

Telco Customer Churn Analysis



Project member:
Contact no:
Email:

Agenda

2

1. Problem Statement

2. Objectives of the Study

3. About Data

4. Data Visualization

5. Data Preprocessing

6. Hypothesis Testing

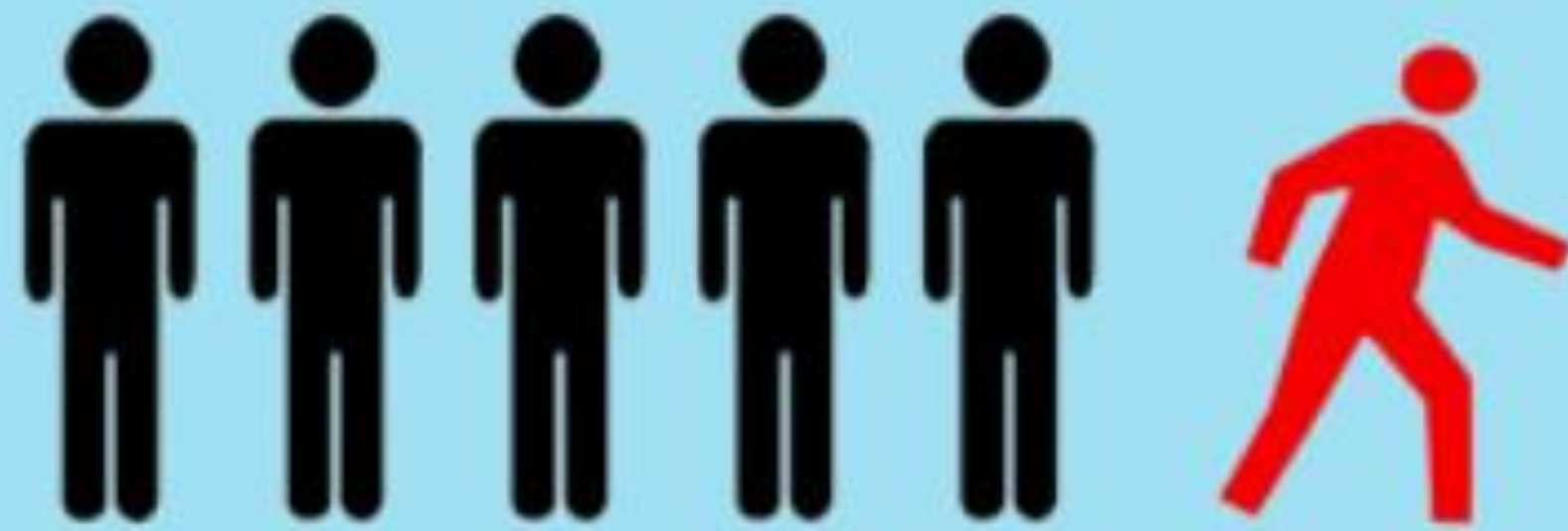
7. Feature Selection

8. Model Training and Evaluation

9. Best Model

10. Limitations, Suggestions and Conclusion

*Did you know that attracting a new customer cost **5 times** as much as keeping an existing one?*



“Annual Churn Rate of Telco Industry is 31%”

Source: Customergauge.com

Problem Statement

3

- ‘Customer churn’ is defined as the process of subscribers (either prepaid or post paid) switching from one service provider.
- With the enormous increase in the number of customers using telephone services, the marketing division for a telco company wants to attract more new customers and avoid contract termination from existing customers (churn rate).
- With proper management of customers, we can minimize the susceptibility to churn and maximize the profitability of the company. A mechanism needs to be established to analyze the attributed of profitability.
- Identifying these potential customers early on who may voluntarily churn and providing them retention incentives in form of discounts & combo offers will help the organization to retain those customers and reduce revenue loss. subscribers.

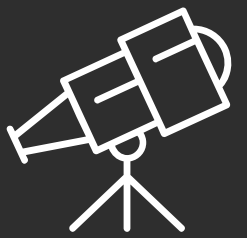
Objectives of the Study

4



Objective 01

Identify the exiting trends and most appropriate factors associate with the customer prediction



Objective 02

Build a high accurate prediction model for predicting the customers who are likely to churn from the network in near future.

If we can archive these goals, the company can easily identify the probable churn customers and they can push attractive campaigns to prevent customer churn.



About Data

5

Data Source

Kaggle

www.kaggle.com/blastchar/telco-customer-churn

Dependent Variable

Churn (Yes, No)

Categorical variable

Independent Variables

19 (16 Categorical & 3 numeric)

Which can be classified into 3 groups. Such as Demographical information, account information and service information

No of Observations

7043

Demographical Information

(all variables are categorical)

- Gender (Female, Male)
- SeniorCitizen (0, 1)
- Partner (Yes, No)
- Dependents

Service Information

(all variables are categorical)

- PhoneService (Yes, No)
- MultipleLines (No phone service, No, Yes)
- InternetServices (DSL, Fiber optic, No)
- OnlineSecurity (No internet service, No, Yes)
- OnlineBackup
- DeviceProtection
- TechSupport
- StreamingTV
- StreamingMovies

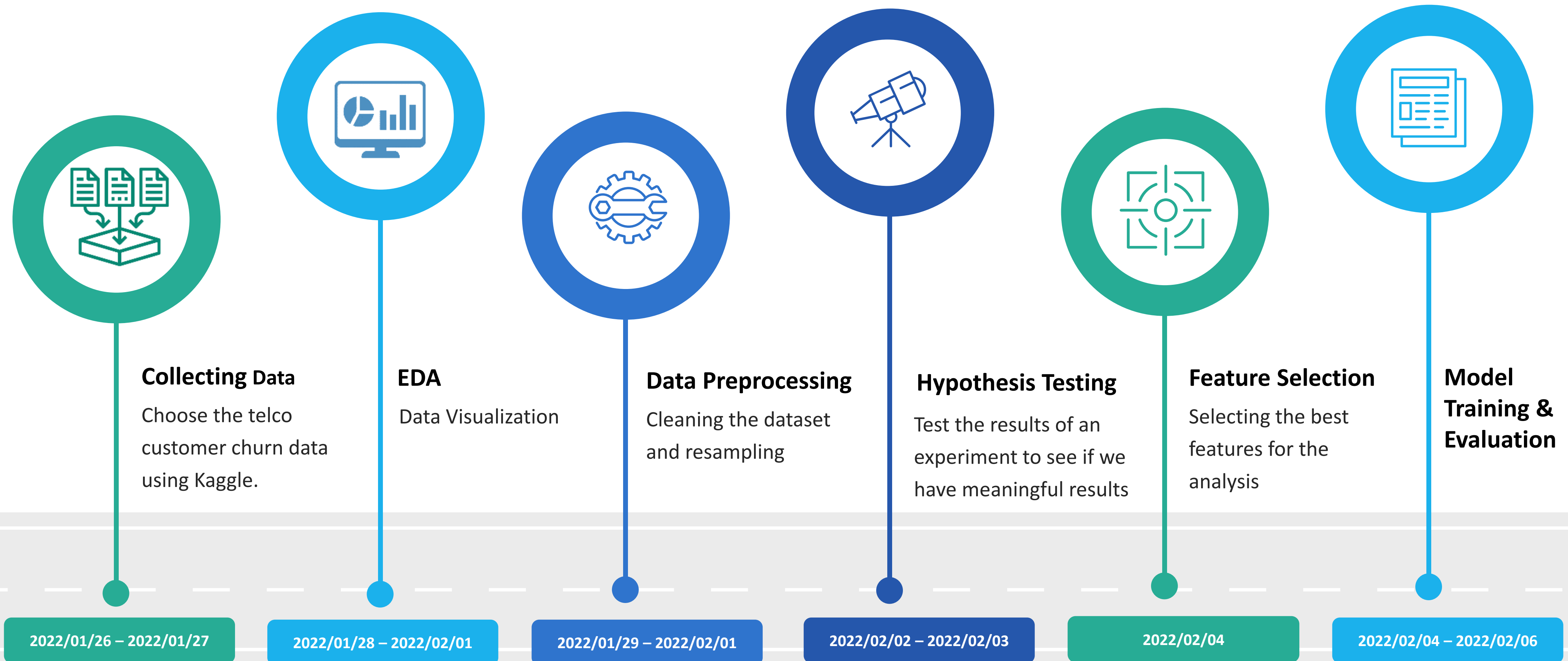
Account Information

(all variables are categorical except Tenure, MonthlyCharge and TotalCharge)

- tenure (numeric values)
- Contract ((Month-to-Month, One year, Two year)
- PaperlessBilling (Yes, No)
- PaymentMethod (Electronic check, Mailed check,...)
- MontlyCharges (numeric)
- TotalCharges (numeric)

