# CUSTOMER CHURN ANALYSIS: TELCO



MUWANI ROBSON

May ,29 2020

## Executive Summary

Telecom companies typically spend most of their effort and resources on customer acquisition, even though the cost of retaining an existing customer is five times lower than acquiring a new one. Customer retention is a measure of how many of your customers continue to buy from you over time and are therefore loyal to your brand. Churn, sometimes known as customer attrition, is at the opposite end of the spectrum, i.e. the number of customers that stop buying from your company. The aim of this research project is to provide business analytics, visualize and discover the root causes of churn in Telecoms based on the Telco data set and implore strategies to minimise deactivations and increase customer lifetime value.

## Introduction

Telecom industry retention surveys show that while price and product/ service are important, most customers leave a service provider because of dissatisfaction with the way they are treated. It costs hundreds of dollars to acquire a new telecom customer. When a customer churns, you not only lose the future revenue from this customer, but also the resources you spent on acquiring the customer in the first place. Researches by Bain company estimates that for a telecom provider with 5 million customers and an average churn of 2 to 2.5%, a reduction in churn by even 50 basis points would be worth $410 million in customer lifetime value over 30 months. When customers leave after poor experiences, they not only will not return, but they often amplify their message of dissatisfaction to others using social media. Key to this churn analysis is that high customer retention means long term customer value for the business, hence the need to manage customer attrition rate in saturated businesses like telecoms.

## Objectives

- How to increase customer lifetime value (CLV)

- Promote actions that drive customer satisfaction, spend and loyalty

- Analyse customer distribution and interactions across service platforms

- Improve customer service delivery

## The Study

The study is based on a subset of a large data set from Telco company. The dataset has 71,047 and 58 variables. Only 11 relevant variables were filtered for project analytics. Missing data classification on churn status being our target variable was excluded from the study.

## The Variables:

The study has 57 selected predictor variables and one target variable. For the purpose of this project we selected the most important predictors and reduced them to 12. The composite variables for the data set are as follows:

CustomerID, Churn, MonthlyRevenue, MonthlyMinutes, TotalRecurringCharge, DirectorAssistedCalls, OverageMinutes, RoamingCalls, PercChangeMinutes, PercChangeRevenues, DroppedCalls, BlockedCalls, UnansweredCalls, CustomerCareCalls, ThreewayCalls, ReceivedCalls, OutboundCalls, InboundCalls, PeakCallsInOut ,OffPeakCallsInOut, DroppedBlockedCalls, CallForwardingCalls, CallWaitingCalls, MonthsInService ,UniqueSubs ActiveSubs ,ServiceArea, Handsets, HandsetModels ,CurrentEquipmentDays, AgeHH1, AgeHH2, ChildrenInHH ,HandsetRefurbished, HandsetWebCapable, TruckOwner, RVOwner ,Homeownership, BuysViaMailOrder, RespondsToMailOffers, OptOutMailings, NonUSTravel ,OwnsComputer, HasCreditCard, RetentionCalls, RetentionOffersAccepted, NewCellphoneUser, NotNewCellphoneUser, ReferralsMadeBySubscriber, IncomeGroup, OwnsMotorcycle ,AdjustmentsToCreditRating, HandsetPrice, MadeCallToRetentionTeam, CreditRating, PrizmCode ,Occupation & MaritalStatus

**The Study Framework**

The study has a total of 10 selected predictor variables, 5 of them grouped as modifiable churn drivers (numeric type), 5 as non modifiable churn drivers(character type).
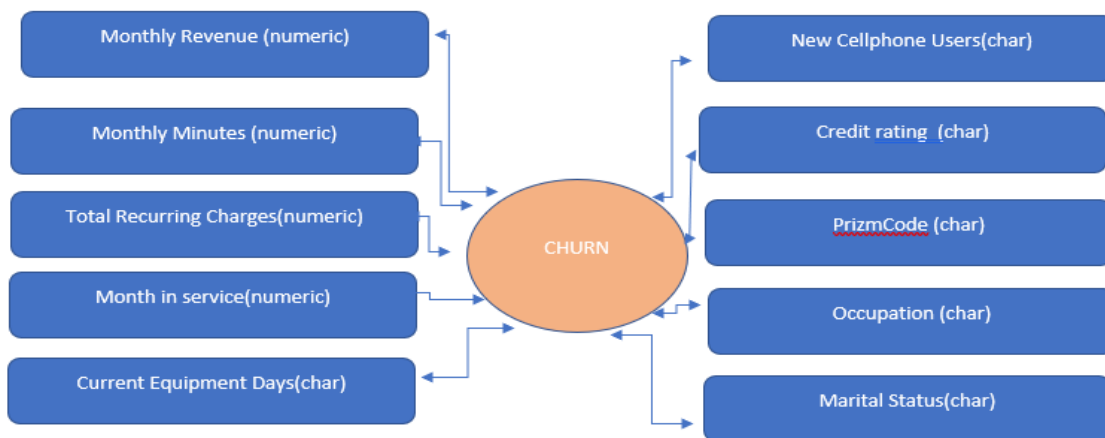
See the flow chart below:



Figure 1, study framework

**Hypotheses**

Using the Telco dataset, the researcher aims to answer the following questions:

- o What are the churn drivers in Telco?

- o What are the telco churn behaviors?

- o Does churn results in significant revenue loss?

- o What are the demographic pattens in churn behavior?

- o What is the disconnection trend/behavior?

- o What is the relationship between credit rating in churn behavior?

- o What is the relationship between equipment days in churn behavior?

- o What is the relationship between occupation in churn behavior?

- o What is the relationship between location in churn behavior?

## Methodology

SAS 9.4 was used in reading, management, analysis and modelling of the data. A total of 15 SAS PROC statements were used in the analysis of the data (see Appendix).

SAS proc import statement was used to import the dataset for maximum control of the final output. Data types were established at import with nine categorical variables and four numerical variables. Custom formats were created for improved analysis and readability of the data. Missing data values on character variables was excluded in the analysis. Missing values on continuous variables was filled with mean.

| Variable | Treatment of missing values |
| --- | --- |
| Monthly Revenue | Filled with mean |
| Monthly Minutes | Filled with mean |
| Total recurring charges | Filled with mean |
| Month in service | No missing values |
| Current Equipment Days | No missing values |
| New cellphone users | No missing values |
| Churn (target variable) | Excluded |
| Credit rating | Excluded |
| Prizm Code | Excluded |
| Occupation | No missing values |
| Churn(target variable) | Excluded |

