# Final Project Data Science

**Fauzia Yumna Ayupuspita**
**Iqbal Muhammad**
**Muhammad Fahmi**
**Siti Rabiatul Adwiyah**
**Yoga Mahardika Sidiq**

**Diabetes Prediction**

# Outline

**1** Overview About Diabetes

**2** Problem Definition, Goal, Methodology

**3** Result Overview

**4** EDA

**5** Machine Learning

**6** Conclusion

# Overview about Diabetes

## What is Diabetes?
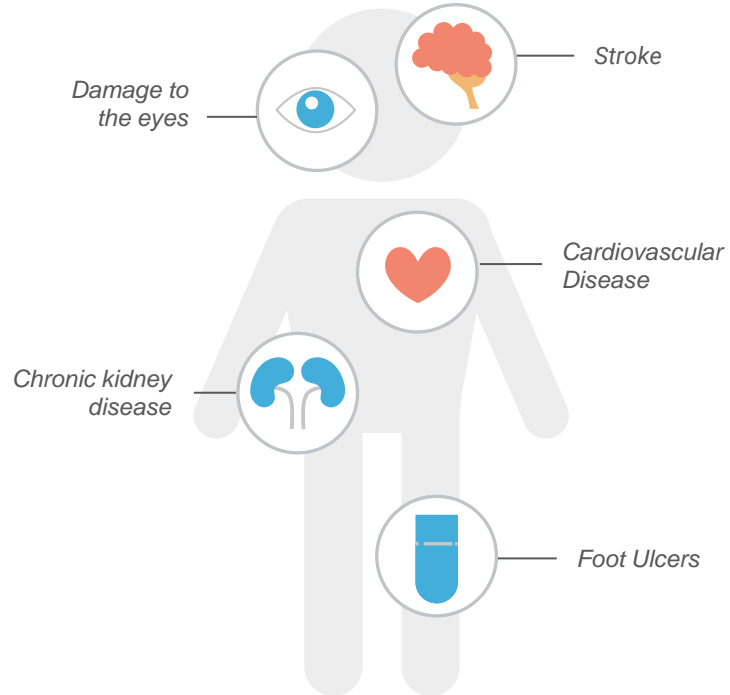
**Disease that occurs when your blood glucose is too high**

## Symptoms

**Thirst**

**Fatigue**

**Weight loss**

## Organs Affected

*Stroke*

*Damage to the eyes*

*Cardiovascular Disease*

*Chronic kidney disease*

*Foot Ulcers*

# Problem, Goals, Methodology

## Problem

- Healthcare sector have large amount databases
- Existing method for diabetes detection is uses lab tests
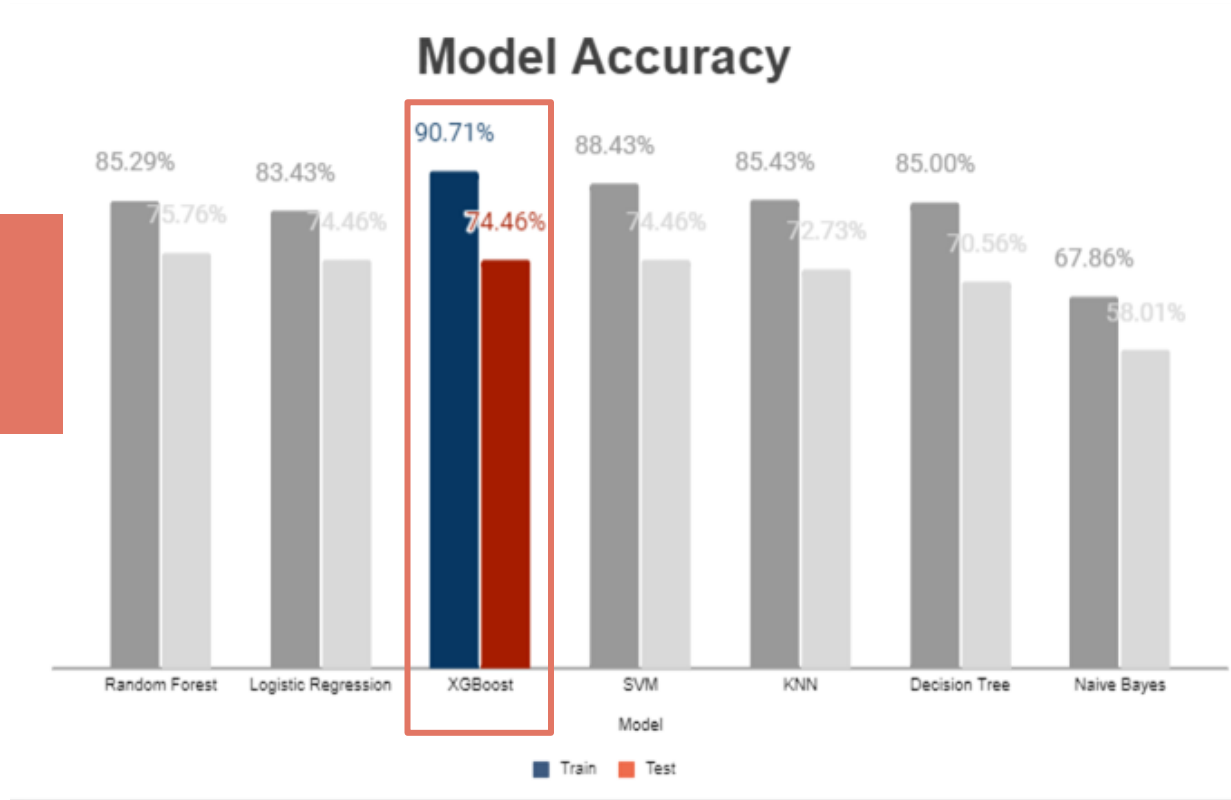- Existing Method is time consuming

