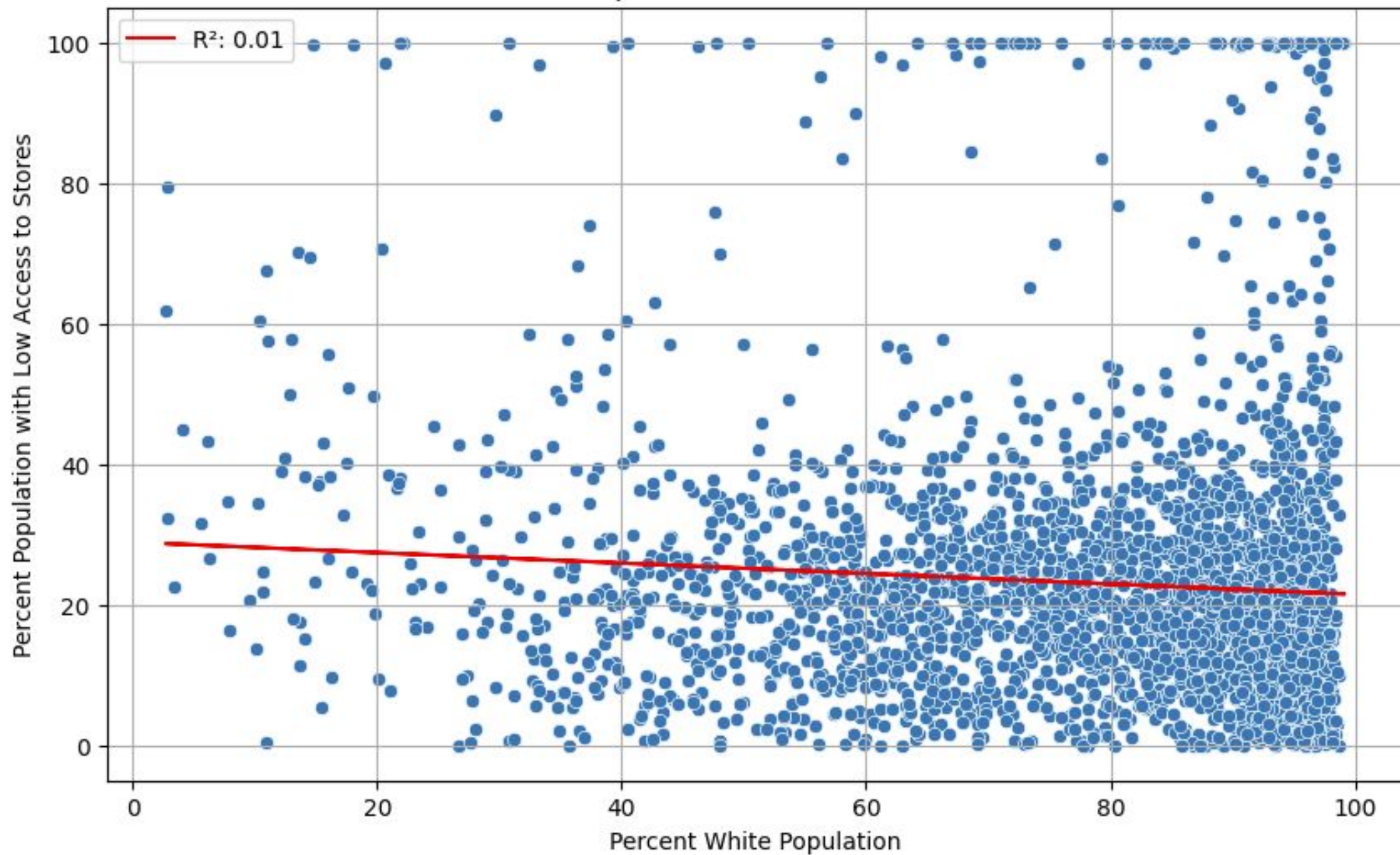


How does socioeconomic status affect access
to healthy food and healthcare?

