

The background of the slide is composed of several overlapping triangles in shades of blue and red. The top-left area is dominated by various shades of blue, while the bottom-left and bottom-right areas feature shades of red. The triangles vary in opacity, creating a layered, geometric effect.

# TELCO CUSTOMER CHURN PREDICTION


BY RAKA R.A PRASETYO



## BACKGROUND

Churn quantifies the number of customers who have left a brand by cancelling their subscription or stopping paying for the services. This is bad news for any business as it costs five times as much to attract a new customer as it does to keep an existing one. A high customer churn rate will hit a company's finances hard.

<https://www.invespcro.com/blog/customer-acquisition-retention/>

- 
- Utilize machine learning for predicting customer churn
  - Compare and choose the best machine learning model for this case
  - Model Evaluation using ROC AUC



## OBJECTIVE

# DATA INTRODUCTION

Source : <https://www.kaggle.com/yeanzc/telco-customer-churn-ibm-dataset>

Data related to a customer who left  
churn\_label, churn\_scores, churn\_value, churn reason

Data related to a customer subscribed service  
phone\_service, multiple\_line, internet\_service, online\_security, online\_backup,  
device\_protection, tech\_support, streaming\_tv, streaming\_movies

Data related to a customer demographic info  
country, state, city, zip\_code, lat\_long, latitude, longitude, gender, senior\_citizen,  
partner, dependents

Data related to a customer account information  
tenure\_month, contract, paperless\_billing, payment\_method, monthly\_charges,  
total\_charges

