

# STRATEGIC THINKING

Banking Industry

Ariel Goldman



# BANK CHURNERS



=





## Business Understanding

Our project is focused on the banking industry, with the main goal of improving and retaining the customers to reduce churn rates.

### Hypothesis

**Analyse why customers will churn or not from the banks based on their demographic features and transaction history.**

### General Goal

**Develop a predictive Machine Learning Model which can accurately classify customers as either churn or not.**

### Success Criteria/Indicators

**The success of our project we will be measured by the accuracy of our Machine Learning Models in predicting customer churn.**



## Technology Used

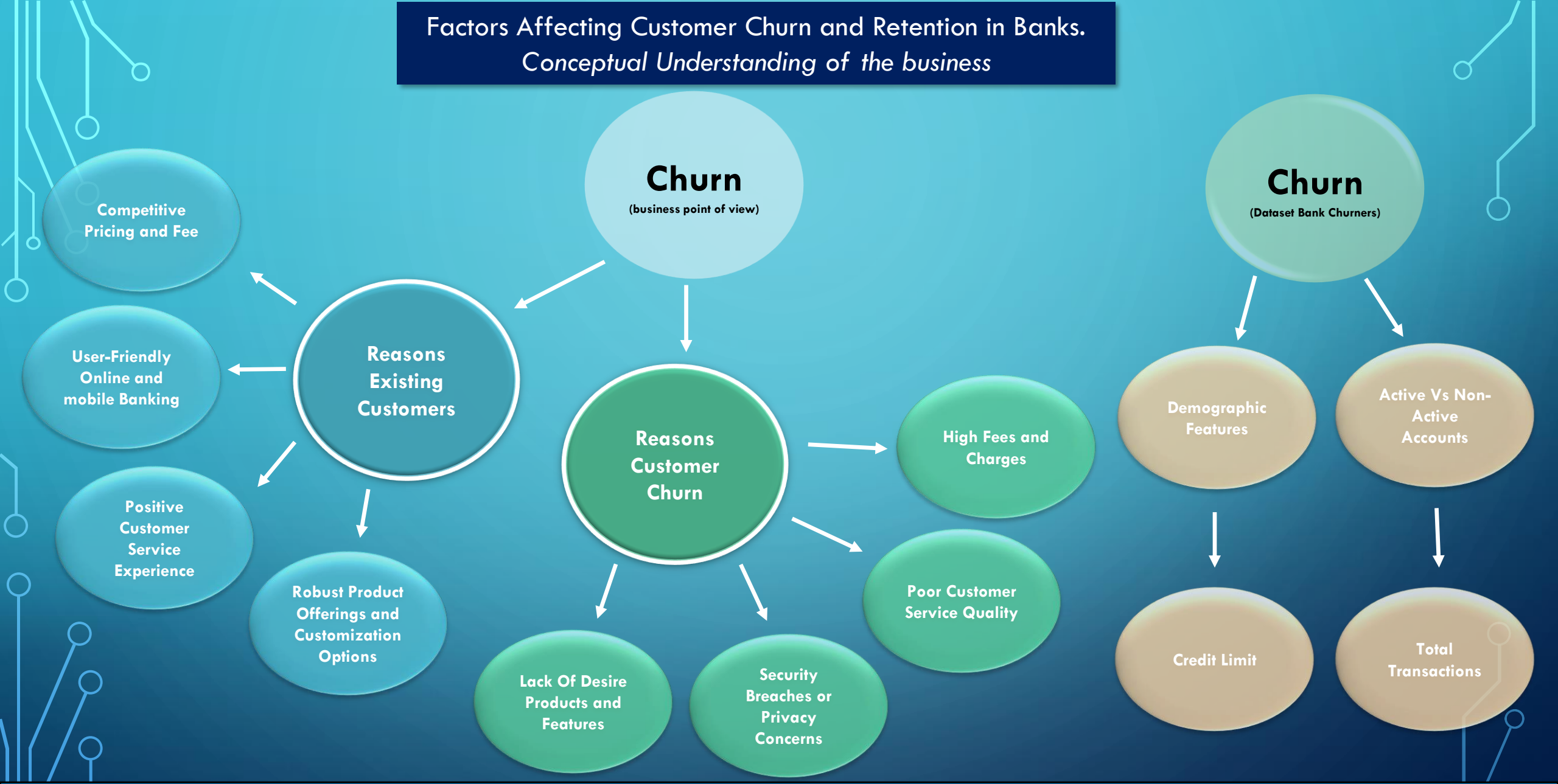
### Machine Learning Models

- **Decision Tree**
- **SupportVector Machine**
- **Logistic Regression**
- **AdaBoost**
- **Random Forest**
- **Gaussian Naïve Bayes**
- **KNeighbours**

### Libraries

- **Scikit-learn**
- **Seaborn**
- **Matplotlib**
- **Pandas**
- **NumPy**
- **Imbelear.over\_sampling.SMOTE**
- **Category\_encoders**
- **Missigno**

Factors Affecting Customer Churn and Retention in Banks.  
Conceptual Understanding of the business

