

Final Report (Lab-Style Redo)

Course Title: Machine Learning

Section: _____

Team Members:

- Mohamed Mostafa : 23101594
- Marwan Khaled : 23101599
- Mohamed Adel : 23101899

Dataset: Top_10000_Movies.csv

Figures exported: 24

Section 1: Problem Domain

We analyze a dataset of movies to study popularity and engagement signals (votes), ratings behavior, and basic relationships between movie attributes.

We then build ML models to predict high engagement and a numeric outcome (when available).

Section 2: Project Summary

Workflow:

- Load + audit + clean the dataset (types, missing values, duplicates)
- EDA with distributions, correlations, and groupby analysis
- Scoring model (weighted rating) to reduce small-sample bias
- ML: 4 lab classifiers (LogReg, DecisionTree, GaussianNB, KNN) + Linear Regression
- Export all plots to PDF and ZIP

Section 3: Source Code

Source code is provided in the notebook:

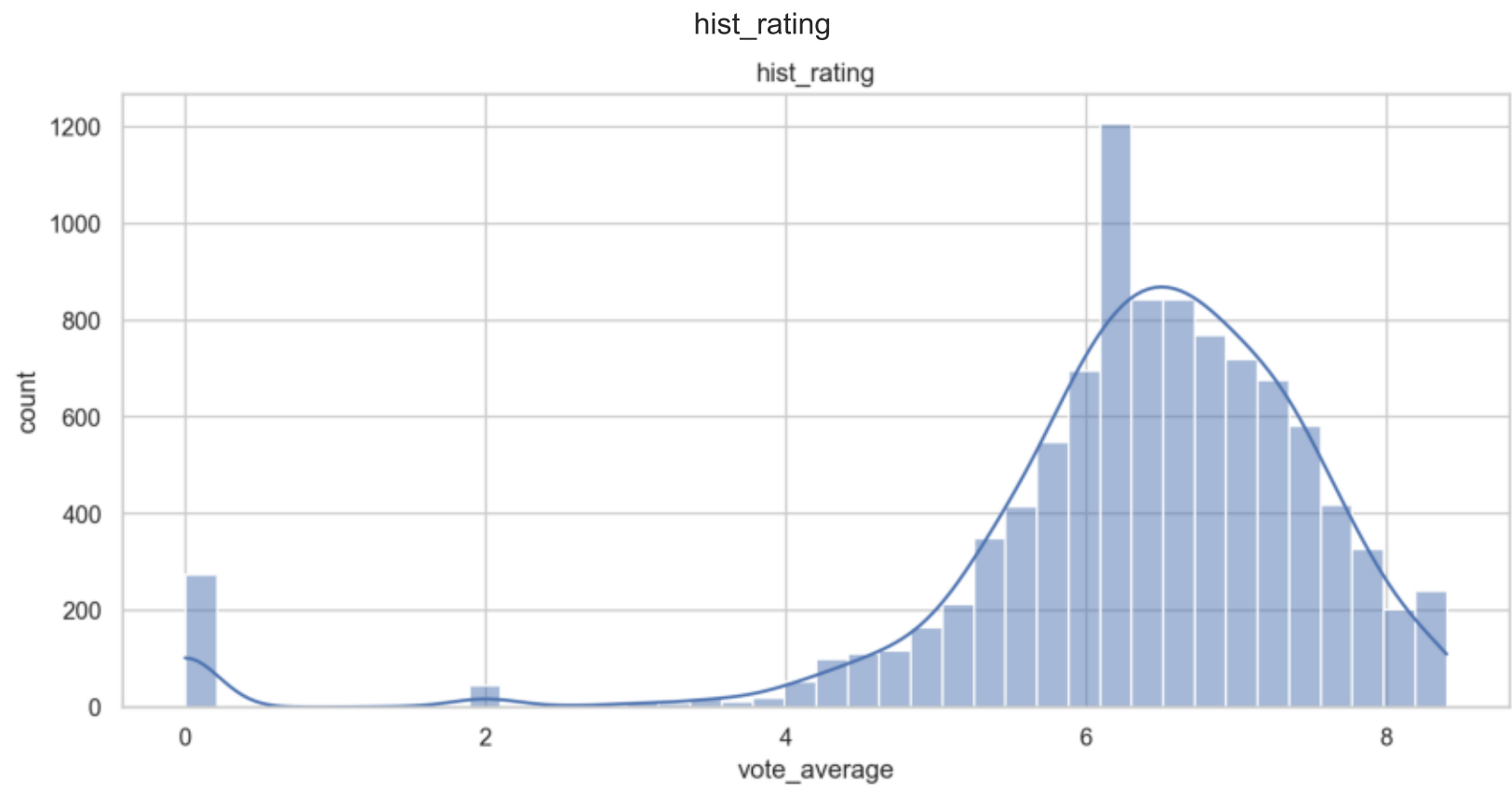
- MI_Project_LabStyle_Redo.ipynb

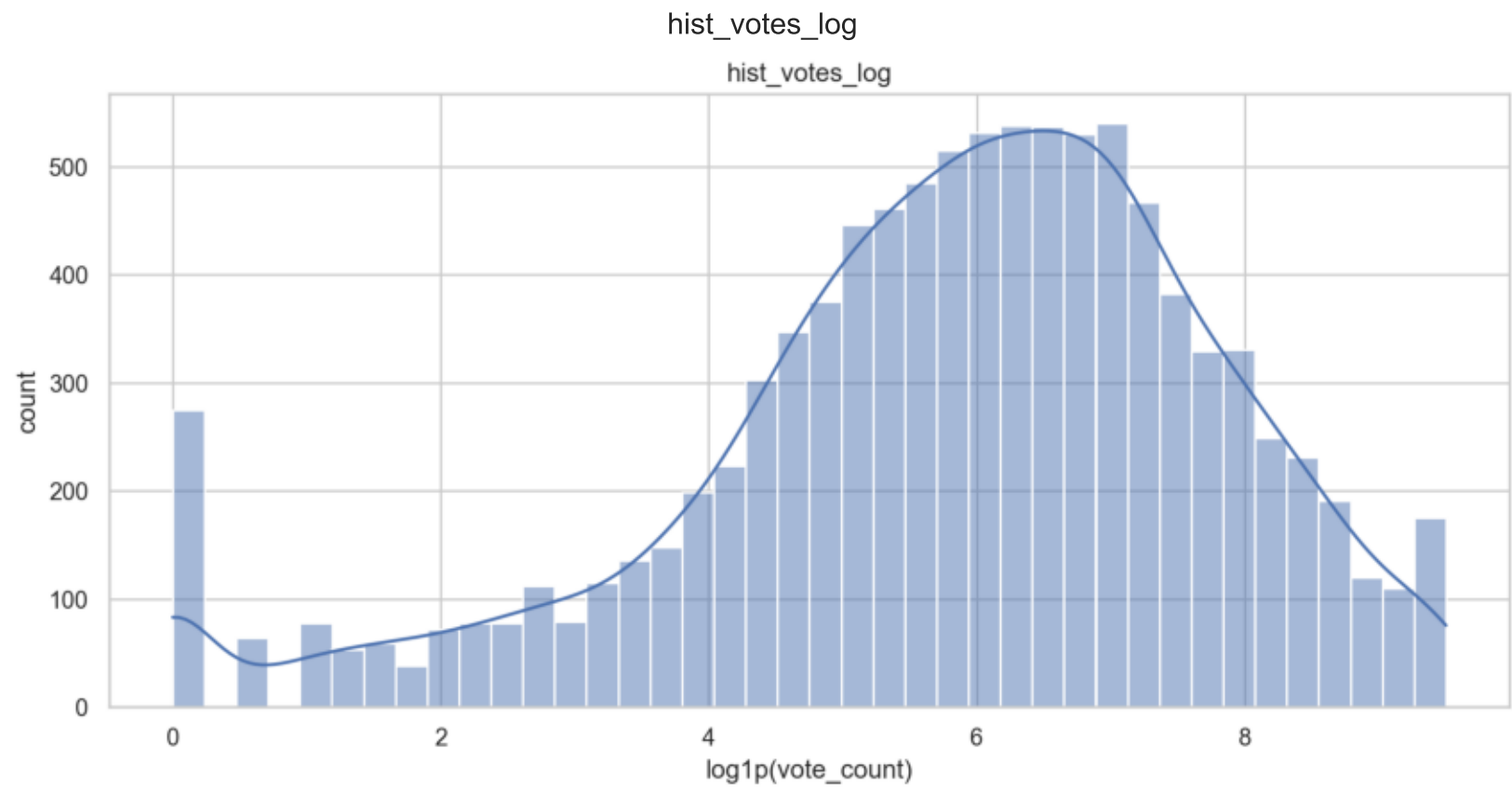
Key implementation points:

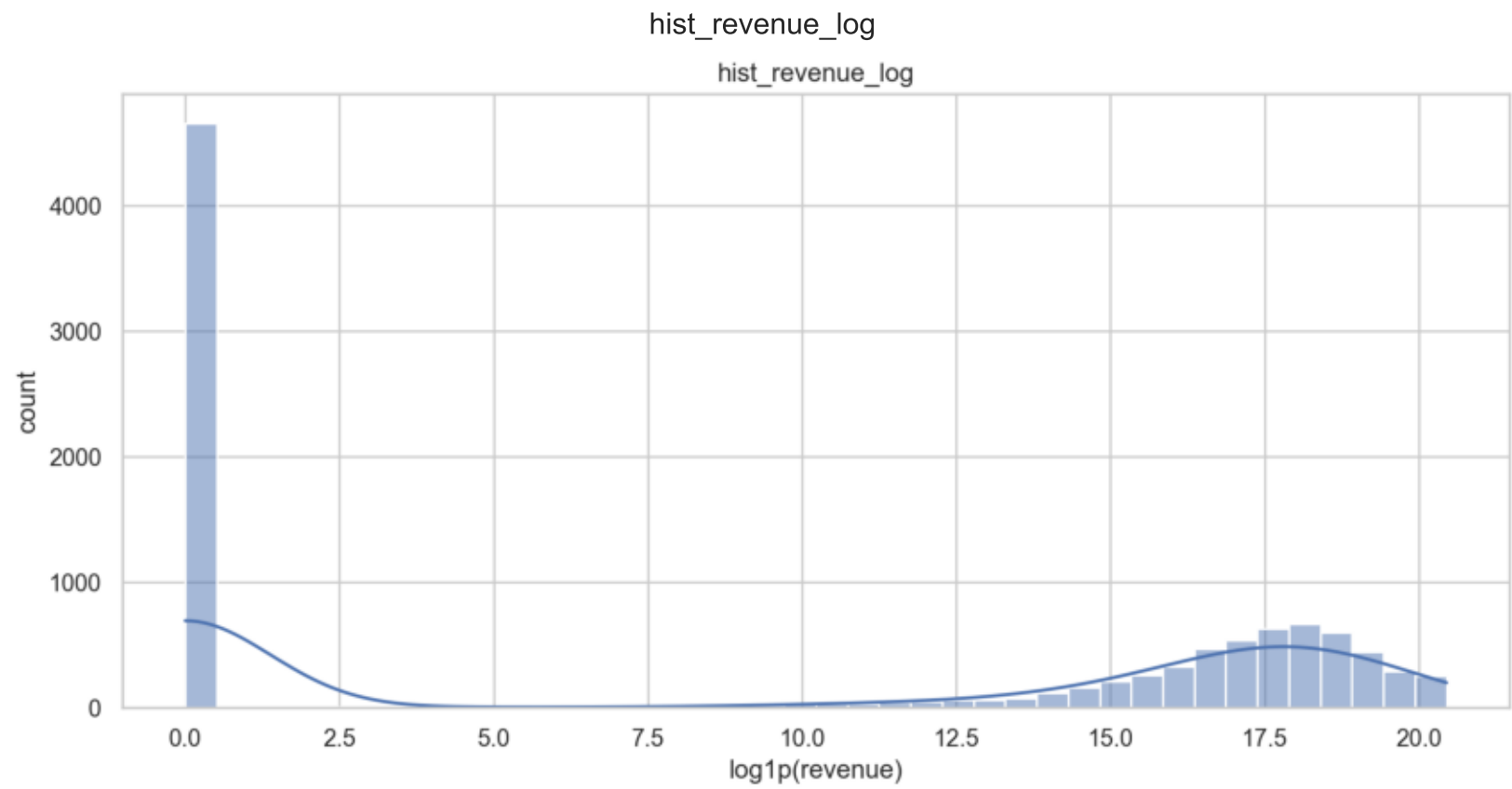
- Cleaning & feature engineering with defensive column detection
- Scikit-learn pipelines + ColumnTransformer preprocessing
- Standard metrics (Accuracy/Precision/Recall/F1, confusion matrix, ROC when available)
- Linear Regression metrics (MAE/RMSE/R2)

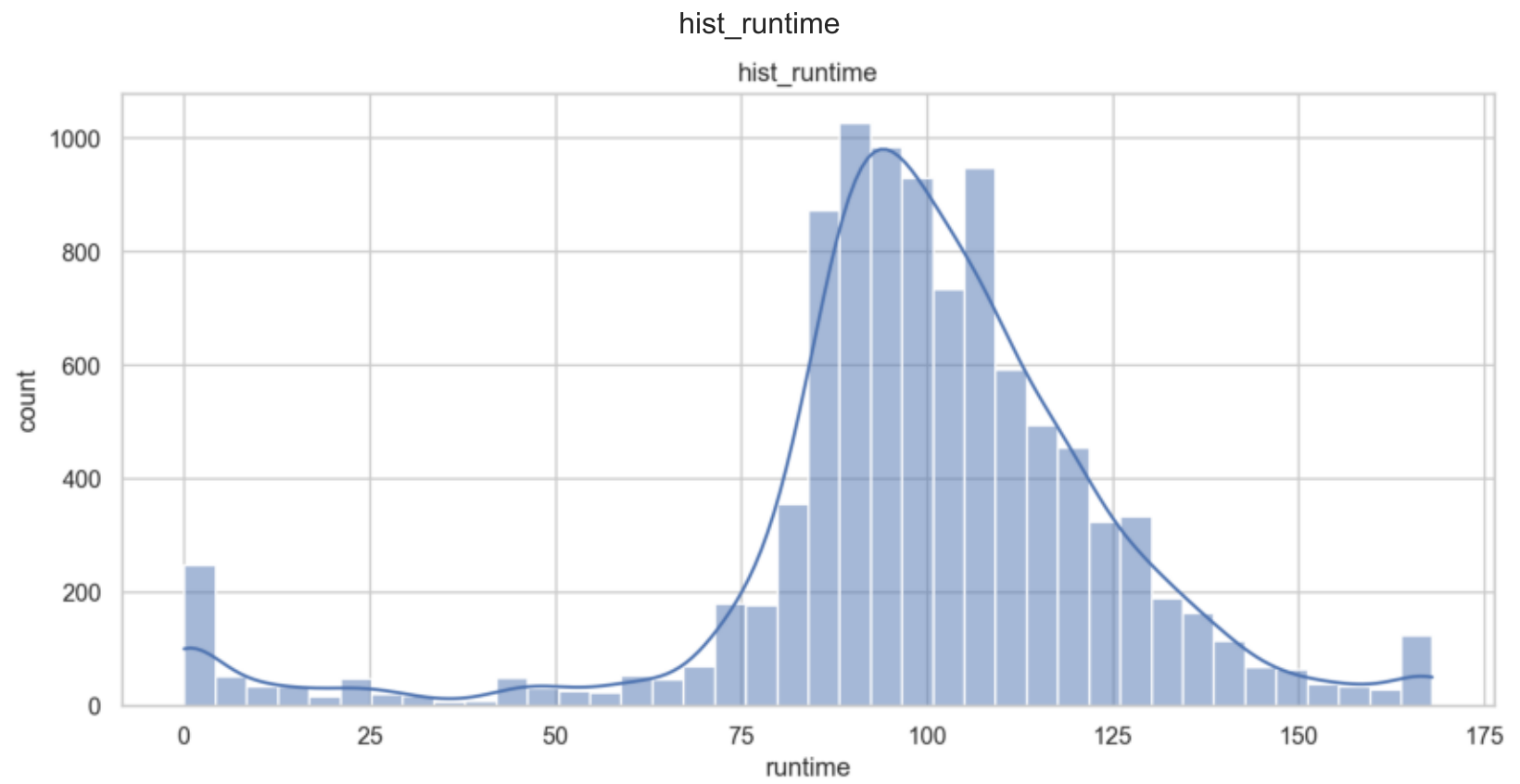
Section 4: Visualization Snapshots

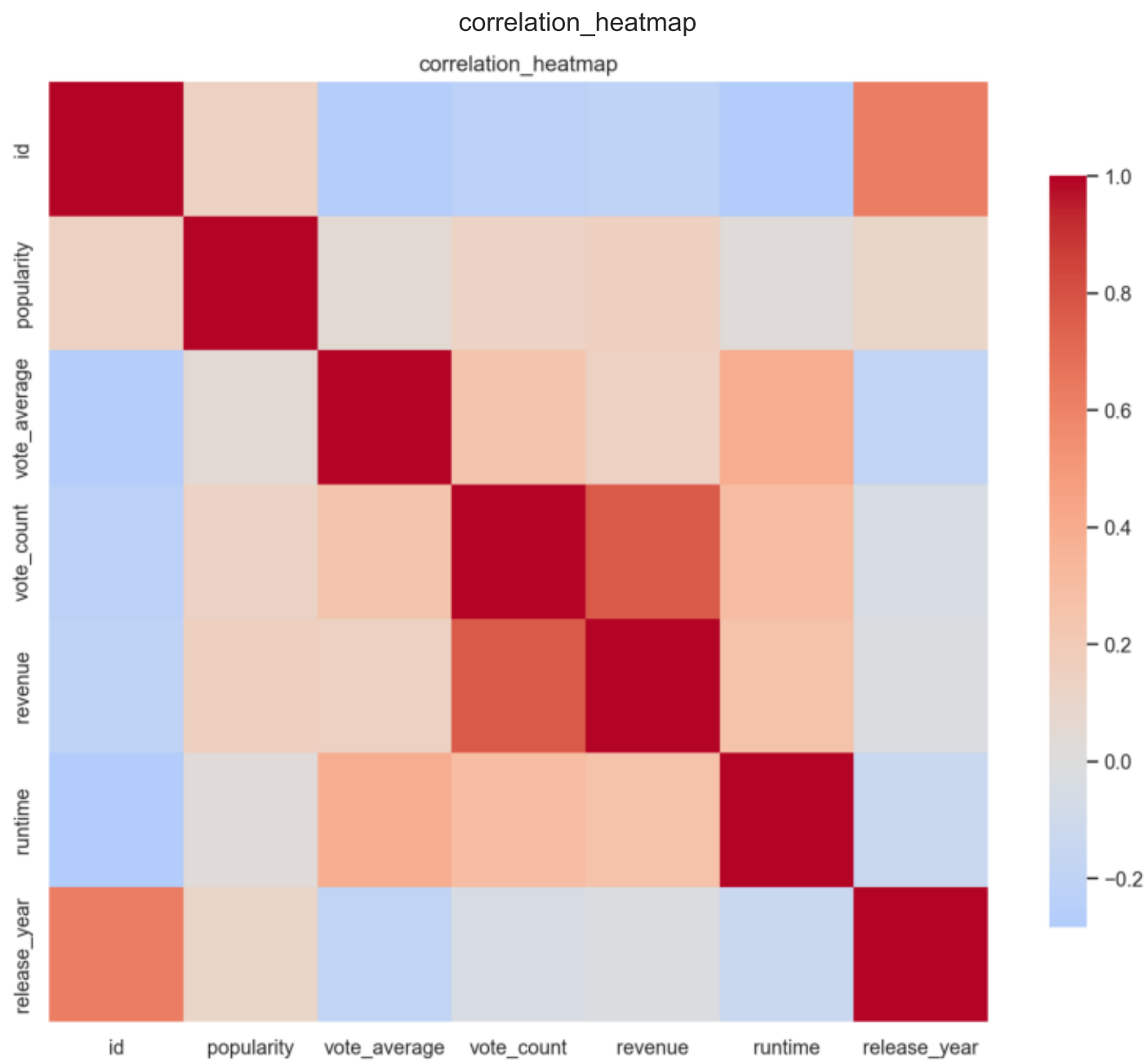
The following pages include every figure saved during EDA + ML evaluation.



