# Cover Page

Course Title: Machine Learning
Team Members:
- Mohamed Mostafa : 23101594
- Marwan Khaled : 23101599
- Mohamed Adel : 23101899

Generated: 2025-12-19 12:42

Project: Movie Metadata Analytics for Streaming Acquisition Decisions

# Section 1: Problem Domain

We model movie performance signals to support streaming acquisition decisions.
Classification target: is_high_rated derived from vote_average.
Regression target: log(1+revenue) for movies with reported revenue.

# Section 2: Project Summary

Problem Domain: Streaming platform decision support using real-world movie metadata.
Dataset: Top_10000_Movies.csv (10014 rows, 13 columns).
Dirty data indicators: missing/blank text, zero revenues, skewed distributions, outliers, and semi-structured categorical genres.
ML tasks: (1) Classification of High Rated movies; (2) Regression for log-revenue.
Models: Logistic Regression, Random Forest, SVC, KNN (classification) and Linear/Ridge, Random Forest Regressor, SVR, KNN (regression).
Evaluation: Cross-validation and test metrics; scaling vs normalization comparisons included.

# Section 3: Source Code

Source Code: This submission includes Ml_Project.ipynb and a template download script.
Preprocessing uses ColumnTransformer with robust imputation and one-hot encoding.
Models are compared under three preprocessing modes: none, StandardScaler (scaling), and MinMax (normalization).

# Section 4: Visualization Snapshots

The following pages include snapshots of model comparisons, EDA, required visualizations, and model evaluation plots.

# Model Comparison (Classification Table)

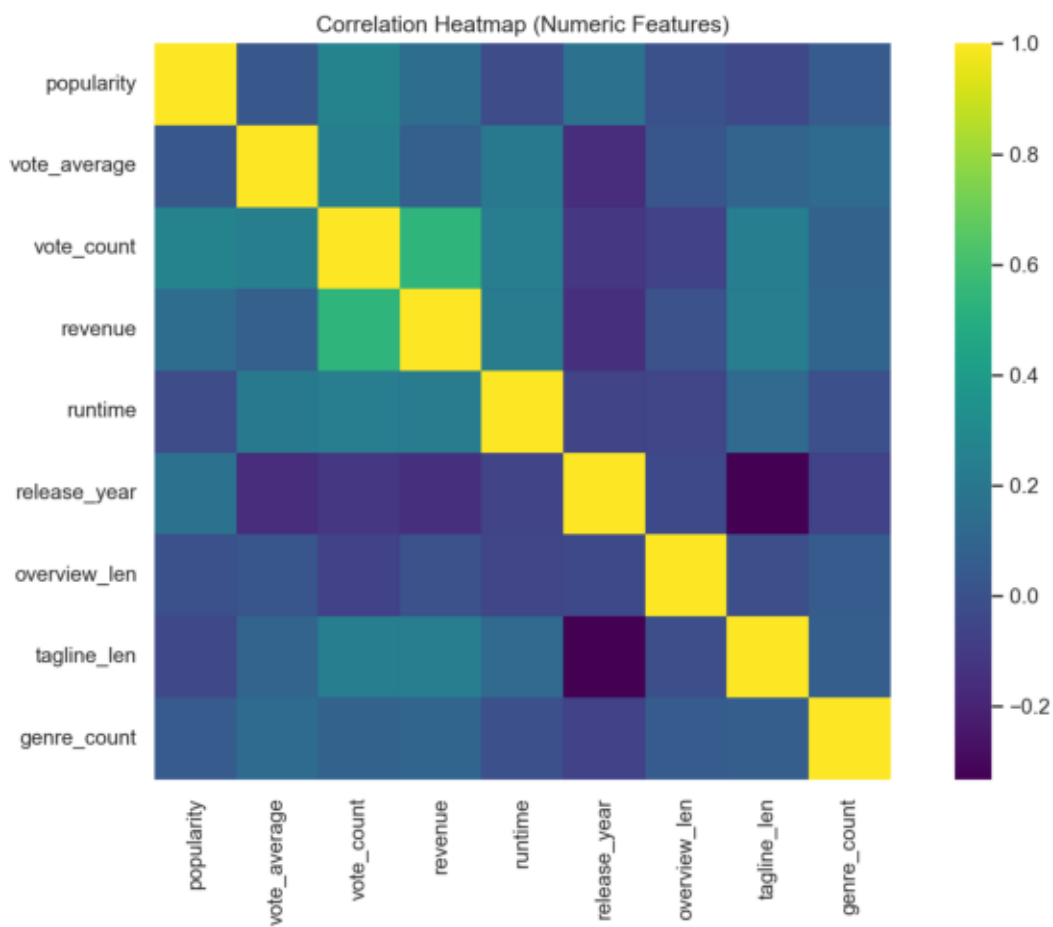**Classification (4 Models): Mean CV ROC-AUC by Preprocessing Mode**