Table 1: Treatment of missing data values

## Descriptive Analysis

The data reported a total of 36336 active accounts and 14711 deactivated accounts, representing 71.18% and 28.82% respectively. Notably, the deactivation rate is extremely high according to Bain telecommunications survey.
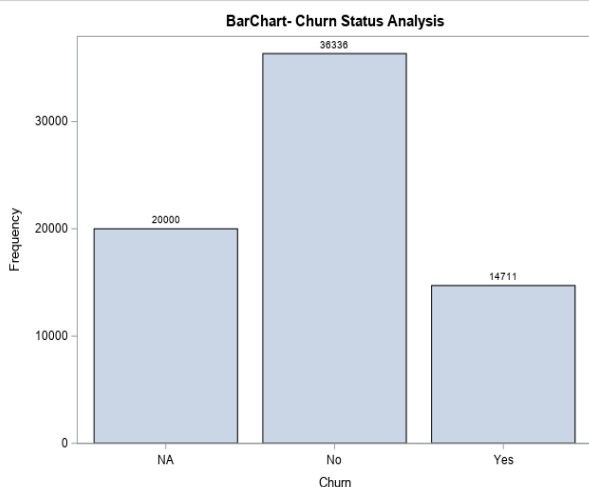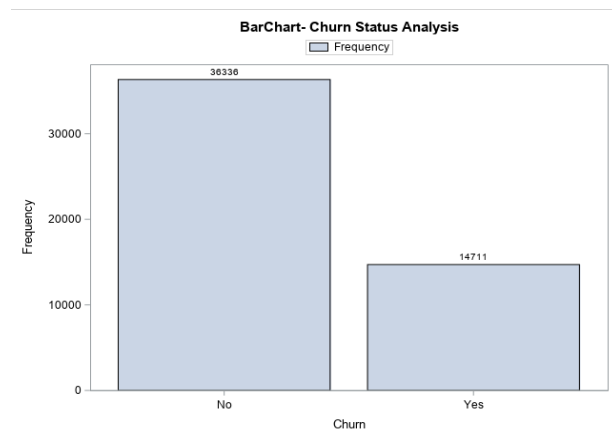


Figure 2 (including unclassified accounts)



Figure 3 (excluding unclassified accounts)

## Univariate Analysis

Univariate analysis can give a better understanding of the distribution of each variable in the framework. The main purpose of univariate analysis was to describe the data and find patterns that exist within it. Additionally, using univariate analysis I was able to observe each of the variables I terms of near normal distribution. This will allow us to use parametric statistics later that require normality as a prerequisite. No transformation was done on the variables as they depicted a near normal distribution, however, monthly revenue had extreme values thus skewed to the right and total recurring charges were more conical in shape.
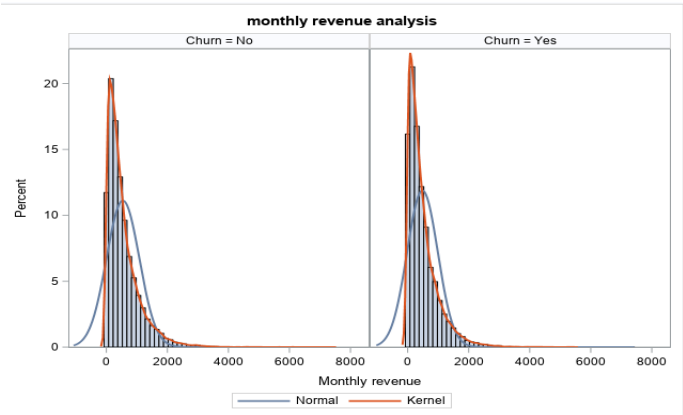


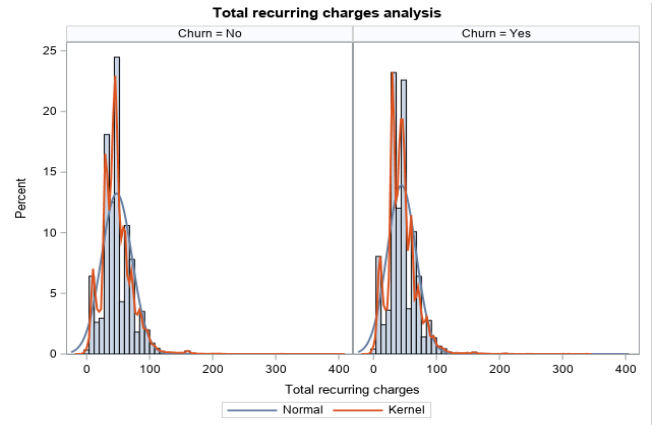Fig 4: distribution of Monthly charges



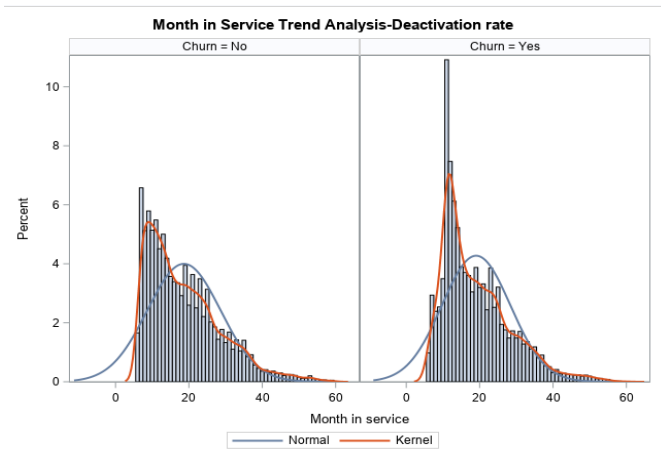Fig 5: distribution of monthly recurring charges
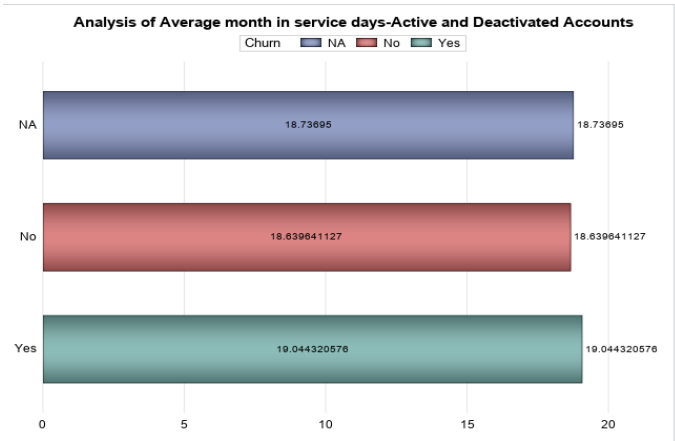


Fig 6: Months in Service Trend



Fig 7: Average days of Active Accounts (active and deactivated)

Most customers deactivations happening around 19mnths and the maximum stay is 60months (5 years)

## Bivariate analysis

Bivariate analysis is one of the statistical analysis where two variables are observed. One variable here is dependent while the other is independent. These variables are usually denoted by X and Y. So, here we analyse the changes occurred between the two variables and to what extent. It is important to compare two variables to each other in the analysis. Boxplots can show the distribution of a variable based on a group that it is a part of. In this case, a comparison of each distribution is seen when CAD is present and when it is not.

## Revenue Analysis



**BarGraph- Recurring Charges by bandwidth Analysis**

**Chi-Square Tests for Inc_bandwidth**

**The FREQ Procedure**

| Frequency Expected Cell Chi-Square Percent | Table of Churn by Inc_bandwidth | | | | | |
|---|---|---|---|---|---|---|
| | | Inc_bandwidth | | | | |
| Churn | $0-$24 | $25-$49 | $75-$99 | $50-$74 | premium | Total |
| Yes | 1840 | 7397 | 1067 | 3995 | 283 | 14582 |
| | 1665.4 | 7109.5 | 1270.7 | 4187.2 | 349.29 | |
| | 18.307 | 11.63 | 32.647 | 8.8214 | 12.581 | |
| | 3.63 | 14.58 | 2.10 | 7.88 | 0.56 | 28.75 |
| No | 3953 | 17333 | 3353 | 10570 | 932 | 36141 |
| | 4127.6 | 17621 | 3149.3 | 10378 | 865.71 | |
| | 7.3866 | 4.6924 | 13.172 | 3.5592 | 5.0763 | |
| | 7.79 | 34.17 | 6.61 | 20.84 | 1.84 | 71.25 |
| Total | 5793 | 24730 | 4420 | 14565 | 1215 | 50723 |
| | 11.42 | 48.76 | 8.71 | 28.71 | 2.40 | 100.00 |
| Frequency Missing = 324 | | | | | | |

