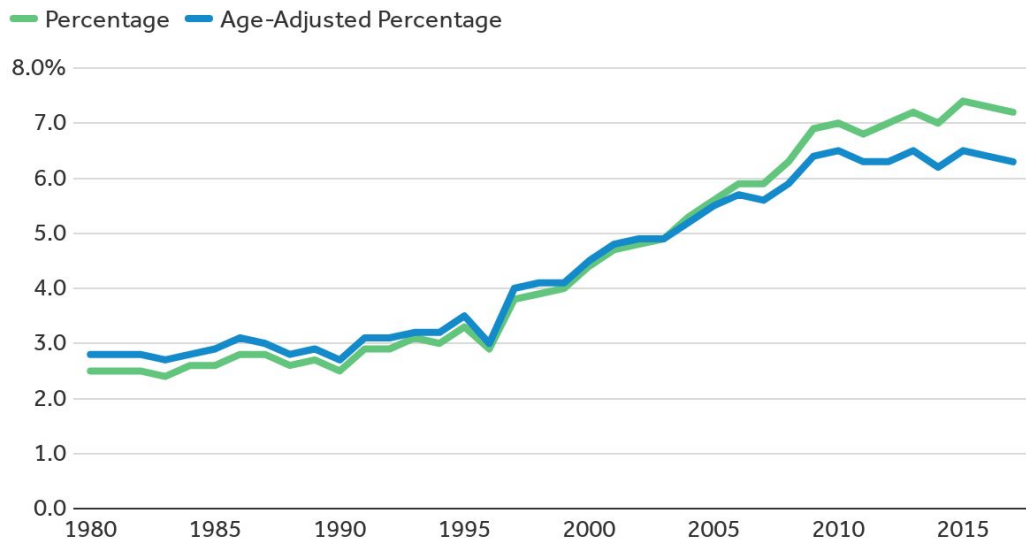# Diabetes Prediction using Medical History

Group 11 - Yunrui Jiang, Jonah Soong, Justin Sin, Anyin Huang, Zeyu Chang

# Motivation

- Diabetes
  a. Chronic condition leading to life threatening issues
  b. Caused by increased resistance to insulin leading to high glucose levels
- 38.4 million people (11.6% of the US)
- 97.6 million adults have prediabetes (38.0% of US adults)

### Share of total population with diagnosed diabetes, 1980-2017

— Percentage  — Age-Adjusted Percentage

Source: US Diabetes Surveillance System

Peterson-KFF
**Health System Tracker**

# Objective

1.  What are some key indicators for diabetes?

2.  How well do existing prediction techniques forecast diabetes?

# Methodology

- Pearson Correlation Coefficients were computed on different subsets of data as main heuristic for strong indicators
    - Categorical variables were binarized
    - Data was stratified to see if coefficients changed

- 5 different prediction ML models were tested for accuracy and precision

# Datasets overview

(https://www.kaggle.com/datasets/iammustafatz/diabetes-prediction-dataset/data)

| | gender | age | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 80.00 | 0 | 1 | never | 25.19 | 6.60 | 140 | 0 |
| 1 | Female | 54.00 | 0 | 0 | No Info | 27.32 | 6.60 | 80 | 0 |
| 2 | Male | 28.00 | 0 | 0 | never | 27.32 | 5.70 | 158 | 0 |
| 3 | Female | 36.00 | 0 | 0 | current | 23.45 | 5.00 | 155 | 0 |
| 4 | Male | 76.00 | 1 | 1 | current | 20.14 | 4.80 | 155 | 0 |