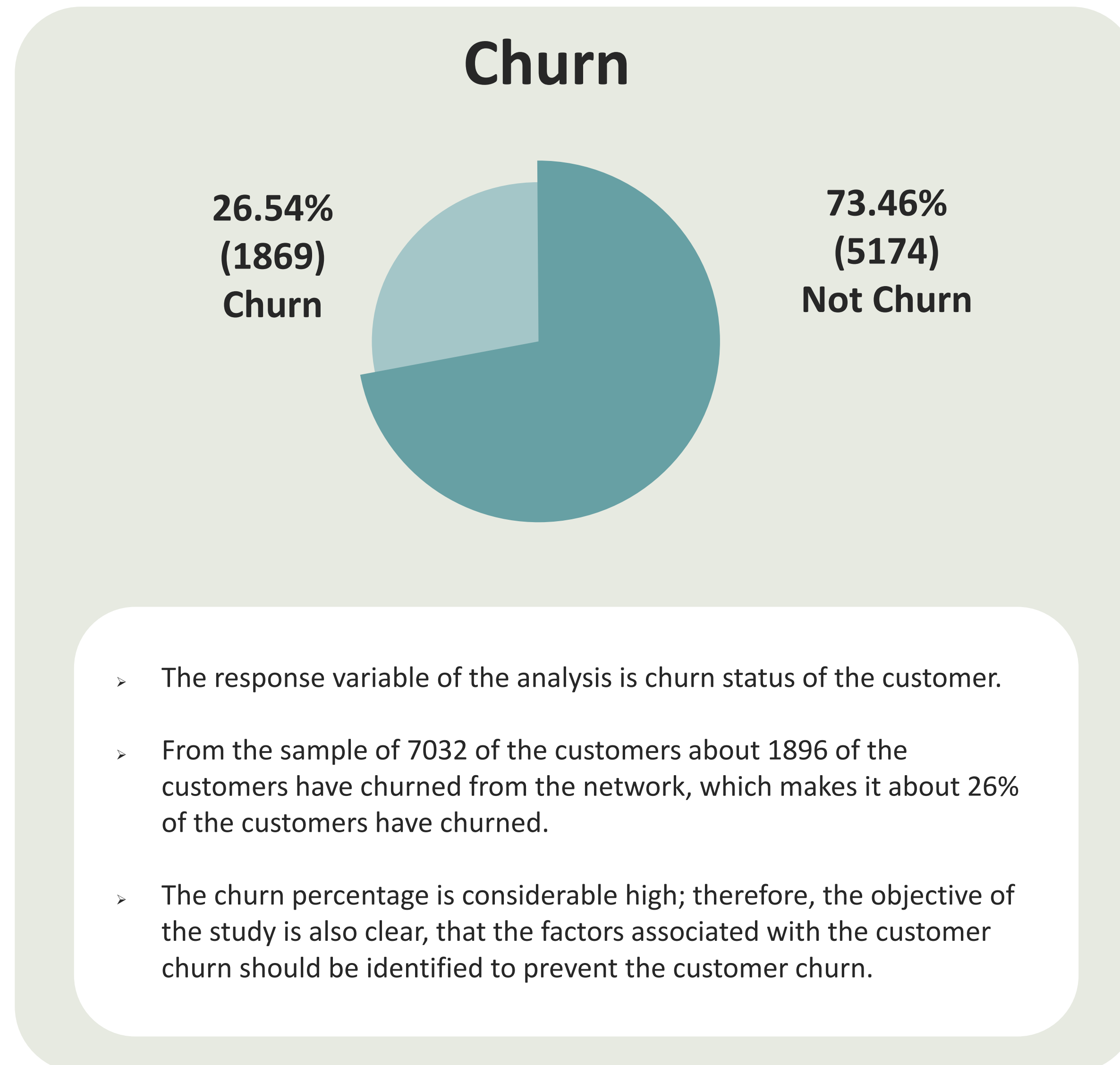
ROADMAP

6



Microsoft Excel, Orange and R Studio will be used for the Analysis

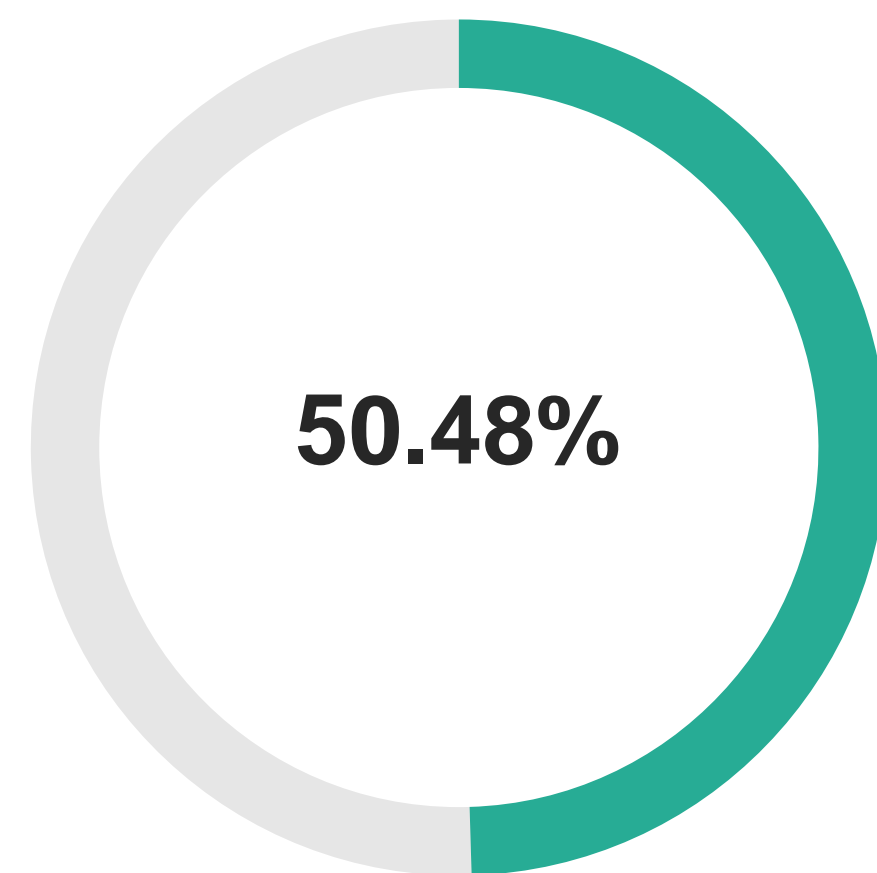
Churn Status



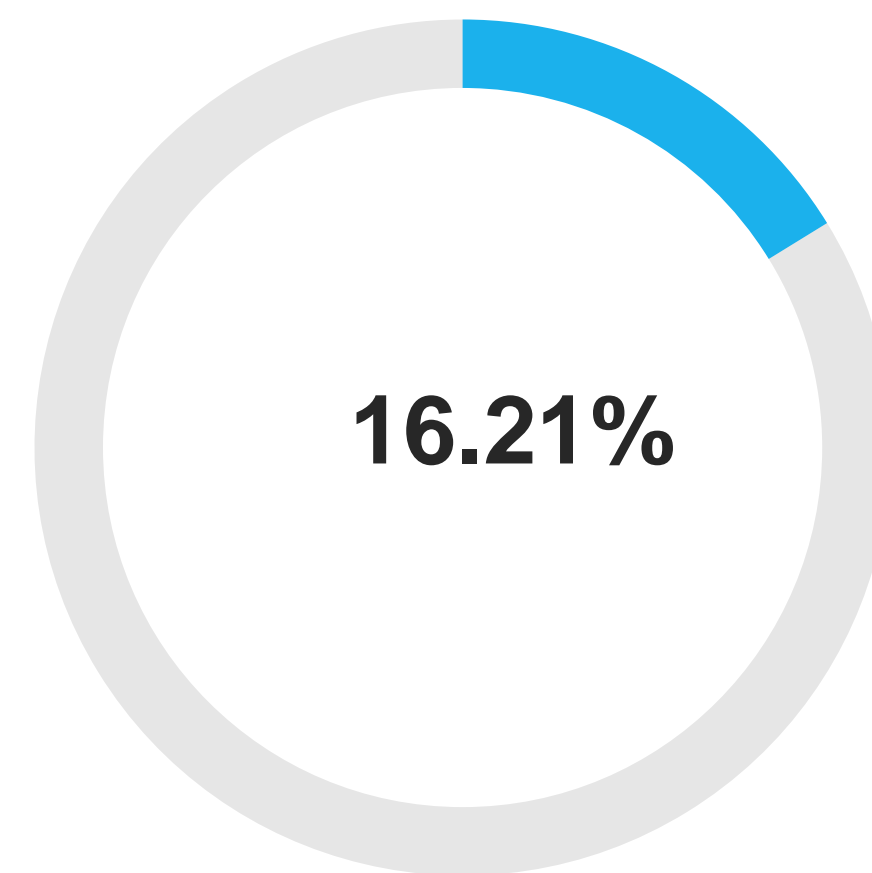
Customer Characteristics

Most of the customers are,

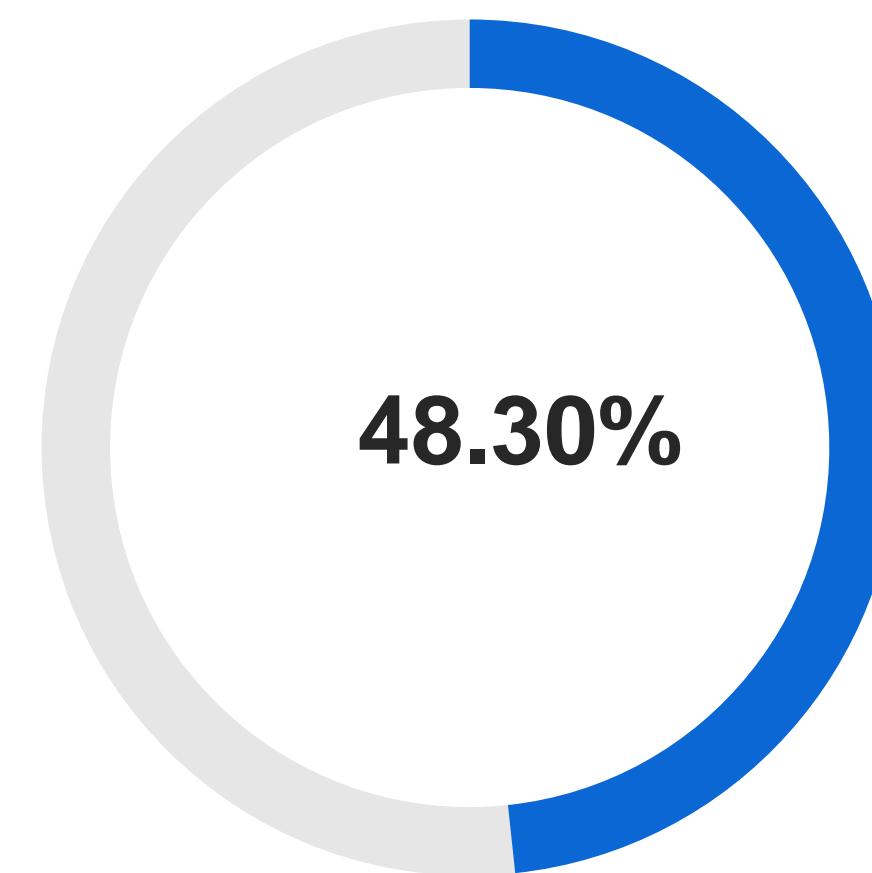
- **Male**
- **Not a senior citizen**
- **Do not a partner and a dependent person**



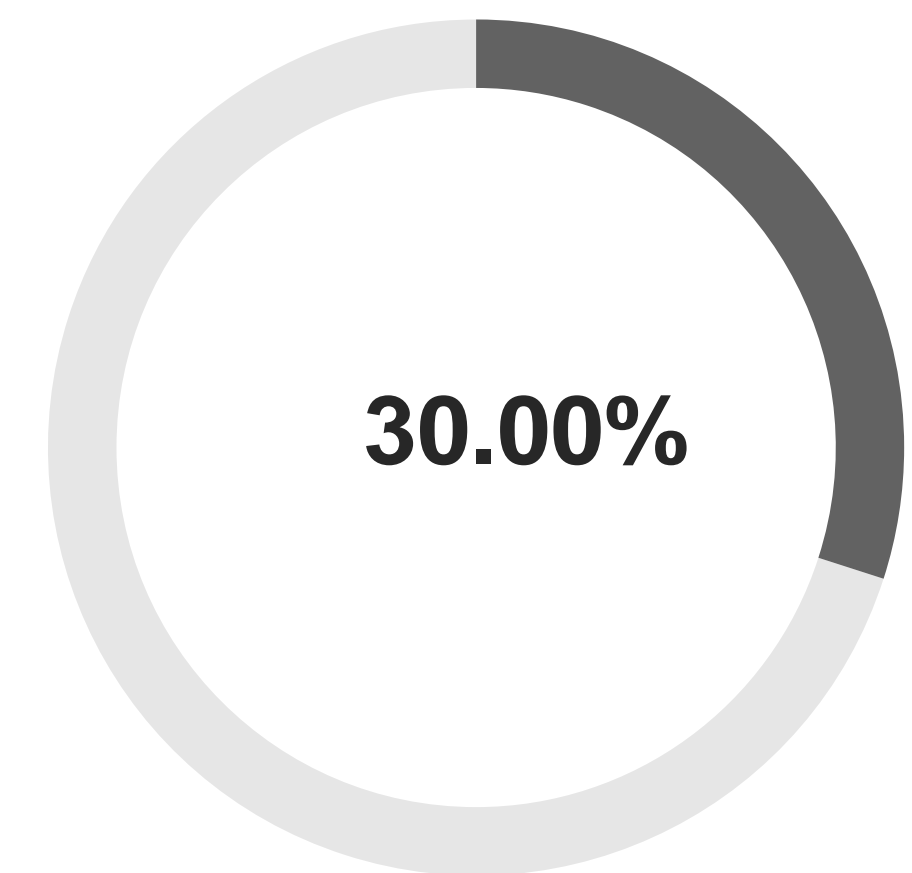
Male



Senior Citizen



Having a Partner



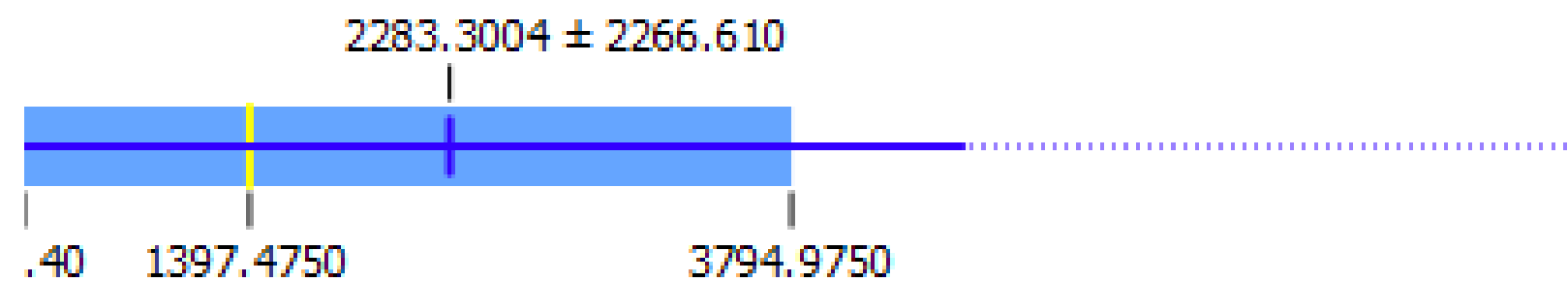
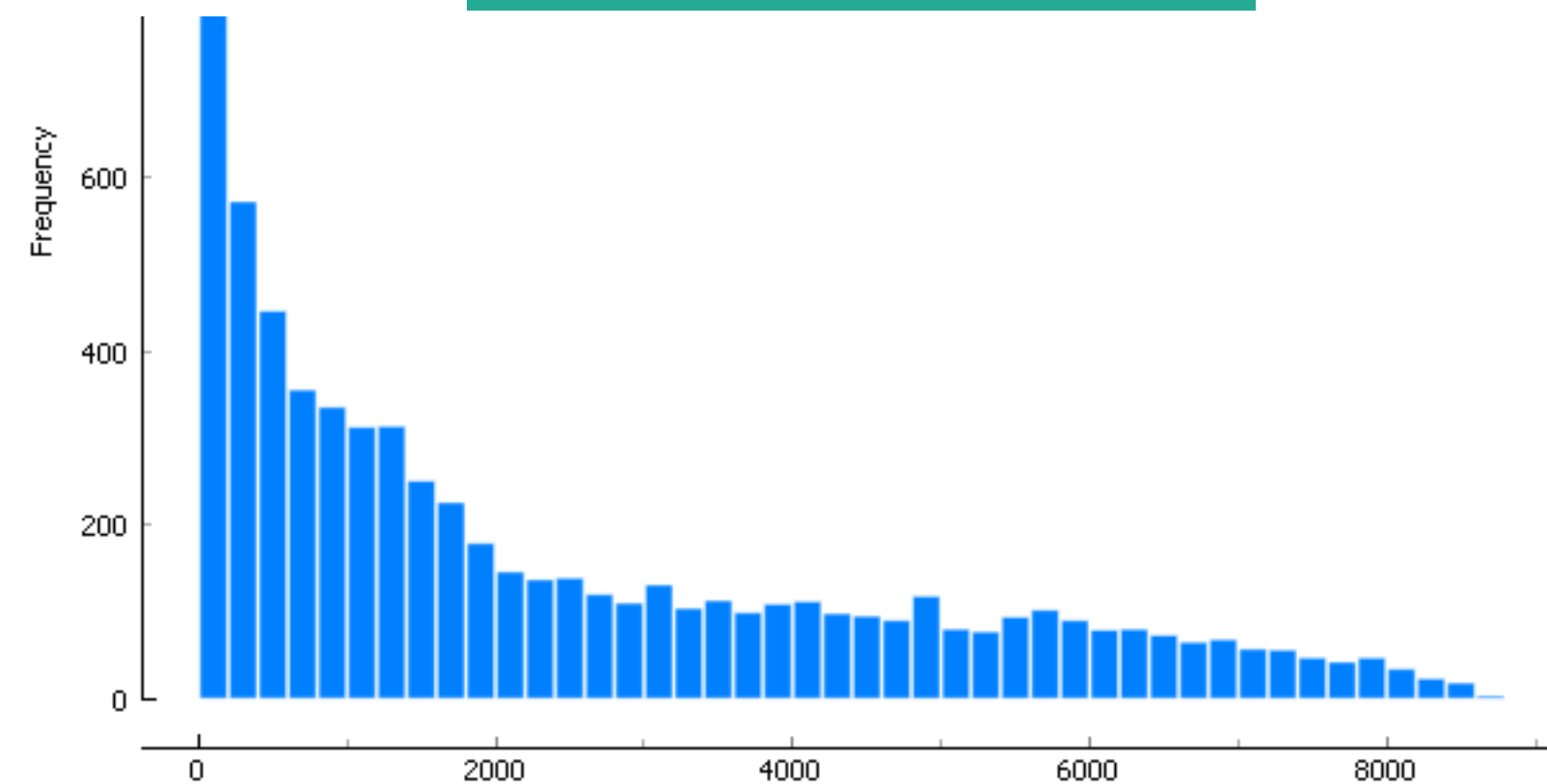
Having a Dependent

