# Employee Attrition predection

By: Muneera Alshunaifi

Statistics

DATA

70%

0.49

# *Outline*

Introduction

Data

Pre-processing

Modeling and evaluation

Conclusions

# *Introduction*

## WHAT IS ATTRITION

The Employee attrition occurs when an employee departs a company for a variety of reasons

## WHY?

Many reasons such as monthly income or job role

## SOLUTION

Build model that Predicts the appropriate reasons that led to the attrition decision in order to reduce the attrition within a company in the future

# *Data Description*

- Obtained from kaggle.com
- Fictional dataset created by IBM data scientist
- 1470 rows, 35 columns

# *Data Description*
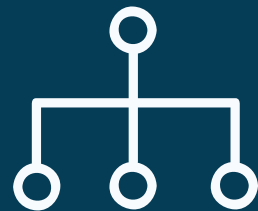
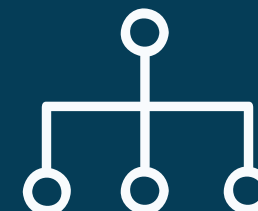| Column Type | Description | Columns names |
|---|---|---|
| **Numeric columns** | Related to personal information | age, distance_from_home, employee_number (id variable) |
| | Related to income | hourly_rate, daily_rate, monthly_rate, monthly_income, percent_salary_hike |
| | Related to time in company | years_at_company, years_in_current_role, years_since_last_promotion, years_with_curr_manager, total_working_years |
| | Others | num_companies_worked, standard_hours(to delete), training_times_last_year, employee_count (to delete) |
| **Categorical columns** | Binary | Attrition (Target variable), gender, over18 (to delete), over_time |
| | Nominal | department, education_field, job_role, marital_status |
| | Ordinal | environment_satisfaction, job_satisfaction, relationship_satisfaction, work_life_balance,job_involvement,performance_rating business_travel, education, job_level, stock_option_level |

# Pre processing steps

EDA

Feature
engineering

Scaling

Sampling

# EDA

# *Using pandas profiling*

## INSIGHT

Part of the analysis using pandas profiling shows overview of the dataset



## Overview

Overview    Reproduction    Warnings **11**

### Dataset statistics

| | |
|---|---|
| **Number of variables** | 32 |
| **Number of observations** | 1470 |
| **Missing cells** | 0 |
| **Missing cells (%)** | 0.0% |
| **Duplicate rows** | 0 |
| **Duplicate rows (%)** | 0.0% |
| **Total size in memory** | 1.0 MiB |
| **Average record size in memory** | 722.8 B |

### Variable types

| | |
|---|---|
| **NUM** | 17 |
| **CAT** | 13 |
| **BOOL** | 2 |

# EDA

# *Correlations*

## INSIGHT

No segnificant corr between attrition and any feature

# *Feature engineering*

- Exclude unnecessary features (over18, standard hours, employee count)
- Convert the attrition column from Yes,No to 1,0
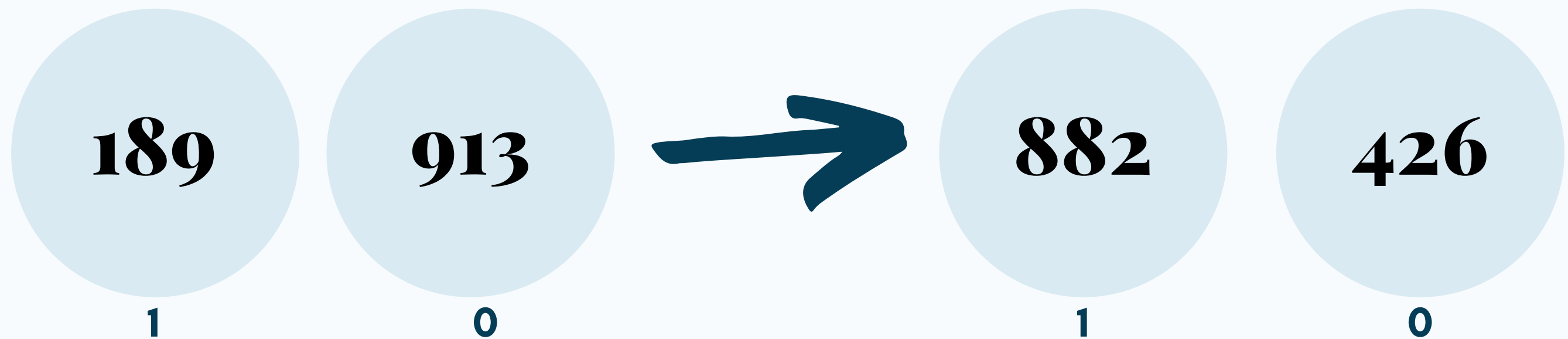- Split the data into 30% test and 70% train

