



أكاديمية سدايا
SDAIA Academy

Comparison of classifier Algorithms on bank marketing Dataset

Presented by :
Arwa AlBassam & Nouf Alsaeed

10/7/2021

OUTLINES

01

Introduction

02

Workflow

03

Dataset Description

04

EDA

05

Baseline Model

06

Pre-processing

07

Modelling

08

Results

09

Conclusion

10

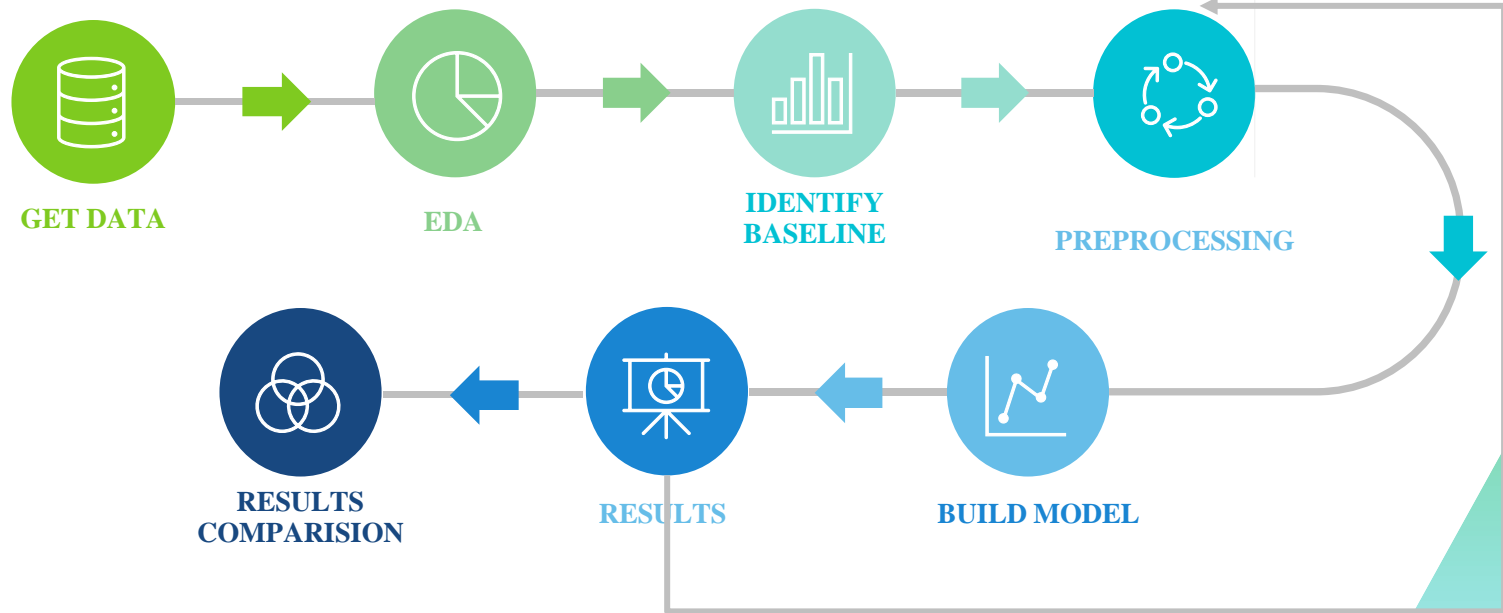
Future work

Introduction

The data is related with direct marketing campaigns (phone calls) of a Portuguese banking institution. The classification goal is to predict if the client will subscribe a term deposit (variable y).