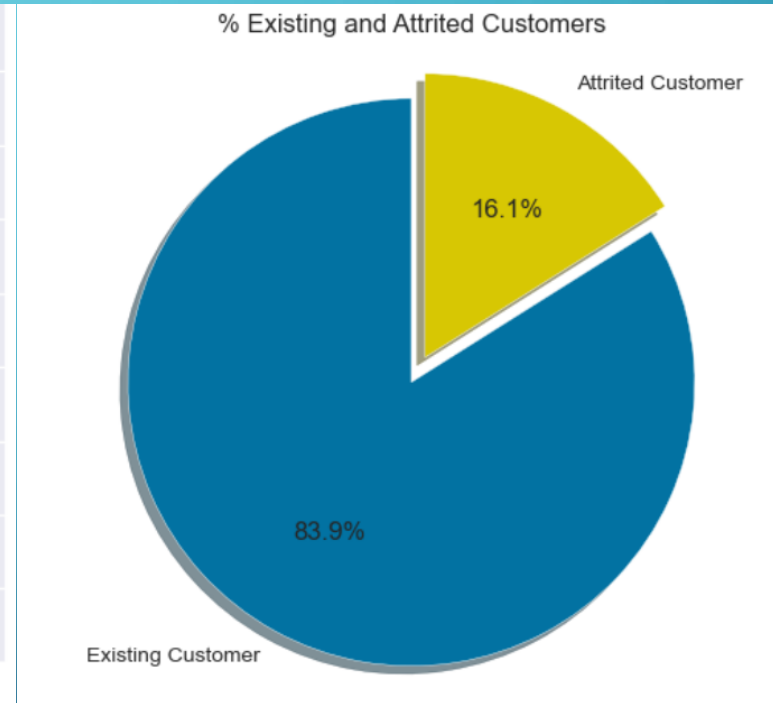




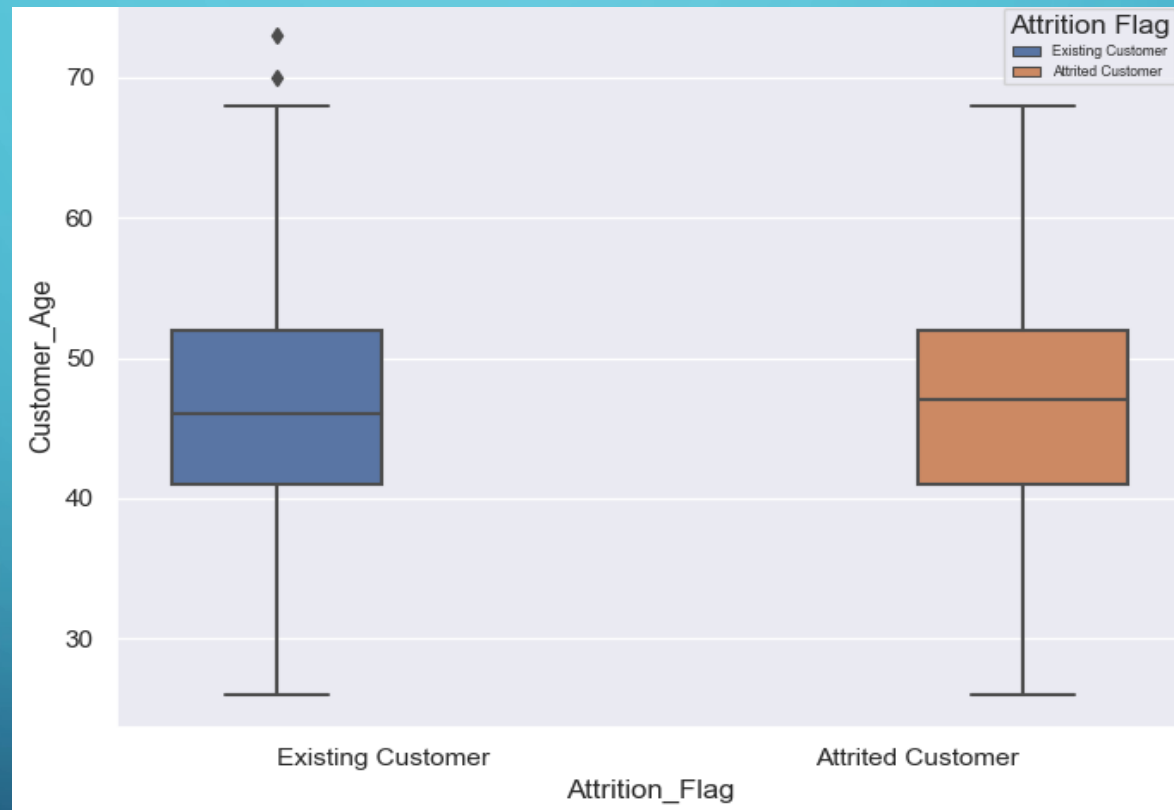
# Exploratory Data Analysis



## Imbalanced Dataset



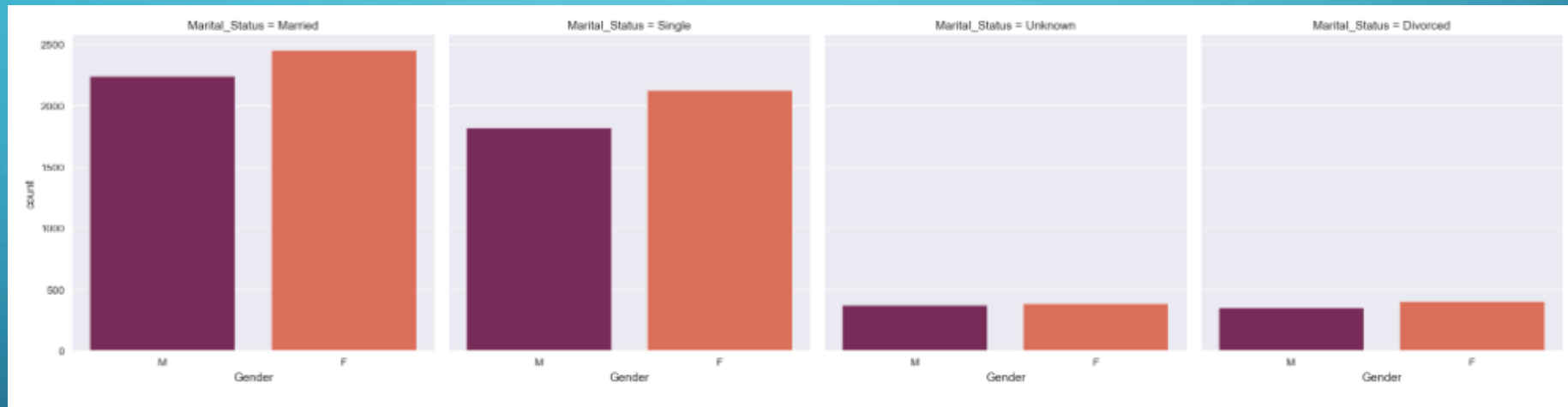
It seems that older customers are more likely to leave the bank.



