# Detecting Diabetes Early
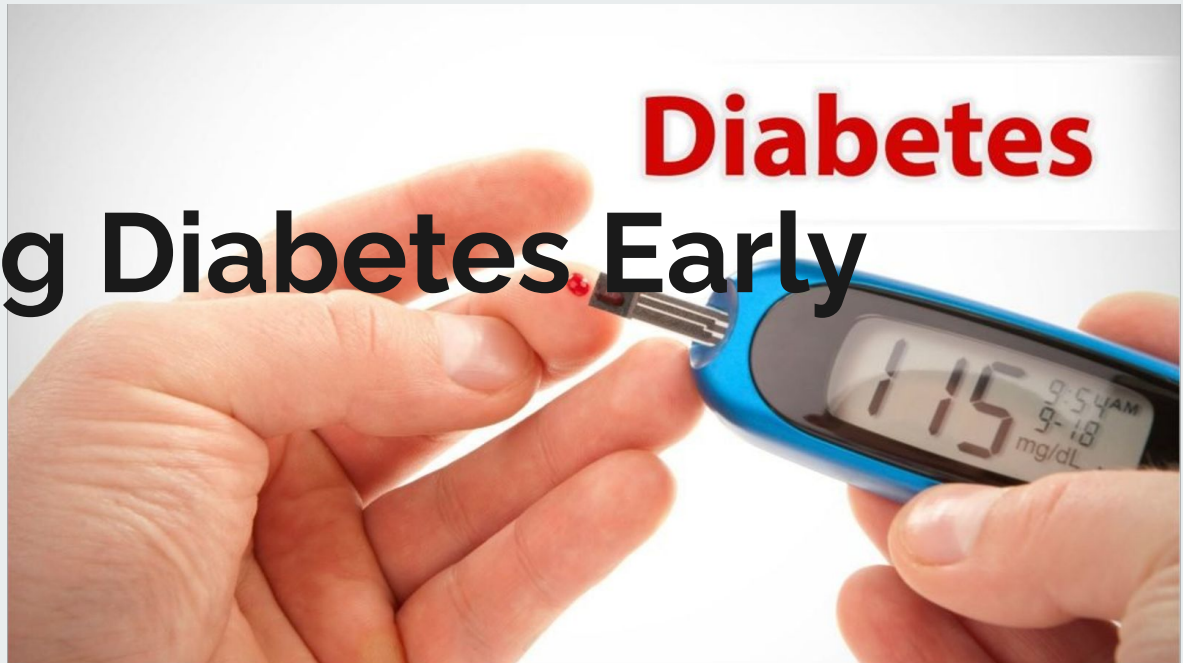
Daniel Kim



**Detecting Diabetes Early**

# DIABETES SYMPTOMS

## Problem Statement & Goal - How to Detect it?

Feeling Hungry

Feeling Thirsty

Blurry Eyesight

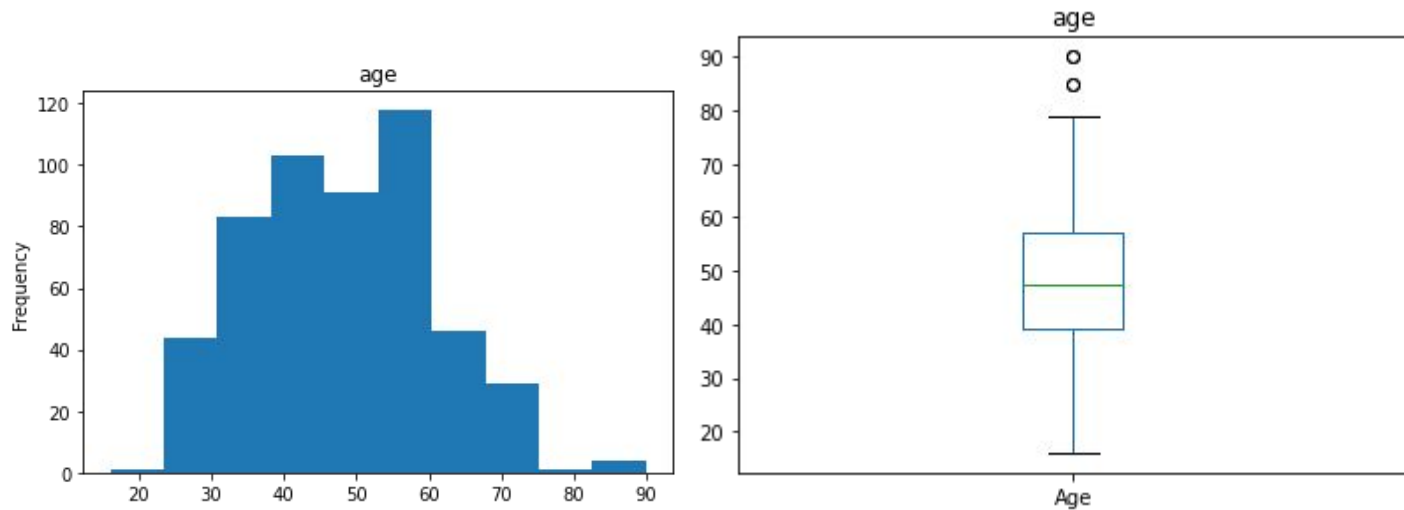Sudden Weight Loss

Frequent Urination

# Data Wrangling

1. Detecting NaN & Redundancy
   a. There were no NaN
   b. Although there were redundancies, it is too early to decide to remove redundancies because the redundancies are not the same but different individuals having same symptoms.
   c. Except for Age, all other features are binary, so it was needed to do one hot encoding to make a machine learning understand the data better.

2. One Hot Encoding

