# Methodology & Architecture

- **Architecture of proposed model:**

# Data exploration & Problem Formulation

Scatter plot between pairs of variables





corr matrix



class distribution

**Important Information Discovered**

- The target classes distribution is not even
- Collinearity among some of the predictors
- No clear distinction between target classes in some predictors

**Problem Formulation**

- Imbalance classification problem
- Distinction between target classes in some predictors need to be amplified
- It is expected that model that won't be significantly affected by collinearity will perform better

# Data preparation

## Data Cleaning

- Purpose:
  - Remove noise from the data to enhance model performance

- Methodology:
  - Outlier removal with condensed nearest neighbors undersampling technique
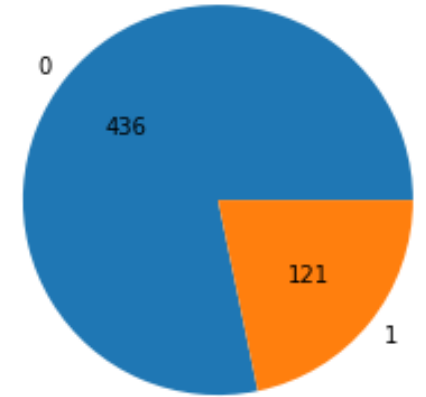
- Description:
  - This method can not only make the class distribution in the training set more balanced, but also remove outlier from both classes in the training set, we implemented a modification of the aforementioned technique called **Tomek Links**, developed by Tomek (1976)
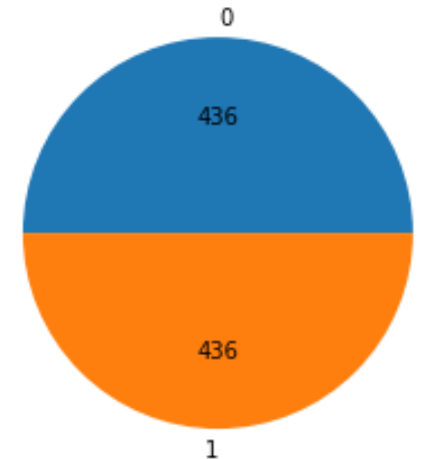


*Class Distribution after Tomek Link*

## Imbalance Handling

- Purpose:
  - Balance the distribution of classes in training set such that the model fitted will demonstrate a higher ability to identity minority class accurately.

- Methodology:
  - Oversampling with synthetic minority over-sampling technique (SMOTE)

- Description:
  - We implemented a modification of the aforementioned technique called **Borderline SMOTE**, developed by Han, H., Wang, WY., Mao, BH. (2005). It can identify noises and signals in minority class, and only oversample the signals.



*Class Distribution after Borderline SMOTE*

# Data Preparation

- **Normalization:**
  - ➤ *StandardScaler()*

$$Z = \frac{x - \mu}{\sigma}$$

  where, μ is the mean of all sample data, and σ denotes the standard deviation.

```
from sklearn.preprocessing import StandardScaler
stdScaler = StandardScaler().fit(X_train)  #return mean and sd
X_train = stdScaler.transform(X_train)  #return sdscaler number
X_test = stdScaler.transform(X_test)
```

  - ➤ Normalize **before** PCA to ensure each feature is **in the same scale** to prevent over-capturing of certain features with large values, and accelerate convergence of gradient descent method.

- **PCA:**
  - ➤ Reduce Dimensionality.
  - ➤ Eliminate redundancy and data noise.
  - ➤ Identify the **most variable direction** of the observations.
  - ➤ Allow us to determine the property of density $f(X)$ when $Y$ is unknown.

```
from sklearn.decomposition import PCA

pca=PCA(n_components=2, svd_solver='auto').fit(X_train)
# Dimensionality reduction
X_train_pca=pca.transform(X_train)
X_test_pca=pca.transform(X_test)
```

