# CUSTOMER CHURN PREDICTION

## Capstone Project – Final Report

## Abstract

This is a binary classification problem to predict customers who would churn

## Swetha Kunapuli

PGP-DSBA June B – 21 Batch

# Contents

## List of tables

# List of figures

# 1. Introduction - Business problem definition

## 1.1.    Defining problem statement

A DTH 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 potential churners. In this company, account churn is a major issue because 1 account can have multiple customers. Hence by losing one account the company might be losing more than one customer.

The current project is aimed at developing a churn prediction model for this company and to provide business recommendations on a campaign focused on retaining customers. The campaign recommendation should be such that it does not entail a huge cost for retention of customers and should remain within a budget earmarked for this purpose.

## 1.2.    Need of the study/project

A DTH provider's biggest cost is cost of acquisition of a new customer. A customer thus acquired, will need to be retained for quite a few years so that the initial cost of acquisition is recovered back and that particular account is profitable1. Due to this reason, customer churn directly impacts the profitability of a DTH operator. DTH providers also are in a constant pressure to increase their customer base to maintain their profitability as most of them have a fixed broadcaster/content provider fee irrespective of the number of customers in their customer base2. So, more the number of customers, greater their profitability. Hence it becomes very important to not only increase the customer base but also protect the current customer base.

Acquiring a new customer can cost five times more than retaining an existing customer. Increasing customer retention by 5% can increase profits from 25-95%.

As customer churn directly impacts both the top-line and bottom-line revenue of the business, existing customer base needs to be protected. Providing all customers with offers to retain them would make a dent in the profitability and hence it is very important to focus only on select set of customers who are at a higher risk of churning.

# 2. EDA and Business Implication

## 2.1.    Data collection

- ➢ The dataset contains 'Account' level data which is master data along with features of account such as gender, marital status, city tier, account user count of primary account holder, whether the account is live or churned, segment the account belongs to. It also possibly contains certain derived features such as tenure (which probably could have been derived from account open date).
- ➢ The dataset also contains information taken/derived from transaction data and rolled up at Account level and given as a feature for the account – for e.g., Number of days since none of the account holders contacted customer care, monthly average cash back for last 12 months, number of complaints made last year, number of times customer care contacted last year, revenue per month in last 12 months, how many

times customers have used coupons to pay in last 12 months, satisfaction score and customer service score.
- ➢ Most of the transaction data roll-up at account level has been done for 12 months (previous year). However, the revenue growth percentage has been taken for last year in comparison with previous one year which implies that 24-months' worth of data has been used to calculate this field.
- ➢ As the dataset has been provided, the methodology used by customer to extract data is not known. The frequency at which this dataset is extracted has also not been specified.

## 2.2. Visual inspection of data

| | AccountID | Churn | Tenure | City_Tier | CC_Contacted_LY | Payment | Gender | Service_Score | Account_user_count | account_seg |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 20000 | 1 | 4 | 3.0 | 6.0 | Debit Card | Female | 3.0 | 3 | S |
| 1 | 20001 | 1 | 0 | 1.0 | 8.0 | UPI | Male | 3.0 | 4 | Regular |
| 2 | 20002 | 1 | 0 | 1.0 | 30.0 | Debit Card | Male | 2.0 | 4 | Regular |
| 3 | 20003 | 1 | 0 | 3.0 | 15.0 | Debit Card | Male | 2.0 | 4 | S |
| 4 | 20004 | 1 | 0 | 1.0 | 12.0 | Credit Card | Male | 2.0 | 3 | Regular |

```
(11260, 19)
```

*Table 1 - Head of the dataset*

- ➢ The dataset has 11260 rows and 19 columns.

```
RangeIndex: 11260 entries, 0 to 11259
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)
```

*Table 2 - Columns and data types*

- ➢ There are 5 columns of float type, 2 columns of integer type and 12 columns of object type.
- ➢ There are several columns that are supposed to be read as numeric, instead they have been read as object type for e.g., Tenure is a numeric field but has been read

as object. Those columns need to be checked for special characters and need to be treated before the column can be changed to numeric for further processing.
- ➢ There are no duplicate rows in the data set. Each account id has one unique row.
- ➢ Several columns have null values.

## 2.3.  Understanding the attributes

The following table shows the attribute names, their description and the kind of values that they contain. Although some of the variable names are slightly long, they do not have blanks or special characters in them. Hence, it has been decided to let the current column names stay as-is as they are self-explanatory and would be easy to understand and interpret when seen in the plots as part of univariate and bivariate analysis. The variable names would be changed later to shorten or make it uniform when one hot encoding is done in a later section.

| S.no | Column | Column Description | Data Description |
|---|---|---|---|
| 1 | AccountID | account unique identifier | Unique ID. Hence, it will not be used in modelling |
| 2 | Churn | account churn flag (Target) | Target variable. Contains 1 for churned and 0 for non-churned |
| 3 | Tenure | Tenure of account | Continuous field. Contains values ranging from 0 to 99 |
| 4 | City_Tier | Tier of primary customer's city | Categorical ordinal - values 1,2,3 |
| 5 | CC_Contacted_LY | How many times all the customers of the account have contacted customer care in last 12months | Continuous field. Contains values ranging from 4 to 132 |
| 6 | Payment | Preferred Payment mode of the customers in the account | Categorical nominal - values Credit card, debit card, E wallet, UPI, Cash on Delivery |
| 7 | Gender | Gender of the primary customer | Categorical nominal - values Male, Female, M and F (M and F need to be converted to Male and Female) |
| 8 | Service_Score | Satisfaction score given by customers of the account on service provided by company | Categorical ordinal - values 0 to 5 |
| 9 | Account_user_count | Number of customers tagged with this account | Limited range. Can be treated as categorical - values 1 to 6 |
| 10 | account_segment | Account segmentation on the basis of spend | Categorical nominal - values HNI, Regular, Regular Plus, Super, Super plus and variations with + |
| 11 | CC_Agent_Score | Satisfaction score given on customer care service provided | Categorical ordinal - values 1 to 5 |
| 12 | Marital_Status | Marital status of primary customer | Categorical nominal - contains values Married, Single and Divorced |
| 13 | rev_per_month | Monthly average revenue from account in last 12 months | Continuous field. Contains values ranging from 1 to 140 |
| 14 | Complain_ly | Complaints raised by account in last 12 months | Categorical - 0 (for no) or 1 (for yes) |
| 15 | rev_growth_yoy | revenue growth percentage of the account | Continuous field. Contains values ranging from 4 to 28 |

| | | (last 12 months vs last 24 to 13 month) | |
|---|---|---|---|
| 16 | coupon_used_for_payment | How many times customers have used coupons to do the payment in last 12 months | Continuous field, but with limited range. Contains values ranging from 0 to 16 |
| 17 | Day_Since_CC_connect | Number of days since no customers in the account has contacted the customer care | Continuous field. Contains values ranging from 0 to 47 |
| 18 | Cashback | Monthly average cash back generated by account in last 12 months | Continuous field. Contains values ranging from 0 to 1997 |
| 19 | Login_device | Preferred login device of the customers in the account | Categorical nominal - contains values Mobile, Computer |

*Table 3 - Attribute Description*

**Descriptive Statistics:**

| | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| Churn | 11260.0 | 0.168384 | 0.374223 | 0.0 | 0.00 | 0.00 | 0.00 | 1.0 |
| Tenure | 11042.0 | 11.025086 | 12.879782 | 0.0 | 2.00 | 9.00 | 16.00 | 99.0 |
| City_Tier | 11148.0 | 1.653929 | 0.915015 | 1.0 | 1.00 | 1.00 | 3.00 | 3.0 |
| CC_Contacted_LY | 11158.0 | 17.867091 | 8.853269 | 4.0 | 11.00 | 16.00 | 23.00 | 132.0 |
| Service_Score | 11162.0 | 2.902526 | 0.725584 | 0.0 | 2.00 | 3.00 | 3.00 | 5.0 |
| Account_user_count | 10816.0 | 3.692862 | 1.022976 | 1.0 | 3.00 | 4.00 | 4.00 | 6.0 |
| CC_Agent_Score | 11144.0 | 3.066493 | 1.379772 | 1.0 | 2.00 | 3.00 | 4.00 | 5.0 |
| rev_per_month | 10469.0 | 6.362594 | 11.909686 | 1.0 | 3.00 | 5.00 | 7.00 | 140.0 |
| Complain_ly | 10903.0 | 0.285334 | 0.451594 | 0.0 | 0.00 | 0.00 | 1.00 | 1.0 |
| rev_growth_yoy | 11257.0 | 16.193391 | 3.757721 | 4.0 | 13.00 | 15.00 | 19.00 | 28.0 |
| coupon_used_for_payment | 11257.0 | 1.790619 | 1.969551 | 0.0 | 1.00 | 1.00 | 2.00 | 16.0 |
| Day_Since_CC_connect | 10902.0 | 4.633187 | 3.697637 | 0.0 | 2.00 | 3.00 | 8.00 | 47.0 |
| cashback | 10787.0 | 196.236370 | 178.660514 | 0.0 | 147.21 | 165.25 | 200.01 | 1997.0 |

*Table 4 – Basic descriptive statistics*

➤ Tenure seems to have a very huge range up to 99. It could be in months in which case those would be valid values.
➤ The maximum limit for many of the variables seems to be very far apart from the 75th percentile for many variables such as cash back, revenue per month and customer contacted last year. There seems to be significant positive skew in these variables. A look at the boxplot and histogram will confirm the presence of outliers.

## 2.4.   Univariate  analysis

Univariate analysis is done for the purpose of observing distribution and spread for every continuous attribute and distribution of data in categories for categorical ones. It has been done by observing:

➤ Box plots and histograms for continuous variables.

➢ Count plots for categorical variables.

## Continuous data – Box plot



*Figure 1 – Box plot for numeric variables*

## Continuous data – Histogram

*Figure 2 – Histogram for numeric variables*

|  | Skewness |
|---|---|
| Tenure | 3.90 |
| CC_Contacted_LY | 1.42 |
| rev_per_month | 9.09 |
| rev_growth_yoy | 0.75 |
| coupon_used_for_payment | 2.58 |
| Day_Since_CC_connect | 1.27 |
| cashback | 8.77 |

*Table 5 – Numeric variables skewness*

**Observations:**

➢ All numeric variables with the exception of rev_growth_yoy have outliers. Some outliers for certain variables are closer to the whisker, whereas there are a group of outliers that are far beyond the whisker with no in between values. For instance, rev_per_month has a huge space between 30 and 100 indicating absence of values in that range. Those outliers in the extreme values do not correlate with corresponding outliers in cashback field. We cannot rule out these outliers as incorrect values, they may belong to hotels with many rooms. But models like logistic regression are sensitive to outliers and may not give good performance if outliers are left untreated.

➢ Hence, we will follow two approaches to modelling – one set of data with outliers treated for outlier sensitive models and other set of data with outliers not treated (left as-is) for outlier resistant models such as Random forest.

➢ Coupon_used_for_payment has a very limited range 0 to 16. Hence, for the purpose of this analysis, the outliers will not be treated (similar to categorical variables).

➢ All numeric variables with the exception of rev_growth_yoy have a high positive skew.