   After one -hot encoding, all the other binary columns ending with such as '_Female', '_Negative', or '_No' were removed.
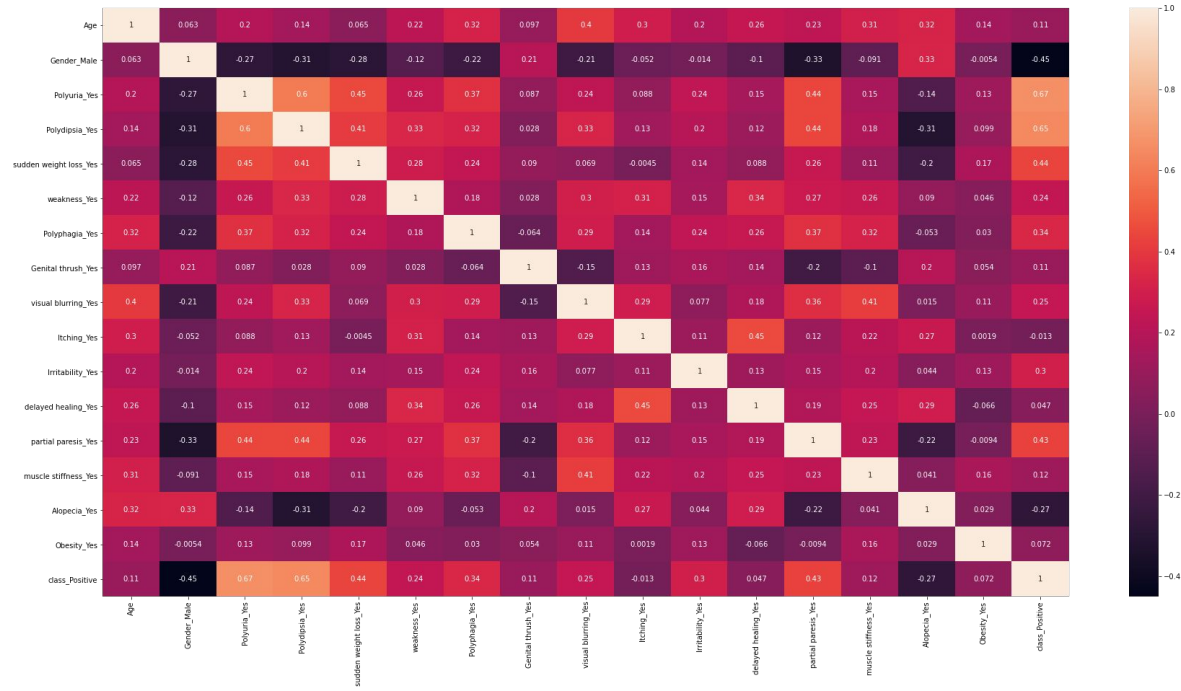
# Exploratory Data Analysis - Histogram & Boxplot



'Age' is the only feature showing normal distribution with a little right-skewness because all other features are binary. By its Boxplot, 'Age' seems having 2 outliers.

# Exploratory Data Analysis - Heat Map

Too many features seem correlated to each other, so it is needed to simplify their correlation by feature importance of several machine learning algorithms.

# Machine Learning - Feature Importance

'Polyuria' showed the strongest feature when it comes to identifying 'class_positive'. So we can conclude that 'Polyuria' plays a key role in detect diabetes.

