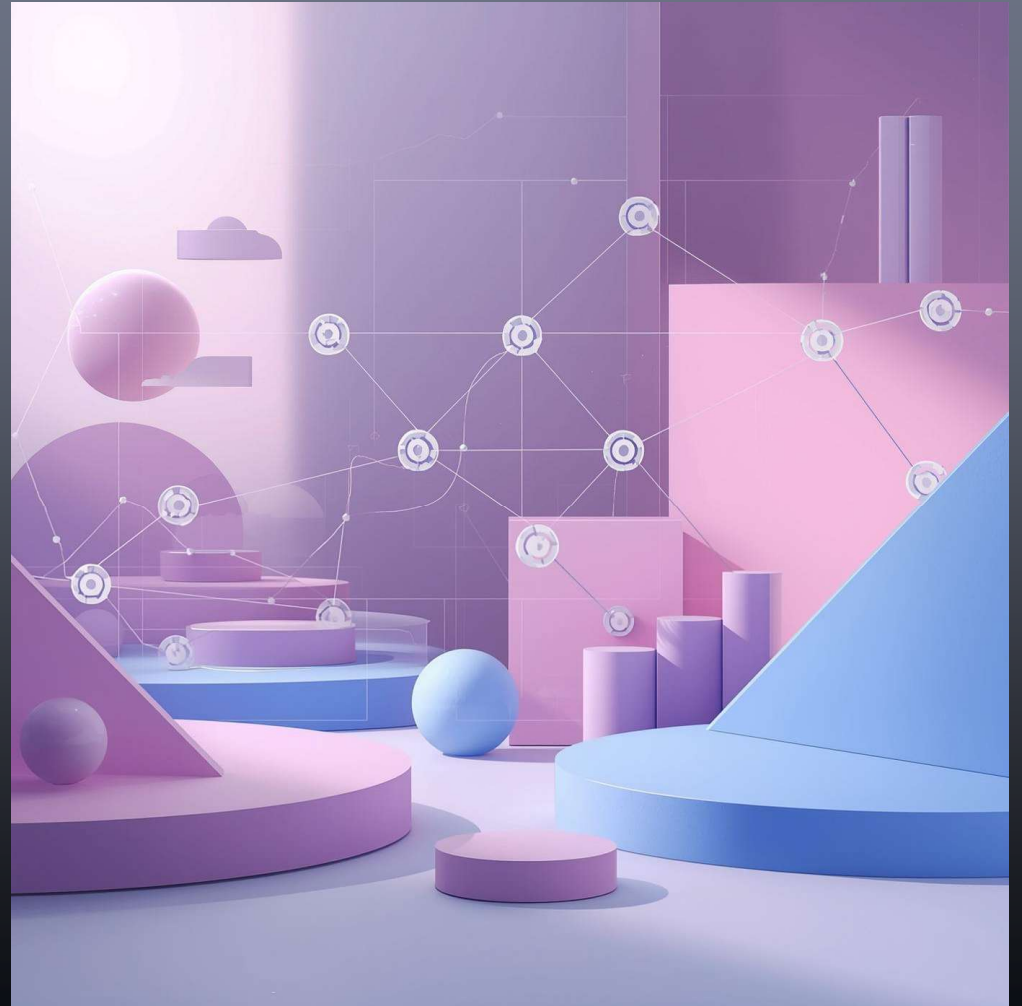


Telecom Churn Prediction

Predicting customer churn using machine learning

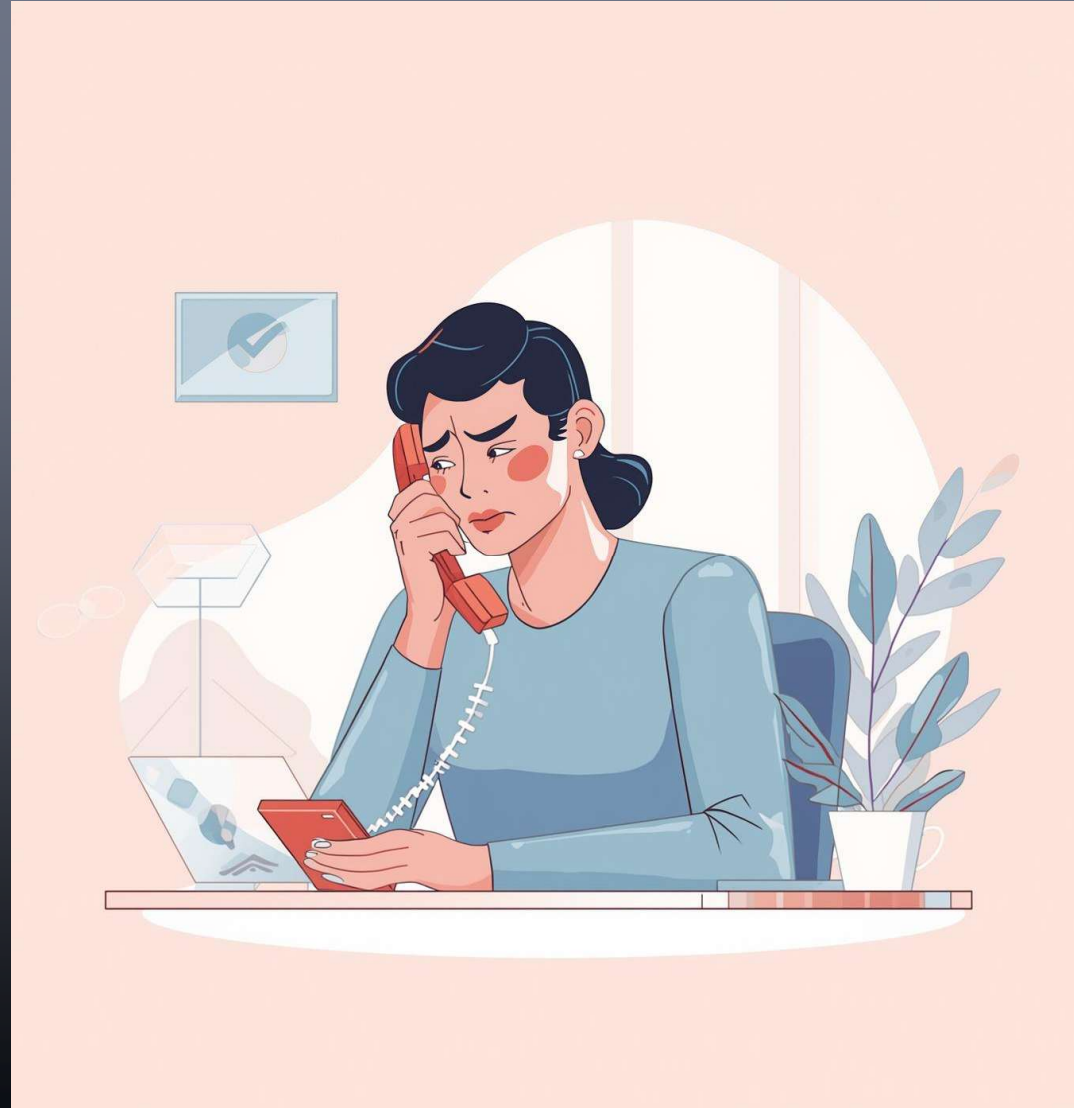
Presented by:

Vrushali Oak
Data Scientist



Understanding Churn

Customer churn refers where clients discontinue their service. Understanding churn is essential as it impacts revenue and growth. Identifying the reasons behind customer dissatisfaction helps businesses **improve services** and develop effective retention strategies.