**Categorical fields – Count plot**

*Figure 3 – Count plot for categorical variables*



```
0    83.2
1    16.8
Name: Churn, dtype: float64
```

*Figure 4– Count plot of target variable*

**Observations:**

➢ This is an imbalanced dataset with target variable containing 16.8% churn.
➢ Tier 1 cities have more accounts followed by Tier 3 cities.

➢ Most of the accounts pay through debit card followed by credit card. UPI ranks last amongst payment methods.
➢ Number of male account holders outnumbers females.
➢ Regular plus and Super are top two account segment types by number.
➢ Top score for both Customer service agent and Service score are 3.
➢ Married customers have the most accounts followed by single.
➢ Most accounts do not have a customer complaint filed last year.
➢ Most account holders use Mobile for logging and using services.

## 2.5.    Bivariate analysis & Multivariate analysis

Bivariate analysis shows the relationship between two variables. Here, the predictor variables have been taken and their relationship with the target variable has been plotted. The influence of the predictor variables on the target variable can be observed in these bivariate plots.

**Continuous predictor variables that show some ability to separate the target variables**



*Figure 5 - Histogram and Box plots for Tenure vs Churn*



*Figure 6 - Histogram and Box plots for Day_since_CC_connect vs Churn*



*Figure 7 - Histogram and Box plots for CC_Contacted_LY vs Churn*

**Continuous predictor variables that isn't able to show a clear separation between target classes**



*Figure 8 - Histogram and Box plots for Coupn_used_for_payment vs Churn*



*Figure 9 - Histogram and Box plots for Cashback vs Churn*



*Figure 10 - Histogram and Box plots for rev_per_month vs Churn*

*Figure 11 - Histogram and Box plots for rev_growth_yoy vs Churn*

**Observations:**

➤ From the above plots we can see that variables such as Tenure, Days_since_CC_connect have some influence on the target variable.

➤ The median lines for churned and Non-churned observations when plotted against these variables show a difference. Whereas, the medians in boxplots for variables such as coupon_used_for_payment, rev_growth_yoy do not show much difference in churned vs. non-churned distributions.

**Continuous variables: Anova (Analysis of Variance)**

Analysis of Variance is a statistical method, used to check if the means of two or more groups that are significantly different from each other. Hypothesis for the test is as follows:

*H0: Means of all groups are equal*
*Ha: At least means of one pair of the groups is differen*t

Stats model library was used to perform the Anova test.

**Results of Anova test for following variables and churn**

| Variable | F - statistic | Probability of > F | Inference at significance level of 5% |
|---|---|---|---|
| Tenure | 634.6 | 3.30E-36 | Reject null hypothesis. The means are different. Variable significant to model building |
| CC_Contacted_LY | 58.25 | 2.49E-14 | Reject null hypothesis. The means are different. Variable significant to model building |
| rev_per_month | 5.32 | 0.021 | Reject null hypothesis. The means are different. Variable significant to model building |
| rev_growth_yoy | 2.17 | 0.141 | Cannot reject null hypothesis. The means are equal. Variable can be dropped |
| coupon_used_for_payment | 2.47 | 0.116 | Cannot reject null hypothesis. The means are equal. Variable can be dropped |
| Day_Since_CC_connect | 243.9 | 2.10E-54 | Reject null hypothesis. The means are different. Variable significant to model building |
| Cashback | 11.32 | 0.001 | Reject null hypothesis. The means are different. Variable significant to model building |

*Table 6 – Results of Anova test*

**Observations:**

At a significance level of 0.05 (5%), the tests for the variables rev_growth_yoy and coupon_used_for_payment have given a p-value of greater than 0.05. In these two cases, H0 cannot be rejected, i.e., the means for the two groups churn=0 and churn=1 for these variables are the same. This implies that since the groups are not too different, these two variables cannot be significant predictors of the target variable. This is in line with what was visually observed using the bivariate box plots for these two variables.

**Bivariate plots for categorical variables vs. churned**



*Figure 12 – Stacked bar chart for categorical variables*

**Observations:**

In the above stacked bar charts, active customers are called current customers. The first bar shows distribution of the categorical predictor variable being analysed within churned customer and the second bar shows distribution within active/current customers. Absolute comparison cannot be done as this is a stacked bar with the height of the bar representing 100% of churned and non-churned/current customers respectively. The 18 respective absolute counts differ. Only interpretation from these plots is the % distribution between various categories of the categorical variable being plotted.

- From the above charts, some of the variables like city_tier, account_segment, Complain_ly, Marital_Status, CC_Agent_score seem to show a difference in distribution when churned vs current customers are considered. These variables may have an influence over the target variable churn and may contribute to the model.
- Other variables such as Gender and Service_Score have more or less similar distributions within Churned and Current/active customer bars. These variables may not have significant contribution towards the model.

**Categorical variables: Chi-squared test of independence at significance level 0.05**

These variables were subjected to Chi square test of independence to decide if they are statistically significant enough to be included in the model or not. The Chi-Square test of independence is used to determine if there is a significant relationship between two categorical variables. The frequency of each category for one categorical variable is compared across the categories of the second categorical variable. The data can be displayed in a contingency table where each row represents a category for one variable and each column represents a category for the other variable.

> **Null hypothesis:** *There is no relationship between the two categorical variables.*
> **Alternative hypothesis:** *There is a relationship between the two categorical variables.*

| | Variable | chi2 | p-value | chi2_output |
|---|---|---|---|---|
| 0 | Gender | 8.983146 | 2.724812e-03 | Reject Ho; Dependent. |
| 1 | Service_Score | 18.414690 | 2.469166e-03 | Reject Ho; Dependent. |
| 2 | City_Tier | 80.288817 | 3.677095e-18 | Reject Ho; Dependent. |
| 3 | Payment | 103.799617 | 1.526348e-21 | Reject Ho; Dependent. |
| 4 | Account_user_count | 154.959445 | 1.173574e-31 | Reject Ho; Dependent. |
| 5 | account_segment | 567.068402 | 2.073937e-121 | Reject Ho; Dependent. |
| 6 | CC_Agent_Score | 139.031565 | 4.549521e-29 | Reject Ho; Dependent. |
| 7 | Marital_Status | 379.808123 | 3.355165e-83 | Reject Ho; Dependent. |
| 8 | Complain_ly | 688.084739 | 1.166239e-151 | Reject Ho; Dependent. |

*Table 7 – Results of Chi-Squared test*

A p-value less than 0.05 (typically ≤ 0.05) is statistically significant. It indicates strong evidence against the null hypothesis, as there is less than a 5% probability the null is correct. Since the p-value returned for all the categorical variables is less than 0.05, the null hypothesis can be rejected. Hence at 5% level of significance, it may be concluded that churn is not independent of these categorical variables. Hence, we will proceed to retain all these categorical predictor variables at this point.

**Pair plots for the numeric variables with hue set as target variable**



*Figure 13 – Pair plot for predictor numeric variables*

**Observations:**

➢ Some of the variables clearly show presence of some clusters for e.g., rev_per_month and Day_Since_CC_connect.
➢ The diagonal kde plot for Tenure shows a slight separation with customers who churned falling on the lower side of Tenure.
➢ There is no linear relationship between any two continuous variables.

**Correlation heat map**

The following heat map shows the correlation (Pearson's) between various numeric predictors in the dataset. The purpose of doing a correlation matrix is to check if any of the variables have a strong correlation that can contribute to multi-collinearity. If there is a strong correlation > 0.80 between two variables, one of them can be dropped as it would not add much to the model's performance.

*Figure 14 – Correlation heat map*

**Observation:** None of the numeric variables show a strong correlation. Hence there is no need to drop any variable.

## 2.6.    Data imbalance

As can be seen from the plot below, the data is imbalanced with respect to the target variable which is the indicator whether customer has churned or not. The distribution given below shows that for every 100 customers acquired by the business, 17 have churned and 83 customers are active. This distribution is skewed towards active/current customers. The objective for this exercise is to be able to predict the customers who would churn i.e., the minority class '1'.

Count plot of target variable - Churn

```
0    83.2
1    16.8
Name: Churn, dtype: float64
```

*Figure 15 – Count plot of target variable*

**Observations:**

- ➢ If a dataset has equal distribution amongst the categorical values taken by the target variable (Churn), it would be called a balanced dataset. In a balanced dataset, the model learns to predict both classes with equal efficiency as there are equal number of observations.
- ➢ In case of customer churn, any business would have less of churned customers and more of active customers. Similarly, this dataset also has close to 17% churned customers. Although this mimics real-life churn data, from a modelling perspective, this may pose a challenge.

**Challenges posed by an imbalanced dataset:**

- ➢ If there are no sufficient observations in the minority class, the model would be unable to learn the patterns in minority class (class of interest) well and may predict majority class better. Using resampling techniques such as oversampling minority class or under-sampling majority class may be beneficial. Compared to under sampling majority class which would result in loss of data, oversampling using techniques like SMOTE (Synthetic Minority Oversampling) may help. This technique would help generate synthetic data for churned customers based on the existing churned customer observations. However, this technique may not always guarantee better model performance. Depending on the algorithm used and the type of SMOTE used, there could be over fitting in the train dataset.
- ➢ Accuracy as an evaluation metric would not be appropriate in imbalanced classification problems. Even if the algorithm predicts all customers as belonging to majority class, it would still result in an accuracy of 83.2%. Using Precision, Recall or F1-score for the minority class may be a better approach for evaluation.

## 2.7. Clustering

> Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar to each other than to those in other groups (clusters).

> Clustering is an unsupervised task and given all the features in the dataset, the clustering algorithm is allowed to group customers such that each group has similar customers and customers of different groups are dissimilar.

> For this purpose, the processed dataset (cleaned, scaled, nulls imputed, outlier treated, categorical features encoded) without the target variable was used.

> As the dataset contained both continuous and categorical predictor variables, kprototypes function from Kprototypes library was used.

> The algorithm was run for 2, 3 and 4 clusters and 3 clusters seemed to have the best separation. The cluster profile was formed by grouping the observations by clusters and finding the mean for all the features.

> Although churn was not part of the features for clustering, it was added part of the cluster profile so that it is possible to appreciate how churn varies for each cluster.

| kproto_3clusters | Churn | Tenure | City_Tier | CC_Contacted_LY | Service_Score | User_Count | CC_Score | Rev_Permonth |
|---|---|---|---|---|---|---|---|---|
| 0 | 0.187416 | 9.499015 | 1.600000 | 13.305577 | 2.855501 | 3.675453 | 3.056379 | 6.323737 |
| 1 | 0.096765 | 13.621100 | 1.715588 | 14.997355 | 2.925544 | 3.669191 | 3.102383 | 6.149150 |
| 2 | 0.223943 | 10.643872 | 1.679526 | 30.488770 | 2.964625 | 3.757564 | 3.039710 | 6.714170 |

| Complain_LY | Days_Since_CC | ... | Payment_UPI | Gender_Male | ACSegment_Regular | ACSegment_Regularplus | ACSegment_Super |
|---|---|---|---|---|---|---|---|
| 0.292164 | 2.219717 | ... | 0.075468 | 0.609535 | 0.026057 | 0.540822 | 0.306312 |
| 0.253474 | 8.873577 | ... | 0.072574 | 0.585835 | 0.101428 | 0.132031 | 0.392888 |
| 0.313060 | 3.826087 | ... | 0.068731 | 0.580060 | 0.013973 | 0.328172 | 0.425604 |

| ACSegment_Superplus | Maritalstatus_Married | Maritalstatus_Single | Logindevice_Mobile | kproto_2clusters |
|---|---|---|---|---|
| 0.044393 | 0.464196 | 0.359390 | 0.671492 | 0.992472 |
| 0.114252 | 0.572136 | 0.270184 | 0.660157 | 0.032644 |
| 0.074018 | 0.563444 | 0.276057 | 0.656344 | 0.720544 |

*Table 8 – 3 cluster profile*

**Business Insights from cluster profiling:**

> Even though this is an unsupervised algorithm without the use of target variable in clustering, the profiling came up with clear separation of groups only for those features that also showed a high F-statistic and high Chi square value in the bivariate analysis.