# Dataset attributes

Age

Gender

Diabetes

Heart Disease

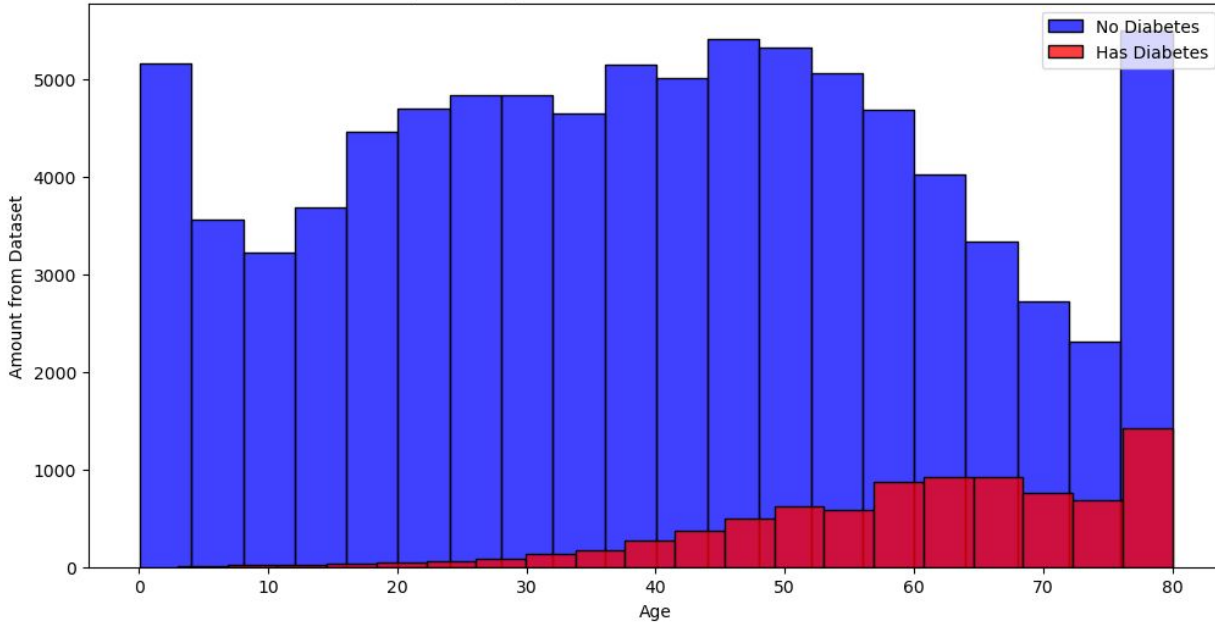BMI (Body mass index)

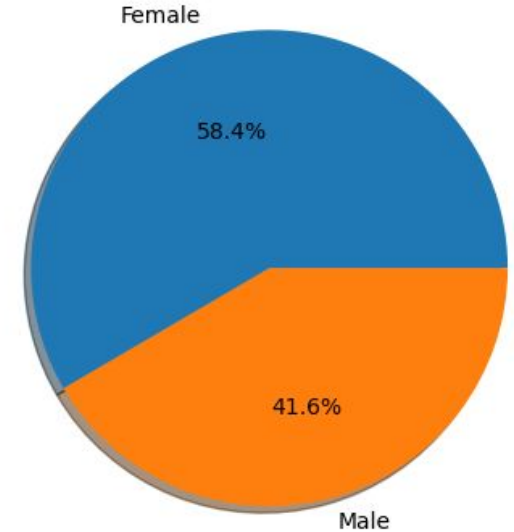Hypertension (high blood pressure)

Smoking history

Blood glucose level

HbA1c level (average blood glucose level over the past 2-3 months)

