

1. Introduction.....	1, 8
2. EDA and Business Implication.....	8, 12
3. Data Cleaning and Pre-processing.....	12, 16
4. Model building.....	17, 27
5. Model validation.....	27, 38
6. Final interpretation/recommendation	38, 44

1. Introduction

The purpose of this whole exercise is to explore the dataset. Do the exploratory data analysis. Explore the dataset and analysis of An E-Commerce Company or DTH, the company wants to develop a model through which they can do churn prediction of the accounts and provide segmented offers to the potential churners. We are using the best module to find religion for customer churn.

Customer Churn is very expensive for any business or organization. A high Churn Rate requires a company to deal with the stress of doubling down to bring in new customers; just to stay afloat. Even a minuscule single-digit increase in the Churn Rate can seriously impede a company's growth rate and what's worse is that high Churn Rates are more likely to compound over time.

This is necessary for every organization to retain their customer to sustain in the market. Customer Churn can lead to financial loss and so we need to study this project.

- **Expose Product Weaknesses:** Churn Analysis plays a crucial role in revealing the patterns that indicate the common motivators for customers to part ways with your company. These could be anything from price productivity to poor product adoption. It is also instrumental in demonstrating the exact way customers engage with your product throughout their lifecycle. You can use this to maximize what your customers already love, and improve upon everything they don't.
- **Uearth Customer Opportunities:** Customer experience improvement comes with a need to understand the customer requirements and expectations at every stage of their journey and moulding your product accordingly. Churn Analysis depicts trends in customer behaviour at every touch point. Personalized engagement preferred by your customers allows you to make your customers feel valued and appreciated.

Executive Summary

An E-Commerce company or DTH (you can choose either of these two domains) provider is facing a lot of competition in the current market and it has become a challenge to retain the existing customers in the current situation. Hence, the company wants to develop a model through which they can do churn prediction of the accounts and provide segmented offers to the potential churners. In this company, account churn is a major thing because 1 account can have multiple customers. Hence by losing one account the company might be losing more than one customer. You have been assigned to develop a churn prediction model for this company and provide business recommendations on the campaign. Your campaign suggestion should be unique and be very clear on the campaign offer because your recommendation will go through the revenue assurance team. If they find that you are giving a lot of free (or subsidized) stuff thereby making a loss to the company; they are not going to approve your recommendation. Hence be very careful while providing campaign recommendations.

Understanding of attributes (variable info, renaming if required)

➤ **Variable info:-**

```
Data columns (total 19 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
 0   AccountID        11260 non-null   int64  
 1   Churn            11260 non-null   int64  
 2   Tenure           11158 non-null   object  
 3   City_Tier        11148 non-null   float64 
 4   CC_Contacted_LY 11158 non-null   float64 
 5   Payment          11151 non-null   object  
 6   Gender           11152 non-null   object  
 7   Service_Score   11162 non-null   float64 
 8   Account_user_count 11148 non-null   object  
 9   account_segment  11163 non-null   object  
 10  CC_Agent_Score  11144 non-null   float64 
 11  Marital_Status  11048 non-null   object  
 12  rev_per_month   11158 non-null   object  
 13  Complain_ly    10903 non-null   float64 
 14  rev_growth_yoy 11260 non-null   object  
 15  coupon_used_for_payment 11260 non-null   object  
 16  Day_Since_CC_connect 10903 non-null   object  
 17  cashback        10789 non-null   object  
 18  Login_device    11039 non-null   object  
dtypes: float64(5), int64(2), object(12)
memory usage: 1.6+ MB
```

fig_1.d

➤ **Objects detail :-**

i.) Login_device

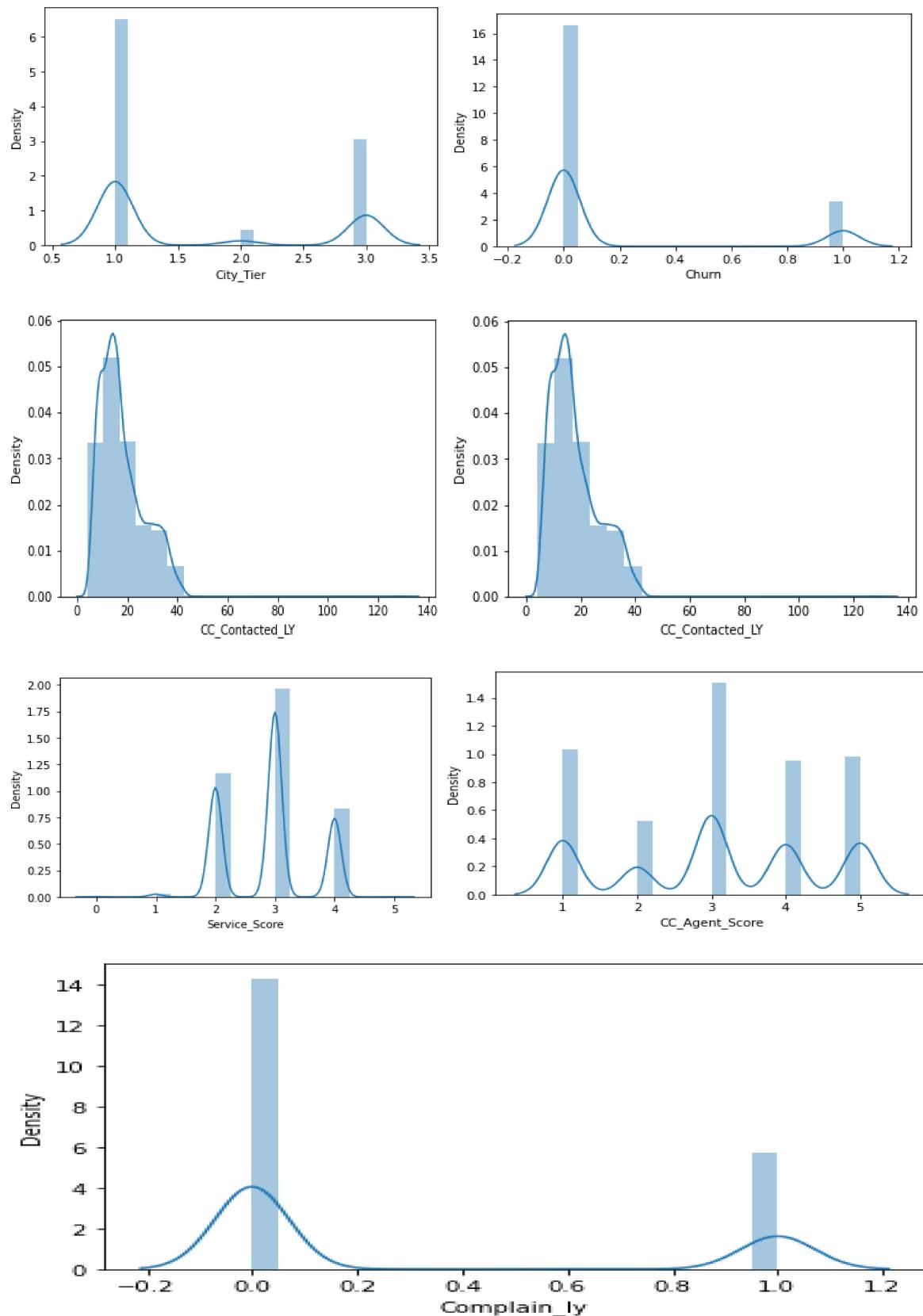
```
Login_device
Mobile      7482
Computer    3018
&&&&       539
dtype: int64
```

ii.) Account_segment

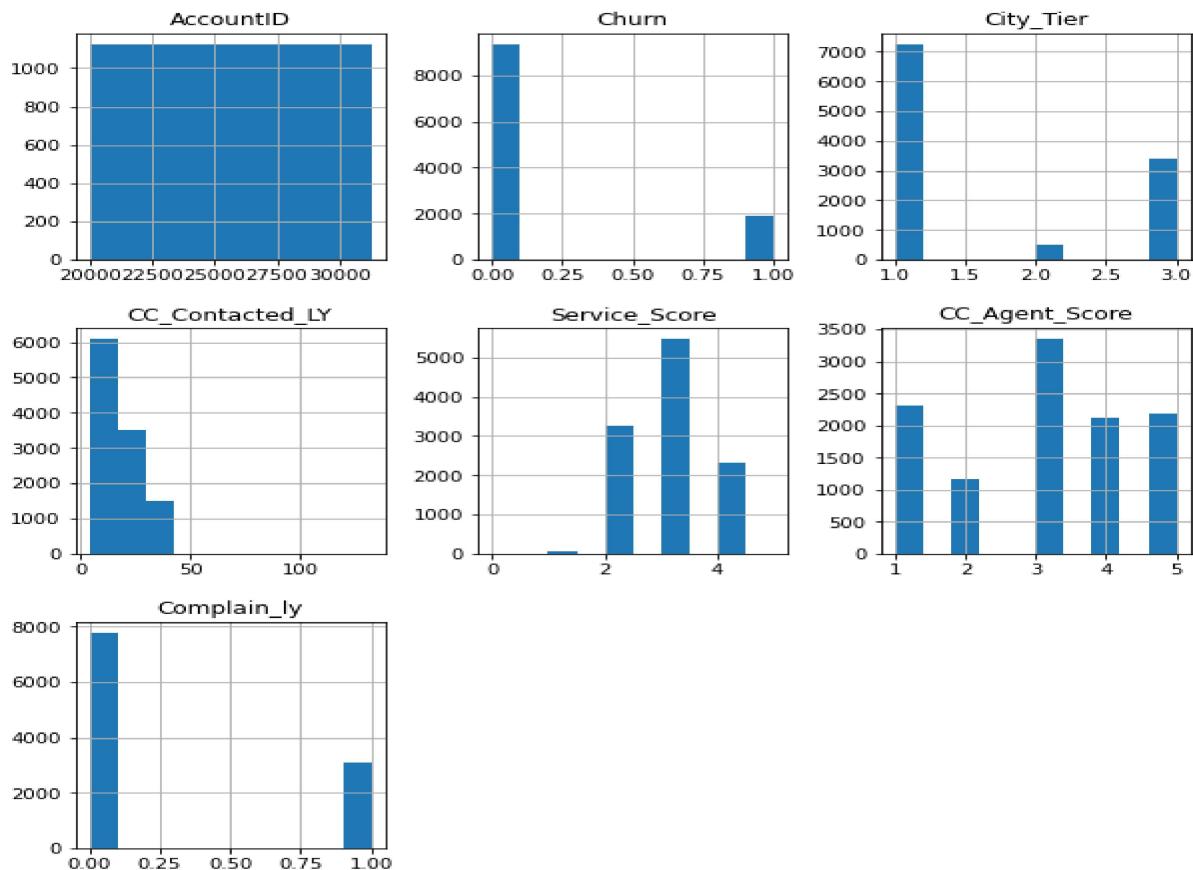
```
account_segment
Super        4062
Regular Plus 3862
HNI          1639
Super Plus   771
Regular      520
Regular +    262
Super +      47
```


