



Citibike Classification

Alexander Sigrist

Found 2013

- Bike Sharing System
- Fleet of 40,000
- Launched e-bikes in 2019



- Identify a model that can effectively detect usertype

Subscribers vs Customers

- What's the difference?
- Why is it useful?

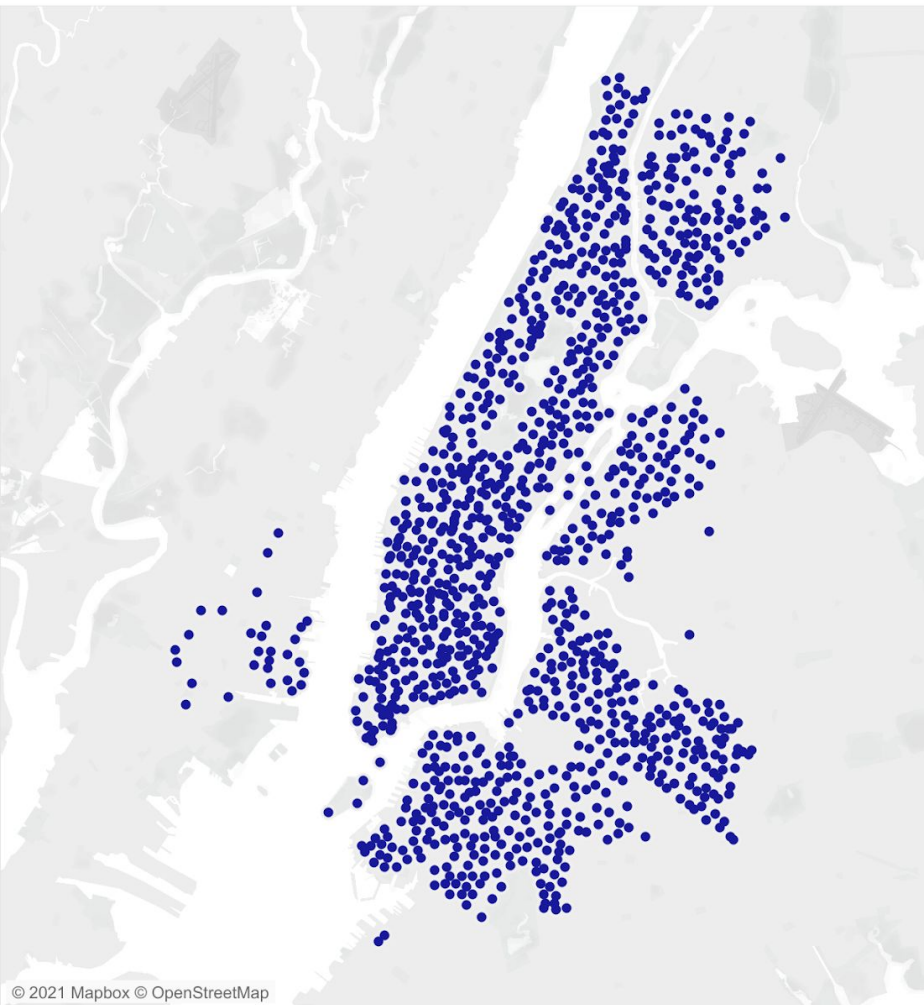
Data Overview

- Citibike Data - Summer 2020
 - Target - Usertype
- Features Anomalies
 - Outliers
 - Class Imbalance