> Higher the tenure, lesser the churn. But the cluster profile also shows that for low tenures (cluster 1 and 0), this is not holding good. To further explore this, it would be good to bin Tenure and check the relationship with Churn through a stacked bar.

> More Regular plus customers have a greater churn compared to least churn clusters.

> High churn clusters had contacted customer care almost twice as the least churn customers.

- The cashback was lower for high churn cluster compared to low churn cluster.
- High churn cluster contacted customer care less times compared to low churn customers.
- The cluster with maximum complaints last year also has the maximum churn. The cluster with minimum complaints last year has the minimum churn. It could be suggested that if any customer files a complaint, the complaint be followed up and resolved until the customer is satisfied. This process needs to be looked at if it can be strengthened.

## 2.8.    Business insights from EDA

- **Customer feedback:** 78% of customers have rated service as 3 or less than 3 (out of a scale of 5). Likewise, 61% of customers have rated customer care agents a score of 3 or less than 3 (out of a scale of 5). This indicates that the customer feedback is pointing to dissatisfaction or bare satisfaction in service and also in customer care engagement.
- **Relationship between Tenure and Churn:** In the bivariate histogram for Tenure vs. Churn, it can be seen that the churn is very high for low tenures. For tenure between 0 and 1, around 51.85% customers have churned. The reason for this churn needs to be drilled down and addressed.
- **Relationship between Account segment and Churn:** More customers in Regular_Plus plan seem to churn. A comparison of this plan and competitors plan with same features and for same pricing range could be done to ascertain if the plan or pricing needs to be changed. Also, customer feedback for customers on this plan could be obtained and analysed to find out why there is more churn in this account segment.
- **Relationship between Monthly revenue and Churn:** The % churn in high revenue customers is slightly more than the churn in lower revenue customers. This is would a cause for concern for the DTH provider that not only is there more churn but more high revenue customers are churning.
- **Relationships between Days since customer care connect and Churn:** Days since customer care connect for churned customers are lesser than for active customers. This shows that churn has happened shortly after the customers have contacted customer care.
- **Relationship between Customer care contacted last year and Churn:** The number of times customer care was contacted previous year was more in churned customers compared to active customers.
- **Relationship between Complaints made last year and Churn:** Proportion of customers who complained is significantly more in churned customers compared to active customers.
- **Relationship between User count and Churn:** Proportion of accounts with user counts 5 and 6 is more in churned customers compared to active customers.
- **Relationship between Payment type and Churn:** Proportion of customers who have paid through E-wallet and Cash on delivery is more within churned customers compared to active customers.
- **Relationship between City tier and Churn:** Proportion of customers who reside in Tier 3 cities is more within churned customers compared to active customers. Whether there is more competition in those cities compared to Tier-1 cities needs to be explored by the DTH provider.
- **Relationship between Marital status and Churn:** Single customers have churned more compared to married or divorced customers.

# 3. Data Cleaning and Pre-processing

## 3.1. Removal of unwanted variables

The following are the checks done to see if any columns can be dropped before modelling exercise:

➢ Any variable that has unique values for each observation – for e.g., AccountID field. This would not contribute to the model as it is just an ID field to tag each observation.
➢ Any variable that remains constant for all or most of the observations as this does not add any strength to prediction. As observed from the histogram (numeric) and count plots/value counts (categorical), there are no variables that have constant value for all observations.
➢ Any variable that has nulls in more than 25 to 30% of observations. The maximum nulls present are in cashback variable and that column contains 4% nulls. Hence no column will be dropped.
➢ Any predictor variable that has a strong correlation with another predictor variable. Then one of the variables can be dropped. As seen in the correlation heat map above, there are no strong correlations, hence no variable needs to be dropped.
➢ Any predictor variable that has a very weak correlation with the target variable. As seen from the Chi square test and Anova test in the bivariate analysis section above, two variables – rev_growth_yoy and coupon_used_for_payment were found to not be significant (at a significance level of 5%). Hence these two variables would be dropped from further processing.

> **AccountID, rev_growth_yoy and coupon_used_for_payment have been removed**

## 3.2. Addition of new variables

Cluster code has been added to the dataset (details about clustering are provided in section 2.7). It may be used when experimenting with model building.

## 3.3. Missing value treatment

**Percentage of nulls**

Percentage nulls or missing values present in the predictor variables of the dataset are as follows:

```
cashback                   4.18
Day_Since_CC_connect       3.17
Complain_ly                3.17
Login_device               1.96
Marital_Status             1.88
CC_Agent_Score             1.03
Account_user_count         0.99
City_Tier                  0.99
Payment                    0.97
Gender                     0.96
rev_per_month              0.91
CC_Contacted_LY            0.91
Tenure                     0.91
Service_Score              0.87
account_segment            0.86
rev_growth_yoy             0.00
coupon_used_for_payment    0.00
dtype: float64
```

*Table 9 – Percentage null values*

*Figure 16 – Visualization of nulls*

**Missing value treatment**

Missing value treatment was done using KNN imputation, a distance-based method. The following treatments were done as they are pre-requisites for missing value treatment using KNN.

➢ All variables need to be numeric. Any object-type categorical variables need to be encoded suitably (label/ one-hot encoding). The following columns were one hot encoded as it contained nominal categorical variables. The first encoded variable was dropped to avoid multi-collinearity.

➢ After encoding, the column names were renamed in order to shorten, make it uniform and remove blank spaces that resulted out of some of the category values from one-hot encoded columns (e.g., Payment_Credit Card).

➢ All variables need to be scaled as KNN is a distance-based algorithm. Scaling was done using Standard Scalar function from SKLearn library for the predictor variables. The target variable was left as-is.

➢ Null imputing was done using Sklearn's KNNImputer function. This algorithm imputed missing values using K-nearest neighbours.

| One hot encoded variables |
|---|
| Payment','Gender','account_segment','Marital_Status','Login_device' |

```
Tenure                    0
City_Tier                 0
CC_Contacted_LY           0
Service_Score             0
User_Count                0
CC_Score                  0
Rev_Permonth              0
Complain_LY               0
Days_Since_CC             0
Cashback                  0
Payment_Creditcard        0
Payment_Debitcard         0
Payment_Ewallet           0
Payment_UPI               0
Gender_Male               0
ACSegment_Regular         0
ACSegment_Regularplus     0
ACSegment_Super           0
ACSegment_Superplus       0
Maritalstatus_Married     0
Maritalstatus_Single      0
Logindevice_Mobile        0
```

*Table 10 – Nulls in predictor variables after KNN imputing*

## 3.4. Variable transformation

➢ **Encoding:** The variables Payment, Gender, Account_Segment, Marital_Status and Login_device are all categorical object type. They need to be converted to numeric variables. The categories in these variables do not have an order. Hence, they were one-hot encoded. When performing one-hot encoding, the variable is split into multiple columns with each column taking binary values corresponding to each category in the original attribute. Additionally, one of the columns in the one-hot encoded variables is dropped in order to avoid multi-collinearity which may cause performance degradation and interpretability issues in some of the models.
➢ **Scaling:** The dataset had been scaled as it is a pre-requisite for any distance-based algorithm like KNN imputer, K-means clustering, KNN and ANN. The scaled data can also be used by all models irrespective of whether they expect scaled input or not.
➢ No other transformation is expected for modelling as of now.

## 3.5. Outlier treatment

The following continuous variables have outliers.



*Figure 17 – Box plot for continuous variables*

**Observations:**

➢ Some outliers for certain variables are closer to the whisker, whereas there are a group of outliers that are far beyond the whisker with no in between values. For instance, rev_per_month has a huge space between 30 and 100 indicating absence of values in that range. Those outliers in the extreme 25 values do not correlate with corresponding outliers in cashback field. We cannot rule out these outliers as incorrect values, they may belong to hotels with many rooms. But models like logistic regression are sensitive to outliers and may not give good performance if outliers are left untreated.

➢ Hence, two approaches to modelling will be performed – one set of data with outliers treated for outlier sensitive models and other set of data with outliers not treated (left as-is) for outlier resistant algorithms such as Random-forest and during tuning/trials of other algorithms.

➢ Coupon_used_for_payment has a very limited range 0 to 16. Hence, for the purpose of this analysis, the outliers will not be treated (similar to categorical variables).

For the outlier treated dataset, outliers beyond upper and lower whiskers were treated by capping to the lower and upper range where

lower_range= 1st quartile - (1.5 * IQR) and

upper_range= 3rd quartile + (1.5 * IQR)

Where, IQR = 3rd quartile – 1st quartile value



*Figure 18 – Box plot for continuous variables post outlier treatment*

| | Skewness |
|---|---|
| **Tenure** | 0.80 |
| **CC_Contacted_LY** | 0.80 |
| **Rev_Permonth** | 0.78 |
| **Complain_LY** | 0.95 |
| **Days_Since_CC** | 0.82 |
| **Cashback** | 0.93 |

*Table 11 – Skewness of numeric variables after outlier treatment*

# 4. Model Building and Tuning

- In this business case, the need is to predict whether a given customer would churn or not. This is a binary classification problem with only two prediction outcomes '0' – will not churn and '1' – will churn.
- Since there is a target variable 'Churn' to be predicted, this is a supervised learning problem.
- There are several algorithms that can be used for classification problems such as these.
    - **Linear classification:** Logistic Regression, Linear Discriminant Analysis.
    - **Non-linear classification algorithms:** SVMs non-linear adaptations, K-nearest neighbour, Artificial Neural Network.
    - **Ensemble models:** Random Forest, Adaboost, Gradient Boost.
- These algorithms have certain assumptions about the data on which they are fit. Depending on the nature of the data, algorithms would give good or poor performance.
- Different treatments of pre-processed data were prepared - with and without outliers, with and without SMOTE resampling, with and without scaling.
- VIF (Variance Inflation Factor) was calculated for all the predictor variables. One by one predictor variables whose VIF was more than 5 were identified. 4 variables – Cashback, Service Score, Clusters and User count had VIF greater than 5. It is to be noted that even though User count had a high VIF value, it was retained as it showed significant Chi square value as part of EDA. The remaining 3 variables were dropped before processing the datasets; hence, none of the models have used these variables as predictors.
- It is to be noted that rev_growth_yoy and coupon_used_for_payment were dropped after Anova and Chi-square test were conducted and double checked against EDA plots. Hence none of the models have used these variables as predictors.
- Scaled data was used for distance-based algorithms such as KNN and ANN. Smote resampled data was also tried for few algorithms.
- Data was split into train and test set in the ratio 70:30. 7882 records were assigned to train dataset and 3378 records were assigned to test dataset. The selection was done in such a manner that both sets had similar distribution of target variable as in the original dataset (i.e., 16.8% churn=1s).
- 8 algorithms were chosen and for each algorithm

    - Base model (with default hyper parameters) was constructed and evaluation on train and test datasets performed.
    - Different data was used and performance measured and recorded.
    - Hyper parameters for the algorithm were tuned using SKlearn library's GridSearchCV and also manually if required.
    - Performance was again measured for the tuned algorithm.
    - Feature importance was extracted from the model through in-built attributes for certain models. Amongst black box models, Sklearn's Permutation feature importance was used as a wrapper function on the model to obtain feature importance.