Fig 8: Recurring charges Analysis by bands                    Table 2: Chi-Square Tests for Recurring charged bands)

## Months in service Analysis



Fig 9 : Month in service  Analysis by segments

**Table of Churn by MonthsInService**

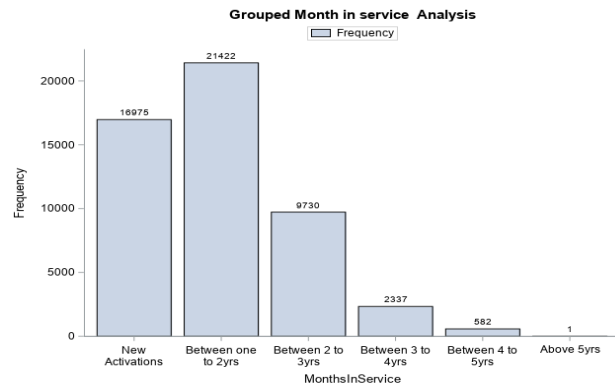| Churn | Above 5yrs | Between 4 to 5yrs | Between 3 to 4yrs | Between 2 to 3yrs | Between one to 2yrs | New Activations | Total |
|---|---|---|---|---|---|---|---|
| Yes | 1 | 152 | 645 | 2802 | 6591 | 4520 | 14711 |
|  | 0.2882 | 167.72 | 673.49 | 2804 | 6173.5 | 4891.9 |  |
|  | 1.7582 | 1.4741 | 1.2051 | 0.0015 | 28.234 | 28.28 |  |
|  | 0.00 | 0.30 | 1.26 | 5.49 | 12.91 | 8.85 | 28.82 |
| No | 0 | 430 | 1692 | 6928 | 14831 | 12455 | 36336 |
|  | 0.7118 | 414.28 | 1663.5 | 6926 | 15248 | 12083 |  |
|  | 0.7118 | 0.5968 | 0.4879 | 0.0006 | 11.431 | 11.449 |  |
|  | 0.00 | 0.84 | 3.31 | 13.57 | 29.05 | 24.40 | 71.18 |
| Total | 1 | 582 | 2337 | 9730 | 21422 | 16975 | 51047 |
|  | 0.00 | 1.14 | 4.58 | 19.06 | 41.97 | 33.25 | 100.00 |

Table 3: Chi-Square Tests for Month in service  Analysis

Pick activations and deactivations between 1 and 2 years, about 41% of the total. customers in the data base.

Most Accounts deactivation between 12 and 24months. (12.91%). New activations terminating service too, about 8.85%.

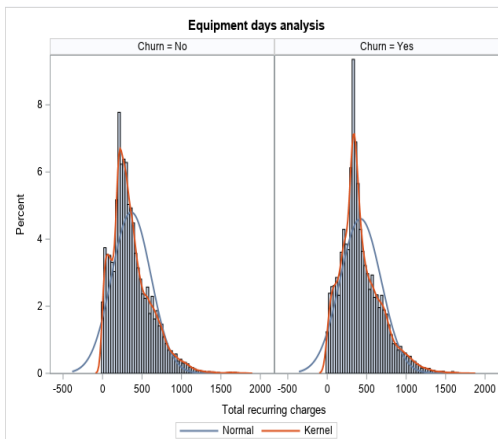## Equipment Days Analysis

| Univariate analysis | Bivarate Analysis | Chi-Square Analysis |
|---|---|---|



Fig 10 : equipment days  Analysis Analysis

**Table of Churn by CurrentEquipmentDays**

| Churn | from 0 days to 1yr | Between 4 to 5yrs | Between one to 2yrs | Between 2 to 3yrs | Between 3 to 4yrs | Total |
|---|---|---|---|---|---|---|
| Yes | 7329 | 28 | 5527 | 1563 | 264 | 14711 |
|  | 8334.2 | 27.955 | 4895.2 | 1277.6 | 176.08 |  |
|  | 121.24 | 0.0001 | 81.54 | 63.779 | 43.894 |  |
|  | 14.36 | 0.05 | 10.83 | 3.06 | 0.52 | 28.82 |
| No | 21590 | 69 | 11459 | 2870 | 347 | 36335 |
|  | 20585 | 69.045 | 12091 | 3155.4 | 434.92 |  |
|  | 49.086 | 299E-7 | 33.013 | 25.822 | 17.771 |  |
|  | 42.30 | 0.14 | 22.45 | 5.62 | 0.68 | 71.18 |
| Total | 28919 | 97 | 16986 | 4433 | 611 | 51046 |
|  | 56.65 | 0.19 | 33.28 | 8.68 | 1.20 | 100.00 |

Table 4: Chi-Square Tests for equipment days

o   **Equipment Days Analysis:** 50% of the churners have equipment days less than one year, Most equipment days up to 2years.

**Demographic Analysis**: Marital Status and churn Behavior

Fig 11 : Marital status in in Churn behavior

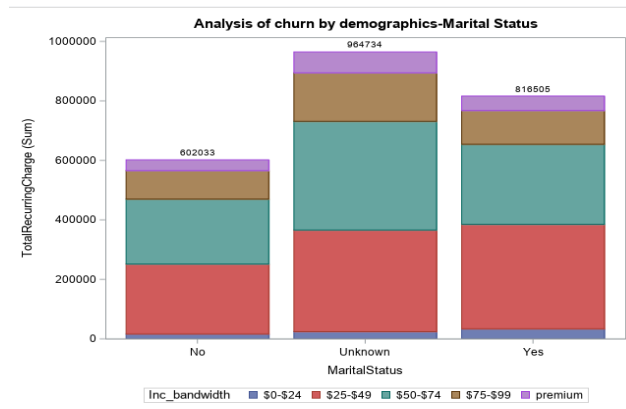Table 5: Chi-Square Tests for Marital status in Churn behavior

### Table of Churn by Marital Status

| Churn | No | Yes | Unknown | Total |
|---|---|---|---|---|
| Yes | 3441 | 5323 | 5947 | 14711 |
| | 3658.8 | 5374.9 | 5677.3 | |
| | 12.965 | 0.502 | 12.817 | |
| | 6.74 | 10.43 | 11.65 | 28.82 |
| No | 9255 | 13328 | 13753 | 36336 |
| | 9037.2 | 13276 | 14023 | |
| | 5.2492 | 0.2033 | 5.189 | |
| | 18.13 | 26.11 | 26.94 | 71.18 |
| Total | 12696 | 18651 | 19700 | 51047 |
| | 24.87 | 36.54 | 38.59 | 100.00 |

**Marital Status**: Churn rate is related to marital status, excluding the unknown category, married people are leaving more than those not married(10.43%) compared to 6.74% for those who did not churn. The rate is proportional to the total, as we have more married clients than not.
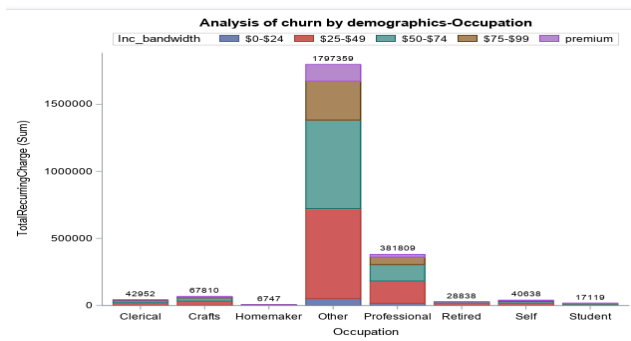
**Demographic Analysis**: Occupation