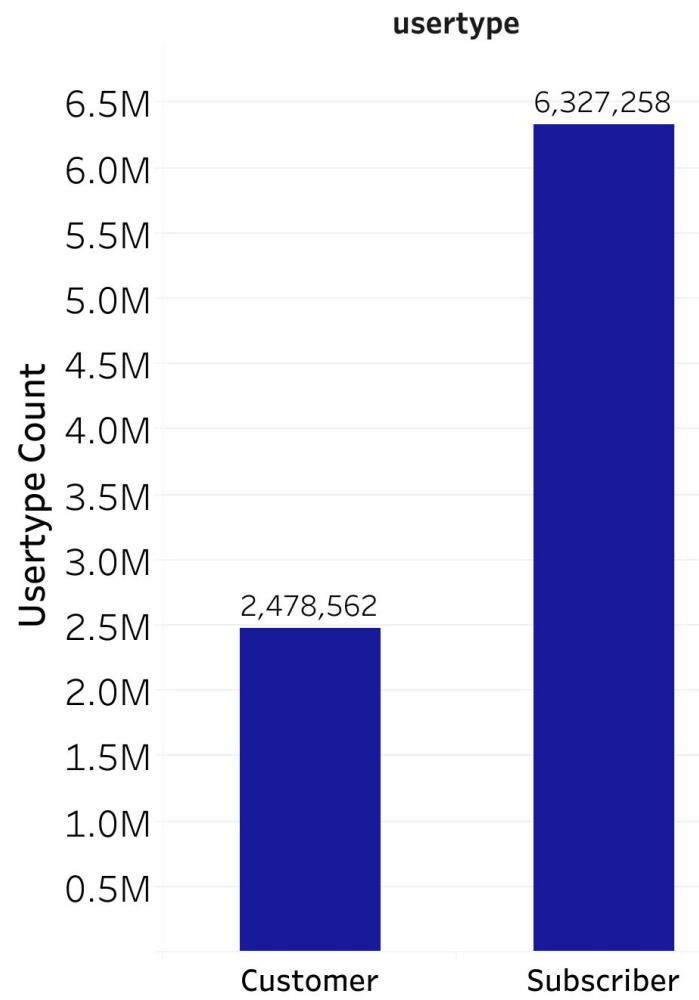
Features

- Trip duration
- Gender
- Month
- Day of the week
- Birth Year
- Day
- Starttime
- Stoptime
- End station latitude & longitude
- Start station latitude & longitude