Workflow



Dataset Description

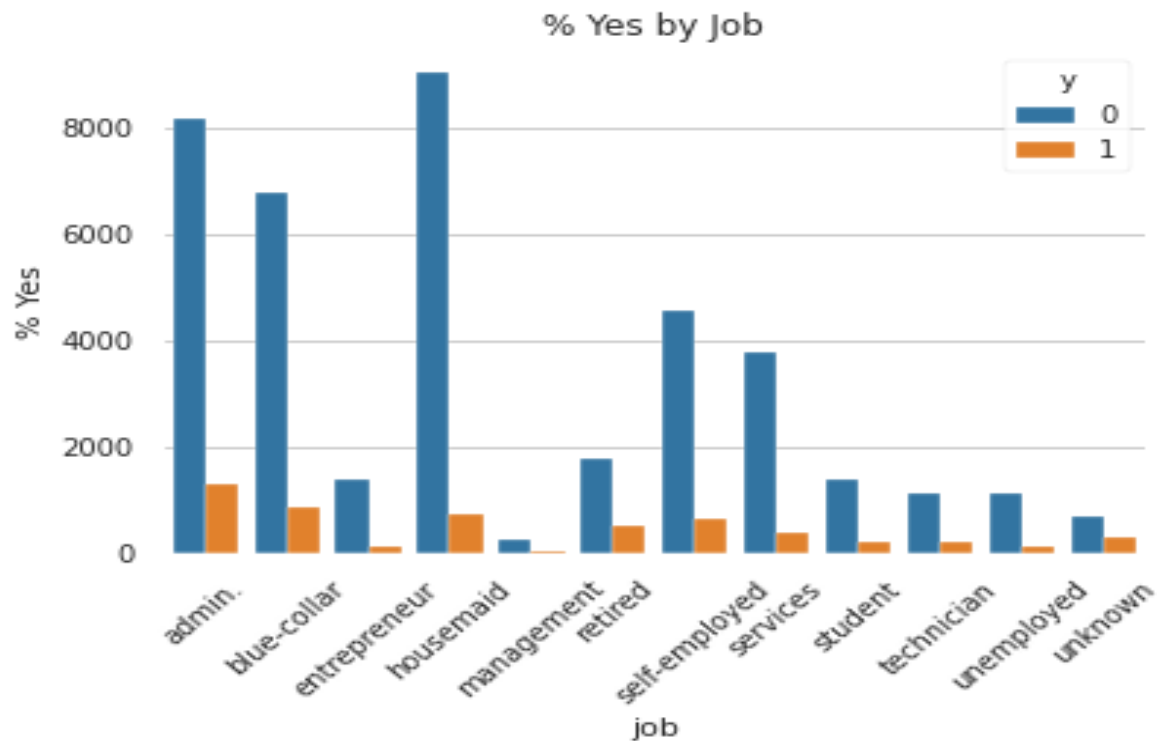
From UCI

45211 Rows X 17 Columns

Feature	Description
Age	Age of clients
Job	Type of job
Marital	Marital status
Education	(Categorical: "unknown", "secondary", "primary", "tertiary")
Balance	Average yearly balance, in euros (numeric)
Housing	Has housing loan? (Binary: "yes", "no")
Loan	Has personal loan? (Binary: "yes", "no")
Poutcome	Outcome of the previous marketing campaign
Y	Has the client subscribed to a term deposit? (Binary: "yes", "no")

EDA

Visualizations



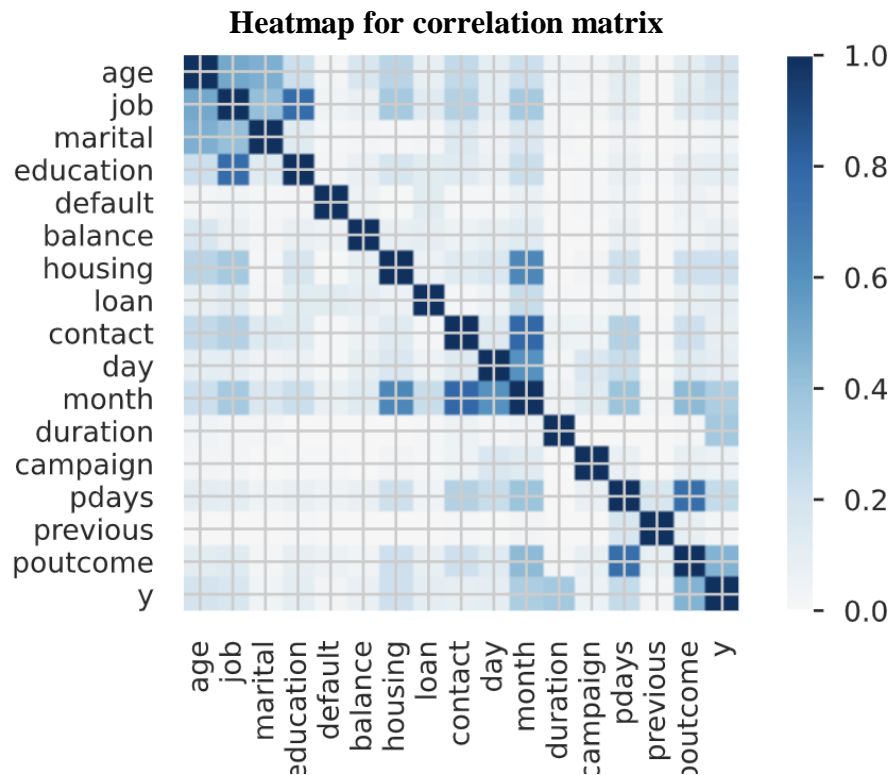
EDA (Cont.)

Dataset statistics

Number of variables	17
Number of observations	45211
Missing cells	0
Missing cells (%)	0.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	29.2 MiB
Average record size in memory	677.2 B

EDA (Cont.)