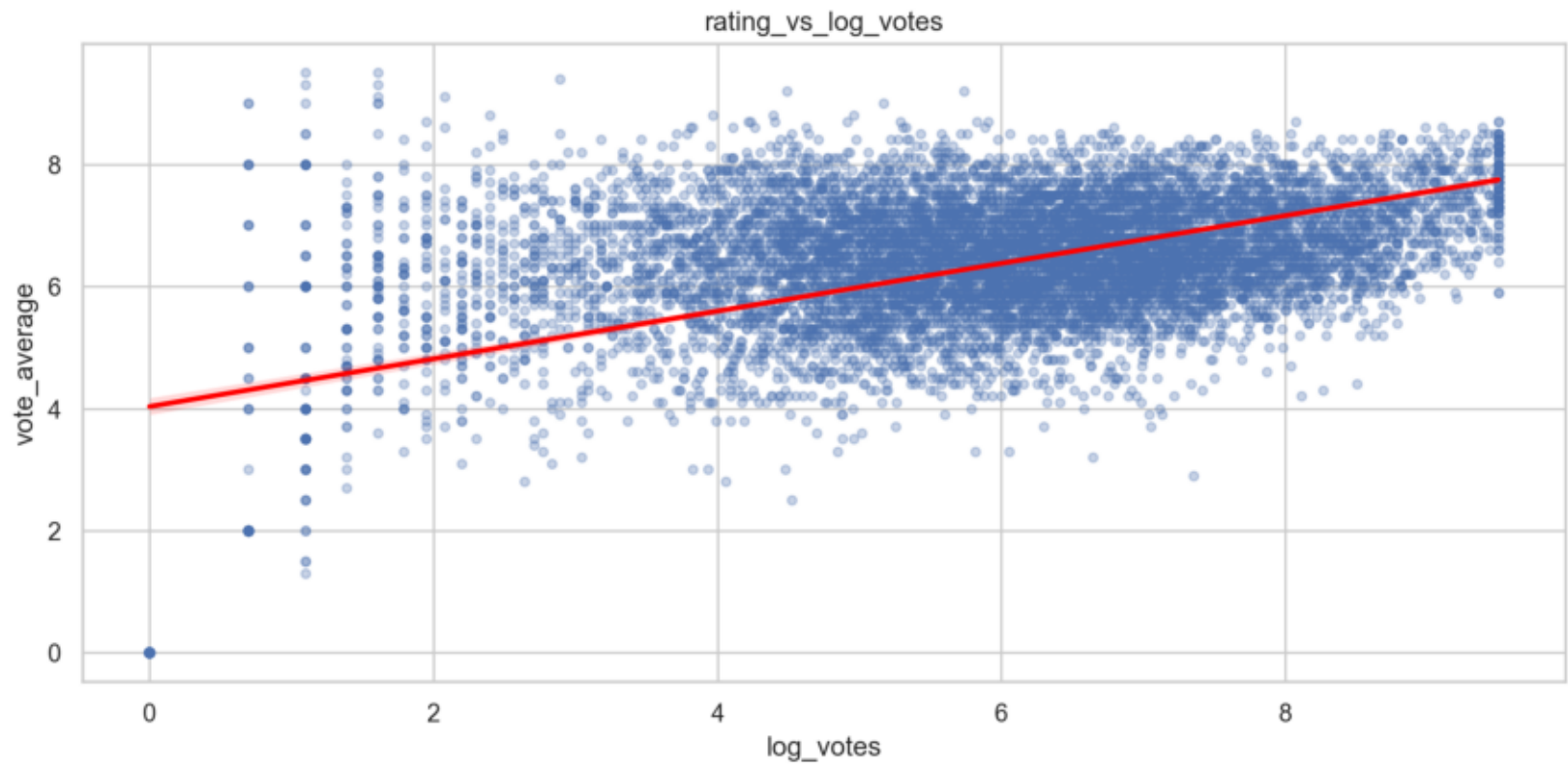




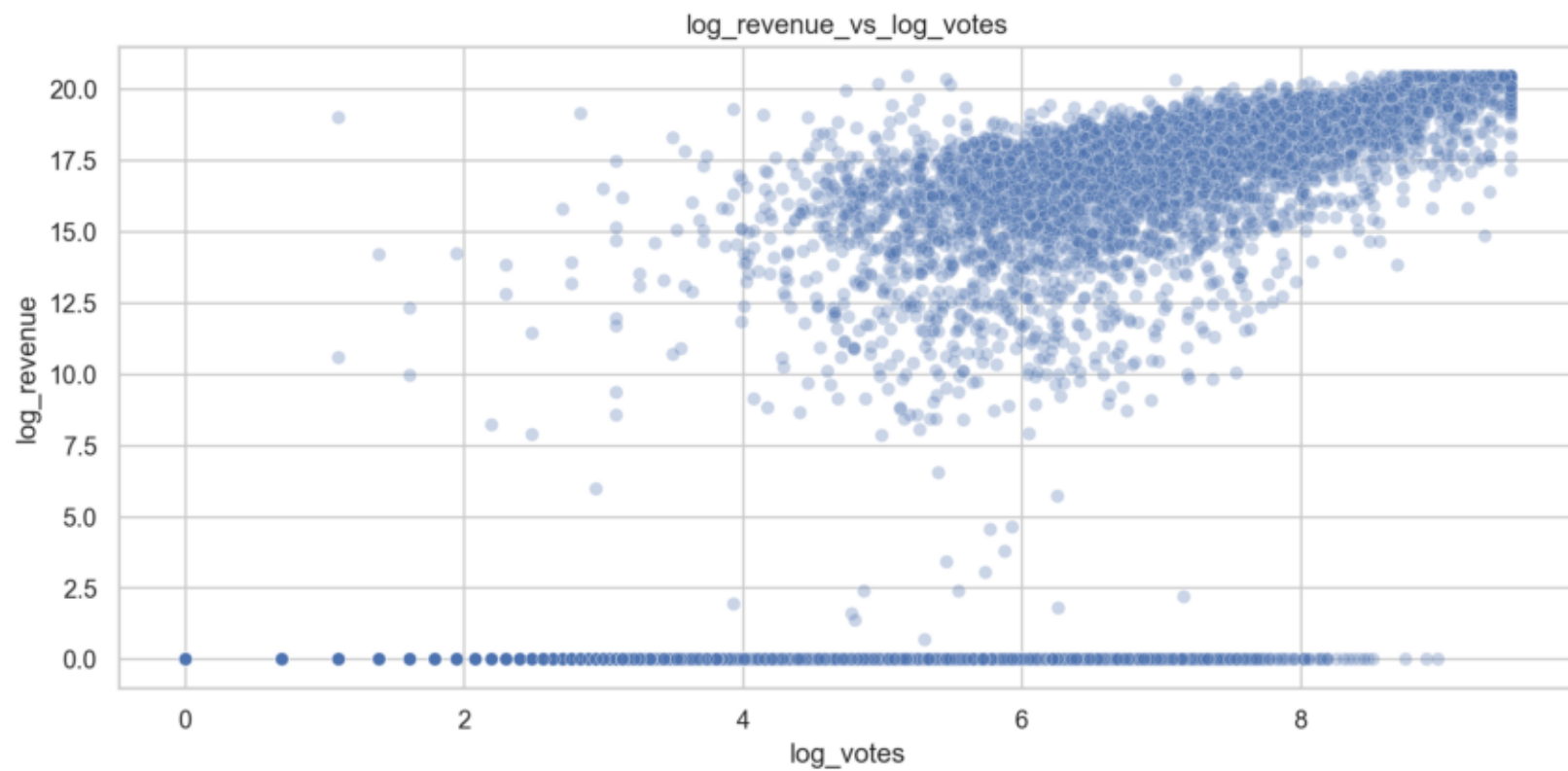




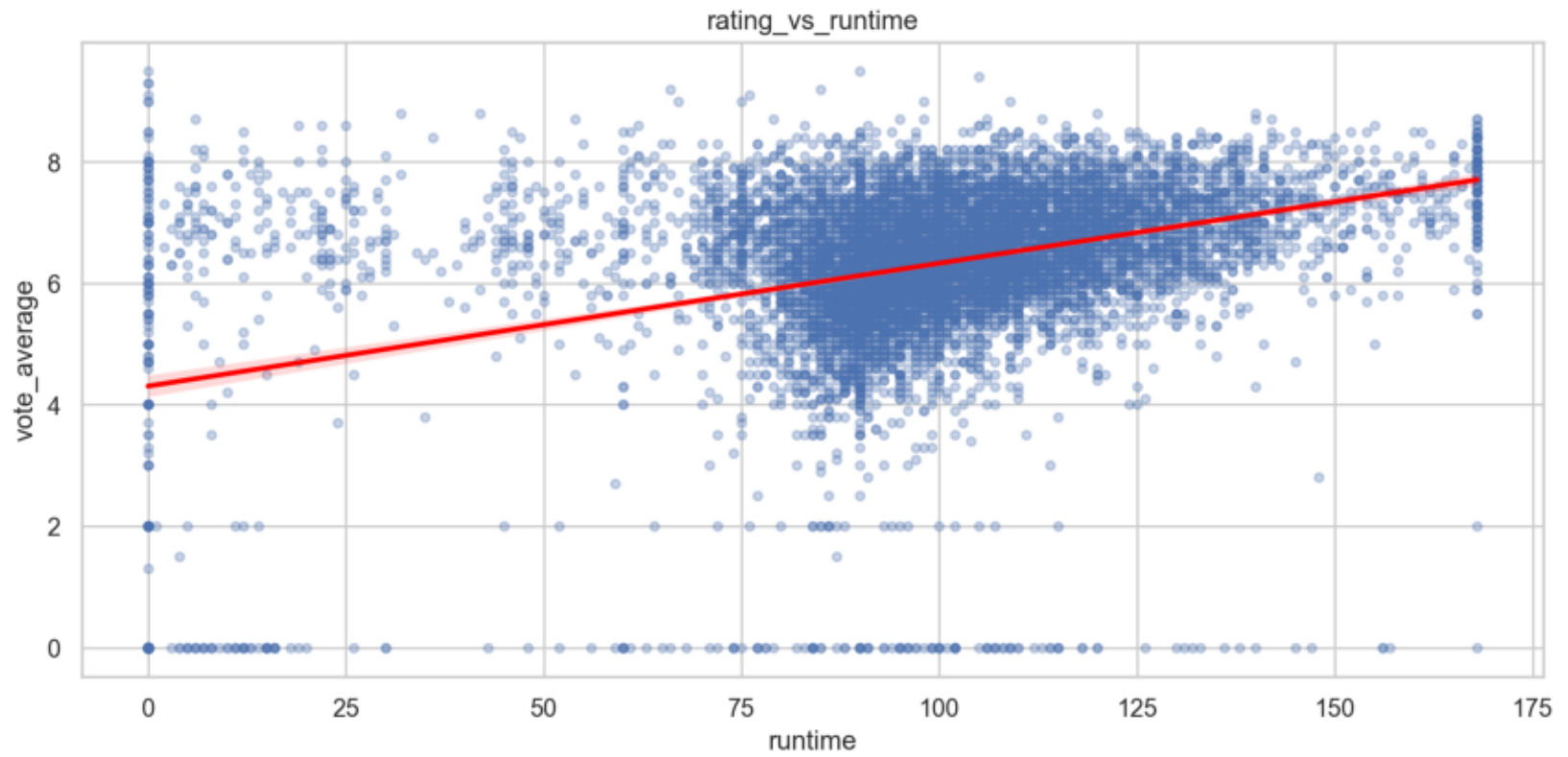
scatter_rating_vs_log_votes



scatter_log_revenue_vs_log_votes