Business Problem

- Customer churn leads to significant revenue loss in telecom
- Retaining existing customers is cheaper than acquiring new ones
- Objective: Predict customers likely to churn so retention teams can act early

Dataset Overview

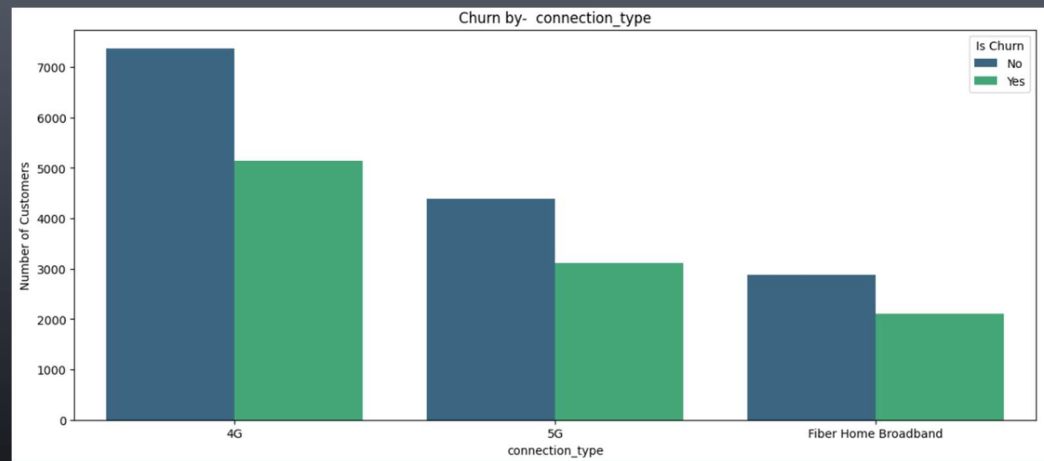
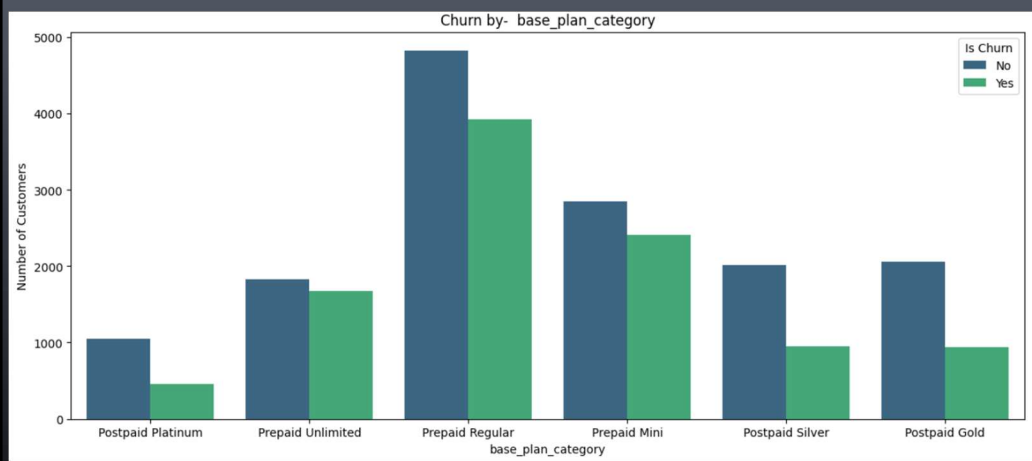
