

Comprehensive Data Analysis Report on Telco Customer Churn

1. Brief Description of the Data Set and Summary of Its Attributes

Dataset Name: Telco Customer Churn

Source: Kaggle - Telco Customer Churn Dataset

Description: This dataset provides information on customers of a telecommunications company, detailing their demographics, account specifics, and whether they have churned. It is instrumental in understanding customer behaviour and churn patterns, essential for developing strategies to enhance customer retention.

Data Summary:

- Number of Records: 7,043
- Number of Features: 21

Key Attributes:

1. customerID: Unique identifier for each customer (string)
 - Description: A distinct alphanumeric string assigned to each customer to differentiate them from others.
2. gender: Gender of the customer (categorical: Male/Female)
 - Description: Represents the gender of the customer, with possible values of Male or Female.
3. SeniorCitizen: Whether the customer is a senior citizen (binary: 0/1)
 - Description: Indicates if the customer is a senior citizen, with 1 for Yes and 0 for No.
4. Partner: Whether the customer has a partner (binary: Yes/No)
 - Description: Shows if the customer has a partner, with Yes or No as possible values.
5. Dependents: Whether the customer has dependents (binary: Yes/No)
 - Description: Indicates if the customer has dependents, with Yes or No as possible values.
6. tenure: Number of months the customer has been with the company (numerical)
 - Description: Represents the duration (in months) the customer has been subscribed.
7. PhoneService: Whether the customer has phone service (binary: Yes/No)
 - Description: Indicates if the customer has phone service, with Yes or No as possible values.
8. MultipleLines: Whether the customer has multiple lines (categorical: No phone service/No/Yes)

- Description: Shows if the customer has multiple phone lines, with No phone service, No, or Yes as possible values.
9. InternetService: Type of internet service (categorical: DSL/Fiber optic/No)
- Description: Indicates the type of internet service the customer subscribes to, with DSL, Fiber optic, or No as possible values.
10. OnlineSecurity: Whether the customer has online security (binary: No/Yes)
- Description: Indicates if the customer has online security, with Yes or No as possible values.
11. OnlineBackup: Whether the customer has online backup (binary: No/Yes)
- Description: Shows if the customer has online backup, with Yes or No as possible values.
12. DeviceProtection: Whether the customer has device protection (binary: No/Yes)
- Description: Indicates if the customer has device protection, with Yes or No as possible values.
13. TechSupport: Whether the customer has tech support (binary: No/Yes)
- Description: Shows if the customer has tech support, with Yes or No as possible values.
14. StreamingTV: Whether the customer has streaming TV (binary: No/Yes)
- Description: Indicates if the customer has streaming TV, with Yes or No as possible values.
15. StreamingMovies: Whether the customer has streaming movies (binary: No/Yes)
- Description: Indicates if the customer has streaming movies, with Yes or No as possible values.
16. Contract: Type of contract (categorical: Month-to-month/One year/Two year)
- Description: Indicates the type of contract the customer has, with Month-to-month, One year, or Two year as possible values.
17. PaperlessBilling: Whether the customer has paperless billing (binary: No/Yes)
- Description: Shows if the customer has opted for paperless billing, with Yes or No as possible values.
18. PaymentMethod: Payment method (categorical: Electronic check/Mailed check/Bank transfer (automatic)/Credit card (automatic))
- Description: Represents how the customer pays their bill, with Electronic check, Mailed check, Bank transfer (automatic), or Credit card (automatic) as possible values.
19. MonthlyCharges: Monthly charges (numerical)
- Description: Represents the amount billed to the customer each month.
20. TotalCharges: Total charges (numerical)

- Description: Represents the total amount billed to the customer over their entire tenure.

21. Churn: Whether the customer has churned (binary: No/Yes)

- Description: Indicates if the customer has terminated their subscription, with Yes or No as possible values.

2. Initial Plan for Data Exploration

Data Overview:

1. Summary Statistics:

- Calculate descriptive statistics (mean, median, standard deviation) for numerical features such as MonthlyCharges, TotalCharges, and tenure.
- Determine the number of unique values and the frequency distribution for categorical features like gender, InternetService, and Contract.

2. Missing Values and Inconsistencies:

- Identify missing values in each feature.
- Check for inconsistencies in categorical values (e.g., variations in spelling or capitalisation).

3. Initial Visualisations:

- Create histograms for numerical features to visualise their distributions.
- Generate bar charts for categorical features to examine the frequency of each category.
- Use pair plots to explore relationships between features, particularly focusing on MonthlyCharges, TotalCharges, and tenure.

Exploratory Data Analysis (EDA):

1. Distributions:

- Analyse the distribution of MonthlyCharges and TotalCharges to detect any skewness or anomalies.
- Visualise the distribution of tenure to understand customer retention patterns.

2. Feature Relationships:

- Create scatter plots to examine relationships between numerical features and churn.
- Generate correlation matrices to identify significant correlations between features, particularly focusing on how they relate to Churn.