2. EDA and Business Implication

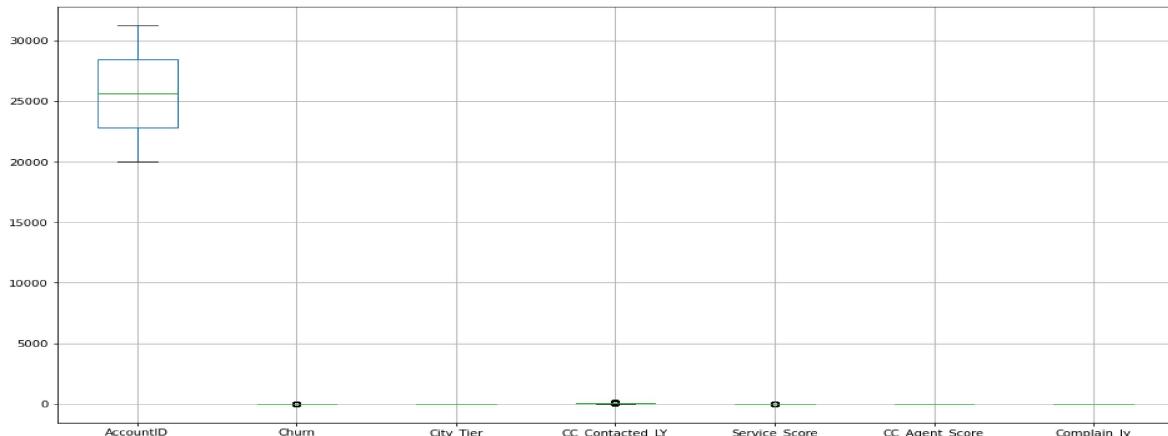
Univariate analysis



fig_2.a



Fig_2.b

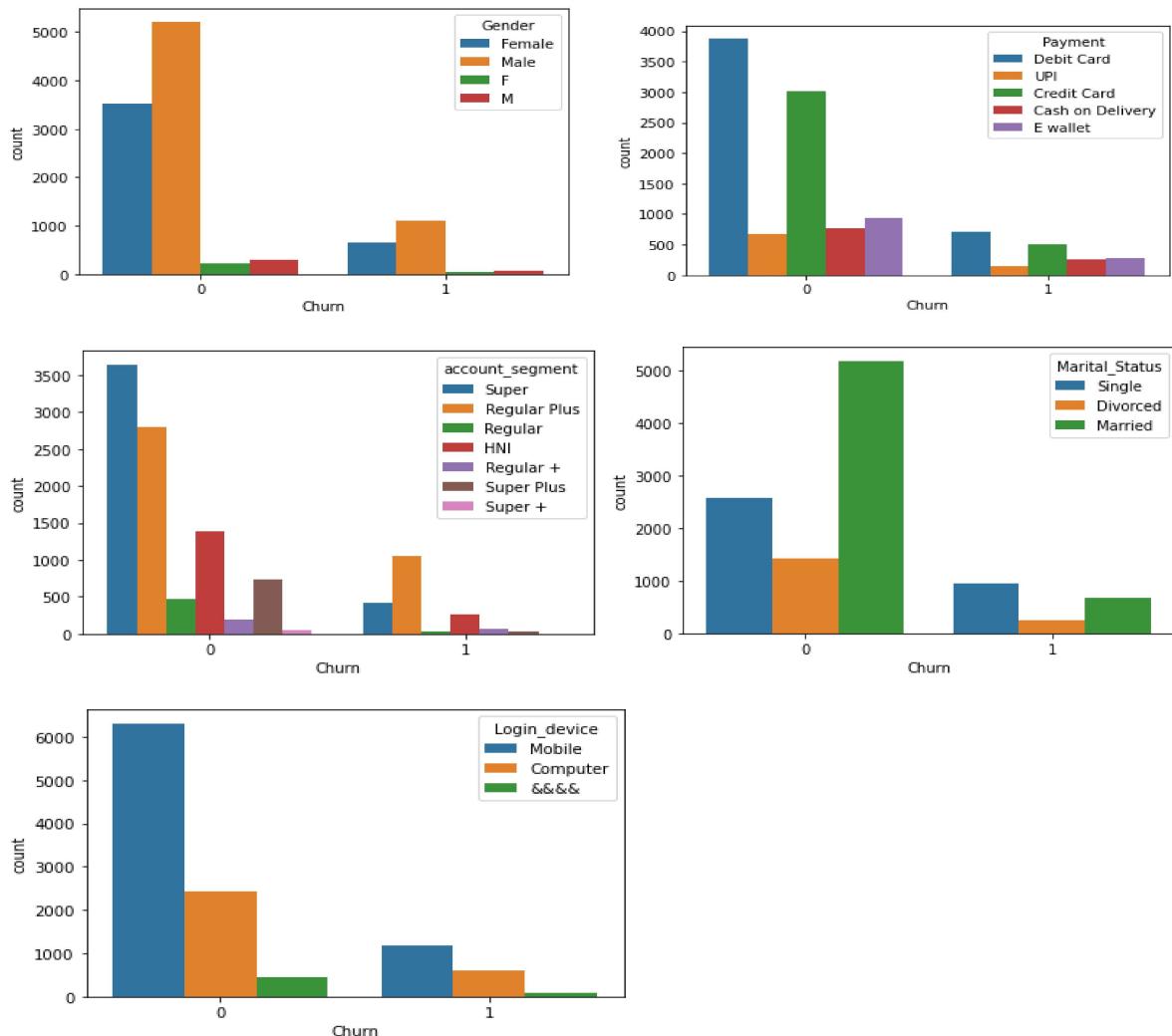


fig_2.c

Conclusion | insight: All columns data single behaviour show in Fig_2.a & Fig_2.b, Fig_2.C.

1. We have a count plot for our target variable churn and we can observe a high no. of churn as compared to non-churn.
2. It is found that customer churn is higher in Tier 1 cities and lowest in tier 2 cities. Tier 1 customer needs more focus to be retained with the organization.
3. It is found that the count for customer care contacted is very high who churn.

Bivariate Analysis



Fig_3.f

Conclusion | insight:

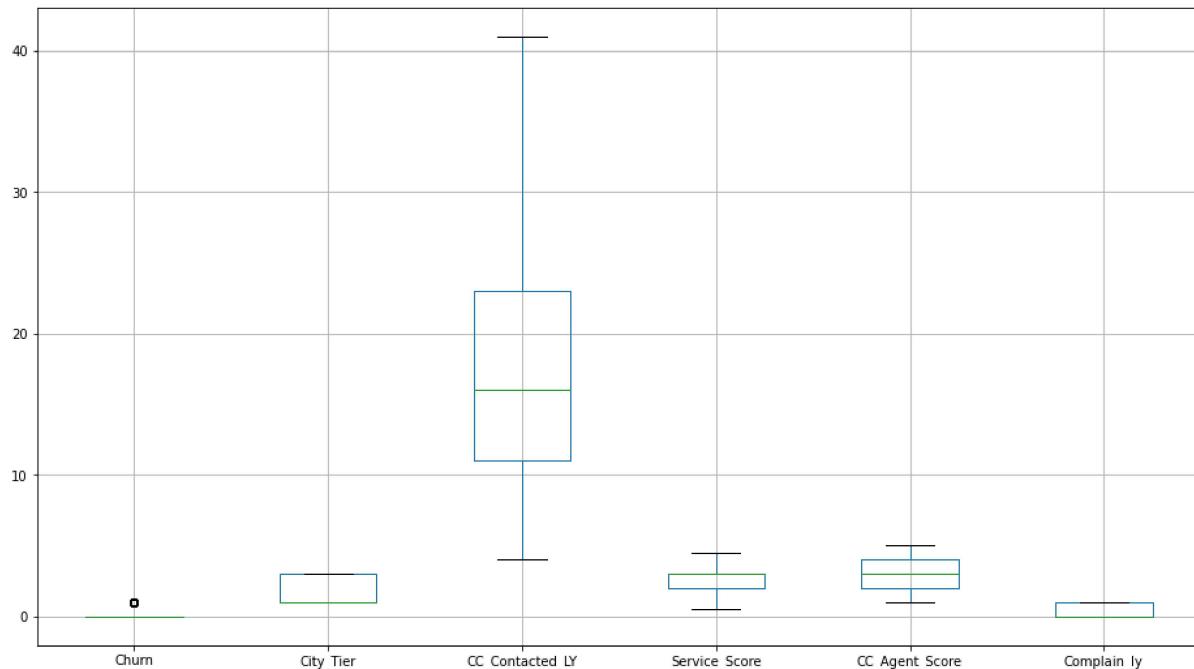
1. Here in Bivariate Analysis, we have obtained churn count with respect to Gender (Male & Female), Payment method (Debit card, UPI, Credit card, Cash on Delivery and E wallet), account segment and marital status (single, Divorced and Married). If we look for different account segment then churn rate is maximum for super account segment and it is minimum in super+ plan. So we can say that Super plan is most disliked and Super + plan is most liked by customer and we should spread it more across.
2. We can clearly see that churn count is higher for male in Gender and churn count is higher for married people in marital status as compared to other categories.

Churn	0
Tenure	0
City_Tier	0
CC_Contacted_LY	0
Payment	0
Gender	0
Service_Score	0
Account_user_count	0
account_segment	0
CC_Agent_Score	0
Marital_Status	0
rev_per_month	0
Complain_ly	0
rev_growth_yoy	0
coupon_used_for_payment	0
Day_Since_CC_connect	0
cashback	0
Login_device	0

Fig_4.d

Conclusion | insight:- all missing value are deleted shown in fig_4.d

Outlier treatment



Fig_4.m

Conclusion | insight:- All outlier treated now no any outliner available show in fig_4.m

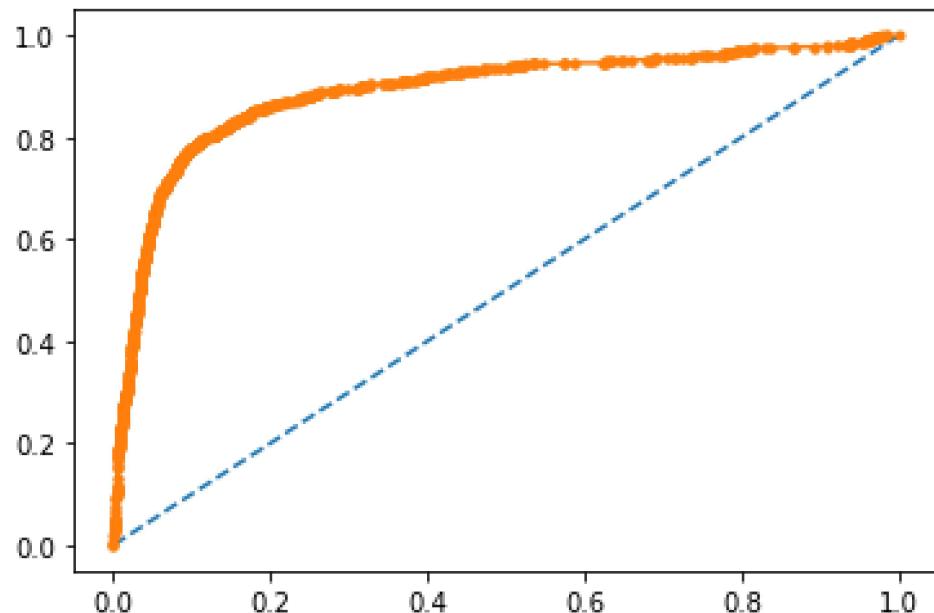
ROC-AUC Graph**AUC:** 0.891

Fig-1.a.1

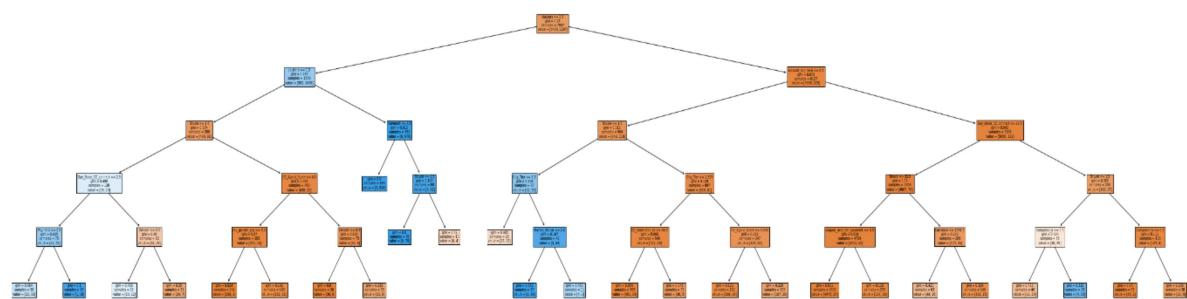
(2) Decision Tree**Train:- Graphviz**

Fig-1.a.2

Train:- confusion_matrix & classification_report , ROC-AUC Graph**confusion_matrix**

```
[[6357,    43],  
 [ 204, 1093]],
```