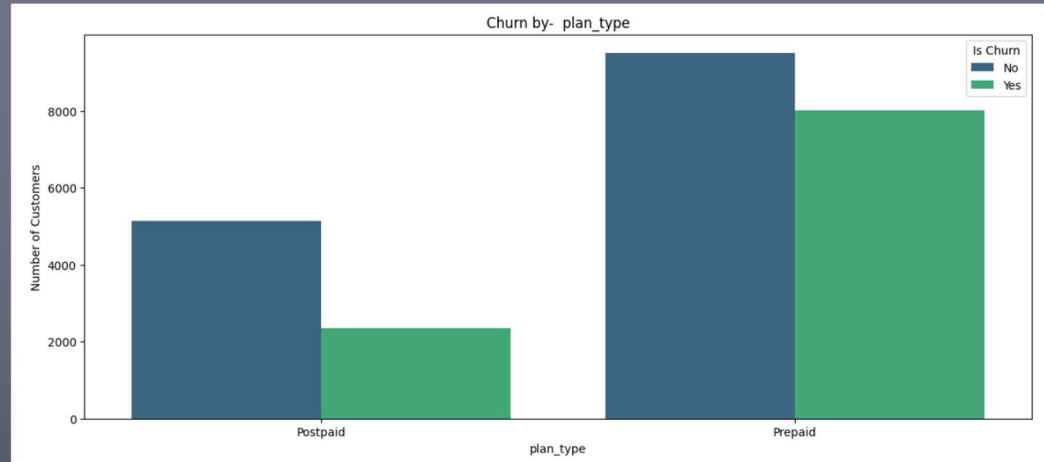
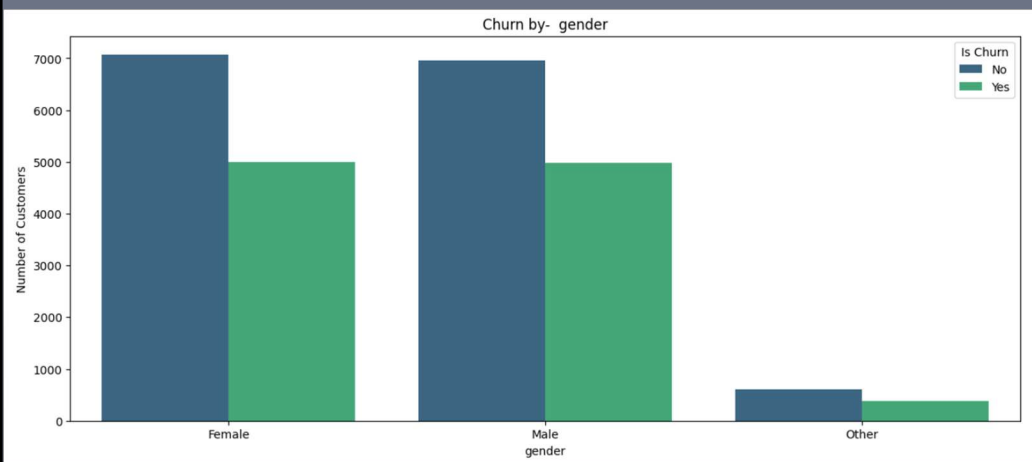
- 25,000 customers, 36 features
- Includes numerical, binary, and categorical variables
- Target variable: is_churn (1 = churned, 0 = retained)
- Churn distribution:
 - 58.6% retained
 - 41.4% churned
 - This is moderate imbalance, not extreme

