# Titanic EDA + Modeling Report

Includes feature engineering and two baseline models
(Logistic Regression, Random Forest)

Source: train.csv

# Dataset overview & missing values

Rows: 891  Columns: 17

Total missing values: 687

Top missing columns:

Cabin: 687

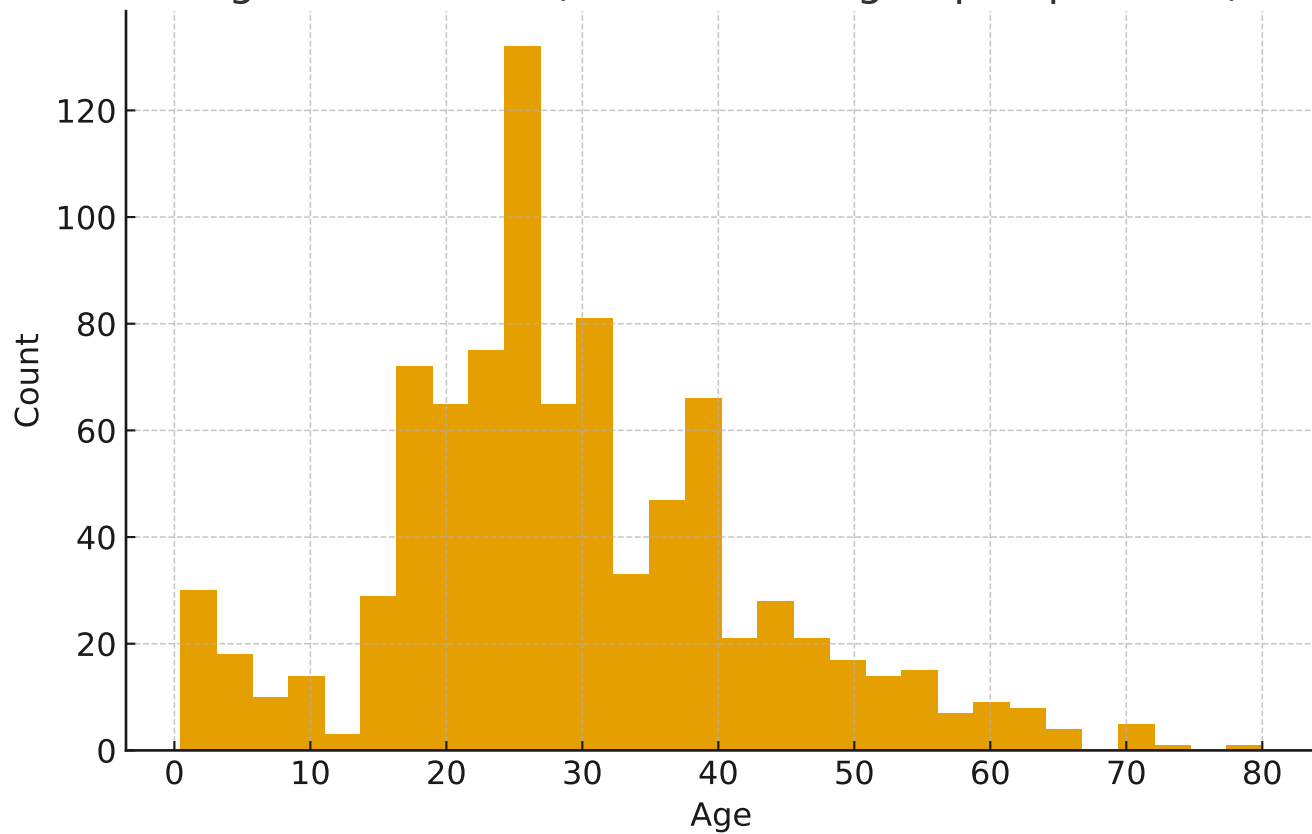PassengerId: 0

Fare: 0

IsAlone: 0
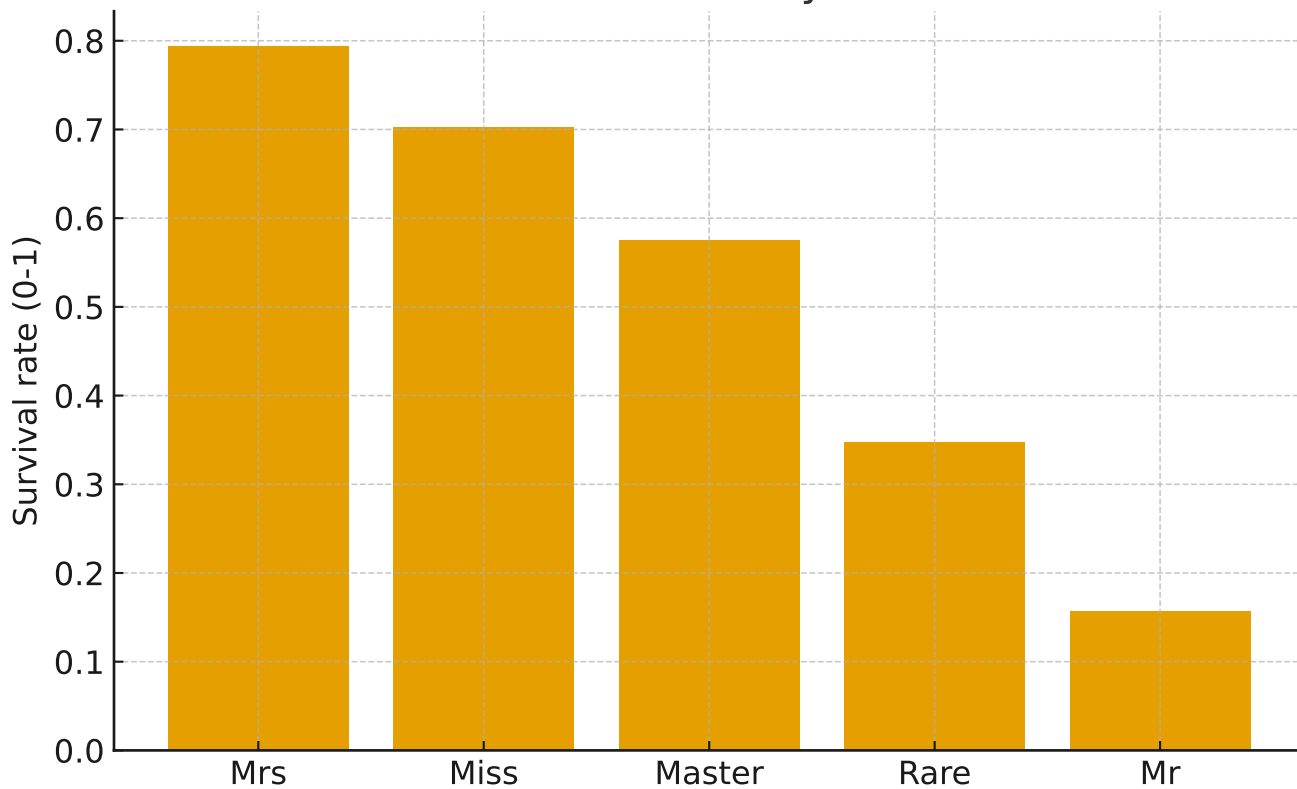
FamilySize: 0

Deck: 0

Title: 0