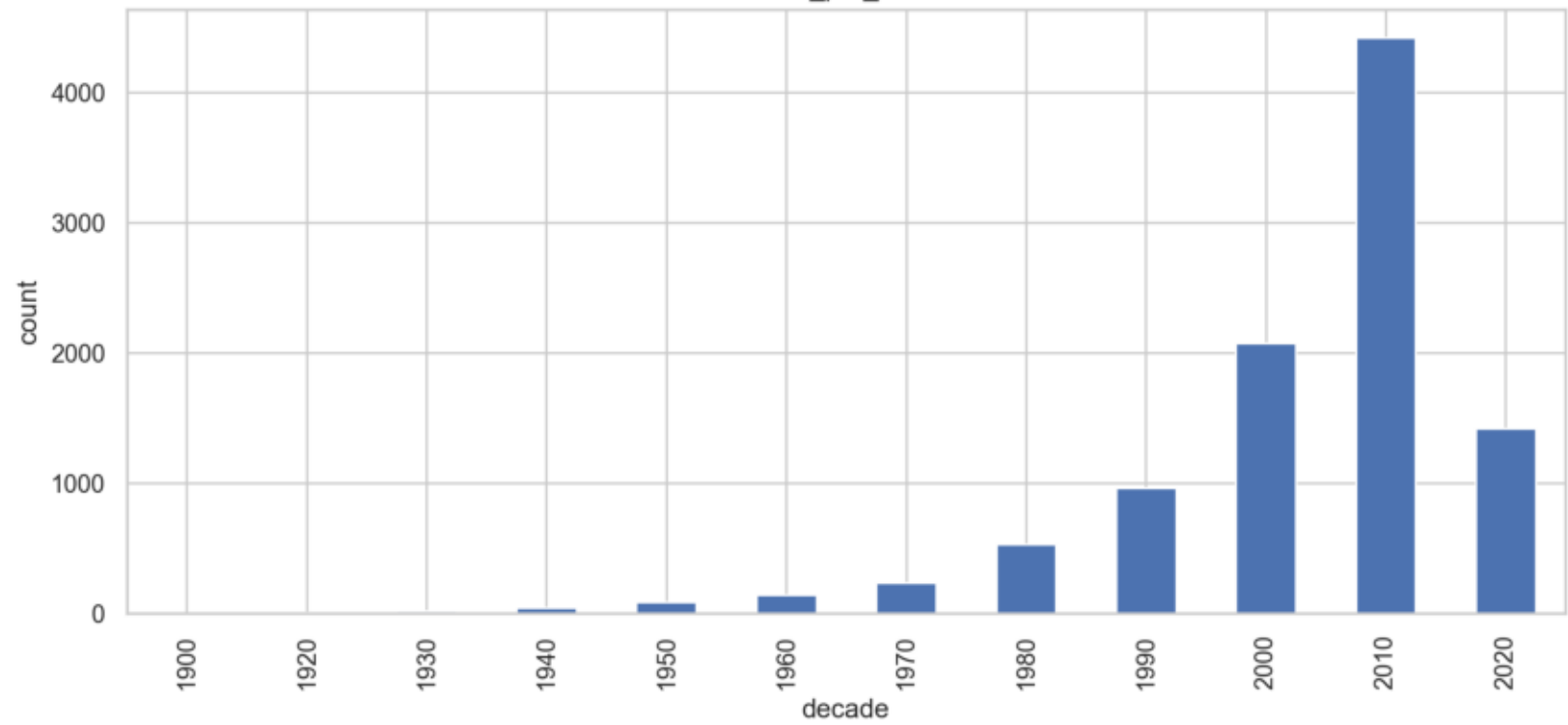


scatter_rating_vs_runtime

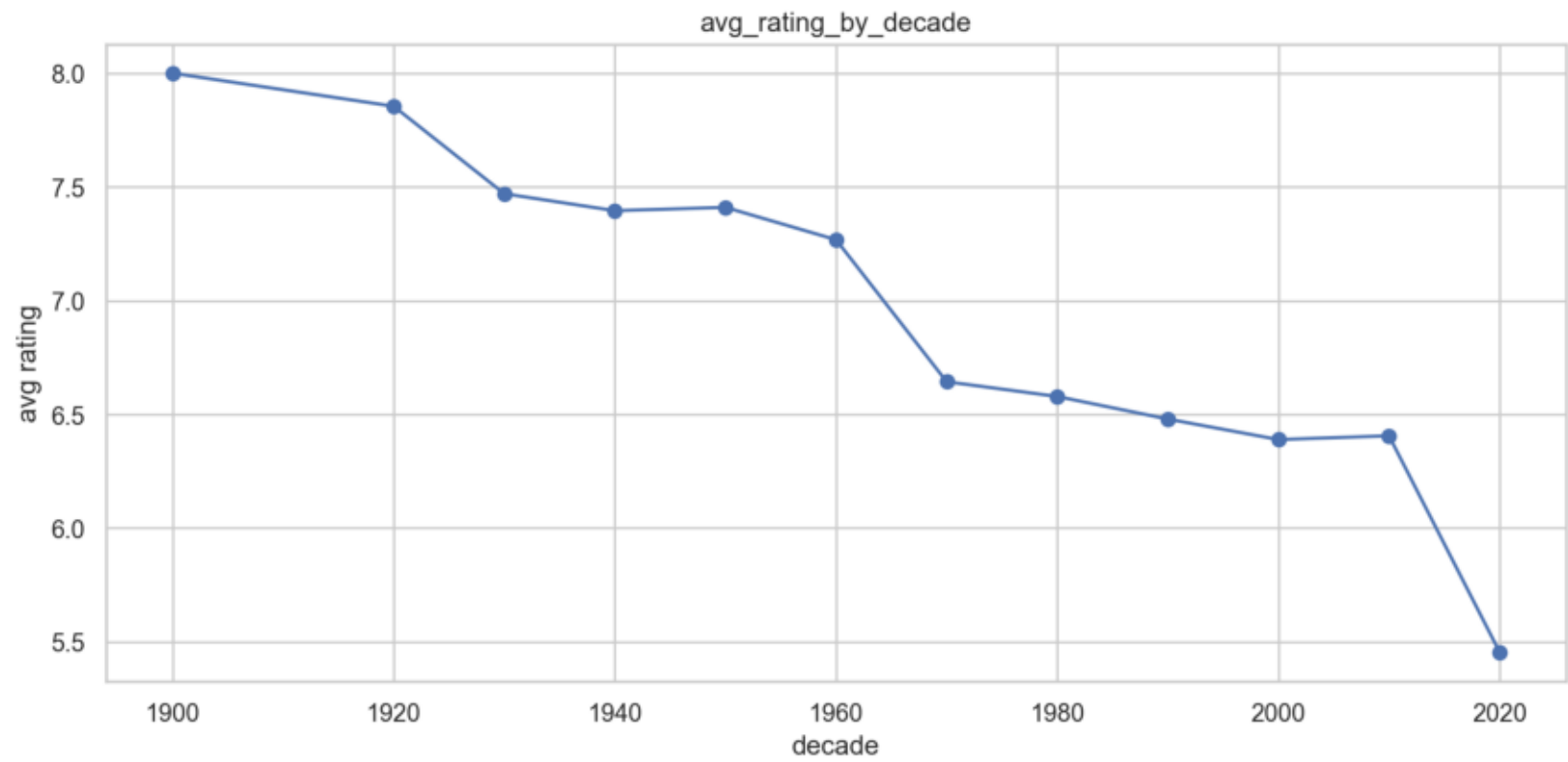


bar_movies_per_decade

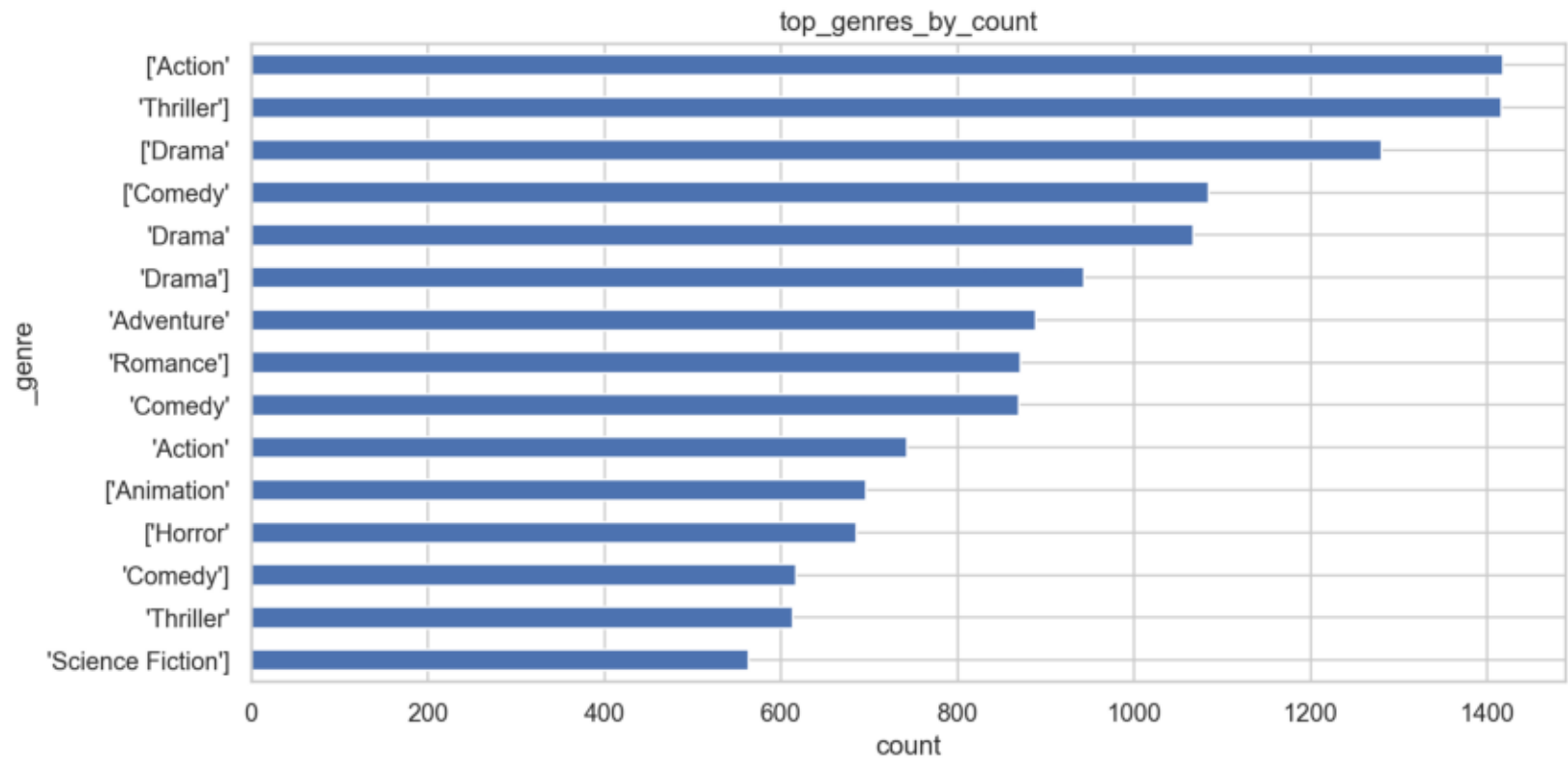
movies_per_decade