Visualizations



Base Model (Random Forest)

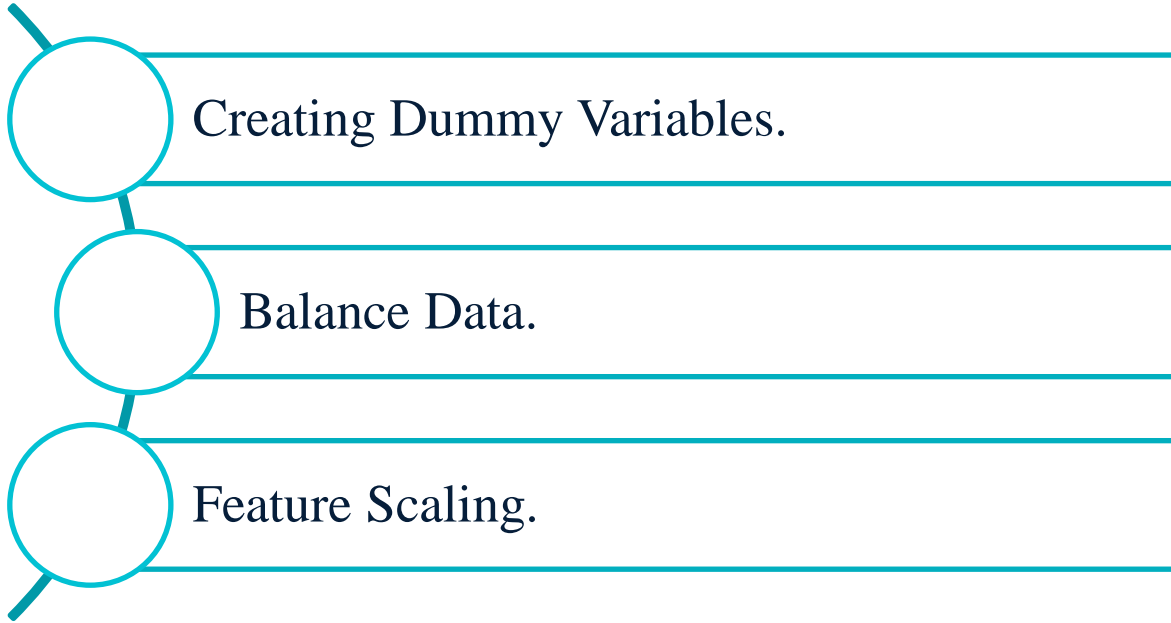


ACCURACY



RECALL

Data Preprocessing

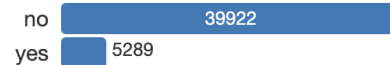


Data Preprocessing (cont.)

Balancing the Model

y
Boolean

Distinct count	2
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Memory size	353.3 KiB



Toggle details

Frequency Table

Value	Count	Frequency (%)
no	39922	88.3%
yes	5289	11.7%

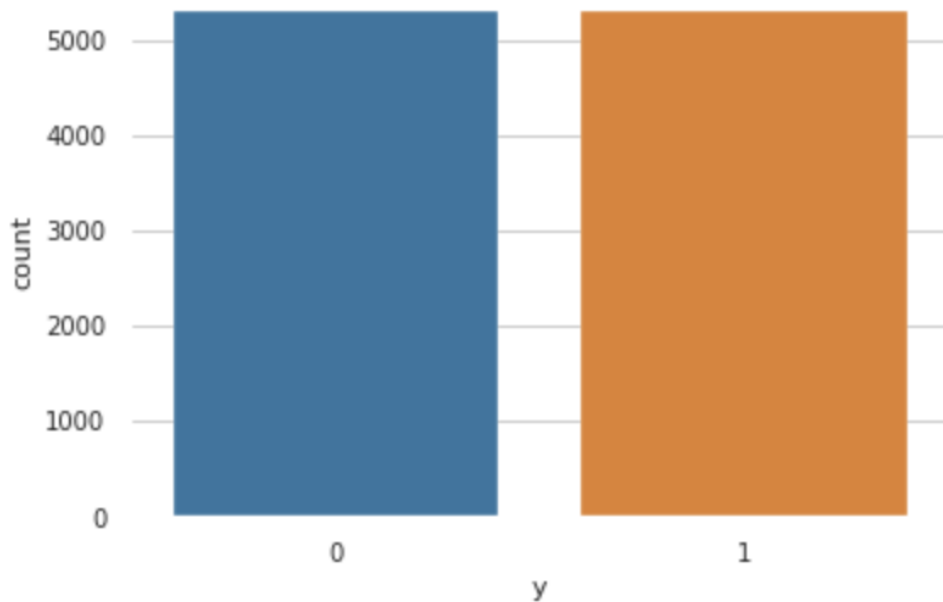
Data Preprocessing (cont.)

