

Telco Customer Churn Analytics & Prediction

Comprehensive Project Report

Executive Summary

This project delivers an end-to-end **data analytics and machine learning solution** to understand, analyse, and predict **customer churn** in the telecommunications industry. Using a dataset of **7,043 customers**, the project applies **exploratory data analysis (EDA)**, **statistical hypothesis testing**, and **machine learning modelling** to uncover churn drivers and predict customer churn risk.

A **Random Forest classification model** achieved **79.1% accuracy** and **0.83 ROC-AUC**, demonstrating strong predictive performance. The solution is deployed as a **multi-page interactive Streamlit dashboard**, enabling **real-time churn prediction, business insights, and decision support** for stakeholders.

1. Project Overview

1.1 Project Objectives

Primary Objectives

1. **Understand Churn Drivers**
Identify customer attributes, services, and pricing factors associated with churn.
2. **Build Predictive Model**
Develop a machine learning model to estimate churn probability for individual customers.
3. **Support Business Decisions**
Translate analytical findings into actionable customer retention strategies.

Target Audience

- Business managers and retention teams
 - Data analysts and BI professionals
 - Customer service leadership
 - Executive stakeholders
-

1.2 Business Problem Statement

Challenge

The telecommunications company experiences a **26.5% customer churn rate**, leading to revenue loss, increased acquisition costs, and operational instability.

Solution Approach

- Analyse historical customer behaviour
 - Identify statistically significant churn predictors
 - Predict at-risk customers before churn occurs
 - Provide insights to support proactive retention strategies
-

2. Data Analytics Methodology

2.1 Data Collection & Preparation

Dataset

- Name: *WA_Fn-UseC_-Telco-Customer-Churn.csv*
- Records: 7,043 customers
- Features: 21 original variables
- Target Variable: **Churn** (Yes / No)

Data Cleaning Process

Step	Action	Result
1	Identified data types	18 categorical, 3 numerical
2	Cleaned <code>TotalCharges</code>	Converted from text, removed 11 blank values
3	Removed non-predictive feature	Dropped <code>customerID</code>
4	Final dataset	7,032 rows × 20 features
5	Data quality check	100% complete dataset

Data Validation

- No duplicate records
- Correct data types
- No missing values in final dataset
- Target variable properly encoded

2.2 Exploratory Data Analysis

Key Dataset Statistics

Metric	Value
Total Customers	7,043
Average Tenure	32.4 months
Average Monthly Charges	\$64.80
Churn Rate	26.54%
Senior Citizens	16.2%

Feature Categories

- **Demographics:** Gender, Senior Citizen, Partner, Dependents
 - **Services:** Internet, Security, Backup, Tech Support, Streaming
 - **Account & Billing:** Tenure, Contract, Payment Method, Charges
-

3. Hypothesis Testing & Statistical Analysis

3.1 Hypothesis 1 – Contract Type Affects Churn

Hypothesis

Customers on longer contracts have lower churn rates.

Test Used

Chi-Square Test of Independence

Results

- Chi-square statistic: **1179.55**
- P-value: **7.33e-257**

Conclusion

 **VALIDATED**

Contract Type	Churn Rate	Customer s
Month-to-month	42.71%	3,875
One year	11.28%	1,472
Two year	2.85%	1,685

Insight

Month-to-month customers are **~15x more likely to churn** than two-year contract customers.

3.2 Hypothesis 2 – Tenure Inversely Affects Churn

Test Used

Independent Samples T-Test

Results

- T-statistic: **31.74**
- P-value: **9.44e-207**
- Cohen's d: **0.857 (large effect)**

Conclusion



VALIDATED

Status	Avg Tenure
Non-Churned	37.65 months
Churned	17.98 months

Insight

Churned customers leave nearly **20 months earlier**, highlighting the importance of early retention.

3.3 Hypothesis 3 – Internet Service Type Affects Churn

Test Used

One-Way ANOVA

Results

- F-statistic: **406.29**
- P-value: **1.06e-167**

Conclusion



Internet Service Churn Rate

Fiber optic	41.89%
DSL	19.00%
No internet	7.43%

Insight

Fiber optic customers churn at more than **2x the rate** of DSL users.