## 4.1.  Build various models

### 4.1.1.  Logistic Regression

Logistic regression is outlier sensitive; hence only outlier treated data has been used for all the models.

- SKlearn's Logistic Regression was used alongside RFE (Recursive feature elimination) to determine number of features that can give best performance as well as the ranking of features that Logistic regression provides.
- Stats model implementation was tried to obtain p-values to determine what features need to be retained in the model and to use the coefficients to understand the relationship of predictor and target variables.

| Model reference | Data treatment | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| LR_model1 | No | No | Default base model | 0.89 | 0.77 | 0.51 | 0.61 | 0.88 | 0.89 | 0.78 | 0.5 | 0.61 | 0.87 |
| LR_model2 | No | No | Default base model, RFE variables=8 | AUC = 0.77 | | | | | AUC = 0.77 | | | | |
| LR_model3 | No | No | Default base model, RFE variables =12 | AUC = 0.78 | | | | | AUC = 0.78 | | | | |
| LR_model4 | No | No | Default base model, RFE vars=16 | 0.89 | 0.78 | 0.49 | 0.6 | 0.88 | 0.89 | 0.8 | 0.47 | 0.59 | 0.87 |
| LR_model5 | No | No | Default base model, RFE vars=18 | 0.89 | 0.77 | 0.49 | 0.6 | 0.88 | 0.89 | 0.79 | 0.48 | 0.6 | 0.87 |
| LR_model6 | No | No | Default base model, RFE vars=19 | 0.89 | 0.77 | 0.49 | 0.6 | 0.88 | 0.89 | 0.8 | 0.48 | 0.6 | 0.87 |
| LR_model7 | No | No | Gridsearch CV, best model for f1 score | 0.89 | 0.77 | 0.5 | 0.61 | 0.88 | 0.89 | 0.78 | 0.49 | 0.6 | 0.87 |
| LR_model8 | Yes | No | Default base model | 0.81 | 0.8 | 0.82 | 0.81 | 0.89 | 0.79 | 0.44 | 0.79 | 0.56 | 0.87 |
| LR_model9 | No | Yes | Default base model | 0.89 | 0.77 | 0.5 | 0.61 | 0.88 | 0.89 | 0.78 | 0.49 | 0.6 | 0.87 |
| LR_model10 Statsmodel | No | No | Iterated 4 times to remove 4 variables | 0.89 | 0.78 | 0.51 | 0.61 | 0.88 | 0.89 | 0.79 | 0.48 | 0.6 | 0.87 |

*Table 12 - Model tuning and performances – Logistic Regression*

**SKLearn Base model with default hyper parameters (Also the best model)**



*Figure 19 - Logistic regression - Base model and best model performance*

## 4.1.2.    Linear Discriminant Analysis

Linear Discriminant Analysis (LDA) is also another machine learning classifier. It works well when there are linearly separable classes in data. It assumes that the underlying data has a Gaussian distribution but can perform well even if assumptions are violated.

| Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| Treated | No | No | Default base model | 0.89 | 0.77 | 0.47 | 0.58 | 0.88 | 0.88 | 0.77 | 0.45 | 0.57 | 0.86 |
| Treated | No | No | Gridsearch CV, best model for f1 score | 0.89 | 0.77 | 0.47 | 0.59 | 0.88 | 0.88 | 0.77 | 0.45 | 0.57 | 0.86 |

*Table 13 - Model tuning and performances – Linear Regression*

- SKlearn's Linear Discriminant Analysis function was used for modelling.
- The base model was run with default hyper parameters and the performance metrics noted.
- The model's hyper parameters were tuned using GridSearchCV function and model.
- Constructed using the best parameters selected by GridSearchCV. This did not provide much improvement over the base model except that f1-score improved by 0.01.
- Data used was outlier treated unscaled dataset.
- Model performance metrics for base model and tuned model have been provided below.

**Linear Discriminant Analysis model with default hyper parameters**



*Figure 20 - LDA - Base model and best model performance*

## 4.1.3.    Support Vector Machine

Support vector machine (SVM) is a popular machine learning algorithm that can be used for classification as well as regression. It works well when there are higher dimensions as well. It is very versatile as there are different kernel functions that can be specified to work well with the given data.

| Model reference | Algorithm | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| SVM_model1 | SVM | Treated | No | Yes | Default base model | 0.94 | 0.93 | 0.71 | 0.8 | 0.97 | 0.93 | 0.9 | 0.64 | 0.75 | 0.94 |
| SVM_model2 | SVM | Treated | No | Yes | Gridsearch CV, best model for f1 score | 0.99 | 0.93 | 1 | 0.96 | 1 | 0.97 | 0.87 | 0.93 | 0.9 | 0.98 |
| SVM_model3 | SVM | Not treated | No | Yes | Gridsearch CV, best model for f1 score | 0.98 | 0.89 | 0.99 | 0.94 | 1 | 0.95 | 0.83 | 0.92 | 0.87 | 0.97 |

*Table 14 - Model tuning and performances – SVM*

In this modelling exercise, SKlearn's Support Vector machine function was used for modelling:

- The model requires scaled data so scaling was done using Sklearn's Standard Scalar.
- The base model was run with default hyper parameters with outlier treated unscaled dataset and the performance metrics noted. Tuning was later done using GridSearchCV.
- Model performance metrics for base model and best model have been provided in the below sections.

**SVM base model with default hyper parameters**



**Train dataset Confusion Matrix**

**Test dataset Confusion Matrix**

**Train dataset Classification report**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.99 | 0.97 | 6555 |
| 1 | 0.93 | 0.71 | 0.80 | 1327 |
| accuracy |  |  | 0.94 | 7882 |
| macro avg | 0.93 | 0.85 | 0.88 | 7882 |
| weighted avg | 0.94 | 0.94 | 0.94 | 7882 |

**Test dataset Classification report**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.93 | 0.99 | 0.96 | 2809 |
| 1 | 0.90 | 0.64 | 0.75 | 569 |
| accuracy |  |  | 0.93 | 3378 |
| macro avg | 0.92 | 0.81 | 0.85 | 3378 |
| weighted avg | 0.93 | 0.93 | 0.92 | 3378 |

AUC-ROC curve for train and test datasets

Train AUC :0.97
Test AUC :0.94

*Figure 21 - SVM - Base model and best model performance*

## 4.1.4.   Artificial Neural Network(ANN)

Artificial Neural Network (ANN) is a powerful machine learning algorithm that can be used for classification as well as regression. This algorithm can learn the complex patterns in underlying data. There is a tendency to over fit, but that can be controlled in the tuning exercise using the hyper parameters.

| Model reference | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Outliers | Smote | Scaling |  | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| ANN_model1 | Treated | No | Yes | Default base model | 1 | 1 | 1 | 1 | 1 | 0.98 | 0.94 | 0.91 | 0.93 | 0.99 |
| ANN_model2 | Treated | No | Yes | Gridsearch CV, best model for f1 score | 0.99 | 0.98 | 0.97 | 0.97 | 1 | 0.97 | 0.92 | 0.89 | 0.91 | 0.99 |
| ANN_model3 alpha = 0.2 | Treated | No | Yes | Default base model with alpha = 0.2 | 0.97 | 0.95 | 0.89 | 0.92 | 0.99 | 0.95 | 0.92 | 0.8 | 0.85 | 0.97 |
| ANN_model3 alpha = 0.1 | Treated | No | Yes | Default base model with alpha = 0.1 | 0.99 | 0.98 | 0.95 | 0.97 | 1 | 0.97 | 0.94 | 0.85 | 0.9 | 0.98 |
| ANN_model3 alpha = 0.05 | Treated | No | Yes | Default base model with alpha = 0.05 | 0.99 | 0.98 | 0.98 | 0.98 | 1 | 0.97 | 0.94 | 0.9 | 0.92 | 0.99 |
| ANN_model3 alpha = 0.04 | Treated | No | Yes | Default base model with alpha = 0.04 | 1 | 0.99 | 0.98 | 0.99 | 1 | 0.97 | 0.94 | 0.9 | 0.92 | 0.99 |
| ANN_model3 alpha = 0.03 | Treated | No | Yes | Default base model with alpha = 0.03 | 1 | 0.99 | 0.98 | 0.99 | 1 | 0.97 | 0.95 | 0.86 | 0.91 | 0.99 |

*Table 15 - Model tuning and performances – ANN*

In this modelling exercise, SKlearn's Multilayer perceptron function was used for modelling:

- ■ The model requires scaled data so scaling was done using Sklearn's Standard Scalar.

- The base model was run with default hyper parameters with outlier treated scaled dataset and the performance metrics noted. Tuning was later done using GridSearchCV.
- Model performance metrics for base model and best model have been provided in the below sections.
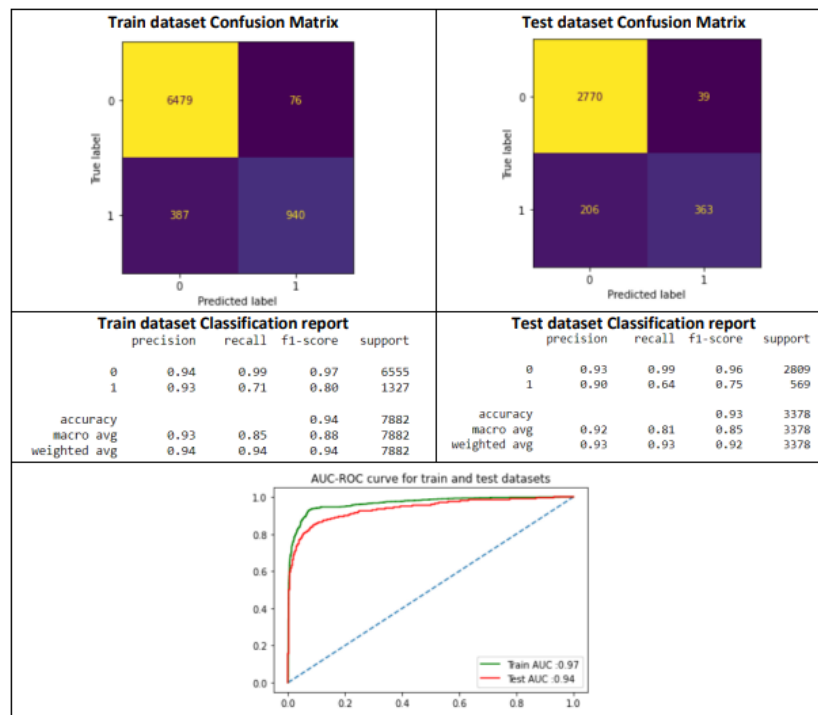
**ANN base model with default hyper parameters**



| Train dataset Confusion Matrix | Test dataset Confusion Matrix |
|---|---|

```
           Train dataset Classification report
              precision    recall  f1-score   support

           0       1.00      1.00      1.00      6555
           1       1.00      1.00      1.00      1327

    accuracy                           1.00      7882
   macro avg       1.00      1.00      1.00      7882
weighted avg       1.00      1.00      1.00      7882
```

```
           Test dataset Classification report
              precision    recall  f1-score   support

           0       0.98      0.99      0.99      2809
           1       0.94      0.91      0.93       569

    accuracy                           0.98      3378
   macro avg       0.96      0.95      0.96      3378
weighted avg       0.97      0.98      0.98      3378
```