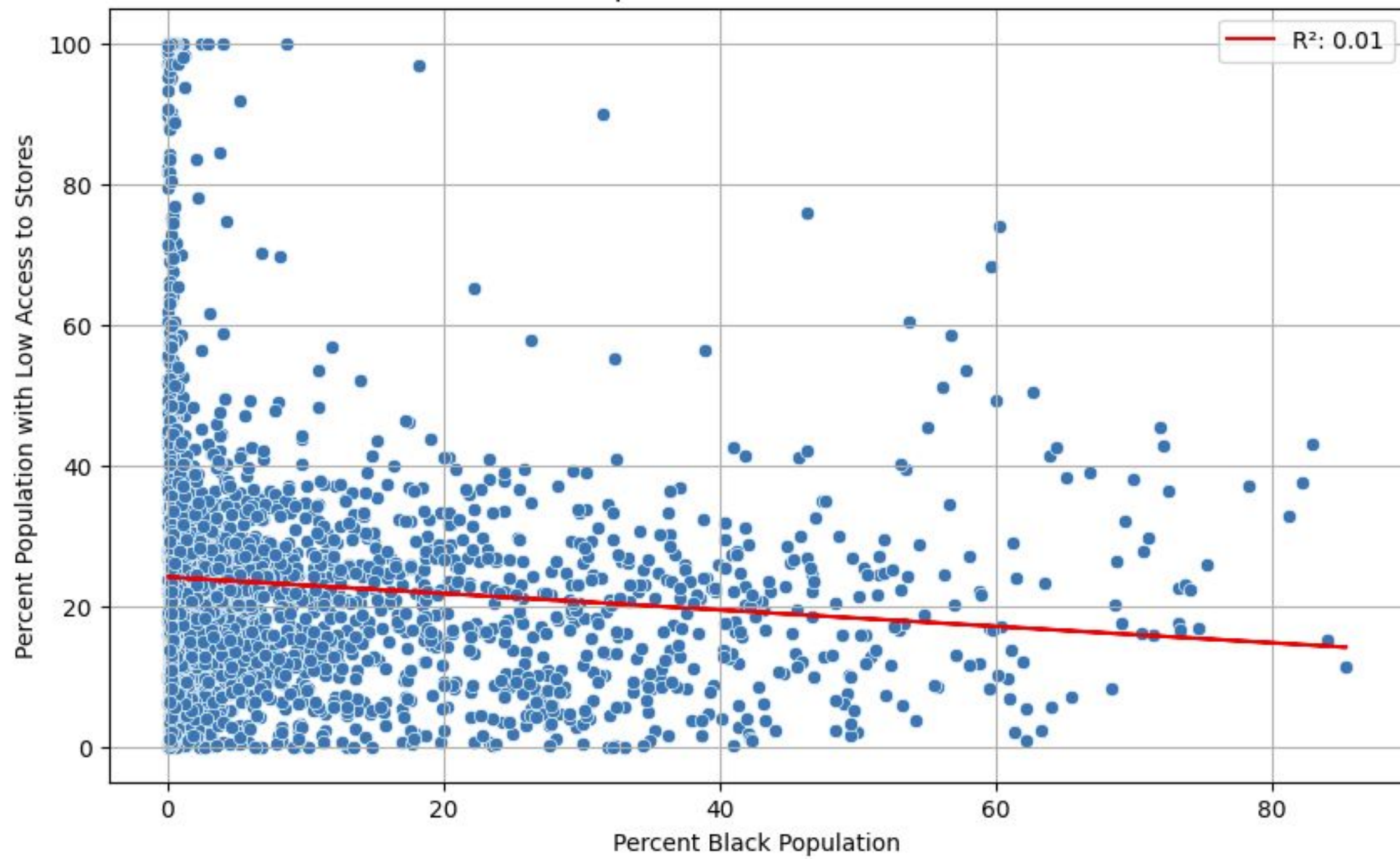
Income vs Low Store Access



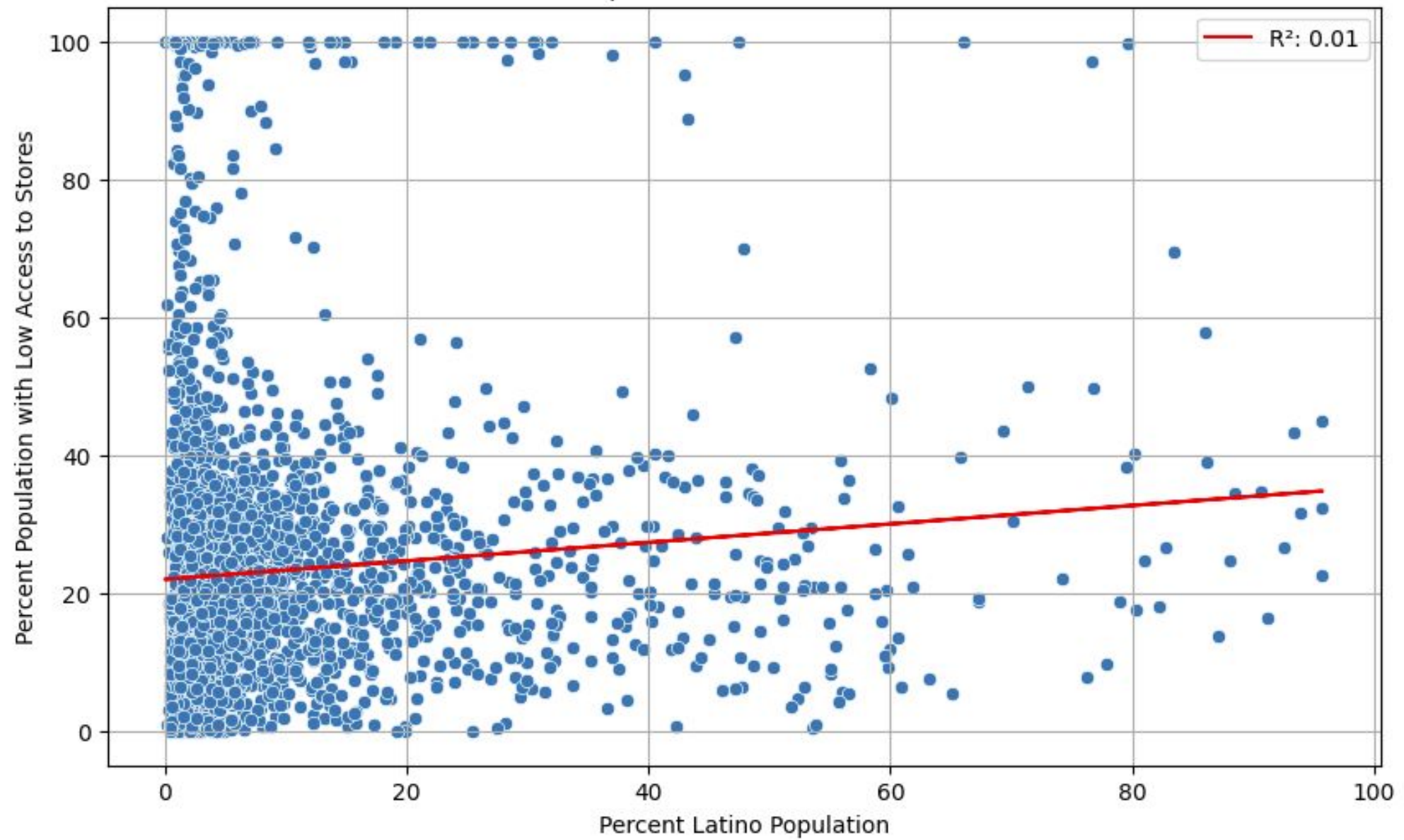
White Population vs Low Store Access



Black Population vs Low Store Access



Latino Population vs Low Store Access