Citibike Locations 09/2020



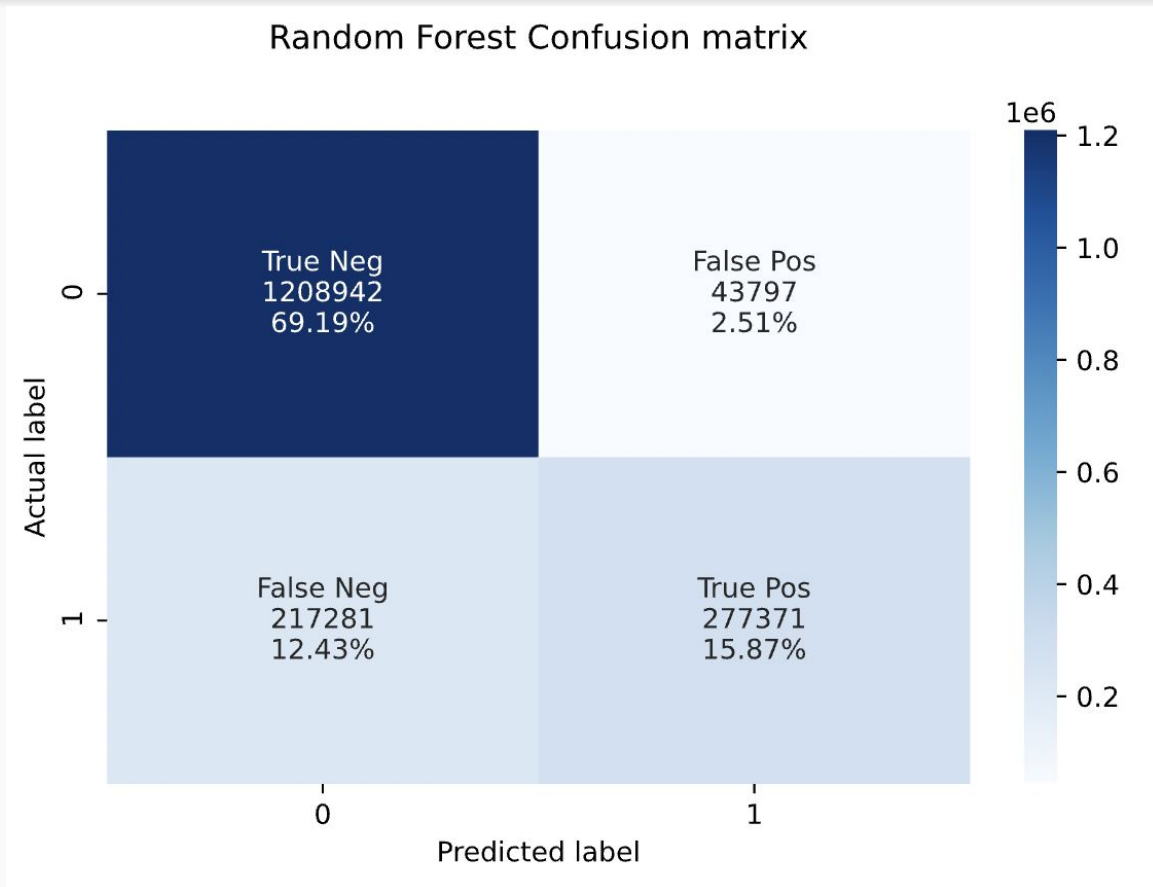
Class Imbalance of 1 : 3 for Usertype



Models & Metrics

Models Tested	Logistic Regression	Naive Bayes Gaussian	Random Forest	XGBoost*
F1 Score:	.523	.573	.68	.662

Models & Metrics



0 = Subscriber

1 = Customer

Models & Metrics

Models Tested	Logistic Regression	Naive Bayes Gaussian	Random Forest	XGBoost*
F1 Score:	.523	.573	.68	.662
Beta = 0.25				

Models & Metrics

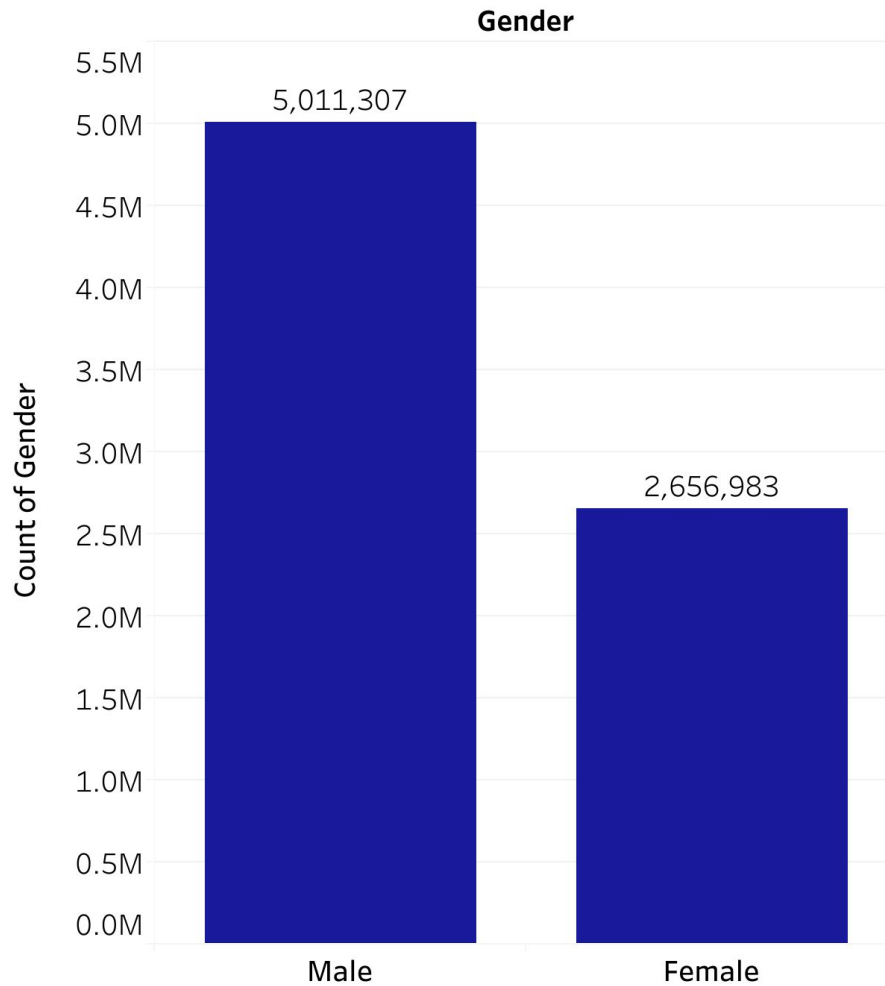
Models Tested	Logistic Regression	Naive Bayes Gaussian	Random Forest	XGBoost*
F1 Score:	.523	.573	.68	.662
Beta = 0.25				
Fbeta Score:	.757	.811	.837	.823

