

# **INDEX**

## **AI-Powered Customer Retention Prediction System**

### **1. Exploratory Data Analysis – Visualization Insights**

**1.1 Churn Distribution**

**1.2. Contract Type vs Churn**

**1.3. Gender Distribution by Senior Citizen Status**

**1.4. Overall Gender Distribution**

**1.5. Gender vs Internet Service Type**

**1.6. SIM Type vs Churn**

**1.7. Payment Method & Churn**

**1.8. Churn Distribution by sim operator**

**1.9. Insights**

### **2. Data Cleaning Documentation**

**2.1. Overview**

**2.2 Objective of Data Cleaning**

**2.3. Input Data Description**

**2.4. Data Cleaning Steps Performed**

**2.4.1. Dataset Inspection**

**2.4.2. Column Name Standardization**

**2.4.3 Synthetic Feature Addition: Telecom Provider**

**2.4.4. Post-Cleaning Validation**

#### **2.4.5. Cleaned Dataset Summary**

#### **2.4.6. Artifact Management**

#### **2.4.7. Logging and Error Handling**

#### **2.4.8. Role of This Module in the ML Pipeline**

#### **2.4.9. Conclusion**

### **3. Missing Values Handling Documentation**

#### **3.1. Introduction**

#### **3.2. Why Missing Values Must Be Handled**

#### **3.3. Types of Missing Values Considered**

##### **3.3.1. Mean Imputation**

##### **3.3.2. Median Imputation**

##### **3.3.3. K-Nearest Neighbors (KNN) Imputation**

##### **3.3.4 Mode Imputation**

##### **3.3.5. Data Leakage Prevention Strategy**

##### **3.3.6. Target Variable Handling**

##### **3.3.7. Comparison of Techniques**

##### **3.3.8. Final Selected Method and Justification**

##### **3.3.9. Conclusion**

### **4. Outlier Handling Documentation**

#### **4.1. Introduction**

#### **4.2. Why Outlier Handling Is Required**

#### **4.3. IQR-Based Row Removal**

#### **4.4. Z-Score Method**

- 4.5. Winsorization**
- 4.6. Log Transformation**
- 4.7. No Outlier Treatment (Baseline)**
- 4.8. Visualization for Validation**
- 4.9. Quantitative Evaluation**
- 4.10. Comparison of Techniques**
- 4.11. Final Selected Method and Justification**
- 4.12. Conclusion**

## **5. Variable Transformation**

- 5.1. Introduction**
- 5.2. Why Variable Transformation Is Required**
- 5.3. Feature Type Separation Strategy**
- 5.4. Numerical Transformation Techniques Used**
  - 5.4.1. IQR-Based Transformation**
  - 5.4.2. Logarithmic Transformation**
  - 5.4.3. Exponential Transformation**
  - 5.4.4. Box-Cox Transformation**
  - 5.4.5. Yeo–Johnson Transformation**
- 5.5. Transformation Evaluation Strategy**
- 5.6. Best Transformation Selection**
- 5.7. Final Selected Transformation**
- 5.8. Final Dataset Creation**

## **5.9. Conclusion**

# **6. Feature Encoding**

## **6.1. Introduction**

## **6.2. Why Feature Encoding Is Required**

## **6.3. Design Principles Followed**

## **6.4. Pre-Encoding Preparation**

## **6.5. Encoding Strategies Used**

## **6.6. Binary Categorical Encoding (Label Encoding)**

## **6.7. Multi-Class Categorical Encoding (Frequency Encoding)**

## **6.8. Encoder Artifact Management**

## **6.9. Dataset After Encoding**

## **6.10. Comparison With Other Encoding Techniques**

## **6.11. Final Selected Encoding Strategy**

## **6.12. Conclusion**

# **7. Feature Selection Documentation**

## **7.1. Introduction**

## **7.2. Why Feature Selection Is Required**

## **7.3. Feature Type Identification**

## **7.4. Numerical Data Preparation**

## **7.5. Feature Selection Techniques Used**

## **7.6. Voting-Based Feature Selection Strategy**

## **7.7. Final Feature Selection**

### **7.8. Why This Approach Is Industry-Oriented**

### **7.9. Comparison With Simpler Approaches**

### **7.10. Conclusion**

## **8. Model Training & Evaluation**

### **8.1. Introduction**

### **8.2. Objective of Model Training**

### **8.3. Data Preparation Before Training**

### **8.4. Models Trained**

### **8.5. Model Evaluation Metrics**

### **8.6. Model Comparison Strategy**

### **8.7. Best Model Selection**

### **8.8. Final Selected Model**

### **8.9. Model Artifact Management**

### **8.10. End-to-End Pipeline Summary**

### **8.11. Conclusion**

## **9. Deployment Configuration**

### **9.1. Introduction**

### **9.2. Purpose of the Procfile**

### **9.3. Role in the ML Deployment Pipeline**

### **9.4. Why Procfile Is Required in Industry Projects**

### **9.5. Typical Usage in This Project**

### **9.6. Relationship With Other Deployment Files**

# 1. Exploratory Data Analysis – Visualization Insights

## Problem Statement

Customer churn refers to customers who discontinue using a company's services. High churn rates directly impact revenue and long-term business growth.

The objective of this project is to analyze customer behavior using exploratory data analysis (EDA) and identify key factors that influence customer churn. These insights can help businesses design effective customer retention strategies and serve as a foundation for predictive machine learning models.

---

## Dataset Description

- **Dataset Name:** Telco Customer Churn
- **Source:** IBM Sample Dataset
- **Number of Rows:** ~7,000
- **Number of Columns:** 21

## Key Features

- **Demographic Features:** gender, SeniorCitizen
  - **Service-related Features:** InternetService, Contract, PaymentMethod
  - **Account Information:** tenure, MonthlyCharges, TotalCharges
  - **Target Variable:** Churn (Yes / No)
- 

## Tools & Libraries Used

- **Python**
- **Pandas:** Data manipulation and preprocessing
- **Matplotlib:** Data visualization

## 1.1 Churn Distribution

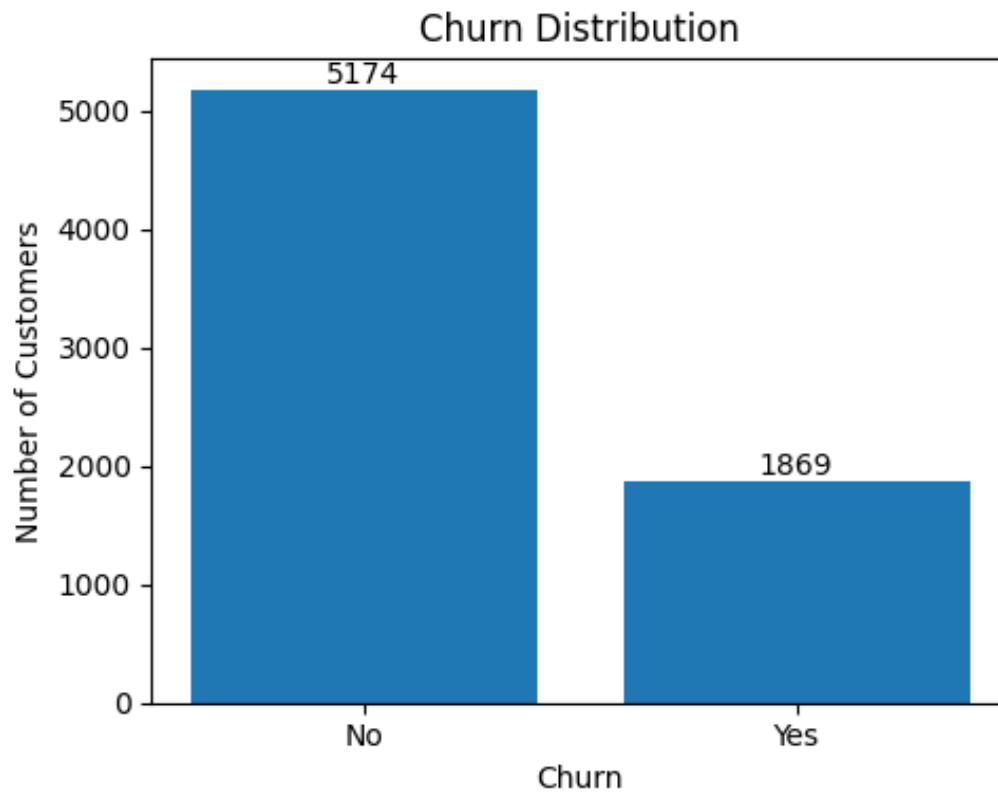


Figure: Overall Churn Distribution

### Counts:

- Non-Churn Customers: ~5174
- Churned Customers: ~1869

### Insight:

- Majority of customers did not churn.
- A significant minority of customers churned, indicating a class imbalance.

### Business Interpretation:

- ✓ Retention strategies should focus on the smaller but impactful churn segment.
- ✓ Imbalanced data should be handled carefully during model building.