# Age and Gender

# Diabetes and Heart Disease



Distribution of diabetes out of 96146 people

HAS 8.8%

Does NOT have 91.2%

Distribution of heart_disease out of 96146 people

HAS 4.1%

Does NOT have 95.9%

8

# BMI

## Body Mass Index - kg/m^2

https://www.cdc.gov/healthyweight/assessing/bmi/adult_bmi/index.html

# Hypertension

High Blood Pressure



HAS 7.8%

92.2%

Does NOT have

# Smoking History

# Blood Glucose

## glucose in bloodstream - mg/dl

https://www.cdc.gov/diabetes/basics/getting-tested.html

# HbA1c

fraction of hemoglobin(blood protein) with glucose attached - %

https://www.cdc.gov/diabetes/basics/getting-tested.html

# What are some key indicators for diabetes?

# How does stratification affect indicators of diabetes?

# Diabetes Correlation Matrix by Gender

- Very similar correlation matrix with non stratified data

# Diabetes Correlation Matrix by Age

- Stronger correlation with blood glucose and HbA1c in older people

- Younger ages don't have strong indicators

Diabetes Correlation Matrix Stratified by Age and Gender

# How well do standard ML models predict diabetes?

# Predictive Modeling

## 5 models are used

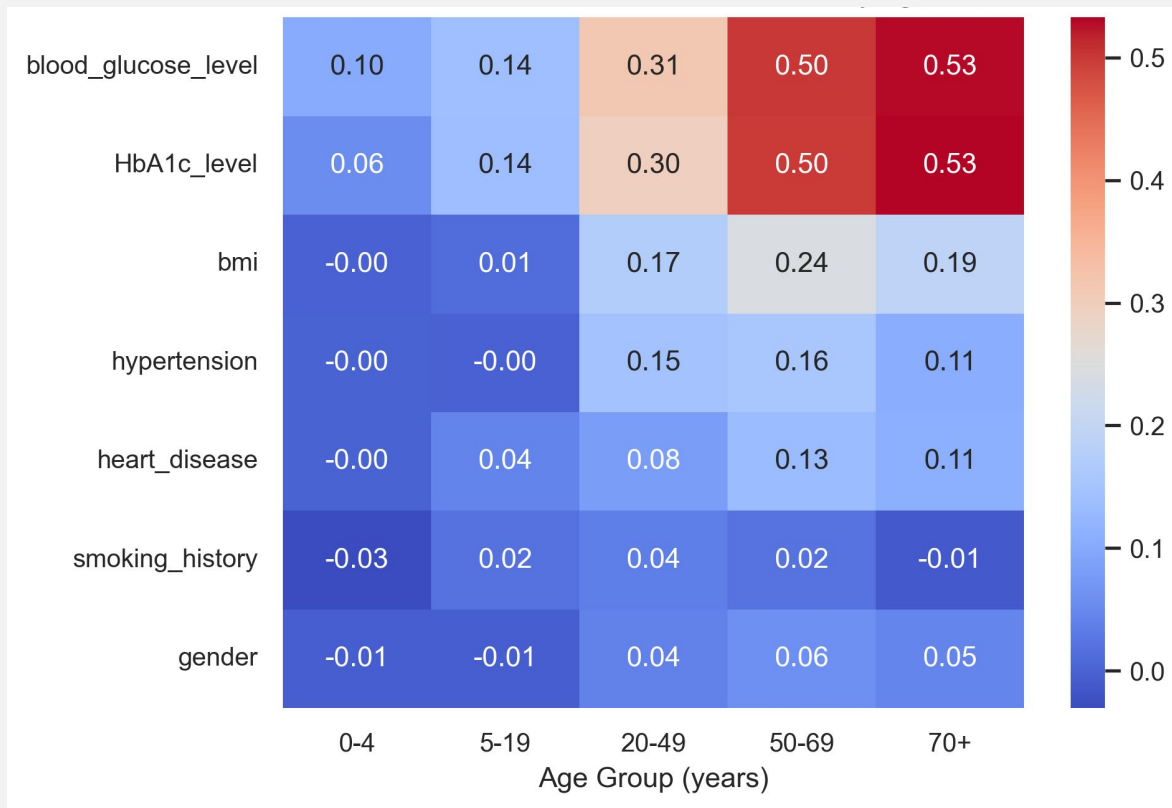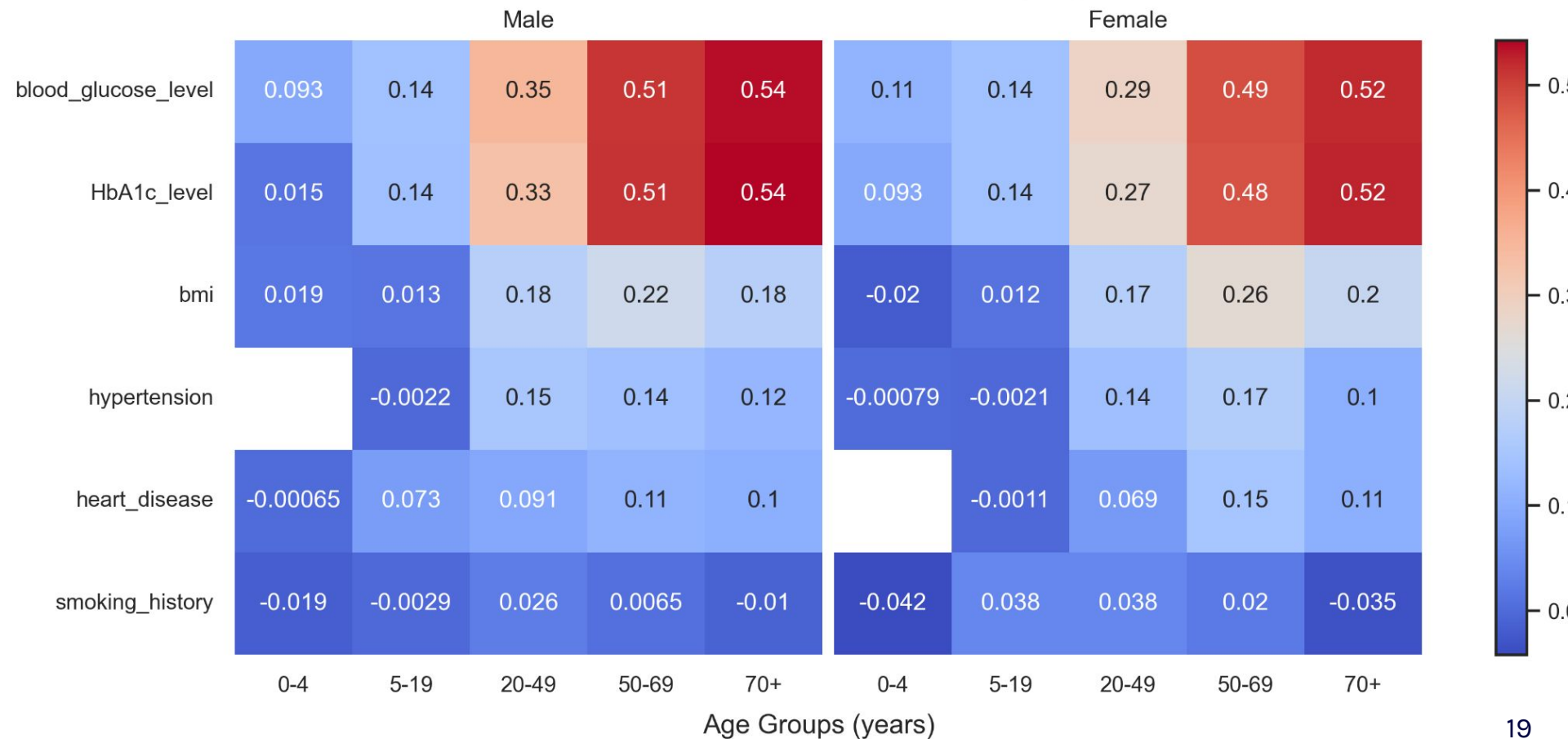- **Logistic Regression** - Predicts the probability of an event's occurrence using a logistic function.
- **Random Forest** - Builds multiple decision trees and combines their predictions to improve accuracy and generalizability.
- **Support Vector Machine (SVM)** - Finds the optimal boundary between classes to maximize the margin between categories.
- **K-Nearest Neighbors (KNN)** - Predicts by finding the K nearest training samples to the test data point.
- **Gradient Boosting** - Incrementally adds weak predictive models (usually decision trees) to minimize the loss function.