# Scaling

- Using MinMax scaler

# Sampling

## USING HYBRID TECHNIQUE (SMOTE+ENN)

- ENN: algorithm for pattern recognition that predicts the pattern of an unknown test sample hinged on the highest gain of intraclass coherence
- ENN+SMOTE helps in doing extensive data cleaning.
- The misclassification by NN's samples from both the classes are removed.
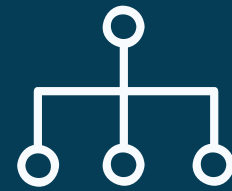
# *Sampling*

## USING HYBRID TECHNIQUE (SMOTE+ENN)

189 (1)  913 (0)  →  882 (1)  426 (0)

# Modeling and evaluation
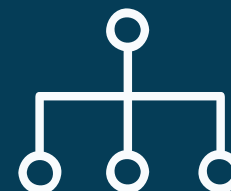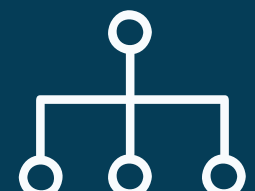
Logistic regression

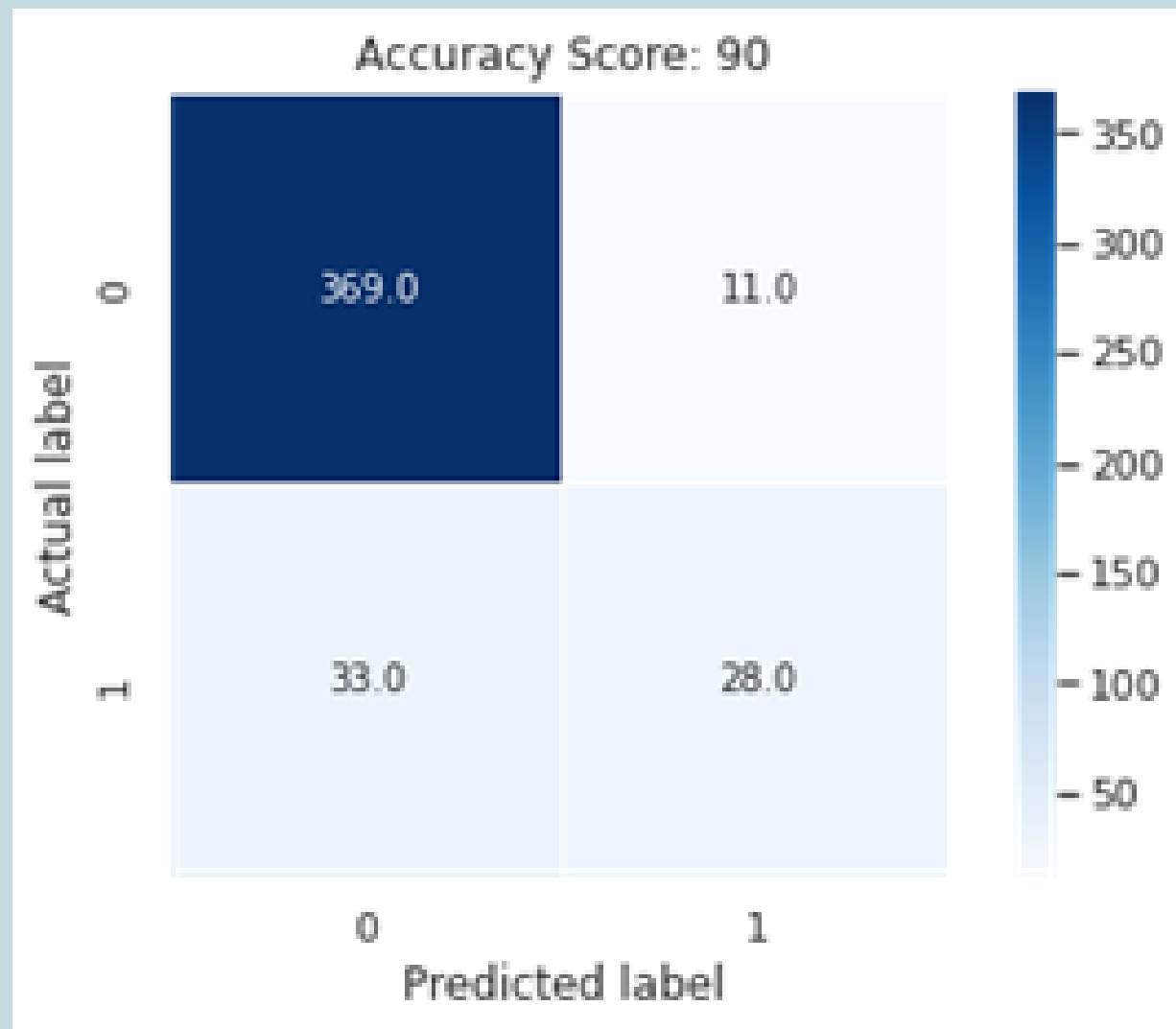Decision Tree
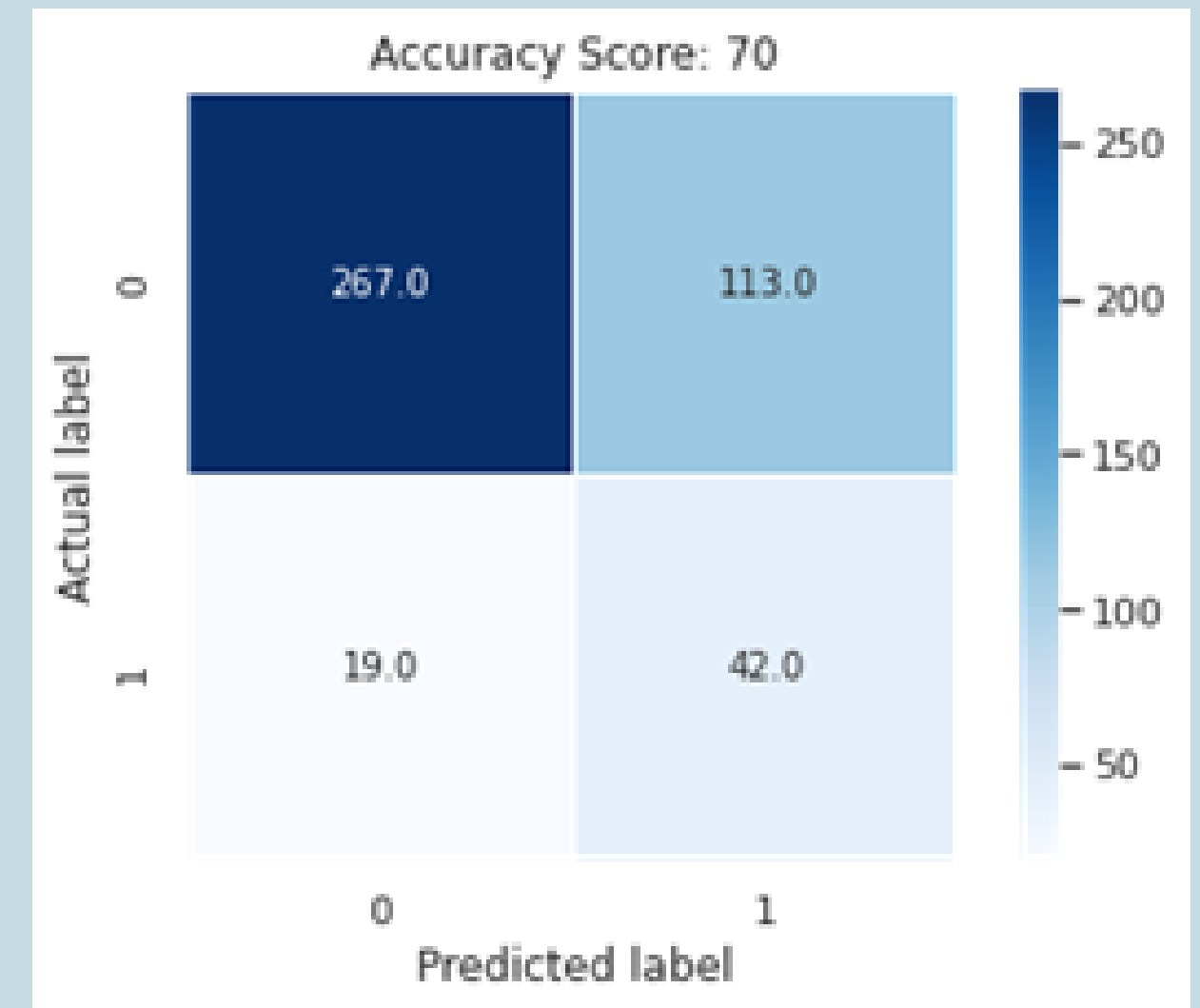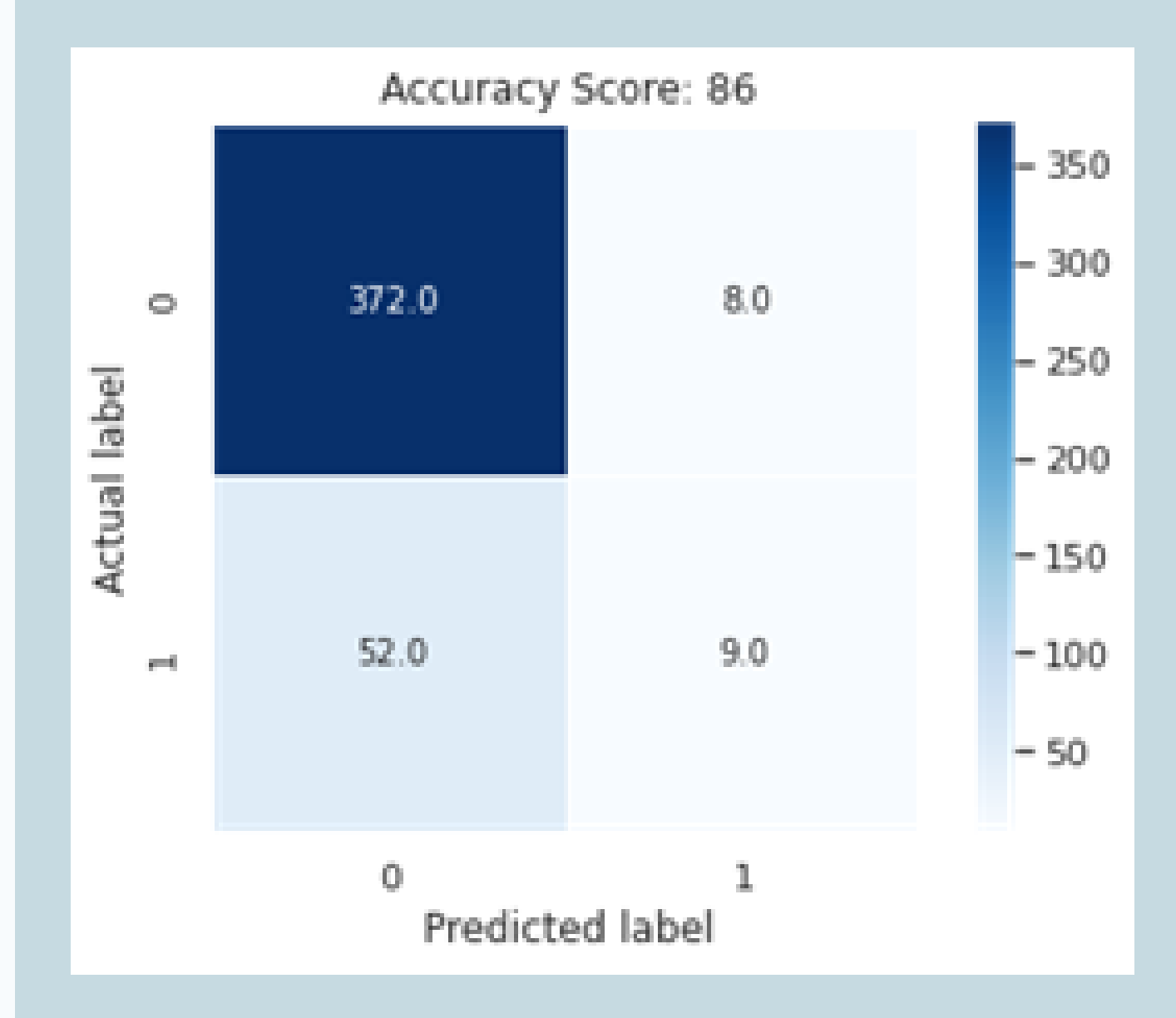
SVM