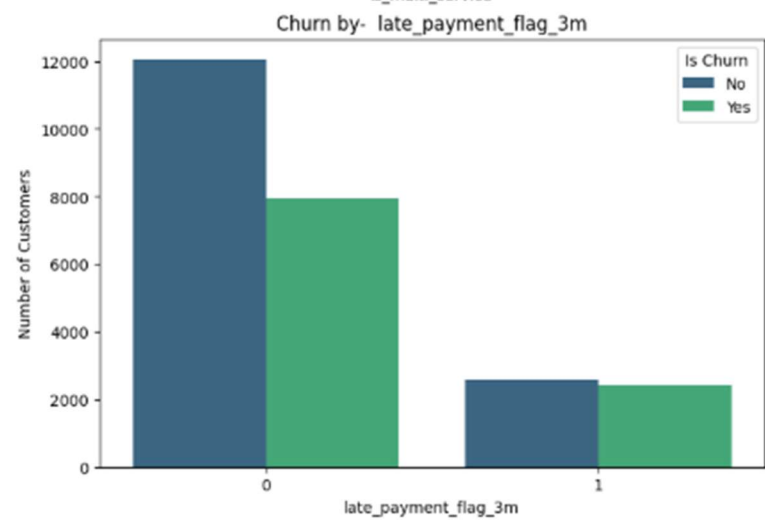
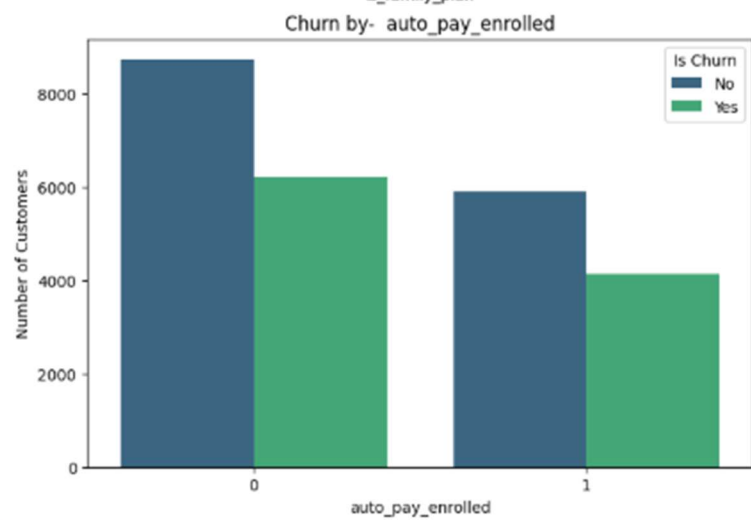
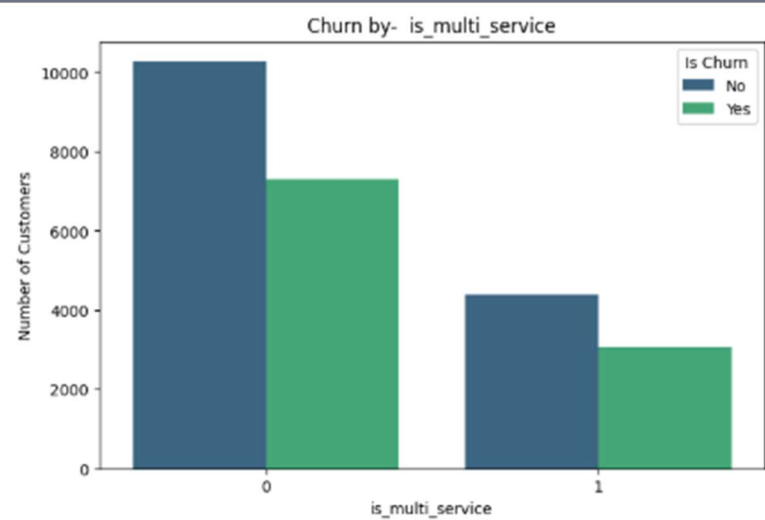
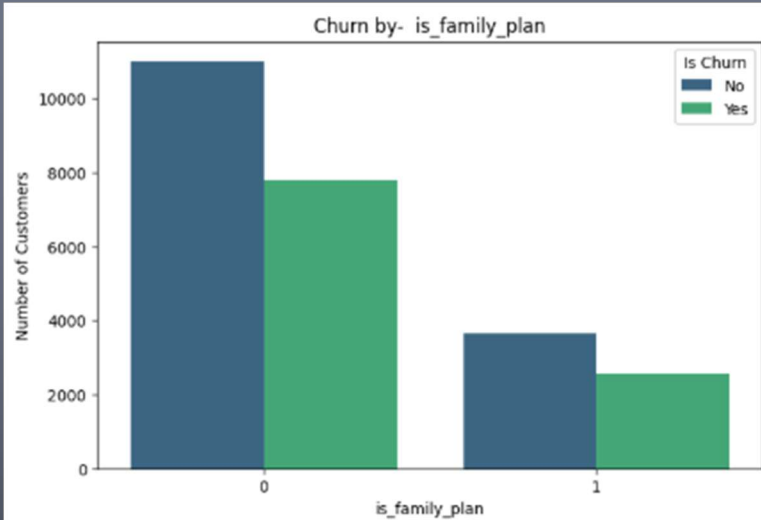
Data Understanding & Cleaning

