Telco Customer Project

Andrew Lee, Hisham Galib

University of Texas at Austin

Table of Contents
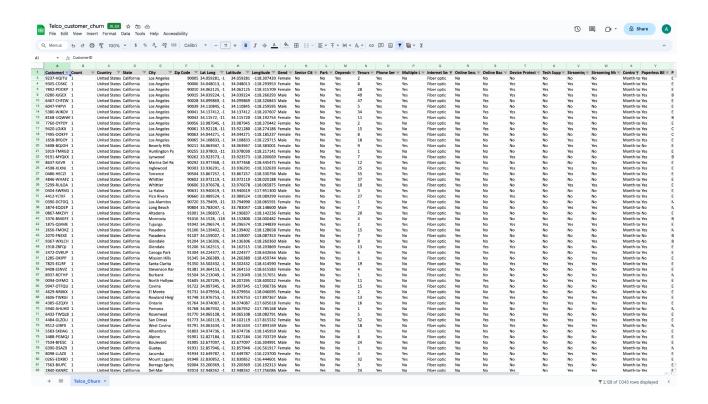
**Introduction**

With growing competition from big companies such as Verizon and AT&T, telecom providers

constantly compete with one another to attract more customers, often taking customers from the

other provider with incentives promised to the customer after switching providers. The situation

described is related to what is referred to as 'churn' in the telecommunication industry. Churn is

the measure of how many customers stop using a product or service. Customer retention has

been an imperative focus from telecom providers to sustain revenue, attempt to foster lasting

relationships, and ensure customer satisfaction. With competition growing constantly and

providers finding new ways to attract customers as well as retain customers, one of the best ways

to go about finding methods to accomplish mitigating churn and potentially finding new

customers is to analyze patterns of churn, uncover insights, and proposing strategies based off

those patterns and insights. Demonstrating these techniques in the Telco Customer Churn dataset,

a synthetic dataset produced by IBM to model a hypothetical telecommunications company. With

the dataset, analyzing the data to find out if the customer is at risk of churn or is undergoing

churn, and finding out how to retain customers through a proposed solution based off the

analysis is the goal. Using key features such as 'Payment Method', 'Monthly Charges', and

'Internet Service' as the leading factors in increasing churn, techniques such as an Exploratory

Data Analysis (EDA) and Synthetic Minority Oversampling Technique (SMOTE) were used to

decide how much these features affected churn and how they can be harnessed to reduce churn as

well as to retain customers.

**Description of Data**

The Telco Customer Churn dataset is a repository of information, providing data on the

behaviors and preferences of customers. The dataset includes key attributes that define

customers, ranging from demographic details to the type of service they subscribe to including

the main variable 'Churn' which indicates whether the customer has chosen to discontinue or

continue their services with Telco. The dataset can be found here.



Variable Explanation

The features we used to determine the main contribution of churn were 'Payment Method' which

would indicate customer's preferred payment method, 'Monthly Charges' which was the cost

factor and how higher charges might lead to higher churn, and 'Internet Service' which indicates

the type of service used which is essential for understanding customer needs. We used them as

the categorical feature and plotted them to understand the variables' relationship to churn.

```
# Explore the distribution of selected categorical variables
categorical_features = ['Payment Method', 'Internet Service']
for col in categorical_features:
    plot_categorical_distribution(df, col, 'Churn')
```
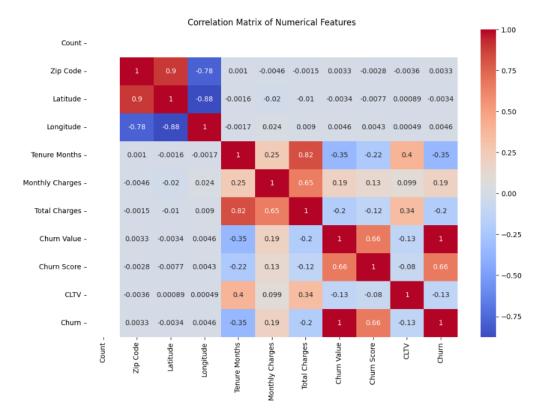
1) Payment method - Payment method can show the preferred payment method chosen by each customer and provide important insight and their transaction behavior and preferences. Different types of payment methods as well as shifts in payment methods might include factors such as financial strain or dissatisfaction which would ultimately lead to churn.