*Figure 22 - ANN - Base model and best model performance*

## 4.1.5.    K-Nearest  Neighbours(KNN)

KNN classifier works by looking at K-Nearest Neighbours to the given data point. It decides the target value based on its neighbours. KNN works on a principle assuming every data point falling near to each other is falling in the same class. It is also a black box model and lacks interpretability. Since it is non-parametric, it may be computationally expensive and require more memory to store training data. It also has a tendency to over fit. Although this model was tried on the given data and tuned extensively, due to the above said reasons, it has been decided not to select this as best model even if model performance is good.

| Model reference | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| KNN_model1 | Treated | No | Yes | Default base model | 0.98 | 0.95 | 0.93 | 0.94 | 1 | 0.96 | 0.88 | 0.85 | 0.87 | 0.98 |
| KNN_model2 | Treated | No | Yes | GridSearchCV, best model for f1 score | 1 | 1 | 1 | 1 | 1 | 0.98 | 0.93 | 0.93 | 0.93 | 0.99 |

*Table 16 - Model tuning and performances – KNN*

As the train data seemed to over fit, a 5-fold cross validation was run on complete data to ensure that the test data performance measure f1-score is holding up. The 5-fold cross

validation gave a mean F1 score of 0.90 across 5 folds. The minimum f1 score in one of the folds was 0.85. Also, KNN as a model is computationally expensive. Hence although KNN seemed to give good performance as far as metrics is concerned, this was not considered to be selected as final model.

**KNN best model performance metrics**

Hyper parameters used: KNeighborsClassifier (algorithm = 'auto', metric= 'minkowski', p= 1, weights= 'distance')
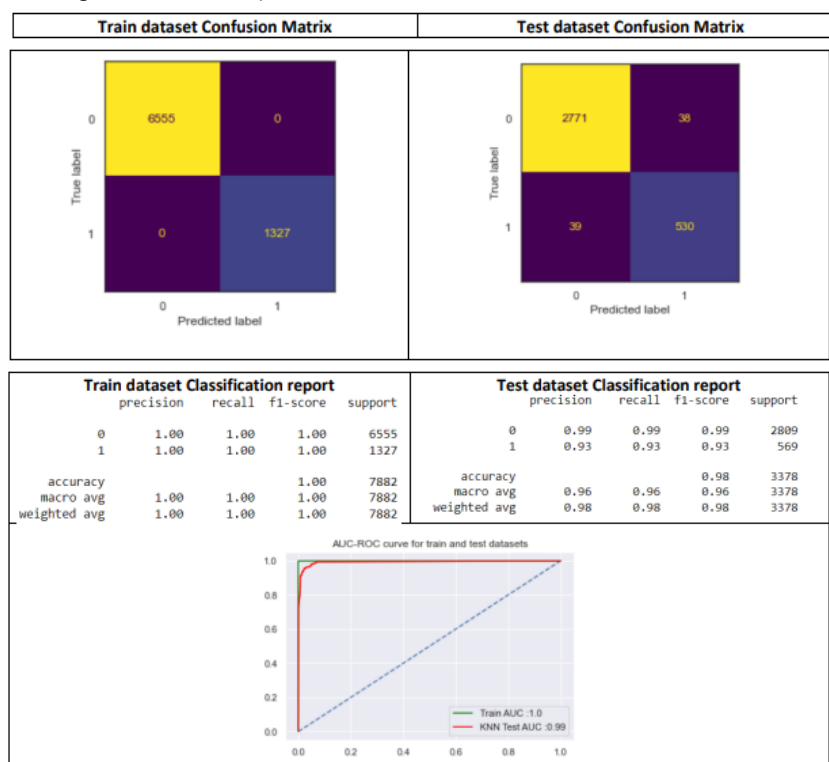


*Figure 23 - KNN - Base model and best model performance*

**Observations:** Train dataset has over fit because of a smaller number of neighbours' selection. Let us try to observe the f1-score and accuracy for different values of K (neighbours) again using the above hyper parameters.
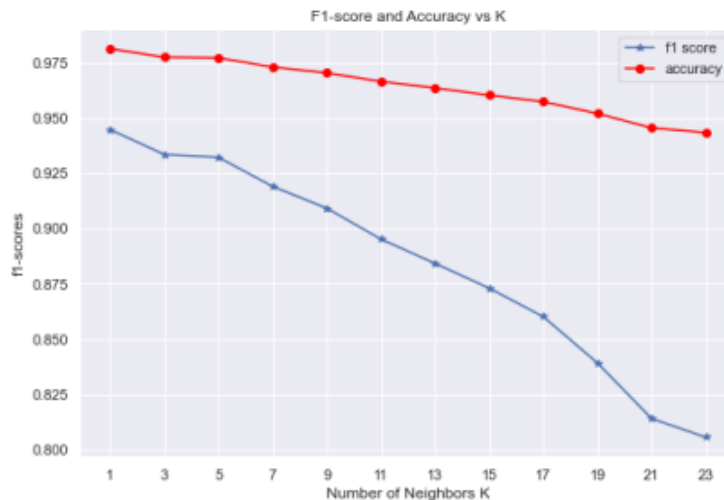
*Figure 24 - KNN – 5 fold cross validation*

**Observations:** 5 seem to be the optimum value but for that, train dataset has fully grown. KNN_model2 has the best grid with best neighbours. Even though the training data performance has fully grown, test data is not far behind and has a difference of 7%. Let's use cross validation on full data to see if the f1-score of 93% holds true. For 5-fold cross validation on full data, the following are the F1-scores.

```
[0.90358127, 0.93530997, 0.93405114, 0.91005291, 0.84656085]
```

**Observations:** The cross-validation scores on train dataset and entire dataset for CV=5 are comparable to test dataset but within the folds, the differences are quite high and there are some inconsistencies within folds. Hence K-means may not always give predictable results as in testing. Increasing neighbours to 7 may reduce variance but the f1-score may drop considerably compared to other models.

## 4.1.6. Ensemble method - Random Forest

Random forest is an ensemble machine learning algorithm that uses bootstrapping to reduce variance in the underlying decision trees. It also selects only a subset of features for each node split decision. Since it is an ensemble of trees with varying features for each node, each tree is different from another. Random forest is resistant to outliers and does not require the data to be scaled.

| Model reference | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| RF_model1 | Treated | No | No | Default base model | 1 | 1 | 1 | 1 | 1 | 0.97 | 0.97 | 0.86 | 0.91 | 0.99 |
| **RF_model2** | **No** | **No** | **No** | **Default base model with outliers** | **1** | **1** | **1** | **1** | **1** | **0.97** | **0.98** | **0.86** | **0.92** | **0.99** |
| RF_model3 | Treated | No | No | GridSearchCV, best model for f1 score | 0.98 | 0.98 | 0.91 | 0.94 | 1 | 0.95 | 0.9 | 0.78 | 0.84 | 0.98 |

*Table 17 - Model tuning and performances – RF*

- The base model was run with default hyper parameters and the performance metrics noted. Tuning was later done using GridSearchCV.

- Model performance metrics for base model and best model have been provided in the below sections.
- The data with outliers gave a slightly better performance compared to outlier treated data. But the difference between Recall in train and test is more than 10%. The model has over fit on train data and hence this was not considered for selection as final model.

**Random Forest base model with default hyper parameters**
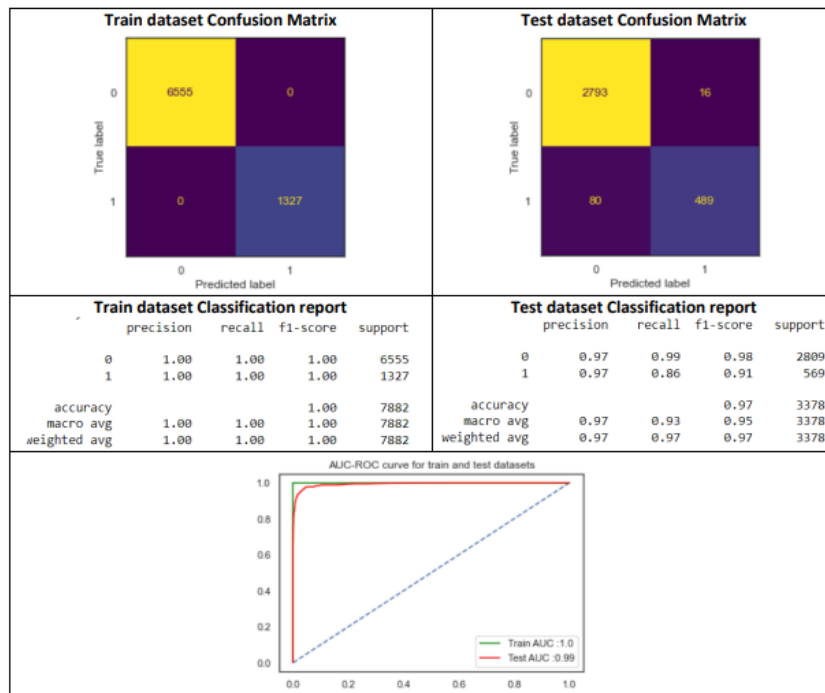


*Figure 25 - RF - Base model and best model performance*

**Observations:** Clearly, the model has over fit on the train dataset as the test dataset shows a recall of only 0.86. The model has to be tuned to reduce variance.

**Random Forest best model – default hyper parameters with outliers**



*Figure 26 - RF - Base model and best model performance with outliers*

## 4.1.7.    Ensemble method – Adaboost

An AdaBoost classifier is a meta-estimator that begins by fitting a classifier on the original dataset and then fits additional copies of the classifier on the same dataset but where the weights of incorrectly classified instances are adjusted such that subsequent classifiers focus more on difficult cases.

| Model reference | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| ADA_model1 | Treated | No | No | Default base model | 0.9 | 0.74 | 0.61 | 0.67 | 0.91 | 0.9 | 0.75 | 0.6 | 0.67 | 0.9 |
| ADA_model2 | Treated | No | No | GridSearchCV, best model for f1 score | 0.9 | 0.75 | 0.6 | 0.67 | 0.92 | 0.9 | 0.76 | 0.6 | 0.67 | 0.91 |

*Table18 - Model tuning and performances – Adaboost*

- The base model was run with default hyper parameters and the performance metrics noted. Tuning was later done using GridSearchCV.
- Model performance metrics for base model and tuned model have been provided in the below sections.