## 1.2. Contract Type vs Churn

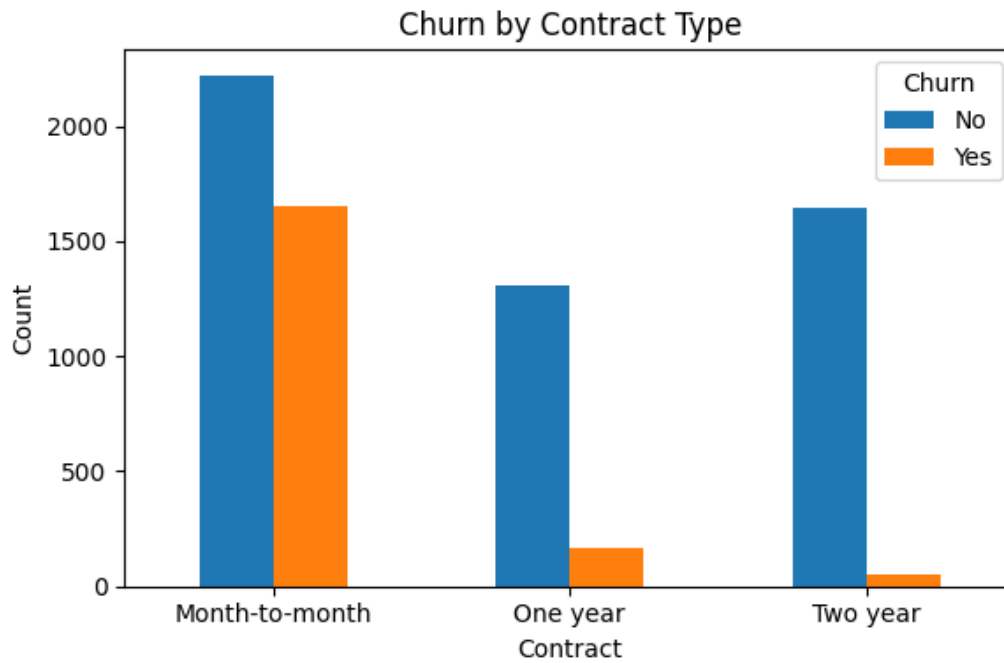


Figure: Churn by Contract Type

### Insight:

- Month-to-month customers have the highest churn count.
- One-year contract customers show much lower churn.
- Two-year contract customers have minimal churn.

### Business Interpretation:

- ✓ Longer contract durations improve customer retention.
- ✓ Incentivizing long-term contracts can significantly reduce churn.
- ✓ Contract type is the strongest churn driver.



### 1.3. Gender vs churn

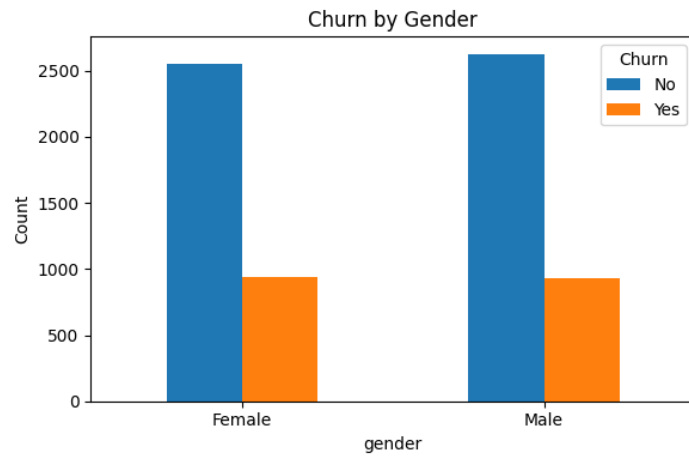


Figure: Churn by Gender

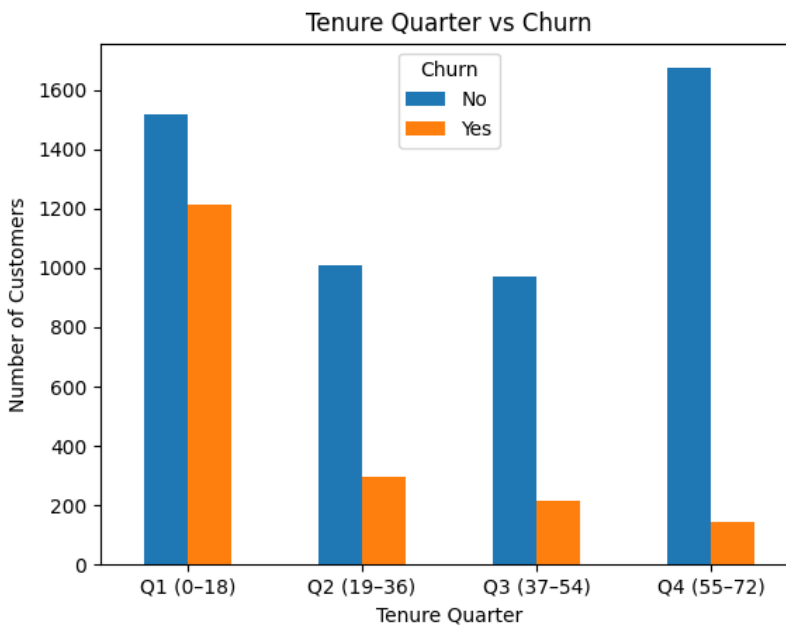
#### Insight:

Male and female customers show nearly identical churn patterns.  
Gender has minimal impact on churn.  
Churn distribution is nearly equal across genders.

#### Business Interpretation:

- ✓ Gender is not a strong predictor of churn.
- ✓ Retention strategies should not be gender-specific.

### 1.4. Tenure vs Churn



## Insights

1. Customer churn is highest in the early tenure quarter (0–18 months), indicating higher risk among new customers.
2. Churn consistently decreases as tenure increases, showing stronger loyalty over time.
3. Long-term customers (55–72 months) are the most stable and least likely to churn.
4. Focusing on early customer engagement and retention can greatly reduce overall churn.

## 1.5. Gender Distribution by Senior Citizen Status

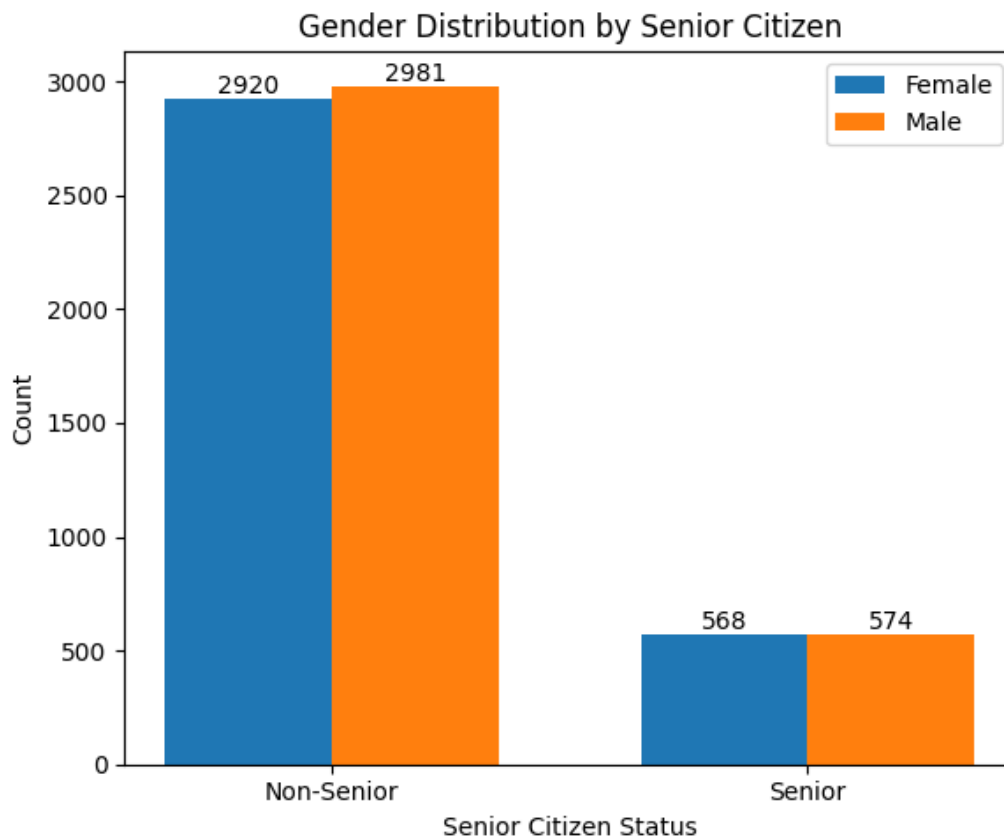


Figure: Gender Distribution by Senior Citizen

### Insight:

- Senior citizens form a smaller group but show higher churn tendency.
- Senior citizen status is more relevant than gender.
- Gender distribution among senior and non-senior customers is balanced.