Ensemble: AdaBoost

Models comparison

Optimization
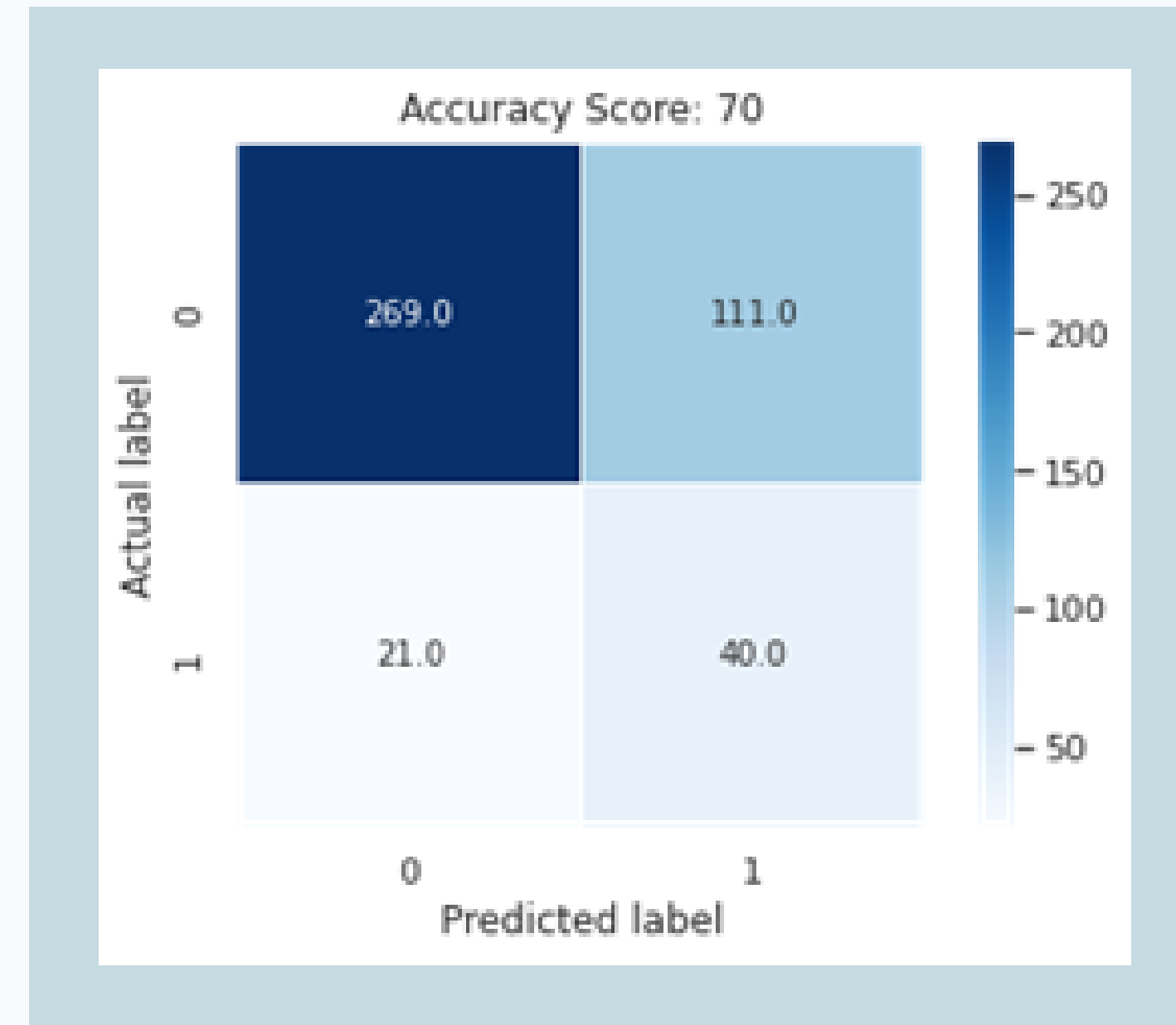
Logistic regression