2) Monthly Charges - Monthly charges essentially show how much customers are paying in a month-to-month basis. Higher monthly charges may lead to increased churn as customers who are cost-sensitive may look for more affordable alternatives. Being able to analyze the correlation between monthly charges and churn provides insight to how pricing and pricing structures impact customer retention. Through the analysis, being able to cater and optimize prices at the threshold where monthly charges significantly affect churn is crucial in meeting customer needs and expectations.

3) Internet Service - Internet services is a key determinant on a customer's specific needs and usage patterns. Being able to understand why a customer would opt into different plans such as basic or premium internet services to meet their needs is crucial as it can help optimize custom plans to meet customer needs. Catering to the varied requirements of their customer base in incredibly important in the world of customer retention.
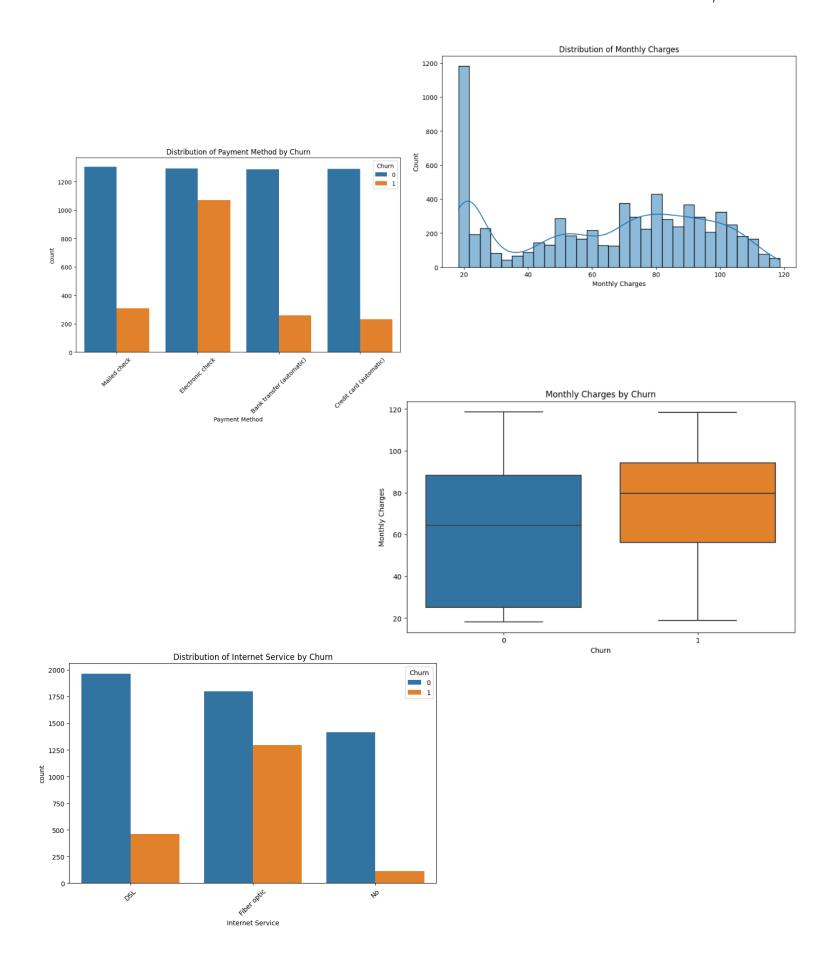
| Feature | Payment Method | Monthly Charges | Internet Service |
|---|---|---|---|
| Data | Object | Float64 | Object |
| Rationale | Indicates customer's preferred payment mode, could affect churn | Cost factor, higher charges might lead to higher churn | Type of service used, essential for understanding customer needs |

## **Methodology**

1) Exploratory Data Analysis (EDA)

After opening up the dataset from kaggle.com using the pandas library, the libraries seaborn and matploblib.pyplot were used to make visualizations in the form of histograms, bar graphs, and box plots.