Data related to IBM company data  
customer\_id, count, cltv

# DATA CLEANSING

ASSIGN DATASET AS DF



CHECK DATA TYPES



ASSIGN NA AS NEW  
UNIQUE DATA



DROP UNUSED DATA

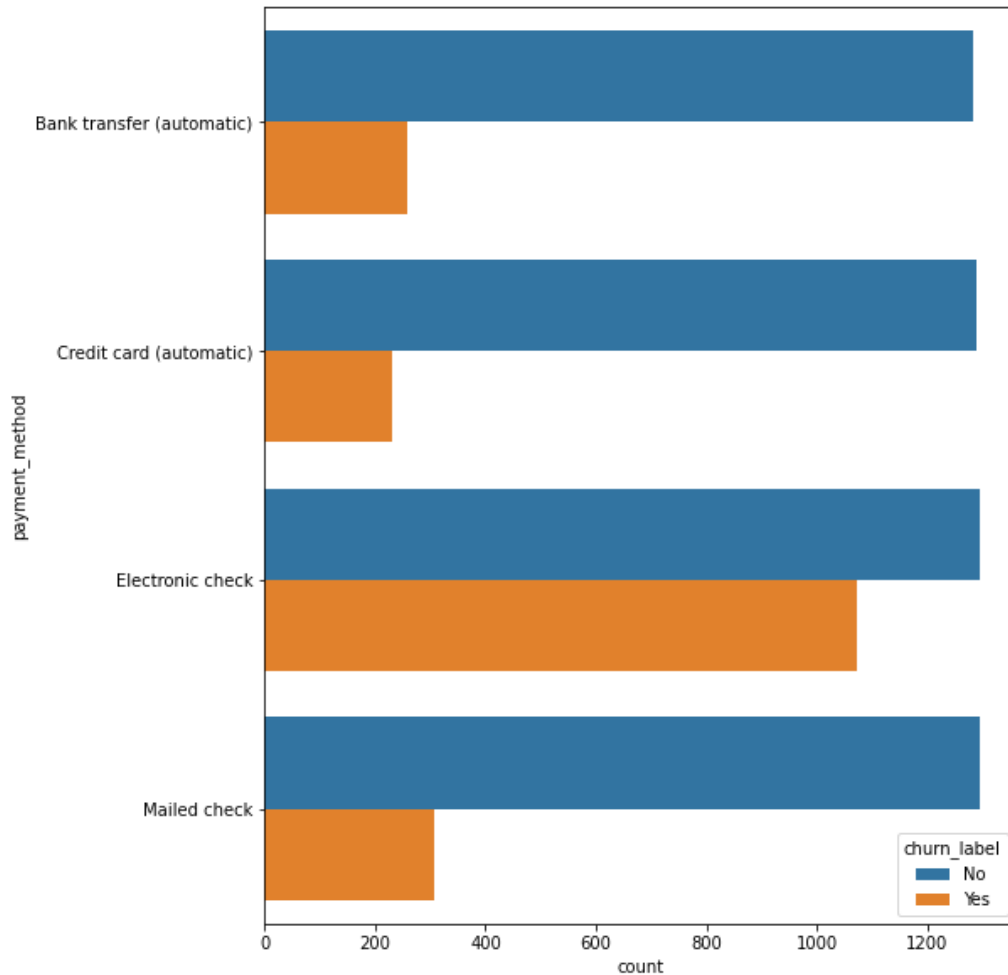


TRANSFORM THE DATA