## Goals

Building predictive model for diabetes prediction so that indicated patient can check further as soon as possible
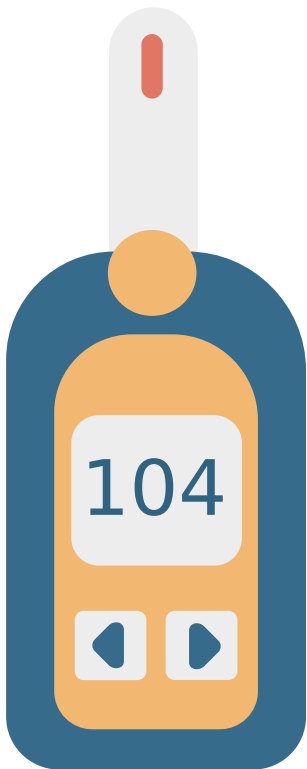
## Methodology

Using Machine Learning - Supervised

5.8

# Result Overview



Model Accuracy

The best model is XGBoost with Accuracy Test is 74.46%

# Methodology



**01**
### Database
Extracting dataset from database *)

**02**
### EDA
Understanding data from univariate, bivariate even multivariate analysis

**03**
### Data Preprocessing
Feature Engineering (Classifying age, Glucose, BMI, Blood Pressure, Insulin and Pregnancies

**04**
### Modelling
Splitting data, model training & model testing, accuracy

**05**
### Result
Diabetes Prediction Model

# Exploratory Data Analysis (EDA)

Data Check

| No | Feature | IsNull |
|----|---------|--------|
| 1 | Pregnancies | 0 |
| 2 | Glucose | 5 |
| 3 | Blood Pressure | 227 |
| 4 | Skin Thickness | 374 |
| 5 | BMI | 11 |
| 6 | Diabetes Pedigree Function | 0 |
| 7 | Age | 0 |
| 8 | Outcome | 0 |