Correlation Matrix of Numerical Features

### Distribution of Monthly Charges



### Distribution of Payment Method by Churn



### Monthly Charges by Churn



### Distribution of Internet Service by Churn

2) Resampling through SMOTE

After our exploratory data analysis, we noticed an imbalance in our dataset. SMOTE was used to create synthetic samples due to its ability to address the issue of class imbalance in our dataset.

**EDA Results**

Distribution of Payment Method by Churn

- Churn distribution demonstrate a larger number of customers not churning (0) to those who do (1).

- Majority of customers prefer an Electronic Check type of payment over other payment methods.

- Many Electronic Check users show a higher churn rate as well, suggesting that there is some sort of dissatisfaction or a need for improved payment processing.

Internet Service

- Fiber Optic users have the highest churn rate of all the other types of services, indicating some dissatisfaction with either the service provided or the price of the service itself.

- Customers with DSL service was shown to lower churn rates, possibly having something to do with perceived value or service satisfaction.

Monthly Charges

- Higher Monthly Charges correlate with increased churn, highlighting price sensitivity among customers.

**SMOTE**

In our analysis, we also used the Synthetic Minority Over-sampling Technique (SMOTE) to address the challenge of class imbalance, particularly in the 'Churn' classification. This technique, introduced by Chawla et al. in their 2002 paper, is pivotal for enhancing the performance of machine learning models in scenarios where the minority class is underrepresented (Chawla et al. 321-357). In our project, SMOTE was applied to the training dataset after the initial phase of preprocessing and feature engineering. By generating synthetic samples of the minority class, SMOTE significantly balanced the class distribution within the dataset. This balancing act was crucial for training the Random Forest Classifier. It ensured that the model did not show bias towards the majority class and improved its ability to generalize well to new unseen data. The effectiveness of this approach was evaluated on an untouched test set, where the model's performance metrics, including the accuracy and classification report, indicated a more robust and equitable handling of both classes.

**Machine Learning Results**

After the exploratory data analysis, the Random Forest Classifier was chosen for its ability to handle unbalanced datasets. The customer churn dataset suffered from class imbalances where other instances of one class significantly tends to outweigh the others. The classifier was trained on a balanced dataset achieved through the application of the Synthetic Minority Over-sampling Technique (SMOTE), ensuring that the minority class was adequately represented.

```
# Model Training: Random Forest Classifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
```

**Evaluation Metrics**

We trained our Random Forest model with hyperparameters to balance bias and variance, thereby optimizing the model's ability to generalize to unseen data. The model's hyperparameters, such as the number of trees in the forest 'n_estimators' and the maximum depth of the trees 'max_depth', were used to manage the complexity of the model and prevent overfitting.

After training, the model was evaluated using a test set that reflected the original distribution of the classes. Using this approach allowed us to assess the model's performance in a more realistic scenario, and ensured the validity and applicability of our findings.

Our evaluation metrics were precision, recall, and the F1-score, each providing insights into different aspects of the model's performance.
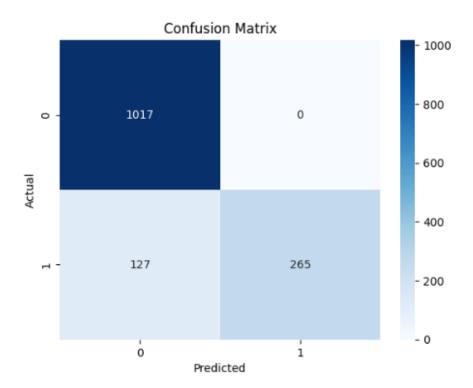
**<u>Model Results</u>**

1) Precision, Recall, F1-Score:

- Precision: For class 0 (no churn), the precision is 89%, indicates the model's accuracy in predicting customers who won't churn.