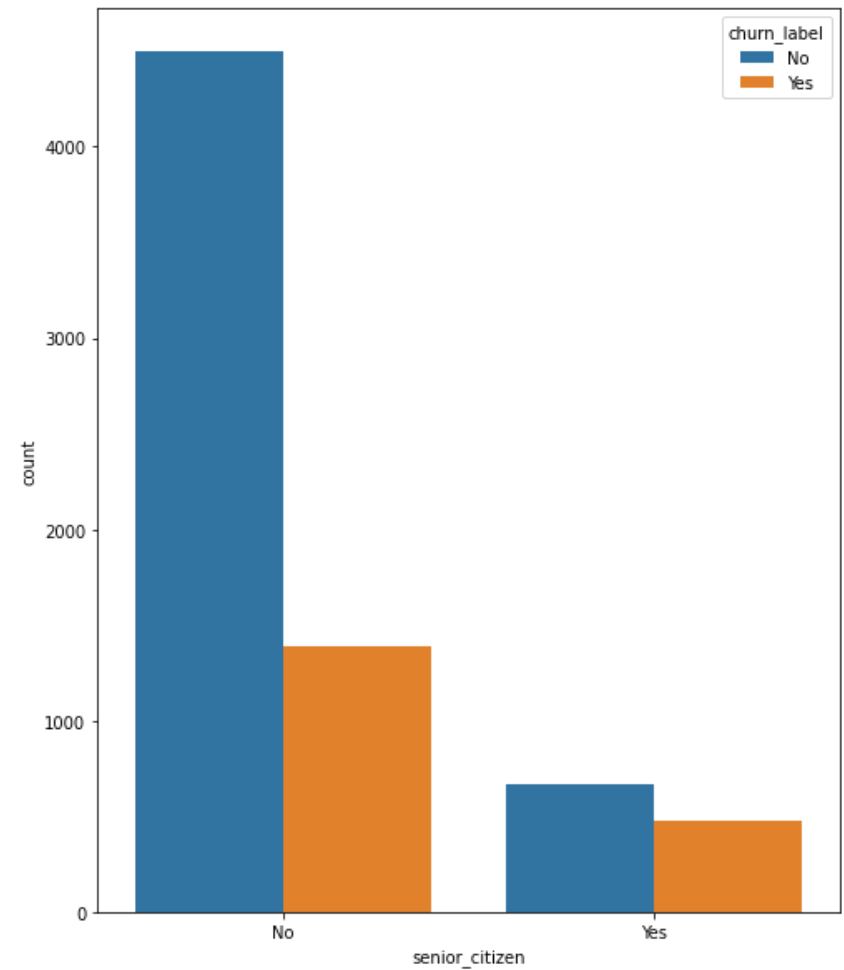
# EXPLORATORY DATA ANALYSIS

## CATEGORICAL VARIABLE

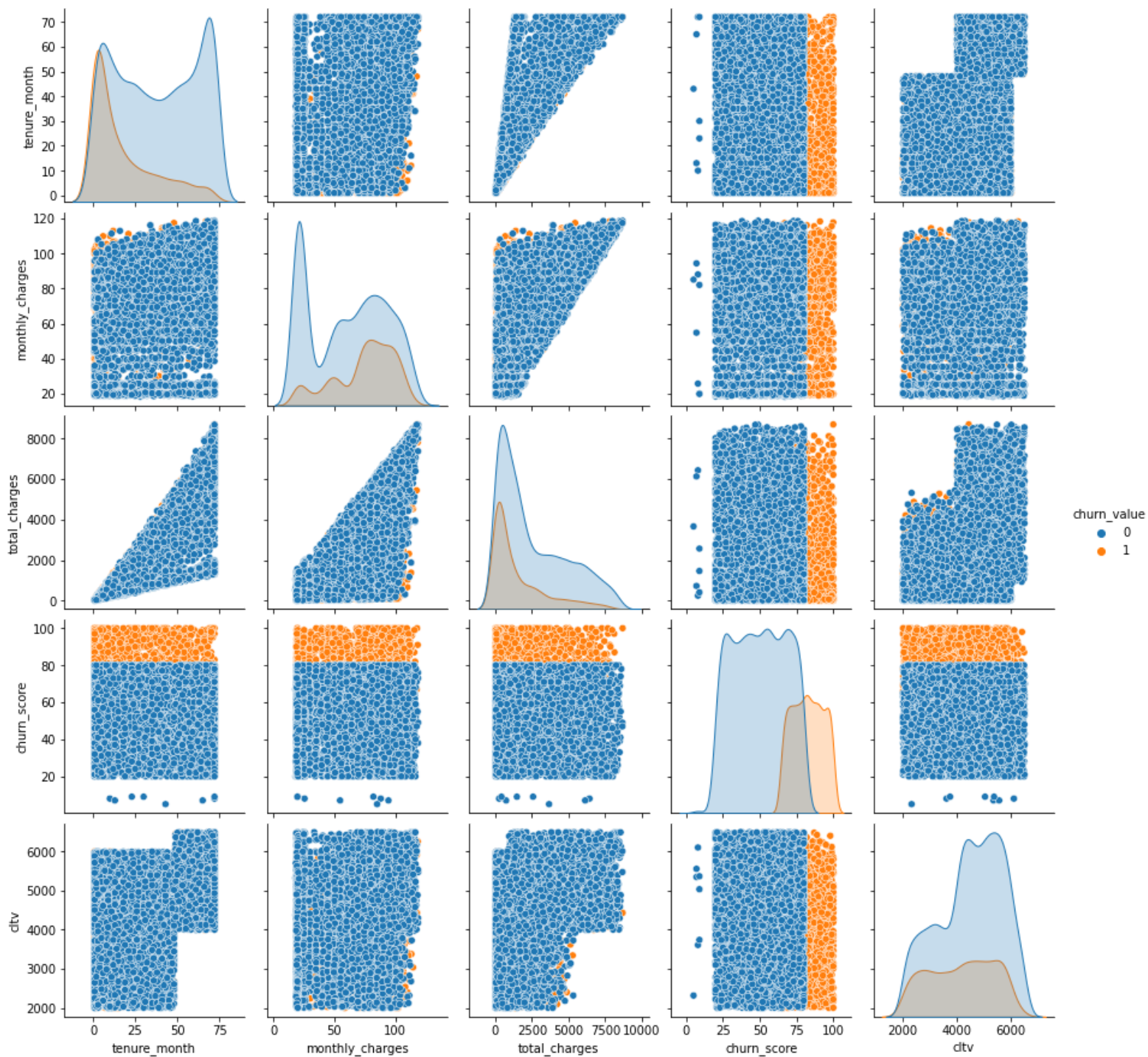
### PAYMENT METHOD



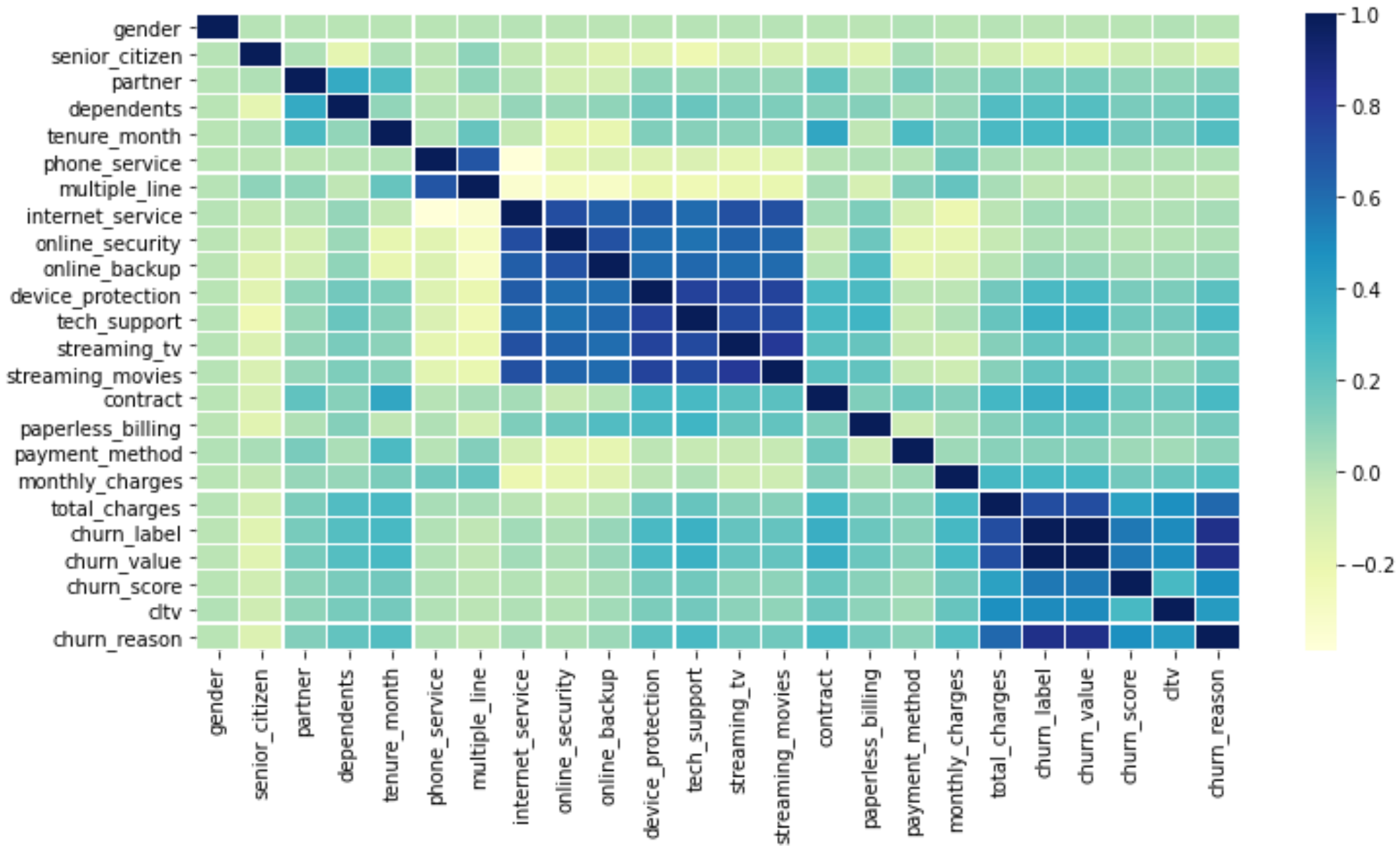
### SENIOR CITIZEN



# NUMERIC COLUMN



# CORRELATION