Fig 12 : Occupation in Churn behavior

Table 6: Chi-Square Tests for Occupation in Churn behavior

### Table of Churn by Occupation

| Churn | Professional | Crafts | Other | Self | Retired | Homemaker | Clerical | Student | Total |
|---|---|---|---|---|---|---|---|---|---|
| Yes | 2467 | 426 | 10932 | 243 | 185 | 51 | 289 | 118 | 14711 |
| | 2523.1 | 437.75 | 10846 | 253.31 | 211.24 | 45.245 | 284.15 | 109.8 | |
| | 1.2457 | 0.3156 | 0.675 | 0.42 | 3.2595 | 0.732 | 0.0828 | 0.6126 | |
| | 4.83 | 0.83 | 21.42 | 0.48 | 0.36 | 0.10 | 0.57 | 0.23 | 28.82 |
| No | 6288 | 1093 | 26705 | 636 | 548 | 106 | 697 | 263 | 36336 |
| | 6231.9 | 1081.2 | 26791 | 625.69 | 521.76 | 111.75 | 701.85 | 271.2 | |
| | 0.5044 | 0.1278 | 0.2733 | 0.1701 | 1.3196 | 0.2964 | 0.0335 | 0.248 | |
| | 12.32 | 2.14 | 52.31 | 1.25 | 1.07 | 0.21 | 1.37 | 0.52 | 71.18 |
| Total | 8755 | 1519 | 37637 | 879 | 733 | 157 | 986 | 381 | 51047 |
| | 17.15 | 2.98 | 73.73 | 1.72 | 1.44 | 0.31 | 1.93 | 0.75 | 100.00 |

**Occupation Analysis**: Most people leaving in the other category, Professionals leaving at a significantly high proportion. Numbers proportionally related to the totals in the subs data base, Shows company not able to handle numbers, a clear sign of dissatisfaction amongst customers.

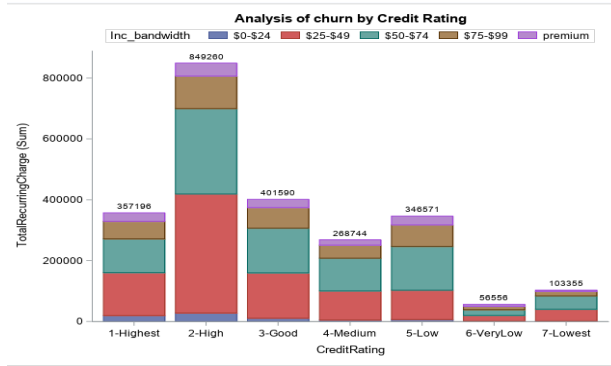MONTHLY CHARGES ANALYSIS OF CHURN AND CREDIT RATING

Fig 13 : Credit Rating in Churn behavior

**Table of Churn by CreditRating**

| | | | | CreditRating | | | | |
|---|---|---|---|---|---|---|---|---|
| Churn | 1-Highest | 4-Medium | 3-Good | 6-VeryLow | 2-High | 5-Low | 7-Lowest | Total |
| Yes | 2628 | 1399 | 2608 | 316 | 5712 | 1436 | 612 | 14711 |
| | 2455.9 | 1543.8 | 2423.6 | 331.99 | 5473.5 | 1872.9 | 609.22 | |
| | 12.058 | 13.583 | 14.024 | 0.7701 | 10.392 | 101.92 | 0.0126 | |
| | 5.15 | 2.74 | 5.11 | 0.62 | 11.19 | 2.81 | 1.20 | 28.82 |
| No | 5894 | 3958 | 5802 | 836 | 13281 | 5063 | 1502 | 36336 |
| | 6066.1 | 3813.2 | 5986.4 | 820.01 | 13519 | 4626.1 | 1504.8 | |
| | 4.8817 | 5.4993 | 5.6777 | 0.3118 | 4.2072 | 41.265 | 0.0051 | |
| | 11.55 | 7.75 | 11.37 | 1.64 | 26.02 | 9.92 | 2.94 | 71.18 |
| Total | 8522 | 5357 | 8410 | 1152 | 18993 | 6499 | 2114 | 51047 |
| | 16.69 | 10.49 | 16.48 | 2.26 | 37.21 | 12.73 | 4.14 | 100.00 |

Table 7: Chi-Square : Credit Rating Churn behavior

**Credit Rating in Churn behavior**: company has fairly good client acquisition with best credit rating (1 and 2) with over 51%. More deactivation also happening from those with good credit over 16%. Company has significant portion of unhealthy customers (19.13%) with bad credit, deactivation does exits in these groups (5 to 7)
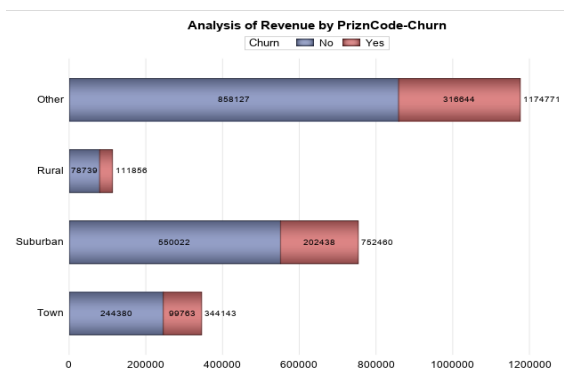
## CHURN &REVENUE ANALYSIS BY AREA CODE



Fig 14 : Location and Churn behavior

**Table of Churn by PrizmCode**

| | | PrizmCode | | | |
|---|---|---|---|---|---|
| Churn | Suburban | Town | Other | Rural | Total |
| Yes | 4609 | 2276 | 7057 | 769 | 14711 |
| | 4719.9 | 2187 | 7105.2 | 698.85 | |
| | 2.6058 | 3.6186 | 0.3271 | 7.0417 | |
| | 9.03 | 4.46 | 13.82 | 1.51 | 28.82 |
| No | 11769 | 5313 | 17598 | 1656 | 36336 |
| | 11658 | 5402 | 17550 | 1726.2 | |
| | 1.055 | 1.465 | 0.1324 | 2.8509 | |
| | 23.06 | 10.41 | 34.47 | 3.24 | 71.18 |
| Total | 16378 | 7589 | 24655 | 2425 | 51047 |
| | 32.08 | 14.87 | 48.30 | 4.75 | 100.00 |

Table 8: Chi-Square : Location and Churn behavior

**Location and Churn behavior**: More churn in suburban areas excluding the others category. The revenue generated and loss is proportional to the clientele data base, more numbers more revenue and more terminations. Need to investigate the other category
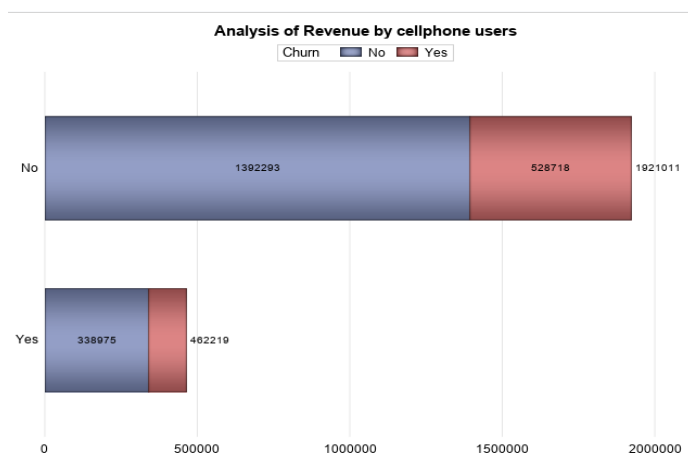
## CELL PHONE USER ANALYSIS