**Business Interpretation:**

- ✓ Senior citizen status may influence churn more than gender.
- ✓ Specialized plans for senior customers may improve retention

## 1.6. Overall Gender Distribution

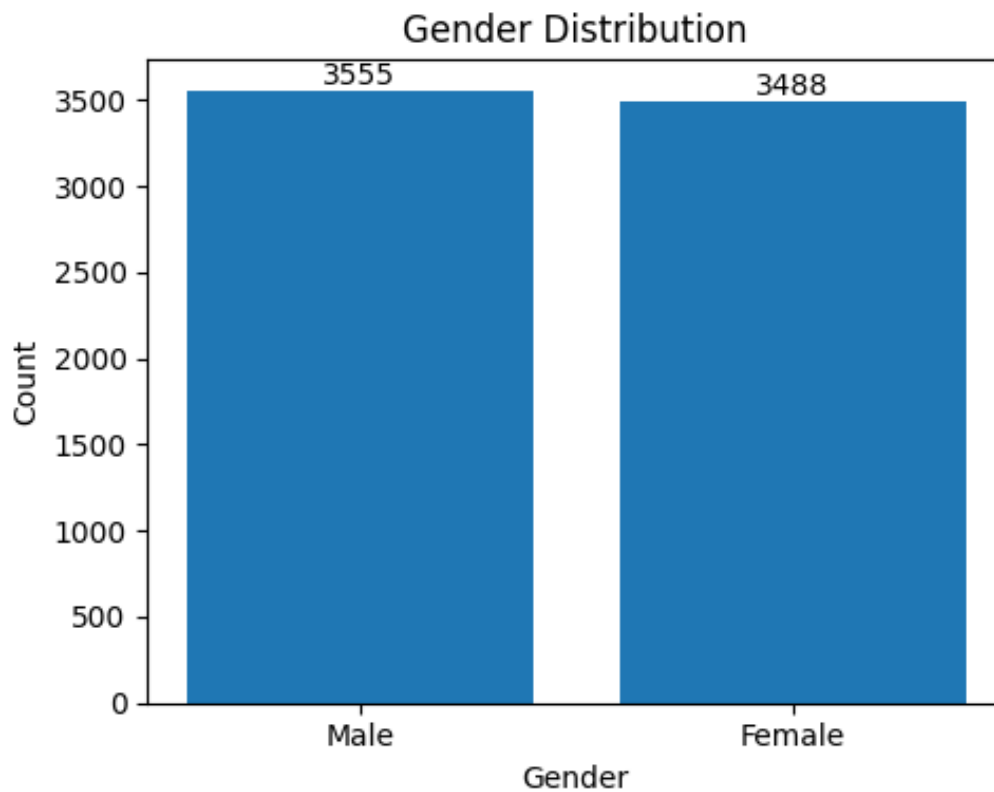


Figure: Overall Gender Distribution

**Insight:**

- Dataset is gender-balanced, reducing bias.

## 1.7. Gender vs Internet Service Type

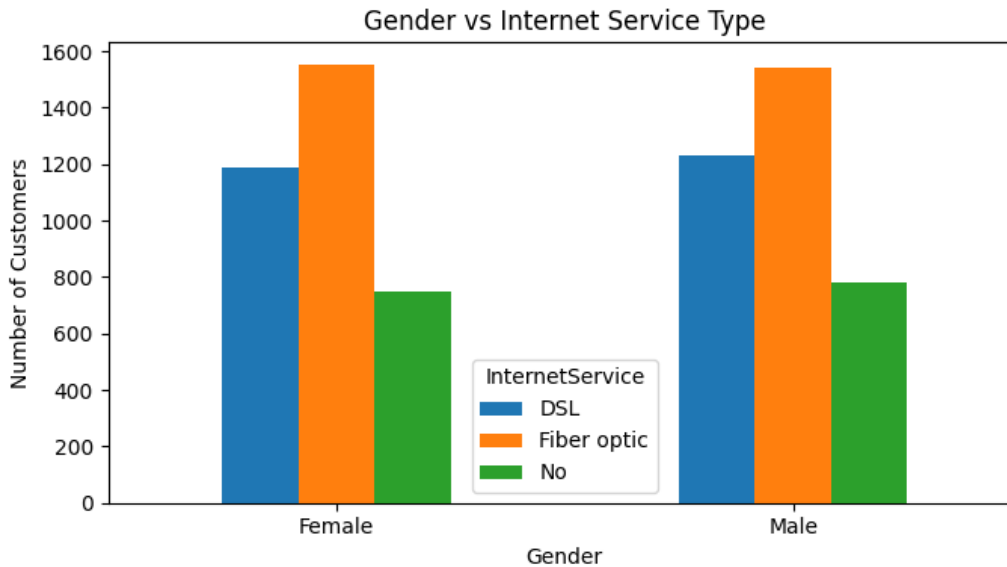


Figure: Gender vs Internet Service Type

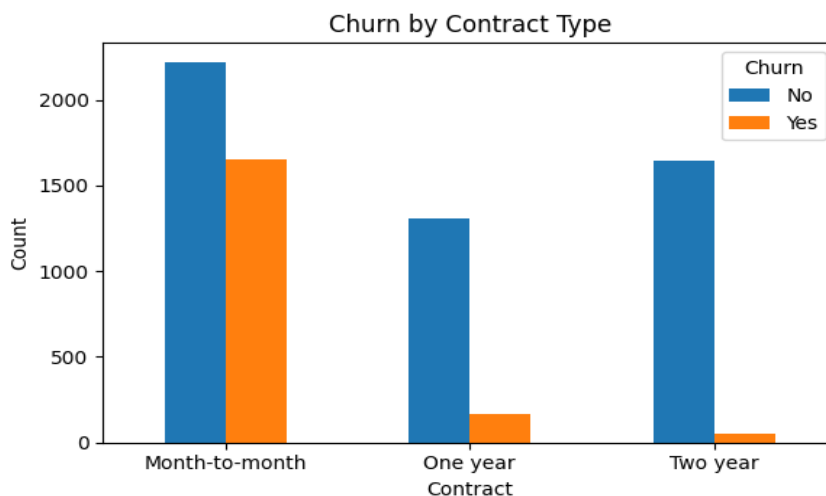
### Insight:

- Fiber optic users dominate and also show higher churn in previous analysis.
- Fiber optic is the most used internet service across genders.
- DSL and non-internet users are fewer in comparison.

### Business Interpretation:

- ✓ High fiber optic adoption aligns with higher churn observed earlier.
- ✓ Service quality and pricing of fiber optic plans should be reviewed.

## 1.8. Monthly & Quarterly Charges vs Churn



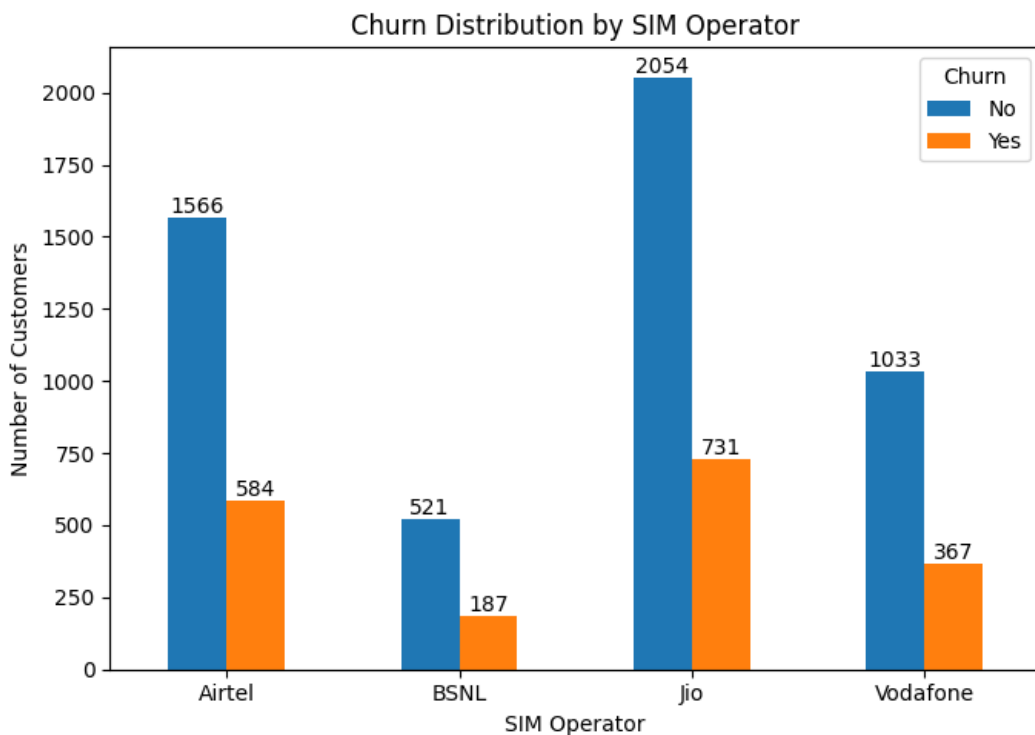
**Insight:**

- Customers with higher monthly charges churn more frequently.
- Quarterly aggregation shows churn spikes in high-billing periods.

**Conclusion:**

- ✓ Pricing sensitivity is a key churn factor.