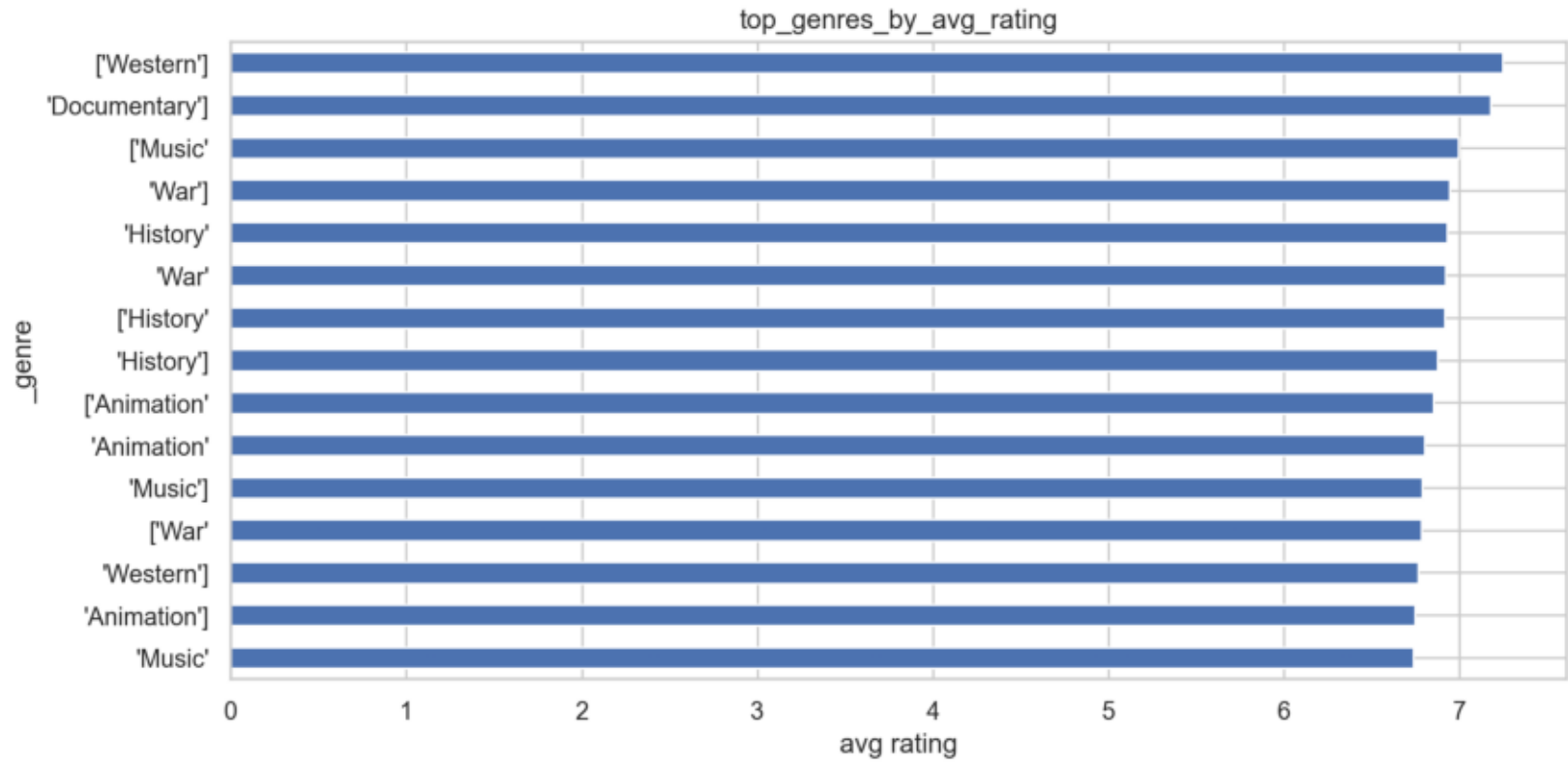
line_avg_rating_by_decade



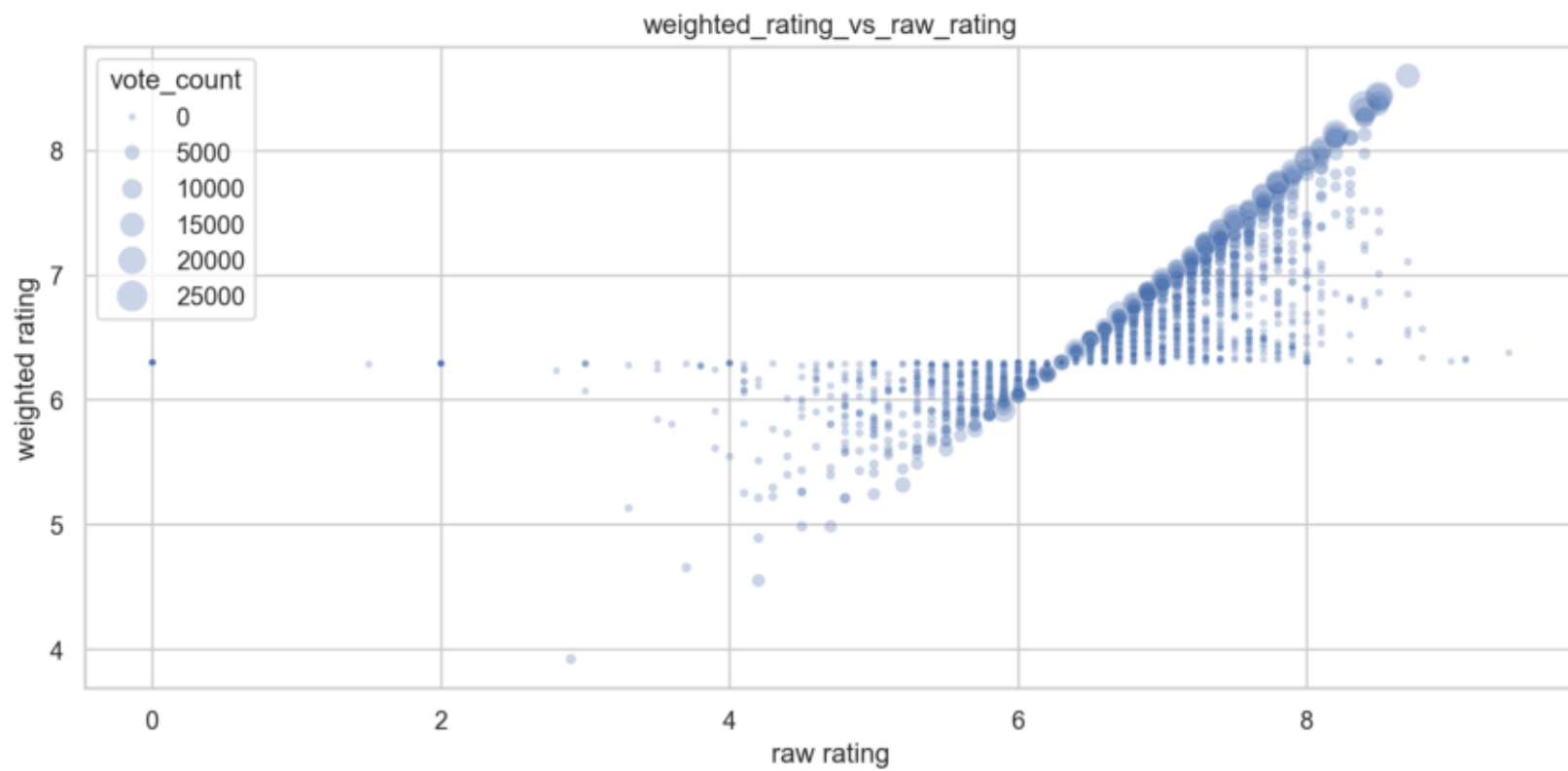
bar_top_genres_by_count



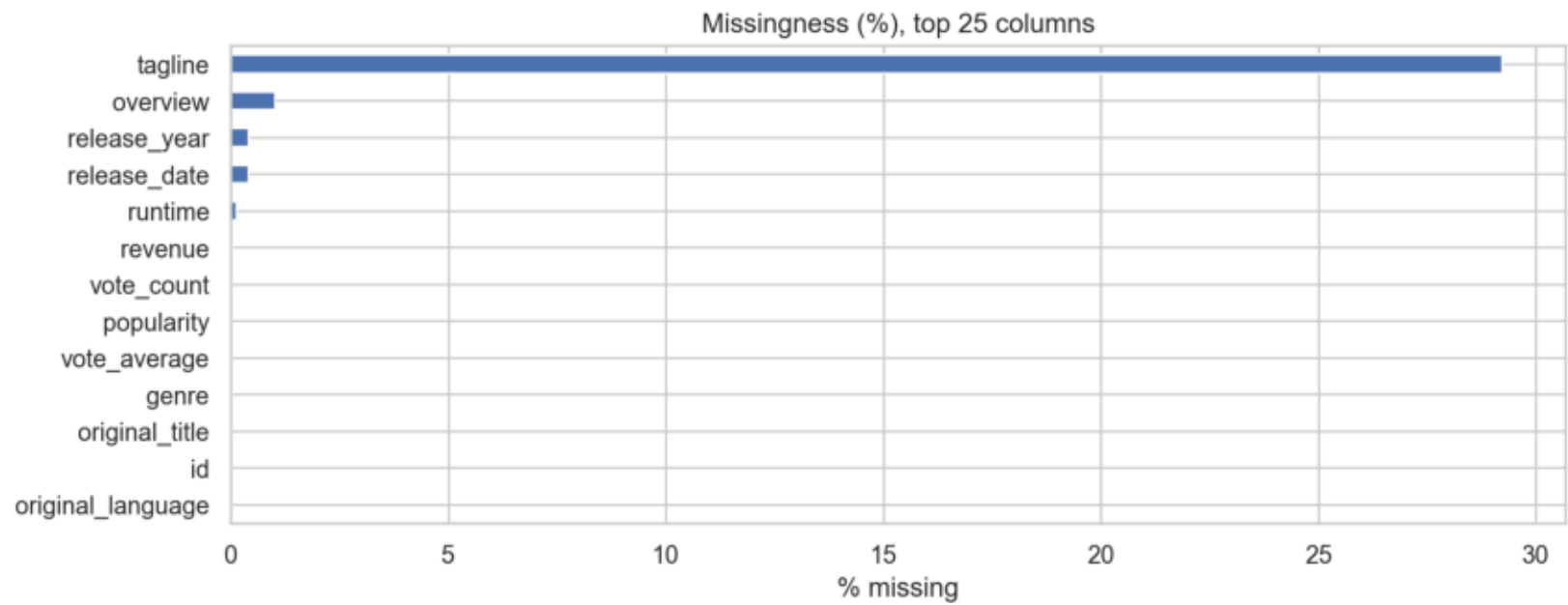
