

Diabetes Classification Model



Snizhana Kurylyuk

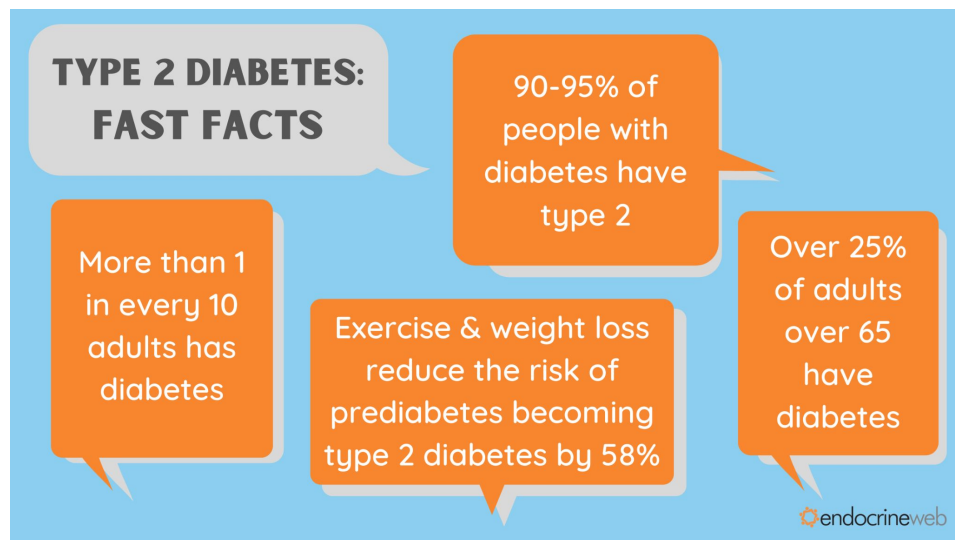
Type 2 Diabetes

An estimated 34.2 million people have diabetes (10.5 percent of the U.S. population).

- An estimated 7.3 million adults ages 18 years or older have diabetes but are undiagnosed (21.4 percent of adults with diabetes).

Data: 250K Data points obtained from the Behavioral Risk Factor Surveillance System (BRFSS) a health-related telephone survey that is collected annually by the CDC.

Target: Build a classification model that can diagnose whether or not an individual has diabetes.



Features



- High Blood Pressure
- High Cholesterol
- Cholesterol Check
- Body Mass Index
- Smoker
- Stroke
- Heart Disease or Attack
- Physical Activity
- Fruits
- Veggies
- Heavy Alcohol Consumption
- Any Healthcare
- No Doctor because of cost
- General Health
- Mental Health
- Physical Health
- Difficulty Walking
- Sex
- Age
- Education
- Income

Can we predict if someone has Diabetes?