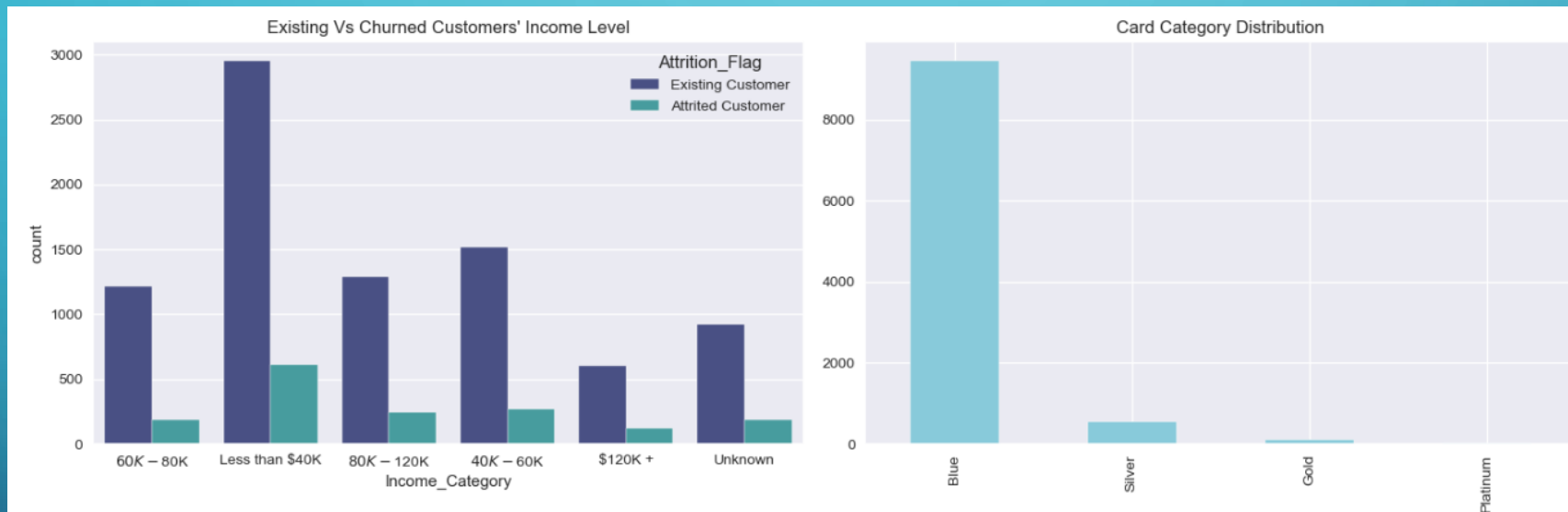


**Customers married and single are the majority in the bank.**

- **Special promotions**
- **Incentives**



- The majority of the customers have an income less than \$40k.
- 93% of the customers have credit card 'Blue'.