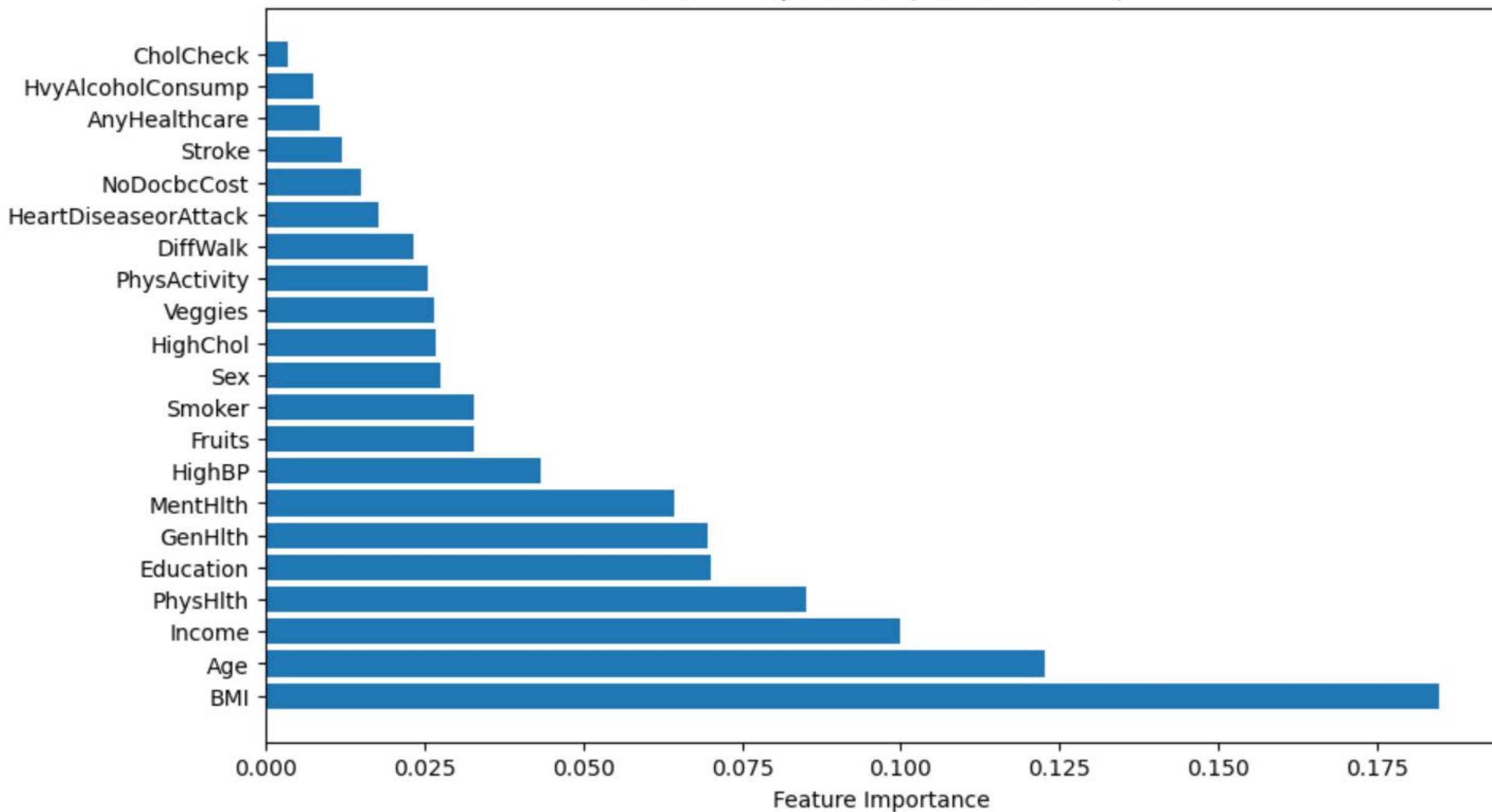


What features are most indicative of someone getting diabetes?

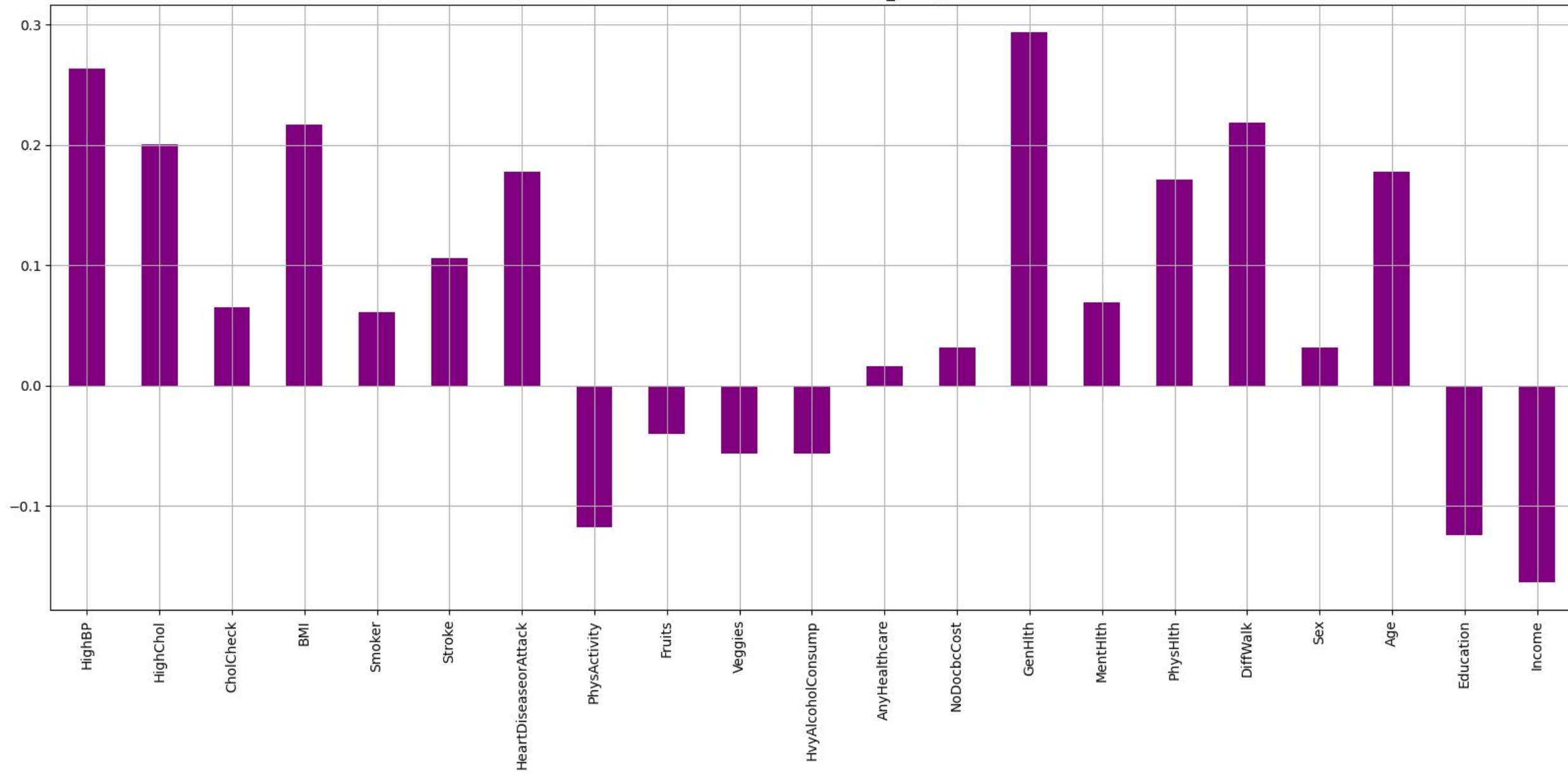
Below are is a comparison of the accuracy and recall for '1s'. Random Forest performed the best.