Next Steps

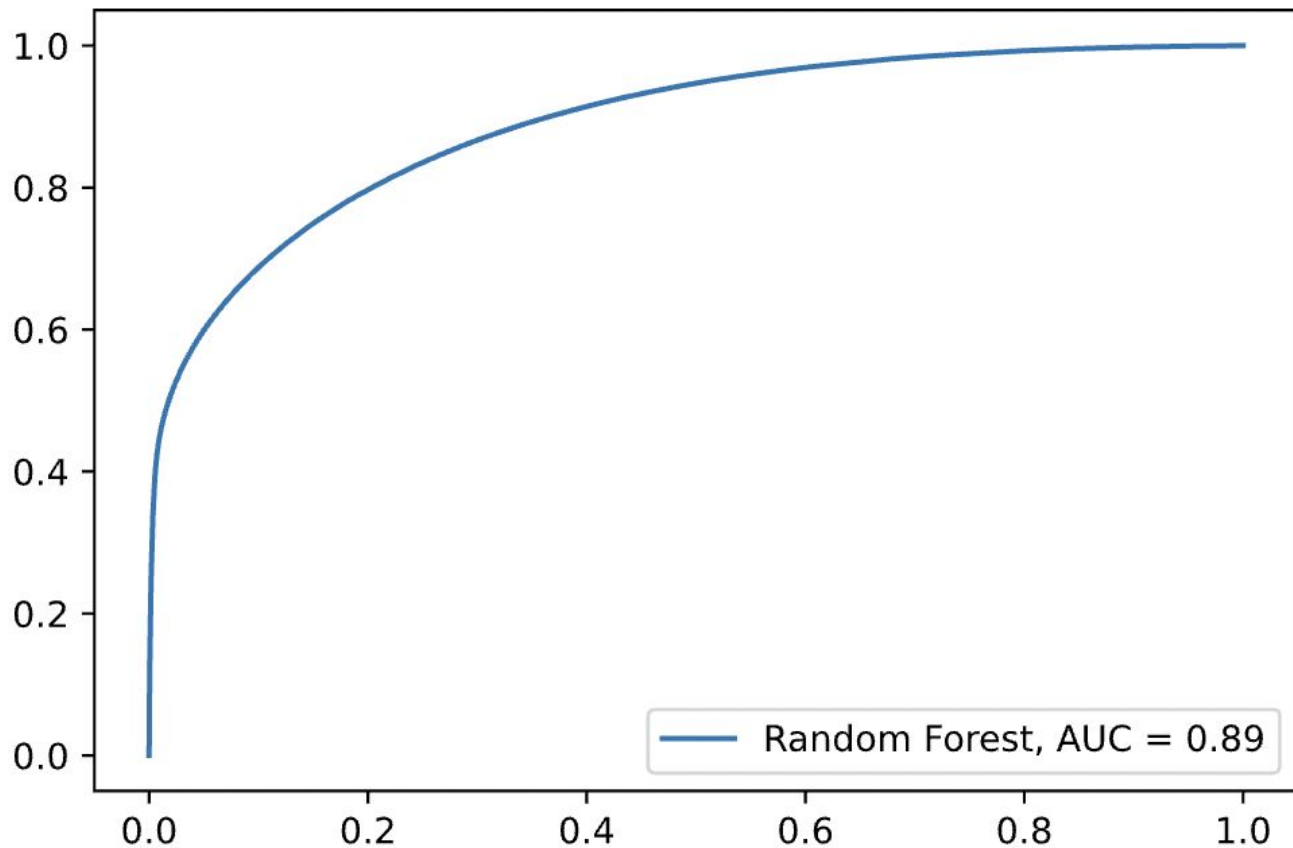
- Examine XGBoost - tune hyperparameters
- Undersample dataset