|        |          | Predicted | |
|--------|----------|----------------|----------------|
|        |          | Negative | Positive |
| Actual | Negative | True Negative | False Negative |
|        | Positive | False Positive | True Positive |

$$Accuracy = \frac{Number\ of\ Correct\ Prediction}{Number\ of\ the\ Results}$$

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

$$F1 = \frac{Precision * Recall}{Precision + Recall}$$

21

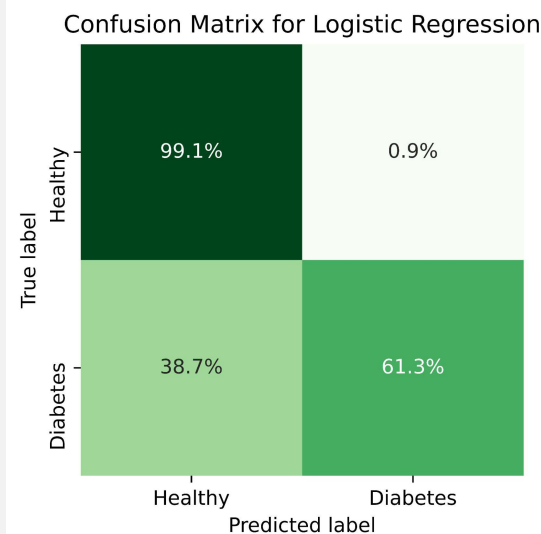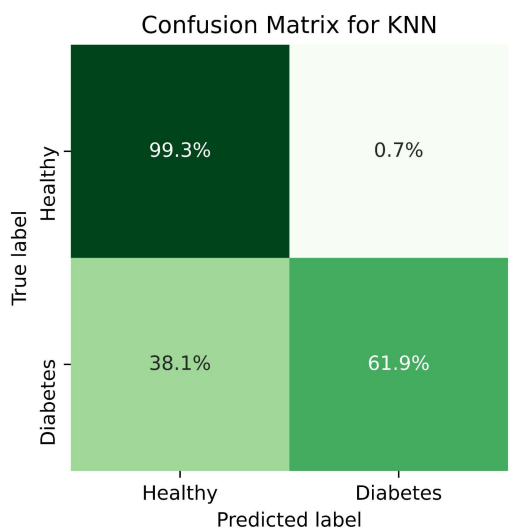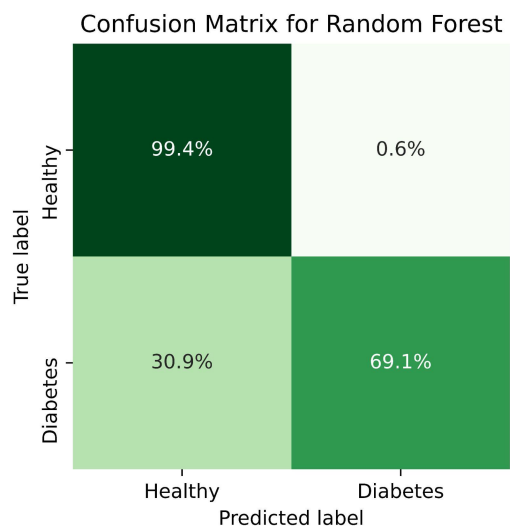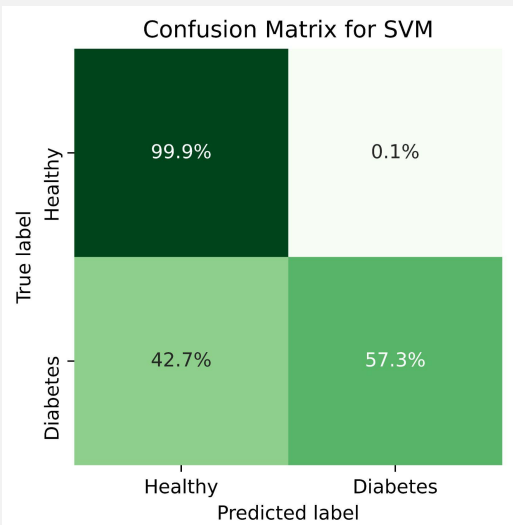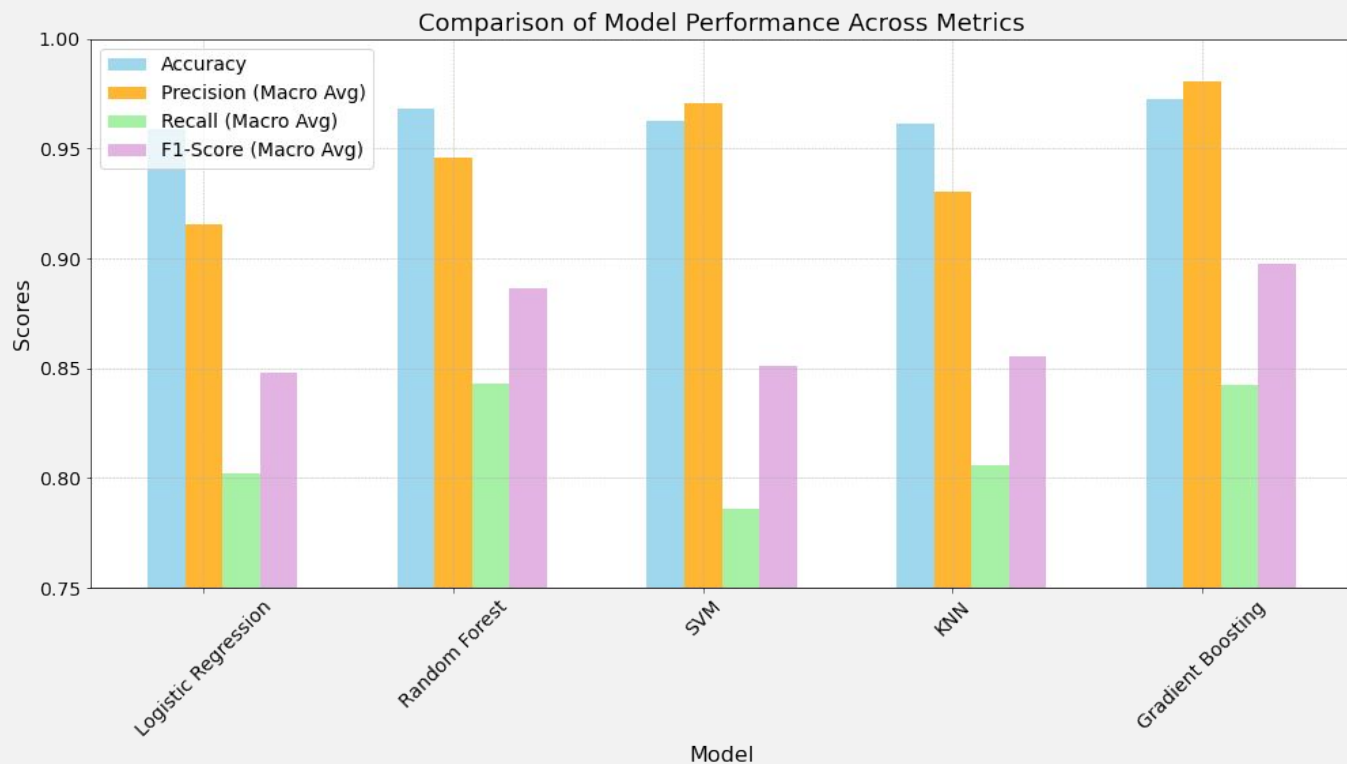|  | Negative | Positive |
|---|---|---|
| **Negative** | True Negative | False Negative |
| **Positive** | False Positive | True Positive |