### Existing Customers

Total Transaction counts Vs Total Transaction Amount

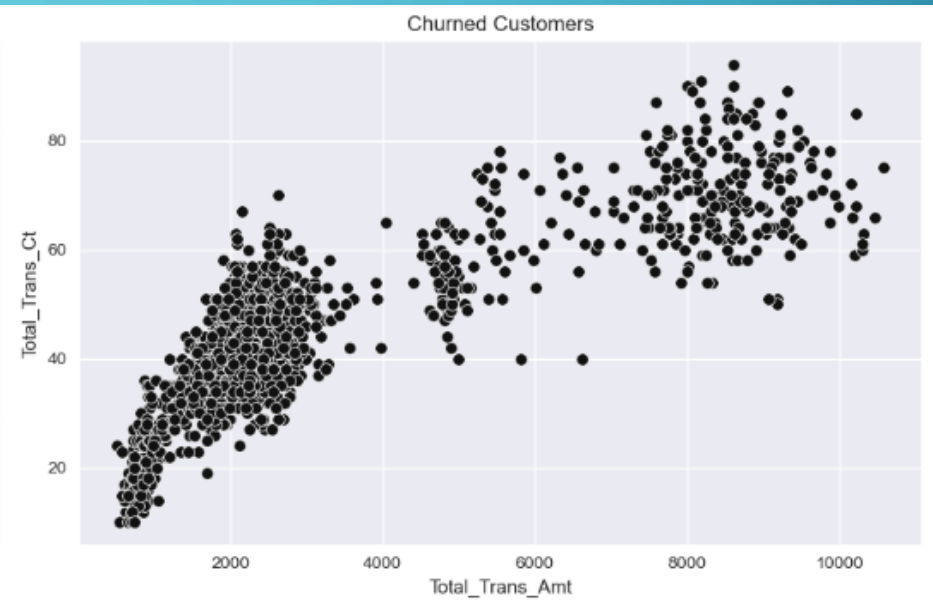
3 distinct clusters

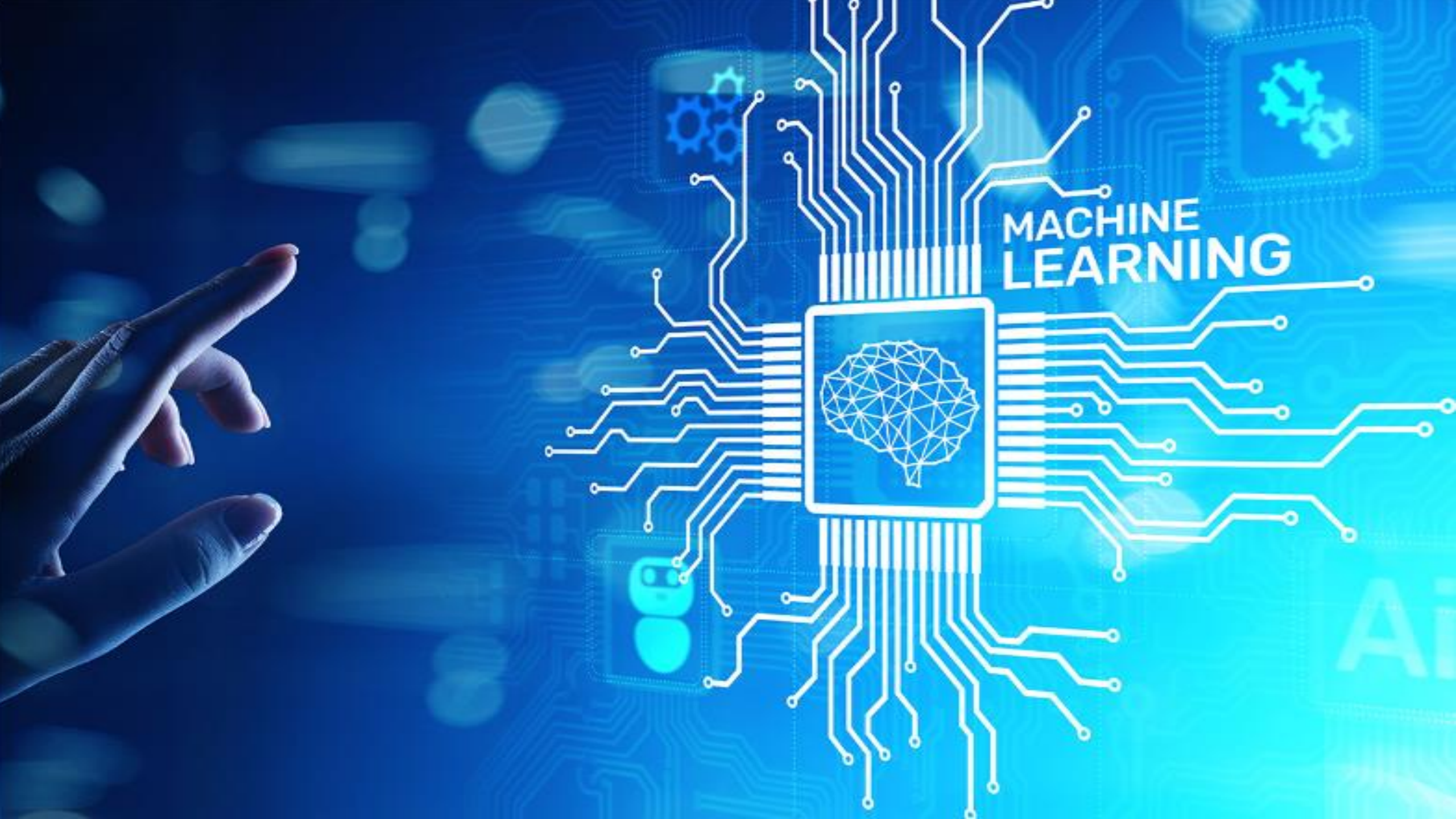


### Churned Customers

Total Transaction Counts Vs Total Transaction Amount

1 distinct cluster





**MACHINE  
LEARNING**





## Accuracy, score and recall based on SVM on balanced dataset

There are 3039 customers in the test set 30%

2225	348
CM	
80	386

There are 2026 customers in the test set 20%

1486	244
CM	
44	252