Customer Characteristics

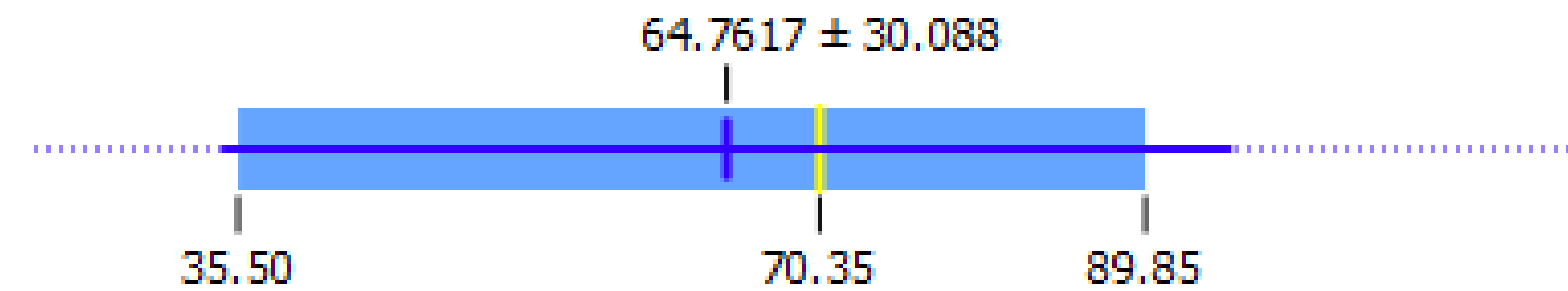
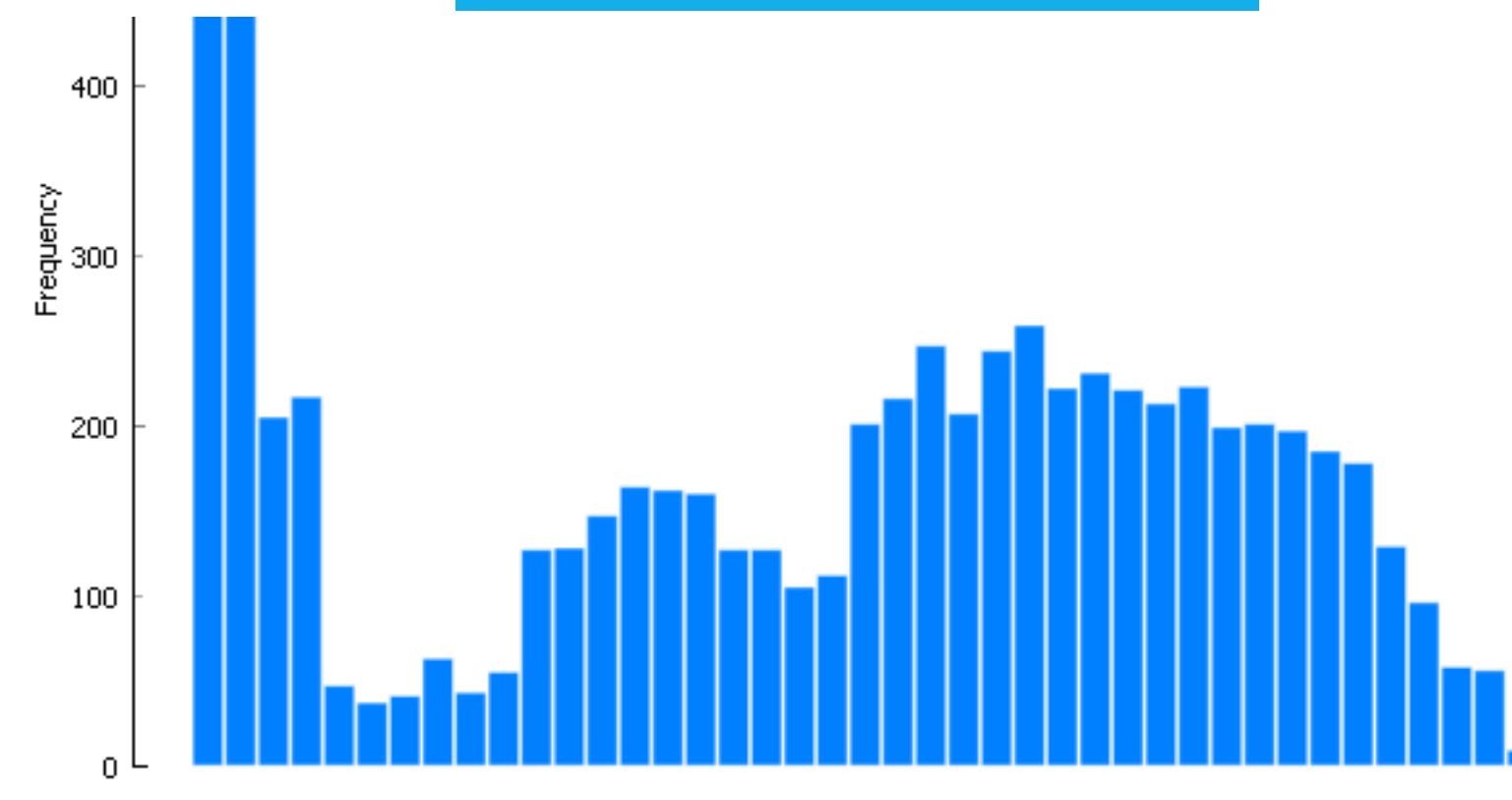
50% of the telco customers,

- Total charge is less than 1397
- Monthly charge is less than 70.35
- Stay less than 29 months in network

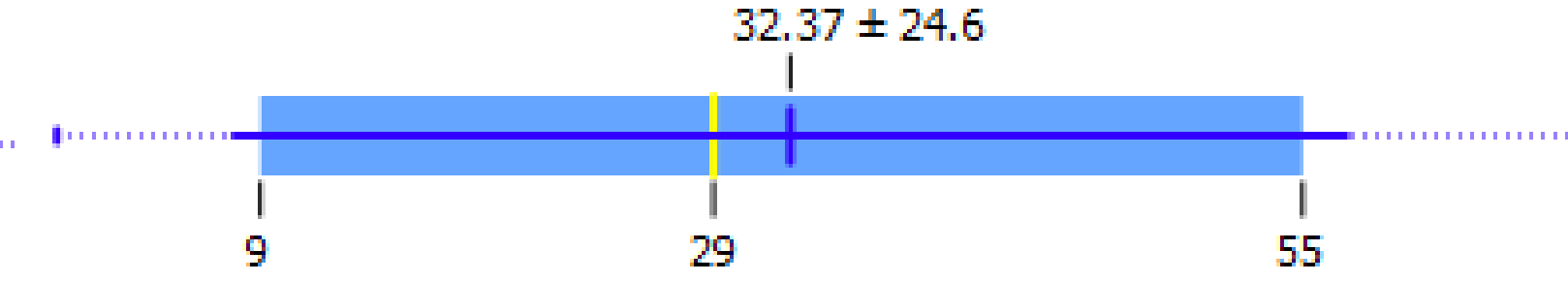
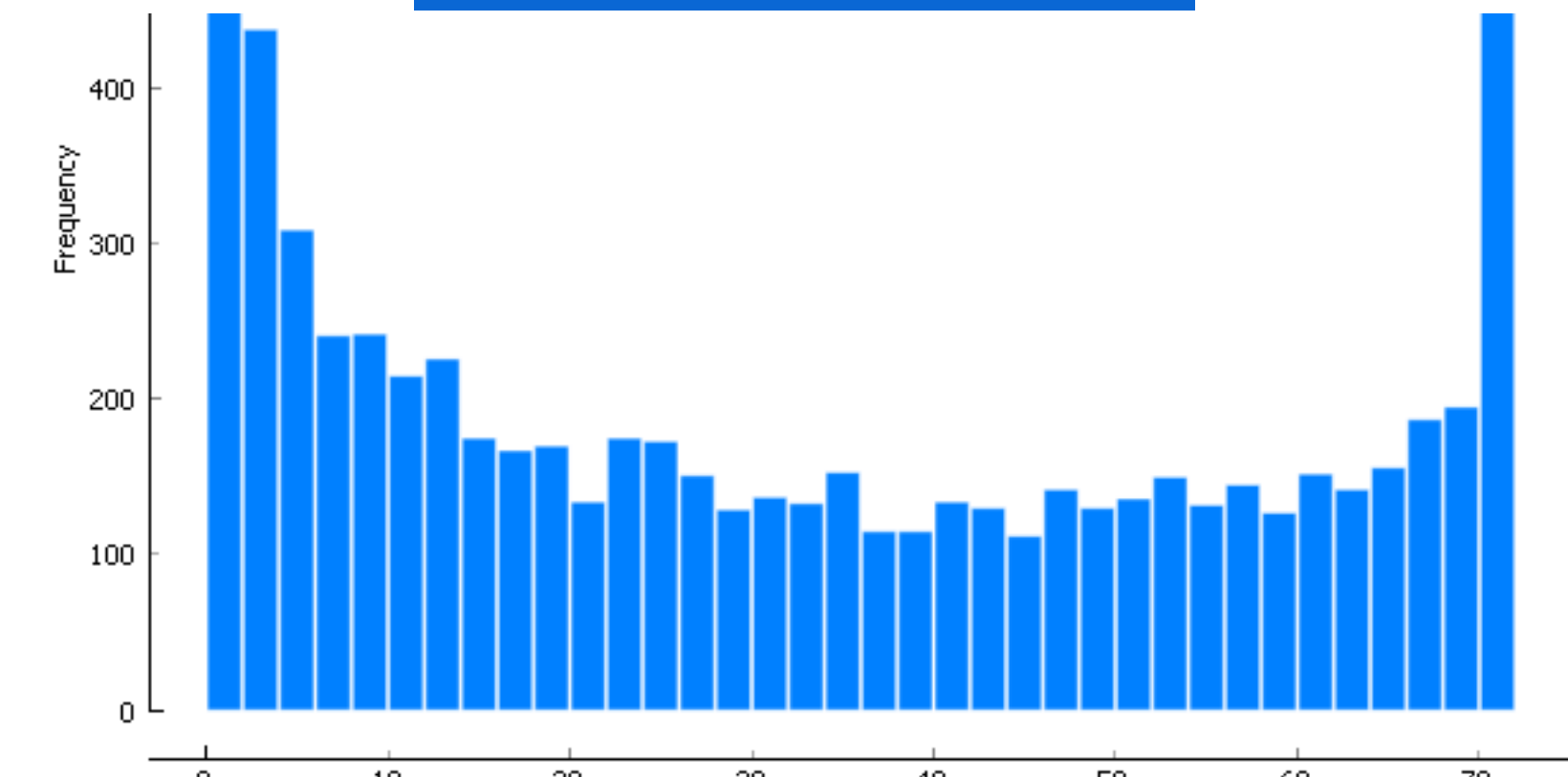
Total Charge