# Exploratory Data Analysis (EDA)

Univariate - Numerical
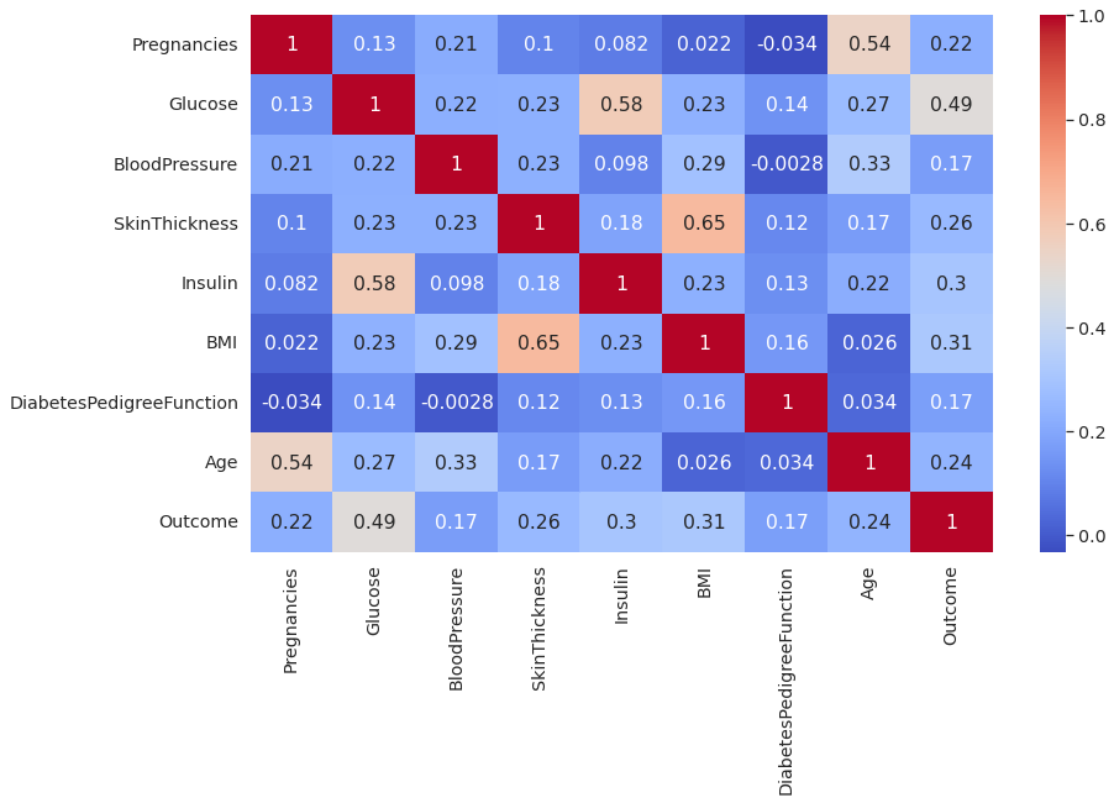


## Normal Distribution

- Glucose and
- Blood Pressure

## Positive Skewed

Most of data are positive skewed except Glucose and Blood Pressure data
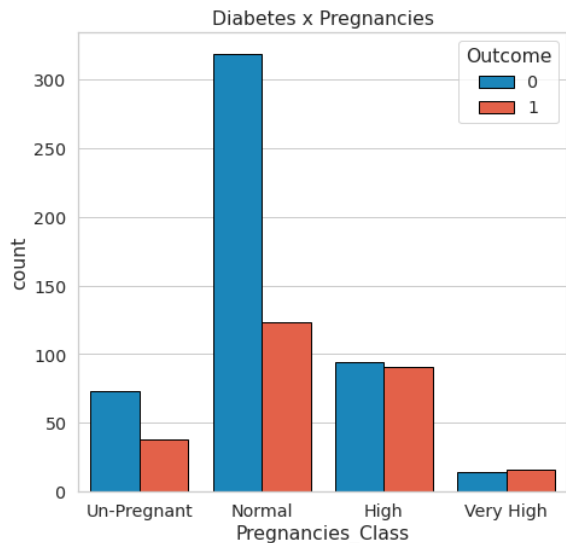
# Exploratory Data Analysis (EDA)
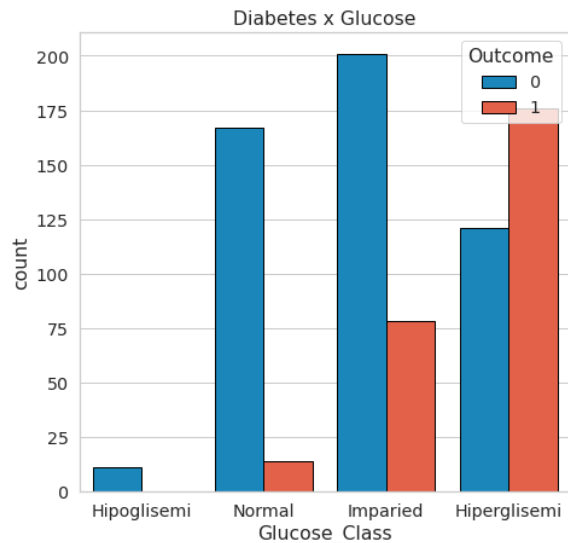
Handling Outlier
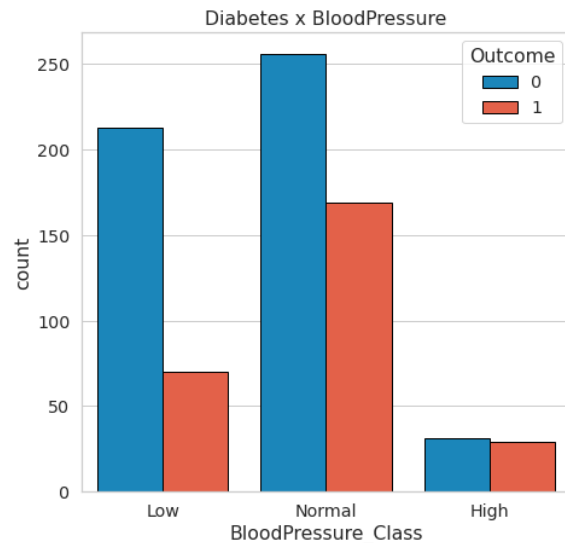
# Exploratory Data Analysis (EDA)