HEATMAP

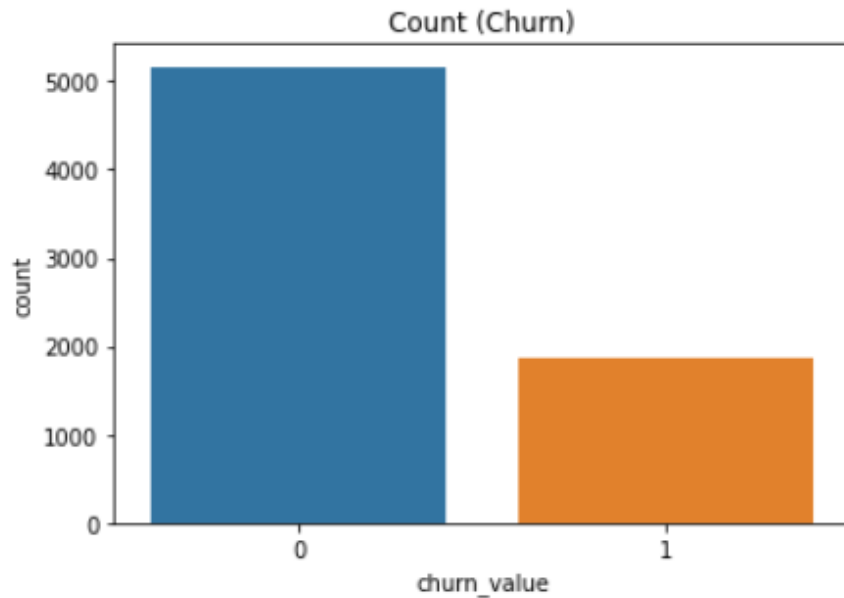


# RANDOM UNDERSAMPLING

Random Undersampling involves randomly selecting examples from the majority class to delete from the training dataset.

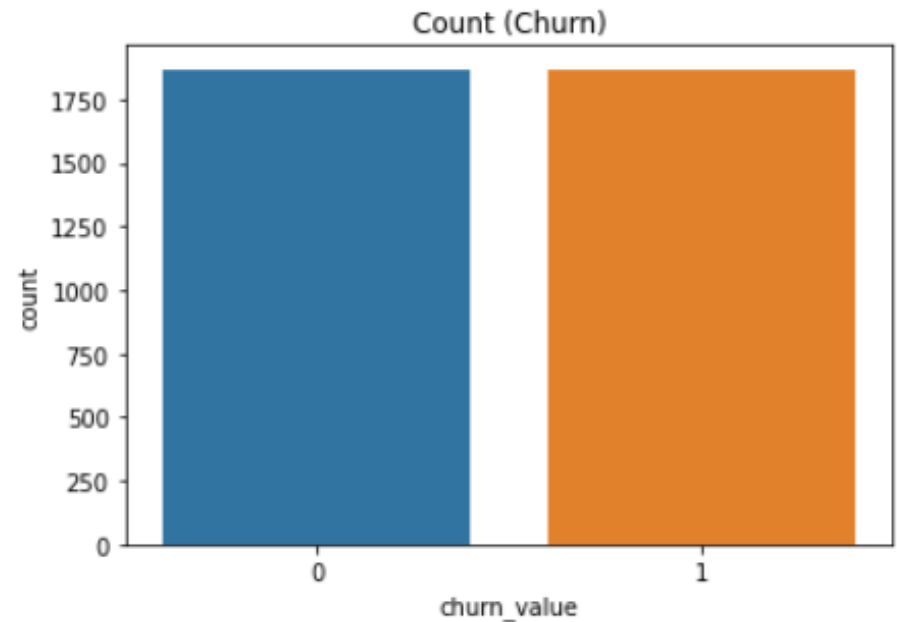
BEFORE

```
Class 0: 5163  
Class 1: 1869  
Proportion: 2.76 : 1  
Text(0.5, 1.0, 'Count (Churn)')
```



AFTER

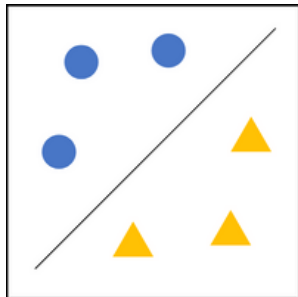
```
1 1869  
0 1869  
Name: churn_value, dtype: int64  
Text(0.5, 1.0, 'Count (Churn)')
```