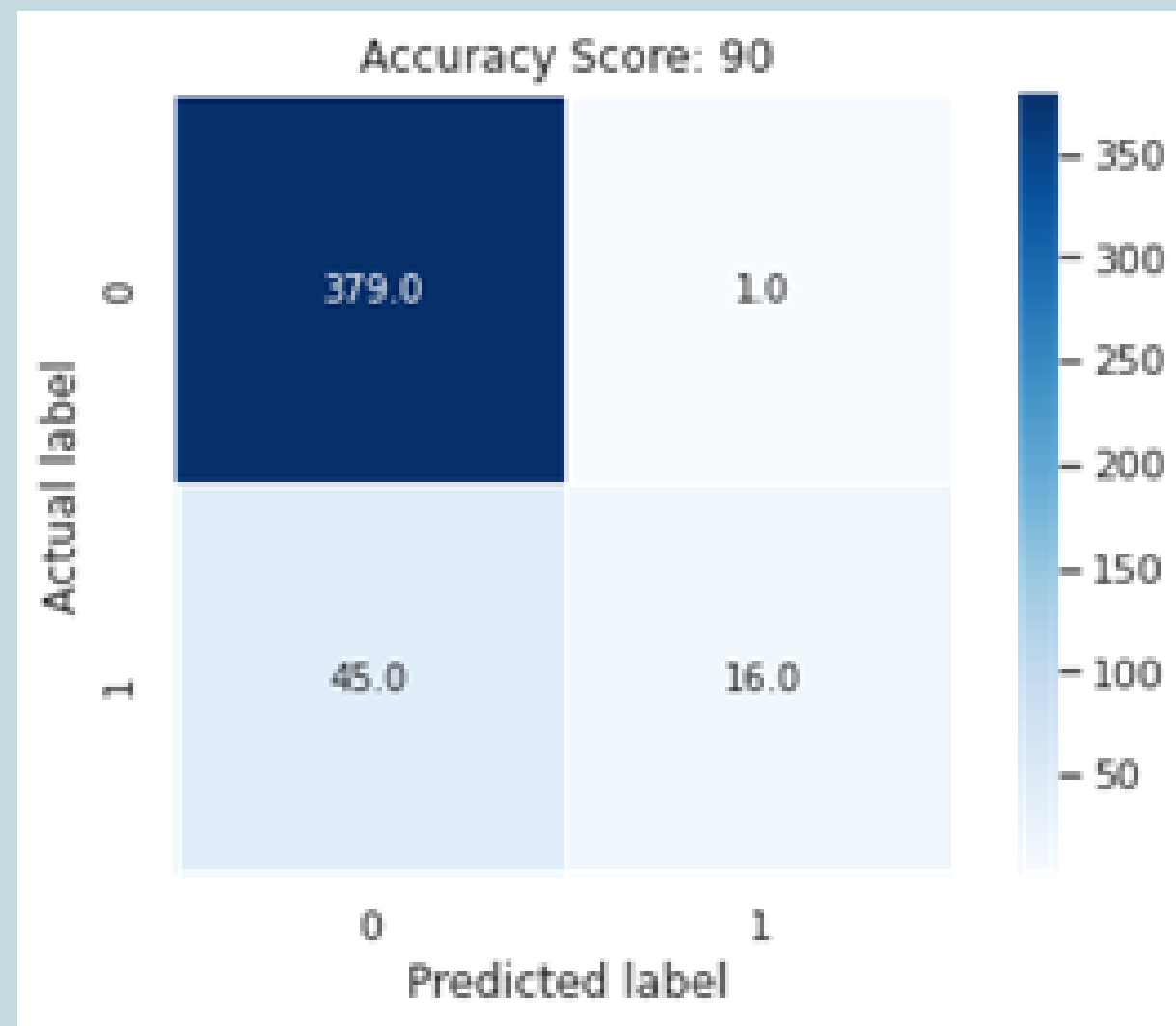Confusion matrix before and after sampling
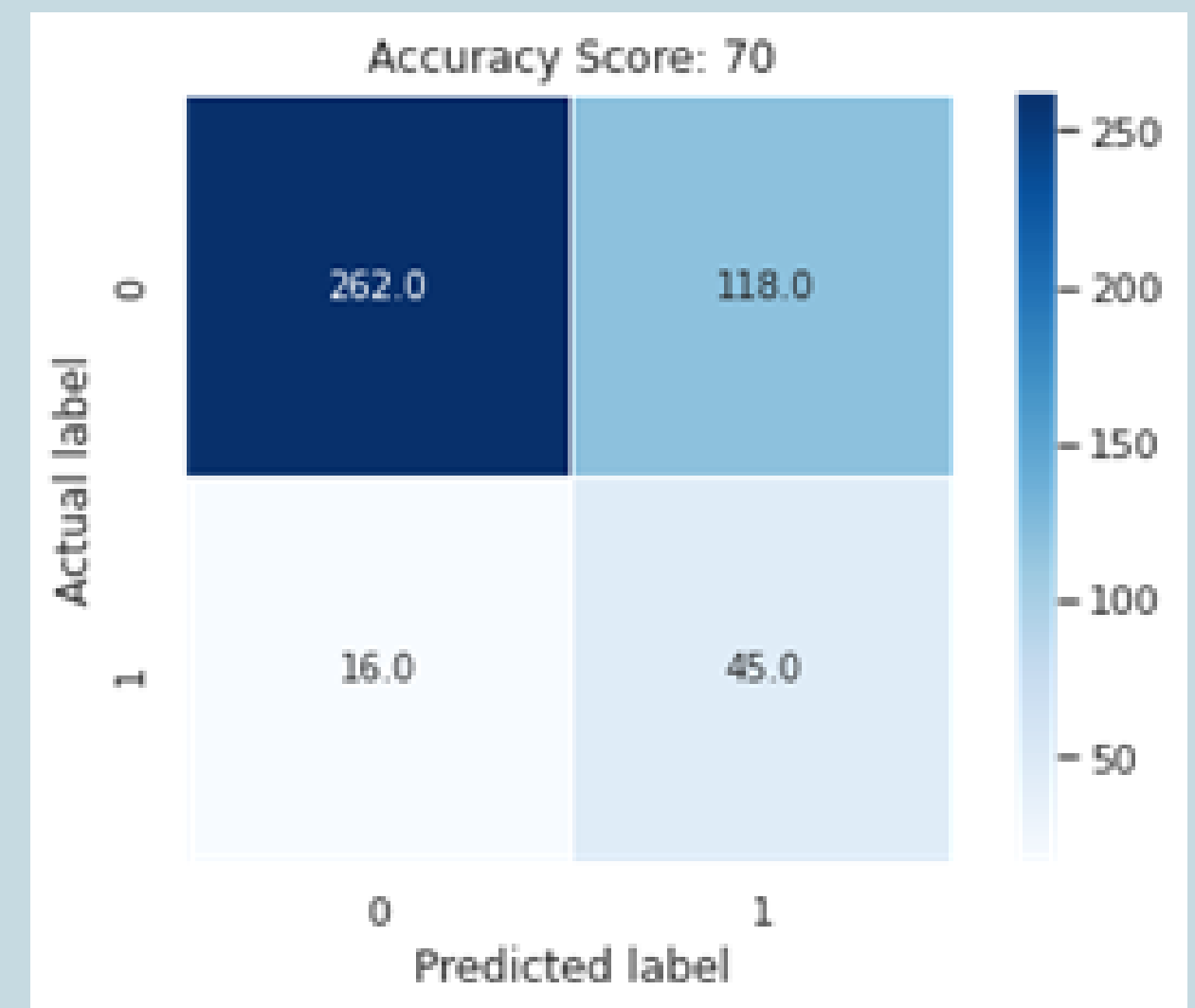
Accuracy Score: 86