Balancing the Model

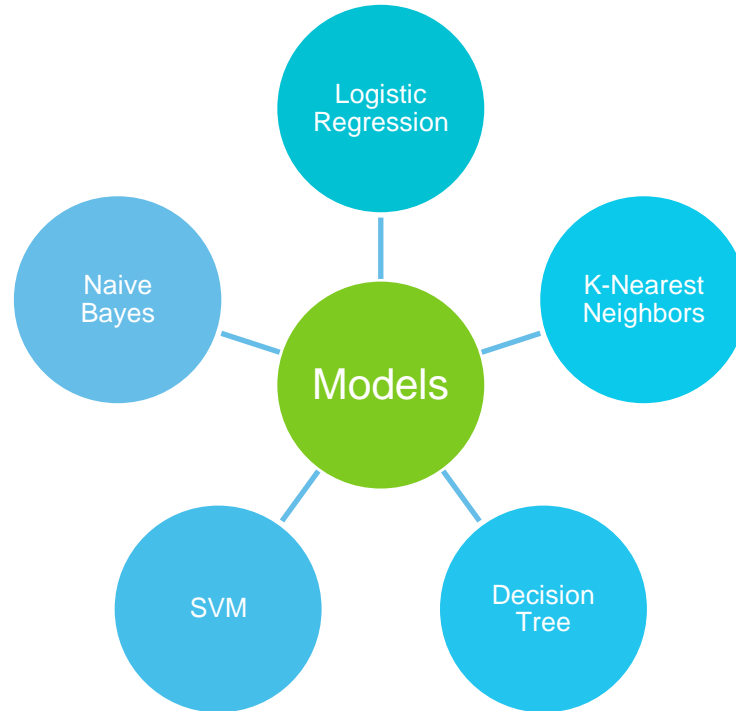
```
1    5289
```

```
0    5289
```

```
Name: y, dtype: int64
```