There are 1013 customers in the test set 10%

761	110
CM	
15	127

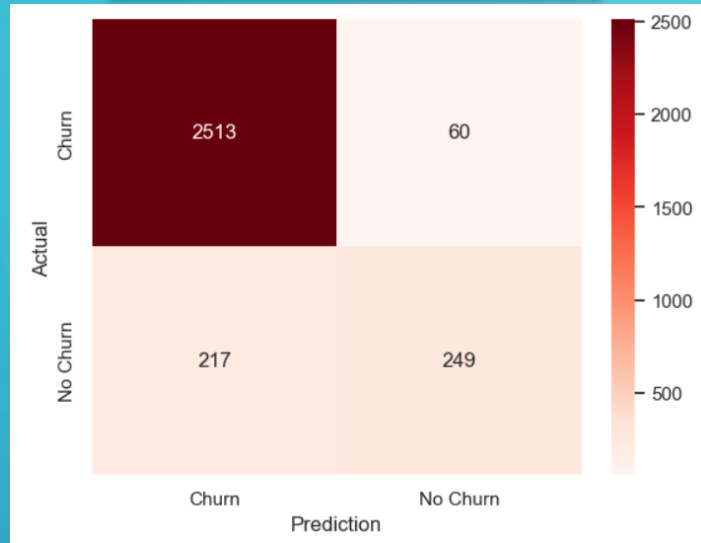
### Model Performance Metrics

<u>Test set 30%</u>	<u>Test set 20%</u>	<u>Test set 10%</u>
Accuracy: 85.92	Accuracy: 85.78%	Accuracy: 87.66%
Precision: 96.53%	Precision: 97.12%	Precision: 98.07%
Recall: 86.47%	Recall: 85.90%	Recall: 87.37%

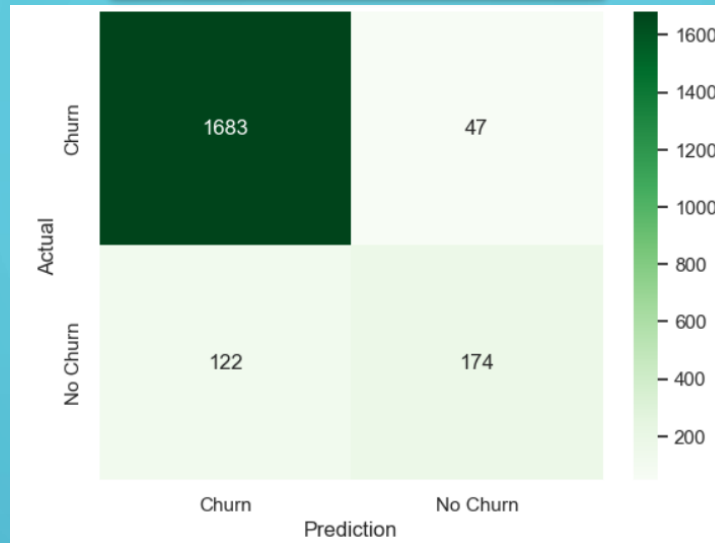


# Accuracy, score and recall based on SVM on imbalanced dataset

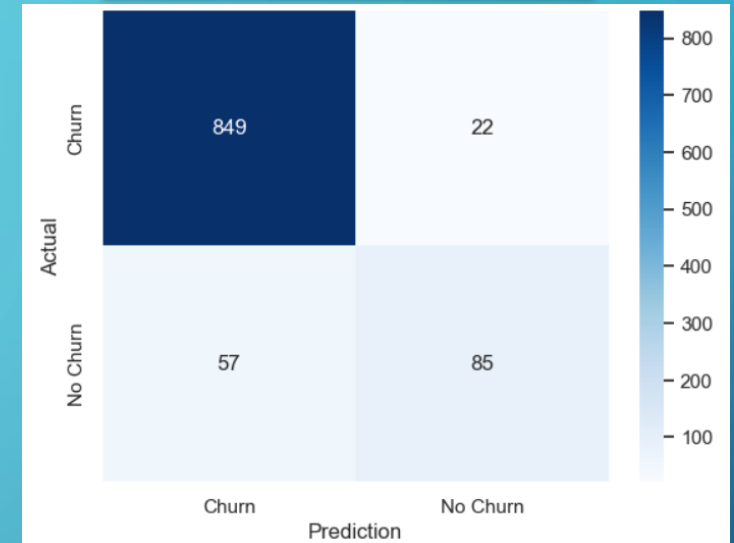
There are 3039 customers in the test set 30%