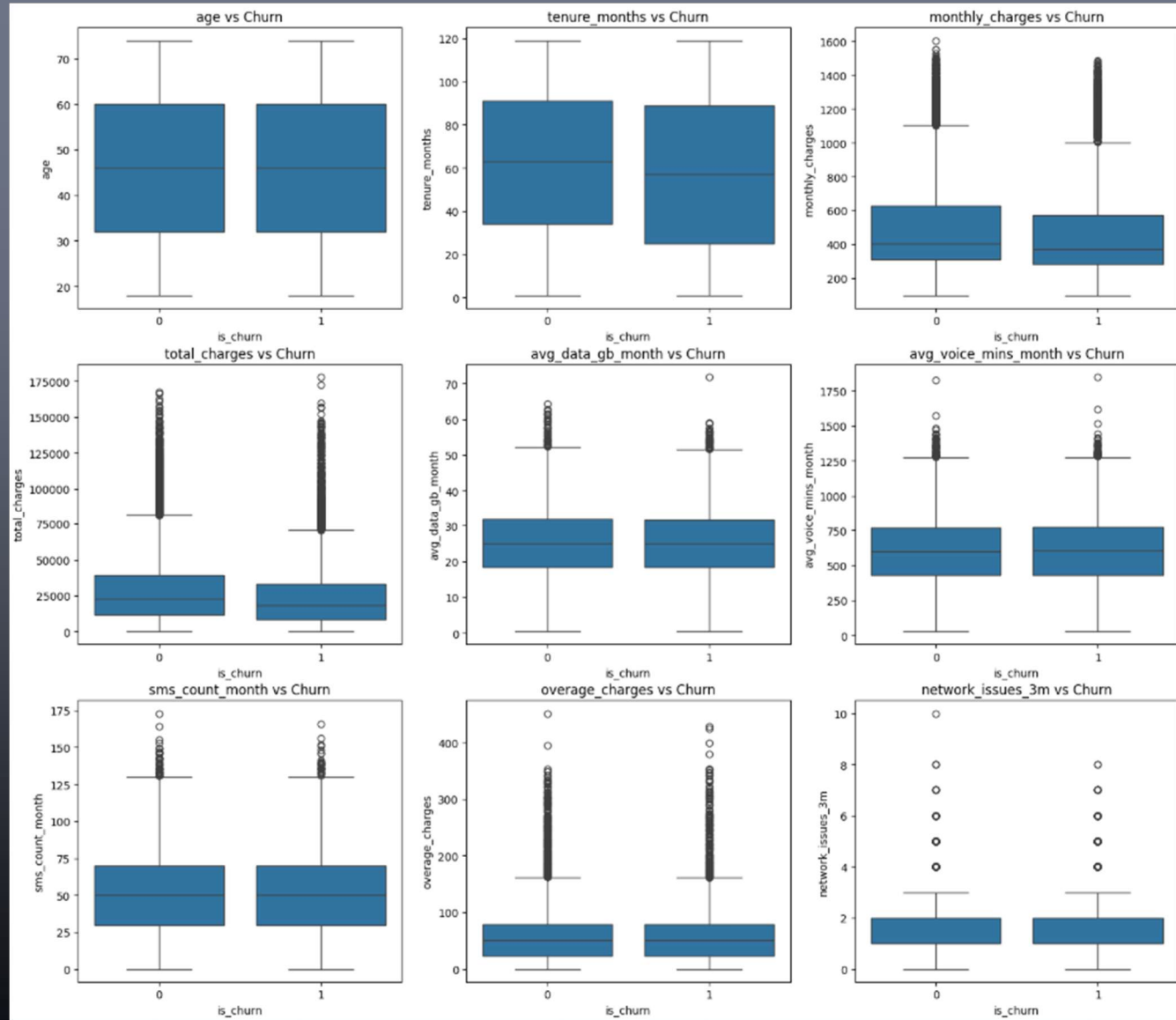
- Removed non-informative and leaky columns:
 - customer_id (identifier)
- Corrected data types
- Verified missing values and distributions- **No missing values** and **no duplicate records** after validation
- Separated numerical, categorical and binary features