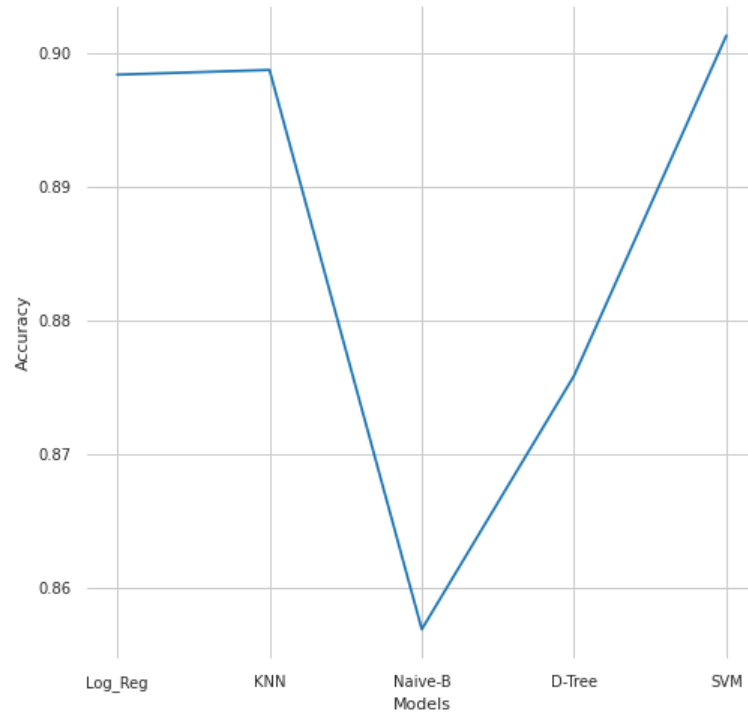


Modelling

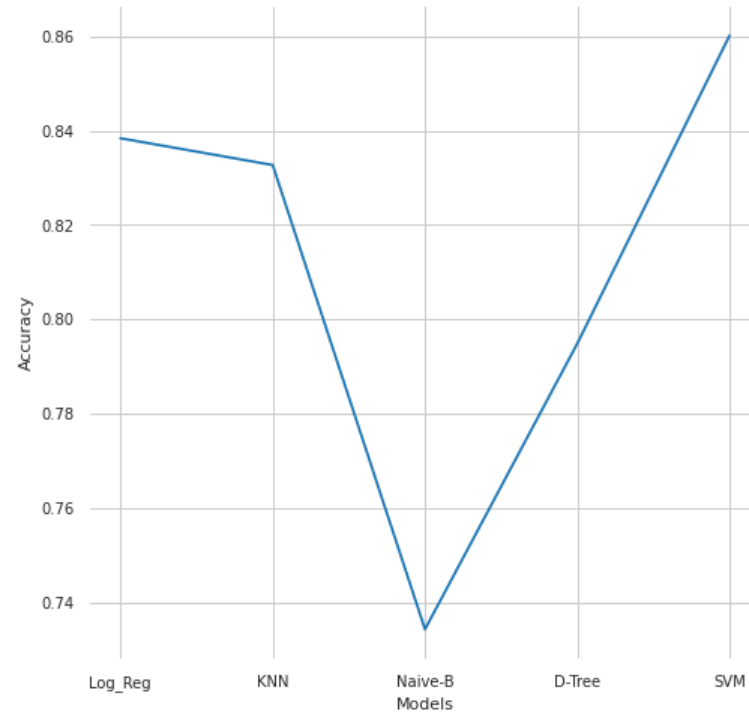


Models Evaluation

Accuracy Before Balancing The Data

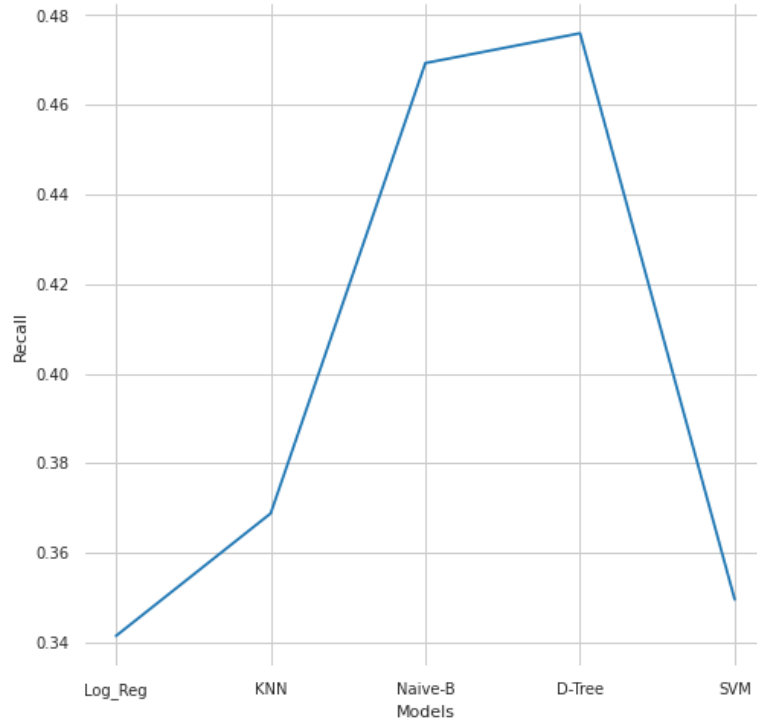