| Model | none | standard | minmax |
|---|---|---|---|
| KNN | 0.576 | 0.774 | 0.791 |
| LogReg | 0.782 | 0.805 | 0.805 |
| RandomForest | 0.847 | 0.846 | 0.846 |
| SVC_RBF | 0.592 | 0.823 | 0.807 |

# Model Comparison (Regression Table)

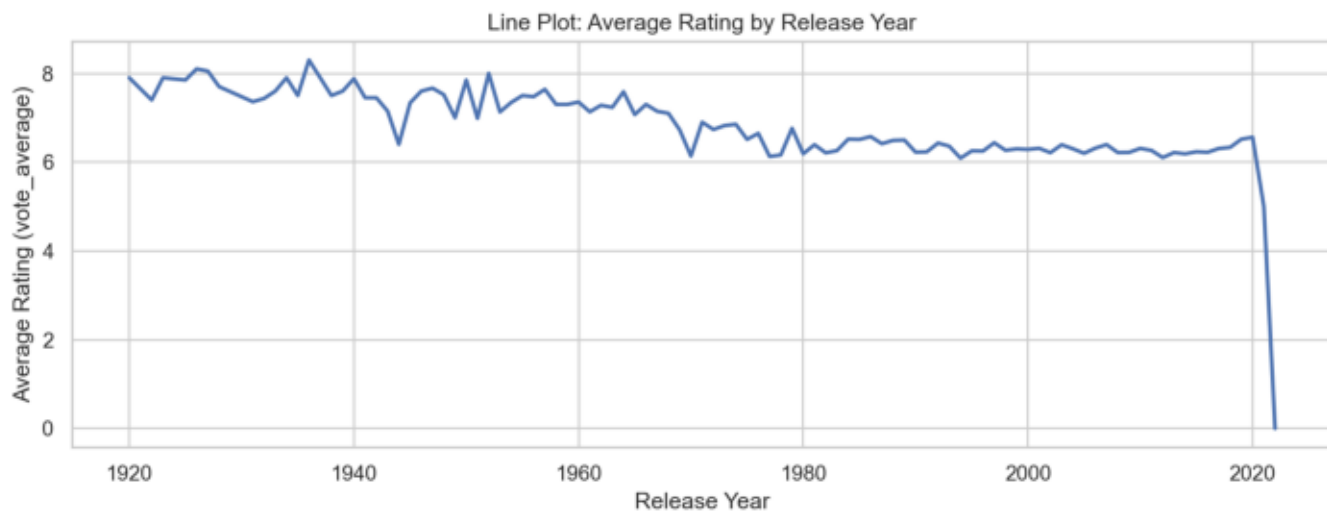**Regression (4 Models): Mean CV R² by Preprocessing Mode (log-revenue)**

| Model | none | standard | minmax |
|---|---|---|---|
| KNN | 0.15 | 0.247 | 0.205 |
| Linear | 0.19 | 0.191 | 0.191 |
| RandomForest | 0.304 | 0.302 | 0.304 |
| Ridge | 0.187 | 0.192 | 0.192 |
| SVR_RBF | 0.102 | 0.288 | 0.202 |

# EDA: Correlation Heatmap



Correlation Heatmap (Numeric Features)