Fig 15 : Phone use/gadgets in Churn behavior



Table 9: Chi-Square : Phone use/gadgets in Churn behavior

**Phone use/gadgets in Churn behavior**: Even new cellphone user leaving (5.41%). More churners on clients using own cell phones

SAMPLE T-TEST FOR STUDY VARIABLE ANALYSIS

| Variable | Statistically significant | p-value |
|---|---|---|
| Monthly Revenue | Yes | 0.0001 |
| Monthly Minutes | Yes | 0.0014 |
| Total recurring charges | Yes | 0.0001 |
| Month in service | Yes | 0.0001 |
| Current Equipment Days | Yes | 0.0001 |
| New cellphone users | No | 0.0821 |
| Credit rating | Yes | 0.0001 |
| Prizm Code | Yes | 0.0003 |
| Occupation | No | 0.1714 |
| Marital Status | Yes | 0.0001 |

Table 1 Two sample t-test between study variables in Churn=Yes and Churn=No

**Inferential Analysis** : Visually, it appears that all variables are key determinants of churn behavior except occupation or whether someone got new cellphone or not.

SCATTERPLOT MATRIX

Plot of continuous data against other continuous variables to assess for multicollinearity between independent variables. High positive correlation between monthly revenue and month in service as well as equipment days was noted from the analysis.

**Outlier Detection**

**Outliers** are extreme values that deviate from other observations on data , they may indicate a variability in a measurement, experimental errors or a novelty. These can be problematic as they can cause a skewness in the data, pulling the mean away from the median in the direction of the outlier.
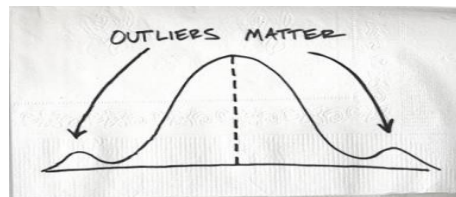


Fig 9: Outlier Analysis by bands

PROC univariate is useful in compiling the five number summary: minimum, Q1, median, Q2, max. Interquartile range can also be reported in this procedures. The outliers have been defined by the following two equations:
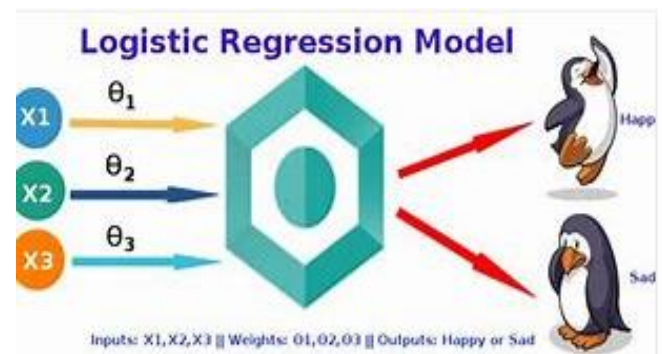
$out\ lier>$ Q3+(**3**\*IQR)(the upper bound) $outlier$ < Q1-(**3**\*IQR)(the lower bound)

The outliers were not removed from the analysis, separate grouped were created upper bound, lower bound and data values within range formulated part of the  bivariate analysis.

**Logistic Regression**

A useful solution will be able to predict the most likely causes of churn and flag any customers at risk. For example, how much of your churn is coming from simple recurrent monthly charges, occupation or location or how much of it is coming from that poorly strategized credit rating process. The summary model results from analysis are set as below:



| Odds Ratio Estimates | | | |
|---|---|---|---|
| Effect | Point Estimate | 95% Wald Confidence Limits | |
| MonthlyRevenue | 1.004 | 1.003 | 1.005 |
| TotalRecurringCharge | 0.994 | 0.993 | 0.995 |
| MonthlyMinutes | 1.000 | 1.000 | 1.000 |
| MonthsInService | 0.991 | 0.989 | 0.994 |
| CurrentEquipmentDays | 1.001 | 1.001 | 1.001 |
| DroppedCalls | 1.007 | 1.004 | 1.009 |
| BlockedCalls | 1.003 | 1.001 | 1.005 |
| AgeHH1 | 0.995 | 0.995 | 0.996 |

- Backward Feature Elimination
  - Above are Features left (from framework)
  - Model concordance = 58.5%
- **Odds ratio =**

$e \wedge (0.434 + 1.004(MR) + 0.994(TRC) + 1.000(MM) + 0.991(MIS) + 1.001(CED) + 1.007(DC) + 1.003(BC) + 0.995(AG))$

For every one unit change in **monthly revenue**, the log odds of churn (versus not churn) increases by 1.004 etc.

**Summary findings**

- High customer disconnections at 28.82% against industry standard of 2.5% per 5million customers
- Monthly revenue has obvious correlation with TotalRecurringCharge & monthly in service.
- Bulk of customers leaving around 20months in service, customer retention less than 2yrs, sign of customer frustration.
- Most contributing income bandwidth/segments are those paying between $25-$49(48.76%) and $50-$74(28.71%)
- More disconnections from these segments that are generating more revenue
- Disconnections highly correlated to equipment's days
- More churners from married customers (10.43%), professionals (4.83%)
- Company losing more of customers with good credit ~16%.
-  In high revenue generating groups, High and Good Credit Rating are more easily to churn.
- Use of gadgets as a hold or strategy to create loyalty proved ineffective as activations with new gadgets are disconnecting from the service.
- More revenue contributions from suburban areas and towns compared to rural communities

**Recommendations**

- o Company to focus on retention policy by conducting various market campaign strategies like running monthly or quarterly promotions, pricing strategies, expand product offering and or improve service delivery.
- o To better service the market, the company must identify needs by segments and link product offer to specific segments, (say age groups, geographical locations, marital status or occupation)
- o Optimise customer loyalty drivers, company must create loyalty groups and document loyalty indicators, say for stable married people, professionals etc.
- o Company must create powerful customer focused value proposition:
- o Relevance (solve problems), Quantifiable value (unique benefits) & differentiations.
- o Company must devise ways for managing total customer experience by conducting service quality surveys and customer support periodically.
- o Continue to focus more on customers with good credit as this will increase sales as revealed by the analysis and thrive to retain them.

APPENDIX

**PROC STATEMENTS USED**

1. PROC FORMAT enables you to define your own informats and formats for variables
2. PROC PRINT  prints the observations in a SAS data set
3. PROC CONTENTS shows the contents of a SAS data set and prints the directory of the SAS library
4. PROC MEANS  provides data summarization tools to compute descriptive statistics for variables across all observations and within groups of observations
5. PROC FREQ produces one-way to n-way frequency and contingency (crosstabulation) tables
6. PROC SORT orders SAS data set observations by the values of one or more character or numeric variables
7. PROC UNIVARIATE provides a variety of descriptive measures, graphical displays, and statistical methods, which you can use to summarize, visualize, analyze, and model the statistical distributions of numeric variables
8. PROC SGPLOT creates one or more plots and overlays them on a single set of axes
9. PROC SGPANEL creates a panel of graph cells for the values of one or more classification variables
10. PROC SGSCATTER creates a paneled graph of scatter plots for multiple combinations of variables
11. PROC TTEST performs t tests and computes confidence limits for one sample, paired observations, two independent samples, and the AB/BA crossover design
12. PROC LOGISTIC investigate the relationship between these discrete responses and a set of explanatories.