Accuracy After Balancing The Data

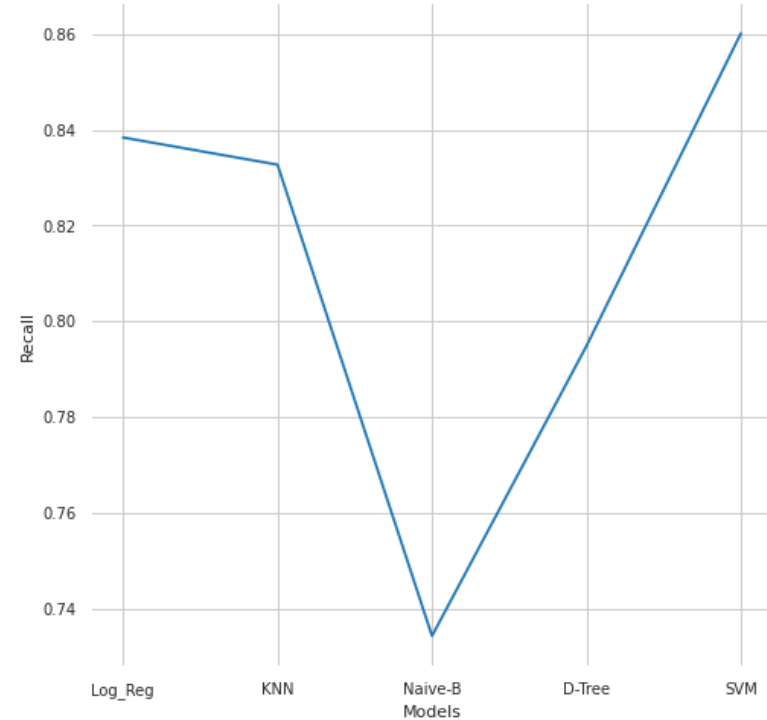


Models Evaluation

Recall Before Balancing The Data

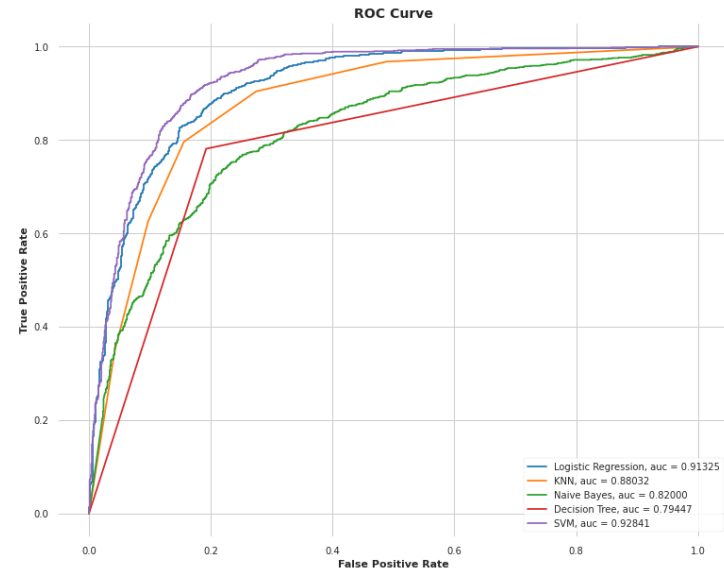
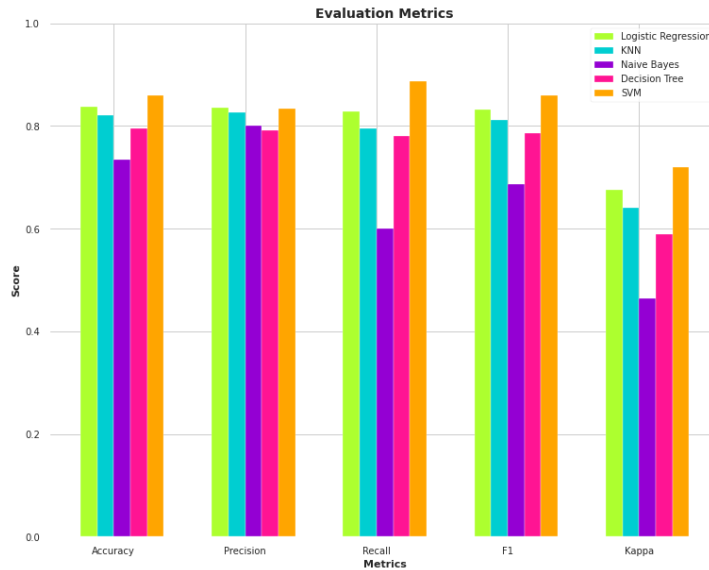