- The model excels in recalling class 0 with a perfect score (100%), showing its effectiveness in creating a customer base who won't churn. While there is a perfect score for recalling class 0, there is a lower recall of 68% for class 1, suggesting different opportunities for identifying customers who do churn.

- F1-Score: The F1-score, a balance between precision and recall, stands at 94% for class 0 and 81% for class 1, showcasing a good overall performance.

```
Random Forest Accuracy: 0.9098651525904897

Random Forest Classification Report:
              precision    recall  f1-score   support

           0       0.89      1.00      0.94      1017
           1       1.00      0.68      0.81       392

    accuracy                           0.91      1409
   macro avg       0.94      0.84      0.87      1409
weighted avg       0.92      0.91      0.90      1409
```

2) Confusion Matrix

- A confusion matrix represents the prediction summary in matrix form and is an indicator for how many predictions are correct and incorrect per class.



Confusion Matrix

- True Positives (TP): 265 instances where the model correctly predicted churn.

- True Negatives (TN): 1017 instances where the model correctly predicted no churn.

- False Positives (FP): 0 instances where the model incorrectly predicted churn.

- False Negatives (FN): 127 instances where the model incorrectly predicted no churn.

This matrix was critical in visualizing the model's predictive power and diagnosing areas where the model could be improved. In particular, it was useful in reducing false negatives to capture more customers at risk of churning.

```
Actual Churn Reduction: 9.013484740951032%
Time Frame for Impact: 12 months
```

```
Actual Estimated Impact:
We anticipate achieving a significant reduction in churn rates by approximately 9.013484740951032% within the first 12 months.
```

**Model Performance, Strengths and Areas for Improvement**

Overall accuracy of the model stood at 90.99%, which is not only impressive, but it demonstrates the classifier's strong predictive capabilities. While accuracy can be used as a valuable metric, the additional context provided by precision, recall, and the F1-score is also essential for having a comprehensive understanding of the model's performance, particularly in an imbalanced class distribution scenario like churn prediction.

Some of the strengths of the model are demonstrated by its high precision in identifying non-churning customers and achieving perfect recall for this group. A high F1-score for class 0 indicated a well-balanced model for the majority class. However, the model displayed a relatively lower recall for class 1, which showed an area that needed to be improved. We would guess that increasing the recall for class 1 would mean enhancing the model's sensitivity to churn risk, thereby providing an opportunity to intervene and retain those customers.

In conclusion, implementing the Random Forest Classifier has worked in our favor. It has shown impressive capabilities in predicting customer churn. Its most notable strengths are its accuracy and its effectiveness in pinpointing true negative cases. However, there is room for improvement in enhancing its sensitivity towards instances of churn. These findings provide a foundation for subsequent model training phases, but they also provide insight. The groundwork for future iterations of model training should have a focus on improving recall for the churn class without compromising the overall accuracy and precision.

**<u>Proposed Solution</u>**

The insights gleaned from our machine learning model make it easy to reduce churn. We have proposed a solution that focuses on leveraging predictive analytics to tailor customer experiences and address the factors contributing to churn.

In terms of our strategic rationale, our Random Forest Classifier has provided us with a clear indication of which customers are at risk of churn. With this predictive power, we can adopt a proactive stance in customer relationship management. The strategic rationale for our solution is described by these points: customer segmentation, personalized engagement plans, and service enhancement initiatives.

For customer segmentation, we propose segmenting the customer base into distinct groups based on their risk of churn, service usage patterns, and billing information. This requires a good understanding of the unique preferences and behaviors of different customer segments which enables the creation of targeted marketing campaigns tailored to specific needs and concerns. This will ensure that interventions are timely and relevant to each customer's unique circumstances.