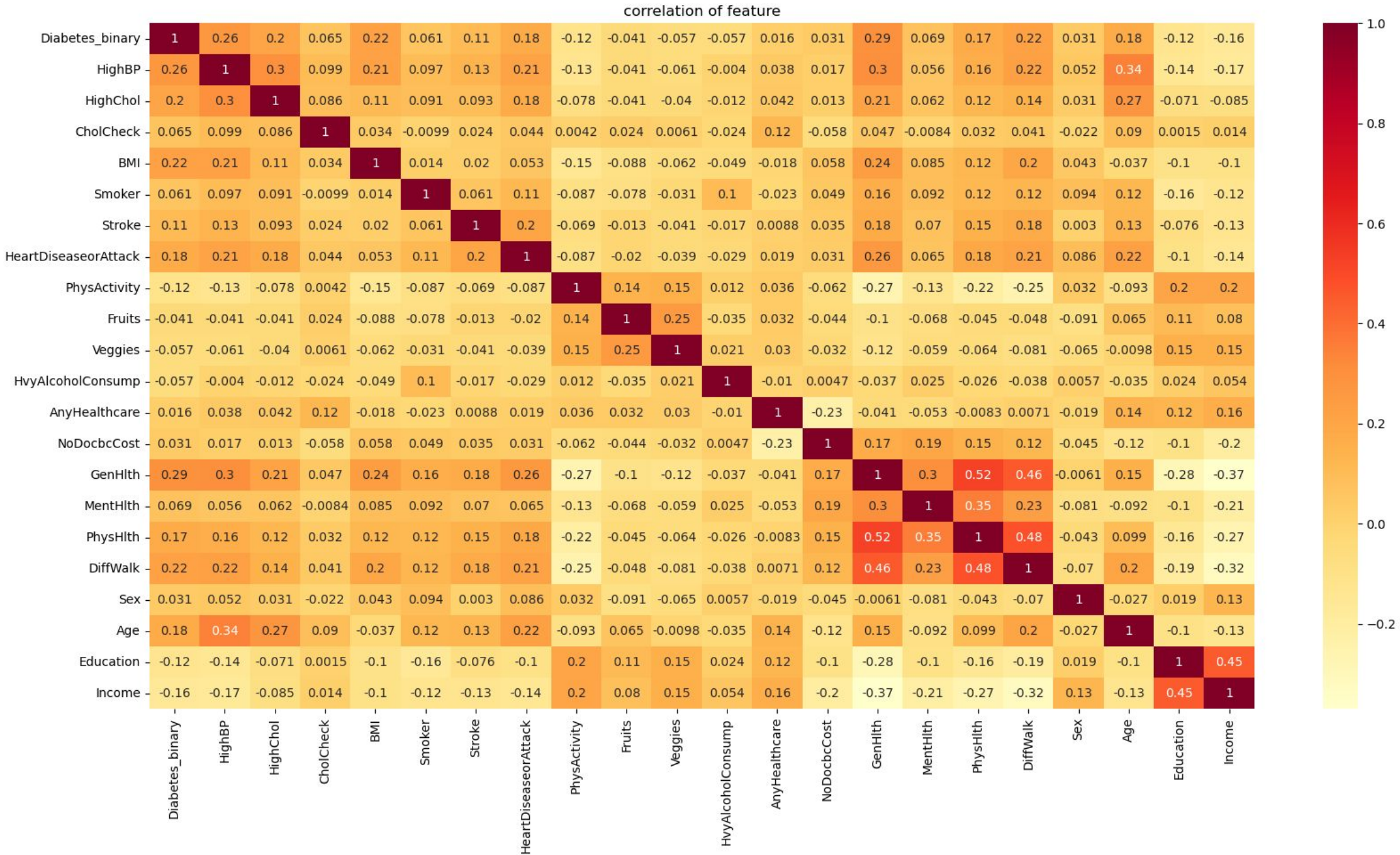
Model	Accuracy	Recall for 1s
Logistic Regression	0.74	0.76
Decision Trees	0.77	0.82
Random Forest	0.78	0.82

Feature Importance (Random Forest)



Correlation with Diabetes_binary





Select Best features using KBest method

```
#using SelectKBest class to extract top 10 best features
```

```
BestFeatures = SelectKBest(score_func=chi2, k=10)
```

```
fit = BestFeatures.fit(X_train_fi,y_train_fi)
```

```
df_scores = pd.DataFrame(fit.scores_)
```

```
df_columns = pd.DataFrame(X_train_fi.columns)
```

```
#concatenating two dataframes for better visualization
```

```
f_Scores = pd.concat([df_columns,df_scores],axis=1)
```

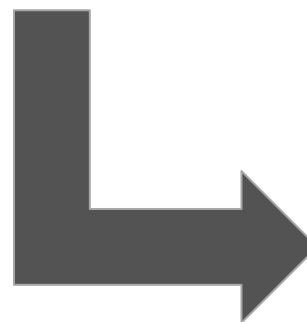
```
f_Scores.columns = ['Feature','Score']
```

```
f_Scores = f_Scores.sort_values('Score', ascending=False)
```

```
f_Scores
```

```
# feature scores
```

KBest was also used to see another method for determining features importance ranking. Results are sorted by highest chi score.