**SAS CODES**

```
/*import the data as csv file*/
```

```
LIBNAME MUW "C:\Users\Admin\Desktop\SAS PROJECT";


PROC IMPORT OUT= MUW.SAS_PROJECT
             DATAFILE= "C:\Users\Admin\Desktop\SAS PROJECT\Project Data F
iles\Telco Churn Data.csv"
             DBMS=CSV REPLACE;
     GETNAMES=YES;
     DATAROW=2;
RUN;

*data preparation -view data set..10obs;
 PROC PRINT DATA = MUW.SAS_PROJECT  (OBS = 10);
TITLE "CHURN DATA SET"; RUN;

/*scan for duplicates and remove them*/
TITLE "Count of Distinct Customer IDs in RSA1";
PROC SQL;  SELECT COUNT(CustomerID)AS TOTAL_COUNT,
COUNT(DISTINCT CustomerID) AS UNIQUE_COUNT  FROM MUW.SAS_PROJECT  ;
QUIT;

/*scan for duplicates and remove them*/

Proc sort data=MUW.SAS_PROJECT  out=MUW.Telcom_data nodupkey;  */sort data
and scan for duplicates/*;
by CustomerID;
run;

/* No duplicates found*/

/*summary data*/
PROC CONTENTS DATA=MUW.Telcom_data;
RUN;

/*data type*/

PROC CONTENTS DATA=MUW.Telcom_data varnum short;
RUN;
/*scan for missing data under var CHURN*/

ods table onewayfreqs=temp;
proc freq data=MUW.Telcom_data;
 table _all_ / missing;
format _numeric_ nmissfmt. _character_ $missfmt.;
run;


/*DESCRIPTIVE ANALYSIS-CHURN*/

PROC FORMAT;
VALUE $CHAR " "="MISSING"
                 OTHER= "NOT MISSING"
RUN;
/* 28.15 missing, label as missing*/
/*label as missing*/

PROC FREQ DATA=MUW.Telcom_data;
```

```
TABLE Churn/MISSING;
FORMAT _CHARACTER_$CHAR.;
RUN;

/*delete missing data on key column-Churn*/

data MUW.Telcom;
set MUW.Telcom_data;
if Churn eq"NA" then delete;
run;

/*descriptive Analysis-Barchart for Churn status; */

proc sgplot data=MUW.Telcom;
   Vbar Churn / datalabel  colormodel=twocolorramp;
   keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
    title"BarChart- Churn Status Analysis (ex_missing";
run;
/*analysis of churn with missing data*/

proc sgplot data=MUW.Telcom_data;
   vbar Churn / datalabel colormodel=twocolorramp;
    title"BarChart- Churn Status Analysis(inc_missing";
run;


*UNIVARIATE ANALYSIS*;

/*scan for missing monthinservice Analysis*/

proc means data=MUW.Telcom N NMISS MIN MEAN STD MAX;
var MonthsInService;
run;

*NO MISSING UNDER MONTHINSERVICE;
ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\monthinservice_graph_1"
style= MoonFlower notoc ;
TITLE "Month in Service Tren Analysis-Deactivation rate";
proc sgpanel data = MUW.Telcom;
panelby Churn / columns = 2;
histogram MonthsInService;
density MonthsInService;  density MonthsInService/ type = kernel;
colaxis label = 'Month in service';
run;
ods pdf close;


/*scan for missing MONTHLY REVENUE*/
TITLE "Missing month in revenue values";
proc means data=MUW.Telcom maxdec=2 N NMISS MIN MEAN STD MAX;
var MonthlyRevenue;
run;

/*imputation of missing: mean=58.83*/
proc sql;
create table Telco_revenue as
select *, coalesce (MonthlyRevenue,58.83) as MonthlyRevenue_updated
```

```sas
from MUW.Telcom;
run;


/*confirm  missing MONTHLY REVENUE values*/
TITLE "missing month in revenue data";
proc means data=Telco_revenue maxdec=2 N NMISS MIN MEAN STD MAX;
var MonthlyRevenue_updated;
run;


ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\monthinservice_graph_1"
style= MoonFlower notoc; ;
TITLE "Monthly Revenue Distribution  Analysis ";
proc sgpanel data = Telco_revenue;
panelby Churn / columns = 2;
histogram MonthlyRevenue_updated;
density MonthlyRevenue_updated / type = kernel;
colaxis label = 'Monthly revenue';
run;
ods pdf close;

Title 'Analysis of  monthly revenue-Active and Deactivated Accounts';
proc sgplot data=Telco_revenue noborder;
hbar Churn / response=MonthlyRevenue_updated stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
/*format MonthsInService 8.0;*/
run;


/*scan for missing MONTHLY REVENUE*/
TITLE "missing month in minutes data";
proc means data=MUW.Telcom maxdec=2 N NMISS MIN MEAN STD MAX;
var MonthlyMinutes;
run;


/*imputation of missing: mean=525.65*/
proc sql;
create table Telco_minutes as
select *, coalesce (MonthlyMinutes,525.65) as MonthlyMinutes_updated
from MUW.Telcom;
run;


/*confirm  missing MONTHLY minutes values*/
TITLE "missing month in revenue data";
proc means data=Telco_minutes maxdec=2 N NMISS MIN MEAN STD MAX;
var MonthlyMinutes_updated;
run;



ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\monthinservice_graph_1"
style= MoonFlower notoc; ;
TITLE "monthly revenue analysis";
proc sgpanel data = Telco_minutes;
panelby Churn / columns = 2;
histogram MonthlyMinutes_updated;
```

```
density MonthlyMinutes_updated;  density MonthlyMinutes_updated/ type =
kernel;
colaxis label = 'Monthly revenue';
run;
ods pdf close;


/*scan for missing MONTHLY RECURRING CHARGES*/
TITLE "missing TotalRecurringCharge values";
proc means data=MUW.Telcom maxdec=2 N NMISS MIN MEAN STD MAX;
var TotalRecurringCharge;
run;

/*imputation of missing: mean=46.83*/
proc sql;
create table Telco_mcharges as
select *, coalesce (TotalRecurringCharge,46.83) as
TotalRecurringCharge_updated
from MUW.Telcom;
run;

/*confirm  missing MONTHLY CHARGES values*/
TITLE "missing total recurring charges values";
proc means data=Telco_mcharges maxdec=2 N NMISS MIN MEAN STD MAX;
var TotalRecurringCharge_updated;
run;


TITLE "Total recurring charges analysis-consolidated";
proc sgplot data=MUW.Telcom;
histogram TotalRecurringCharge;
density TotalRecurringCharge;
title"Monthly total revenue distribution";
run;


ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\monthinservice_graph_1"
style= MoonFlower notoc; ;
TITLE "Total recurring charges analysis";
proc sgpanel data = Telco_mcharges;
panelby Churn / columns = 2;
histogram TotalRecurringCharge_updated ;
density TotalRecurringCharge_updated;density TotalRecurringCharge_updated/
type = kernel;
colaxis label = 'Total recurring charges';
run;
ods pdf close;

/* "Total recurring charges analysis*/

TITLE "Total recurring charges analysis";
proc sgplot data = Telco_mcharges ;
histogram TotalRecurringCharge_updated / showbins;
density TotalRecurringCharge_updated;
density TotalRecurringCharge_updated / type = kernel;
yaxis grid;  xaxis label = 'total charges';
keylegend / location = inside    position = topright;
```

```sas
  title 'Distribution of total charges';
run;


/*scan for missing dropped calls*/
TITLE "missing TotalRecurringCharge values";
proc means data=MUW.Telcom maxdec=2 N NMISS MIN MEAN STD MAX;
var DroppedCalls
;
run;

/* no missing for dropped calls*/

/*distribution of monthly revenue*/

proc sgplot data=MUW.Telcom;
histogram TotalRecurringCharge;
density TotalRecurringCharge;
title"Monthly total revenue distribution";
run;


/*scan for missing under equipment days*/
TITLE "missing equipment days values";
proc means data=MUW.Telcom maxdec=2 N NMISS MIN MEAN STD MAX;
var CurrentEquipmentDays;
run;

/*delete missing data values*/

data MUW.Equip_days;
set MUW.Telcom_data;