Monthly Charge



Tenure



1000.00 2000.00 3000.00 4000.00 5000.00 6000.00 7000.00

0 40.00 60.00 80.00 100.00 120.00

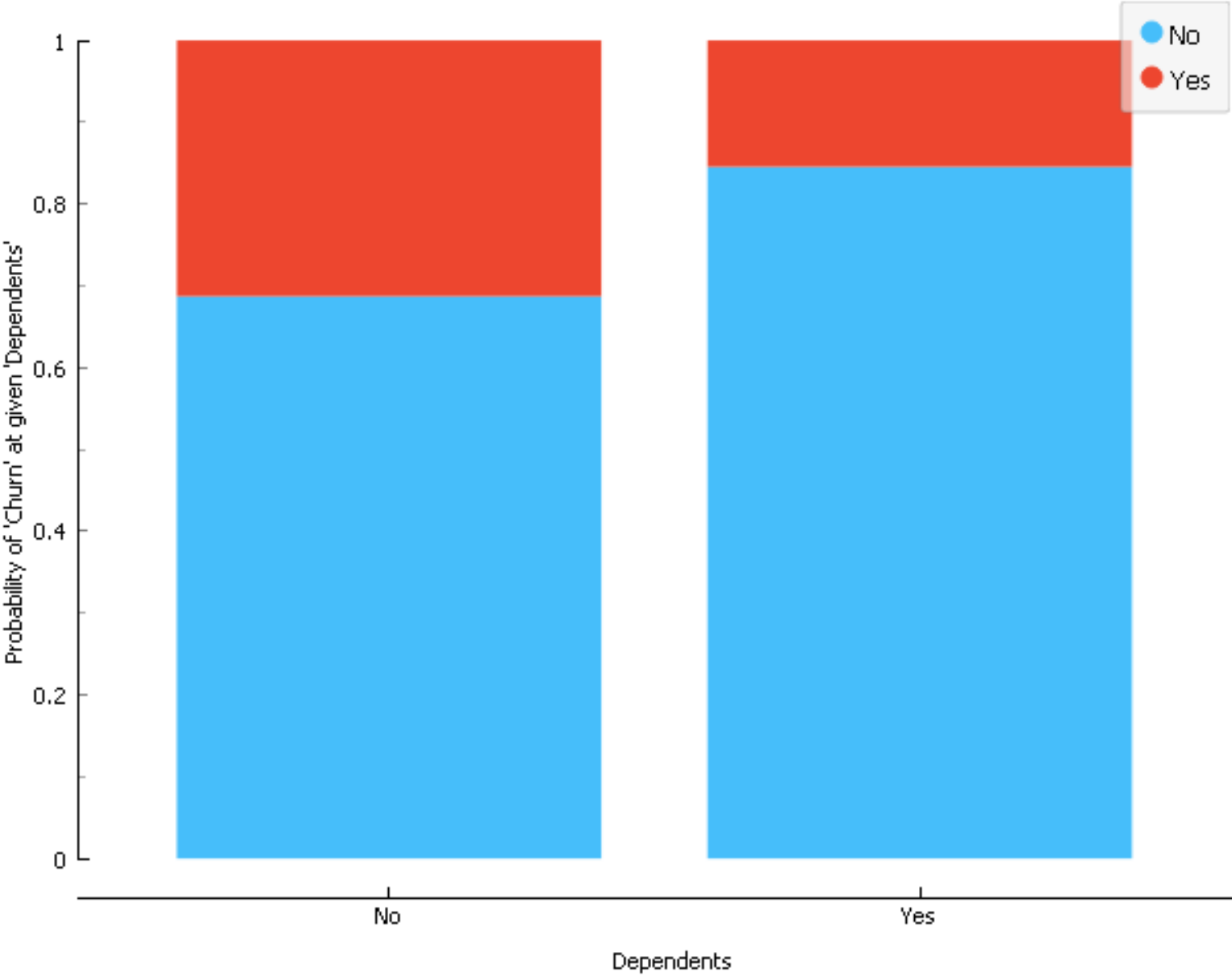
0 10 20 30 40 50 60

Customers who are more likely to Churn

Demographical Information

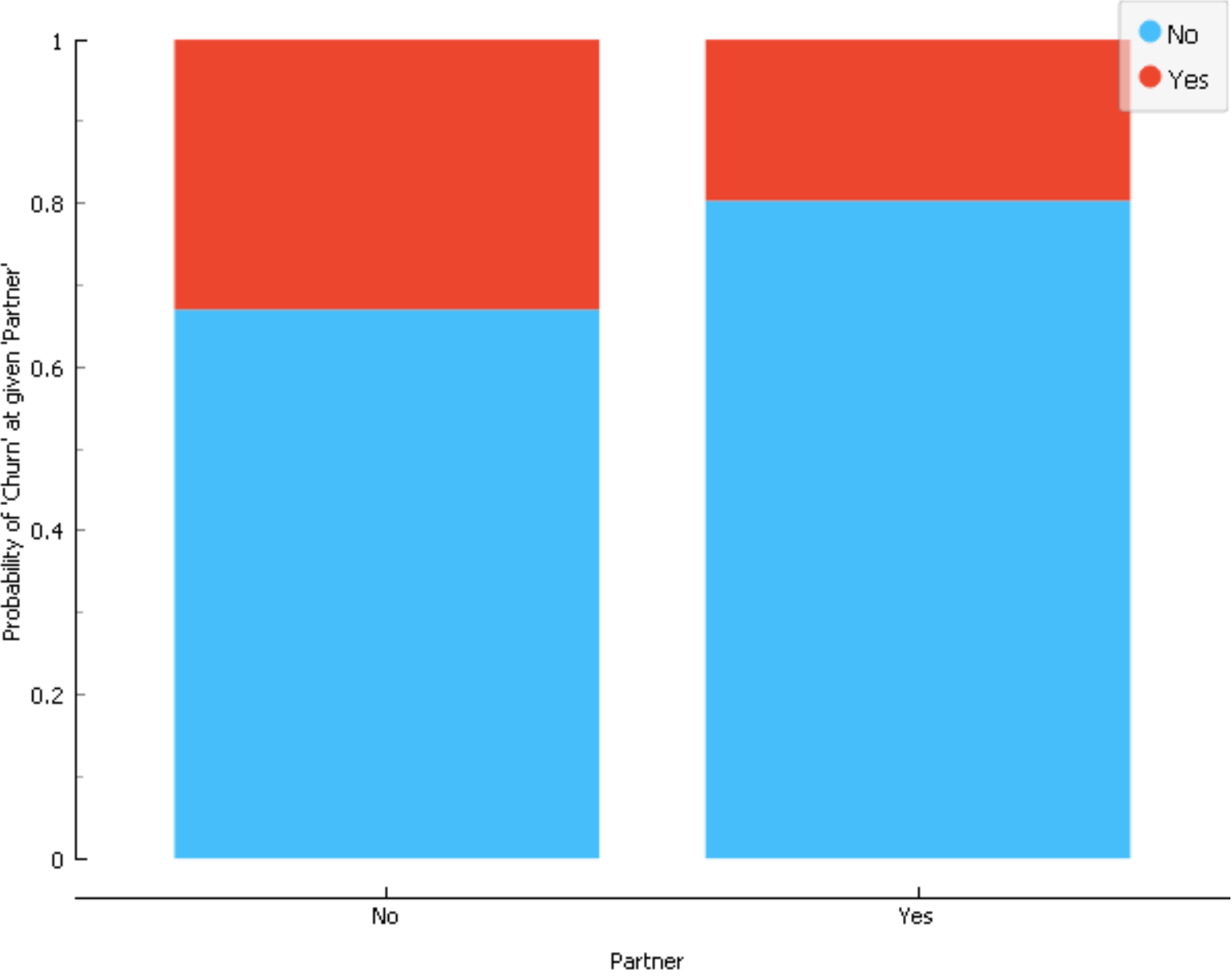
Dependents

Does not Have



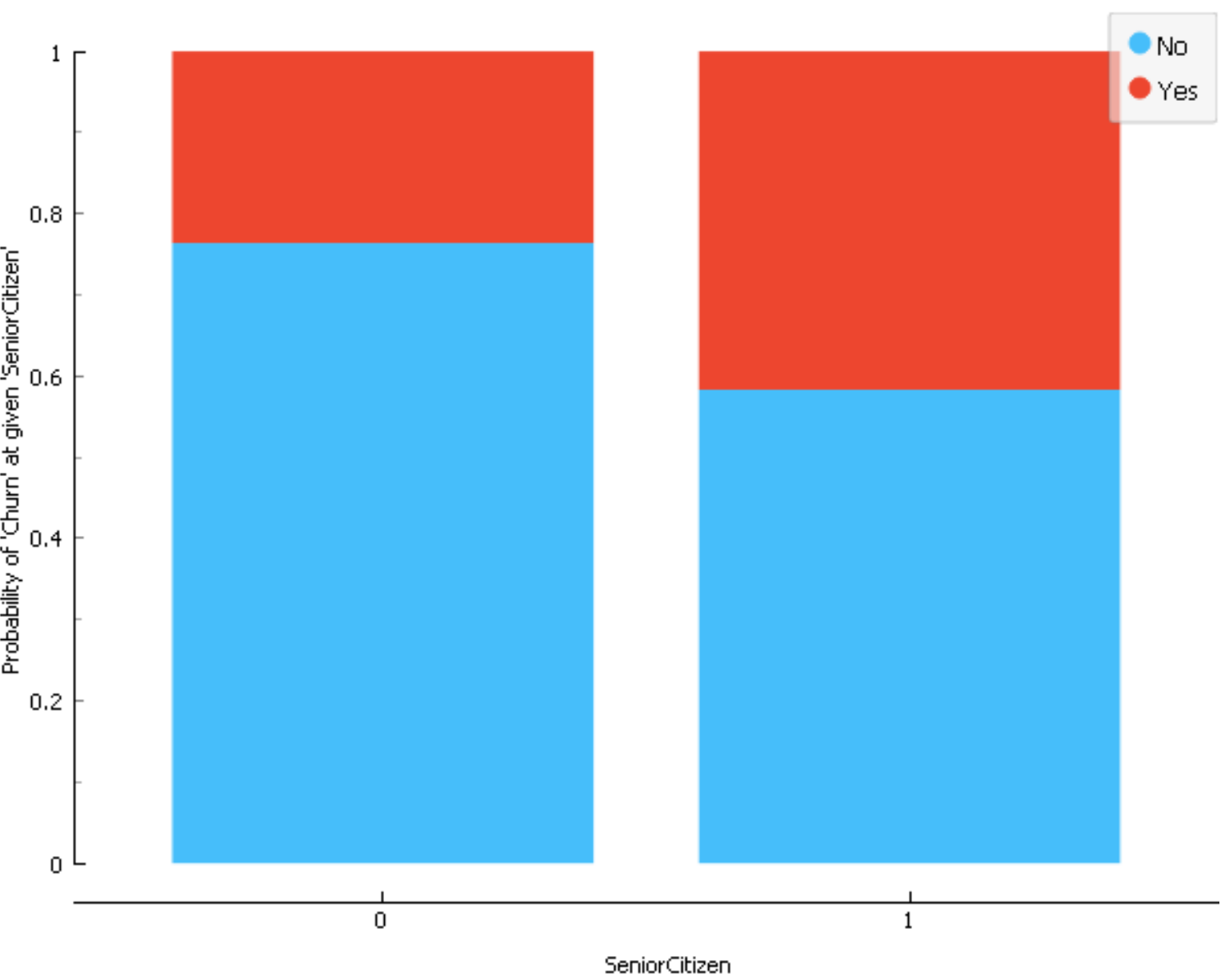
Partner

Does not Have



Senior Citizen

Yes



Customers who are more likely to Churn

Account Information & Internet Service

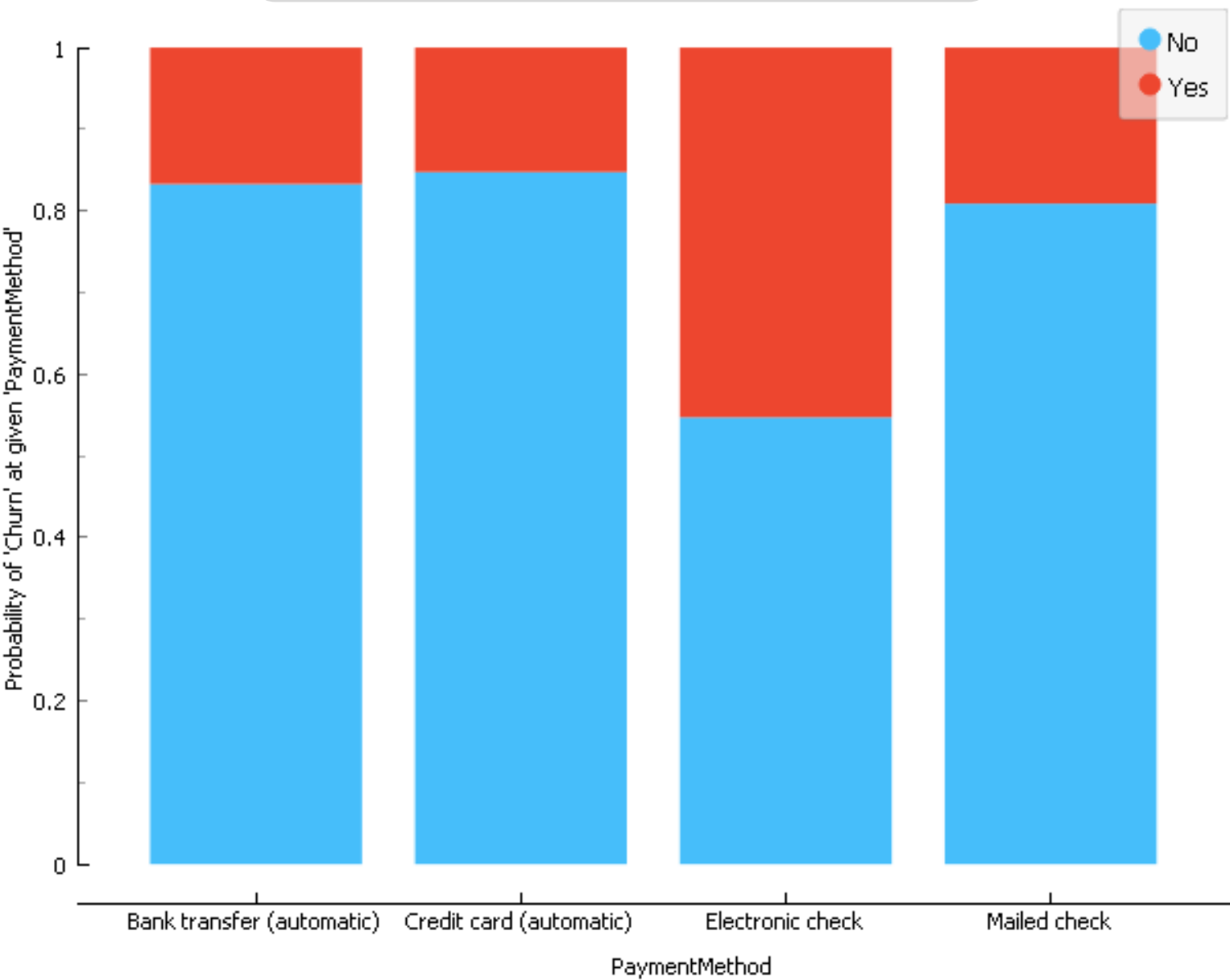
Most of the customers,

- use electronic check method as their payment method
- subscribed for paperless billing
- use fiber optics