**Adaboost base model with default hyper parameters**



**Train dataset Confusion Matrix**

**Test dataset Confusion Matrix**

**Train dataset Classification report**

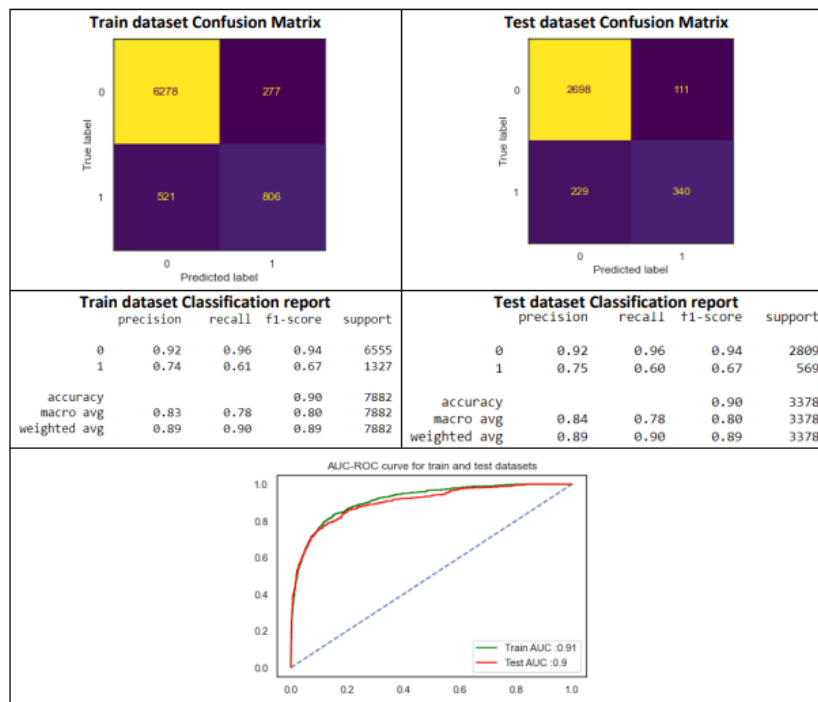|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.92 | 0.96 | 0.94 | 6555 |
| 1 | 0.74 | 0.61 | 0.67 | 1327 |
| accuracy |  |  | 0.90 | 7882 |
| macro avg | 0.83 | 0.78 | 0.80 | 7882 |
| weighted avg | 0.89 | 0.90 | 0.89 | 7882 |

**Test dataset Classification report**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.92 | 0.96 | 0.94 | 2809 |
| 1 | 0.75 | 0.60 | 0.67 | 569 |
| accuracy |  |  | 0.90 | 3378 |
| macro avg | 0.84 | 0.78 | 0.80 | 3378 |
| weighted avg | 0.89 | 0.90 | 0.89 | 3378 |

*Figure 27 - Adaboost - Base model and best model performance*

**Observations:** There is no over fit or under fit issues with the model, but the model has not performed well on predicting minority class.

## 4.1.8.    Ensemble method – Gradient Boost

Gradient boost is an ensemble machine learning algorithm that trains underlying models in a gradual, additive and sequential manner. In this modelling exercise, SKlearn's Gradient Boost classifier function was used for modelling:
The base model was run with default hyper parameters and the performance metrics noted. Tuning was later done using GridSearchCV.

| Model reference | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Outliers | Smote | Scaling |  | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| GB_model1 | Treated | No | No | Default base model | 0.91 | 0.82 | 0.61 | 0.69 | 0.93 | 0.91 | 0.82 | 0.57 | 0.67 | 0.91 |
| GB_model2 | Treated | No | No | GridSearchCV, best model for f1 score | 1 | 1 | 1 | 1 | 1 | 0.99 | 0.99 | 0.94 | 0.96 | 1 |

*Table 19 - Model tuning and performances – Gradient Boost*

The model is robust (no over fit or under fit) but the recall is quite poor in both train and test data. Hyper parameter tuning was done to see if performance can be improved on the model.

The model's hyper parameters were tuned using GridSearchCV function from Sklearn library and model constructed using the best parameters selected.

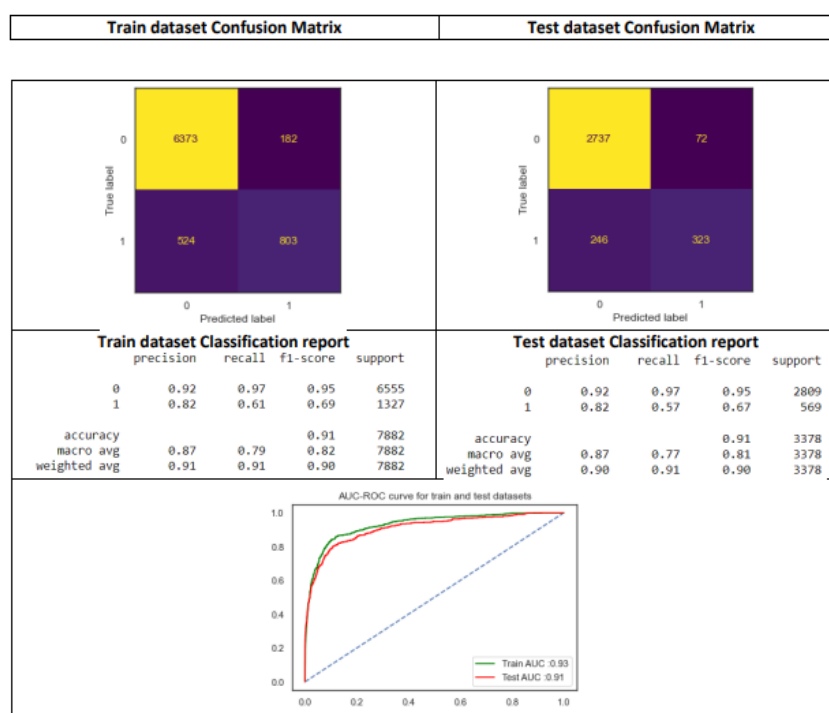**Gradient boost base model with default hyper parameters**



| Train dataset Confusion Matrix | Test dataset Confusion Matrix |

*Figure 28 - Gradient Boost - Base model and best model performance*

## 4.2.  Model Tuning

**Effort for model tuning**

Model performance improvement was accomplished using the following methods:

Around 8 different algorithms were tried across linear, non-linear and ensemble methods.

- The first model for each algorithm was the base model with the default hyper parameters. Model tuning to achieve better performance was done by tuning the hyper parameters using Grid Search CV, a function offered by Sklearn to automatically try out multiple parameter options for each hyper parameter.
- The underlying data was changed (Smote resampled/ non-resampled, outlier treated/not treated) and the effect of model performance observed for some of the models.
- Ensemble methods were used (Random Forest, Adaboost, Gradient Boost) and their hyper parameters were also tuned.
- Hyper parameters are part of the code and hence not included here. The best model for each algorithm has been highlighted in green in all tables under Section 4.1. The model reference can be used to locate the model in the python code.

**Logistic regression - model tuning**

- The base model (Section 4.1.1) was further tuned by using RFE to change the number of predictors used. It was found that reducing variables did not yield better results. Hence all 20 predictors were retained.
- GridSearchCV function from Sklearn library was used to tune the hyper parameters. It was found that tuned hyper parameters did not perform better than the base model.
- Further, Scaled data and Smote resampled data was used with the base model. It was found that Smote resulted in over fitting precision of class 1. Scaling also did not improve performance compared to non-scaled data.
- The best model is the same as the base model mentioned in Table 12 under Section 4.1.1 which is LR_model1 (reference to the python code) highlighted in green in the table.

**Stats model Logistic regression - model tuning**

- Stats model provides a model summary when a model is fit on the train data. The p-value of the predictor variables in the summary was used for deciding the significance of each predictor variable to the model.
- One by one the predictor variables whose p-value was > 0.05 were removed from the model and model rebuilt. Over multiple iterations (as shown below), the variables were eliminated one by one such that only variables with p-value < 0.05 remained in the model. 16 variables were significant at a level of 0.05.
  - o Iteration 1: ACSegment_Superplus removed.
  - o Iteration 2: Payment_Ewallet removed.
  - o Iteration 3: ACSegment_Regularplus removed.
  - o Iteration 4: Maritalstatus_Married removed.

- The F1-score of the final model built was close but not greater than the SkLearn's best model (described in Table 12 under Section 4.1.1). Hence this was not selected as the best model in logistic regression.

**Linear Discriminant Analysis - model tuning**

- GridSearchCV function from Sklearn library was used to tune the hyper parameters on the base model described in Section 4.1.2. It was found that tuned hyper parameters performed almost similar to the base model.
- The best model is highlighted in green in the table and the metrics for this model is given in Table 13 under Section 4.1.2 (hyper parameter tuned model)

**Support vector machines - model tuning**

- The model's hyper parameters were tuned using GridSearchCV function and model constructed using the best parameters selected by GridSearchCV. This provided significant improvement over the base model's performance simply by changing the kernel to 'poly'. A penalty parameter ('C') had to be included as the model was over fitting on training data.
- The best grid model was also run with data that was not outlier treated. Although this was better than base model with default parameters, the previous model using outlier treated data had the best performance in this algorithm.
- The best model is highlighted in green in the table and the metrics for this model is given in Table 14 under Section 4.1.3 (hyper parameter tuned model).

**Artificial Neural Network - model tuning**

- The model's hyper parameters were tuned using GridSearchCV function and model constructed using the best parameters selected by GridSearchCV.
- As can be seen from the base model, the model had over fit on train dataset (all 1s). During tuning, higher alpha was provided in order to add a penalty that will regularize the model by reducing weights.
- Lower hidden layer sizes were also provided to GridsearchCV in order to reduce over fitting.
- The best model is highlighted in green in the table and the metrics for this model is given in Table 15 under Section 4.1.4 (hyper parameter tuned model).
- The base model has a 0.01 better recall than ANN_model3, but ANN_model3 has been chosen as the best model as the train and test performances are more comparable whereas in base model, the difference is more. Also, the train dataset in base model has over fit and completely learnt even the noise in the train dataset. Hence ANN_model3 which introduces a constraint through penalty is chosen to the best model.

**Random Forest - model tuning**

- The model's hyper parameters were tuned using GridSearchCV function and model constructed using the best parameters selected by GridSearchCV.
- As can be seen from the base model, the model had over fit on train dataset (all 1s). During tuning, the tree depth was contained, the minimum samples in each leaf increased to reduce over fit. The metrics for this model is given in Table 17 under Section 4.1.6.
- One other model using data that was not outlier treated (i.e., outliers left as-is) was built using default hyper parameters. That model showed a 0.01 improvement over outlier treated base model for precision and f1- score, but recall was still an over fit.

**Adaboost - model tuning**

- The model's hyper parameters were tuned using GridSearchCV function and model constructed using the best parameters selected by GridSearchCV.
- The best model is highlighted in green in the Table 18 under Section 4.1.7. Tuning has not resulted in noticeable improvement in the model.

**Gradient Boost - model tuning**

- The model's hyper parameters were tuned using GridSearchCV function and model constructed using the best parameters selected by GridSearchCV.
- Tuning improved the model performance considerably and the best model is highlighted in green in the Table 19 under Section 4.1.8.

# 5. Model Validation

The given data was split into train and test dataset in the ratio 70:30. The test dataset was held out and kept for validation purpose only.

All models were trained only on the train dataset. The trained model was used to predict train dataset target variable. The performance metrics such as Accuracy, F1-score, and Precision, Recall, confusion matrix, ROC curve and AUC was observed and recorded on train dataset.