There are 2026 customers in the test set 20%



There are 1013 customers in the test set 10%



## Model Performance Metrics

<u>Test set 30%</u>	<u>Test set 20%</u>	<u>Test set 10%</u>
Accuracy: 90.89%	Accuracy: 91.66%	Accuracy: 92.20%
Precision: 92.05%	Precision: 93.24%	Precision: 93.71%
Recall: 97.67%	Recall: 90.63%	Recall: 97.47%

## Classification Report Random Forest

Test Set 30%

	precision	recall	f1-score	support
1	0.97	0.98	0.97	2573
2	0.86	0.86	0.86	466
accuracy			0.96	3039
macro avg	0.92	0.92	0.92	3039
weighted avg	0.96	0.96	0.96	3039

Test Set 20%

	precision	recall	f1-score	support
1	0.97	0.98	0.97	1703
2	0.87	0.86	0.87	323
accuracy			0.96	2026
macro avg	0.92	0.92	0.92	2026
weighted avg	0.96	0.96	0.96	2026

Test Set 10%

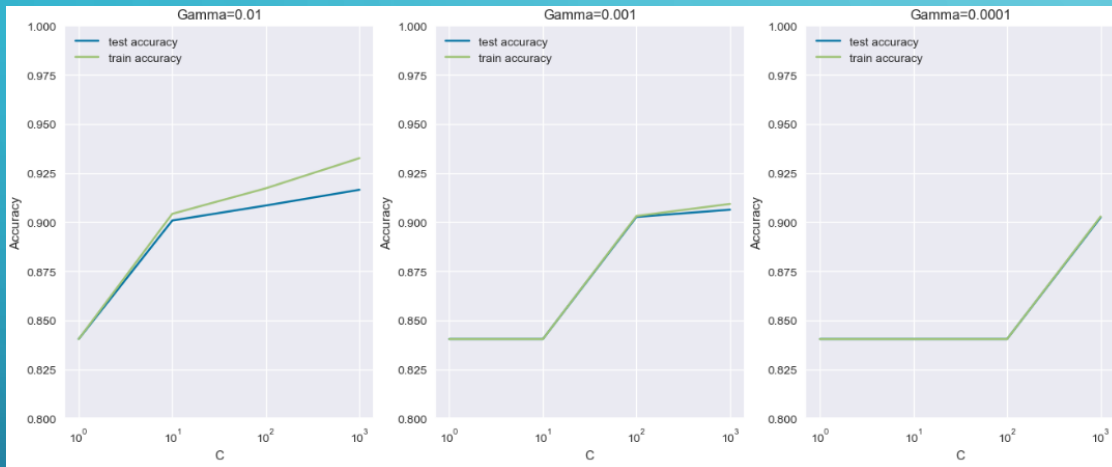
	precision	recall	f1-score	support
1	0.87	0.98	0.92	871
2	0.35	0.06	0.11	142
accuracy			0.85	1013
macro avg	0.61	0.52	0.51	1013
weighted avg	0.79	0.85	0.81	1013

# Deployment of Machine Learning Model

## SVM and Random Forest

### SVM

- Supervised learning
- Classify customers as churn or non-churn based on banking behavior
- GridSearchCV to find optimal hyperparameter ('Kernel=RBF')



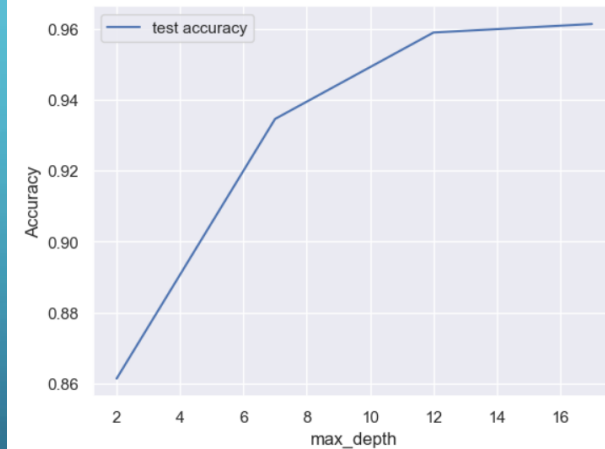
### RBF Linear Kernel Model

- The model has a good balance between being too simple (underfitting) and being too complex (overfitting).
- The best test score is 92%, it means the model is able to predict the correct output (generated by the model ) based on the provided data.