Among these customers, who are more likely to churn,

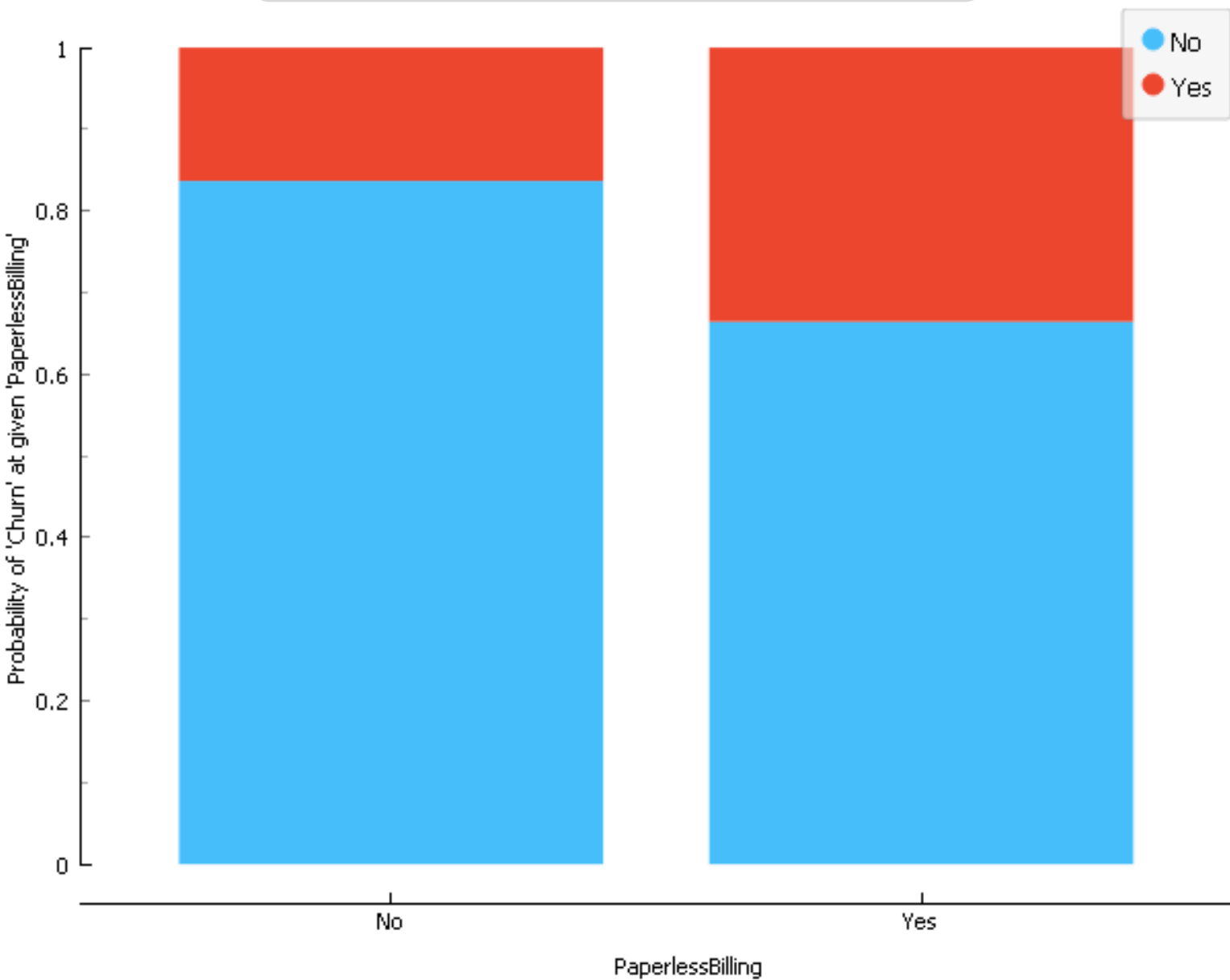
Payment Method

Electronic Check



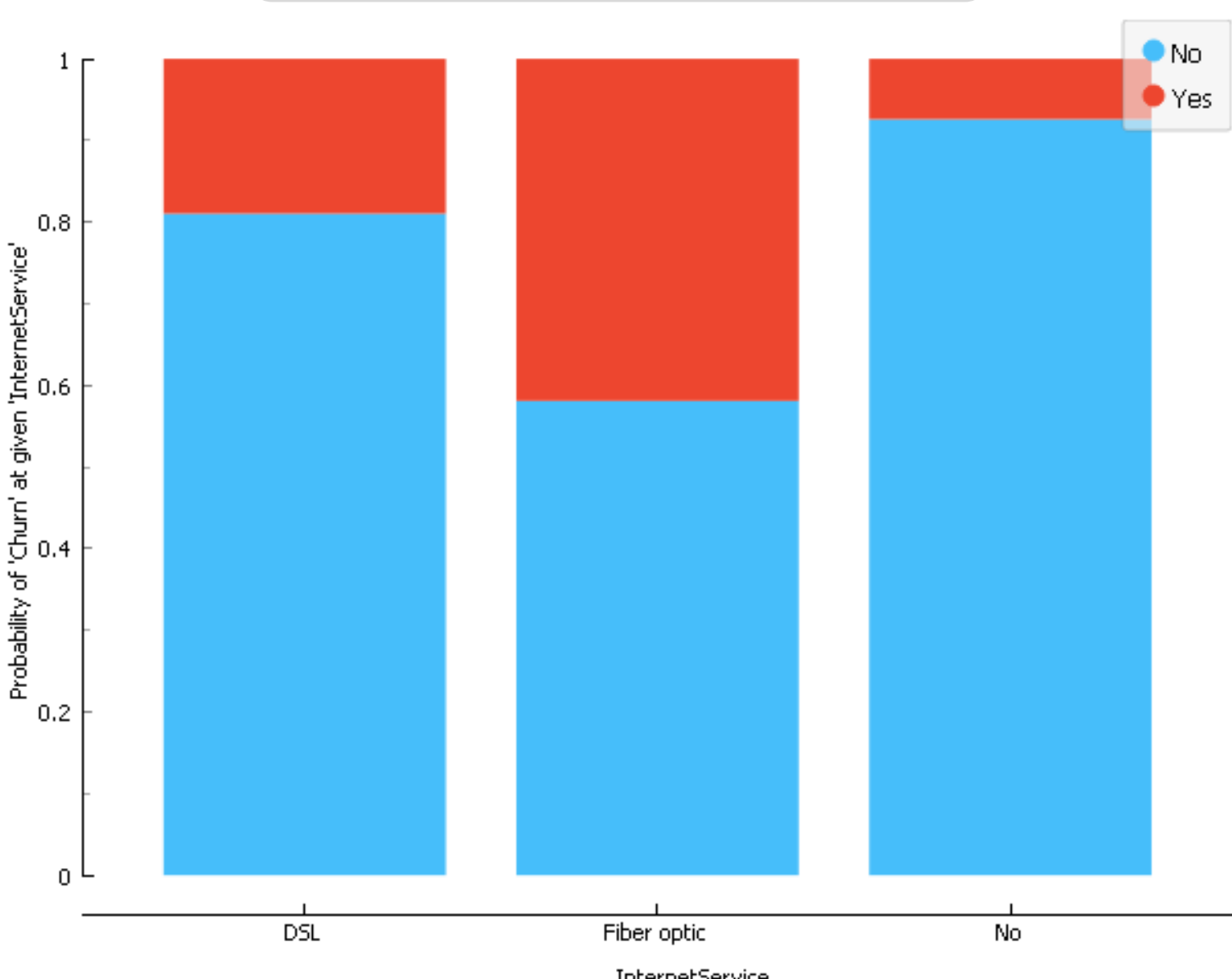
Paperless Bill

Yes (Subscribed)