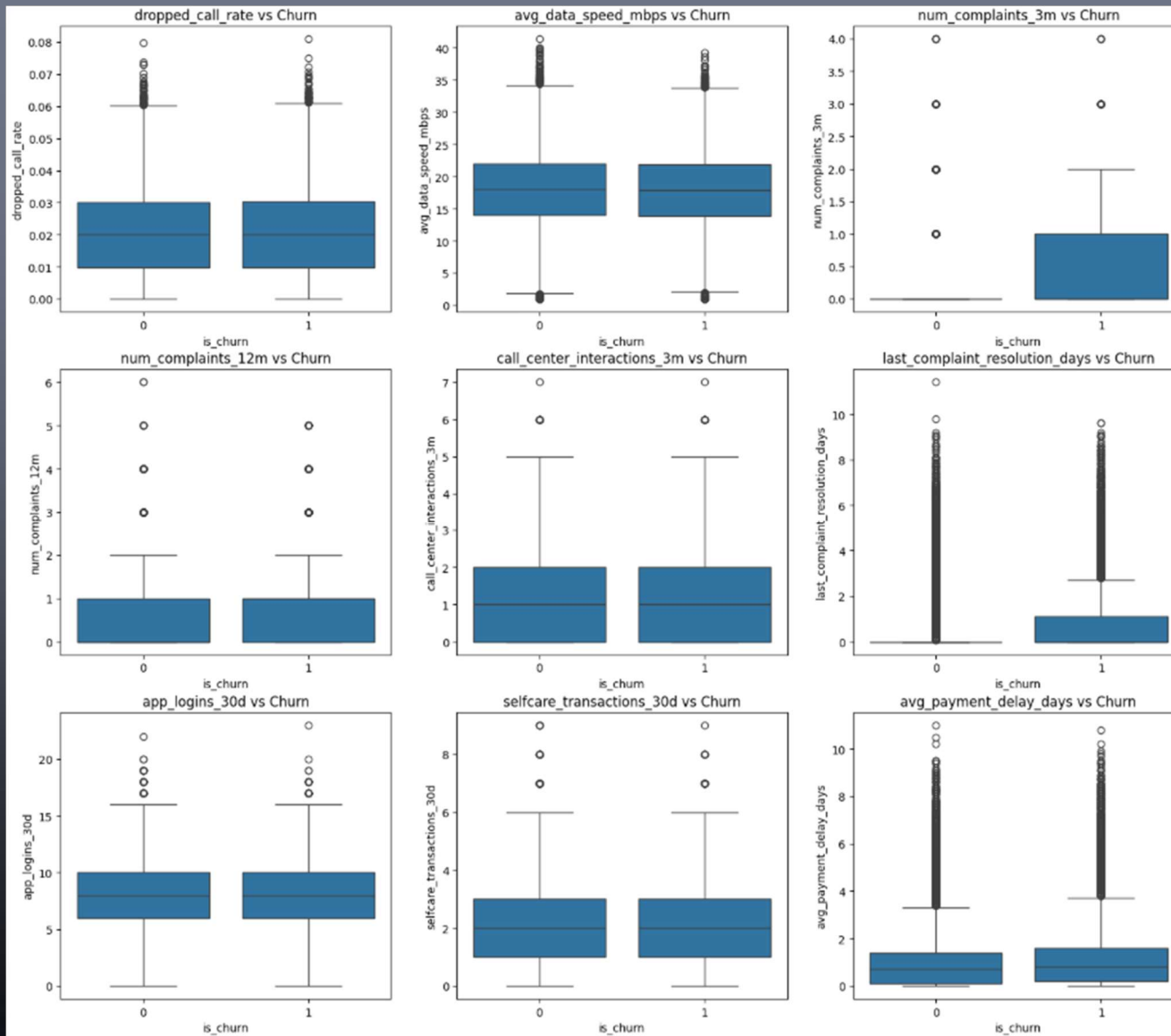
Exploratory Data Analysis (EDA)