classification_report

	precision	recall	f1-score	support
0	0.97	0.99	0.98	6400
1	0.96	0.84	0.90	1297
accuracy			0.97	7697
macro avg	0.97	0.92	0.94	7697
weighted avg	0.97	0.97	0.97	7697

ROC-AUC Graph

AUC: 0.977

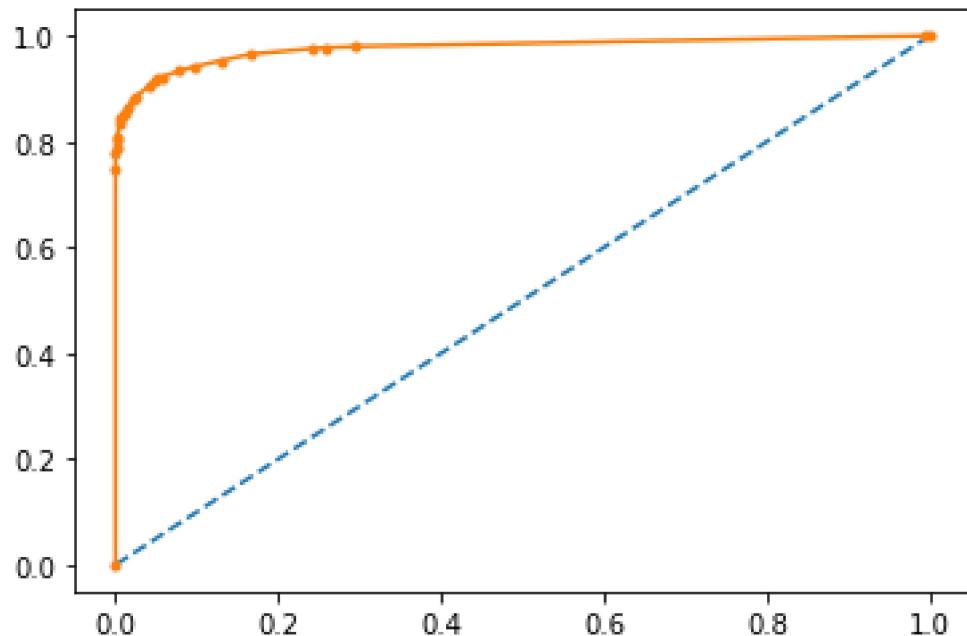


Fig-1.a.3

(3) Random Forest

Train:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[6357,    43],
 [ 204, 1093]],
```

classification_report

	precision	recall	f1-score	support
0	0.97	0.99	0.98	6400
1	0.96	0.84	0.90	1297
accuracy			0.97	7697
macro avg	0.97	0.92	0.94	7697
weighted avg	0.97	0.97	0.97	7697

ROC-AUC Graph

AUC: 0.9931

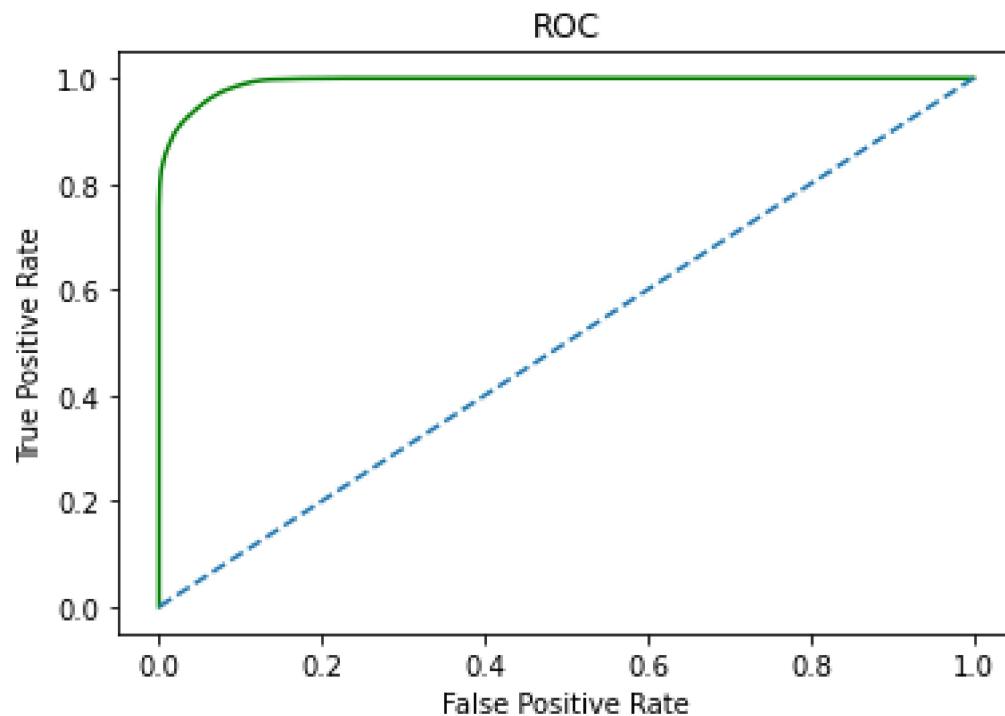


Fig-1.a

(4) Linear Discriminant Analysis

Train:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[5966, 434],  
 [ 495, 802]]
```

classification_report

	precision	recall	f1-score	support
0	0.84	0.98	0.91	2745
1	0.45	0.06	0.11	554
accuracy			0.83	3299
macro avg	0.65	0.52	0.51	3299
weighted avg	0.77	0.83	0.77	3299

ROC-AUC Graph

AUC: 0.8406

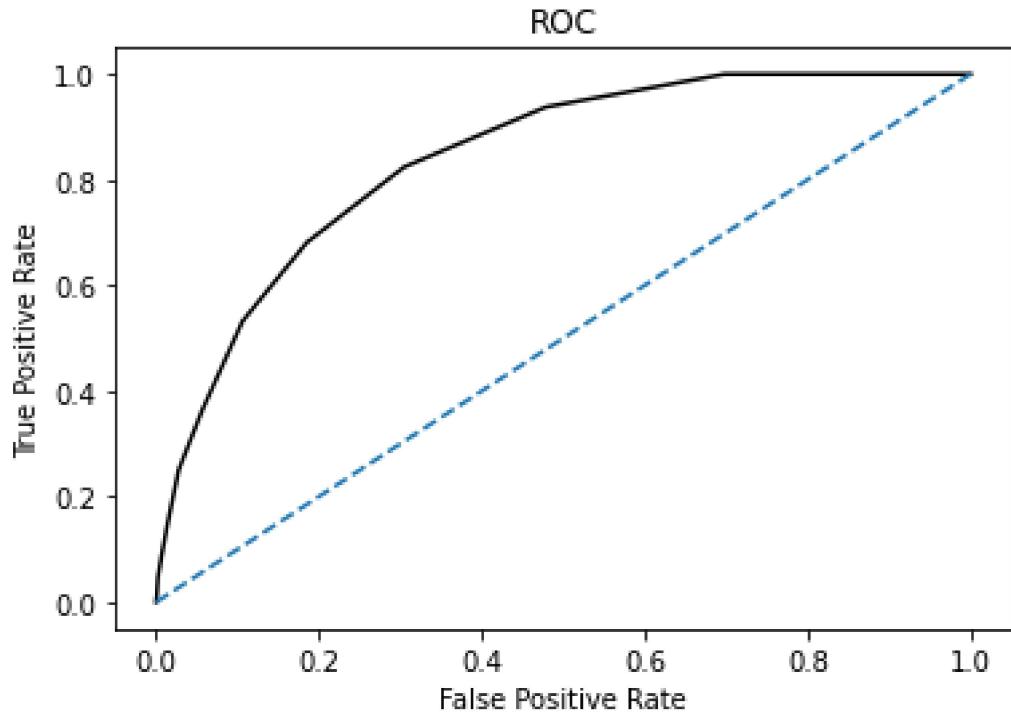


Fig-1.a.6

(6) Naive Bayes

Train:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[6014, 386],  
 [ 398, 899]],
```

classification_report

	precision	recall	f1-score	support
0	0.94	0.94	0.94	6400
1	0.70	0.69	0.70	1297
accuracy			0.90	7697
macro avg	0.82	0.82	0.82	7697
weighted avg	0.90	0.90	0.90	7697

ROC-AUC Graph

AUC: 0.868

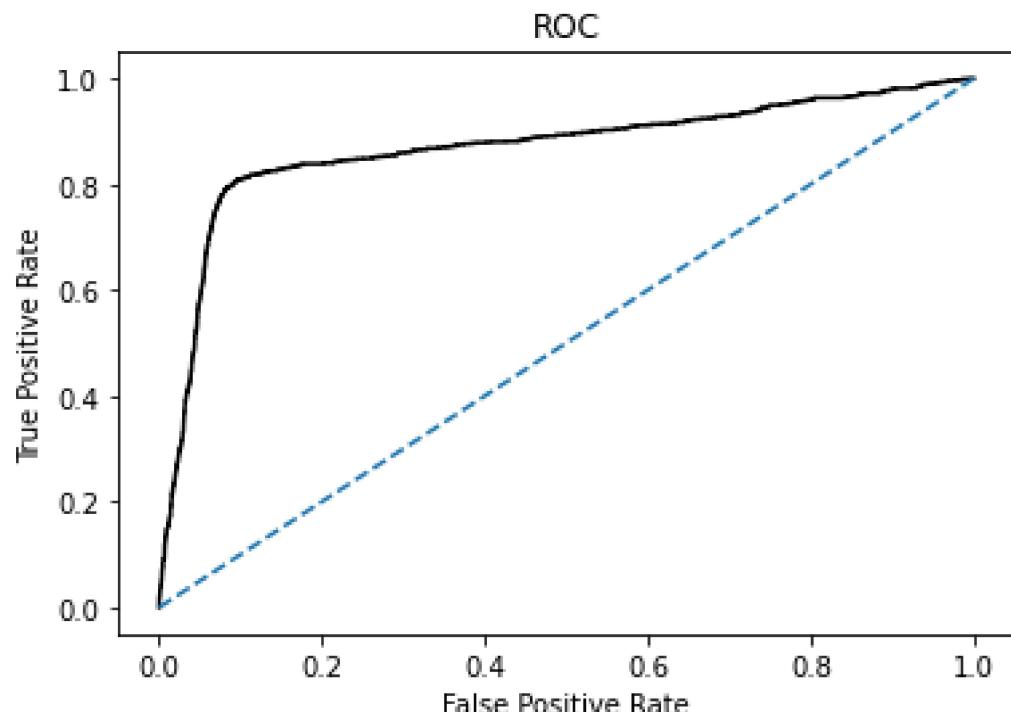


Fig-1.a.7

(7) Gradient Boosting

Train:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[6400,     0],
[1297,     0]]
```

classification_report

	precision	recall	f1-score	support
0	0.83	1.00	0.91	6400
1	0.00	0.00	0.00	1297
accuracy			0.83	7697
macro avg	0.42	0.50	0.45	7697
weighted avg	0.69	0.83	0.75	7697