### Random Forest

- Supervised learning
- Classify customers as churn or non-churn based on banking behavior
- Hyperparameter Tuning - Max\_depth

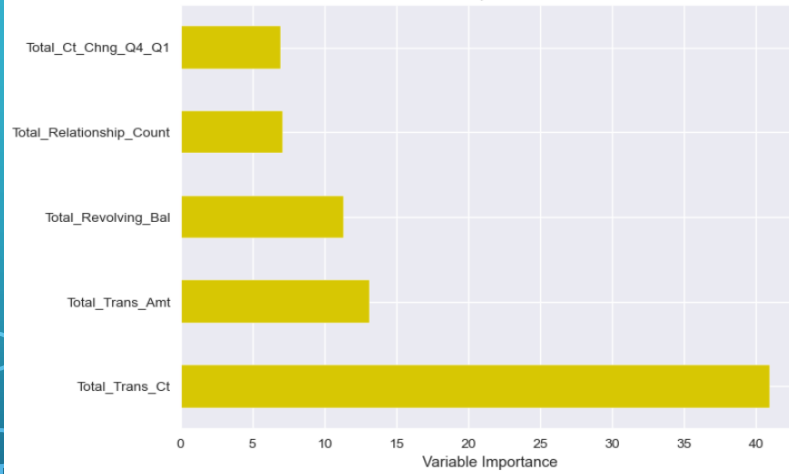
params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	split4_test_score	mean_test_score	std_test_score	rank_test_score
{'max_depth': 2}	0.86	0.87	0.88	0.87	0.87	0.87	0.00	4
{'max_depth': 7}	0.94	0.93	0.94	0.94	0.94	0.94	0.00	3
{'max_depth': 12}	0.96	0.96	0.96	0.95	0.96	0.96	0.00	2
{'max_depth': 17}	0.96	0.96	0.96	0.96	0.97	0.96	0.00	1



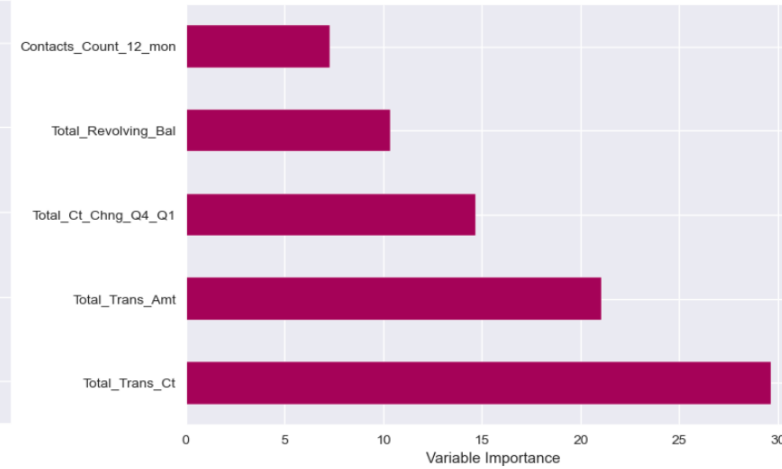
After fitting the model with the values of 'max\_depth', the plot indicates that the test accuracy 96% is high when 'max\_depth' is set in a range of 2 to 20 while splitting the model into 5.