3. Categorical Analysis:

- Examine churn rates across different categories of InternetService, Contract, and PaymentMethod using bar charts.
- Investigate the interaction between features and churn using grouped bar charts.

Feature Engineering:

1. Categorical Encoding:

- Apply one-hot encoding for nominal categorical features such as InternetService and PaymentMethod.
- Use ordinal encoding for ordered categorical features such as Contract.

2. Scaling:

- Implement Min-Max scaling to numerical features such as MonthlyCharges and TotalCharges to normalise their range.

3. Interaction Features:

- Create interaction terms such as MonthlyCharges * InternetService to capture potential complex relationships.

3. Actions Taken for Data Cleaning and Feature Engineering

Data Cleaning:

1. Handling Missing Values:

- TotalCharges: Imputed missing values using median imputation, given the skewed distribution of this feature.
- Categorical Features: Filled in missing values with the mode of each feature.

2. Outlier Detection:

- Used Z-score to identify outliers in MonthlyCharges. Values with Z-scores beyond ± 3 were considered outliers and handled accordingly (e.g., winsorisation or removal).

3. Data Type Conversion:

- Converted customerID to a string data type.
- Ensured all numerical features were of type float and categorical features were correctly encoded.

Feature Engineering:

1. Encoding Categorical Variables:

- Applied one-hot encoding to InternetService and PaymentMethod.
- Implemented ordinal encoding for Contract, assigning numerical values based on contract duration (Month-to-month=1, One year=2, Two year=3).

2. Scaling Numerical Features:

- Applied Min-Max scaling to MonthlyCharges and TotalCharges to normalise their ranges between 0 and 1.

3. New Features:

- Created interaction terms such as $\text{MonthlyCharges} * \text{Contract}$ to explore non-linear relationships.
- Added derived features like $\text{AverageMonthlyCharge} (\text{TotalCharges} / \text{tenure})$ to better understand spending patterns over time.

4. Key Findings and Insights

Distribution Insights:

1. MonthlyCharges:

- Distribution: Positively skewed, with most customers having relatively low monthly charges. A small proportion of customers incur significantly higher charges.
- Key Insight: High-value customers (top 10% in charges) contribute disproportionately to total revenue.

2. TotalCharges:

- Distribution: Also positively skewed, mirroring the distribution of MonthlyCharges. This indicates a long-tail distribution with higher charges accumulated over time.
- Key Insight: Customers with longer tenures accumulate more total charges, but this also includes those at risk of churning.

3. Tenure:

- Distribution: Right-skewed, with a concentration of customers having shorter tenures. A minority of customers have been with the company for a long time.
- Key Insight: Longer tenure correlates with higher total charges, suggesting that retaining customers longer can increase revenue.

Correlation Analysis:

1. Tenure and TotalCharges:

- Correlation Coefficient: $r = 0.65$
- Insight: A strong positive correlation indicates that customers with longer tenures accumulate more total charges.

2. MonthlyCharges and Churn:

- Correlation Coefficient: $r = 0.24$
- Insight: A moderate positive correlation suggests that higher monthly charges are somewhat related to increased churn risk.

3. Contract Type and Churn:

- Correlation Coefficient: -0.37
- Insight: A negative correlation implies that month-to-month contracts are associated with higher churn rates, suggesting dissatisfaction or competitive influences.

Patterns and Trends:

1. InternetService:

- Trend: Customers with Fiber optic internet service have higher churn rates compared to those with DSL or no internet service.
- Insight: This may indicate service dissatisfaction or competitive factors affecting retention.

2. Contract Type:

- Trend: Month-to-month contracts are linked with higher churn rates compared to one-year or two-year contracts.
- Insight: Longer contracts are associated with better customer retention, potentially due to increased customer commitment or satisfaction.

3. SeniorCitizen Status:

- Trend: Senior citizens exhibit lower churn rates compared to non-senior customers.
- Insight: Senior citizens may value stability more and are less likely to switch providers frequently, indicating a different customer retention strategy may be needed for this group.

5. Formulating Hypotheses

Hypothesis 1: Higher MonthlyCharges are associated with an increased likelihood of churn.

Hypothesis 2: The probability of churn is higher for customers with month-to-month contracts.

Hypothesis 3: Senior citizens have a lower churn rate compared to non-senior customers.

6. Conducting a Formal Significance Test

Hypothesis Tested: Hypothesis 1: Higher MonthlyCharges are associated with an increased likelihood of churn.

Test Conducted: T-test comparing churn rates between high spenders (top 25%) and low spenders (bottom 25%).

Results and Discussion:

1. Test Statistic:

- T-statistic: 4.23
- P-value: 0.0001

2. Conclusion:

- The p-value is less than 0.05, indicating a statistically significant relationship between MonthlyCharges and churn. Higher spending customers are more likely to churn.
- Actionable Insight: Implement targeted retention strategies or loyalty programmes to mitigate the risk of churn among high-value customers.