ROC-AUC Graph

AUC: 0.5

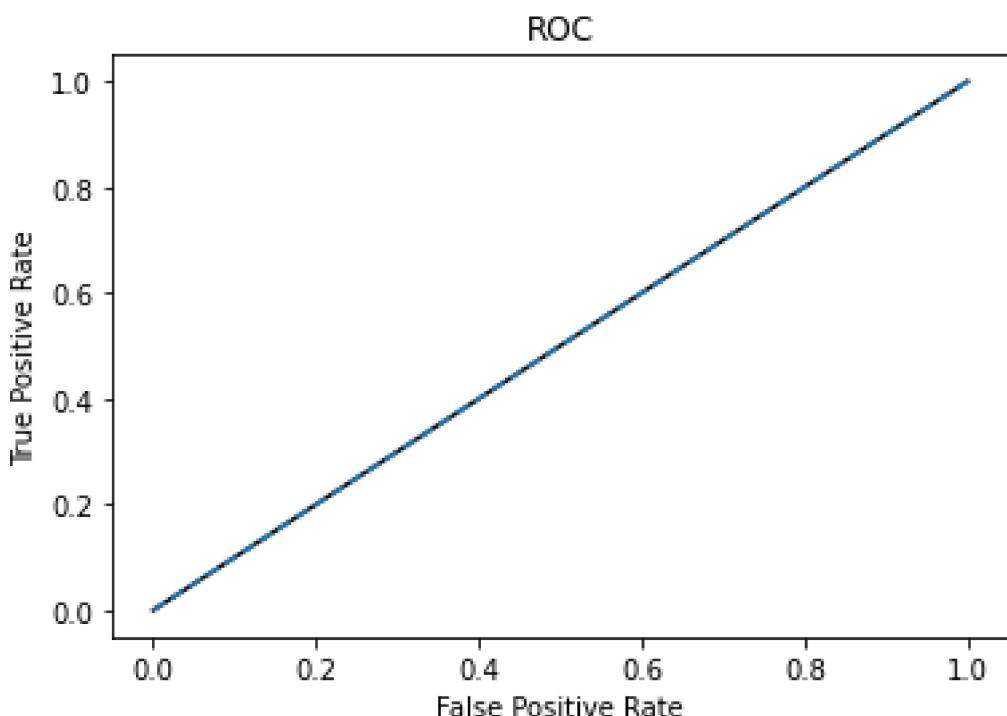


Fig-1.a.8

(8) Extreme Gradient Boosting

Train:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[6400, 0],
 [18, 1279]]
```

classification_report

	precision	recall	f1-score	support
0	1.00	1.00	1.00	6400
1	1.00	0.99	0.99	1297
accuracy			1.00	7697
macro avg	1.00	0.99	1.00	7697
weighted avg	1.00	1.00	1.00	7697

ROC-AUC Graph

AUC: 0.999

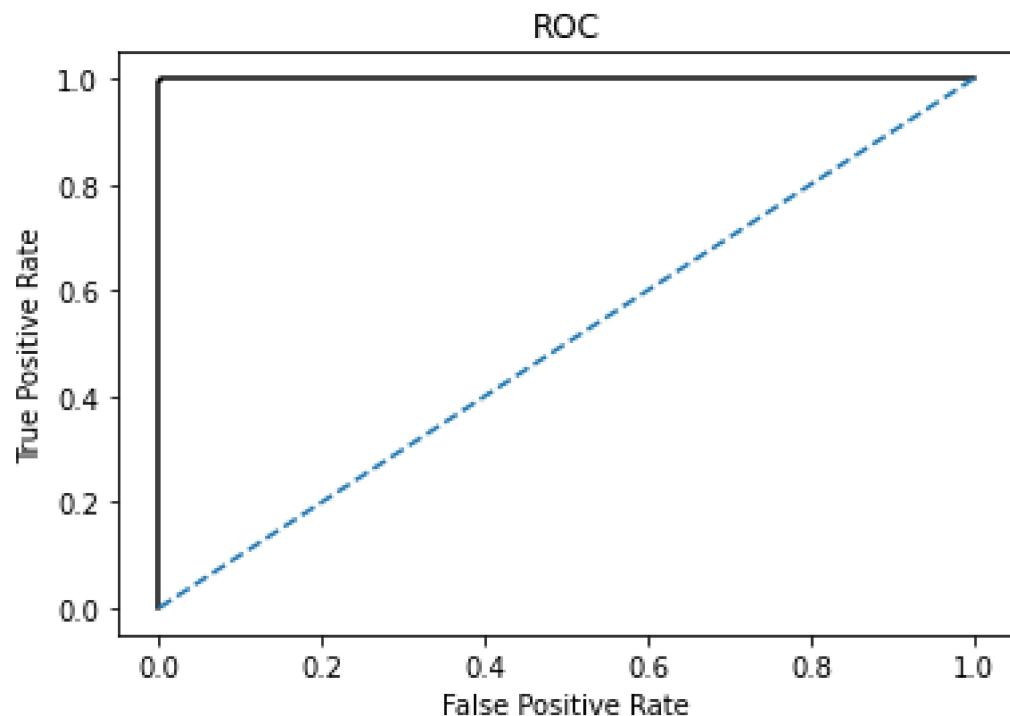


Fig-1.a.9

(9) Extra Tree Classifier

Train:- `confusion_matrix & classification_report , ROC-AUC Graph`

`confusion_matrix`

```
[6400, 0,
[1297, 0]]
```

`classification_report`

	precision	recall	f1-score	support
0	0.83	1.00	0.91	6400
1	0.00	0.00	0.00	1297
accuracy			0.83	7697
macro avg	0.42	0.50	0.45	7697
weighted avg	0.69	0.83	0.75	7697

ROC-AUC Graph

AUC: 0.5

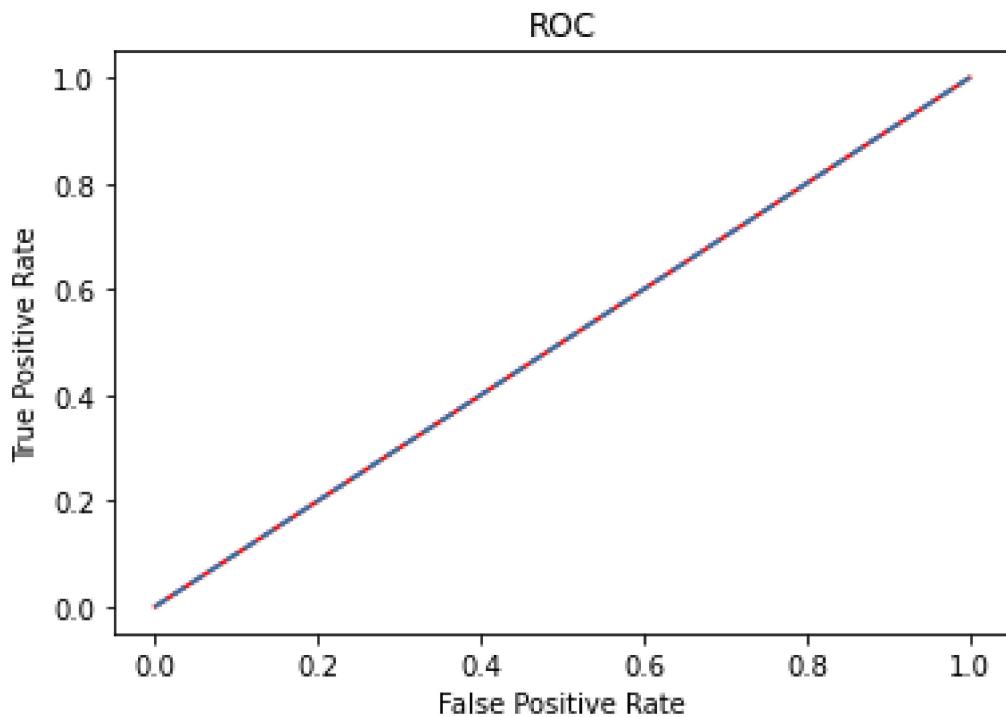


Fig-1.a.10

Conclusion | insight: all modal names and train modal confusion matrix and classification and accuracy are mentioned in all models.

1. The most important feature for determining customer churn for the random forest model is tenure i.e. for how long the customer is attached with the organization.

Why Random Forest chosen ?

Random Forest was chosen because this model outperforms rest all the models in all the performance metrics either it is accuracy, precision, f1-score, recall or AUC Score.

That's why the Random Forest model is chosen for the given dataset.

Efforts to improve model performance

1. We performed cross-validation 5 to find out the best parameters for the model as shown above to improve model performance.
2. We used the best parameters to build models and get the best performance metrics for the same.

Q 5. Model validation

Test your predictive model against the test set using various appropriate performance metrics.

Conclusion | insight:

(1) Logistics Regression

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2645, 100]
[ 270, 284]
```

classification_report

	precision	recall	f1-score	support
0	0.91	0.96	0.93	2745
1	0.74	0.51	0.61	554
accuracy			0.89	3299
macro avg	0.82	0.74	0.77	3299
weighted avg	0.88	0.89	0.88	3299