# Model 1. Logistic Regression and QDA

**Step 1:**
Input features and the response to train

X1,X2,X3,X4,X5,X6,X7,X8, X9,X10,X11,X12    &    Y
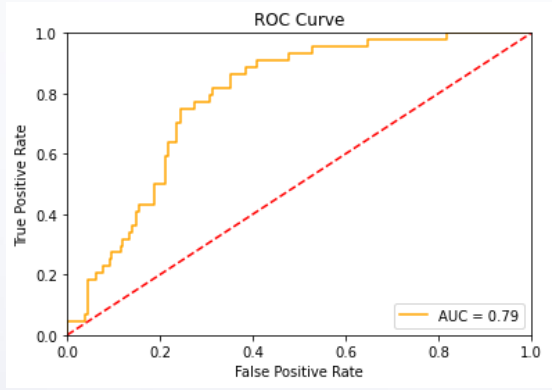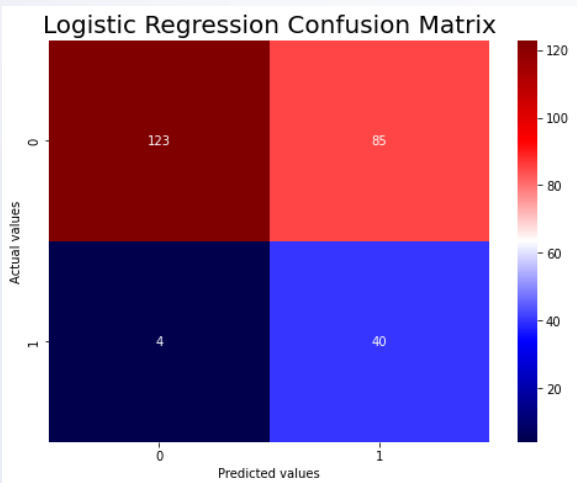
**Step 2:**
Model specification

**(1)Logistic Regression**
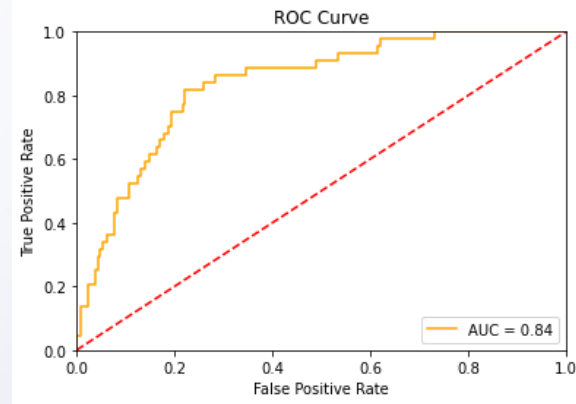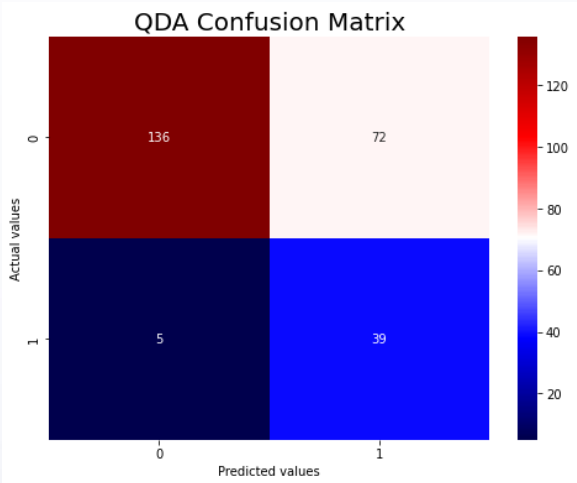
**(2) QDA**

**Step 3:**
Model evaluation

AUC=0.79

AUC=0.84



|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.97 | 0.59 | 0.73 | 208 |
| 1 | 0.32 | 0.91 | 0.47 | 44 |
| accuracy |  |  | 0.65 | 252 |
| macro avg | 0.64 | 0.75 | 0.60 | 252 |
| weighted avg | 0.86 | 0.65 | 0.69 | 252 |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.96 | 0.65 | 0.78 | 208 |
| 1 | 0.35 | 0.89 | 0.50 | 44 |
| accuracy |  |  | 0.69 | 252 |
| macro avg | 0.66 | 0.77 | 0.64 | 252 |
| weighted avg | 0.86 | 0.69 | 0.73 | 252 |