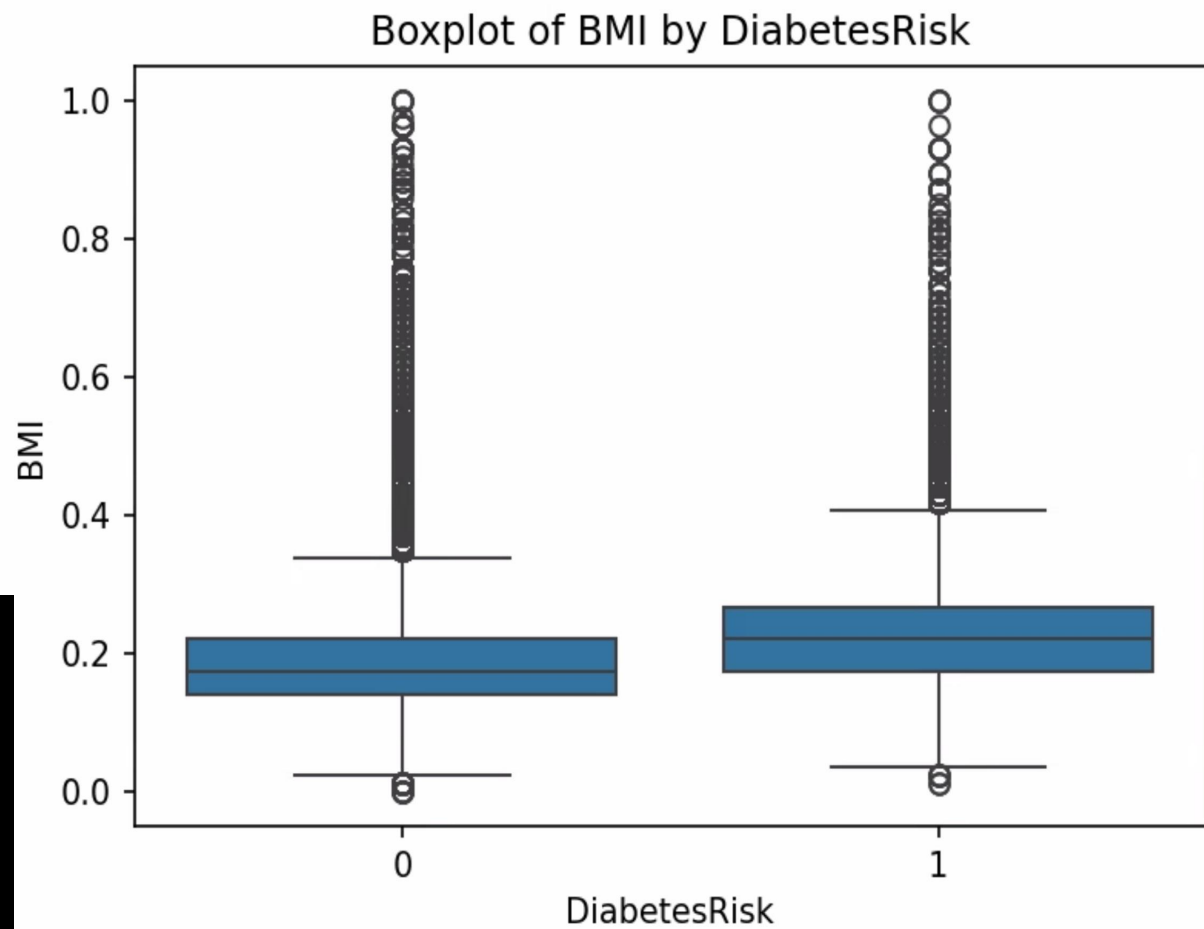
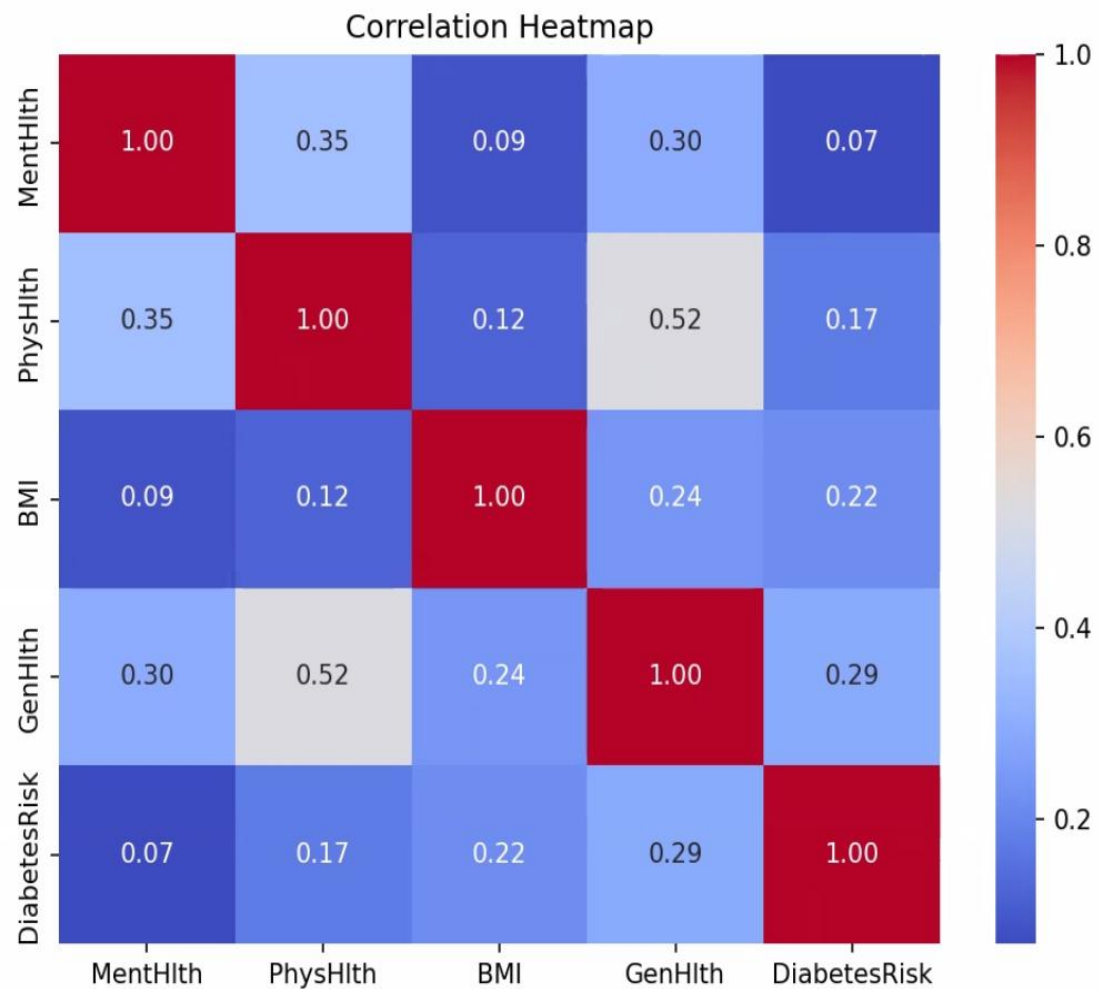
	Feature	Score
15	PhysHlth	133424.406534
14	MentHlth	21029.632228
3	BMI	18355.166400
16	DiffWalk	10059.506391
0	HighBP	10029.013935
13	GenHlth	9938.507776
18	Age	9276.141199
6	HeartDiseaseorAttack	7221.975378
1	HighChol	5859.710582
20	Income	4829.816361
5	Stroke	2725.225194
7	PhysActivity	861.887532
10	HvyAlcoholConsump	779.424807
19	Education	756.035496
4	Smoker	521.978858
12	NoDocbcCost	229.542412