Appendix

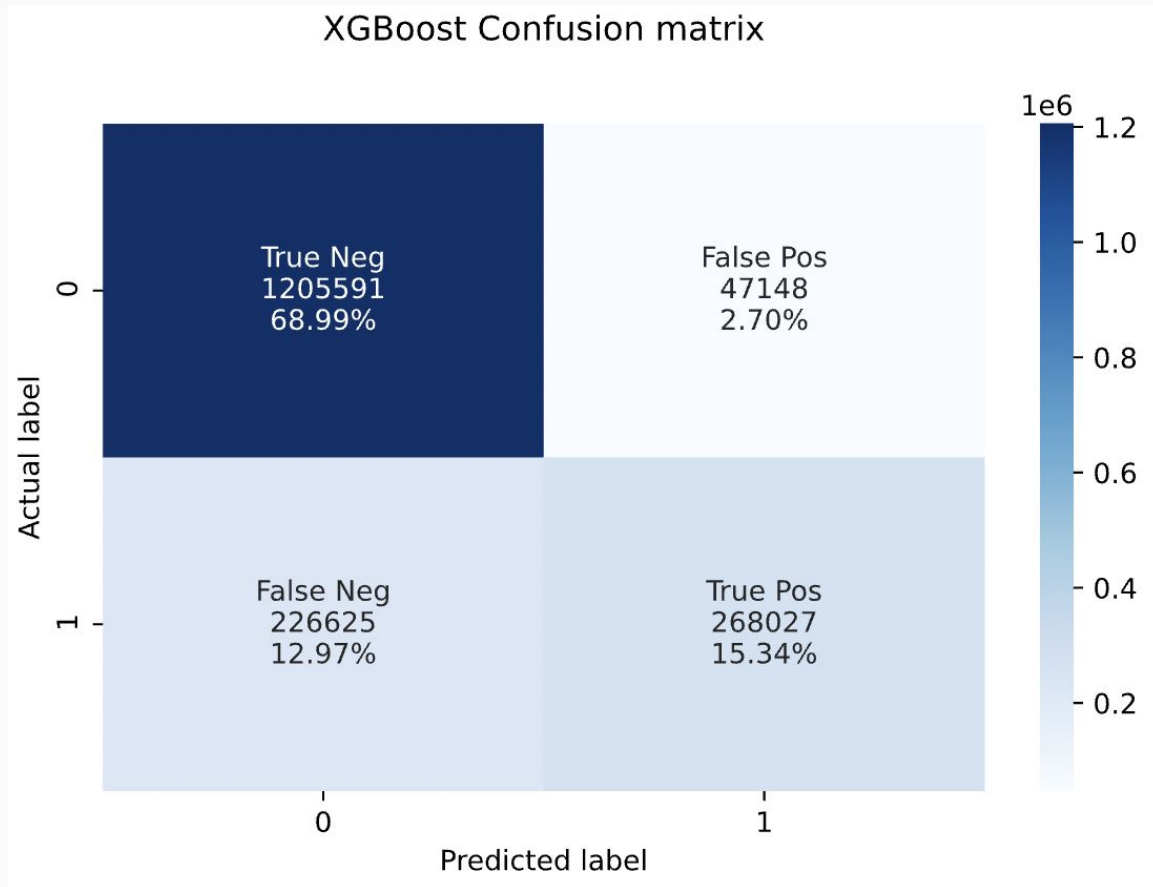
Gender Distribution is 2 : 1