Accuracy Score: 70

**Decision Tree**

**Confusion matrix before and after sampling**

**Confusion matrix before and after sampling**

Accuracy Score: 88

Accuracy Score: 77

*AdaBoost*

**Confusion matrix before and after sampling**

# Models Comparison

| Model | Resample | Precision | Recall | F1-score |
|---|---|---|---|---|
| Logistic Regression | actual | 0.717949 | 0.459016 | 0.560000 |
| Logistic Regression | smote+enn | 0.270968 | 0.688525 | 0.388889 |
| Decision Tree | actual | 0.529412 | 0.147541 | 0.230769 |
| Decision Tree | smote+enn | 0.264901 | 0.655738 | 0.377358 |
| SVM | actual | 0.941176 | 0.262295 | 0.410256 |
| SVM | smote+enn | 0.276074 | 0.737705 | 0.401786 |
| AdaBoost | actual | 0.578947 | 0.360656 | 0.444444 |
| AdaBoost | smote+enn | 0.343511 | 0.737705 | 0.468750 |

# Models Comparison



**INSIGHT:** SVM before sampling has best precision score , however we cannot decide yet if it's the best model  or not

# Models Comparison



ROC-AUC

Legend:
- Logistic, AUC=0.7863244176013805
- SVM, AUC=0.783304572907679
- AdaBoost, AUC=0.8057377049180329
- Decesion Tree, AUC=0.7169542709232097