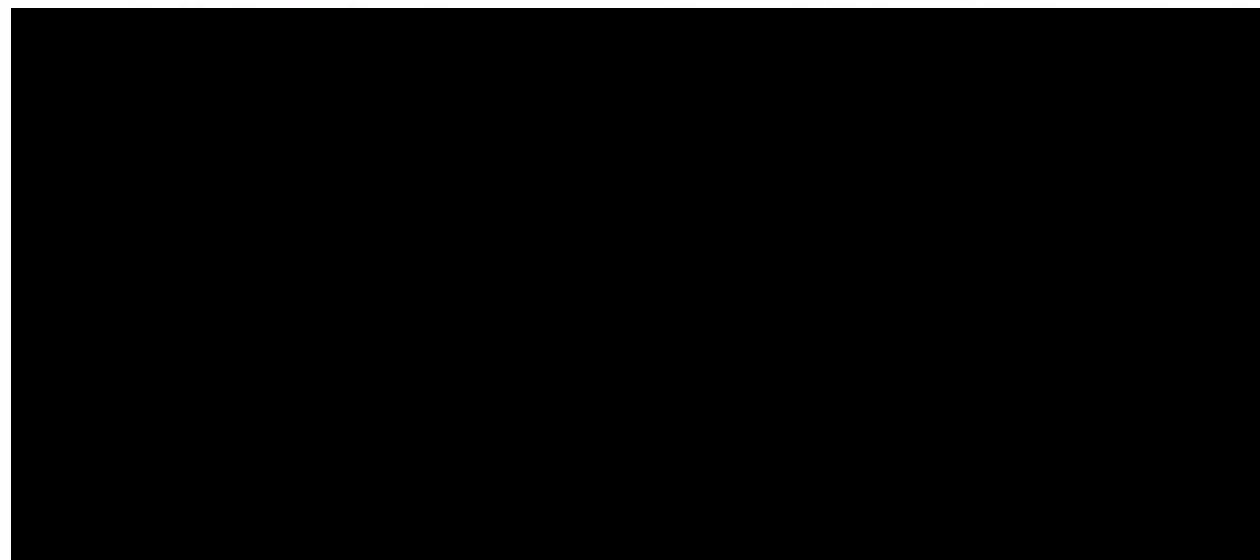
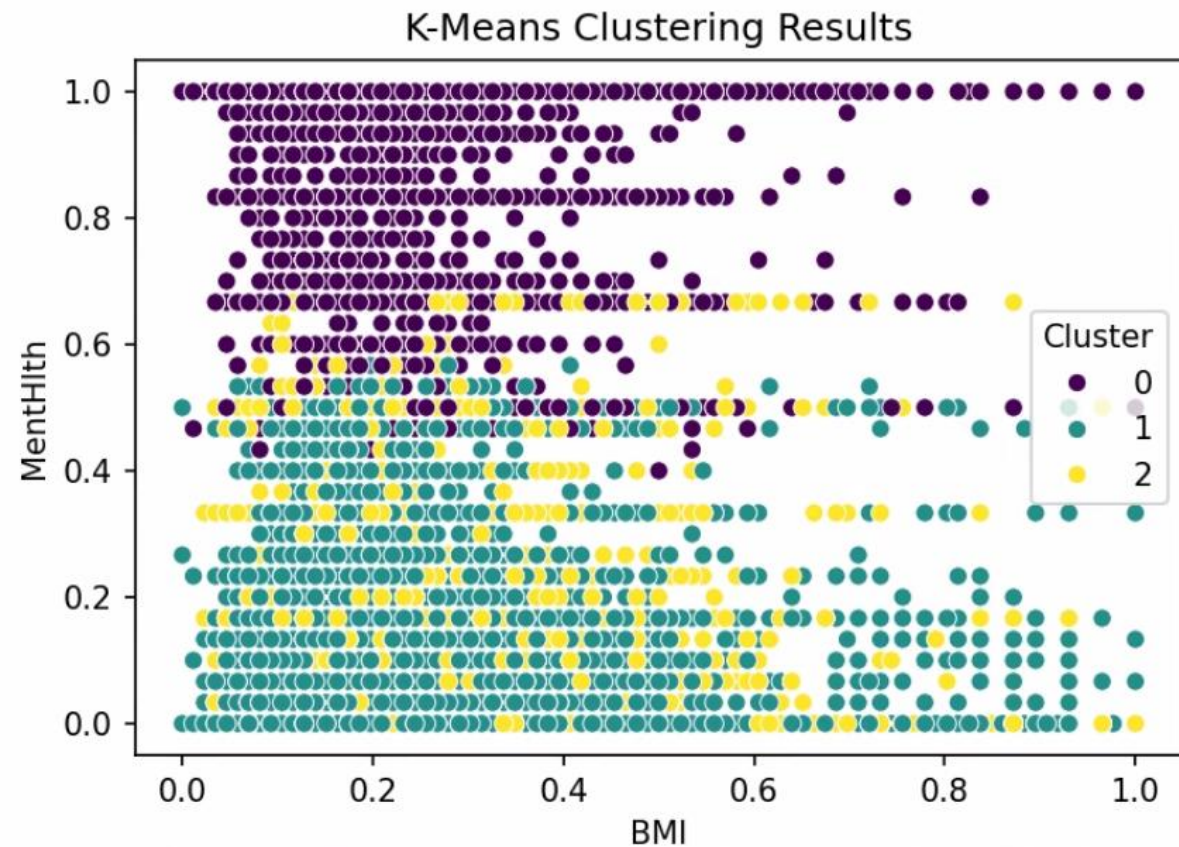
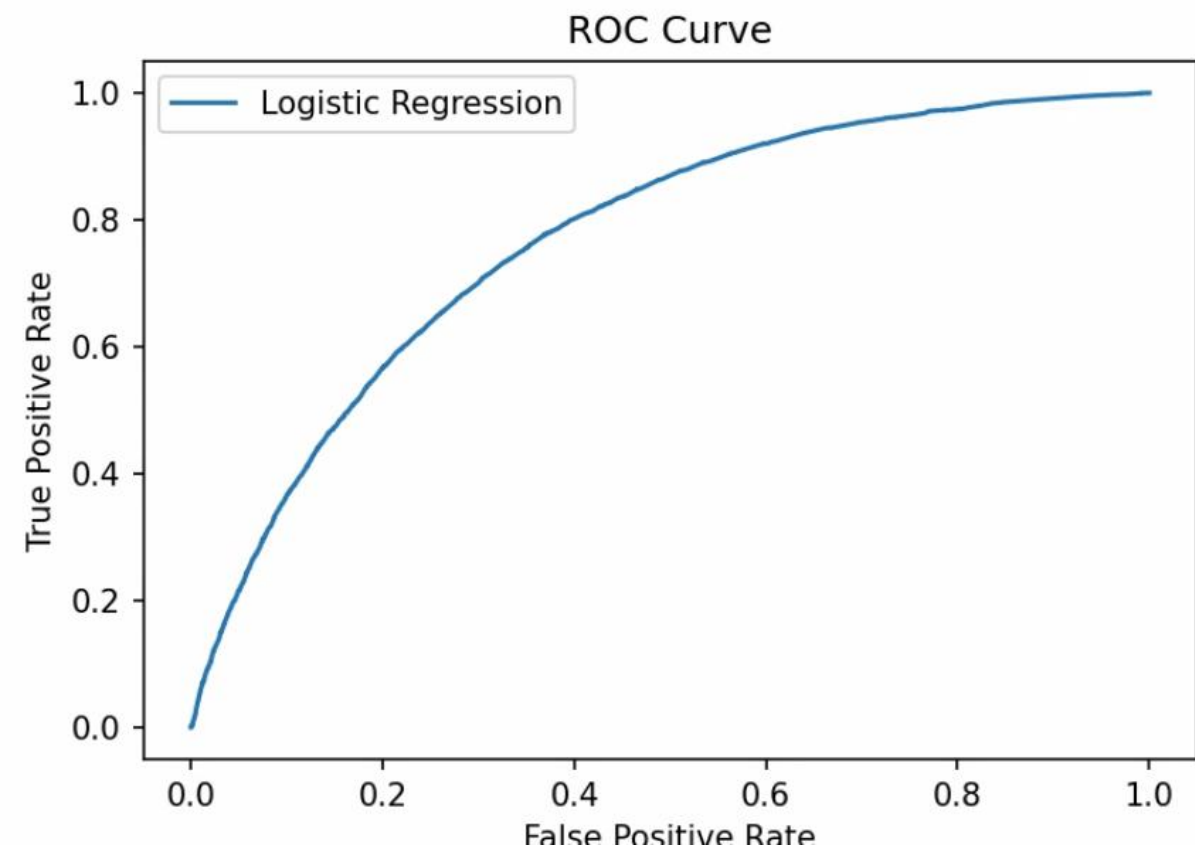