if CurrentEquipmentDays eq . then delete;
run;

/*confirm deletion*/
proc means data=MUW.Equip_days maxdec=2 N NMISS MIN MEAN STD MAX;
var CurrentEquipmentDays;
run;

/*Equipment days analysis*/

ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\monthinservice_graph_1"
style= MoonFlower notoc; ;
TITLE "Equipment days analysis";
proc sgpanel data = Telco_mcharges;
panelby Churn / columns = 2;
histogram CurrentEquipmentDays ;
density CurrentEquipmentDays;density CurrentEquipmentDays/ type = kernel;
colaxis label = 'Total recurring charges';
run;
ods pdf close;

/*Analysis month in equipments days */

proc format;
value equip_days low-365 ="from 0 days to 1yr"
```

```
                              365-730      ="Between one to 2yrs"
                              730-1095     ="Between 2 to 3yrs"
                              1095-1460    ="Between 3 to 4yrs"
                              1460-1825    ="Between 4 to 5yrs"
                              1825-high    ="Above 5yrs";
                              run;
/*Bar graph Analysis for Equipment days*/

proc sgplot data=MUW.Equip_days noborder;
vbar CurrentEquipmentDays /datalabel;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format CurrentEquipmentDays equip_days.;
title "Grouped equipment days  Analysis";
run;

/*Chi square Analysis for equipment days*/

proc freq data=MUW.Telcom order=data;
   tables Churn*CurrentEquipmentDays / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   format CurrentEquipmentDays equip_days.;
   title 'Chi-Square Tests for equipment days';
run;

/*Analysis of equipment days*/
Title 'Analysis of equipment days';
proc sgplot data=MUW.Telcom noborder;
hbar CurrentEquipmentDays / response=CurrentEquipmentDays stat=freq
group=CurrentEquipmentDays displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
/*keylegend / location=outside position=top fillheight=10 fillaspect=2 ;*/
format CurrentEquipmentDays equip_days.;
run;



/*BIVARIATE & CHI-SQUARE ANALYSIS*/


/*INCOME BANDWITH ANALYSIS*/

data MUW.Telcom_Income; /*formating dates-use this method*/
length Inc_bandwidth $15;
set Telco_mcharges;
     if 0.1< TotalRecurringCharge <25 then Inc_bandwidth ='$0-$24';
     else if     24<TotalRecurringCharge<50 then Inc_bandwidth ='$25-$49';
     else if    49< TotalRecurringCharge<75 then Inc_bandwidth ='$50-$74';
     else if 74<TotalRecurringCharge<100 then Inc_bandwidth ='$75-$99';
     else if TotalRecurringCharge >99 then Inc_bandwidth='premium';
     run;

proc sgplot data=MUW.Telcom_Income noborder;
vbar Inc_bandwidth /datalabel;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
title "BarGraph- Recurring Charges by bandwidth Analysis";
```

```sas
run;

proc freq data=MUW.Telcom_Income order=data;
    tables Churn*Inc_bandwidth / expected cellchi2 norow nocol chisq;
    output out=ChiSqData n nmiss pchi lrchi;
    /*weight Count;*/
    title 'Chi-Square Tests for Inc_bandwidth';
run;

/*Analysis month in service groups*/

proc format;
value months_grp low-12 ="New Activations"
                        12-24 ="Between one to 2yrs"
                        24-36 ="Between 2 to 3yrs"
                        36-48       ="Between 3 to 4yrs"
                        48-60 ="Between 4 to 5yrs"
                        60-high     ="Above 5yrs";
                        run;

proc sgplot data=MUW.Telcom noborder;
vbar MonthsInService /datalabel;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService months_grp.;
title "Grouped Month in service  Analysis";
run;

proc freq data=MUW.Telcom_Income order=data;
    tables Churn*MonthsInService / expected cellchi2 norow nocol chisq;
    output out=ChiSqData n nmiss pchi lrchi;
    format MonthsInService months_grp.;
    title 'Chi-Square Tests for month in service grps';
run;


Title 'Analysis of Average month in service days-Active and Deactivated
Accounts';

proc sgplot data=MUW.Telcom_data noborder;
hbar Churn / response=MonthsInService stat=mean group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
/*format MonthsInService 8.0;*/
run;

ods graphics on;
proc anova data = MUW.Telcom;
    class Churn;
    model  MonthlyRevenue= Churn;
    means Churn/scheffe;
    title "customer monthly revenue analysis";
run;
ods graphics off;
```

```
Title 'Analysis of  dropped calls-Active and Deactivated Accounts';

proc sgplot data=MUW.Telcom noborder;
hbar Churn / response=DroppedCalls stat=sum group=Churn displaybaseline=auto
barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
/*format MonthsInService 8.0;*/
run;


ods graphics on;
proc anova data = MUW.Telcom;
   class Churn;
   model  DroppedCalls= Churn;
   means Churn/scheffe;
   title "customer monthly revenue analysis";
run;
ods graphics off;


proc sgplot data=MUW.Telcom;
   hbar MaritalStatus / colormodel=twocolorramp;
   keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
run;



Title 'Analysis of churn by demographics-Marital Status';
proc sgplot data=MUW.Telcom_Income;
vbar MaritalStatus / datalabel response=TotalRecurringCharge
Group=Inc_bandwidth;
run;



proc freq data=MUW.Telcom order=data;
   tables Churn*MaritalStatus / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   title 'Chi-Square Tests for Marital status';
run;



/*analysis of occupation vs churn*/

Title 'Analysis of churn by demographics-Occupation';
proc sgplot data=MUW.Telcom_Income noborder;
vbar Occupation /datalabel response=TotalRecurringCharge
Group=Inc_bandwidth;;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
/*format MonthsInService 8.0;*/
run;

proc freq data=MUW.Telcom order=data;
   tables Churn*Occupation / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
```

```sas
      title 'Chi-Square Tests for Occupation';
run;


/*analysis of CreditRating vs churn*/

Title 'Analysis of churn by Credit Rating';
proc sgplot data=MUW.Telcom_Income noborder;
vbar CreditRating /datalabel response=TotalRecurringCharge
Group=Inc_bandwidth;;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
/*format MonthsInService 8.0;*/
run;



proc freq data=MUW.Telcom order=data;
   tables Churn*CreditRating / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   title 'Chi-Square Tests for Credit Rating';
run;

/*primcode area Analysis*/
title "Analysis of churn by PriznCode";
proc freq data=MUW.Telcom order=data;
   tables Churn*PrizmCode
 / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/

run;

title "Analysis of Revenue by PriznCode-Churn";
proc sgplot data=MUW.Telcom noborder;
hbar PrizmCode / response=TotalRecurringCharge stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;

/*New/Not cellphone users  Analysis*/
title "Analysis of Not New cellphone users";
proc freq data=MUW.Telcom order=data;
   tables Churn*NotNewCellphoneUser / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   run;

title "Analysis of New cellphone users";