# Visualization 1: Line Plot

Line Plot: Average Rating by Release Year

# Visualization 2: Area Plot



Area Plot: Total Reported Revenue by Release Year

# Visualization 3: Histogram



Histogram: Runtime Distribution
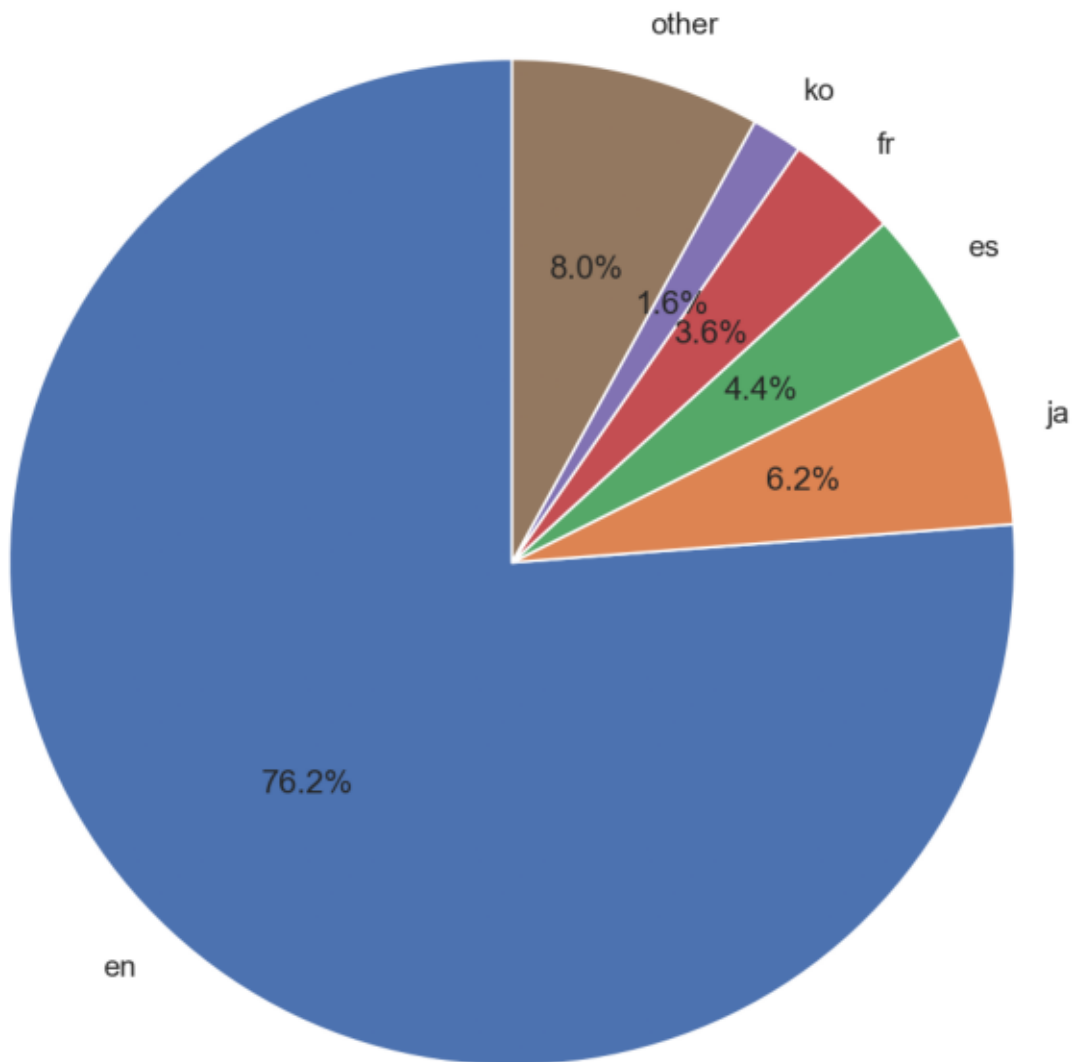
# Visualization 4: Bar Chart



Bar Chart: Top 10 Main Genres by Count

# Visualization 5: Pie Chart

Pie Chart: Original Language Distribution (Top 5 + Other)



other

ko

fr

es

8.0%

1.6%

3.6%

4.4%

ja

6.2%

76.2%

en

# Visualization 6: Box Plot



Box Plot: Rating Distribution by Main Genre (Top 6)

# Visualization 7: Scatter Plot



Scatter Plot: Popularity vs Rating

# Visualization 8: Bubble Plot



Bubble Plot: Popularity vs Revenue (Bubble Size = Vote Count)

# Classification: CV ROC-AUC Comparison



Classification Model Comparison (5-fold CV ROC-AUC)

# Classification: Confusion Matrix



Confusion Matrix (Test Set)

# Classification: ROC Curve



ROC Curve (Test Set)

# Regression: CV R2 Comparison



Regression Model Comparison (5-fold CV $R^2$ on log-revenue)

# Regression: Predicted vs Actual



Regression: Predicted vs Actual (log-revenue)

# Regression: Residuals



Regression Residuals vs Predicted