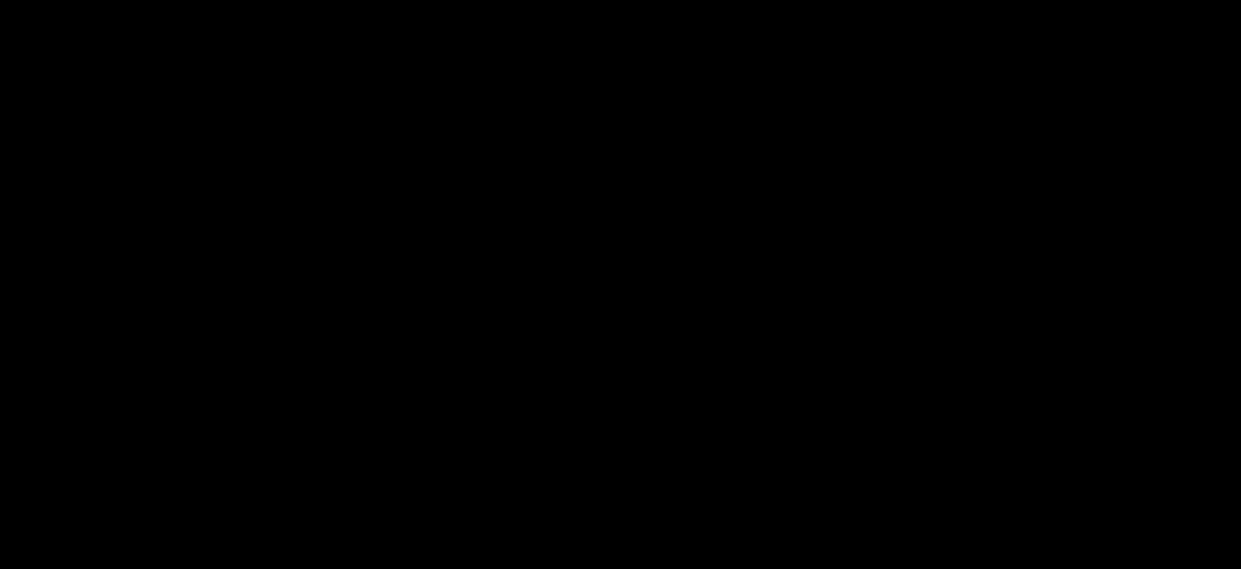
How does Mental Health and Insulin resistance strongly influence diabetes risk?