Besides 'Polyuria', there are the following common strong features to detect diabetes even if their importance rates vary: Age, Gender_Male, and Polydipsia.

| Algorithms | Decision Tree Classifier | Random Forest Classifier | Gradient Classifier |
|---|---|---|---|
| Feature Importance % | Age: 12%<br>Gender_Male: 10%<br>**Polyuria: 42%**<br>Polydipsia: 6%<br>Irritability: 9%<br>muscle stiffness: 8%<br>Alopecia: 7% | Age: 10%<br>Gender_Male: 9%<br>**Polyuria: 21%**<br>Polydipsia: 18%<br>sudden weight loss: 6%<br>partial paresis: 5% | Age: 8%<br>Gender_Male: 11%<br>**Polyuria: 31%**<br>Polydipsia: 25%<br>Irritability: 6%<br>Alopecia: 7% |

# Machine Learning - Detection Accuracy

The Decision Tree Classifier showed the best score (97%). Other algorithms showed quite remarkable accuracy scores, but their scores are still lower than the Decision Tree Classifier. The reason for the difference seems due to 'max_depth' per each algorithm; the Decision Tree Classifier had more depths than others: 30 vs. 5.

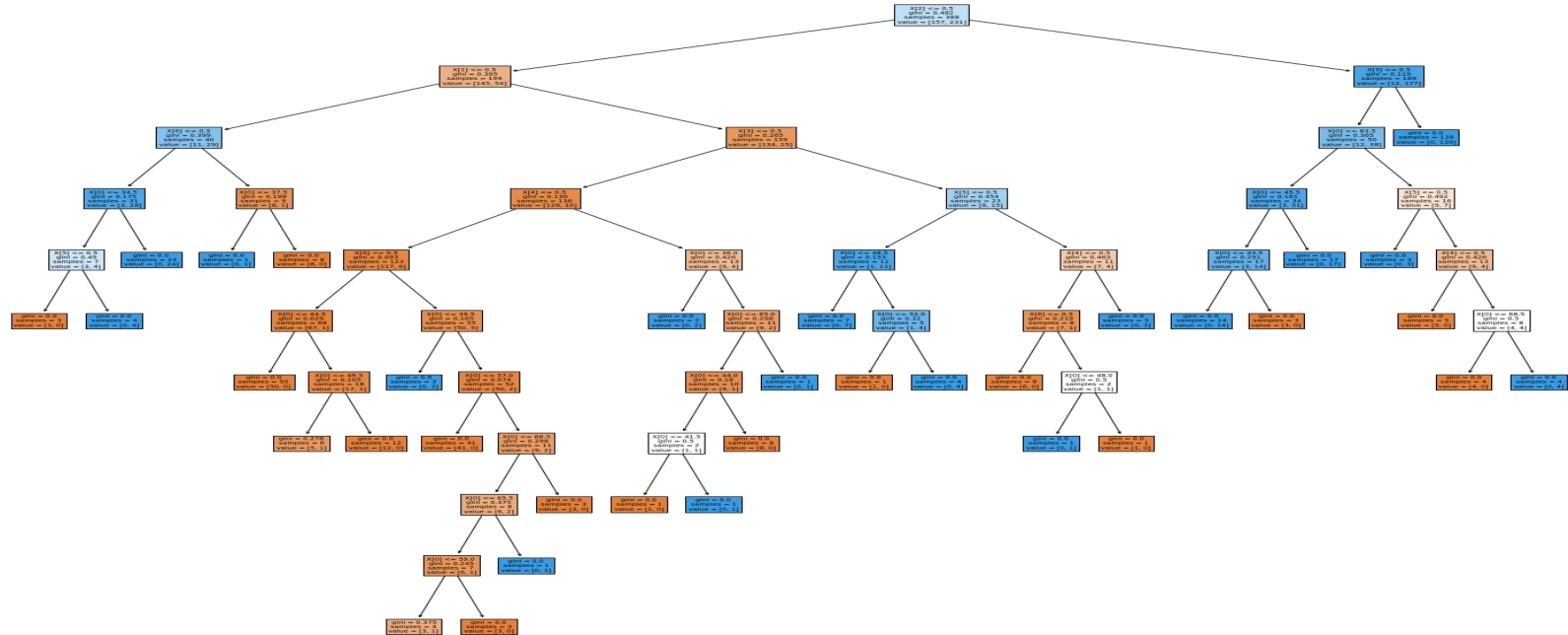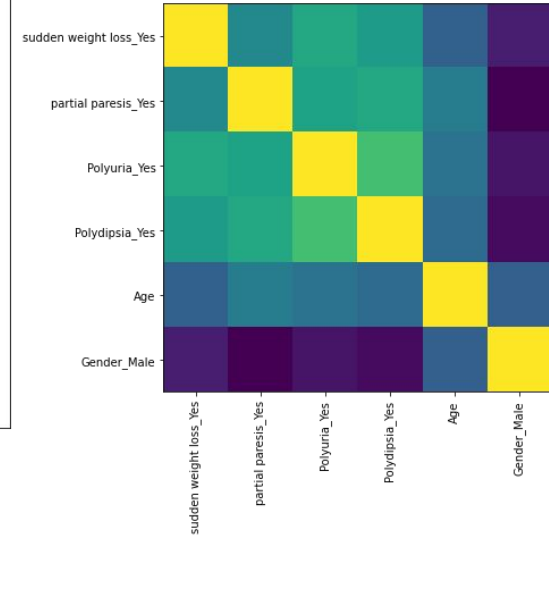| Algorithms | **Decision Tree Classifier** | Random Forest Classifier | Gradient Boost Classifier |
|---|---|---|---|
| Accuracy (R-Squared) | **97%** | 93% | 95% |
| Best Parameters by Randomized Search CV | **'splitter': 'random', 'random_state': 42, 'min_samples_split': 2, 'min_samples_leaf': 1, 'max_features': 'auto', 'max_depth': 30, 'criterion': 'entropy', 'class_weight': 'balanced'** | 'n_estimators': 1200, 'min_samples_split': 2, 'min_samples_leaf': 2, 'max_features': None, 'max_depth': 5, 'criterion': 'gini', 'class_weight': 'balanced' | 'n_estimators': 500, 'min_samples_split': 15, 'min_samples_leaf': 1, 'max_features': 'sqrt', 'max_depth': 5, 'loss': 'exponential', 'criterion': 'mae' |