- Performed correlation heatmap on continuous numeric features only
- Identified:
High correlation between monthly_charges and arpu
Redundant pricing variables
- Dropped columns -
'monthly_charges',
'num_complaints_3m',
'total_charges'
- Used
Bar plots for- Binary and Categorical features
Box plots for- Numerical features









Key Insights from EDA

- High overage charges → higher churn probability
- Frequent complaints and call center interactions → churn risk
- Payment delays strongly associated with churn
- Service quality issues (low speed, dropped calls) impact churn

Feature Selection Decisions

- Dropped 'monthly_charges' due to high correlation with arpu
- Kept arpu as it reflects actual realized customer cost
- Dropped cumulative and leakage-prone features (monthly_charges', 'num_complaints_3m')
- Retained behavior-driven features impacting churn

Feature Engineering

- **Created High Overage Risk Flag to capture billing shock behavior**
- **Engineered behavioral signals instead of raw extremes**
- **Focused on business meaning, not just statistical transformation**

Outlier Handling Strategy

- Did NOT blindly remove outliers
- Applied soft capping only where outliers likely represent noise:
 - Billing spikes
 - Measurement errors
- Preserved extreme values where they represent real churn behavior

Data Preparation

- One-hot encoded categorical variables
- Feature scaling applied for Logistic Regression
- Stratified train-test split to maintain churn ratio

Models Used

- Logistic Regression (baseline & interpretable)
- Random Forest (non-linear, ensemble)
- XGBoost (gradient boosting)
- Purpose:
Compare simplicity vs performance vs generalization

Evaluation Metrics

- ROC-AUC (primary metric – ranking churn risk)
- Precision, Recall, F1-score
- Train vs Test performance to detect overfitting

Model Performance

	Model	ROC-AUC
0	Logistic Regression (unscaled)	0.647326
1	Logistic Regression (scaled)	0.647234
2	Random Forest	0.647614
3	XGBoost (baseline)	0.635404
4	XGBoost (tuned)	0.654467

Model Performance Comparison

- Logistic Regression: stable but limited non-linear learning
- Random Forest: better than LR, moderate generalization
- XGBoost: highest CV and test ROC-AUC

XGBoost performed best in ranking high-risk churn customers

Why XGBoost Performed Best


- Captures non-linear churn behavior
- Learns feature interactions automatically
- Focuses on hard-to-predict churn cases
- Handles mixed feature types and imbalance effectively

Final Model Selection

- Selected XGBoost as primary model
- Logistic Regression retained as interpretable baseline
- Acknowledged mild overfitting and discussed regularization improvements

Business Impact

- Customers ranked by churn probability
- Enables targeted retention campaigns
- Improves marketing efficiency
- Reduces revenue leakage due to churn

A white sunburst graphic with many thin lines radiating from a central point, positioned behind the word 'Conclusion'.

Conclusion

- Built an end-to-end churn prediction pipeline
- Combined business understanding with ML modeling
- Delivered a practical, deployable churn scoring solution