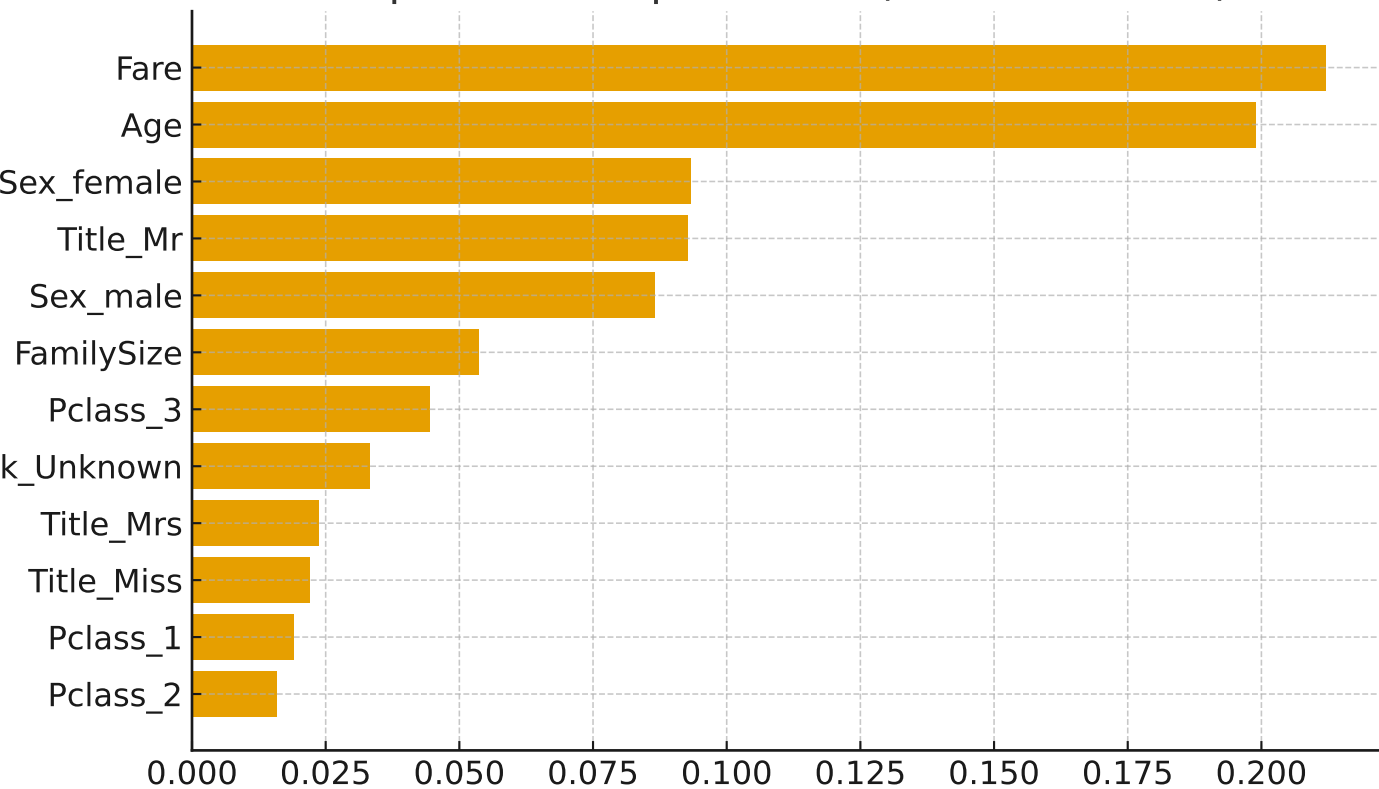
Embarked: 0

Ticket: 0

Survived: 0

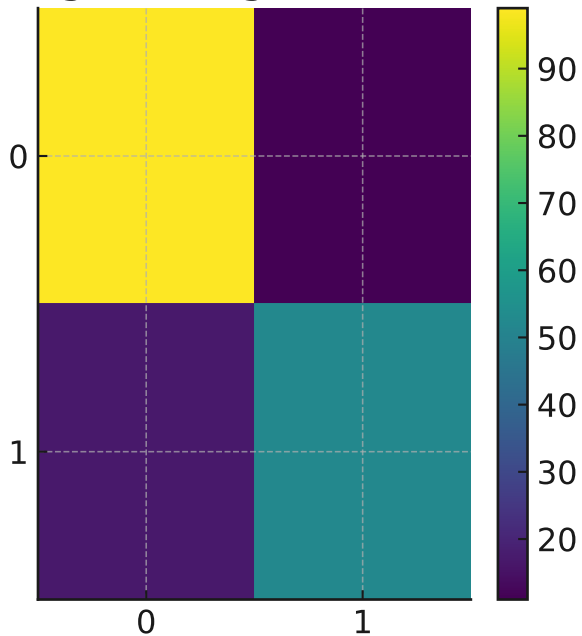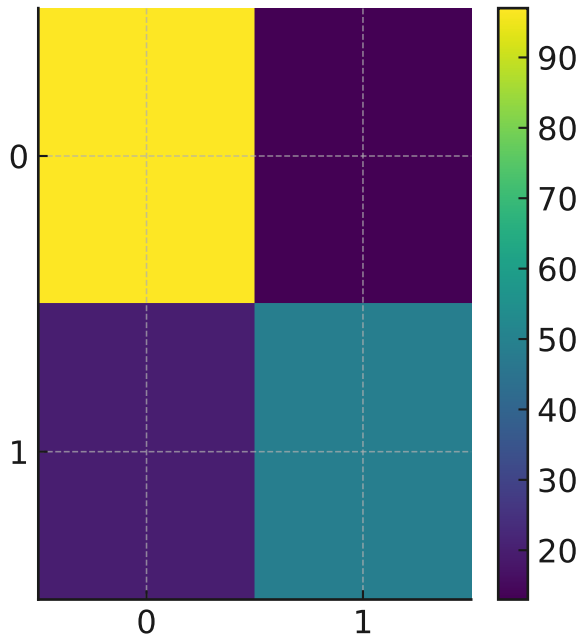Age distribution (after median group imputation)