Bivariate – Categorical



In proportion, **Very High** pregnancies class (>10 times pregnant) is the highest class that have diabetic
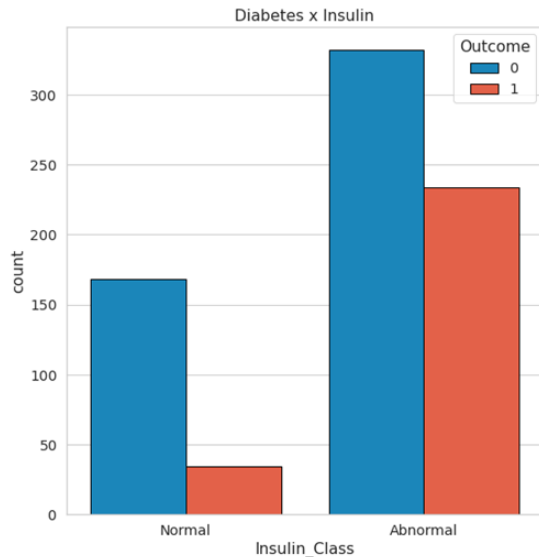
- Hipoglisemi Glucose is non diabetic
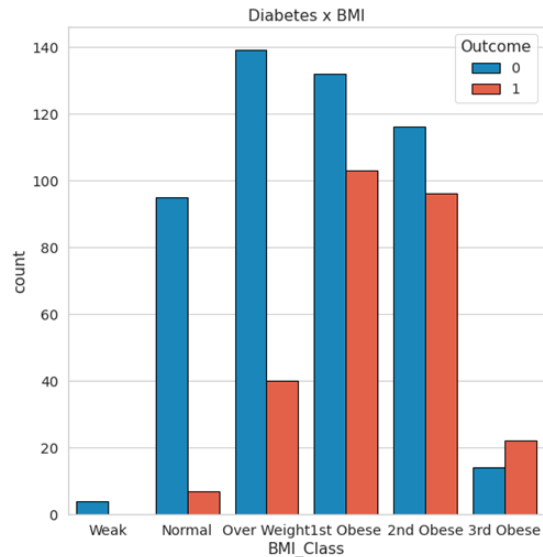- Majority Hiperglisemi Class have diabetic

In proportion, **High Blood Pressure Class** (>90) is the highest class that have diabetic

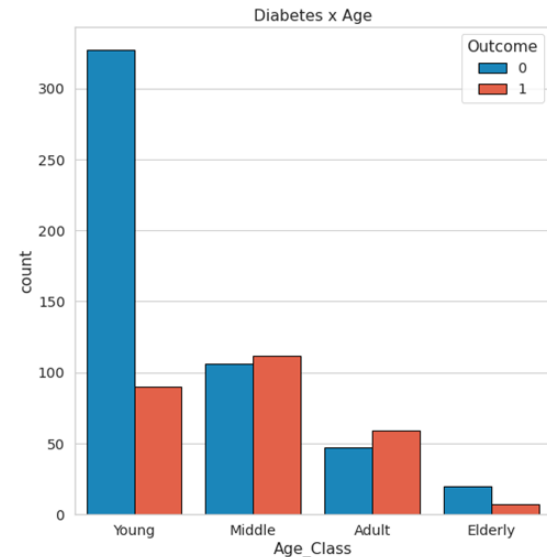# Exploratory Data Analysis (EDA)

Bivariate – Categorical



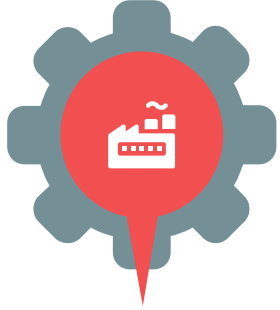In Average, patient with **Abnormal Insulin** class have diabetic