Internet Service

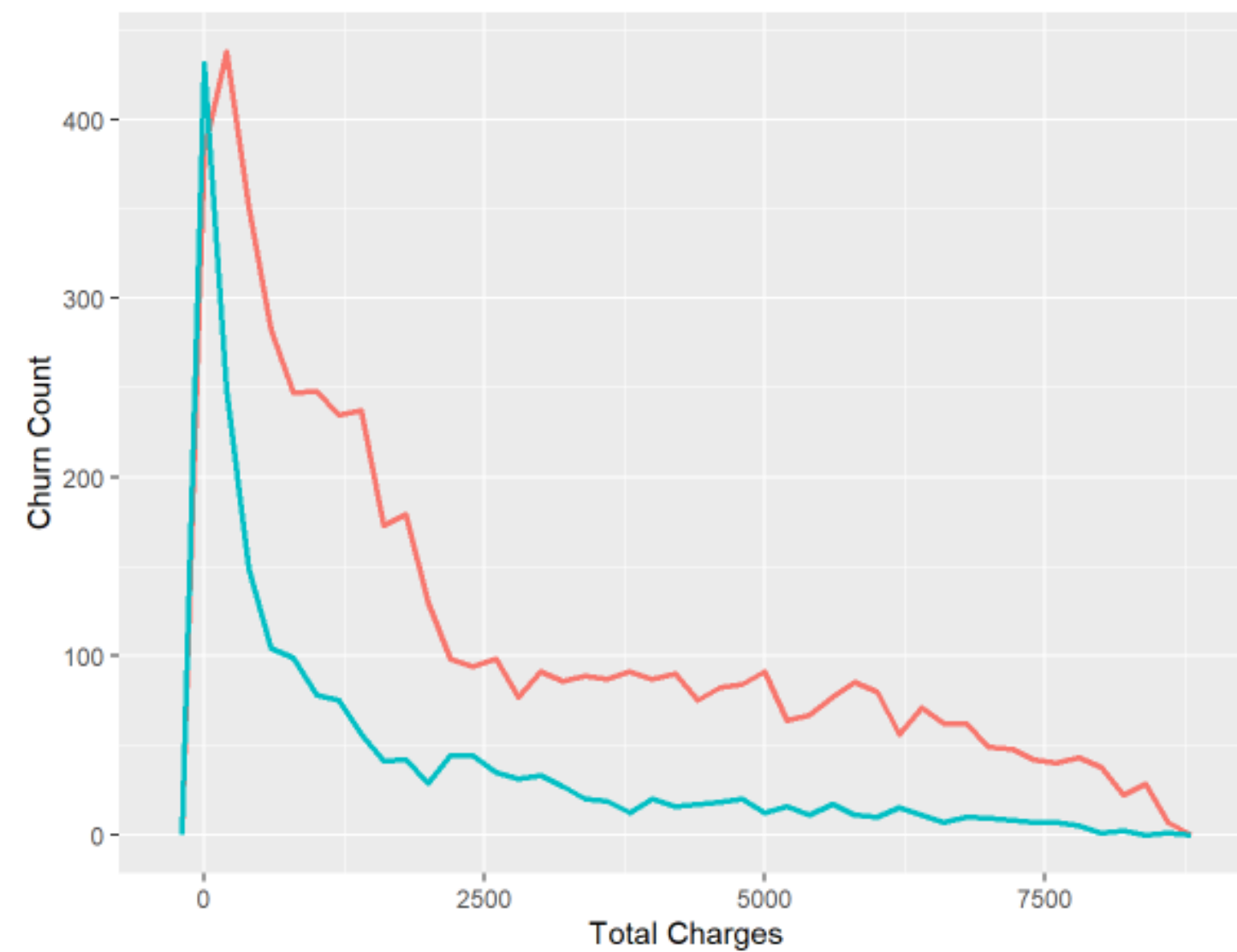
Fiber Optics



Customers Churn Counts

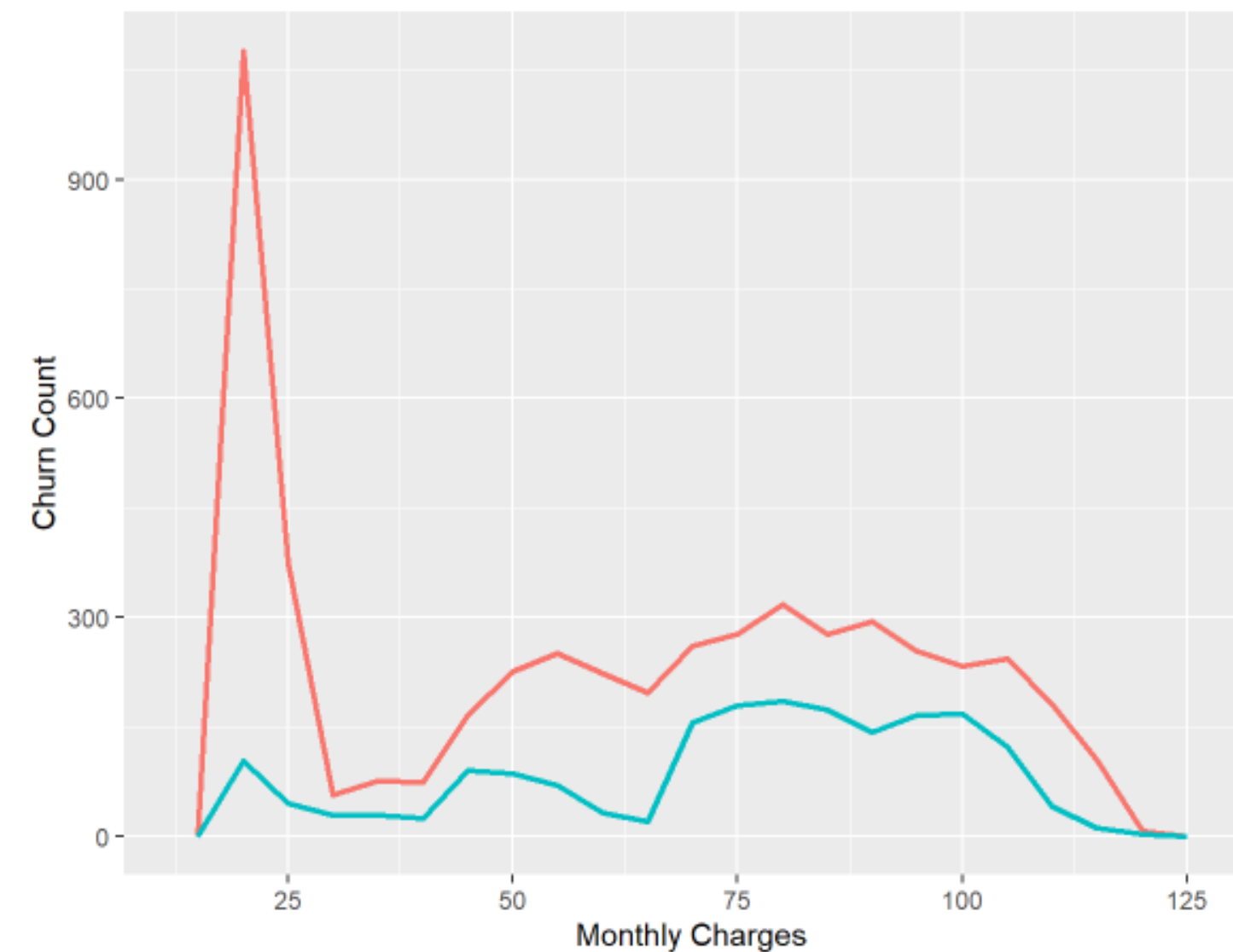
12

Churn Vs Total Charge



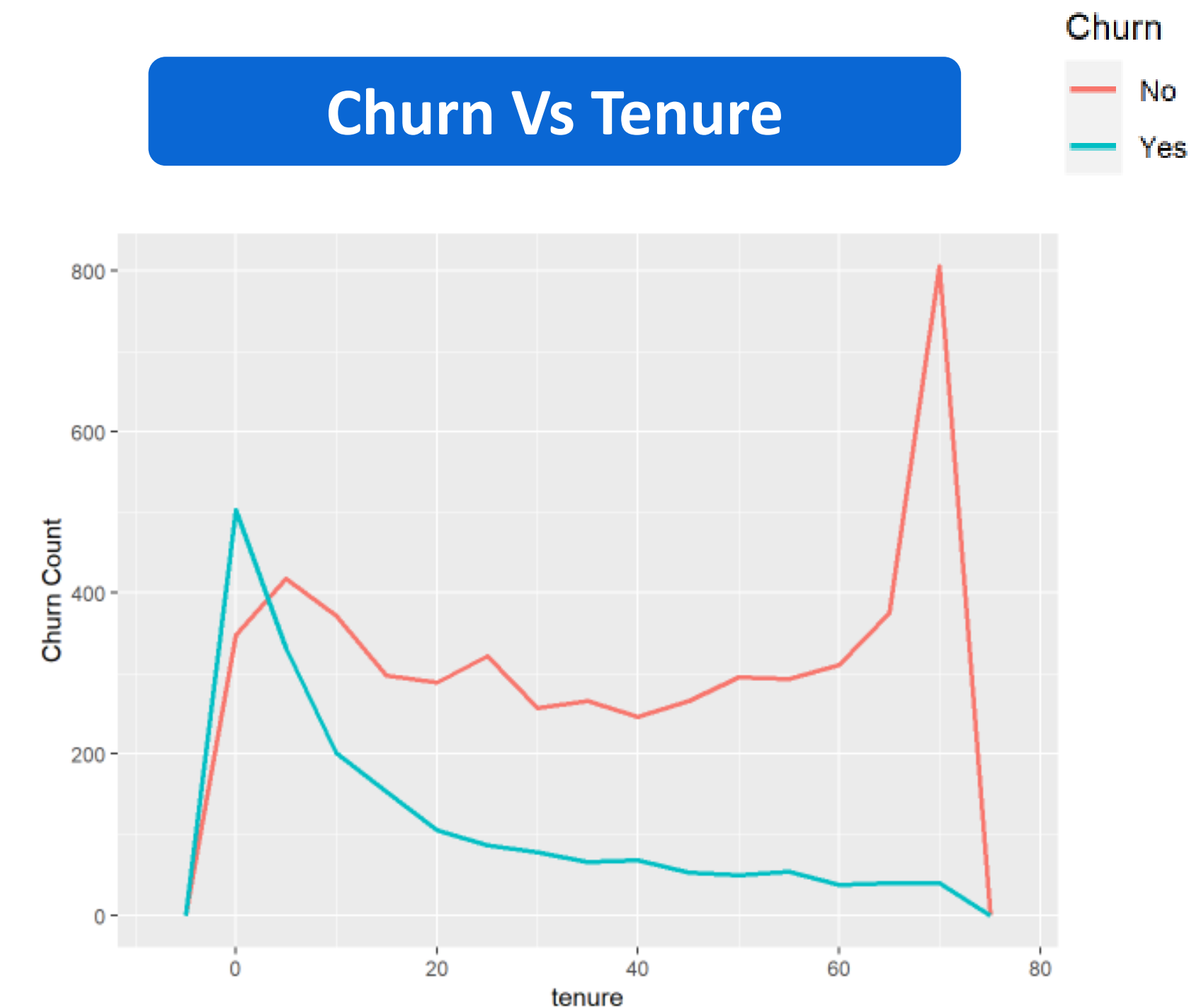
- The total charge two-line charts follow the same pattern for both churned and current customers.
- Although the count of churned customers is low in the dataset, still the count who use below about 250, are same for both categories.

Churn Vs Monthly Charge



- The pattern is different for the customers who have churned from the network.
- The highest frequency of customers belonged to below 25 category.
- More customers who have churned, have used between monthly charges between 75-100.

Churn Vs Tenure

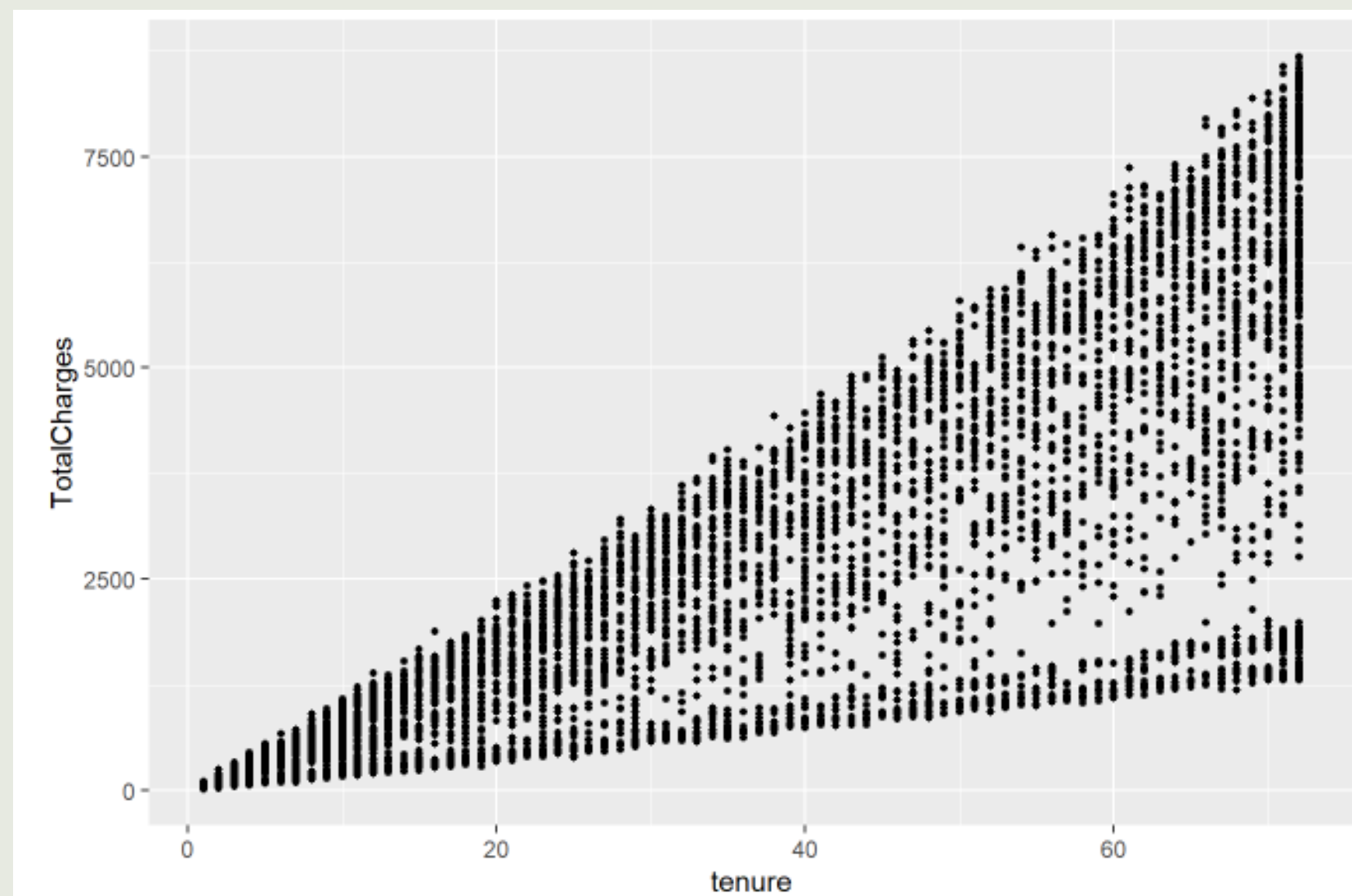


- As the tenure increases count of churn customers are decreasing, while the non churn customers are increasing

Correlation Analysis

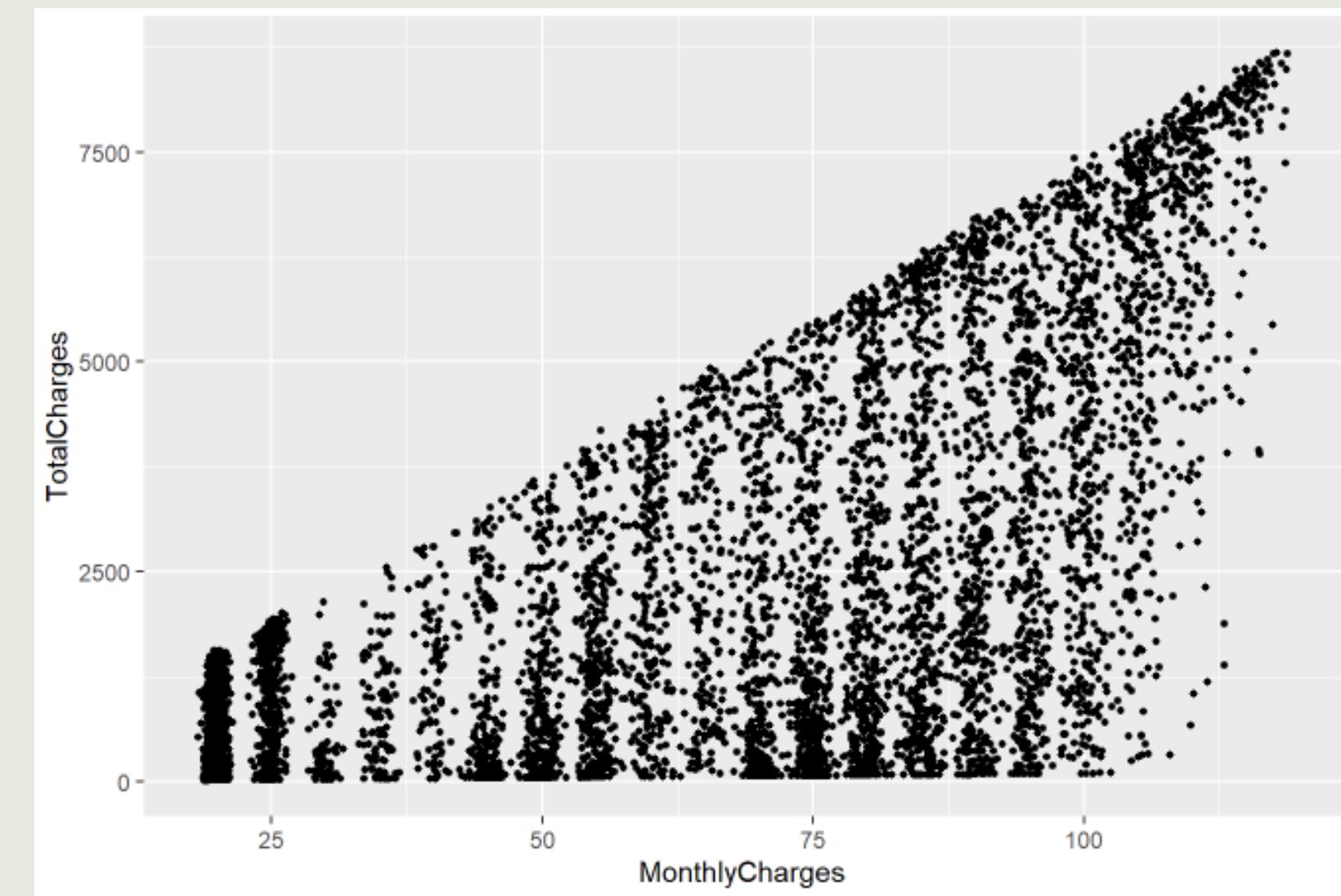
Pearson Correlation Coefficients

Total Charge Vs Tenure



- The correlation coefficient is 0.826.
- There is a linear positive strong relationship between the two factors.