Logistic Regression Confusion Matrix

ROC Curve — AUC = 0.79
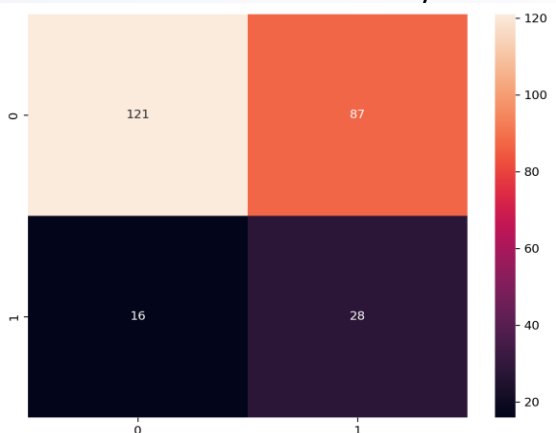
QDA Confusion Matrix

ROC Curve — AUC = 0.84

# Model 3. SVM

- **A popular method used to solve the binary classification problem and implement the prediction of binary data, offering high accuracy.**
- **It can easily handle multiple continuous and categorical variables.**



**Step 1:**
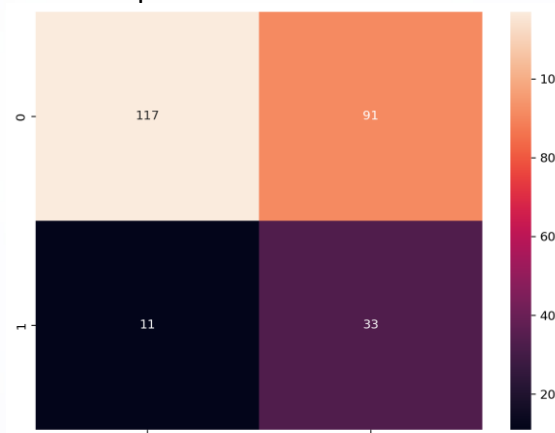Fit the model and obtain the confusion matrix and accuracy score.

*GridSearchCV* **method**

**Parameter Tuning**

**Step 2:**
Fit the optimal model with the best parameter combination.

**Evaluation**

**Step 3:**
Obtain **ROC/AUC**.

**Results**

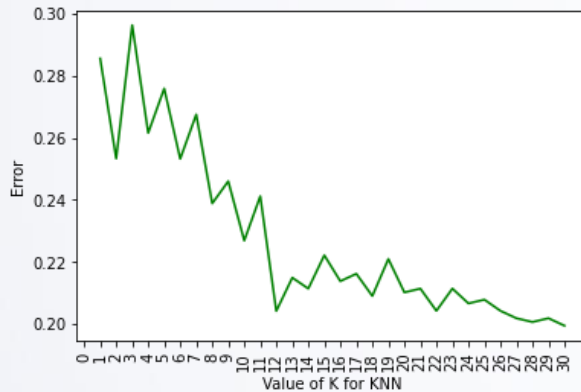| Evaluation Metrics | Value |
|---|---|
| Accuracy (Training Set) | 0.72 |
| Accuracy (Test Set) | 0.59 |
| AUC | 0.72 |

# Model 4. KNN

- **One of the simplest and most common algorithms.**
- **New data can be easily classified into a well-suited category.**

**Step 1:**
Select the optimal **K** by applying **Cross-validation** method.



Select **K** with the **minimum error, but K is not very large to shorten training time.**

**Step 2:**
Let **K = 12** and fit the optimal model. Obtain the accuracy and confusion matrix.



**Evaluation**

**Step 3:**
Obtain **ROC/AUC**.



|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.91 | 0.66 | 0.76 | 208 |
| 1 | 0.30 | 0.68 | 0.41 | 44 |
| accuracy |  |  | 0.66 | 252 |
| macro avg | 0.60 | 0.67 | 0.59 | 252 |
| weighted avg | 0.80 | 0.66 | 0.70 | 252 |

**Results**

| Evaluation Metrics | Value |
|---|---|
| Accuracy (Training Set) | 0.75 |
| Accuracy (Test Set) | 0.67 |
| AUC | 0.74 |