proc freq data=MUW.Telcom order=data;
   tables Churn*NewCellphoneUser / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   run;
```

```sas
title "Analysis of Revenue by cellphone users";
proc sgplot data=MUW.Telcom noborder;
hbar NewCellphoneUser / response=TotalRecurringCharge stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;




/*Revenue Analysis by Credit Rating, grouped by Churn*/

title "Analysis of Revenue by Credit Rating-Churn";
proc sgplot data=MUW.Telcom noborder;
hbar CreditRating / response=TotalRecurringCharge stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;




/*Revenue Analysis by IncomeGroup, grouped by Churn status*/

title "Analysis of Revenue by IncomeGroup-Churn";
proc sgplot data=MUW.Telcom noborder;
hbar IncomeGroup / response=TotalRecurringCharge stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;




/*Revenue Analysis by IncomeGroup, grouped by Churn status*/

title "Analysis of Revenue by marital status-Churn";
proc sgplot data=MUW.Telcom noborder;
hbar MaritalStatus / response=TotalRecurringCharge stat=sum group=Churn
displaybaseline=auto barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;

title "Analysis of blocked calls by Churn";
```

```
proc sgplot data=MUW.Telcom noborder;
hbar Churn / response= BlockedCalls stat=sum group=Churn displaybaseline=auto
barwidth=0.6
seglabel datalabel dataskin=pressed;
yaxis display=(noline noticks nolabel);
xaxis display=(noline noticks nolabel) grid;
keylegend / location=outside position=top fillheight=10 fillaspect=2 ;
format MonthsInService 8.0;
run;

ods graphics on;
proc anova data = MUW.Telcom;
    class Churn;
    model  BlockedCalls= Churn;
    means Churn/scheffe;
    title "customer BlockedCalls analysis";
run;
ods graphics off;


/*OUTLIER DETECTION*/

       *monthly revenue;
       *total recurring charges
       *monthly minutes;
       *dropped calls;
       *blocked calls;


/*calculate quartiles and inter quartiles for Monthly revenue*/


proc sgplot data= Telco_revenue;
vbox   MonthlyRevenue_updated;
run;


proc means data=Telco_revenue maxdec=2;
var MonthlyRevenue_updated;
output out =revenue p25=Q1 p75=Q3 qrange =IQR;
run;

data revenue_01;
set revenue;
lower_limit =Q1-(3*IQR);
upper_limit=Q3+(3*IQR);
drop _TYPE_ _FREQ_;
run;
proc print data=revenue_01; run;

/*create catesian product*/

proc sql;
create table revenue_02 as
select A.*,B.*
from Telco_revenue as A, revenue_01 as B
;
```

```
quit;


data revenue_03;
set revenue_02;
if MonthlyRevenue_updated le lower_limit then range ="below lower limit";
else if  MonthlyRevenue_updated ge upper_limit then range ="above upper
limit";
else range ="within range"
;
run;

/*bivariate analysis for monthly revenue range*/

proc freq data=revenue_03 order=data;
    tables Churn*range / expected cellchi2 norow nocol chisq;
    output out=ChiSqData n nmiss pchi lrchi;
    /*weight Count;*/
    title 'Chi-Square Tests for monthly revenue range';
run;


/*calculate quartiles and inter quartiles for monthly recurring charges */

proc means data=Telco_mcharges maxdec=2;
var TotalRecurringCharge_updated;
output out =charges p25=Q1 p75=Q3 qrange =IQR;
run;

data charges_01;
set charges;
lower_limit =Q1-(3*IQR);
upper_limit=Q3+(3*IQR);
drop _TYPE_ _FREQ_;
run;
proc print data=charges_01; run;

/*create catesian product*/

proc sql;
create table charges_02 as
select A.*,B.*
from Telco_mcharges as A, charges_01 as B
;
quit;


data charges_03;
set charges_02;
if TotalRecurringCharge_updated le lower_limit then charge_range ="below
lower limit";
else if  TotalRecurringCharge_updated ge upper_limit then charge_range
="above upper limit";
else charge_range ="within range"
;
run;
```

```sas
proc freq data=charges_03 order=data;
   tables Churn*charge_range / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   title 'Chi-Square Tests for monthly charge range';
run;


/*calculate quartiles and inter quartiles for monthly recurring charges */

proc means data=Telco_mcharges maxdec=2;
var TotalRecurringCharge_updated;
output out =charges p25=Q1 p75=Q3 qrange =IQR;
run;

data charges_01;
set charges;
lower_limit =Q1-(3*IQR);
upper_limit=Q3+(3*IQR);
drop _TYPE_ _FREQ_;
run;
proc print data=charges_01; run;

/*create catesian product*/

proc sql;
create table charges_02 as
select A.*,B.*
from Telco_mcharges as A, charges_01 as B
;
quit;


data charges_03;
set charges_02;
if TotalRecurringCharge_updated le lower_limit then charge_range ="below
lower limit";
else if  TotalRecurringCharge_updated ge upper_limit then charge_range
="above upper limit";
else charge_range ="within range"
;
run;


proc freq data=charges_03 order=data;
   tables Churn*charge_range / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   title 'Chi-Square Tests for monthly charge range';
run;


/*calculate range for monthly minutes */

proc univariate data=Telco_minutes ;
```

```
var MonthlyMinutes_updated;
run;


proc means data=Telco_minutes maxdec=2;
var MonthlyMinutes_updated;
output out =call_minutes p25=Q1 p75=Q3 qrange =IQR;
run;

data minutes_01;
set call_minutes;
lower_limit =Q1-(3*IQR);
upper_limit=Q3+(3*IQR);
drop _TYPE_ _FREQ_;
run;
proc print data=minutes_01; run;

/*create catesian product*/

proc sql;
create table minutes_02 as
select A.*,B.*
from Telco_minutes as A, minutes_01 as B
;
quit;


data minutes_03;
set minutes_02;
if MonthlyMinutes_updated le lower_limit then call_range ="below lower
limit";
else if  MonthlyMinutes_updated ge upper_limit then call_range ="above upper
limit";
else call_range ="within range"
;
run;


proc freq data=minutes_03 order=data;
   tables Churn*call_range / expected cellchi2 norow nocol chisq;
   output out=ChiSqData n nmiss pchi lrchi;
   /*weight Count;*/
   title 'Chi-Square Tests for call range in minutes';
run;


*Scatterplot Matrix;
ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\matrix_graph_1" style=
MoonFlower notoc;
proc sgscatter data = MUW.Telcom;
matrix MonthlyRevenue TotalRecurringCharge MonthsInService
CurrentEquipmentDays;
/*label MonthsInService = 'transformed';*/
title 'Scatterplot Matrix of churn  Risk Factors';
run;
ods pdf close
```

```
/*Logisic Regression model for prediction, odds ratio*/
ods pdf file ="C:\Users\Admin\Desktop\SAS PROJECT\matrix_graph_1" style=
MoonFlower notoc;
proc logistic data = MUW.Telcom desc plots(only) = oddsratio plots(only) =
roc;
class MaritalStatus CreditRating PrizmCode Occupation;
model Churn = MonthlyRevenue TotalRecurringCharge MonthlyMinutes
MonthsInService
            CurrentEquipmentDays MaritalStatus CreditRating PrizmCode
Occupation
 / expb selection = backward;
output out = outdata p = pred_prob lower = low upper = up;
title 'Logistic Regression for Churn';
 run;
ods pdf close;



                /*************end of code*******************/
```