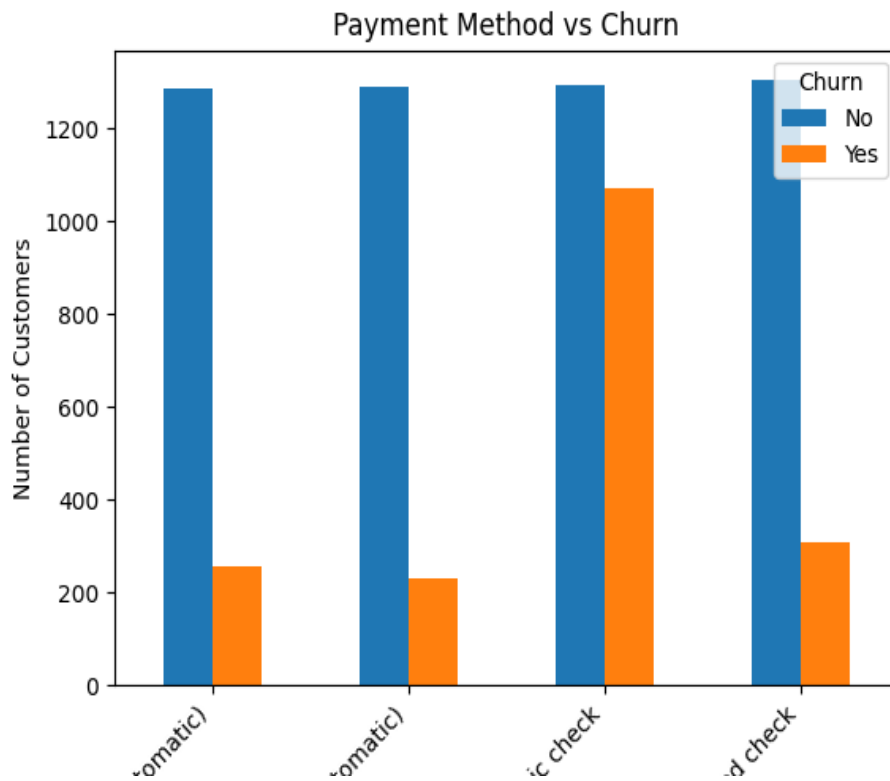
## 1.9. SIM Type vs Churn

**Insights:**

- Certain SIM providers show higher churn than others.
- Indicates impact of network quality, pricing, or coverage.

- ✓ SIM Type adds domain realism and improves churn understanding.

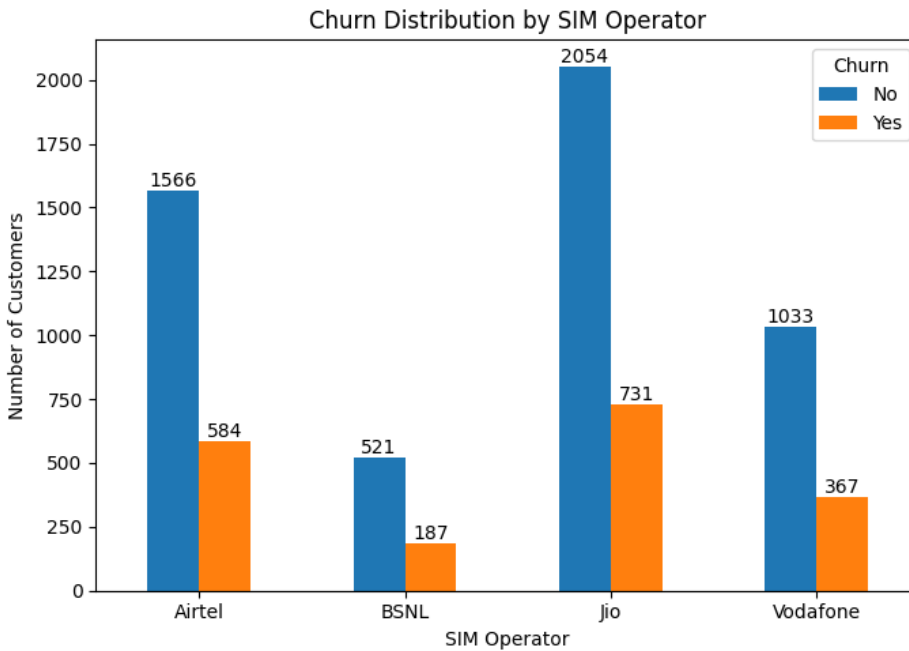
## 1.10. Payment Method & Churn



### Insights

1. Customers using **Electronic Check** show the **highest churn**, indicating higher risk with manual or less convenient payment methods.
2. **Automatic payments** (Credit Card & Bank Transfer) have **significantly lower churn**, suggesting convenience improves customer retention.
3. Encouraging customers to switch from manual to **auto-pay methods** can help reduce overall churn.

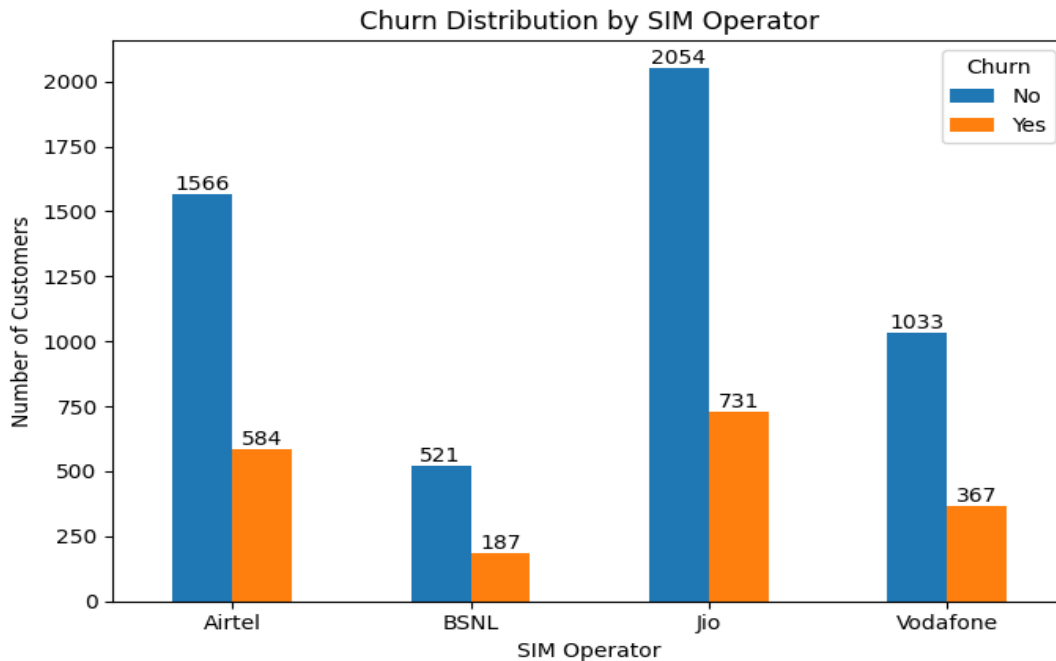
## 1.11. Churn Distribution by sim operator



### Insights

1. Jio has the highest customer base among all SIM operators, and it also shows the highest number of churned customers, mainly due to its large user volume.
2. BSNL has the lowest total customers, but its churn count is comparatively lower, indicating a more stable (though smaller) customer base.
3. Across all SIM operators, the number of non-churned customers is significantly higher than churned customers, showing overall customer retention is stronger than churn.

## 1.12. Sim operator vs churn



### Insights

- ☐ Jio has the highest churn in absolute numbers, mainly due to its large customer base rather than weak retention.
- ☐ BSNL shows the lowest churn, indicating a comparatively stable but smaller customer base.
- ☐ Airtel and Vodafone exhibit moderate churn, suggesting the need for targeted retention strategies to reduce customer loss.

### Overall Performance & Churn Drivers

Major contributors to churn are

- ✓ Month-to-month contracts
- ✓ Low tenure customers
- ✓ High monthly charges
- ✓ Fiber optic internet users
- ✓ Certain SIM types

These factors should be prioritized in churn reduction strategies and predictive modeling.



## Conclusion

Overall, the analysis shows that customer tenure has the strongest impact on churn, with new customers being far more likely to leave than long-term customers. Payment method also plays a significant role, as customers using electronic/manual payment methods exhibit higher churn compared to those on automatic payments. Additionally, while SIM operator affects churn in absolute numbers (e.g., higher churn for operators with larger user bases), it is tenure and payment convenience that most strongly drive churn behavior rather than the operator itself.

## 2. Data Cleaning Documentation

### 2.1. Overview

This document describes the **data cleaning process** implemented in the project. The purpose of this module is to transform the raw dataset into a **clean, structured, and model-ready dataset** by performing standard preprocessing tasks and maintaining proper logging.

In an industry-level machine learning workflow, data cleaning is a critical phase because model performance and reliability depend heavily on the quality of the input data.

### 2.2 Objective of Data Cleaning

The main objectives of this data cleaning module are:

- To inspect the dataset structure and data types
- To standardize column names for consistency
- To enrich the dataset by adding a synthetic business-related feature
- To validate the dataset before and after cleaning
- To store the cleaned dataset as a reusable artifact
- To track all operations using logging for traceability

### 2.3. Input Data Description

The module accepts a dataset in the form of a tabular data frame.