The higher the class of BMI then the more patients have diabetic

Middle Age Class is the highest class have diabetic

# Machine Learning
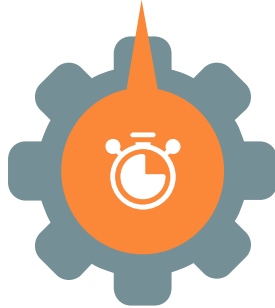
**Encoder**

Change categorical variables to numeric

**Train-Test Split**

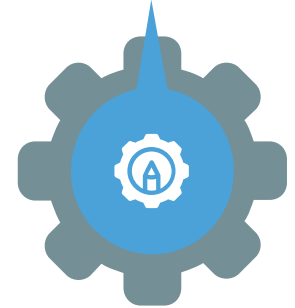Separating data for learning and testing

**Resampling**

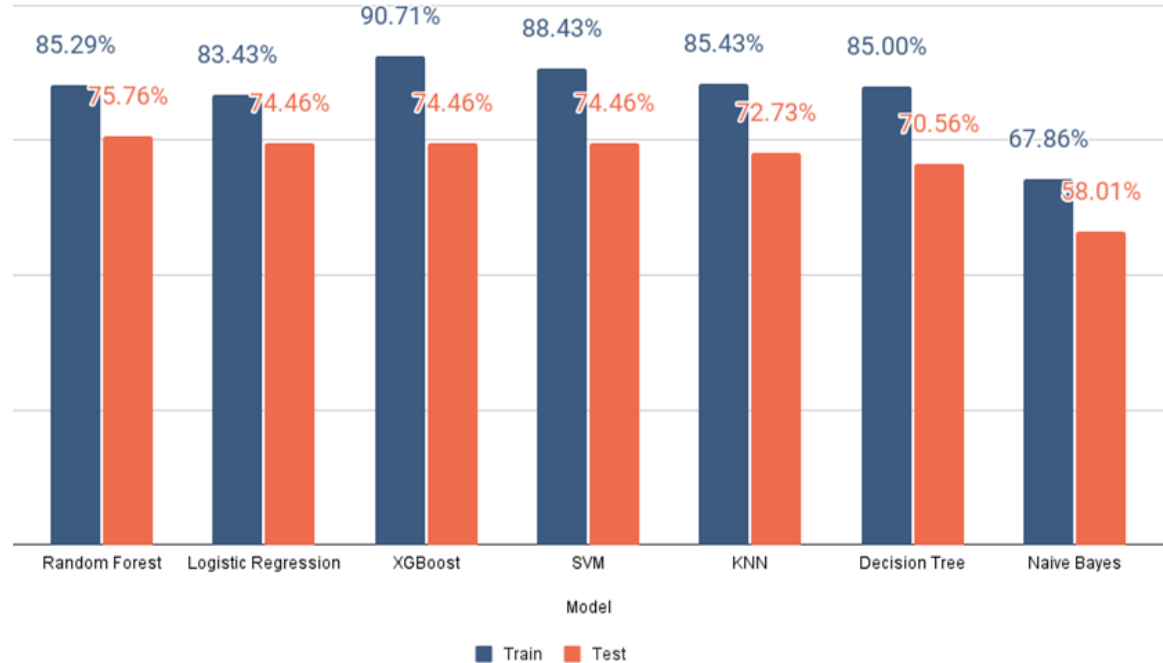Perform SMOTE on diabetes data to create balanced data