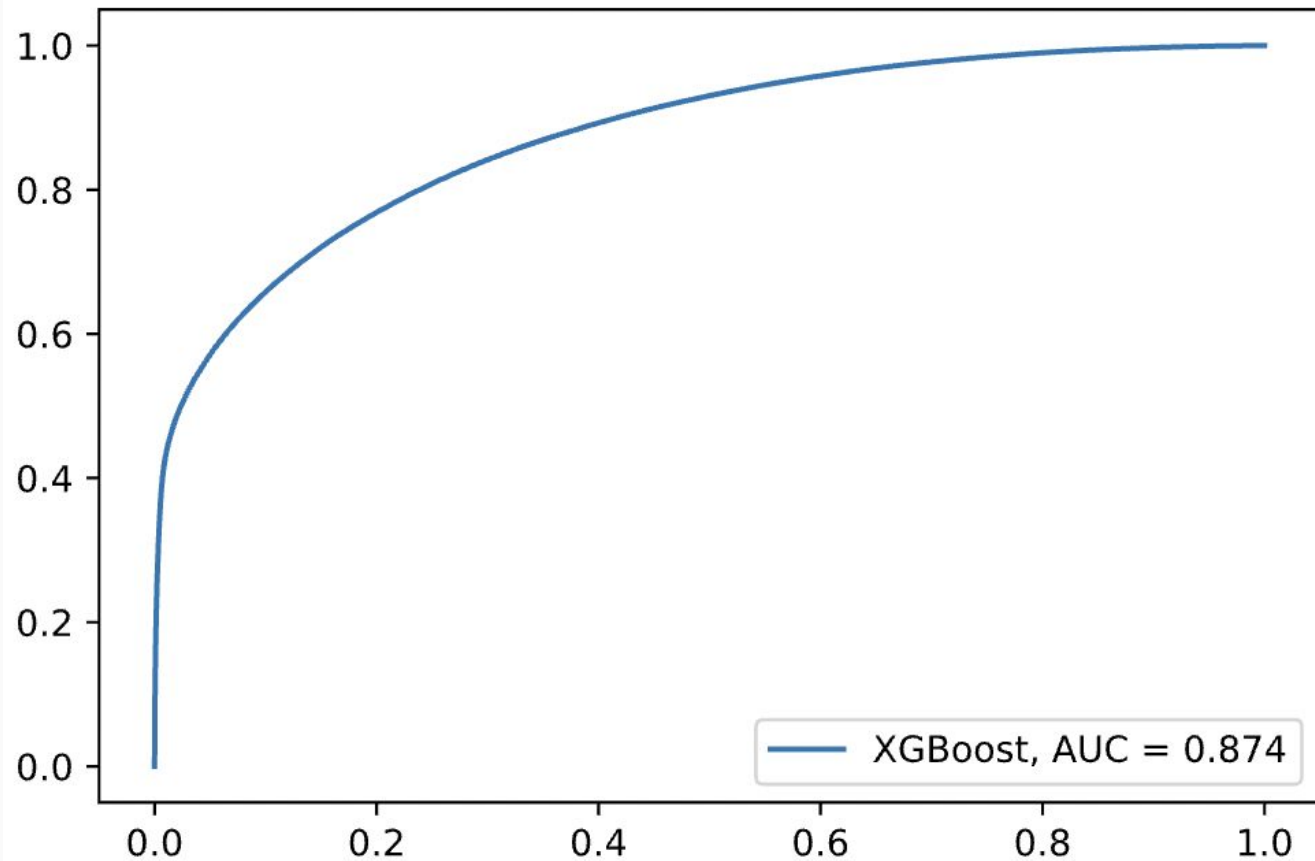
Appendix



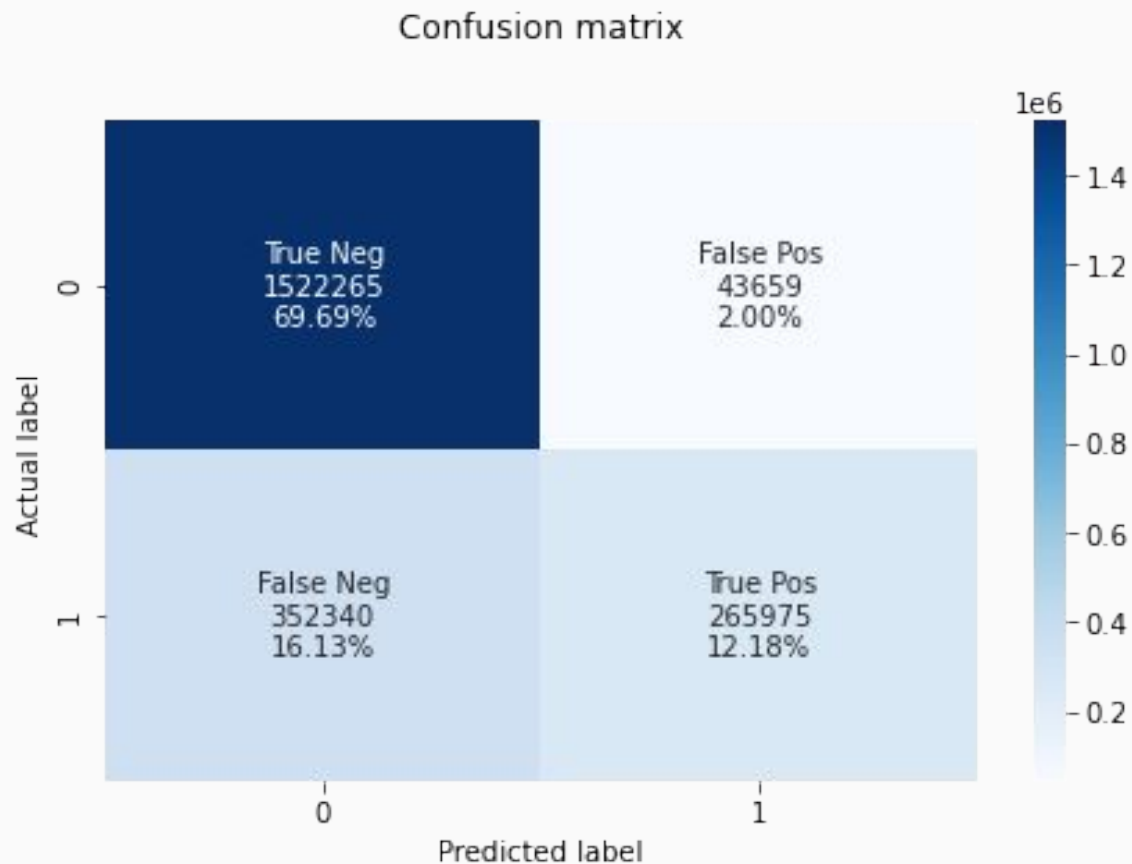