# *Optimization*

## OPTIMIZING SVM

- Using GridSearchCV
- Best parameters were: C=1, Gamma=1, Kernel=RBF

# Models Comparison: after optimizing



Model = SVM-Optimized

actual | smote+enn

**INSIGHT:** the precision decreased in actual and increased in sampled
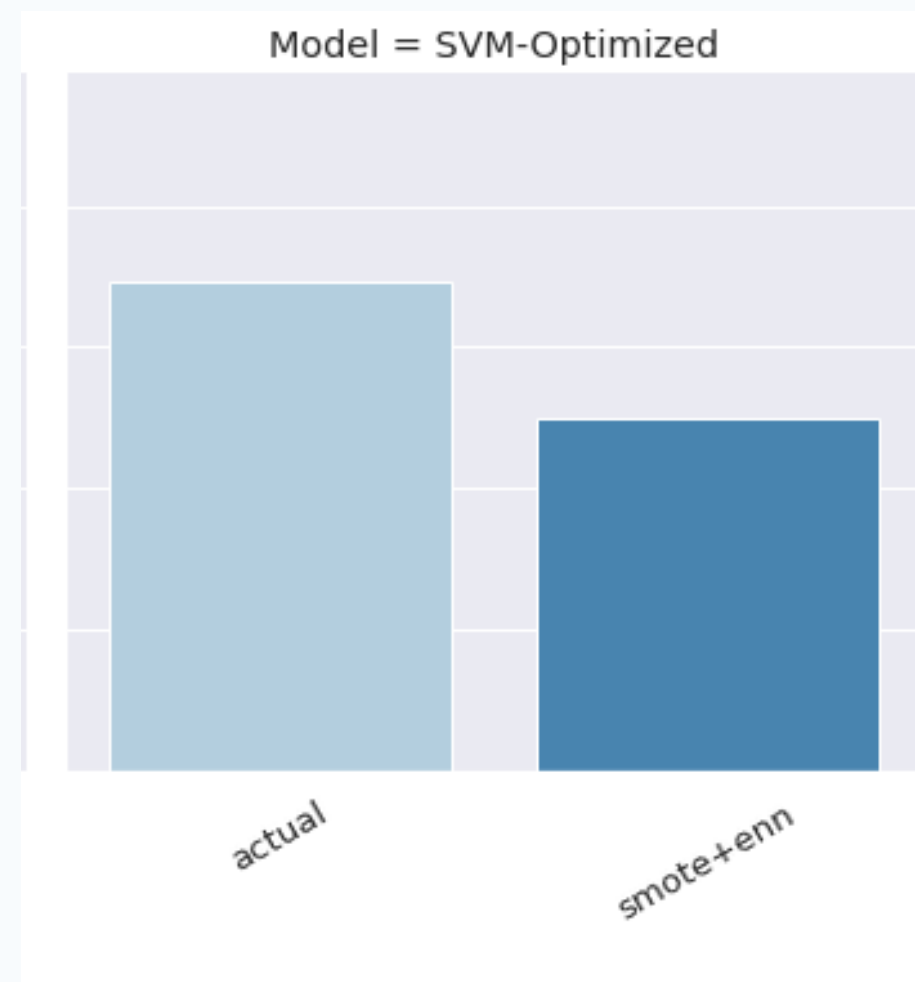
# Models Comparison: after optimizing

| Model | Resample | Precision | Recall | F1-score |
|---|---|---|---|---|
| Logistic Regression | actual | 0.717949 | 0.459016 | 0.560000 |
| Logistic Regression | smote+enn | 0.270968 | 0.688525 | 0.388889 |
| Decision Tree | actual | 0.529412 | 0.147541 | 0.230769 |
| Decision Tree | smote+enn | 0.264901 | 0.655738 | 0.377358 |
| SVM | actual | 0.941176 | 0.262295 | 0.410256 |
| SVM | smote+enn | 0.276074 | 0.737705 | 0.401786 |
| AdaBoost | actual | 0.578947 | 0.360656 | 0.444444 |
| AdaBoost | smote+enn | 0.343511 | 0.737705 | 0.468750 |
| SVM-Optimized | actual | 0.692308 | 0.442623 | 0.540000 |
| SVM-Optimized | smote+enn | 0.500000 | 0.311475 | 0.383838 |

# *Conclusions*

- SVM AFTER OPTIMIZATION AND LOGISTIC REGRESSION MAY BE GOOD MODELS FOR OUR CASE

- ADABOOST CLASSIFIER HAS THE LARGEST AUC AMONG ALL MODELS WE CAN SAY THAT ITS ALSO A GOOD CLASSIFIER

- DECESION TREE SEEM TO HAVE BEEN OVERFITTED DUE TO LOW VOLUME DATA IF THE DATA WAS LARGER TREES WOULD PERFORM BETTER

# Thank you!

ANY QUESTIONS?