**Modelling**

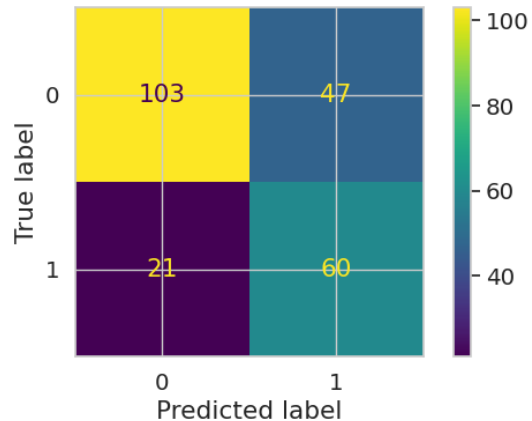Create 7 supervised modeling scenarios
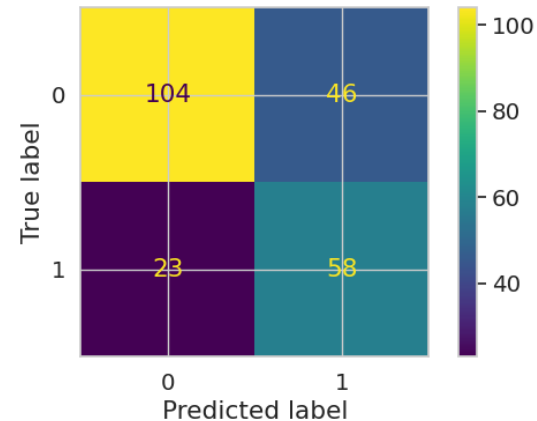
Result Machine Learning Model
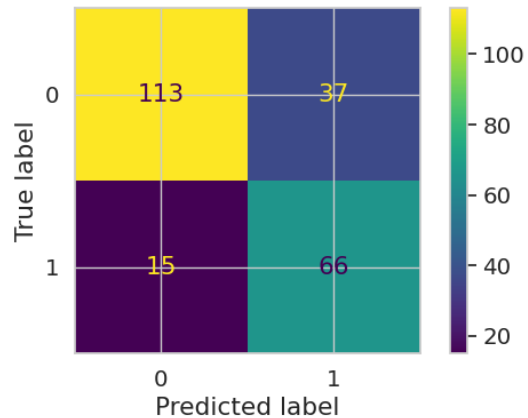
# Confusion Matrix Models
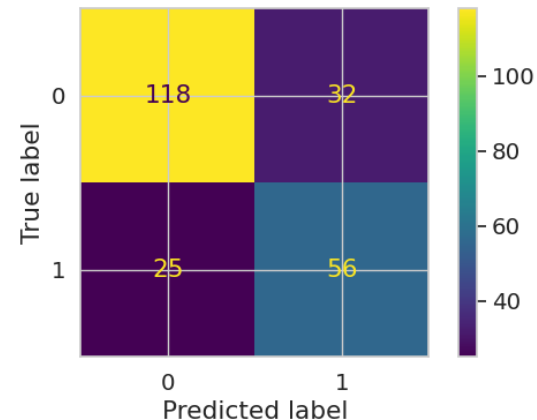


Confusion Matrix for Decision Tree Model

Confusion Matrix for KNN Model

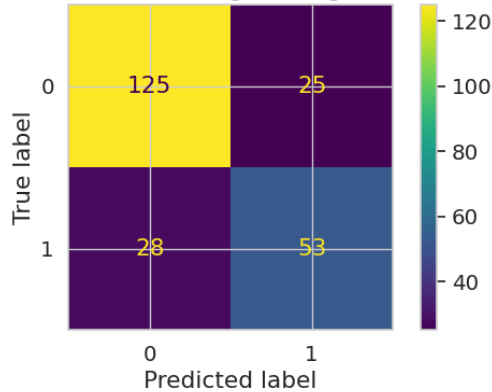Confusion Matrix for Random Forest Model
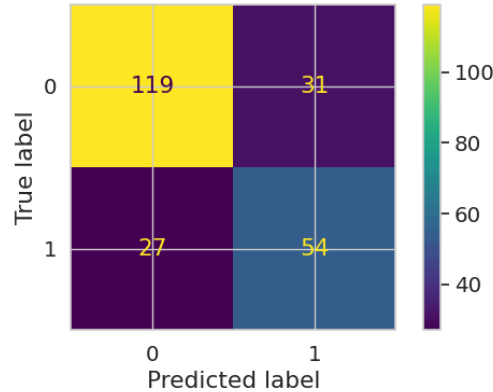
Confusion Matrix for SVM Model
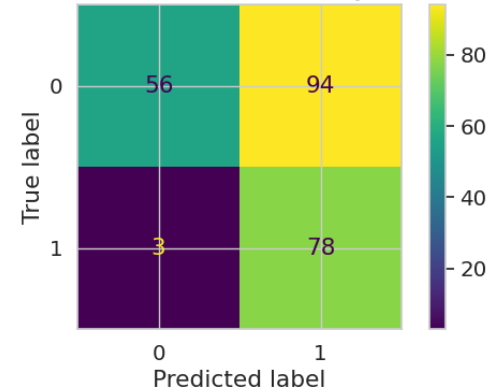
# Confusion Matrix Models



Confusion Matrix for Logistic Regression Model

Confusion Matrix for XGBoost Model

Confusion Matrix for Naive Bayes Model

After looking at the confusion matrix for 7 models, we can draw the conclusion that **the best model to use is XGBoost**, considering the highest proportion of positive true & negative true values

# THANK YOU|