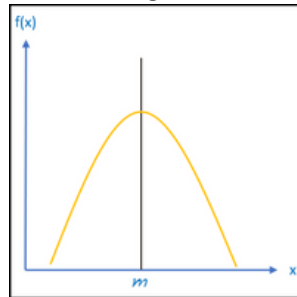


# MODELLING

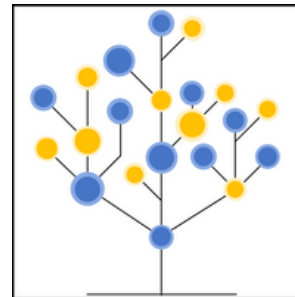
Support Vector  
Machine



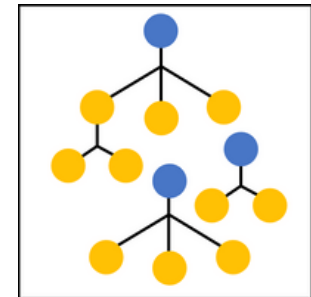
Gaussian Naive  
Bayes



Decision Tree  
Classifier



Random Forest  
Classifier



## Train Set

80 % Total Data

$X = 2990, 19$

$y = 2990, 1$

## Test Set

20 % Total Data

$X = 748, 19$

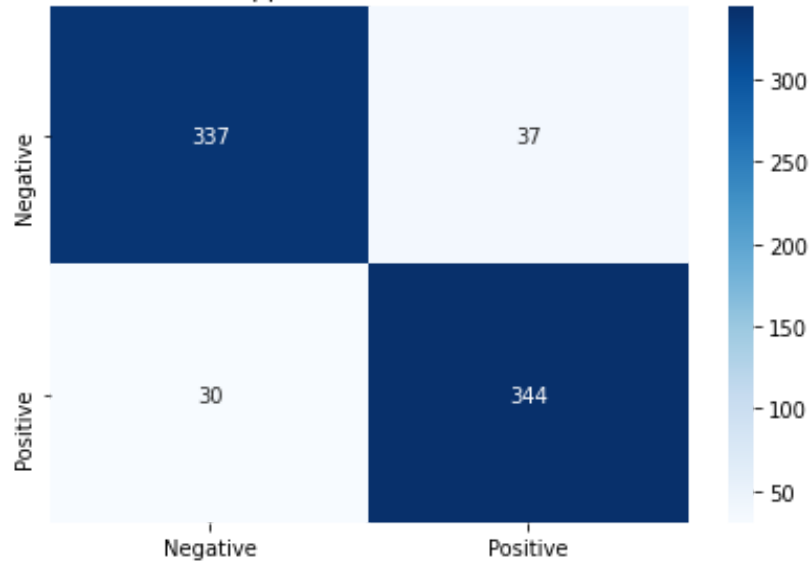
$y = 748, 1$

# MODEL COMPARISON

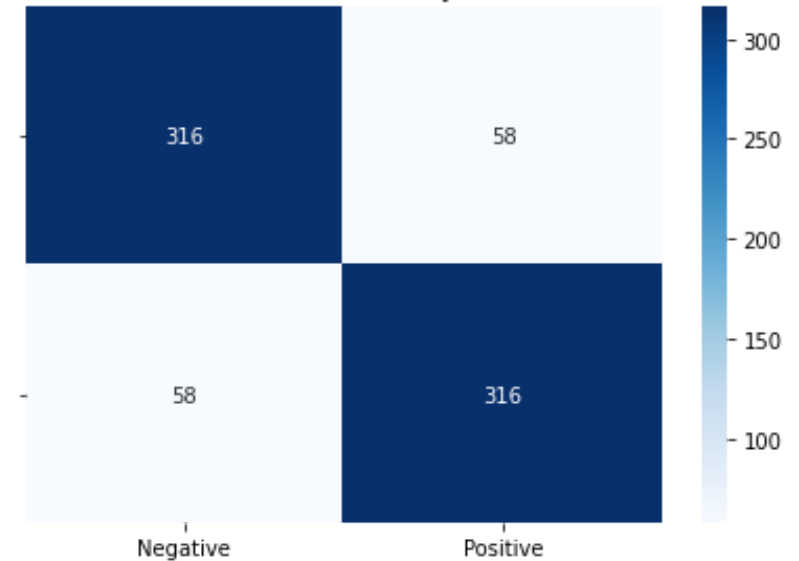
MODEL	ACCURACY	RECALL	PRECISION	F1 SCORE
Support Vector Machine	0.910	0.919	0.902	0.911
Gaussian Naive Bayes	0.844	0.844	0.844	0.844
Decision Tree Classifier	0.899	0.911	0.890	0.900
Random Forest Classifier	0.919	0.919	0.919	0.919