ROC-AUC Graph

AUC: 0.878

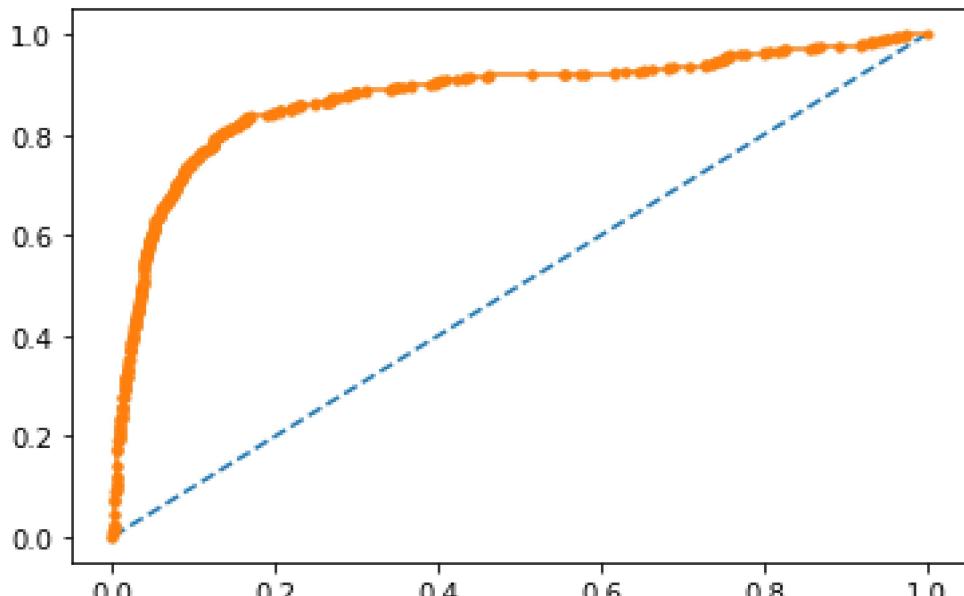


Fig-1.b

(2) Decision Tree

Test:- `confusion_matrix & classification_report , ROC-AUC Graph`

`confusion_matrix`

```
[2719,   26],
[ 106,  448]
```

`classification_report`

	precision	recall	f1-score	support
0	0.96	0.99	0.98	2745
1	0.95	0.81	0.87	554
accuracy			0.96	3299
macro avg	0.95	0.90	0.92	3299
weighted avg	0.96	0.96	0.96	3299

`ROC-AUC Graph`

AUC: 0.971

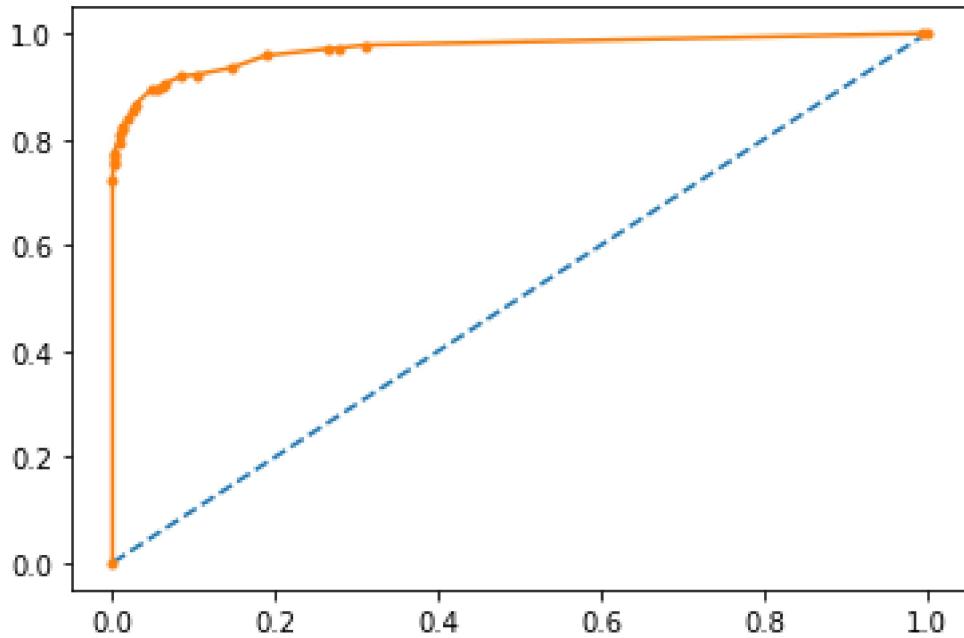


Fig-1.b.2

(3) Random Forest

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2719,   26]
[ 106, 448]
```

classification_report

	precision	recall	f1-score	support
0	0.96	0.99	0.98	2745
1	0.95	0.81	0.87	554
accuracy			0.96	3299
macro avg	0.95	0.90	0.92	3299
weighted avg	0.96	0.96	0.96	3299

ROC-AUC Graph

AUC: 0.978

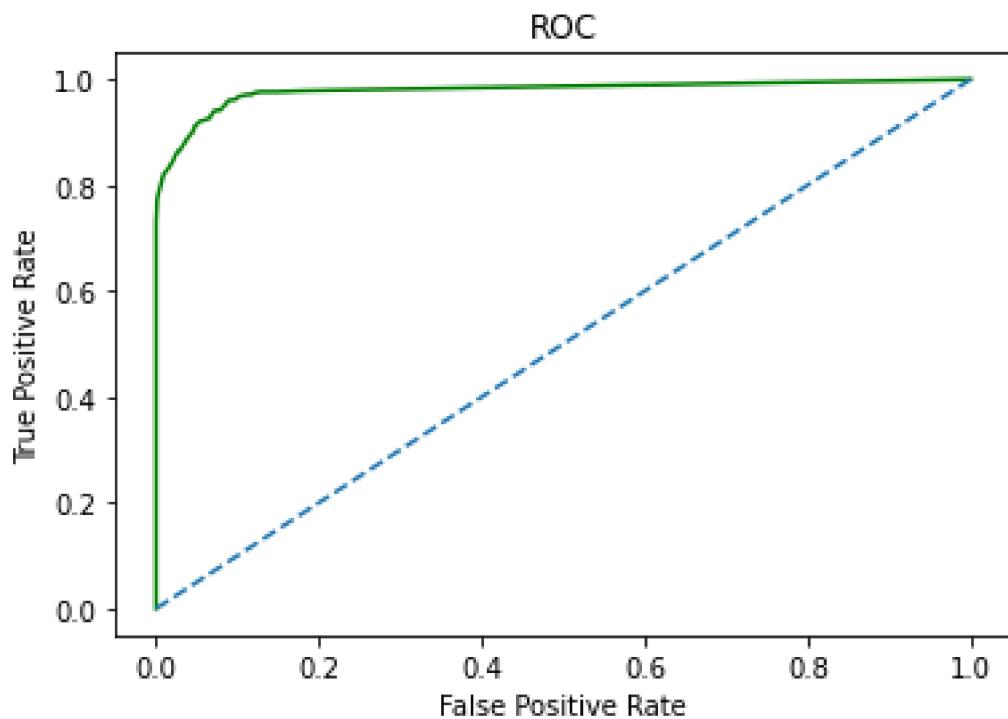


Fig-1.b.3

(4) Linear Discriminant Analysis

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[2567, 178],
 [232, 322]]
```

classification_report

	precision	recall	f1-score	support
0	0.92	0.94	0.93	2745
1	0.64	0.58	0.61	554
accuracy			0.88	3299
macro avg	0.78	0.76	0.77	3299
weighted avg	0.87	0.88	0.87	3299

ROC-AUC Graph

AUC: 0.971

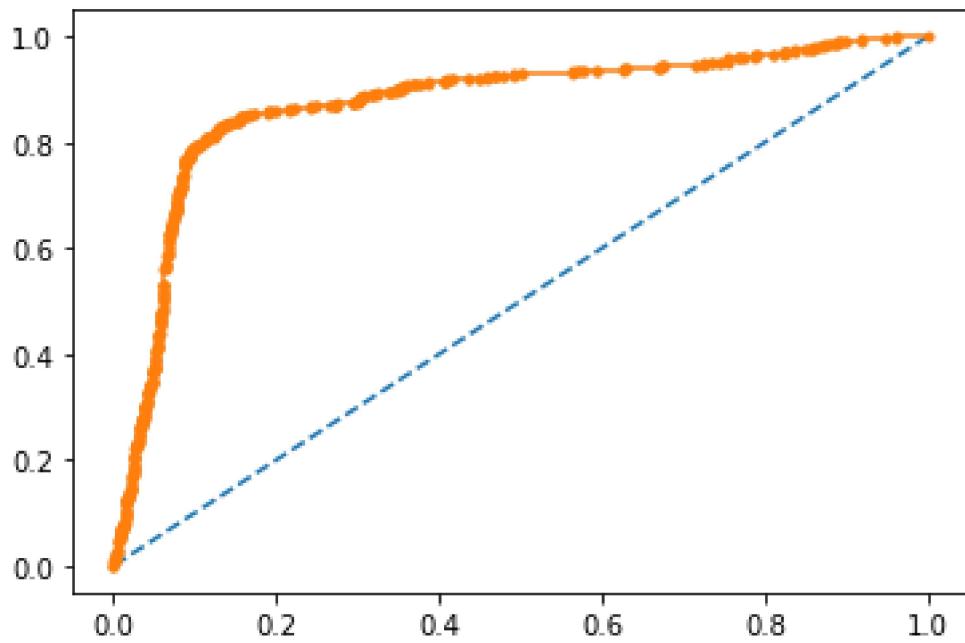


Fig-1.b.4

(5) K Nearest Neighbours

Test:- `confusion_matrix & classification_report , ROC-AUC Graph`

`confusion_matrix`

```
[2703, 42],  
[ 519, 35]]
```

`classification_report`

	precision	recall	f1-score	support
0	0.84	0.98	0.91	2745
1	0.45	0.06	0.11	554
accuracy			0.83	3299
macro avg	0.65	0.52	0.51	3299
weighted avg	0.77	0.83	0.77	3299