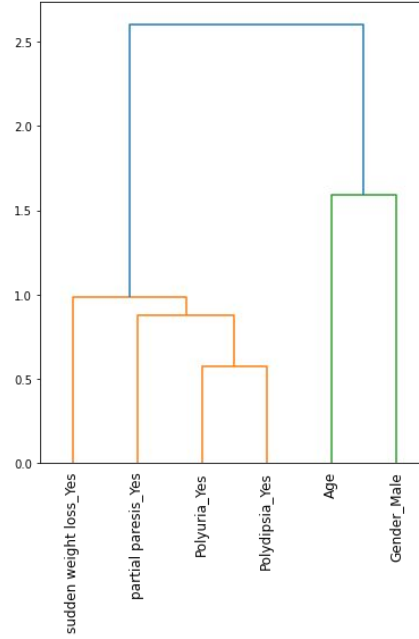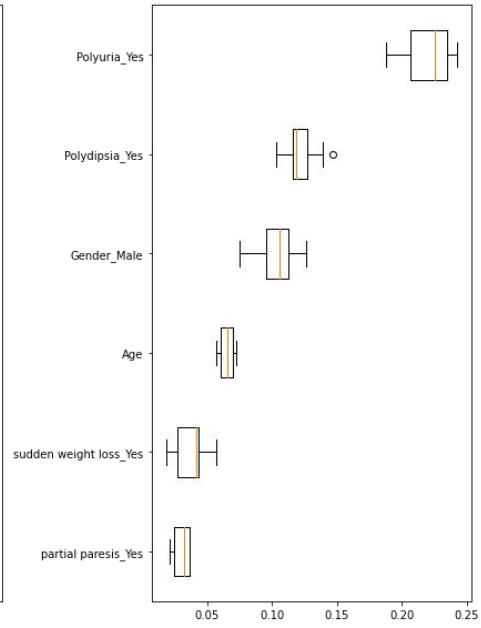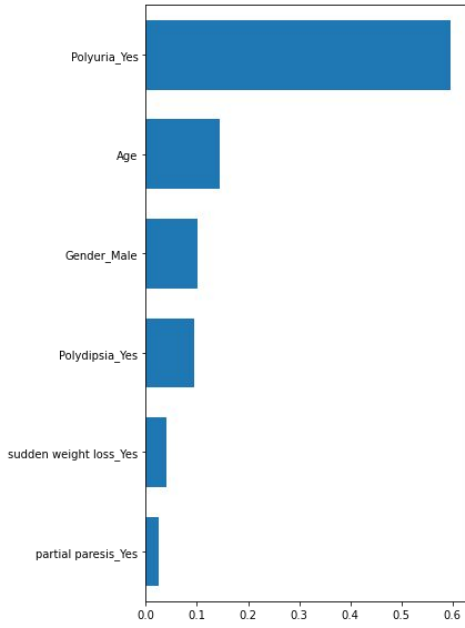# Recommendation - For Better Detection of Diabetes

In order to identify diabetes, you may conduct either of the following testing methods on a trade-off basis for both efficiency and effectiveness.

- Combinations of the above-mentioned features can give you optimal diabetes procedures. For example, a test of 'Polyuria' should be an independent variable for detecting diabetes. Such following features should be dependant variables with higher weights than others: 'Age', 'Gender_Male', and 'Polydipsia'. And then, the other else features should be optional test variables with lower weights.

- You can just conduct a diabetes test by the Decision Tree Classifier method.

# In-Depth Analysis - Decision Tree

# In-Depth Analysis - Random Forest Classifier

# In-Depth Analysis - Gradient Boost Classifier