Monthly Charge Vs Tenure



- The correlation coefficient is 0.651.
- There is a linear positive moderate relationship between the two factors.

Correlation Analysis

Overall Correlation Matrix

	<i>tenure</i>	<i>Monthly Charges</i>	<i>Total Charges</i>
tenure	1.000		
Monthly Charges	0.247	1.000	
Total Charges	0.826	0.651	1.000

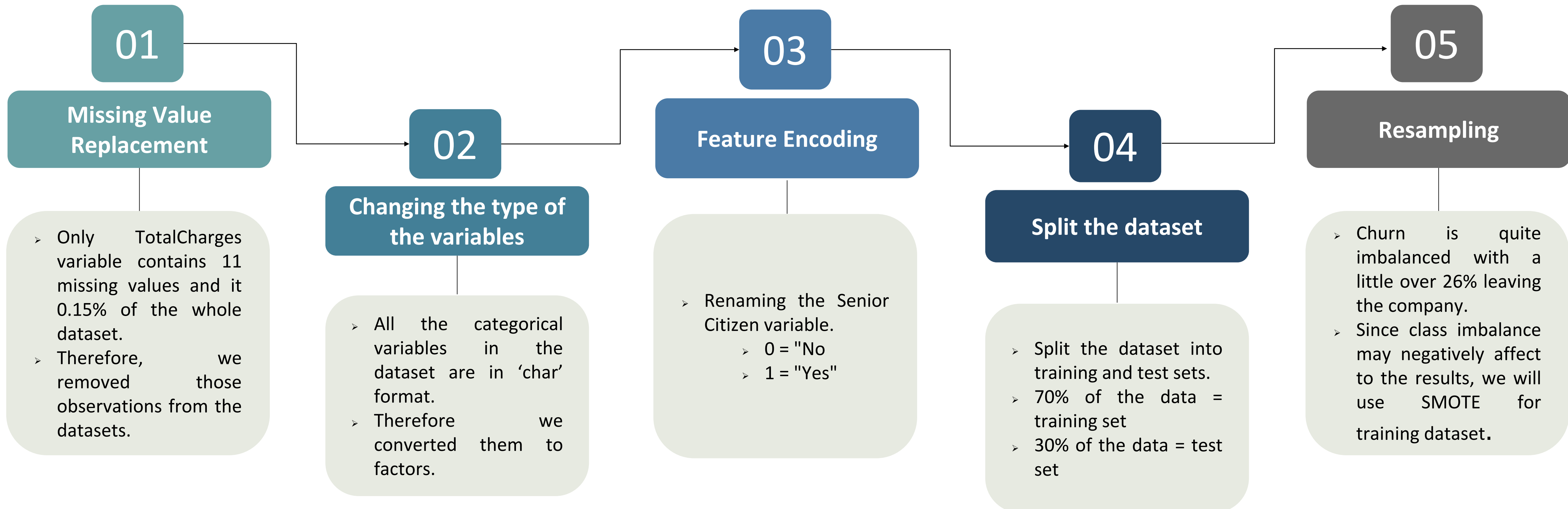
- The strongest positive relationship is observed between tenure and total charge ($r = 0.826$).
- Total charge and Monthly charge have a positive moderate relationship with correlation coefficient of 0.651.
- Tenure variable has positive weak relationships with monthly charge ($r = 0.247$).

Data Preprocessing Workflow

We have identified few data quality issues in the data.

- **TotalCharge variable contains 11 missing values.**
- **All the categorical variables in the dataset are in 'char' format**
- **Senior Citizen variable has labelling issue**
- **Data imbalance issue**

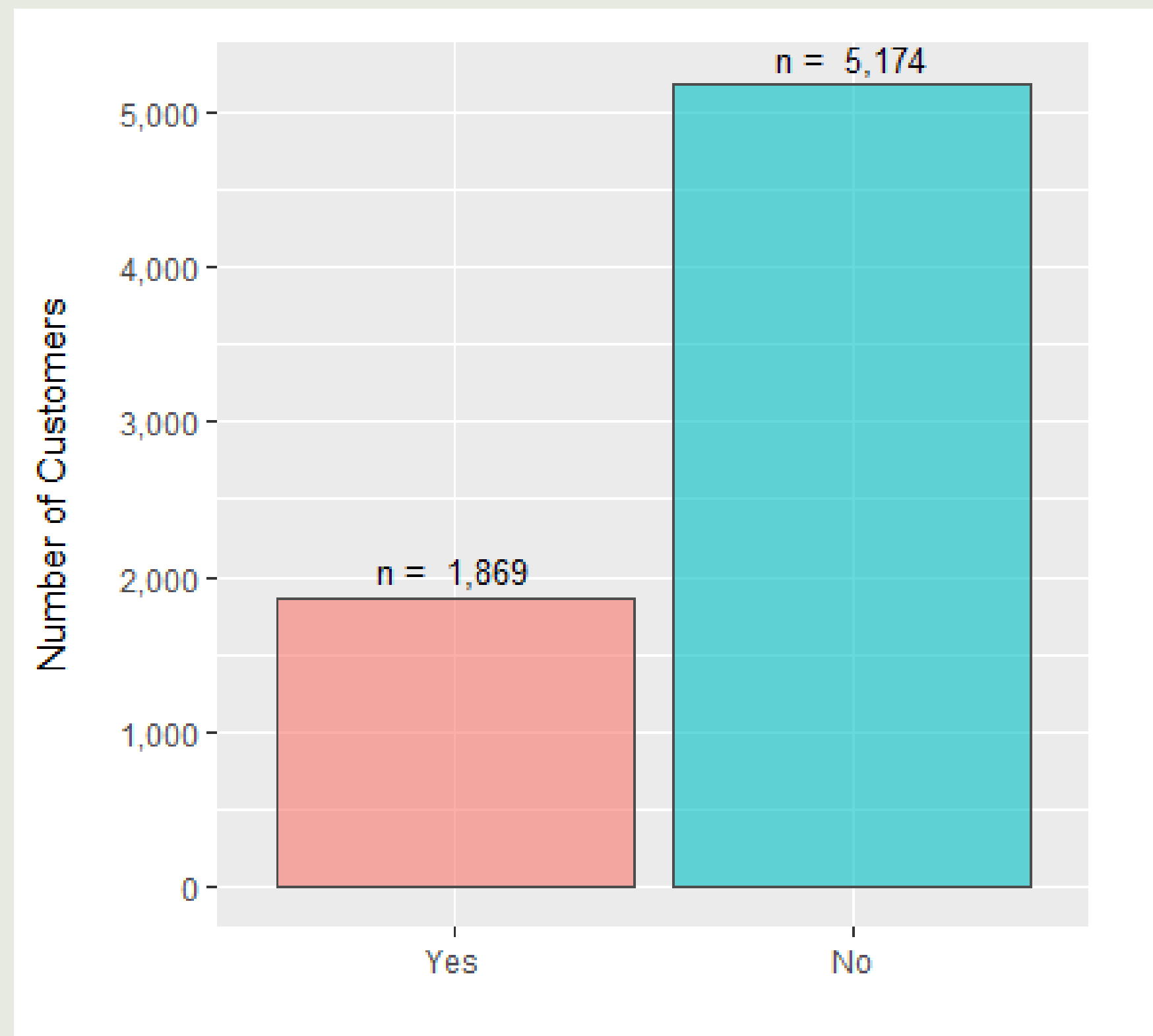
We have applied following methods to overcome these issues,



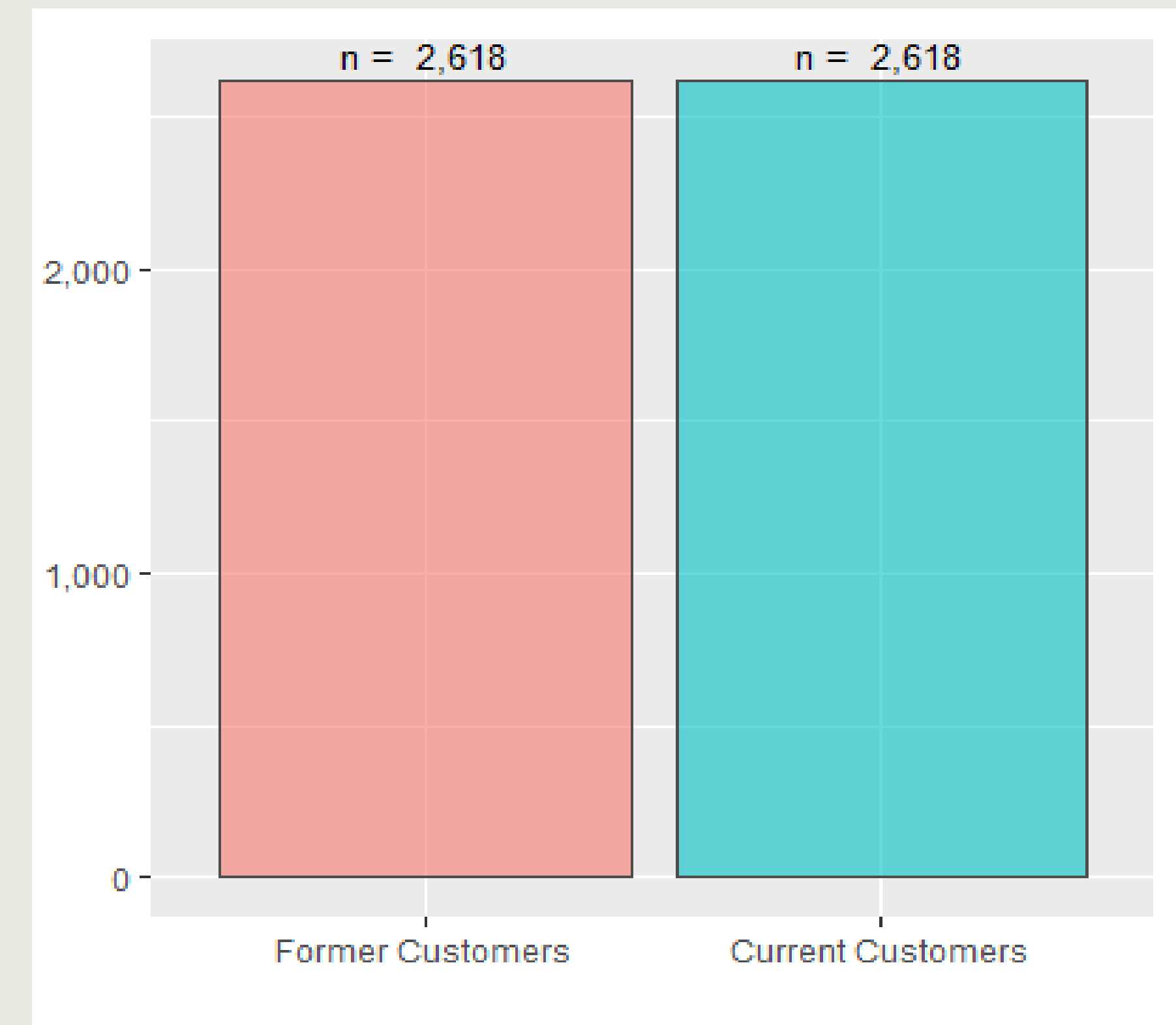
Original Dataset vs Resampled Dataset

16

Original Dataset



After Resampling



Testing if the average monthly charges differ between males and females

H_0 : The average monthly charge does not differ between males and females ($\mu_1 = \mu_2$).
 H_1 : The average monthly charge does not differ between males and females ($\mu_1 \neq \mu_2$).

Normality Check

```
Shapiro-Wilk Normality Test (alpha = 0.05)
```

```
data : df$MonthlyCharges and gender
```

	Level	Statistic	p.value	Normality
1	Female	0.9215085	2.615068e-39	Reject
2	Male	0.9201174	7.152397e-40	Reject

- The shapiro wilk test also shows a p-value less than 0.001 for both males and females, which indicates the significance of the normality test.

**Conclusion: Data are not normally distributed
(non parametric)**

Wilcoxon rank sum test

```
Wilcoxon rank sum test with continuity correction
```

```
data: MonthlyCharges by gender
```

```
W = 6298265, p-value = 0.249
```

```
alternative hypothesis: true location shift is not equal to 0
```

Significance level (=0.05)

P-value = 0.249

As the p-value (=0.249) > 0.05

Fail to reject the null hypothesis.

Conclusion: There is not sufficient evidence to conclude that the average monthly charge for males and females differ.

Testing if the total charge differs by Payment Method

$$H_0: \mu_1 = \mu_2 = \mu_3.$$
$$H_1: \text{At least one mean is different}$$

Normality Check

```
Shapiro-Wilk Normality Test (alpha = 0.05)
-----
data : df$TotalCharges and PaymentMethod

      Level Statistic      p.value Normality
1 Bank transfer (automatic) 0.9228526 2.123728e-27 Reject
2  Credit card (automatic) 0.9175288 5.073903e-28 Reject
3      Electronic check 0.8496637 9.382990e-43 Reject
4      Mailed check 0.7103853 3.209453e-46 Reject
-----
```

. The Shapiro wilk test shows the p-values for all the categories are less than 0.05 and the tests are significant.

Conclusion: The normality assumption is not met for the data.

Testing for homogeneity

```
Levene's Test for Homogeneity of Variance (center = median)
      Df F value      Pr(>F)
group   3  218.93 < 2.2e-16 ***
      7028
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
-----
```

From the Levene’s test it is clear that the variances among the four groups are not similar as the p-value is less than 0.001.
The homogeniety assumprion is also not met for this data.