1

Step 1

- EDA & Data Engineering

2

Step 2

- Baseline & Feature Engineering

3

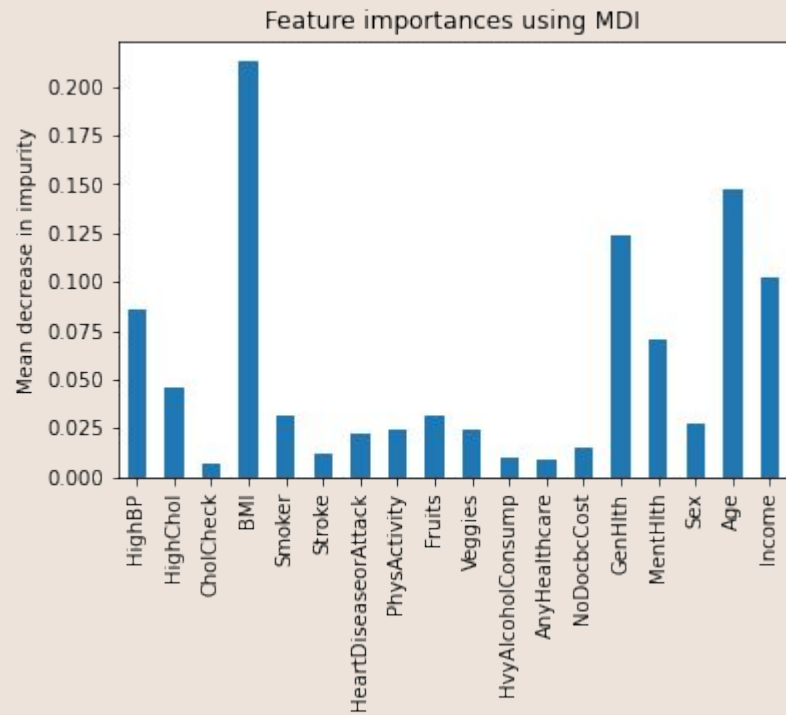
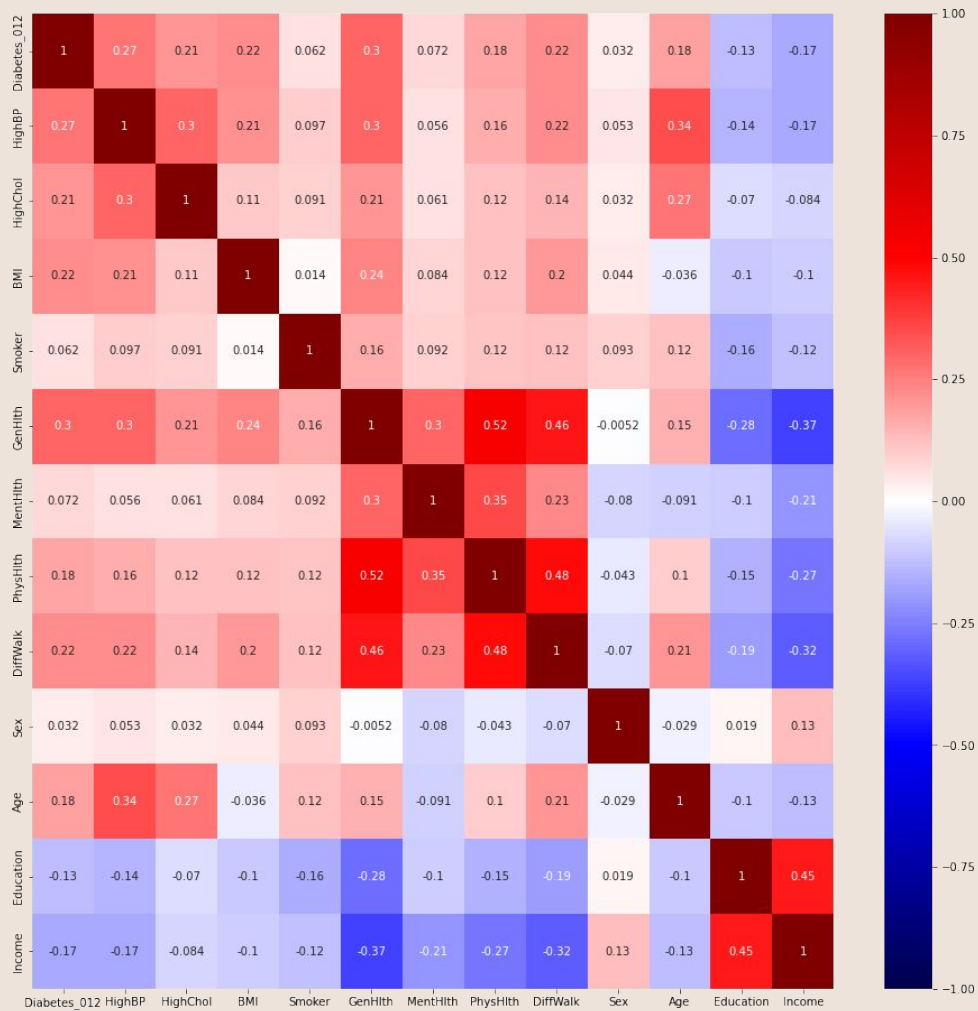
Step 3

- Model Comparison

4

Step 4

- Best Performance Model Selection



Classification Algorithms Validation Metrics

K Nearest Neighbors

- Accuracy: 0.7673
- Precision: 0.3128
- Recall: 0.5360
- F1: 0.3951
- Cross Validation Accuracy Score : 0.77854

Random Forest:

- Accuracy: 0.7698
- Precision: 0.3488
- Recall: 0.7195
- F1: 0.4698
- Cross Validation Accuracy Score: 0.80256