Recall After Balancing The Data

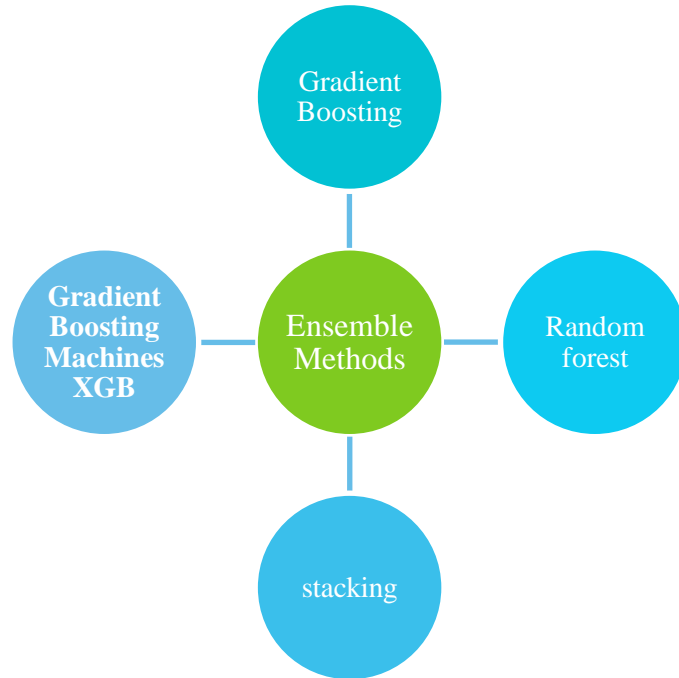


Models Evaluation

Model Comparison



Ensemble Methods



Stacking Ensemble Family

Average Voting



ACCURACY: 86.48 %

Max Voting



ACCURACY: 86.24%

Weighted Average Ensemble



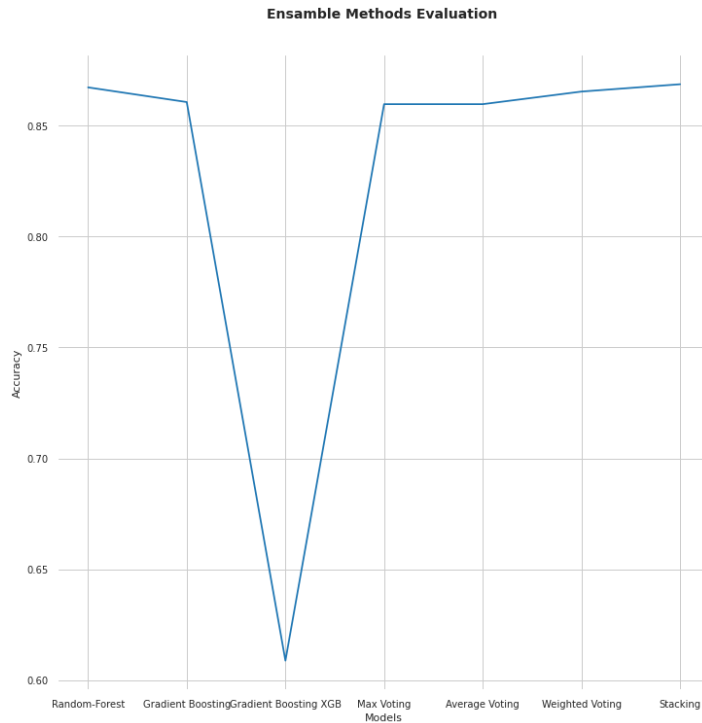
ACCURACY: 86.06%

Stacking



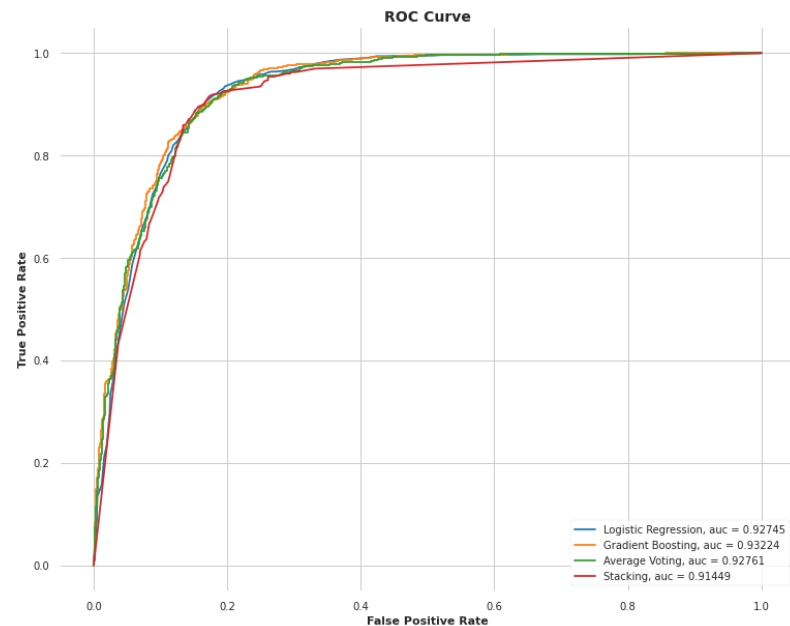
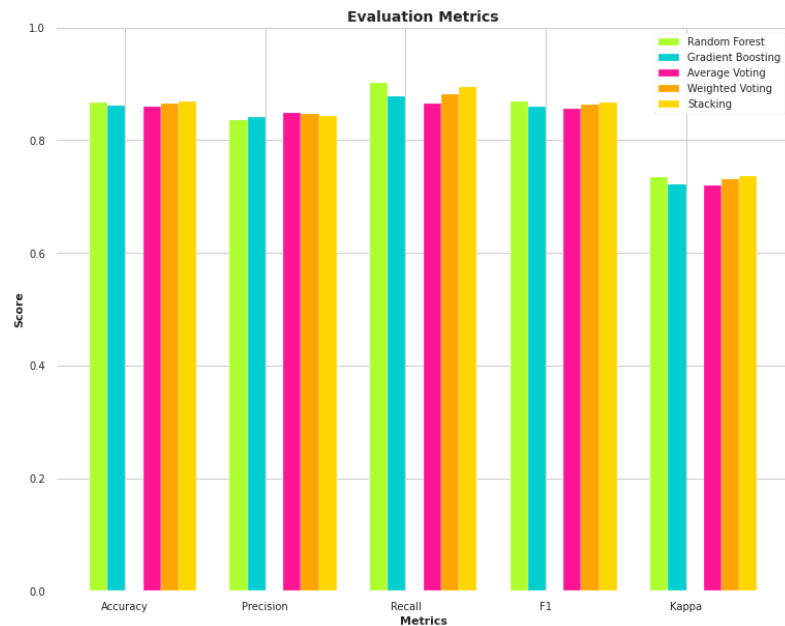
ACCURACY: 86.7 %