To explore these relationships

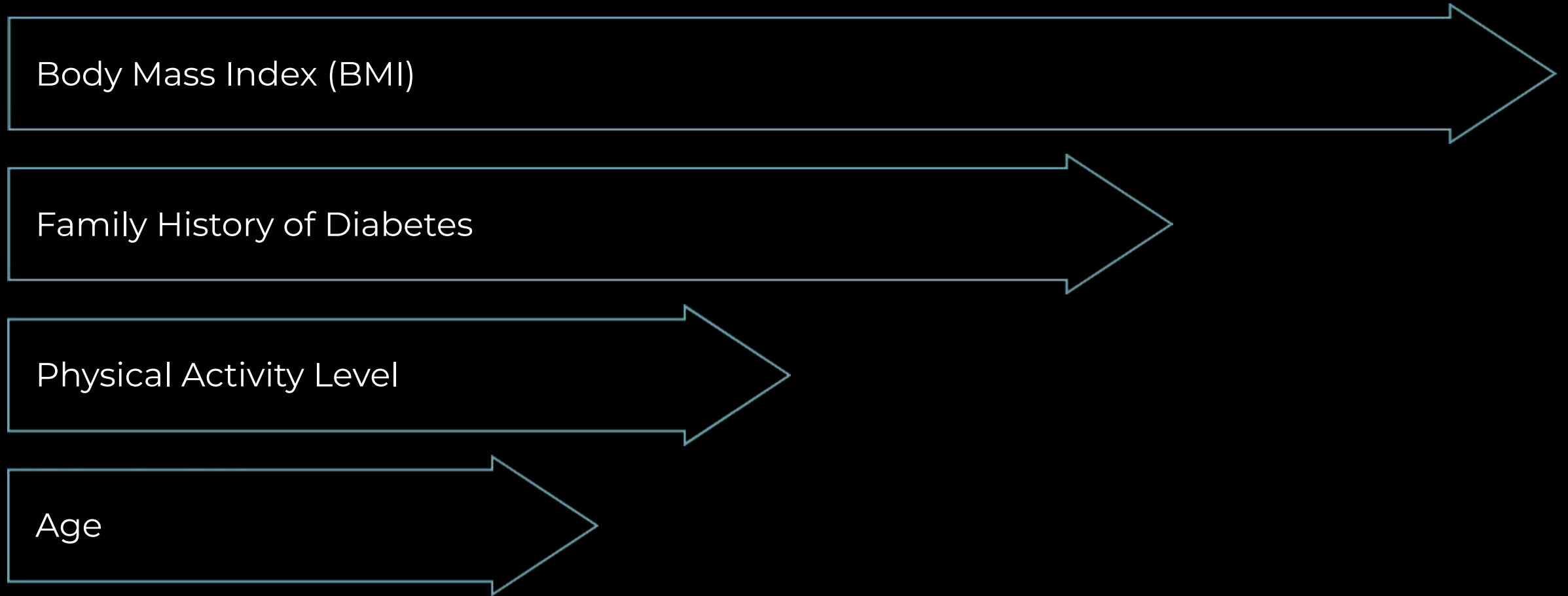
- Correlation analysis,
- Applied clustering methods
- Evaluated predictive performance using logistic regression.







Feature Importance

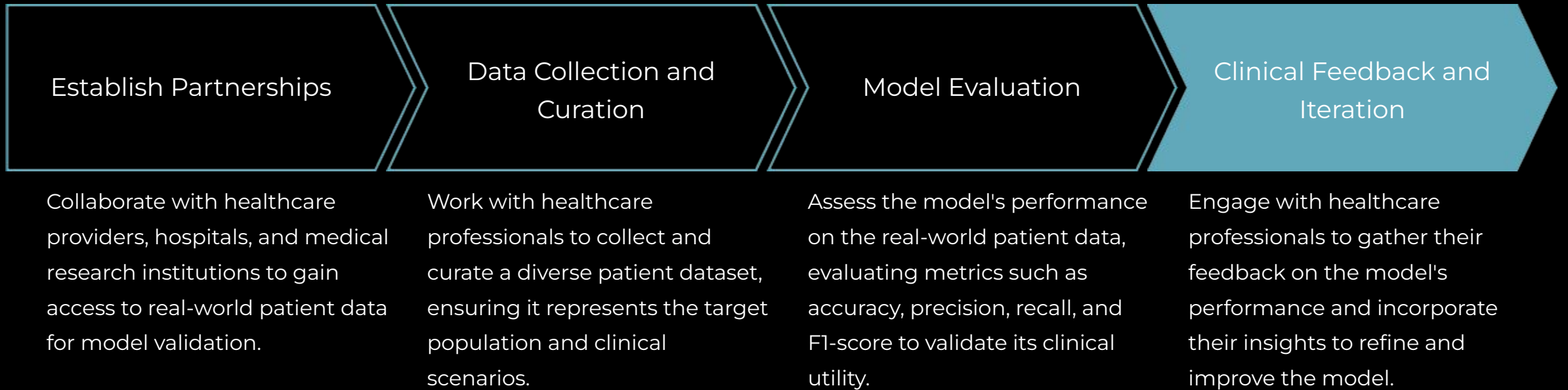


Can we use a subset of the risk factors to accurately predict whether an individual has diabetes?

The results of the logistic regression model shows that we can correlate BMI with diabetes with a sufficient degree of accuracy

- The random forest results also indicates that BMI, income and age are important factors
- Therefore creating a predictive model is possible and is a good avenue for future work.

Clinical Validation



Future Considerations