For our personalized engagement plans, we will be leveraging the segmentation to implement

them. Customers identified with a higher risk of churn can be targeted with specialized campaigns, which may include personalized discount offers, service upgrades, or loyalty rewards. For those with high monthly charges or dissatisfaction with payment methods, we can offer customized billing solutions or payment plans to alleviate financial concerns.

For service enhancement initiatives, our data points to specific service types, such as Fiber Optic internet, correlating with higher churn rates. Addressing this, we propose service enhancement initiatives where we can refine the quality of service, adjust pricing strategies, or introduce new features that align with customer expectations and needs.

**Estimated Impacts**

The implementation of our solution is expected to achieve a significant reduction in churn rates, increase customer lifetime value, and enhance overall customer satisfaction. By aligning our efforts with predictive insights, we aim to foster a more loyal customer base, reduce acquisition costs, and bolster the company's market position.

**Conclusion**

The Telco Customer Churn project was quite the journey of data exploration, insightful analysis, and predictive modeling to confront the perennial challenge of customer churn. The synthesis of exploratory data analysis (EDA) and advanced machine learning techniques helped lead us to use a Random Forest Classifier, which would help identify customers at risk of discontinuing their services.

Our accuracy of 90.99% was impressive as it showed its effectiveness in distinguishing between customers who will churn and those who will remain. However, the machine learning model's recall for predicting churn (class 1) showed room for improvement. It is within this context that we made proposed solution to use targeted marketing campaigns and personalized

customer engagement strategies to reduce churn. Our model's predictive strength informed us to use strategic customer segmentation, enabling personalized engagement plans and service enhancement initiatives tailored to individual customer profiles. We expect that this would shift the churn paradigm significantly.

The project's ending is not merely just the use of statistics but it is a testament to the power of data-driven decision-making. The anticipated churn reduction of 9.01% within the first 12 months is a powerful insight, as it can lead to new customer loyalty and satisfaction. As we advance, the project serves as a foundation for future innovations in predictive analytics and customer relationship management, showcasing a new era where customer retention can be redefined through data.

**Appendix: Q & A Session**

Following the presentation on our data science project, a brief Q&A session was held. This appendix documents the key question raised and the response provided.

**Question 1**: "Why did you use SMOTE in your analysis, given it was not covered in our course?"

**Answer**:

We decided to use SMOTE (Synthetic Minority Over-sampling Technique) in our analysis because we needed to address the issue of class imbalance in our dataset. We chose SMOTE for its effectiveness in generating synthetic samples for the minority class, thereby creating a more balanced dataset and improving the performance of our predictive models. Moreover, this approach ensures more accurate and equitable predictions across all classes, reflecting our commitment to addressing dataset biases and ensuring the reliability and fairness of our models.

**Question 2:** "Are you going to make a key feature table with the variables used?"

**Answer**: We said we would. Initially, our report and presentation did not include a key feature table. However, recognizing its importance for clarity and understanding, we committed to incorporating it in our final presentation and report. The key feature table we subsequently added provided a clear overview of the variables used, their descriptions, data types, and roles in our analysis, thereby enhancing the transparency and comprehensibility of our project.

**Question 3**: "You had too much text and not enough graphics in your presentation, are you going to use more graphics?"

**Answer**: We said we would. This feedback was crucial as we realized we were lacking in this regard of communication in data science. The necessary balance between text and visual elements. We significantly increased the use of visual elements, including charts, graphs, and infographics, to better illustrate our findings and methodology. These graphical enhancements not only made the presentation more engaging but also facilitated a clearer and more intuitive understanding of our analysis. The additional visuals provided a more dynamic and accessible way for the audience to grasp complex data and concepts, ensuring that our presentation was not only informative but also visually compelling.

**<u>References</u>**

Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research, 16, 321-357.

**<u>Repositories and other links</u>**

https://github.com/leea36/I310DProjectCode.git

Open up the TelcoProject.ipynb and run the entire code from the beginning. Make sure to have the dataset in your drive to run the dataset. The dataset is here for convenience.

🟨 Telco Customer Project