`ROC-AUC Graph`

AUC: 0.711

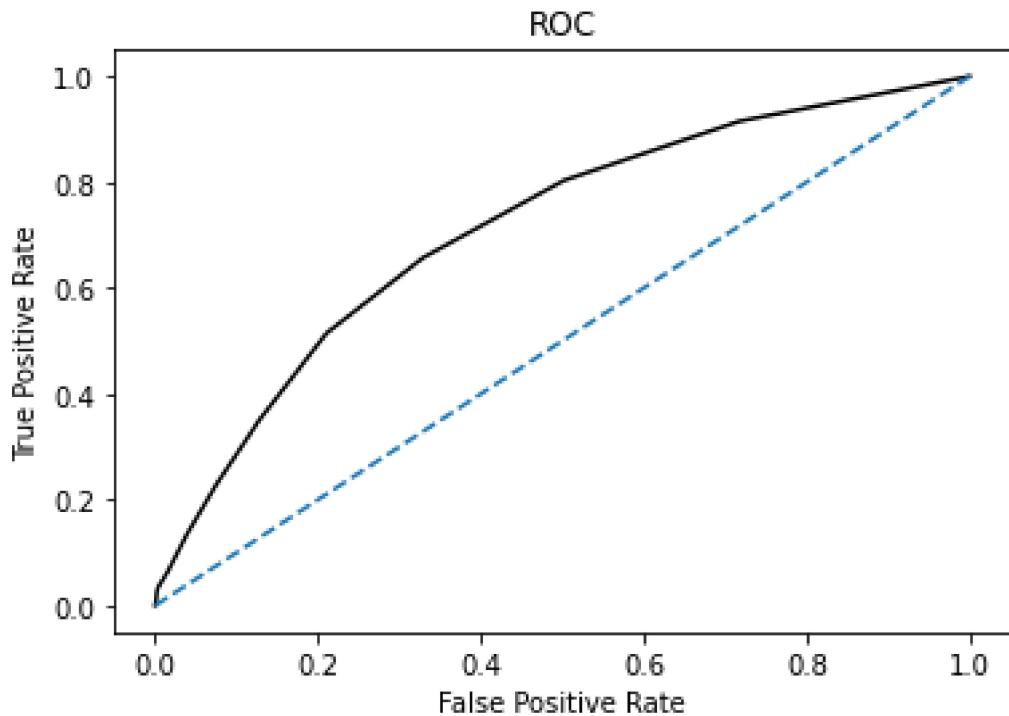


Fig-1.b.5

(6) Naive Bayes

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[ [2597, 148],
  [ 183, 371]],
```

classification_report

	precision	recall	f1-score	support
0	0.93	0.95	0.94	2745
1	0.71	0.67	0.69	554
accuracy			0.90	3299
macro avg	0.82	0.81	0.82	3299
weighted avg	0.90	0.90	0.90	3299

ROC-AUC Graph

AUC: 0.8504

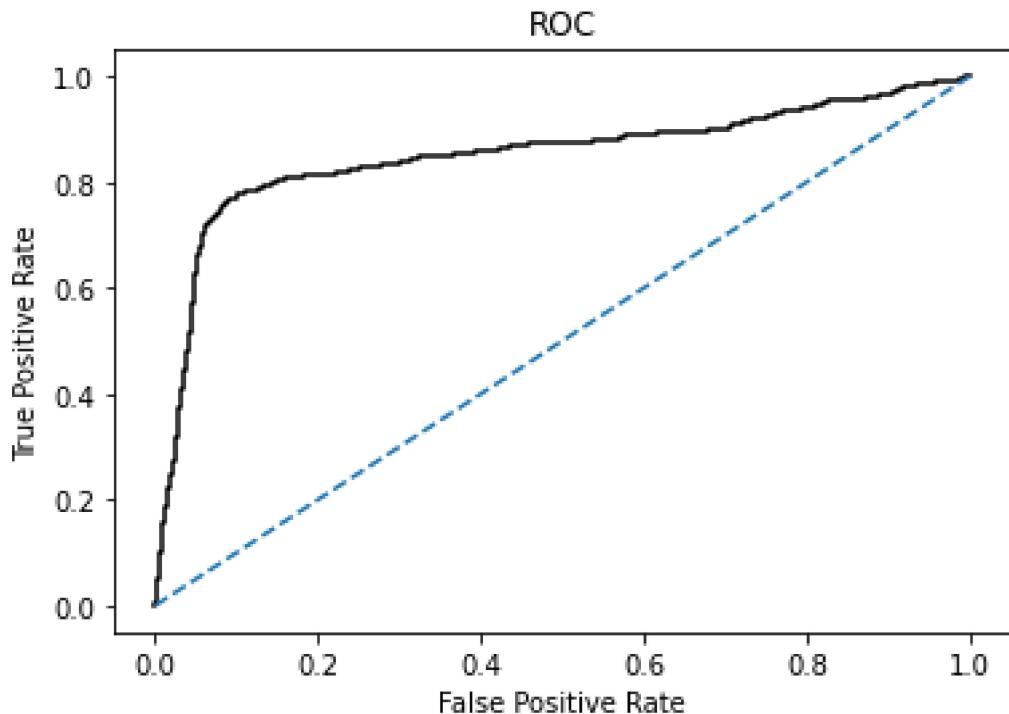


Fig-1.b.6

(7) Gradient Boosting

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2745,    0],
[ 554,    0]]
```

classification_report

	precision	recall	f1-score	support
0	0.83	1.00	0.91	2745
1	0.00	0.00	0.00	554
accuracy			0.83	3299
macro avg	0.42	0.50	0.45	3299
weighted avg	0.69	0.83	0.76	3299

ROC-AUC Graph

AUC: 0.5

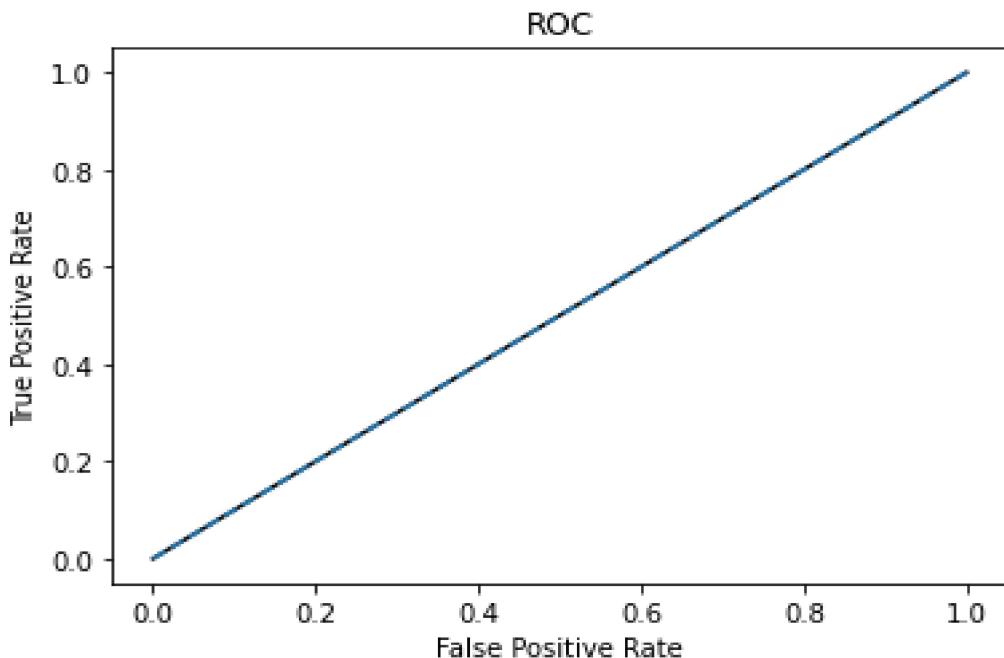


Fig-1.b.7

(8) Extreme Gradient Boosting

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2726,    19],
[   58,  496]]
```

classification_report

	precision	recall	f1-score	support
0	0.98	0.99	0.99	2745
1	0.96	0.90	0.93	554
accuracy			0.98	3299
macro avg	0.97	0.94	0.96	3299
weighted avg	0.98	0.98	0.98	3299

ROC-AUC Graph

AUC: 0.995

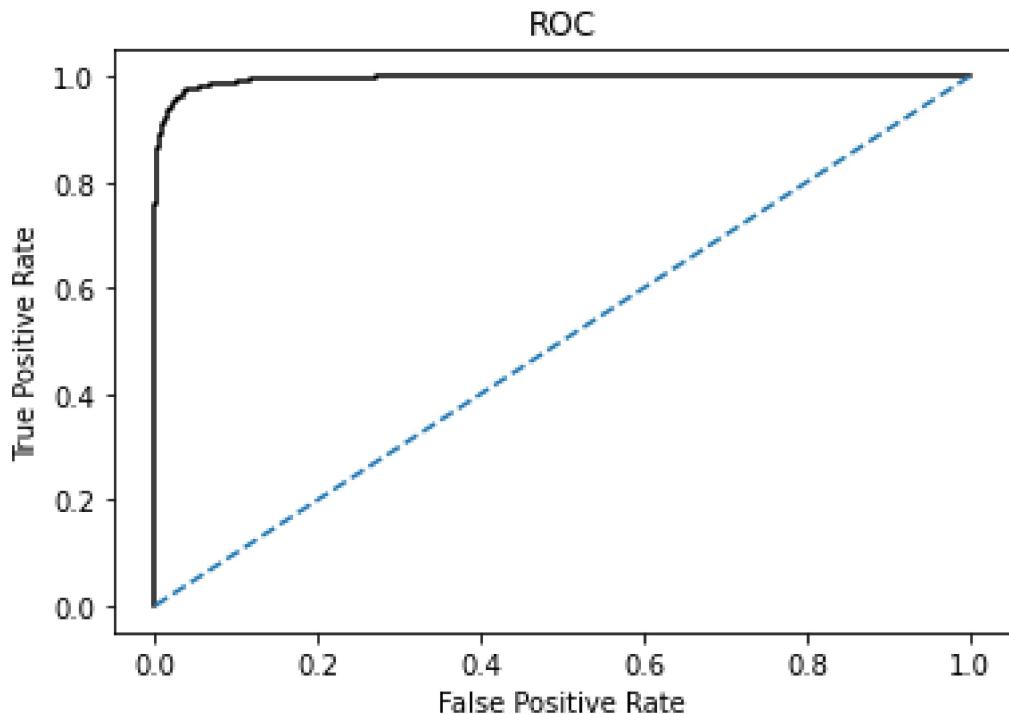


Fig-1.b.8

(9) Extra Tree Classifier

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[[2745,     0]
 [ 554,     0]]
```

classification_report

	precision	recall	f1-score	support
0	0.83	1.00	0.91	2745
1	0.00	0.00	0.00	554
accuracy			0.83	3299
macro avg	0.42	0.50	0.45	3299
weighted avg	0.69	0.83	0.76	3299

ROC-AUC Graph

AUC: 0.5

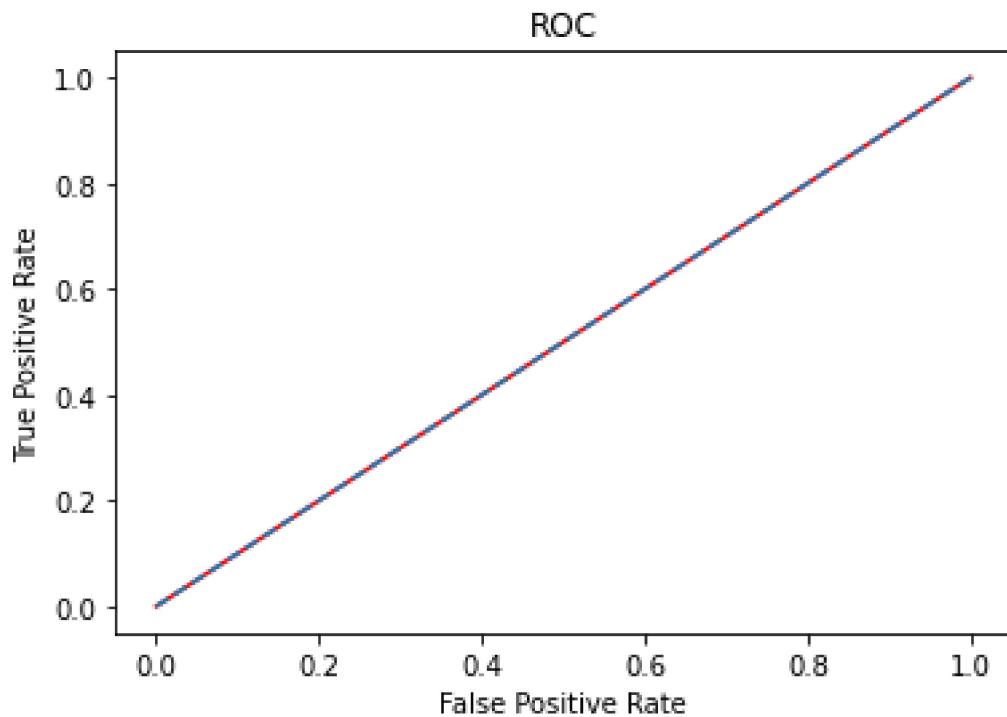


Fig-1.b.9

Conclusion | insight: All test modal are performed best modal results are RF and CART.