3.4 Hypothesis 4 – Monthly Charges Correlate with Churn

Test Used

Pearson Correlation

Results

- Pearson r: **0.193**
- P-value: **6.76e-60**

Conclusion

 **VALIDATED**

Insight

Customers with higher monthly charges are more likely to churn, indicating price sensitivity.

4. Machine Learning Model Development

4.1 Model Architecture

- Algorithm: **Random Forest Classifier**
- Trees: 100
- Max Depth: 15
- Min Samples Split: 10

Rationale

- Handles mixed data types
 - Captures non-linear relationships
 - Provides feature importance for business interpretation
-

4.2 Data Preparation

- Categorical encoding applied to 15 features
 - Numerical features scaled
 - Final features: **19 predictive variables**
 - Train/Test Split: **80% / 20% (stratified)**
-

4.3 Model Performance

Metric	Score
Accuracy	79.10%
Precision	63.89%
Recall	49.20%
F1-Score	0.556
ROC-AUC	0.832

Confusion Matrix Highlights

- True Negatives: 929
- True Positives: 184
- False Negatives: 190

4.4 Cross-Validation

- Mean ROC-AUC: **0.8436**
- Std Dev: **0.0141**

 **Model stability confirmed**

4.5 Feature Importance

Top Predictors

1. Tenure (16.8%)
 2. Total Charges (14.5%)
 3. Contract Type (14.5%)
 4. Monthly Charges (13.3%)
 5. Online Security (8.3%)
-

5. Streamlit Dashboard Development

5.1 Dashboard Architecture

Multi-page Streamlit application:

- **Home**
- **Analytics**
- **Feature Importance**
- **Prediction**

- **Model Performance**
 - **About**
-

5.2 Key Features

- Interactive EDA visualisations
 - Real-time churn prediction
 - Business-friendly KPIs
 - Model performance metrics
 - Clear storytelling for non-technical users
-

6. Key Insights & Business Recommendations

Critical Findings

- Contract type is the strongest churn driver
- Early-stage customers are most at risk
- Fiber optic service shows high churn
- Pricing sensitivity impacts churn

Strategic Actions

- Incentivise long-term contracts
- Improve fiber service quality
- Strengthen onboarding (0–12 months)

- Introduce value-based pricing bundles
-

7. Model Strengths & Limitations

Strengths

- High accuracy and ROC-AUC
- Stable cross-validation performance
- Interpretable feature importance

Limitations

- Moderate recall for churn class
 - Class imbalance
 - No temporal features
-

8. Ethical Considerations & Data Governance

- Dataset anonymised
 - No discriminatory attributes used
 - Predictions support service improvement
 - GDPR-aligned analytical practices
-

9. Deployment & Maintenance

Current Status

- Fully functional local Streamlit dashboard
- Real-time predictions enabled

Future Deployment

- Streamlit Cloud / Heroku
 - Automated retraining pipeline
 - Monitoring for data drift
-

10. Project Outcomes & Metrics

- 4 hypotheses validated
 - 79.1% ML accuracy achieved
 - 5-page dashboard delivered
 - Business-ready insights generated
-

11. Learnings & Future Development

Key Learnings

- Importance of interpretability
- Business framing of ML outputs
- Dashboard-driven storytelling

Future Enhancements

- SHAP explainability

- Time-series churn modelling
 - Retention recommendation engine
-

12. Conclusion

This project successfully demonstrates the **full data analytics lifecycle**—from data exploration and hypothesis testing to machine learning deployment and business storytelling.

Recommendation:

Proceed with production deployment and pilot retention strategies targeting high-risk customers. Expected churn reduction: **10–15% within 6 months**.

Appendix: Project Artifacts

- 5 Jupyter notebooks
- Trained ML model
- Multi-page Streamlit dashboard
- Business-ready documentation

Technologies Used

Python, Pandas, NumPy, Scikit-learn, Matplotlib, Streamlit, Joblib

Project Status:

Complete & Production-Ready