Appendix



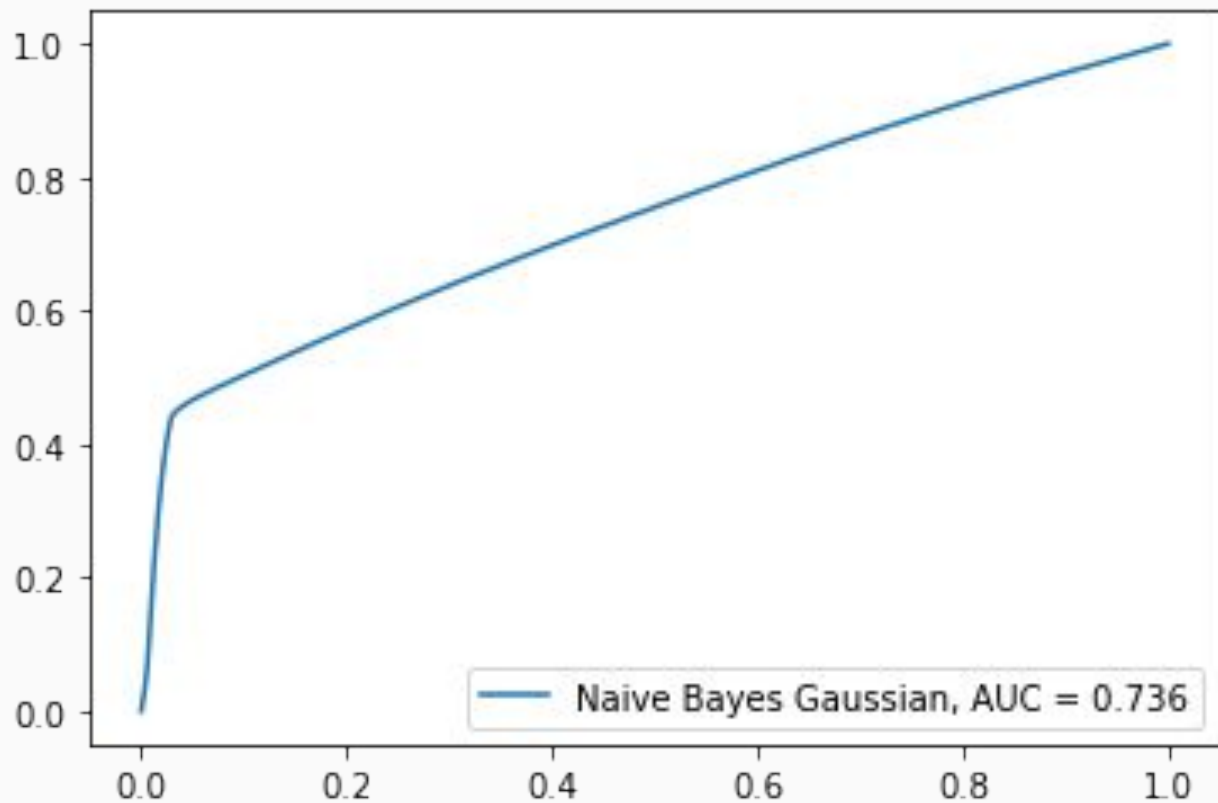
Appendix