Extreme Gradient Boosting

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2726,    19],  
[  58,  496]]
```

classification_report

	precision	recall	f1-score	support
0	0.98	0.99	0.99	2745
1	0.96	0.90	0.93	554
accuracy			0.98	3299
macro avg	0.97	0.94	0.96	3299
weighted avg	0.98	0.98	0.98	3299

ROC-AUC Graph

AUC: 0.995

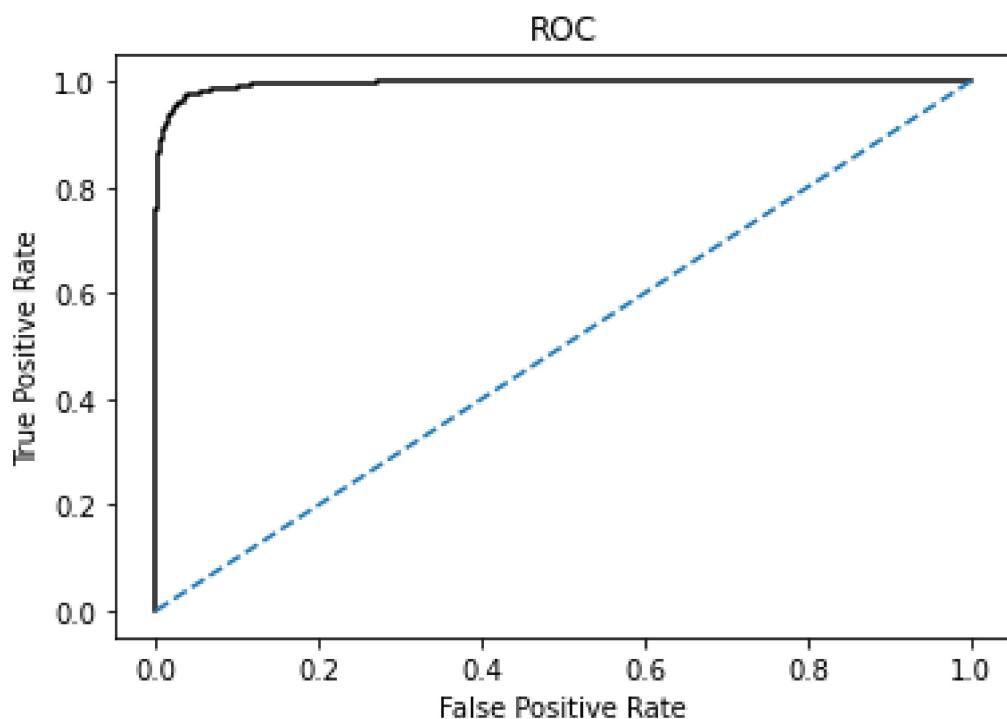


Fig-2.a.1

Gradient Boosting

Test:- confusion_matrix & classification_report , ROC-AUC Graph

confusion_matrix