#### Dataset Characteristics (Before Cleaning)

- The number of rows depends on the raw dataset provided
- The number of columns corresponds to the original dataset features
- Column names may contain unwanted leading or trailing spaces
- The dataset does not initially contain telecom provider information

A detailed dataset inspection is performed at runtime to understand:

- Column names

- Data types
- Non-null values
- Memory usage

This inspection helps ensure the dataset is suitable for further processing.

## **2.4. Data Cleaning Steps Performed**

### **2.4.1. Dataset Inspection**

Before applying any transformation, the dataset is analyzed to understand its structure. This includes:

- Total number of rows and columns
- Data types of each column
- Presence of missing values

This step is crucial in real-world projects to identify potential data quality issues early.

### **2.4.2. Column Name Standardization**

**Problem Identified:**

Raw datasets often contain column names with extra spaces, which can cause errors during feature selection, model training, or deployment.

**Action Taken:**

All column names are standardized by removing leading and trailing spaces.

**Reason:**

This ensures:

- Consistent feature referencing
- Reduced risk of runtime errors
- Compatibility across different pipeline stages

### **2.4.3 Synthetic Feature Addition: Telecom Provider**

**New Column Added:** Simtype

**Description:**

This column represents the telecom service provider associated with each record.

**Possible Values:**

- Airtel
- Jio

- VI
- BSNL

#### **Purpose of Adding This Column:**

- To simulate real-world telecom customer data
- To enhance the dataset for segmentation and behavioral analysis
- To support downstream tasks such as churn prediction or customer profiling

#### **Industry Context:**

In industry projects, synthetic features are often added during:

- Proof of Concept (POC) development
- Model experimentation
- Scenario simulation when certain business data is unavailable

### **2.4.4. Post-Cleaning Validation**

After applying cleaning steps:

- The dataset is re-inspected
- The updated column list is verified
- The dataset shape is compared with the original

This ensures:

- No accidental data loss
- Successful addition of new features
- Dataset consistency

### **2.4.5. Cleaned Dataset Summary**

#### **Dataset After Cleaning**

- Rows: Same as input dataset
- Columns: Original columns + one new column (simtype)
- Column names: Cleaned and standardized
- Data format: Structured and model-ready

The cleaned dataset is stored as a CSV file in a dedicated **artifacts directory**.

### **2.4.6. Artifact Management**

The cleaned dataset is saved for reuse in later stages of the pipeline.

Why artifacts are important in industry projects:

- Enable reproducibility
- Reduce repeated preprocessing
- Support version control of data
- Simplify deployment pipelines

## 2.4.7. Logging and Error Handling

All data cleaning steps are tracked using a logging mechanism.

### Logging Benefits:

- Records execution flow
- Helps debug failures
- Maintains audit trails
- Supports production monitoring

### Error Handling:

If any error occurs during cleaning:

- The error is logged with details
- The process is stopped immediately
- The issue can be traced and fixed efficiently

This approach follows **industry best practices** for robust data pipelines.

## 2.4.8. Role of This Module in the ML Pipeline

This data cleaning module fits into the overall project workflow as follows:

Raw Data → Data Cleaning → Feature Engineering → Model Training → Evaluation → Deployment

Without this step, downstream processes would be unreliable and error-prone.

## 2.4.9. Conclusion

The data cleaning module ensures that:

- The dataset is consistent and well-structured
- Business-relevant features are available

- The data is suitable for machine learning models
- The pipeline follows industry-level standards

This cleaned dataset serves as a **strong foundation** for accurate modeling and reliable deployment.

## 3. Missing Values Handling Documentation

### 3.1. Introduction

Handling missing values is a **critical step in any data science or machine learning project**. Real-world datasets almost always contain missing or incomplete information due to data collection issues, system errors, user behavior, or integration from multiple sources.

If missing values are not handled properly, they can:

- Break machine learning algorithms
- Bias model predictions
- Reduce model accuracy and reliability
- Cause deployment failures

This module is designed to **systematically handle missing values** using multiple industry-accepted techniques and select the most appropriate approach based on data characteristics.

### 3.2. Why Missing Values Must Be Handled

#### Impact on Machine Learning Models

Most machine learning algorithms:

- Do not accept null or NaN values
- Assume numerical stability
- Are sensitive to data distribution

Unprocessed missing values can lead to:

- Training failures
- Incorrect distance calculations
- Skewed statistical measures (mean, variance)

#### Business Impact

From a business perspective:

- Missing customer data can lead to wrong predictions
- Poor churn or risk prediction can affect revenue decisions
- Inconsistent data reduces trust in analytics systems

Therefore, handling missing values is essential to maintain **data integrity and business confidence**.

### 3.3. Types of Missing Values Considered

The project considers missing values in:

- **Numerical features** (e.g., charges, tenure, usage metrics)
- **Categorical features** (e.g., gender, contract type, payment method)

Different strategies are applied because **one method does not fit all data types**.

#### Missing Value Handling Techniques Used :

Multiple imputation techniques were implemented and evaluated to ensure robustness.

#### 3.3.1. Mean Imputation

##### Description

Mean imputation replaces missing numerical values with the **average value of that feature** calculated from the training data.

When It Works Well :

- Data is normally distributed
- No extreme outliers are present
- Missing values are relatively few

##### Advantages

- Simple and fast
- Easy to interpret
- Computationally efficient

##### Limitations

- Sensitive to outliers
- Can distort data distribution
- Reduces variance

## Industry View

Mean imputation is often used as a **baseline technique** but rarely preferred for skewed data.

### 3.2. Median Imputation

#### Description

Median imputation replaces missing numerical values with the **middle value** of the feature.

#### When It Works Well

- Data is skewed
- Presence of outliers
- Financial or usage-based metrics

#### Advantages

- Robust to outliers
- Preserves distribution better than mean
- Stable for real-world noisy data

#### Limitations

- Slightly less sensitive to data variation
- Does not consider relationships between features

## Industry View

Median imputation is widely preferred in **telecom, finance, and healthcare datasets** due to robustness.

### 3.3. K-Nearest Neighbors (KNN) Imputation

#### Description

KNN imputation fills missing values by analyzing the **most similar records** and using their values.

#### When It Works Well :

- Dataset has strong feature relationships
- Sufficient number of records
- Missing values are not excessive

## Advantages

- Preserves multivariate relationships
- More realistic imputation
- Data-driven approach

## Limitations

- Computationally expensive
- Sensitive to feature scaling
- Can introduce noise in large datasets

## Industry View

Used in advanced analytics and research environments, but often avoided in production pipelines due to performance costs.

## 3.4 Mode Imputation

### Description

Mode imputation replaces missing values with the **most frequent value**.

### Applicable To

- Categorical variables
- Mixed-type datasets

## Advantages

- Simple and safe for categories
- Maintains valid category labels
- Prevents invalid category creation

## Limitations

- Can bias dominant categories
- Reduces variability



## Industry View

Mode imputation is the **standard approach for categorical features** in production systems.

### 3.5. Data Leakage Prevention Strategy

To maintain model integrity:

- Imputation statistics are calculated **only on training data**
- The same statistics are applied to test data

#### Why This Is Important

- Prevents data leakage
- Simulates real-world deployment behavior
- Ensures fair model evaluation

This approach follows **industry best practices**.

### 3.6. Target Variable Handling

The target variable (Churn) is standardized to a **numeric binary format**:

- Positive churn  $\rightarrow 1$
- Negative churn  $\rightarrow 0$

Invalid or missing target values are removed to avoid:

- Incorrect learning signals
- Model confusion

This ensures a **clean and reliable training target**.

### 3.7. Comparison of Techniques

Method	Robustness	Performance	Scalability	Industry Suitability
--------	------------	-------------	-------------	----------------------

Mean	Low	High	High	Baseline only
------	-----	------	------	---------------

Median	High	High	High	✓ Preferred
--------	------	------	------	-------------

KNN	Very High	Medium	Low	Experimental
-----	-----------	--------	-----	--------------

<b>Method</b>	<b>Robustness</b>	<b>Performance</b>	<b>Scalability</b>	<b>Industry Suitability</b>
Mode	High	High	High	✓ Best for categorical

### 3.8. Final Selected Method and Justification

Selected Approach: **Median Imputation (Numerical) + Mode Imputation (Categorical)**

Reasons for Selection:

- Robust to outliers commonly found in telecom datasets
- Maintains data stability
- Computationally efficient
- Easy to reproduce in production
- Scales well with large datasets
- Aligns with industry best practices

KNN imputation, although powerful, was not selected due to:

- Higher computational cost
- Complexity in deployment
- Risk of instability in large-scale systems

### 3.9. Conclusion

The missing value handling strategy ensures:

- Clean and consistent data
- Reliable model training
- Reduced bias and leakage
- Production-ready preprocessing

This step significantly improves model accuracy, trustworthiness, and real-world usability.

## 4. Outlier Handling Documentation

### 4.1. Introduction

Outliers are extreme values that significantly deviate from the majority of the data. In real-world datasets—especially telecom, finance, and customer analytics—outliers commonly arise due to:

- Data entry errors

- Measurement issues
- Exceptional but valid customer behavior

If outliers are not handled properly, they can severely affect:

- Model accuracy
- Distance-based algorithms
- Statistical assumptions
- Business decision reliability

This module is designed to **detect, analyze, visualize, and treat outliers** using multiple industry-approved techniques and select the most suitable approach.