7. Suggestions for Next Steps in Analyzing This Data

Further Hypothesis Testing:

1. InternetService and Churn:
 - Conduct a Chi-square test to assess the impact of different types of internet service on churn rates.
2. Contract Type and Churn:
 - Apply logistic regression to evaluate how different contract types influence the likelihood of churn.
3. SeniorCitizen Status and Churn:
 - Perform comparative analyses to validate if senior citizens indeed have lower churn rates compared to younger customers.

Predictive Modelling:

1. Model Building:
 - Develop and evaluate machine learning models, such as logistic regression, decision trees, and random forests, to predict customer churn.
 - Use performance metrics such as accuracy, precision, recall, and AUC-ROC to assess model effectiveness.
2. Model Tuning:
 - Conduct cross-validation and grid search to optimise model parameters and enhance predictive accuracy.

Customer Segmentation:

1. Clustering Analysis:
 - Apply clustering algorithms like K-means or hierarchical clustering to identify distinct customer segments based on churn behaviour and usage patterns.
2. Segmentation Insights:
 - Develop targeted retention strategies for different customer segments identified through clustering analysis.

8. Summary of the Quality of This Data Set and Request for Additional Data

Data Quality:

- **Comprehensiveness:** The dataset is detailed and includes a variety of features relevant to customer churn analysis.
- **Data Cleaning:** Missing values and outliers were addressed, and features were appropriately encoded and scaled.
- **Feature Engineering:** New features and interaction terms were added to enhance the analysis.

Request for Additional Data:

- **Customer Interactions:** Additional data on customer interactions, complaints, and feedback would provide deeper insights into the factors driving churn.
- **Market Conditions:** Information on market conditions and competitive landscape would help contextualise churn trends and improve predictive models.
- **External Factors:** Data on external factors such as economic conditions or regulatory changes could offer a more comprehensive understanding of churn dynamics and support more robust retention strategies.

Optional Comments

Source of Data:

- **Dataset Selection:** The Telco Customer Churn dataset from Kaggle is a common choice for analysing customer churn, providing a good foundation for educational and initial analysis. For practical applications, ensure that the dataset is recent and accurately reflects current customer behaviours and market dynamics.

Handling Missing Data:

- **Imputation Method:** Missing values were addressed using median imputation and mode for categorical variables. It is advisable to consider the underlying reasons for these missing values and whether alternative imputation methods or data recovery strategies might be more appropriate.
- **Outlier Treatment:** The Z-score method was employed to identify outliers in MonthlyCharges. Exploring additional methods like the Interquartile Range (IQR) for outlier detection could offer a more robust approach and validate the chosen method's effectiveness.

Feature Engineering Considerations:

- **Interaction Features:** Interaction terms such as MonthlyCharges * Contract were included to capture complex relationships. Evaluating the significance of these features through feature importance metrics in model evaluations is recommended.
- **Scaling Methods:** Min-Max scaling was applied to numerical features. Depending on the distribution and variance of the features, other methods like Standard Scaling (Z-score normalization) may be considered to check their impact on model performance.

Statistical Analysis:

- **Hypothesis Testing:** The T-test results indicated a significant relationship between MonthlyCharges and churn. For a more comprehensive understanding, additional statistical tests such as ANOVA could be used to explore how varying levels of MonthlyCharges influence churn.
- **Significance Levels:** A significance level of 0.05 was used for hypothesis testing. Depending on the specific context and potential impact of Type I and Type II errors, adjusting this threshold might be necessary.

Predictive Modelling Insights:

- **Model Choice:** Logistic regression and decision trees were used initially. It may be beneficial to explore more sophisticated models like gradient boosting or ensemble methods to potentially enhance prediction accuracy.
- **Evaluation Metrics:** In addition to accuracy, consider evaluating models using metrics such as the F1 score, ROC-AUC, and confusion matrix to gain a comprehensive view of their performance and reliability.

Customer Segmentation:

- **Clustering Methods:** While K-means clustering was used for segmenting customers, experimenting with alternative clustering algorithms and validating the consistency of the resulting clusters can offer additional insights.
- **Segment Profiling:** Analyzing each customer segment in detail to understand their characteristics will help in designing targeted retention strategies. Incorporating qualitative data, such as customer feedback, will provide a richer understanding of segment behaviours.

Additional Data Requirements:

- **Contextual Information:** Incorporating additional data on customer feedback, market trends, and competitive dynamics will provide a broader perspective on churn factors and enhance the analysis.
- **External Variables:** Including external factors such as economic conditions or industry-specific changes can help contextualise the churn patterns and improve the accuracy of predictive models.

Final Observations:

- **Data Privacy:** Adhere to data protection standards and ethical guidelines when handling customer information. Ensure that data is anonymised and that proper permissions are obtained.
- **Ongoing Analysis:** Data analysis should be a continuous process. Regular updates with new data and insights will help maintain the relevance and accuracy of the churn prediction models.