Survival rate by Title

Top feature importances (Random Forest)

Logistic Regression CM | Random Forest CM
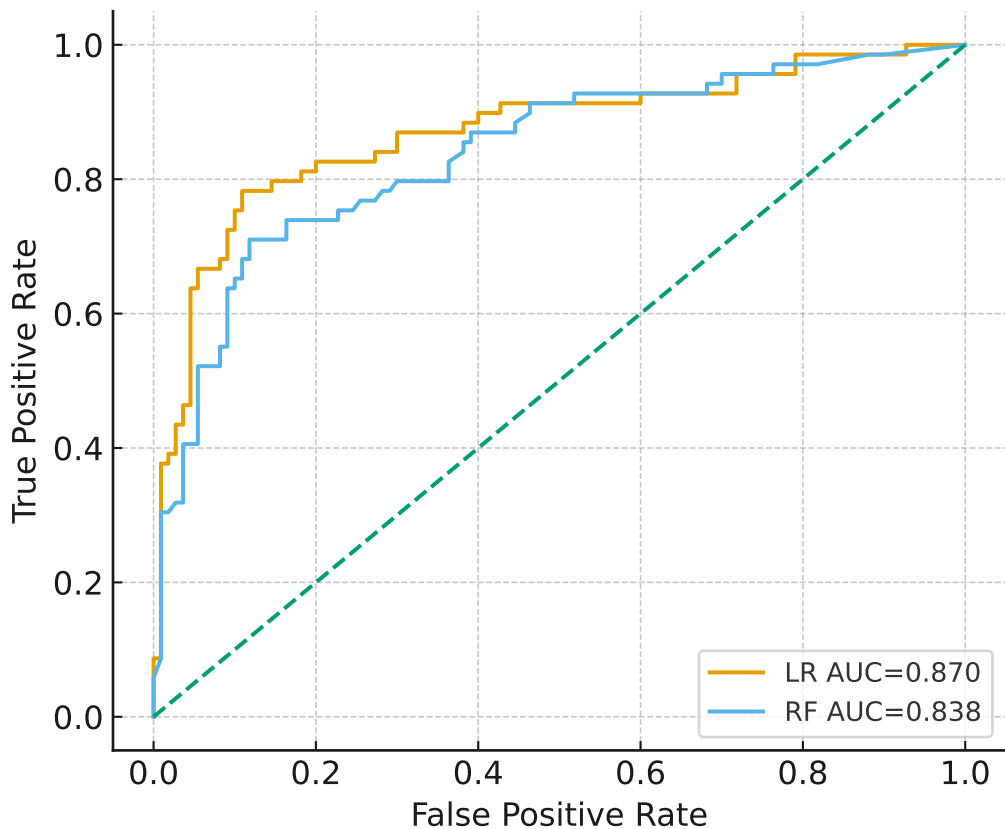
ROC Curves

LR AUC=0.870
RF AUC=0.838

Observations and Recommendations:

- Logistic Regression accuracy: 0.844, ROC AUC: 0.870

- Random Forest accuracy: 0.816, ROC AUC: 0.838

- Important features (from RF): Fare, Age, Sex_female, Title_Mr, Sex_male, FamilySize

- Next steps: tune hyperparameters, cross-validate, test more feature engineering (ticket groups, inter

- For final model, consider K-fold CV and calibration of probabilities.