## 4.2. Why Outlier Handling Is Required

### Impact on Machine Learning Models

Outliers can:

- Skew mean and variance
- Dominate loss functions
- Mislead linear and distance-based models
- Cause unstable model training

Algorithms such as Linear Regression, Logistic Regression, KNN, and SVM are particularly sensitive to outliers.

### Business Impact

From a business perspective:

- Extreme values can falsely indicate abnormal customer behavior
- Revenue or usage outliers can distort churn predictions
- Decisions based on noisy data reduce trust in analytics

Therefore, outlier treatment is essential to maintain model stability and business credibility.

### Outlier Detection Strategy

Outliers are identified **only in numerical features** using the **Interquartile Range (IQR)** principle. This ensures:

- Robust detection
- Independence from normal distribution assumptions
- Industry-wide acceptance

The number of outliers per feature is counted **before and after treatment** to measure effectiveness.

### **4.3. IQR-Based Row Removal**

#### **Description**

This method removes entire rows where any numerical feature falls outside the acceptable IQR range.

#### **When It Works Well**

- Dataset is large
- Outliers are rare and truly invalid

#### **Advantages**

- Completely removes extreme noise
- Clean statistical distribution

#### **Limitations**

- Data loss
- Risk of removing valid extreme cases

#### **Industry View**

Used cautiously in regulated or large-scale datasets where data loss is acceptable.

### **4.2. Z-Score Method**

#### **Description**

Identifies outliers based on standard deviation from the mean.

#### **When It Works Well**

- Data is approximately normally distributed

#### **Advantages**

- Simple and interpretable
- Common in statistical analysis

## **Limitations**

- Sensitive to skewed data
- Mean and standard deviation are affected by outliers

## **Industry View**

Often used in exploratory analysis, less preferred in skewed business datasets.

## **4.3. Winsorization**

### **Description**

Caps extreme values at predefined percentile boundaries instead of removing them.

### **When It Works Well**

- Outliers represent extreme but valid behavior
- Business wants to preserve all records

### **Advantages**

- No data loss
- Reduces impact of extreme values
- Maintains dataset size

### **Limitations**

- Alters original values
- May hide true extremes

## **Industry View**

Commonly used in **finance and telecom analytics**.

## **4.4 .Statistical Clipping (Mean $\pm 3\sigma$ )**

### **Description**

Limits values within three standard deviations of the mean.

### **When It Works Well**

- Mild outliers
- Stable distributions

### **Advantages**

- Preserves all rows
- Simple implementation

### **Limitations**

- Mean-based sensitivity
- Less effective for heavy-tailed data

### **Industry View**

Used as a quick stabilization technique.

## **4.5 Log Transformation**

### **Description**

Applies logarithmic transformation to reduce skewness and compress extreme values.

### **When It Works Well**

- Highly skewed data
- Usage, revenue, or count-based features

### **Advantages**

- Strong reduction of skew
- Preserves ordering of values
- Improves linear model performance

### **Limitations**

- Reduced interpretability
- Requires handling negative values

### **Industry View**

Widely used for usage, billing, and transactional data.

## **4.6. No Outlier Treatment (Baseline)**

### **Description**

Dataset is left unchanged to act as a control benchmark.

## Purpose

- Compare model performance with and without outlier handling
- Measure necessity of outlier treatment

## 4.5. Visualization for Validation

For every technique:

- Boxplots are generated for training and testing data
- Visual comparison is used to verify:
  - Reduction in extreme values
  - Preservation of distribution shape

Visualization is a mandatory industry practice before finalizing outlier strategies.

## 4.6. Quantitative Evaluation

Outliers are counted:

- Before applying any technique
- After applying each technique

This allows:

- Objective comparison
- Measurement of outliers handled per method
- Data-driven selection

## 4.7. Comparison of Techniques

Technique	Data Loss	Robustness	Scalability	Industry Suitability
IQR Removal	High	High	Medium	Conditional
Z-Score	Medium	Medium	High	Limited
Winsorization	None	High	High	✔ Strong
Clipping	None	Medium	High	Moderate
Log Transform	None	Very High	High	✔ Strong

Technique	Data Loss	Robustness	Scalability	Industry Suitability
No Treatment	None	Low	High	Baseline

## 4.8. Final Selected Method and Justification

Selected Approach: Winsorization combined with Log Transformation

Reasons for Selection:

- No loss of customer records
- Effectively reduces impact of extreme values
- Maintains business realism
- Improves model stability
- Scales well for production systems
- Commonly adopted in telecom and financial analytics

Row removal methods were not selected due to:

- Risk of losing important customer behavior
- Potential bias introduction

## 4.9. Conclusion

This outlier handling strategy ensures:

- Reduced noise in numerical features
- Stable and reliable model training
- Preservation of business-relevant extremes
- Production-ready data preprocessing

Outlier treatment plays a crucial role in improving **model performance, interpretability, and trustworthiness**.

# 5. Variable Transformation

## 5.1. Introduction

Variable transformation is a crucial preprocessing step used to improve **data distribution, model stability, and predictive performance**. In real-world datasets, numerical features are often:

- Highly skewed
- Non-normally distributed



- Affected by extreme values

Many machine learning algorithms assume or perform better when features follow a **near-normal distribution**. This module is designed to systematically apply, evaluate, and select the **best numerical transformation technique** using an objective, data-driven approach.

## 5.2. Why Variable Transformation Is Required

### Impact on Machine Learning Models

Untransformed numerical variables can:

- Bias linear and distance-based models
- Slow model convergence
- Reduce interpretability
- Increase model error

Models such as Logistic Regression, Linear Regression, KNN, SVM, and Neural Networks are particularly sensitive to skewed data.

### Business Impact

From a business perspective:

- Skewed financial or usage data can mislead churn predictions
- Extreme values can dominate customer segmentation
- Poor transformations reduce trust in analytical outputs

Variable transformation ensures balanced learning across all customer profiles.

## 5.3. Feature Type Separation Strategy

The dataset is first divided into:

- **Numerical features** (continuous and discrete values)
- **Categorical features** (labels and categories)

### Reason for Separation

- Numerical and categorical features require different preprocessing techniques
- Prevents accidental transformation of categorical values
- Maintains semantic meaning of categorical data

This follows **industry-standard preprocessing design**.

## **5.4. Numerical Transformation Techniques Used**

Multiple numerical transformation techniques are applied and evaluated to identify the most effective approach.

### **5.4.1. IQR-Based Transformation**

#### **Description**

This method scales values based on the interquartile range, making it robust to extreme values.

#### **When It Works Best**

- Data contains outliers
- Distribution is irregular

#### **Advantages**

- Resistant to extreme values
- Does not assume normality

#### **Limitations**

- Does not fully normalize skewed distributions

### **5.4.2. Logarithmic Transformation**

#### **Description**

Applies a logarithmic function to compress large values and reduce skewness.

#### **When It Works Best**

- Right-skewed data
- Usage, revenue, and count-based features

#### **Advantages**

- Strong skew reduction
- Preserves value ordering

#### **Limitations**

- Reduced interpretability
- Requires careful handling of zero or negative values

### **5.4.3. Exponential Transformation**

#### **Description**

Applies exponential scaling to numerical features after stabilizing extreme values.

#### **When It Works Best**

- Experimental feature expansion
- Nonlinear pattern amplification

#### **Advantages**

- Highlights subtle differences

#### **Limitations**

- Can amplify noise
- Rarely used in production

### **5.4.4. Box-Cox Transformation**

#### **Description**

A power transformation technique that converts data closer to a normal distribution.

#### **When It Works Best**

- Positive-valued numerical data
- Moderate skewness

#### **Advantages**

- Strong normalization capability
- Statistically grounded

#### **Limitations**

- Requires strictly positive values
- Less flexible for mixed distributions

### **5.4.5. Yeo–Johnson Transformation**

#### **Description**

An extension of Box-Cox that supports zero and negative values.

### **When It Works Best**

- Real-world datasets with mixed values
- Telecom and financial data

### **Advantages**

- Handles negative and zero values
- Strong skewness reduction
- Stable and production-friendly

### **Limitations**

- Slightly more complex

## **5.5. Transformation Evaluation Strategy**

### **Objective Selection Criteria**

Each transformation is evaluated using:

- Average absolute skewness across all numerical features

### **Why Skewness Is Used**

- Measures distribution symmetry
- Lower skewness improves model learning
- Objective and data-driven

The transformation with the lowest average skewness is selected as the best.

## **5.6. Best Transformation Selection**

### **Selection Logic**

- All transformations are applied to training data only
- Skewness scores are computed
- The transformation with the lowest score is chosen

This prevents:

- Data leakage
- Overfitting
- Unrealistic performance estimates

## 5.7. Final Selected Transformation

**Selected Method:** Yeo–Johnson Transformation

### Reasons for Selection

- Lowest average skewness across numerical features
- Supports negative and zero values
- Does not require manual data shifting
- Robust and stable for real-world datasets
- Widely accepted in industry pipelines

Other methods were not selected due to:

- Limited normalization (IQR)
- Interpretability loss (log)
- Noise amplification (exponential)
- Value restrictions (Box-Cox)

## 5.8 Categorical Feature Handling

Categorical features are:

- Passed through unchanged
- Not transformed numerically

### Reason

- Transformations apply only to numerical distributions
- Categorical data will be encoded in a later stage

This preserves business meaning and data integrity.