Logistic Regression:

- Accuracy: 0.7896
- Precision: 0.3639
- Recall: 0.6476
- F1: 0.4659
- Cross Validation Accuracy Score : 0.74455

Stacking Ensembling:

- Accuracy: 0.8623
- Precision: 0.5388
- Recall: 0.1988
- F1: 0.2905
- Cross Validation Accuracy Score: 0.92498

Decision Tree:

- Accuracy: 0.7699
- Precision: 0.3330
- Recall: 0.6218
- F1: 0.4337
- Cross Validation Accuracy Score: 0.71748

Voting Ensembling:

- Accuracy: 0.7974
- Precision: 0.3560
- Recall: 0.5296
- F1: 0.4258
- Cross Validation Accuracy Score: 0.8623

Hyperparameter Tuning

First Model

Accuracy: 0.7698
Precision: 0.3488
Recall: 0.7195
F1: 0.4698
Cross Validation Accuracy Score: 0.80256

Second Model

Accuracy: 0.8201
Precision: 0.3440
Recall: 0.2969
F1: 0.3187
Cross Validation Accuracy Score: 0.92541

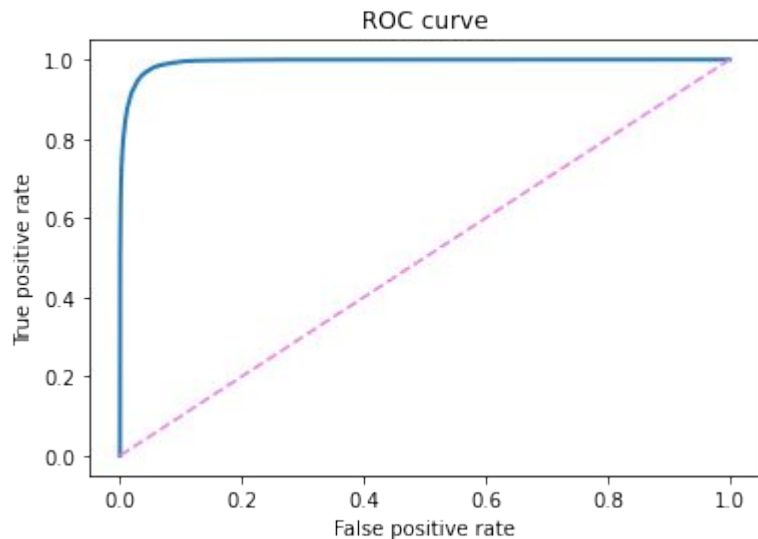
Third Model

Accuracy: 0.8312
Precision: 0.3881
Recall: 0.3310
F1: 0.3572
Cross Validation Accuracy Score: 0.91561

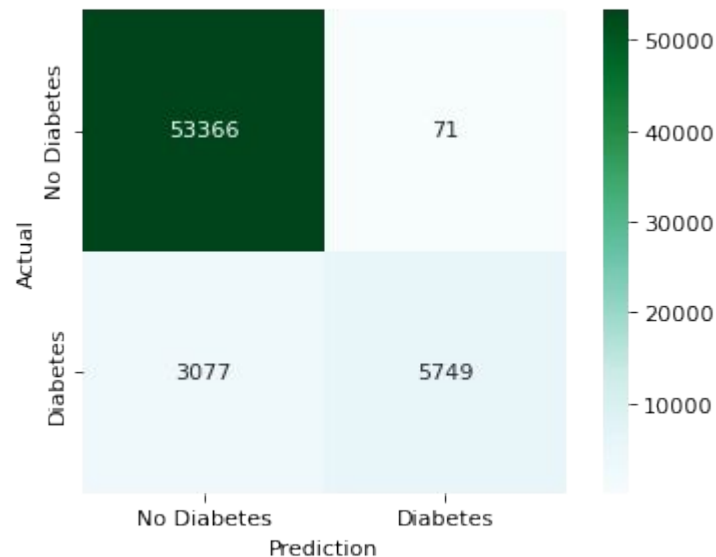
Best Model

Accuracy: 0.8295
Precision: 0.3963
Recall: 0.3869
F1: 0.3915
Cross Validation Accuracy Score: 0.91611

Best Performance Model: Random Forest



ROC AUC Score = 0.99427



Questions