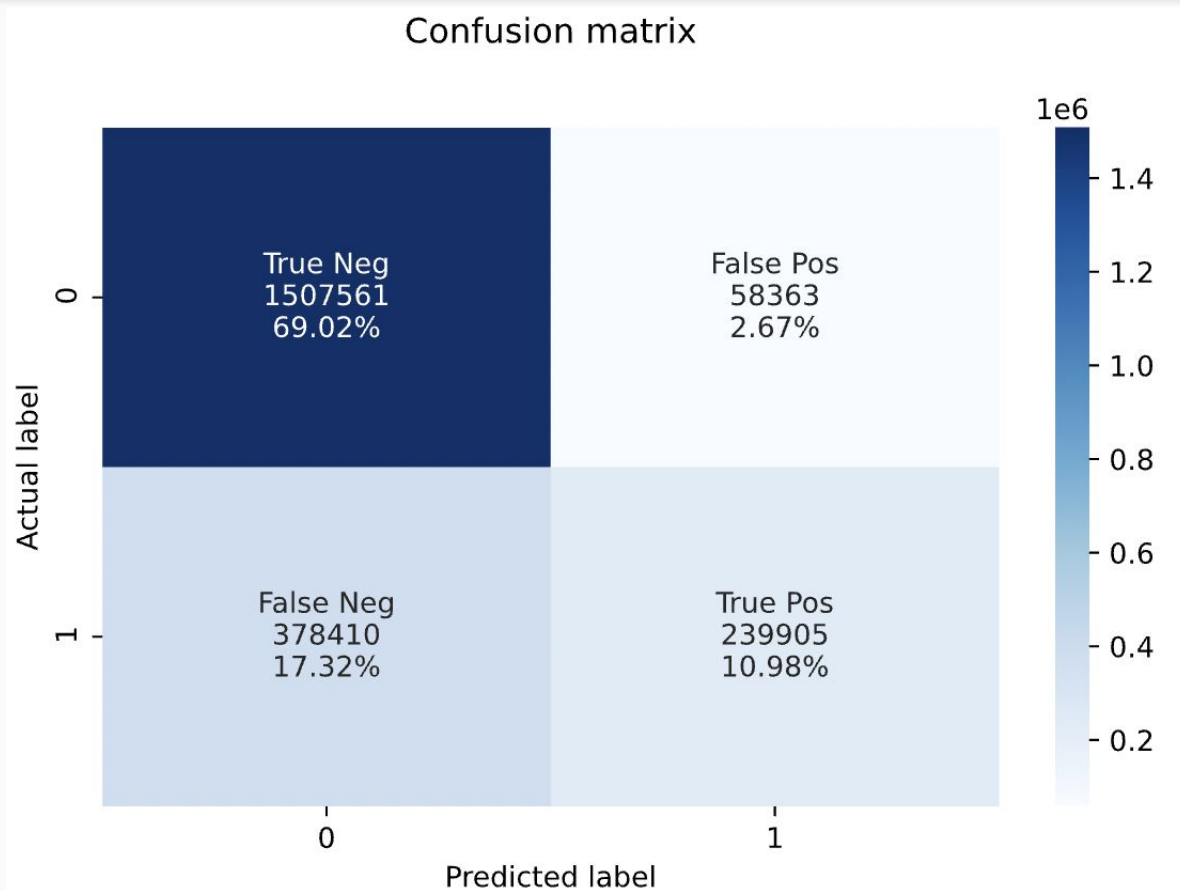
Appendix



Appendix



Appendix



Appendix