## 5.1. Best performing model

**Criteria for the best performing model**

**Primary criteria:** **Precision, Recall & F1-score for 1s:** This is the case of class imbalance as the dataset has 16.8% churns. In this case study 'Precision' and 'Recall' of class 1 or the minority class is most important. Combining these 2 metrics, F1-score for class 1 is also used in the comparison.

- **Precision** is defined as True positive/ (True Positive + False Positive). It answers the question - Out of all customers that the model identifies as churning customers, how many are actual churners? This is the most important factor in this case as budget for offers to retain customers would be limited and more false positives would mean spending that budget on customers who would not churn anyway. The problem statement states that the revenue assurance team is very stringent about providing freebies where it is not required. Translated into metric, this would mean that the precision for 1's/churns should be highest.
- **Recall** is defined as True Positive / (True Positive + False Negative). It answers the question - Out of all actual churning customers, how many does the model correctly identify as churners? This is very important as the purpose of the project is to identify as many churners as possible in order to give special offers in order to retain them. For the DTH provider, customer acquisition cost is very high and hence retention is of utmost importance. This is the whole reason why the project exists and hence recall for 1s is also important in this case.
- **F1-score** is harmonic mean of Precision and Recall. Where both Precision and Recall are important, this metric can be used as a single metric that needs to be optimized for.

**Secondary criteria**

- **Accuracy:** This is a classification problem and the dataset has class imbalance. That is, the proportion of churn and non-churn customers is not equal. With imbalanced classes, it's easy to get a high accuracy without actually making useful predictions. So, accuracy as an evaluation metric makes sense only if the class labels are uniformly distributed. We are concerned with correct prediction of churn customers (class 1). Hence 'Accuracy' is not a correct metric to compare various models but for the sake of completeness and to ensure that 0s (majority class) are not overlooked, it is still recorded in the comparison matrix.
- **AUROC:** In addition to the above metrics, the Area under curve of ROC curve is also used to evaluate model performance. An ROC curve (or receiver operating characteristic curve) is a plot that summarizes the performance of a binary classification model on the positive class. It is a curve that is constructed by evaluating true positives and false positives for different threshold values. As visualizing ROC curve is difficult for actual comparison, the Area under Curve (AUC) metric helps with a numeric comparison. The closer the AUC is to 1, the better the model. However, like accuracy this also works well for balanced dataset. For the sake of completeness, this is also recorded in comparison matrix.

The following table shows performance of the best model from each algorithm built so far. The metrics given in the below table are all for minority class/churners (1s) which is the class of interest. The best performer has been highlighted in green. The figure below that shows the ROC curve for all models and the best performer is the blue line.

| Model reference | Algorithm used | Data Used | | | Hyper parameters | Train data | | | | | Test data | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Outliers | Smote | Scaling | | Accuracy | Precision | Recall | F1 | AUC | Accuracy | Precision | Recall | F1 | AUC |
| LR_model1 | Logistic Regression | Treated | No | No | Default base model | 0.89 | 0.77 | 0.51 | 0.61 | 0.88 | 0.89 | 0.78 | 0.5 | 0.61 | 0.87 |
| LDA_model1 | LDA | Treated | No | No | Gridsearch CV, best model for f1 score | 0.89 | 0.77 | 0.47 | 0.59 | 0.88 | 0.88 | 0.77 | 0.45 | 0.57 | 0.86 |
| SVM_model2 | SVM | Treated | No | Yes | Gridsearch CV, best model for f1 score | 0.99 | 0.93 | 1 | 0.96 | 1 | 0.97 | 0.87 | 0.93 | 0.9 | 0.98 |
| ANN_model3 alpha = 0.05 | ANN | Treated | No | Yes | Default base model with alpha = 0.05 | 0.99 | 0.98 | 0.98 | 0.98 | 1 | 0.97 | 0.94 | 0.9 | 0.92 | 0.99 |
| RF_model2 | RandomForest | No | No | No | Default base model with outliers | 1 | 1 | 1 | 1 | 1 | 0.97 | 0.98 | 0.86 | 0.92 | 0.99 |
| ADA_model2 | Adaboost | Treated | No | No | GridSearchCV, best model for f1 score | 0.9 | 0.75 | 0.6 | 0.67 | 0.92 | 0.9 | 0.76 | 0.6 | 0.67 | 0.91 |
| GB_model2 | Gradient Boost | Treated | No | No | GridSearchCV, best model for f1 score | 1 | 1 | 1 | 1 | 1 | 0.99 | 0.99 | 0.94 | 1 | 1 |

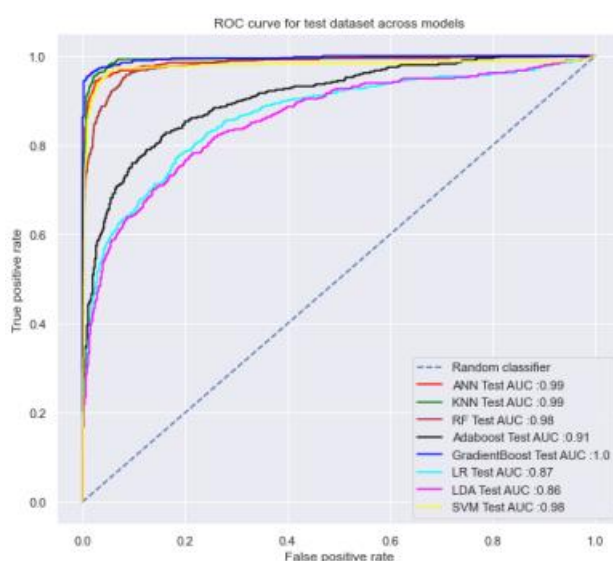*Table 20 - A comparison of best model from all algorithms*



*Figure 29 - AUC-ROC curve – All models*

## 5.2. Why Gradient boost is the best model

- The hyper parameter tuned Gradient Boost model (GB_model2) has given the best performance in terms of test data precision, recall and f1-score which are the primary evaluation metrics. The Accuracy and AUC are also highest amongst all the models. As the model looked like it over fit on train dataset, a further 5-fold and 10-fold cross-validation was done on complete dataset in which the F1-score on all folds was either comparable or greater than test data f1-score (the test data performance held true or was better for all folds).
- The difference between train data metrics and test data metrics is within 10%.
- The model is interpretable. Sklearn provided feature importance for the model.

## 5.3. How can business use these metrics?

- A precision of 0.99 implies that out of 100 customers that the model has identified as churned, 99 would actually churn and 1 would not. Any marketing budget allocated

for a targeted campaign for retention of these customers would be most optimally utilized as only 1/100 customers would be incorrectly identified as churn.

▪ A recall of 0.94 implies that for 100 customers who actually churn, the model would have identified 94 as churn and 6 as not-churn. This would mean that the campaign would target these 94 customers and there is scope for retaining these customers and missing out on 6 customers.

▪ Based on the customer base for which prediction needs to be done, business can use the above to project the numbers of customers who would churn and how many the model would identify and miss.

▪ That could be further used to come up with per customer budget (if total budget for retention campaign is known) so that appropriate offers can be designed for each customer. If per customer budget is known, cost projections for retention campaign can be calculated. If limited budget is available and not all customers can be covered, the model can also provide the probability of churning so that high probability customers can be targeted first. Also, a different perspective to this problem would be to do a segmentation and target high value customer segment.

# 6. Final Interpretation and Business Recommendations

## 6.1. Model Interpretation from best models

**Gradient Boost model interpretation**

Gradient boost model has given very good performance after tuning. The model is robust (no over fit or under fit). The figures below show the feature importance given by the Gradient Boost model and two other top scoring models (F1-score) – Random Forest and ANN. A comparison of top 6 features that contributed to these models is provided in a table.
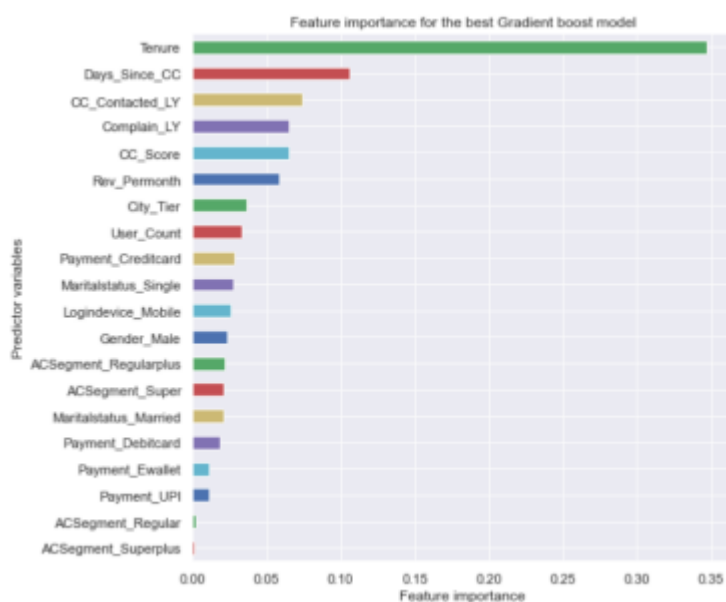


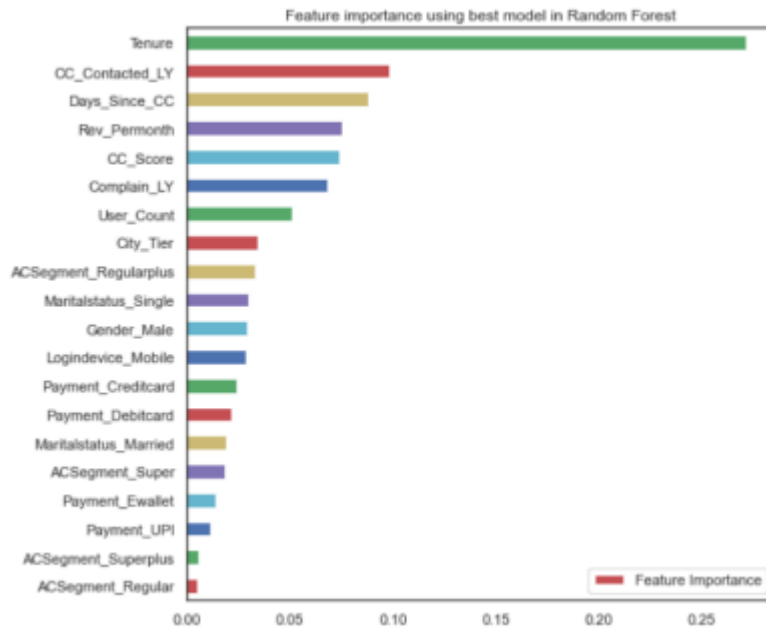*Figure 30 - Feature importance – Gradient Boost*

*Figure 31 - Feature importance – RF*



*Figure 32 - Feature importance – ANN*

* ANN is a black box model and feature importance is not directly available. Hence Sklearn's permutation feature importance function was used with this model. The permutation feature importance is defined to be the decrease in a model score when a single feature value is randomly shuffled. This procedure breaks the relationship between the feature and the target, thus the drop in the model score is indicative of how much the model depends on the feature.

| | Gradient Boost | Random Forest | ANN |
|---|---|---|---|
| Feature1 | Tenure | Tenure | Tenure |
| Feature2 | Days since customer care last contacted | Number of times customer care was contacted last year | Was complaint made last year |
| Feature3 | Number of times customer care was contacted last year | Days since customer care last contacted | Customer care agent score |
| Feature4 | Was complaint made last year | Revenue per month | City tier |
| Feature5 | Customer care agent score | Customer care agent score | Marital status single |
| Feature6 | Revenue per month | Was complaint made last year | Payment credit card |

*Table 21 - Comparison of feature importance of top 3 models*

**Observations:** The top 5 features that have influenced Gradient boost model are Tenure, Days since last customer connect, number of times customer contacted last year, complaint last year and customer care score. Together, they add up to almost 66% of the total feature importance.