```
[2745,    0],
[ 554,    0]]
```

classification_report

	precision	recall	f1-score	support
0	0.83	1.00	0.91	2745
1	0.00	0.00	0.00	554
accuracy			0.83	3299
macro avg	0.42	0.50	0.45	3299
weighted avg	0.69	0.83	0.76	3299

ROC-AUC Graph

AUC: 0.5

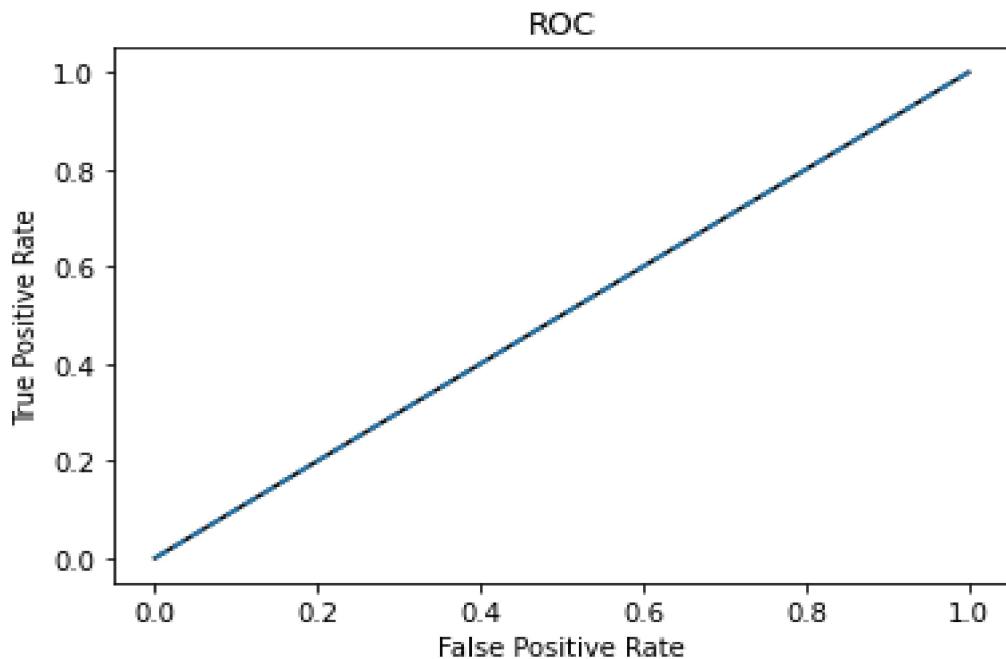


Fig-2.a.2

Conclusion | insight:

Model validation is done on the test dataset which is 30% of the complete data set. Not just accuracy but all the performance metrics are kept in mind before concluding the Random Forest as Best Model. Using Gradient Boosting & Extreme Gradient Boosting both models are performed but accuracy is low as compared to another model so not the best solution using boosting techniques.

Q6. Final interpretation / recommendation

Interpretation of the most optimum model and its implication on the business Conclusion | insight:

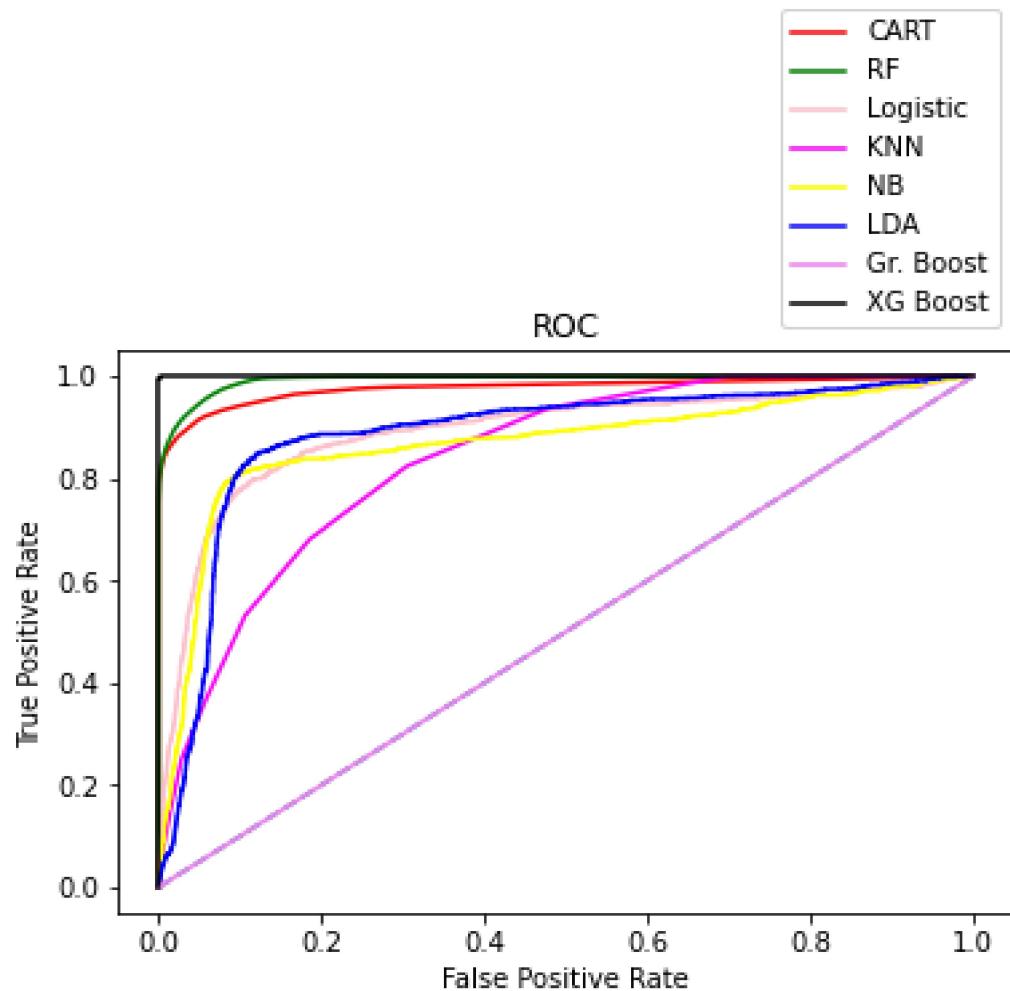
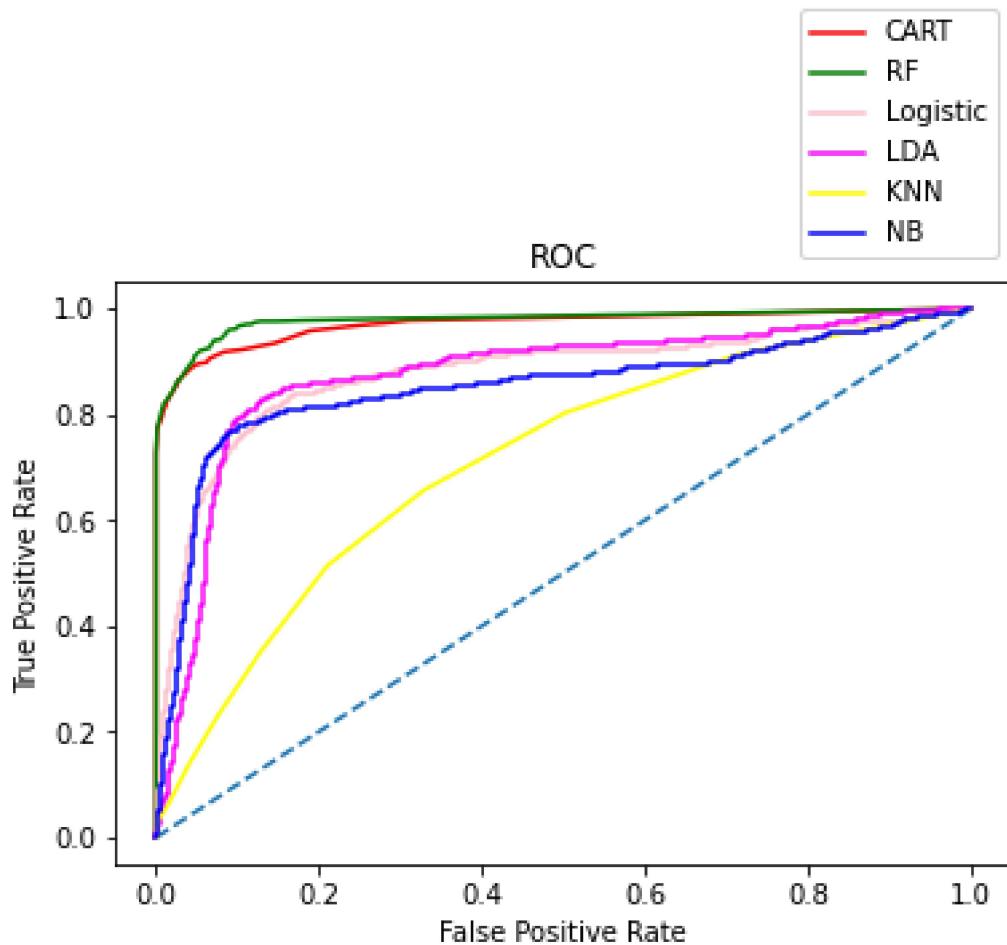


Fig-6.a

These predictions all Train Eight models need to achieve high AUC values. The sample data is divided into 70% for training and 30% for testing. We have applied feature engineering, effective feature transformation, and selection approach to make the features ready for machine learning algorithms. Best XG Boost, RF & CART AUC curve is best shown in campier all models.

**Fig-6.b**

This prediction all Test six models need to achieve high AUC values. The sample data is divided into 70% for training and 30% for testing. We have applied feature engineering, effective feature transformation, and selection approach to make the features ready for machine learning algorithms. Best Model RF & CART AUC curves are best shown in campier all models.

Train

	CART Train	RF Train	Log Train	LDA Train	KNN Train	NB Train	Gr.Boost Train	XG Boost Train	Extra Tree Train
Accuracy	0.97	0.97	0.89	0.88	0.84	0.90	0.83	1.00	0.83
AUC	0.98	0.99	0.89	0.89	0.84	0.87	0.50	1.00	0.50
Recall	0.84	0.84	0.54	0.62	0.15	0.69	0.00	0.99	0.00
Precision	0.96	0.96	0.75	0.65	0.68	0.70	0.00	1.00	0.00
F1 Score	0.90	0.90	0.63	0.63	0.24	0.70	0.00	0.99	0.00

Fig-6.c

Test

	CART Test	RF Test	Log Test	LDA Test	KNN Test	NB Test	Gr.Boost Test	XG Boost Test	Extra Tree Test
Accuracy	0.96	0.96	0.89	0.88	0.83	0.90	0.83	0.98	0.83
AUC	0.97	0.98	0.88	0.87	0.71	0.85	0.50	1.00	0.50
Recall	0.81	0.81	0.51	0.94	0.06	0.67	0.00	0.90	0.00
Precision	0.95	0.95	0.74	0.92	0.45	0.71	0.00	0.96	0.00
F1 Score	0.87	0.87	0.61	0.93	0.11	0.69	0.00	0.93	0.00

Fig-6.d

These algorithms are Logistics Regression, Decision Tree, Random Forest, Linear Discriminant Analysis, K Nearest Neighbours, Naive Bayes, Gradient Boosting, Extreme Gradient Boosting, and Extra Tree Classifier. The method of preparation and selection of features and entering the EDTH features had the biggest impact on the success of this model since the value of AUC reached 96.0%. Random Forest & CART model achieved the best results in all measurements. The AUC value was 96.0%. This comes in second place and the LDA & KNN and Decision Tree came third and fourth regarding AUC values. We have evaluated the models by fitting a new dataset related to different periods and without any proactive action from marketing, RF also gave the best result with 96.0% AUC. The decrease in result could be due to the non-stationary data model phenomenon, so the model needs training each period of time.

What we have understood from assessing the tenure is that there are quite a lot of customers who are new to the

Platform. This is a good opportunity for the company to bring in new strategies and work upon increasing Customer retention.

The number of people that have contacted customer care is high. The number of times CC has been contacted by a single account is high as 132, this is an alarming high. We have noted that comparatively 3 to 4 users per account are higher in number. This is helping us to understand. That the products and the offers we need to plan to offer should be crafted keeping in mind that it is a for widest age distribution. Better offers around less-priced subscription plans will help us to retain new customers, increase customer numbers. The lower-priced subscription plans are more favoured by the customers. The revenue growth is seeming bright with mean and median around 15%.

Impurity-based Importance train value

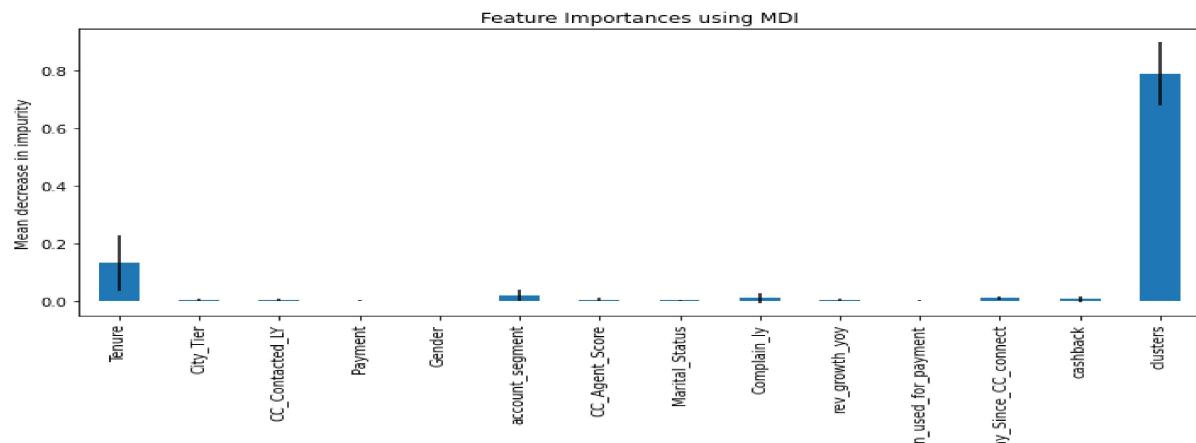


Fig-2.c.1

Permutation-based Importance test value

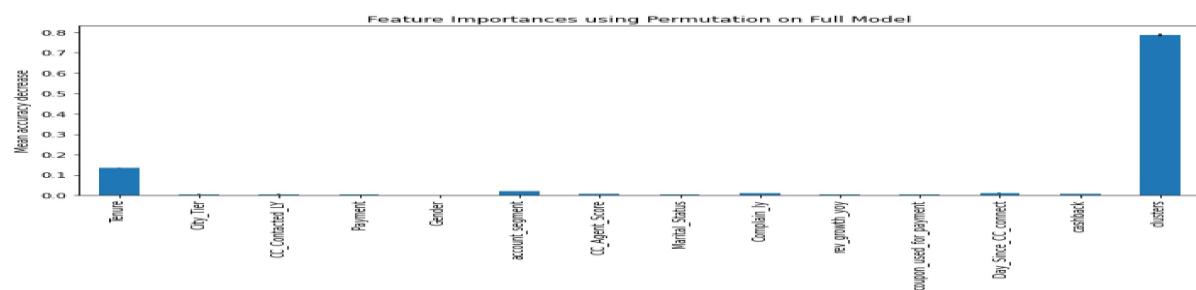


Fig-2.c.2

SHAP Value-based Importance

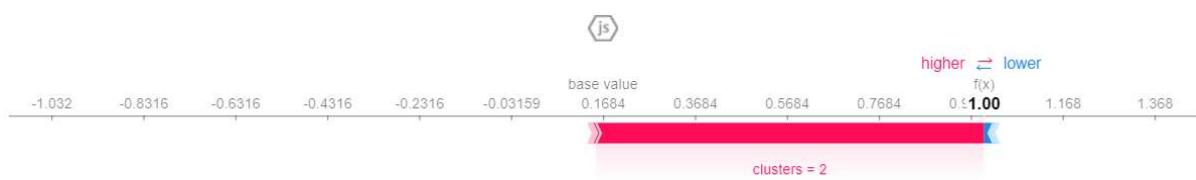


Fig-2.c.3

EL15 Feature Importance

Weight	Feature	Contribution?	Feature
0.7865 ± 0.2209	x13	+0.832	<BIAS>
0.1334 ± 0.1917	x0	+0.158	x0
0.0198 ± 0.0370	x5	+0.010	x8
0.0119 ± 0.0124	x11	+0.008	x2
0.0112 ± 0.0335	x8	+0.006	x6
0.0071 ± 0.0170	x12	+0.005	x1
0.0065 ± 0.0103	x6	+0.004	x11
0.0054 ± 0.0065	x2	+0.002	x3
0.0050 ± 0.0084	x1	+0.001	x10
0.0039 ± 0.0061	x9	-0.000	x7
0.0033 ± 0.0068	x7	-0.004	x5
0.0027 ± 0.0057	x3	-0.006	x12
0.0025 ± 0.0047	x10	-0.036	x13
0.0009 ± 0.0023	x4	-0.043	x9

y=0 (probability 0.934) top features

Contribution?	Feature	Value
+0.832	<BIAS>	1.000
+0.158	Tenure	32.000
+0.010	Complain_ly	0.000
+0.008	CC_Contacted_LY	8.000
+0.006	CC_Agent_Score	3.000
+0.005	City_Tier	1.000
+0.004	Day_Since_CC_connect	8.000
+0.002	Payment	2.000
+0.001	coupon_used_for_payment	1.000
-0.000	Marital_Status	0.000
-0.004	account_segment	5.000
-0.006	cashback	1913.000
-0.036	clusters	1.000
-0.043	rev_growth_yoy	16.000

y=0 (probability 0.934) top features

Contribution?	Feature
+0.832	<BIAS>
+0.158	x0
+0.010	x8
+0.008	x2
+0.006	x6
+0.005	x1
+0.004	x11
+0.002	x3
+0.001	x10
-0.000	x7
-0.004	x5
-0.006	x12
-0.036	x13
-0.043	x9

y=0 (probability 1.000) top features

Contribution?	Feature	Value
+0.832	<BIAS>	1.000
+0.116	clusters	3.000
+0.024	Tenure	11.000
+0.016	account_segment	4.000
+0.006	Day_Since_CC_connect	2.000
+0.002	cashback	398.000
+0.002	Complain_ly	0.000
+0.002	coupon_used_for_payment	0.000
+0.001	Gender	1.000
+0.001	Payment	1.000
+0.000	rev_growth_yoy	7.000
+0.000	CC_Agent_Score	3.000
+0.000	City_Tier	1.000
+0.000	CC_Contacted_LY	16.000
-0.000	Marital_Status	0.000

Conclusion | insight:

Fig-2.c.4

Recommendations

We can suggest the management segment customers based on 'City-Tiers'. The reasons for this suggestion are, Clear demographic divide, Very distinct behavior within each tier across multiple parameters. Each tier can be approached with local and custom offers. Marketing can be done both on online and offline platforms across clearly defined tiers.

Insights based on EDA

The importance of this type of research in the telecom market is to help companies make more profit. It has become known that predicting churn is one of the most important sources of income for telecom companies. Hence, this research aimed to build a system that predicts