Conclusion: Instead of ANOVA parametric test, Kruskal Wallis non parametric test will be Used.

Testing if the total charge differs by Payment Method

Applying the Kruskal-Wallis test

```
Kruskal-Wallis rank sum test
```

```
data: df$TotalCharges by df$PaymentMethod
Kruskal-Wallis chi-squared = 1077, df = 3, p-value < 2.2e-16
```

Significance level (=0.001)
P-value = 2.2e-16
As the p-value < 0.001
Reject the null hypothesis.

Conclusion: at least one mean total charge is different from others..

Posthoc Analysis

```
Pairwise comparisons using Wilcoxon rank sum test with continuity correction
```

```
data: df$TotalCharges and df$PaymentMethod
```

	Bank transfer (automatic)	Credit card (automatic)
Credit card (automatic)	0.7	-
Electronic check	<2e-16	<2e-16
Mailed check	<2e-16	<2e-16
	Electronic check	
Credit card (automatic)	-	
Electronic check	-	
Mailed check	<2e-16	

```
P value adjustment method: BH
```

Conclusion: The post hoc analysis shows that except between bank transfers and credit card payments, all the mean total charges are different for all categories.

Testing if the payment method and the contract type are independent

H_0 : payment method and the contract type are independent.
 H_1 : payment method and the contract type are not independent.

Applying Chi-Square Test

```
Pearson's Chi-squared test  
  
data: obs  
X-squared = 1001.6, df = 6, p-value < 2.2e-16
```

- The p-value obtained from the chi squared test for independence is less than 0.001.
- Therefore the null hypothesis is rejected.

Assumption Checking

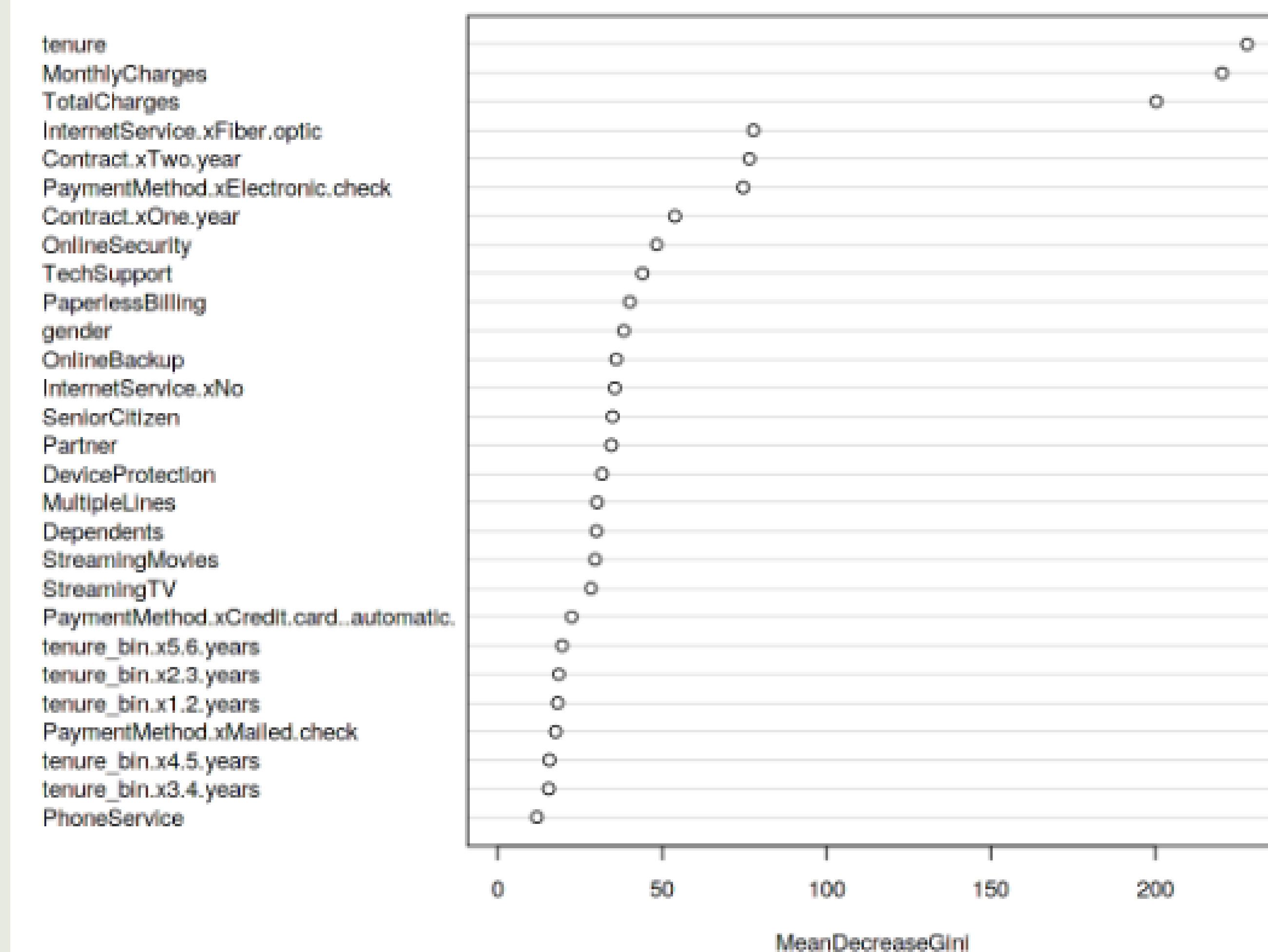
	Month-to-month	One year	Two year
Bank transfer (automatic)	849.4960	322.9181	371.5860
Credit card (automatic)	837.3917	318.3169	366.2914
Electronic check	1301.2033	494.6252	569.1715
Mailed check	886.9090	337.1399	387.9512

The expected values are greater than 5 for each cell. Therefore, the expected value at least be 5, assumption is met for the data.

Conclusion: There is sufficient evidence to conclude that the payment method and the contract type are independent.

Feature Selection

Feature Importance

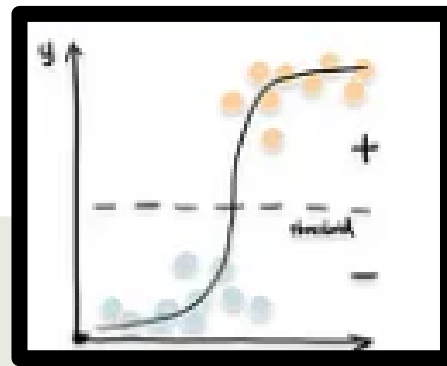


- Most important variables are Tenure, MonthlyCharge and TotalCharge.
- Due to significant correlation (multicollinearity) between tenure and totalcharge, only selected tenure for the analysis.
- All the features used for the analysis except **totalCharge** variable.

Model Training

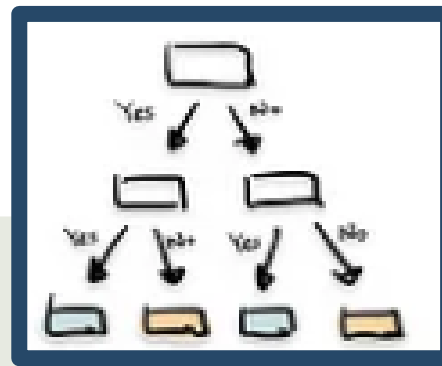
The binary classification technique was used for the resampled dataset.

Following models are applied to the resampled dataset.



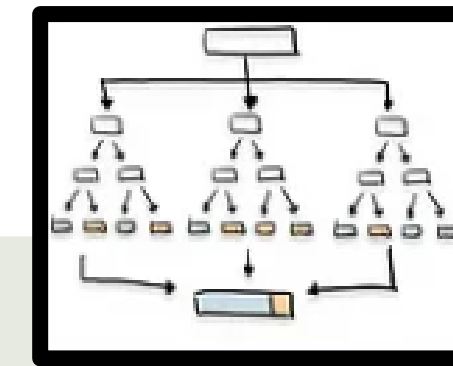
Logistic regression is a process of modeling the probability of a discrete outcome given an input variable.

Logistic Regression



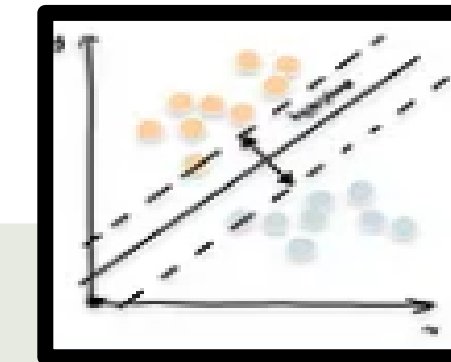
A decision tree is a graph that uses a branching method to illustrate every possible output for a specific input.

Decision Tree



The random forest is a supervised learning algorithm that randomly creates and merges multiple decision trees into one "forest."

Random Forest



Support vector machines (SVMs) are supervised learning models that analyze data and recognize patterns, used for classification and regression analysis

Support Vector Machines (SVM)

Model Evaluation

- Performances of each model are mainly measured using **Accuracy, Precision, Recall and F1-score**.
- But the performance can not directly be assessed by using the accuracy measure as the imbalances of the dataset.
- Hence, the F1-score, recall and precision are the most appropriate measurement to evaluate the model.
- Among these measures, high priority is given to **the recall and f1 score**. Therefore the process is recall-oriented.

	Accuracy	Precision	Recall	F1-score
Logistic Regression	0.7663	0.5358	0.8143	0.6464
Decision Tree	0.7381	0.8375	0.5043	0.6295
Random Forest	0.7533	0.5238	0.7857	0.6286
Support Vector Machines (SVM)	0.6997	0.4644	0.8500	0.6006



- Logistic regression has the best performance on the test set and below summary shows, the coefficients and the performances on the test set.

Logistic Regression Model

$\widehat{churn} = -0.13573 - 0.78385 \text{ tenure} + 0.75276 \text{ MonthlyCharges} + 0.19104 \text{ InternetserviceFiberoptic} + 0.03143 \text{ InternetServiceno} - 0.06496 \text{ PaymentMethodcreditcardauto} + 0.21314 \text{ PaymentMethodecheck} - 0.03878 \text{ PaymentMethodMailedcheck} - 0.31086 \text{ ContractoneYear} - 0.56591 \text{ ContractTwoYear} - 0.16972 \text{ onlinesecurityves} - 0.21416 \text{ TechsupportYes} + 0.14521 \text{ Paperlessillingres} - 0.04852 \text{ genderMale} - 0.26178 \text{ PhoneServiceYes}$

How will you gauge if this model is better able to predict customer churn ?

Model Performance

Accuracy	0.7633	Recall	0.8143
Precision	0.5358	F1 Score	0.6464

Confusion matrix

Reference			
Prediction		No	Yes
	No	1148	120
	Yes	400	440

Logistic Regression Model

$\widehat{churn} = -0.13573 - 0.78385 \text{ tenure} + 0.75276 \text{ MonthlyCharges} + 0.19104 \text{ InternetserviceFiberoptic} + 0.03143 \text{ InternetServiceno} - 0.06496 \text{ PaymentMethodcreditcardauto} + 0.21314 \text{ PaymentMethodecheck} - 0.03878 \text{ PaymentMethodMailedcheck} - 0.31086 \text{ ContractoneYear} - 0.56591 \text{ ContractTwoYear} - 0.16972 \text{ onlinesecurityves} - 0.21416 \text{ TechsupportYes} + 0.14521 \text{ Paperlessbillingres} - 0.04852 \text{ genderMale} - 0.26178 \text{ PhoneServiceYes}$

Interpret the coefficients

➤ -0.78385 tenure

Every one month increase in tenure, expected churn decreases by 0.78385

➤ +0.75276 MonthlyCharge

Every one dollar increase in monthly charge, expected churn increase by 0.75276

Importance of the Best model

26

The model can be used to:

Identify customers who are likely to churn

Can develop this model to identify the high-risk customer

Segmenting the customers based on the churn propensity (high risk/medium risk/low risk)



Identify the most important features of customer churn

Based on the results company can offer special promotions campaigns to high risk/ identified potential churn customers

Can include this predictions base as an input variable for other ongoing machine learning projects

Limitations of our Best model

May not be able to capture complex non-linear relationships between variables

Assumes that the relationships between the independent and dependent variables are always linear

May not be accurate for a huge amount of data

Assume variables followed normal distribution and independent

The model can be prone to overfitting

Our model may not be well-suited for predicting customer churn in cases where the data is time-series in nature



The limitations of the dataset:

- ▶ No of observations in the dataset is relatively small.
- ▶ Lack of variables in the dataset
- ▶ Do not have time series factors/features in the dataset.
- ▶ Imbalanced data for some classes
- ▶ Outdated data



Suggestions on how the dataset can be improved:

- ▶ Increasing the sample size by collecting more data
- ▶ Adding more features to the dataset like customer's network usages, competitor's details, etc.)
- ▶ Adding time series data to the dataset
- ▶ Balancing the class distribution
- ▶ Regularly collecting and updating the data



Conclusion

30

Dataset consists 7043 observations and among those observations, most of them are male, that are not senior citizens and do not have a dependent or partner.

Telco company should target the customer who is a senior citizen, does not have a partner or dependent, use fiber optics and payment done by electronic check

Tenure, MonthlyCharge, Contract, PaymentMethod and **OnlineSecurity** are the top 5 most important features of the analysis,

Dataset is imbalanced for some classes and there is a strong positive relationship between **tenure** and **TotalCharge** variables (multicollinearity occurred)

Logistic regression model gave the best performances among the other algorithms