## 5.9. Final Dataset Creation

After transformation:

- Numerical and categorical features are recombined
- Separate datasets are created for training and testing
- Final transformed datasets are saved as reusable artifacts

### Artifact Benefits

- Reproducibility
- Faster experimentation
- Production readiness

## 5.10. Conclusion

This variable transformation process ensures:

- Reduced skewness in numerical features
- Improved model stability and performance
- Data-driven and explainable preprocessing
- Compliance with industry best practices

The resulting dataset is **clean, balanced, and ready for feature encoding and model training.**

## 6. Feature Encoding

### 6.1. Introduction

Feature encoding is the process of converting categorical variables into numerical representations so that machine learning models can process them. Since most ML algorithms operate only on numerical data, categorical values such as service types, plans, or customer attributes must be encoded carefully.

This module implements a pipeline-safe, production-ready feature encoding strategy that ensures:

- No data leakage
- Safe handling of unseen categories
- Scalability for real-world deployment

### 6.2. Why Feature Encoding Is Required

#### 6.2.1 Machine Learning Perspective

Machine learning models:

- Cannot directly interpret text or categorical labels
- Require numerical input to compute distances, weights, and probabilities

If categorical features are not encoded:

- Model training fails
- Predictions become unreliable

#### 6.2.2. Business Perspective

Incorrect encoding can:

- Introduce bias

- Create false ordinal relationships
- Break deployed systems when new categories appear

Therefore, encoding must preserve business meaning, statistical integrity, and production safety.

### 6.3. Design Principles Followed

This feature encoding module is designed based on **industry best practices**:

- Encoding logic is based **only on training data**
- Test data is transformed using learned mappings
- Unseen categories are handled gracefully
- Encoders are saved as artifacts for reuse during inference
- Identifier columns are excluded from learning

### 6.4. Pre-Encoding Preparation

#### 6.4.1 Identifier Column Removal

Certain columns (such as customer IDs) are:

- Unique identifiers
- Non-informative for prediction

These columns are removed before encoding to:

- Prevent noise
- Avoid misleading patterns
- Improve model generalization

#### 6.4.2 Categorical Feature Identification

Categorical columns are automatically detected based on data type:

- Object
- Category

This dynamic identification ensures:

- Flexibility across datasets
- Reduced manual errors

#### 6.4.3. Encoding Strategies Used

Different encoding techniques are applied based on the nature of categorical features.

## 6.5. Binary Categorical Encoding (Label Encoding)

### Description

Binary categorical features (those with exactly two unique values) are encoded using label encoding.

Example:

- Yes / No
- Male / Female

### Why This Method Is Used

- Simple and efficient
- No artificial ordering issues for binary variables
- Minimal memory footprint

### Handling Unseen Categories

If an unseen category appears in test or production data:

- It is assigned a safe default value
- Prevents system failure

### Industry View

Label encoding for binary features is a widely accepted industry standard.

## 6.6. Multi-Class Categorical Encoding (Frequency Encoding)

### Description

Multi-class categorical features are encoded using **frequency encoding**, where each category is replaced with its occurrence frequency in the training data.

### Why Frequency Encoding Is Used

- Avoids high dimensionality caused by one-hot encoding
- Prevents curse of dimensionality
- Maintains statistical relevance

### Advantages

- Scales well for high-cardinality features
- Preserves relative importance of categories
- Efficient for large datasets



## Handling Unseen Categories

- Unseen categories are assigned a frequency of zero
- Ensures robustness during deployment

## Industry View

Frequency encoding is commonly used in telecom, fintech, and large-scale ML systems.

## Data Leakage Prevention

All encoding rules are:

- Learned strictly from training data
- Applied consistently to test data

This prevents:

- Inflated performance metrics
- Unrealistic evaluation
- Deployment mismatches

This approach aligns with production ML pipeline standards.

## 6.7. Encoder Artifact Management

The learned encoders are saved as reusable artifacts.

### Why This Is Important

- Ensures consistent encoding during inference
- Enables model reproducibility
- Supports CI/CD and MLOps workflows

Encoders can be reused without retraining, which is essential for real-time systems.

## 6.8. Dataset After Encoding

After encoding:

- All categorical features are numeric
- Dataset shape remains consistent
- No loss of records
- Dataset is fully model-ready

The encoded dataset can now safely be used for:

- Feature scaling
- Model training
- Model evaluation
- Deployment

## 6.9. Comparison With Other Encoding Techniques

Encoding Method	Reason Not Selected
One-Hot Encoding	High dimensionality, memory expensive
Target Encoding	Risk of data leakage
Ordinal Encoding	Introduces false order
Hash Encoding	Reduced interpretability

The selected methods provide the best balance of performance, safety, and scalability.

## 6.10. Final Selected Encoding Strategy

### Final Strategy

- **Binary features:** Label Encoding
- **Multi-class features:** Frequency Encoding

### Why This Is the Best Choice

- Production-safe
- Handles unseen categories
- Scales efficiently
- Prevents data leakage
- Industry-aligned

## 6.11. Conclusion

This feature encoding module ensures:

- Reliable numerical representation of categorical data
- Robustness during model training and deployment
- Compliance with industry ML standards

By carefully selecting encoding techniques based on feature characteristics, the pipeline achieves accuracy, stability, and scalability.

## **7. Feature Selection Documentation**

### **7.1. Introduction**

Feature selection is the process of identifying the most relevant features that contribute significantly to predicting the target variable. In real-world machine learning systems, datasets often contain redundant, noisy, or irrelevant features that can negatively affect model performance.

This module implements a robust, multi-method feature selection framework focused on numerical features, following industry best practices to improve:

- Model accuracy
- Generalization
- Training efficiency
- Interpretability

### **7.2. Why Feature Selection Is Required**

#### **7.2.1. Machine Learning Perspective**

Using all available features can:

- Increase overfitting
- Add noise to the model
- Increase computational cost
- Reduce model stability

Many algorithms perform better when trained on a compact and informative feature set.

#### **7.2.2. Business Perspective**

From a business standpoint:

- Too many features make models harder to explain
- Irrelevant features reduce stakeholder trust
- Lean models are easier to deploy and maintain

Feature selection helps create efficient and explainable models suitable for production environments.

#### **7.2.3. Feature Type Identification**

The dataset is first divided into:

- Numerical features

- Categorical features

Only numerical features are considered in this module because:

- The applied statistical and model-based techniques require numeric input
- Categorical features are handled separately during encoding

This separation ensures methodological correctness.

### **7.3. Numerical Data Preparation**

After identifying numerical columns:

- Training and testing datasets are prepared using only numerical features
- A voting structure is initialized to track feature importance across methods

This setup allows multiple independent techniques to collectively influence the final decision.

### **7.4. Feature Selection Techniques Used**

To ensure robustness, multiple complementary techniques are applied.

#### **7.4.1 Constant and Quasi-Constant Feature Removal**

##### **Description**

Features with very low variance carry little to no information and do not contribute meaningfully to prediction.

##### **Why This Is Done**

- Removes features with near-constant values
- Reduces noise
- Simplifies the model

##### **Industry Relevance**

This is a standard initial screening step in enterprise ML pipelines.

#### **7.4.2. Variance Threshold Method**

##### **Description**

Features with variance below a defined threshold are removed.

### **Why This Is Done**

- Low-variance features rarely influence predictions
- Helps reduce dimensionality early

### **Strengths**

- Fast and scalable
- Model-agnostic

## **7.4.3. Correlation-Based Feature Elimination**

### **Description**

Highly correlated features provide redundant information. One of them can be safely removed.

### **Why This Is Done**

- Reduces multicollinearity
- Improves model stability
- Prevents coefficient inflation in linear models

### **Industry View**

This method is widely used in finance, telecom, and healthcare analytics.

## **7.4.4. ANOVA F-Test**

### **Description**

Measures the statistical relationship between each feature and the target variable.

### **Why This Is Done**

- Identifies features with strong discriminatory power
- Focuses on relevance to the target

### **Strengths**

- Statistically sound
- Effective for classification problems

## **7.4.5. LASSO-Based Feature Selection**

### **Description**

LASSO applies L1 regularization, shrinking less important feature coefficients to zero.

### **Why This Is Done**

- Automatically performs feature selection
- Balances model complexity and performance

### **Strengths**

- Embedded method
- Strong theoretical foundation

### **Limitation**

- Sensitive to scaling (handled during implementation)

## **7.4.6. Tree-Based Feature Importance**

### **Description**

Decision Trees and Random Forests are used to identify features based on split importance.

### **Why This Is Done**

- Captures nonlinear relationships
- Works well with complex feature interactions

### **Additional Robustness**

- Feature importance is also evaluated using robust scaling
- Ensures stability against outliers

### **Industry View**

Tree-based methods are highly trusted in enterprise ML systems.

## **7.5. Voting-Based Feature Selection Strategy**

### **Why Voting Is Used**

No single feature selection technique is universally optimal. Each method has its own bias and strengths.

