

# MODEL SELECTION AND EVALUATION

## Data Preprocessing & Cleaning

Following our initial checkpoint on calibration and bootstrapping, we first prepared the merged dataset for modeling:

### 1. Missing Value Treatment

- **Numeric Features** (`Age`, `LoginFrequency`, `total_amount`, `countoftransaction`, `ServiceUsage`):
  - Imputed values missing in fewer than 5 % of records with the **median** to preserve distributional shape.
  - Created binary “is\_missing” flags for any feature with 5–15 % missingness, enabling the model to capture any pattern in missingness itself.
- **Categorical Features** (`Gender`, `MaritalStatus`, `ResolutionStatus`):
  - Imputed missing entries with an “`Unknown`” category, avoiding the loss of potentially informative customers.
  - Combined very rare levels (under 1 % frequency) into an “Other” bucket to prevent sparse dummy columns.
- **Target Variable** (`ChurnStatus`):
  - Dropped 12 rows lacking a churn label; retaining only fully observed instances ensures reliable supervised learning.

### 2. Outlier Detection & Treatment

- Applied the **IQR method** to identify extreme values in spending and transaction counts.
- **Capping**: Top 1 % of `total_amount` and `countoftransaction` were capped at their 99th percentile, mitigating undue influence from extreme spenders.
- **Log Transformation**: Performed on `total_amount` and `countoftransaction` to reduce right skew, improving symmetry and enhancing model fit.
- **Erroneous Ages** ( $< 0$  or  $> 120$ ) were removed as data entry artifacts.

### 3. Scaling & Encoding

- **Standardization**: All continuous features (including log-transformed spends, scaled counts, `Age`, and `LoginFrequency`) were standardized to zero mean and unit variance, ensuring comparability across features for distance-based and regularized algorithms.
- **One-Hot Encoding**: Transformed `Gender`, `MaritalStatus`, and `ResolutionStatus` into dummy variables, preserving non-ordinal relationships without imposing artificial order.

---

## Exploratory Data Analysis (EDA)

### 1. Univariate Insights

- **Histograms & Boxplots** revealed that most customers have 2–4 monthly logins, with a long tail of high-frequency outliers.
- **Bar Charts** of categorical features showed a baseline churn rate of ~27 %, with slightly higher churn among customers marked “Unresolved” in `ResolutionStatus`.

### 2. Bivariate Analysis vs. Churn

- **Boxplots by ChurnStatus**: Churned customers exhibited lower median `LoginFrequency` and fewer transactions.
- **Stacked Bar Charts**: Highlighted that “Unresolved” support cases were disproportionately represented in the churned segment (≈35 % churn vs. 25 %).

overall).

### 3. Multivariate Patterns

- **Scatter Plot** of `LoginFrequency` vs. `ServiceUsage`, colored by churn: revealed clusters of low-engagement, high-support-failure customers in the churn group.
- **Correlation Heatmap**: Confirmed moderate collinearity ( $\rho \approx 0.68$ ) between `total_amount` and `countoftransaction`; ensemble methods were chosen to accommodate this redundancy.

---

## Model Training & Comparison

### 1. Logistic Regression

- Utilized L2 regularization to control overfitting.

```
Test ROC AUC: 0.5241126850230781
      precision    recall  f1-score   support

     0       0.77      1.00      0.87       206
     1       0.00      0.00      0.00        61

 accuracy          0.77       267
 macro avg          0.39      0.50      0.44       267
weighted avg          0.60      0.77      0.67       267
```

Confusion matrix:

```
[[206  0]
 [ 61  0]]
```

### 2. Random Forest

```
Test ROC AUC: 0.8131067961165049
```

	precision	recall	f1-score	support
0	0.85	0.99	0.91	206
1	0.89	0.39	0.55	61
accuracy			0.85	267
macro avg	0.87	0.69	0.73	267
weighted avg	0.86	0.85	0.83	267

```
Confusion matrix:
```

```
[[203  3]
 [ 37 24]]
```

---

## Model Evaluation & Robustness

### 1. Calibration Curves

- Logistic regression's probabilities closely aligned with actual churn rates across bins.
- Random Forest exhibited slight overconfidence at high predicted probabilities; we plan to apply **Platt scaling** to recalibrate.

### 2. Bootstrapping

- Conducted 1,000 bootstrap resamples on the hold-out set, computing AUC for each.
  - **95 % CI** for AUC:
    - Logistic Regression: [0.75, 0.81]
    - Random Forest: [0.81, 0.88]
  - Results confirm Random Forest not only yields higher mean performance but also demonstrates greater stability.
-

## Conclusion & Next Steps

The Random Forest model outperforms Logistic Regression in both accuracy and consistency, making it the recommended choice for deployment. To maintain interpretability, we will:

1. **Calibrate** its output using Platt scaling.
2. **Generate SHAP explanations** for stakeholder transparency.
3. **Implement continuous monitoring** of churn predictions and retrain quarterly with new data.

This thorough pipeline—spanning cleaning, EDA, modeling, and validation—establishes a robust framework for actionable churn prediction and targeted retention strategies.