Ensemble Methods Evaluation

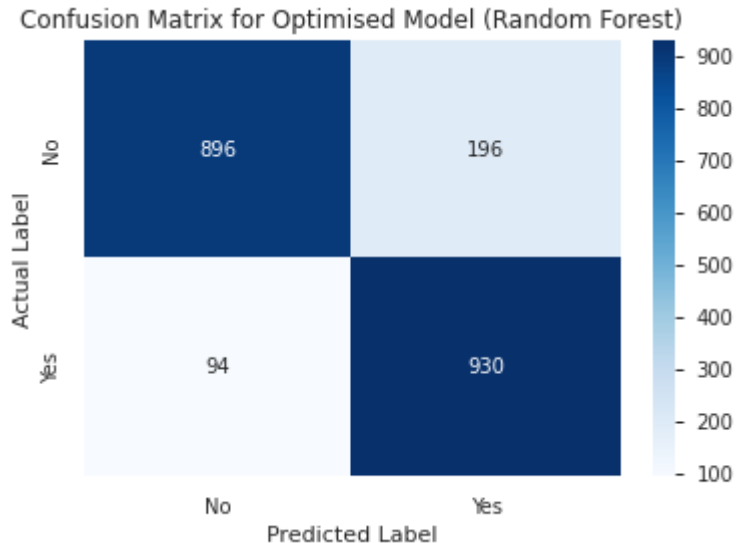


Ensemble Methods Evaluation

Model Comparison



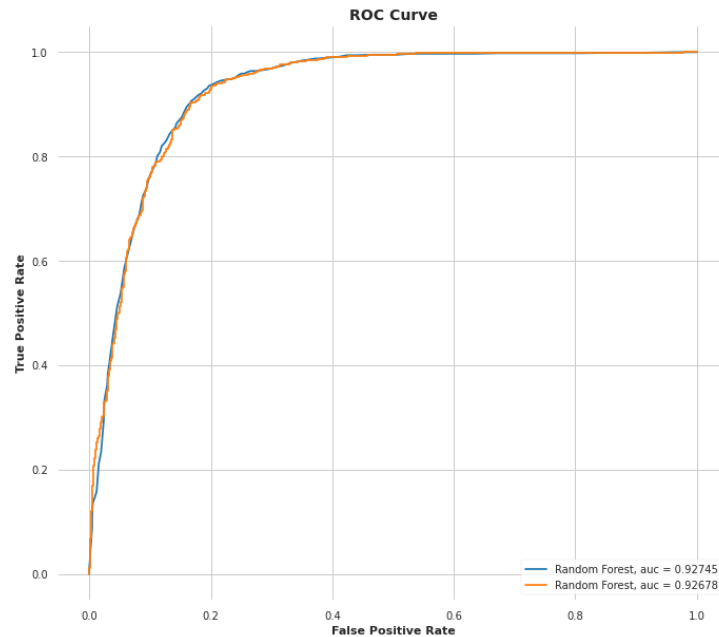
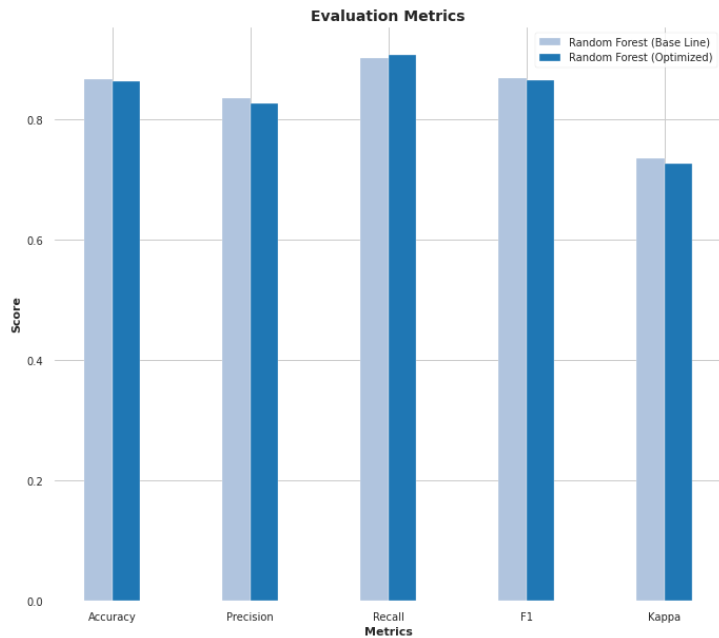
Best model for our project (Random Forest)



ACCURACY: 86.8%
RECALL: 90.23%

Final results (Model Optimization)

Model Comparison



Conclusion

The result is not that much different after optimizing the model using GridSearchCV which can mean that we hit our limit with this model.

Future Work

Try to use other balancing techniques.

Try to work with Deep Learning models.

Increase number of observations.

Feature engineering.



Thank you!