# Model 5. Random forest and XGBoost



**Random Forest**

Training Accuracy = 100%

**XGBoost**

Training Accuracy = 100%

Step 1) Input Processed train data (x,y), fit model, and obtain train accuracy

Step 2) Get prediction and prediction probability with fitted model

Step 3) Get model evaluation with confusion matrix, classification report, and ROC curve
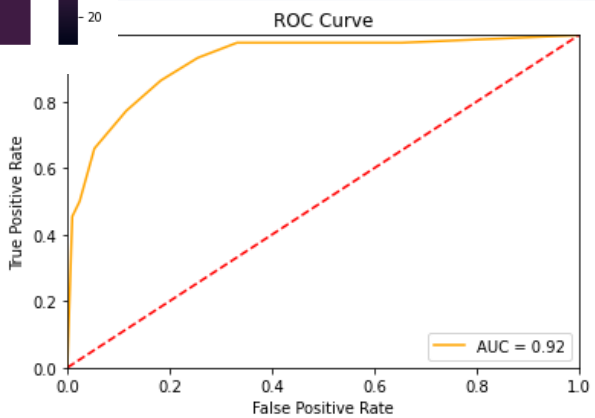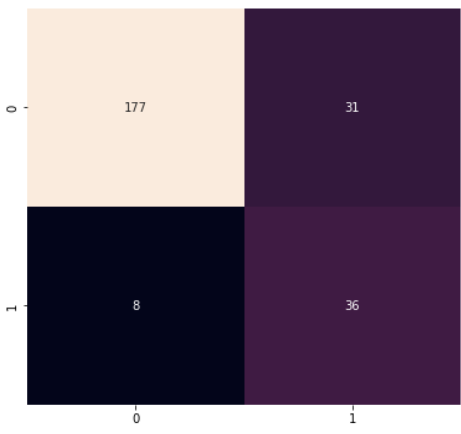
|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.96 | 0.85 | 0.90 | 208 |
| 1 | 0.54 | 0.82 | 0.65 | 44 |
| accuracy |  |  | 0.85 | 252 |
| macro avg | 0.75 | 0.83 | 0.77 | 252 |
| weighted avg | 0.88 | 0.85 | 0.86 | 252 |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.96 | 0.82 | 0.89 | 208 |
| 1 | 0.50 | 0.84 | 0.63 | 44 |
| accuracy |  |  | 0.83 | 252 |
| macro avg | 0.73 | 0.83 | 0.76 | 252 |
| weighted avg | 0.88 | 0.83 | 0.84 | 252 |

# Model comparison by predictability performance

|  | Logistic Regression | QDA | Random Forest | XGBoost | SVM | KNN |
|---|---|---|---|---|---|---|
| Accuracy (train) | 0.72 | 0.81 | 1 | 1 | 0.72 | 0.75 |
| Accuracy (test) | 0.65 | 0.69 | 0.85 | 0.83 | 0.59 | 0.67 |
| F1 score (0) | 0.73 | 0.78 | 0.90 | 0.89 | 0.69 | 0.77 |
| F1 score (1) | 0.47 | 0.50 | 0.65 | 0.63 | 0.39 | 0.42 |
| AUC | 0.79 | 0.84 | 0.92 | 0.89 | 0.72 | 0.74 |

**ROC - AUC**

➢ The area under the ROC curve.

➢ The **larger** the **better, 1** is the **ideal state.**

➢ We use class 1 as the true class when calculating ROC – AUC in our implementation

| Metrics | Formula |
|---|---|
| Accuracy | $\frac{TP+TN}{TP+FP+TN+FN}$ |
| F1-Score | $\frac{2 \times Precision \times Recall}{Precision + Recall}$ |
| Precision | $\frac{TP}{TP+FP}$ |
| Recall | $\frac{TP}{TP+FN}$ |

| Machine learing \ Manual counting | True | False |
|---|---|---|
| True | True Positive (TP) | False Positive (FP) |
| False | False Negative (FN) | True Negative (TN) |

*Rationale to use as core metric:*

- Focus on evaluating the ability of identifying minority class from other class, an important attribute under the context of imbalance classification