### **Voting Mechanism**

- Each technique votes for selected features
- Features selected by multiple methods gain higher importance
- Final features are chosen based on highest votes

This ensemble-style approach:

- Reduces method-specific bias
- Increases reliability
- Reflects consensus importance

## 7.6. Final Feature Selection

### Selection Criteria

- Features with the highest total votes across all methods
- Focus on consistency rather than single-method dominance

### Outcome

- A reduced and optimized numerical feature set
- Same feature subset applied to both training and testing data

This ensures:

- No data leakage
- Consistent model behavior

## 7.7. Why This Approach Is Industry-Oriented

- Combines statistical, linear, and tree-based methods
- Prevents over-reliance on a single technique
- Produces explainable and stable results
- Scales well to large datasets

This mirrors how feature selection is handled **in** real-world enterprise ML systems.

## 7.8. Comparison With Simpler Approaches

Approach	Limitation
Single-method selection	Biased results
Manual selection	Subjective and error-prone
Model-only importance	Overfitting risk
PCA	Loss of interpretability

The chosen multi-method voting approach offers the best balance of performance and interpretability.

## **7.9. Conclusion**

This feature selection module ensures:

- Removal of redundant and noisy features
- Improved model performance and stability
- Reduced computational cost
- Clear and explainable feature choices

By using multiple complementary techniques and a voting mechanism, the pipeline achieves robust, production-ready feature selection.

## **8. Model Training & Evaluation**

### **8.1. Introduction**

Model training is the core phase of a machine learning project where cleaned and engineered features are used to learn patterns that predict the target variable. In this project, multiple classification models are trained, evaluated, and compared using standardized performance metrics to select the best-performing and most reliable model for deployment.

This module follows an industry-grade model evaluation framework, ensuring fairness, reproducibility, and robustness.

### **8.2. Objective of Model Training**

The objectives of this module are:

- Train multiple classification algorithms
- Evaluate models using consistent metrics
- Compare model performance visually and quantitatively
- Select the best model based on objective criteria
- Save trained artifacts for deployment

### **8.3. Data Preparation Before Training**

Before training any model:

- Only numerical features are used
- Infinite values are removed
- Missing values are safely handled



This ensures:

- Model compatibility
- Stable numerical computation
- Reduced runtime errors

This preprocessing step mirrors production inference conditions.

## **8.4. Models Trained**

Multiple industry-standard classification algorithms are trained to ensure a comprehensive comparison.

### **8.4.1 Logistic Regression**

**Purpose:**

Acts as a strong baseline model and provides interpretability.

**Strengths:**

- Simple and fast
- Interpretable coefficients
- Works well with standardized data

**Use Case:**

Preferred when model explainability is critical.

### **8.4.2. K-Nearest Neighbors (KNN)**

**Purpose:**

Captures local data patterns using distance-based learning.

**Strengths:**

- Non-parametric
- Simple logic

**Limitations:**

- Sensitive to scaling and noise
- Computationally expensive for large datasets

### 8.4.3 Naive Bayes

**Purpose:**

Probabilistic classifier based on Bayes theorem.

**Strengths:**

- Fast training
- Works well with independent features

**Limitations:**

- Assumes feature independence
- Limited performance on complex patterns

### 8.4.4 . Decision Tree

**Purpose:**

Captures nonlinear relationships using rule-based splits.

**Strengths:**

- Easy to interpret
- Handles mixed feature types

**Limitations:**

- Prone to overfitting

### 8.4.5. Random Forest

**Purpose:**

An ensemble of decision trees to improve generalization.

**Strengths:**

- High accuracy
- Robust to noise
- Handles nonlinear relationships well

**Industry Usage:**

Widely used in production ML systems.

### 8.4.6. AdaBoost

**Purpose:**

Boosting-based ensemble model that focuses on difficult samples.

**Strengths:**

- Improves weak learners
- Effective on structured data

**Limitations:**

- Sensitive to noisy data

### 8.4.7. XGBoost

**Purpose:**

Advanced gradient boosting algorithm optimized for performance.

**Strengths:**

- Excellent predictive power
- Handles feature interactions
- Highly optimized

**Industry Usage:**

Frequently used in competitive and enterprise ML solutions.

## 8.5. Model Evaluation Metrics

To ensure fair comparison, all models are evaluated using the same metrics.

### 8.5.1 Accuracy

Measures the proportion of correctly predicted instances.

**Limitation:**

Can be misleading for imbalanced datasets.

### 8.5.2 ROC-AUC Score

Measures the model's ability to distinguish between classes across all thresholds.

**Why It Is Preferred:**

- Threshold-independent

- Robust for imbalanced classification problems

### 8.5.3 Classification Report

Provides:

- Precision
- Recall
- F1-score

Used for deeper performance understanding.

### 8.6. Model Comparison Strategy

- All models are trained on the same training data
- Tested on the same unseen test data
- ROC curves are plotted together for visual comparison

This ensures:

- Objective evaluation
- No bias across models

### 8.7. Best Model Selection

#### Selection Criterion

The model with **the** highest ROC-AUC score is selected as the best model.

#### Why ROC-AUC Is Used

- Handles class imbalance well
- Reflects true predictive ability
- Commonly used in industry churn and risk models

### 8.8. Final Selected Model

🏆 **Best Model:** XGBoost Classifier

#### Reasons for Selection:

- Highest ROC-AUC score among all models
- Strong handling of nonlinear relationships
- Robust performance on complex feature interactions
- Proven industry reliability

While simpler models offer interpretability, XGBoost provides the best balance of accuracy and robustness for this problem.

## 8.9. Model Artifact Management

The following artifacts are saved:

- Best trained model
- Feature scaler (if applicable)

Why This Matters

- Enables consistent predictions during deployment
- Supports reproducibility
- Aligns with MLOps best practices

## 8.10. End-to-End Pipeline Summary

The complete pipeline follows this structure:

Data Cleaning → Missing Value Handling → Outlier Treatment → Feature Transformation → Encoding → Feature Selection → Model Training → Model Selection → Deployment

## 8.11. Conclusion

This model training framework ensures:

- Fair comparison across multiple algorithms
- Robust evaluation using industry-standard metrics
- Reliable selection of the best-performing model
- Deployment-ready artifacts

By using multiple models and objective selection criteria, the project achieves high performance, scalability, and production readiness.

# 9. Deployment Configuration

**File Name:** Procfile

## 9.1. Introduction

The Procfile is a deployment configuration file used to define how an application should be executed in a production environment. It is commonly used in cloud platforms such as Render, Heroku, and similar PaaS services.

This file tells the deployment platform:

- Which process to run
- How to start the application
- Which command acts as the entry point

Without a Procfile, the platform may not know how to correctly launch the application.

## 9.2. Purpose of the Procfile

The primary purpose of the Procfile is to **standardize application startup behavior** across environments.

It ensures that:

- The correct application file is executed
- The correct runtime command is used
- The service starts automatically after deployment

This is essential for continuous deployment and production stability.

## 9.3. Role in the ML Deployment Pipeline

In a machine learning project, the Procfile plays a critical role during deployment:

- Starts the web application (API or UI)
- Loads trained models and preprocessing artifacts
- Exposes prediction endpoints to users
- Enables real-time or batch inference

It acts as the bridge between trained ML models and end users.

## 9.4. Why Procfile Is Required in Industry Projects

### 9.4.1. Platform Compatibility

Cloud platforms require a clear definition of:

- Process type (web, worker, etc.)
- Startup command

The Procfile provides this information explicitly, avoiding ambiguity.

### 9.4.2. Automation & CI/CD

The Procfile enables:

- Automated builds
- Repeatable deployments
- Zero manual intervention during release

This aligns with DevOps and MLOps best practices.

### **9.4.3. Reliability & Consistency**

Using a Procfile ensures:

- Same startup behavior in staging and production
- Reduced deployment failures
- Faster rollback and redeployment

### **9.4.4. Typical Usage in This Project**

In this project, the Procfile is used to:

- Launch the machine learning application server
- Ensure the trained model is loaded at startup
- Make the prediction service accessible via the web

This allows the churn prediction system to be:

- Deployed on a cloud platform
- Accessible through a browser or API
- Scalable and production-ready

### **9.4.5. Relationship With Other Deployment Files**

The Procfile works together with:

- Requirements file (dependency management)
- Application entry file (API or UI)
- Saved model artifacts
- Environment variables

Together, these components form a complete deployment setup.

### **9.4.6. Common Industry Scenarios Where Procfile Is Used**

- Deploying ML models as REST APIs
- Hosting Streamlit or Flask applications
- Running background workers for batch prediction
- Managing multiple services in a single application

## 10. Conclusion

Overall, the analysis shows that customer tenure has the strongest impact on churn, with new customers being far more likely to leave than long-term customers. Payment method also plays a significant role, as customers using electronic/manual payment methods exhibit higher churn compared to those on automatic payments. Additionally, while SIM operator affects churn in absolute numbers (e.g., higher churn for operators with larger user bases), it is tenure and payment convenience that most strongly drive churn behavior rather than the operator itself.

**Live deployment link:** <https://customer-churn-prediction-y73h.onrender.com>