# Future Importances

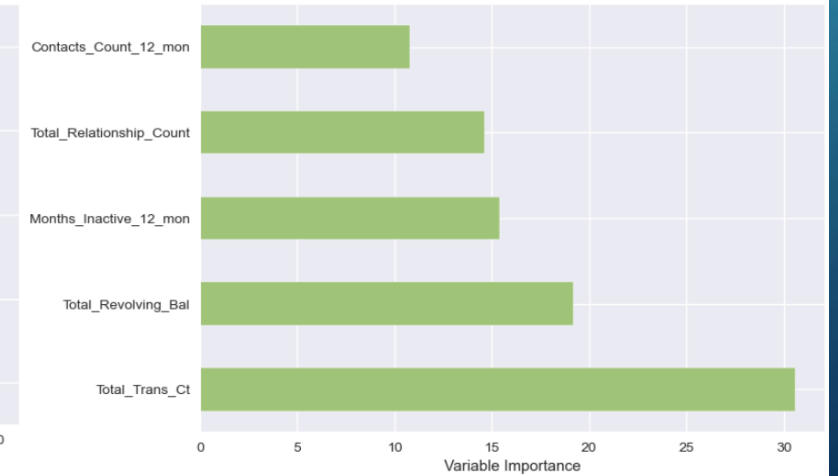
Feature Importances DT



Feature Importances RF

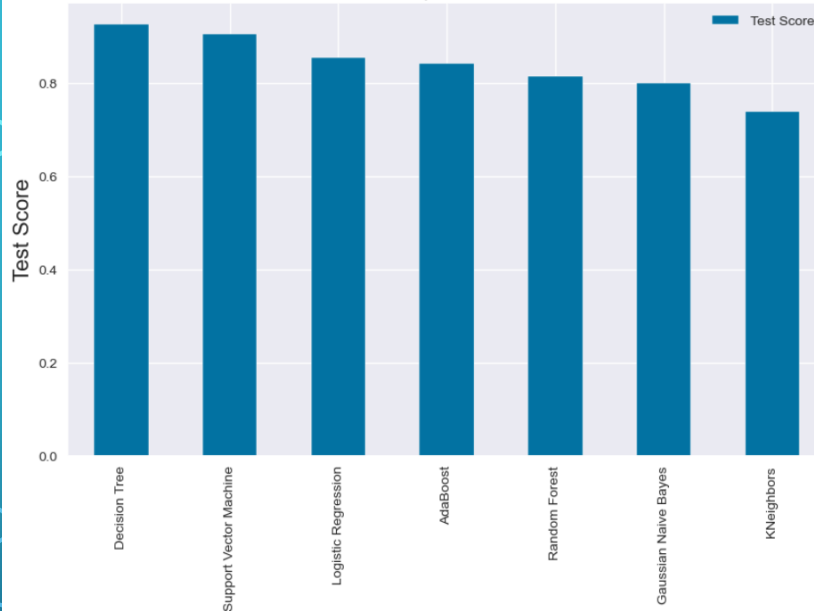


Feature Importances AB

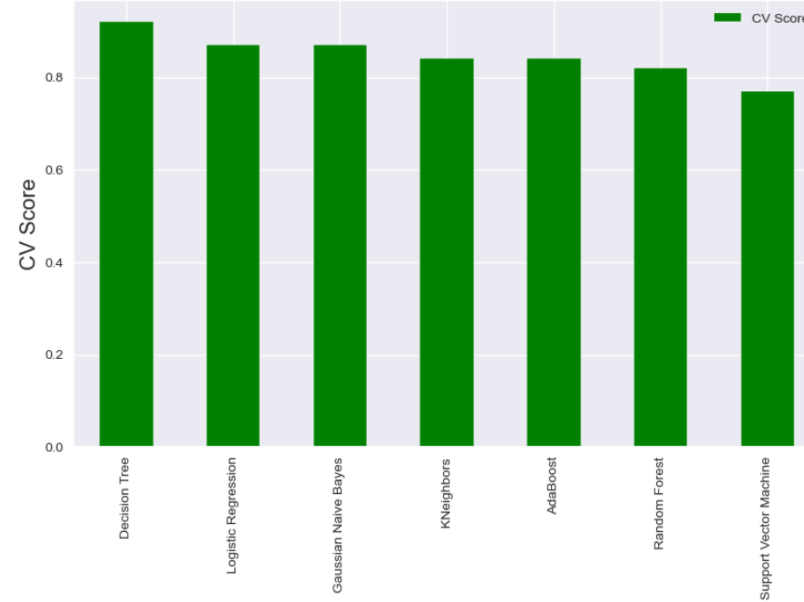


# Model Comparisons

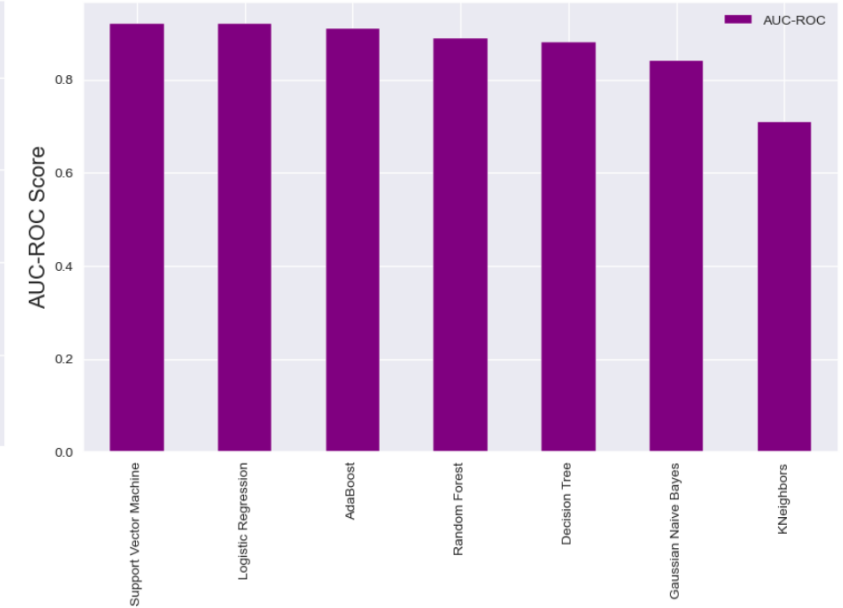
### Test Score Comparison - ML Models



### CV Score Comparison - ML Models



### AUC-ROC Score Comparison - ML Models



## Conclusion

Demographic factors like age, gender, income, marital status, and education are not reliable indicators for predicting customer churn, as shown by Random Forest, Decision Tree, and Adaboost models. Even after an extensive Grid Search, no clear correlation was established. On the other hand, transaction history features consistently proved highly significant. In conclusion, further investigation is needed to identify the most effective model for real-world deployment.



# QUESTIONS?



# Thank you !!!