Top 5 features from Gradient Boost

- From EDA, we can observe that for lower tenures especially within the first month, the churn is higher. Hence, once a customer has been acquired, the first two months is very important to keep the customer satisfied.
- Churned customers had contacted customer care more recently before churning than active customers. Per EDA, median days since last customer connect is higher for active customers. Churned customers had contacted customer care recently before churning.
- The next important parameter to predict customer churn per this model is number of times customer care was contacted by the customer. Per EDA, the median and third quartile of number of times customer care was contacted previous year is higher for churned customers compared to active/current customers.

According to ranking of important features, except tenure, the other features in top 5 are related to customer care or complaints. This may be indicative of need to monitor and improve the customer care processes.

## 6.2. Business Recommendations
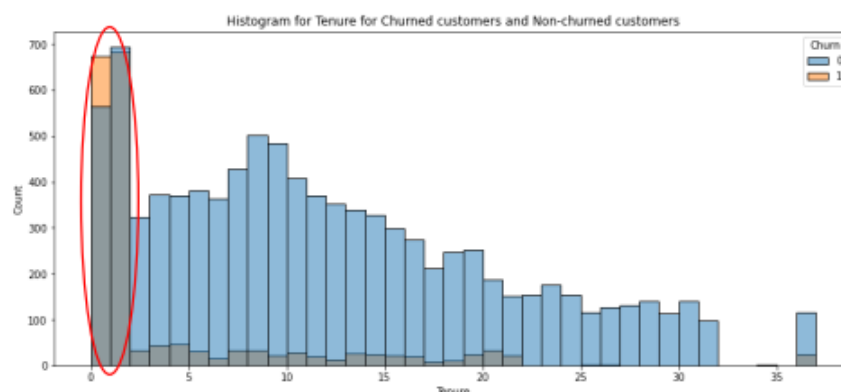
**High churn rate in low tenure customers**



*Figure 33 - Tenure and Churn*

***Insight: Churn is highest between the tenure periods 0 to 2***

**Possible Reasons:** Bad first experience or Trial periods/prepaid accounts that expire automatically if no top-up is done within a predefined period. Important to determine between the above two reasons. Based on high customer care calls, complaints registered and low cashback and coupons for low tenure customers, it points to the first reason.

**Business recommendations:**

- Activation/On boarding team could extend support beyond the initial setup until customers settle down with the service.
- Activation team proactively engages customers for the first month or two.
- Customer care takes a feedback survey about the process so that any hiccups can be understood and sorted out.
- To increase response rates for feedback, gift cards/coupons can be given.

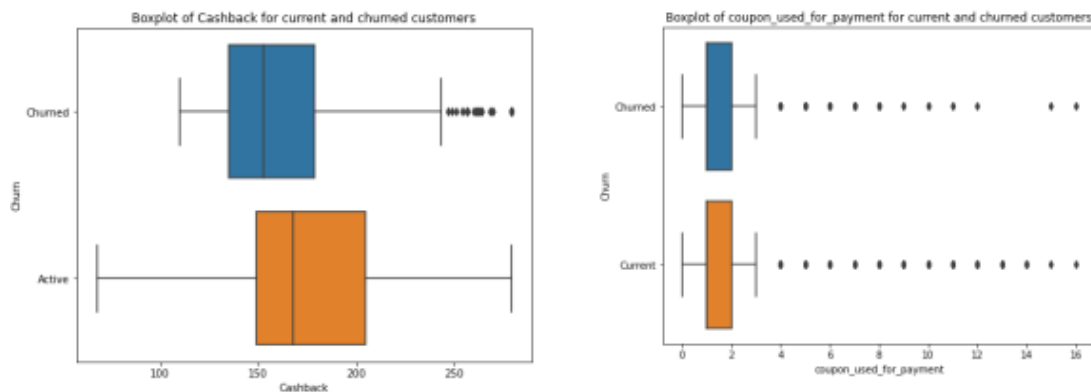**Existing retention programs – Cashback and Coupons**



*Figure 34 - Churn and existing retention programs*

- The churned customers as shown in first boxplot have lesser cashback Figure 8-1 Tenure and Churn Figure 8-2 Churn and existing retention programs 39.
- The churned and active customers have almost used the same number of coupons for payment.

***Insight: The current retention programs do not seem to be focusing on the customers with higher risk of churn.***

**Recommendation:** Review whether existing cashback and coupon programs are still relevant given the current churn model. If they are not relevant, design new retention programs to address current high risk customer group.

**Churn and Customer care service**



*Figure 35 - Churn and Customer care service*

Churned customers seem to have contacted customer care more recently before churning

- The number of times churned customers contacted customer care in the year is higher than number of times active customers contacted customer care.
- 31% of customers who registered complaint churned vs. 11% of customers who have not registered complaint in last year.

_**Insight:**_ _**These indicate behavioural changes in customer before churn happens**_

**Recommendation:** Analyse Complaints & Customer care contact reasons

- Perform Root cause analysis, identify and fix top reasons.
- Establish Service level agreements (if not already present).

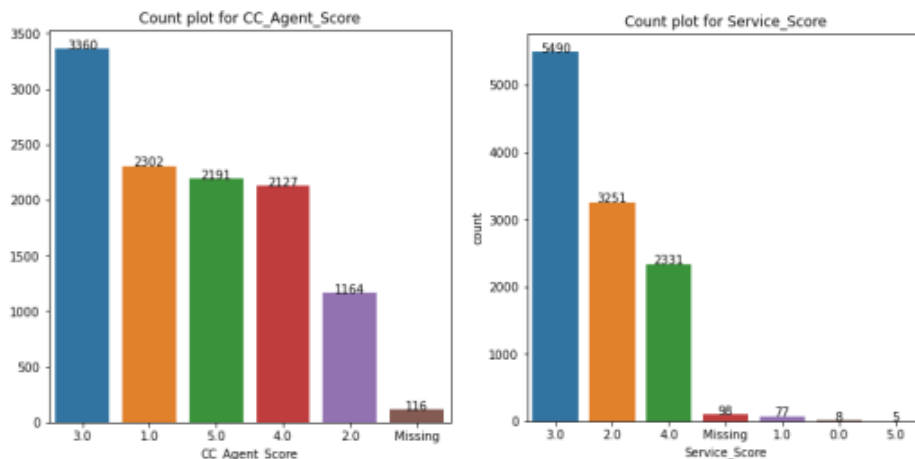**Customer care & Service – Customer perspective**



*Figure 36 – Count plot customer care & service*

***Insights: 78% of customers have rated service as 3 or less than 3 and 61% of customers have rated customer care agents a score of 3 or less than 3.***

**Recommendation:** Analyse customer feedback. Perform Sentiment analysis of the feedback (if any) that went along with scores. Identify top reasons that have resulted in low scores; if subjective feedback not captured, capture that as well.

**Revenue per month and Churn**



*Figure 37 - High churn in high revenue customers*

The % churn of customers in higher revenue group, for revenue >=7 per month is higher than low revenue group customers.

***Insight: More proportion of high revenue customers are leaving compared to less revenue customers.***

**Recommendation:** Create segmented offers for high revenue customers. An illustrative segmentation based on RFM (Recency – Frequency – Monetary) has been shown below along with sample targeted recommendations for each segment. This illustration is based on the test data. A similar segmentation can be done in discussion with client based on what metrics they feel are important to segment based on.

**Segmentation basis**

Recency (how recently they onboarded): High Recency - Tenure <=2
Frequency (how frequently they have contacted customer care): High Frequency - CC_contacted last year >=17
Monetary (Revenue): High Monetary – Revenue per month >=7
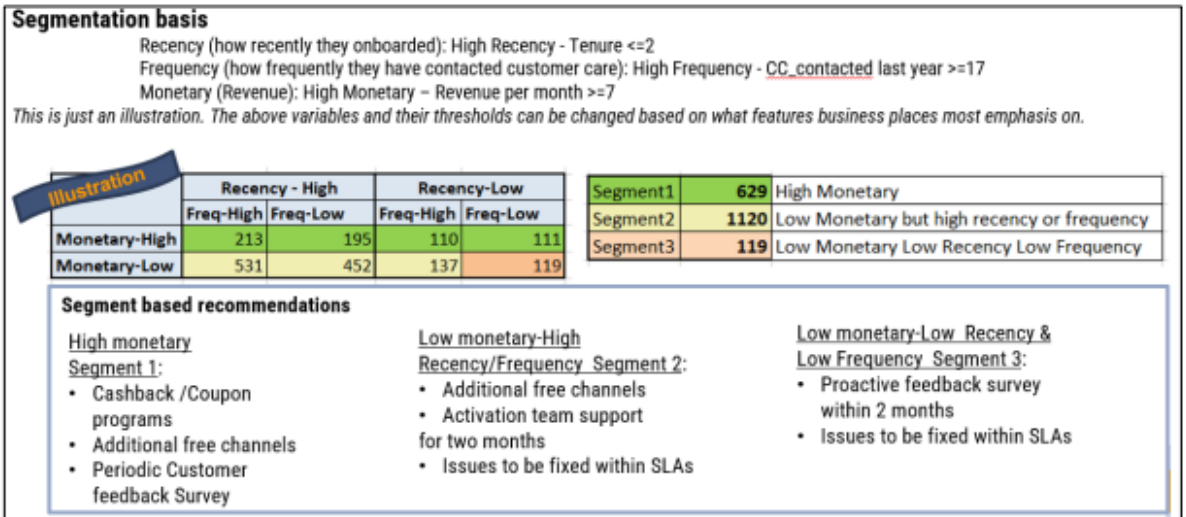*This is just an illustration. The above variables and their thresholds can be changed based on what features business places most emphasis on.*

| Illustration | Recency - High | | Recency-Low | |
|---|---|---|---|---|
| | Freq-High | Freq-Low | Freq-High | Freq-Low |
| Monetary-High | 213 | 195 | 110 | 111 |
| Monetary-Low | 531 | 452 | 137 | 119 |

| | | |
|---|---|---|
| Segment1 | 629 | High Monetary |
| Segment2 | 1120 | Low Monetary but high recency or frequency |
| Segment3 | 119 | Low Monetary Low Recency Low Frequency |

**Segment based recommendations**

High monetary
Segment 1:
• Cashback /Coupon programs
• Additional free channels
• Periodic Customer feedback Survey

Low monetary-High Recency/Frequency  Segment 2:
• Additional free channels
• Activation team support for two months
• Issues to be fixed within SLAs

Low monetary-Low  Recency & Low Frequency  Segment 3:
• Proactive feedback survey within 2 months
• Issues to be fixed within SLAs

*Figure 38 - Segmented offers illustration*

# THE END OF THE REPORT