**Confusion Matrix for Random Forest**

| True label \ Predicted label | Healthy | Diabetes |
|---|---|---|
| Healthy | 99.4% | 0.6% |
| Diabetes | 30.9% | 69.1% |

**Confusion Matrix for Logistic Regression**

| True label \ Predicted label | Healthy | Diabetes |
|---|---|---|
| Healthy | 99.1% | 0.9% |
| Diabetes | 38.7% | 61.3% |

**Confusion Matrix for SVM**

| True label \ Predicted label | Healthy | Diabetes |
|---|---|---|
| Healthy | 99.9% | 0.1% |
| Diabetes | 42.7% | 57.3% |

**Confusion Matrix for KNN**

| True label \ Predicted label | Healthy | Diabetes |
|---|---|---|
| Healthy | 99.3% | 0.7% |
| Diabetes | 38.1% | 61.9% |

**Confusion Matrix for Gradient Boosting**

| True label \ Predicted label | Healthy | Diabetes |
|---|---|---|
| Healthy | 99.9% | 0.1% |
| Diabetes | 31.4% | 68.6% |

| Models<br>Norms | Logistic Regression | Random Forest | SVM | KNN | Gradient Boosting |
|---|---|---|---|---|---|
| Accuracy | 0.802 | 0.843 | 0.786 | 0.806 | 0.842 |
| F1 Score | 0.756 | 0.815 | 0.728 | 0.761 | 0.813 |



Comparison of Model Performance Across Metrics

# Feature importance

# Questions?