bar_top_genres_by_avg_rating



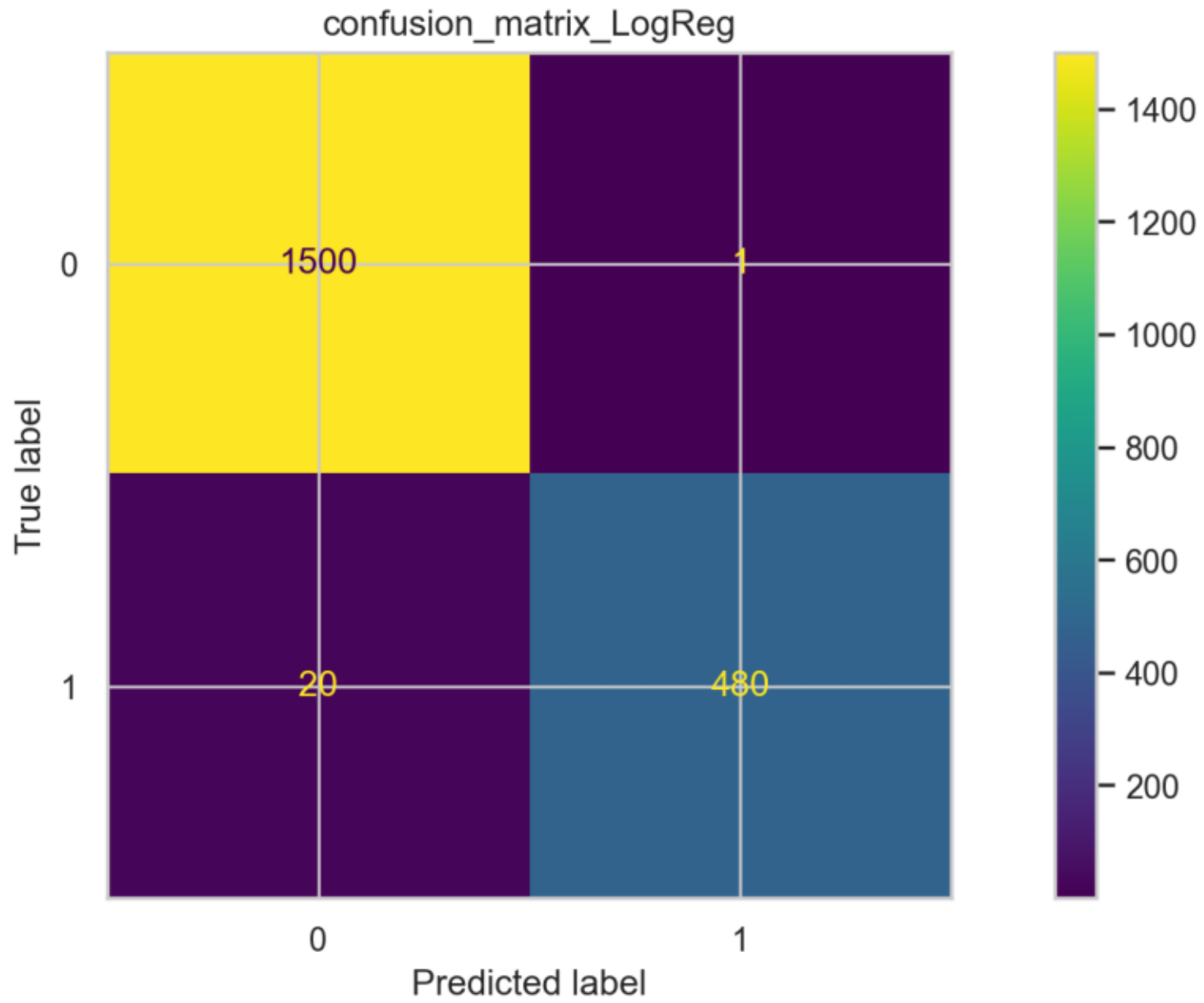
scatter_weighted_vs_raw_rating

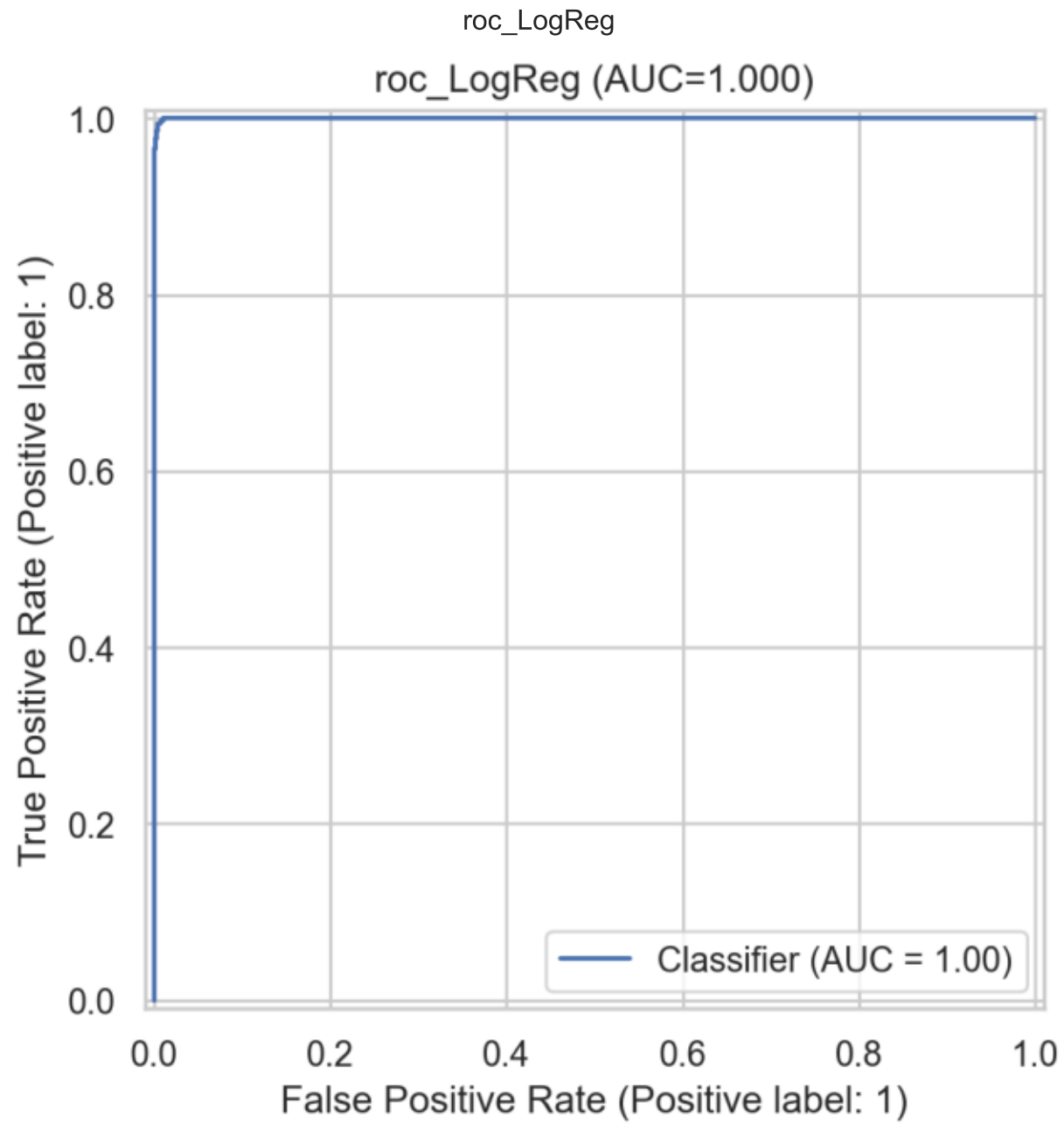


bar_missingness_top25

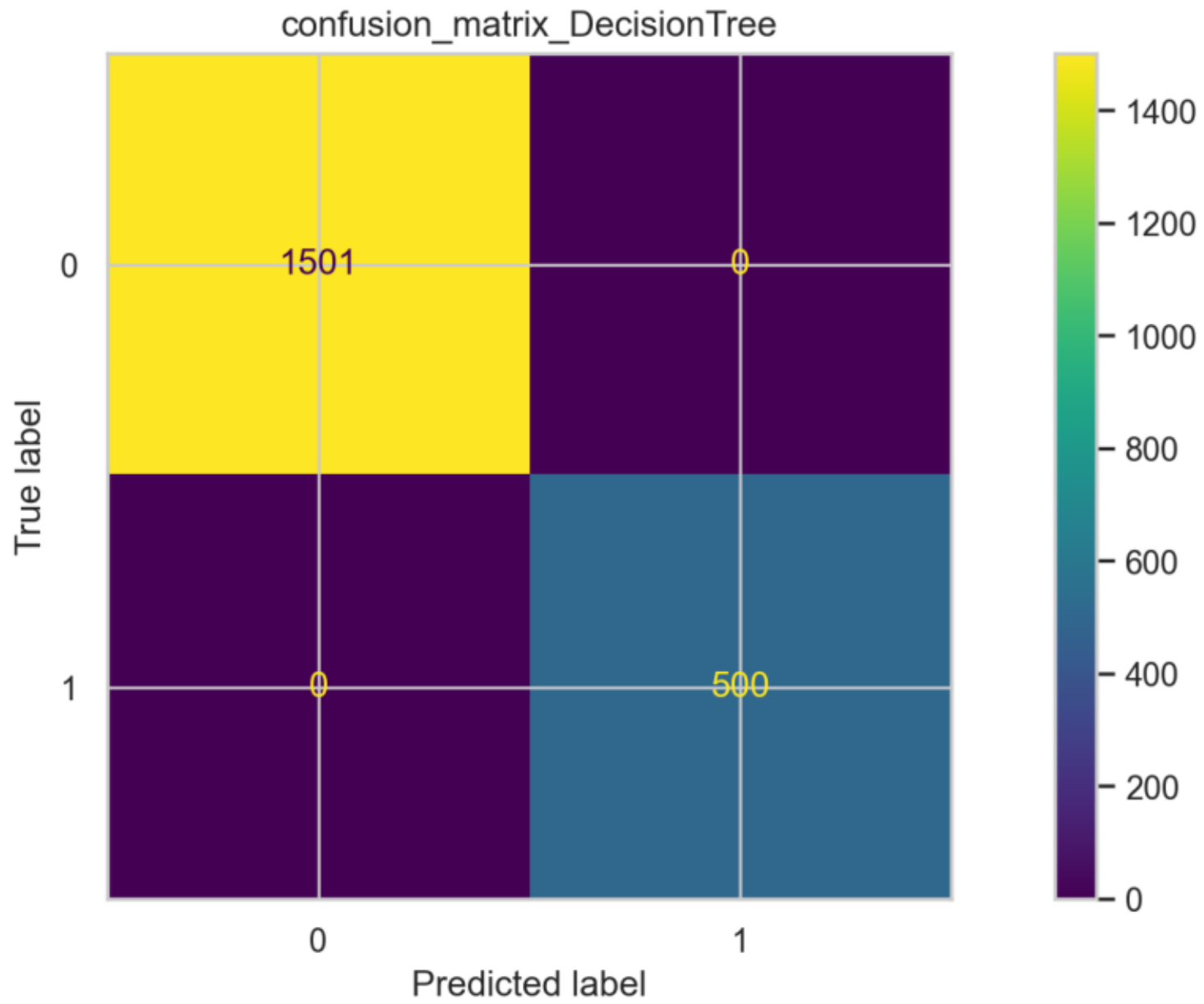