# CONFUSION MATRIX

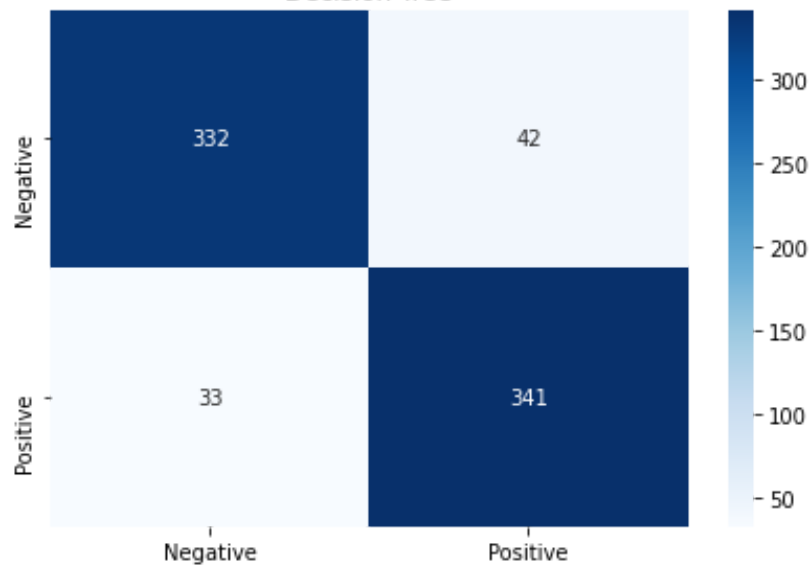
Support Vector Machine



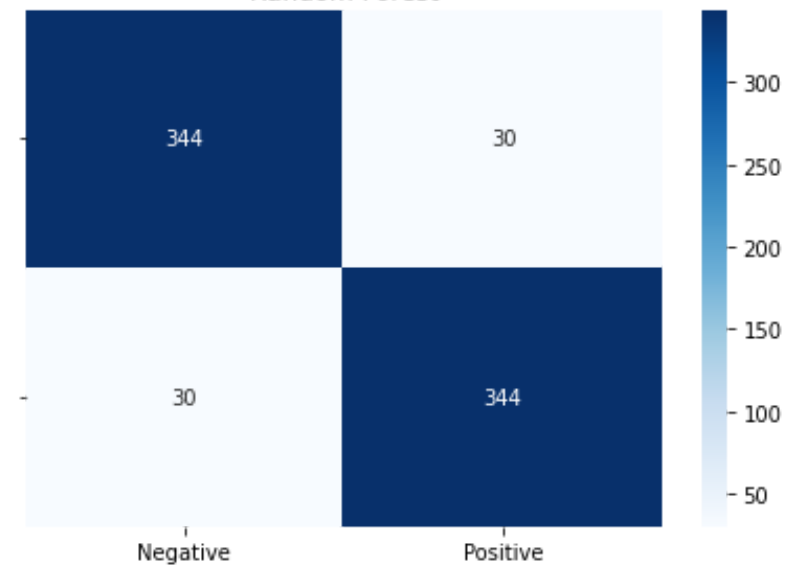
Gaussian Naive Bayes



Decision Tree



Random Forest



# ROC AUC

Area Under Curve

AUC Train & Test : 99.9% & 96.57%

Confusion Matrix Evaluation

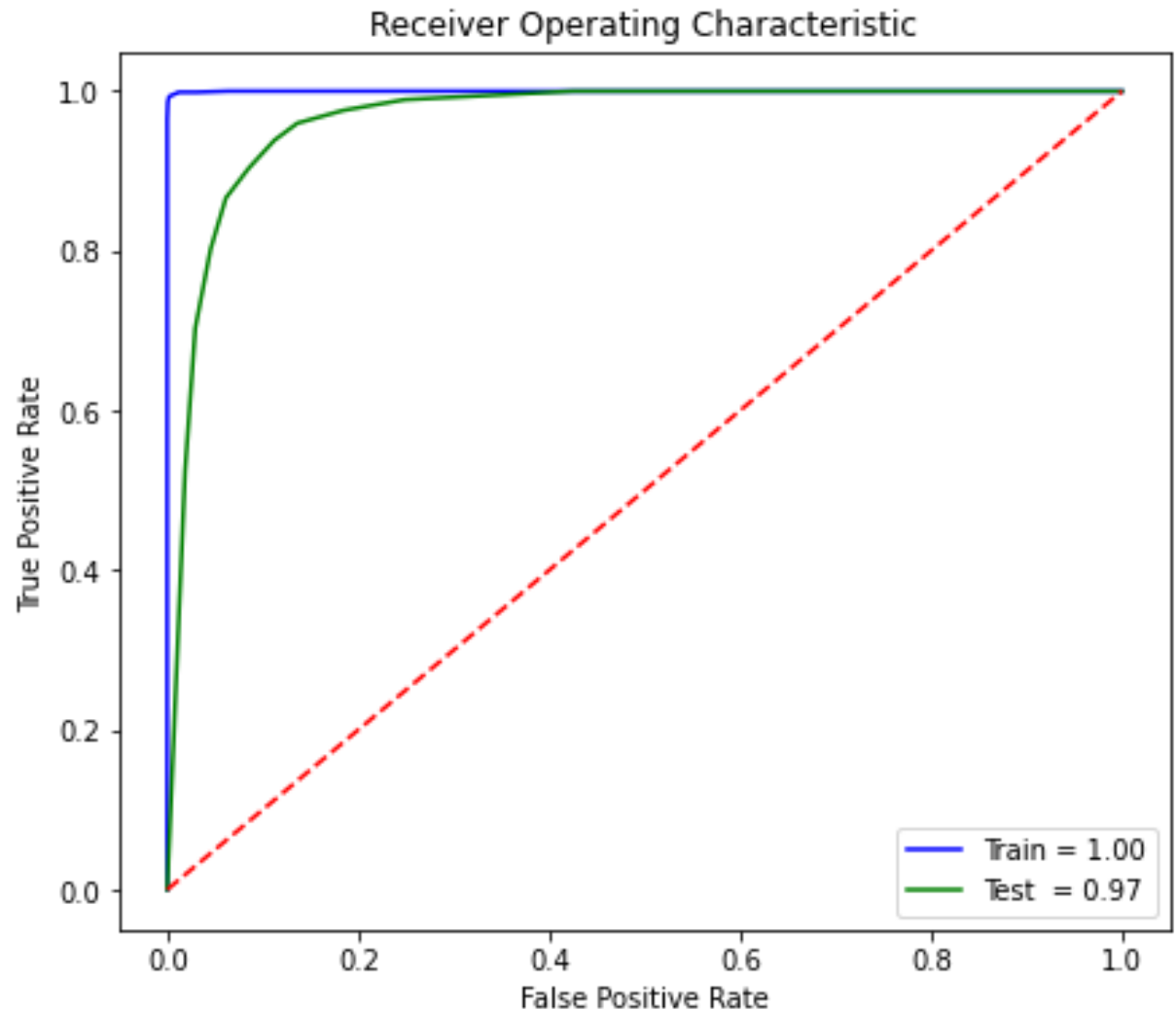
Accuracy Train & Test : 99.36% & 91.98%

Recall Train & Test: 99.13% & 91.98%

Precision Train & Test: 99.60% & 91.98%

F1 Score Train & Test: 99.36% & 91.98%

Log Loss Train & Test: 0.1502 & 3.1399



# SUMMARY

- Some column was dropped because they can cause the model to overfit.
- Random Undersampling was chosen to handle the imbalance data.
- Random Forest Classifier was chosen as it is the best model compared to the other model in this case.
- Model were validated using ROC AUC.

# THANK YOU



<https://www.linkedin.com/in/raka-raprast/>



<https://github.com/raka-raprast>