confusion_matrix_LogReg



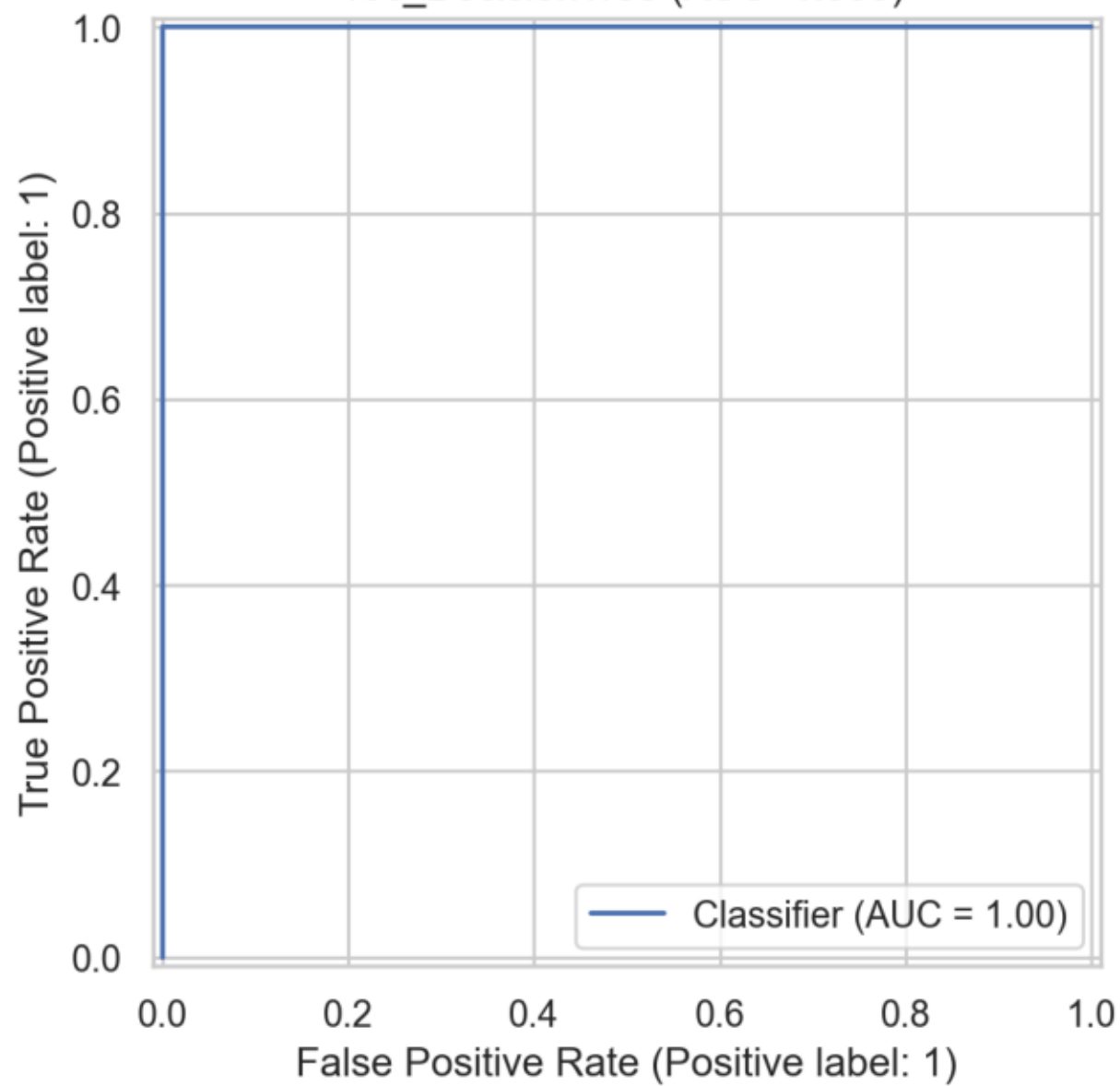


confusion_matrix_DecisionTree

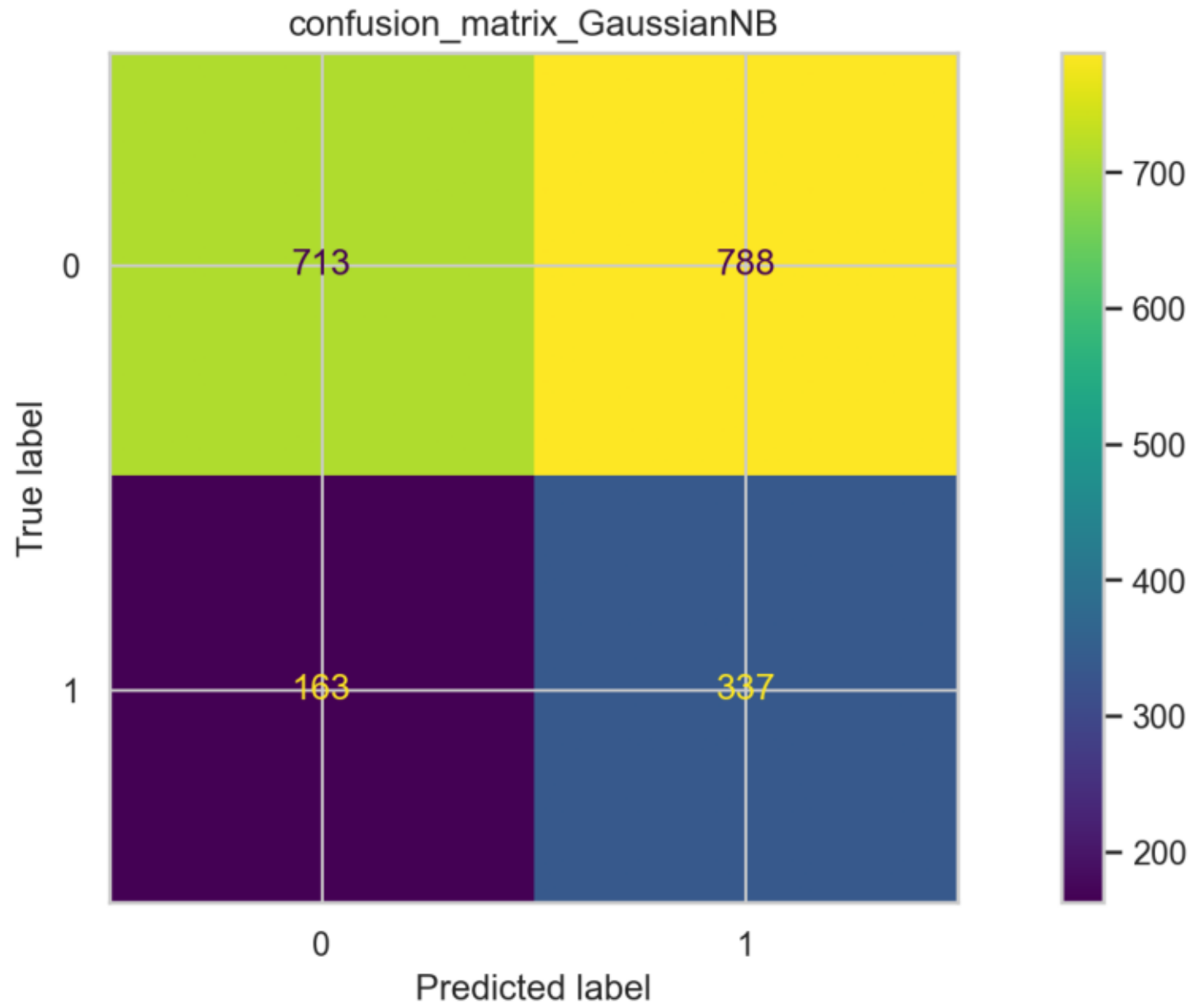


roc_DecisionTree

roc_DecisionTree (AUC=1.000)

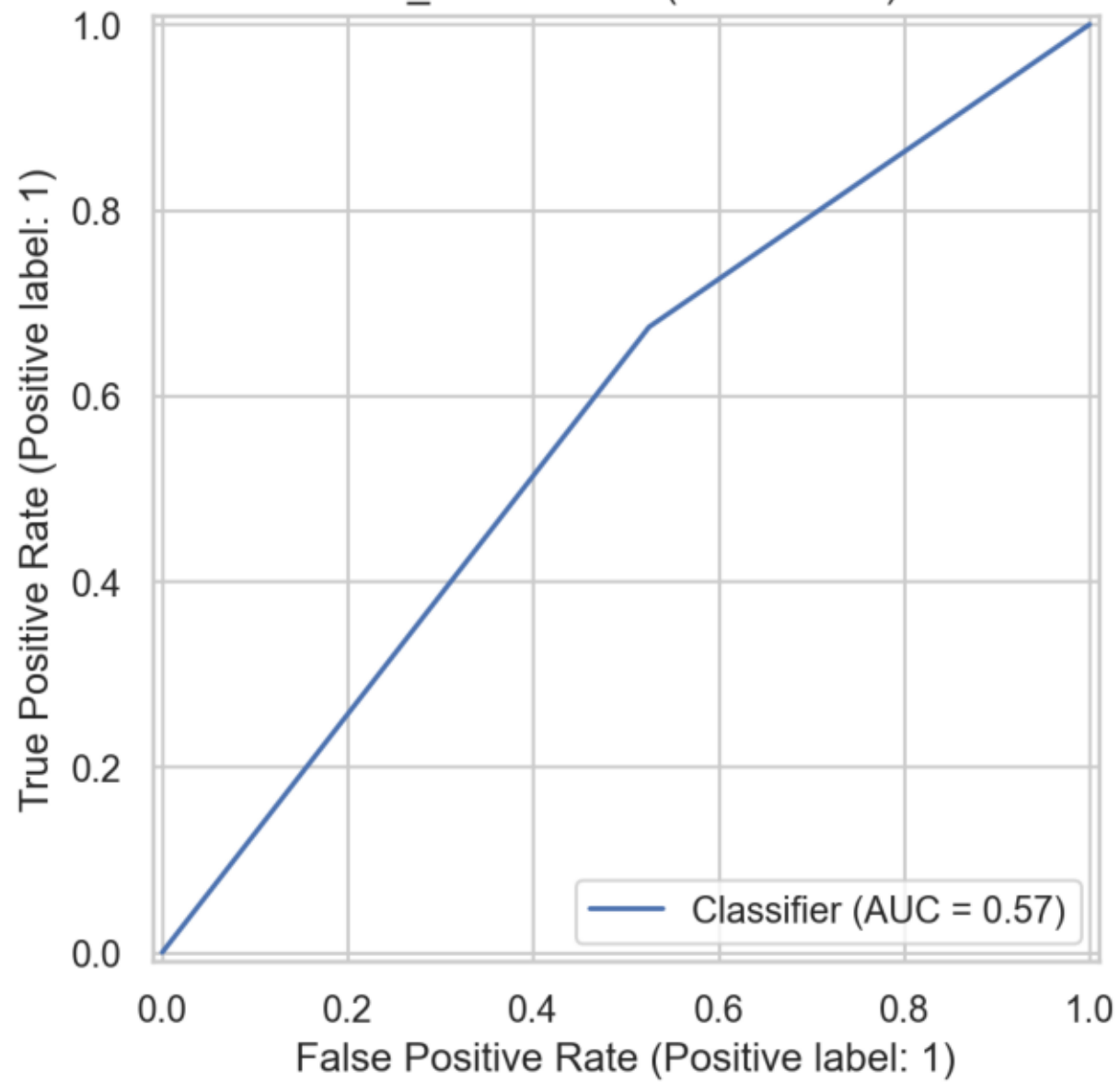


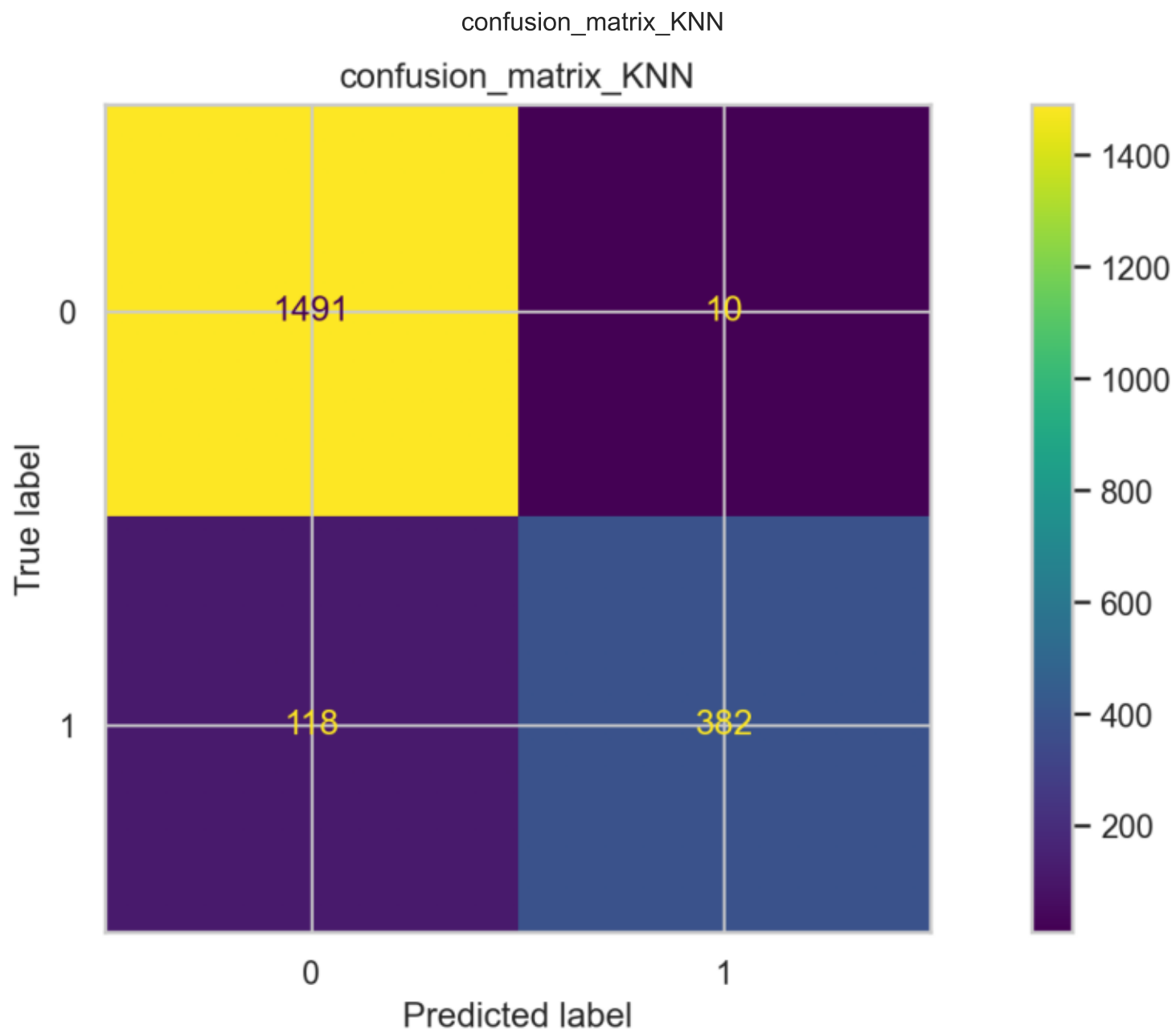
confusion_matrix_GaussianNB

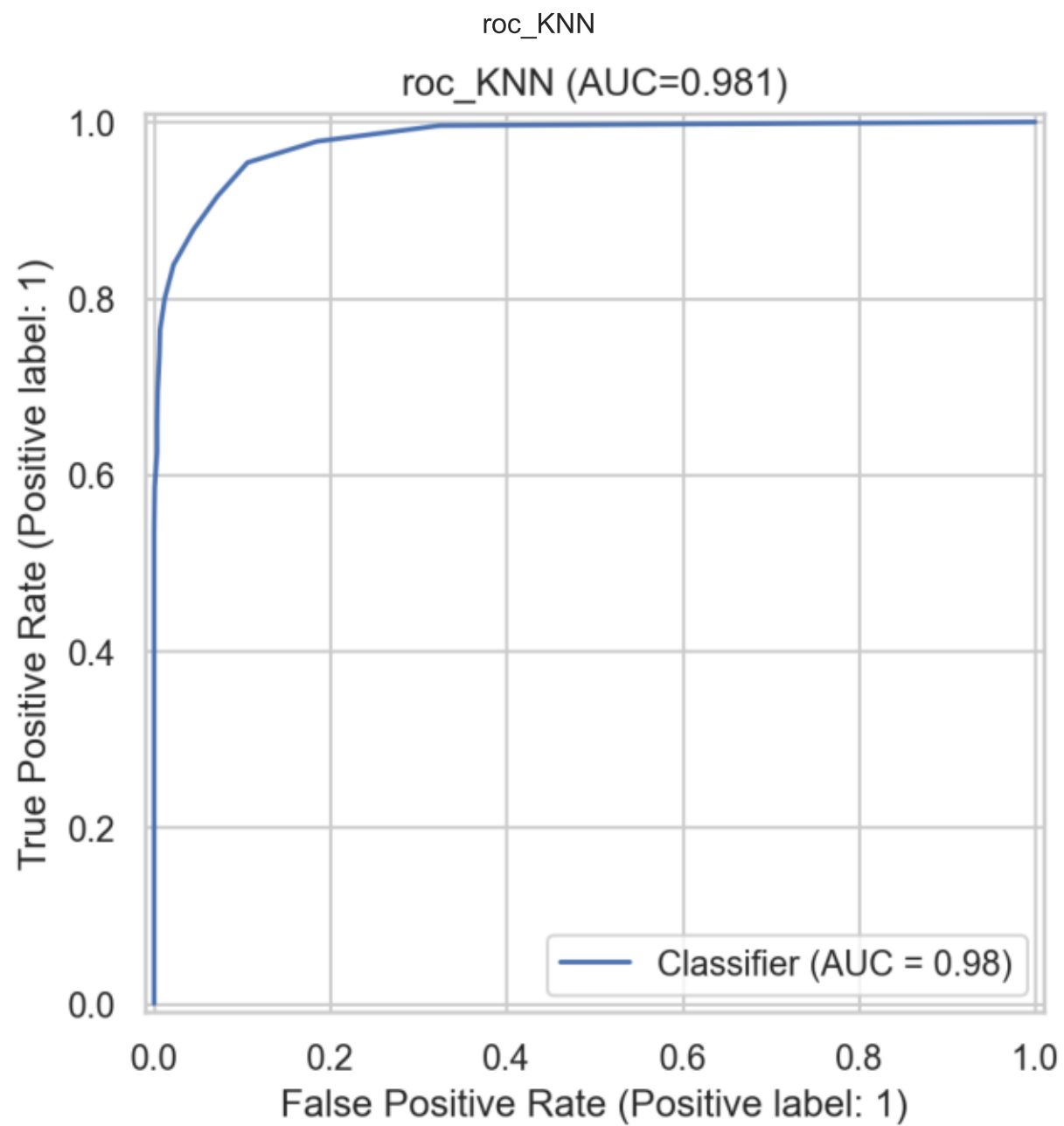


roc_GaussianNB

roc_GaussianNB (AUC=0.575)

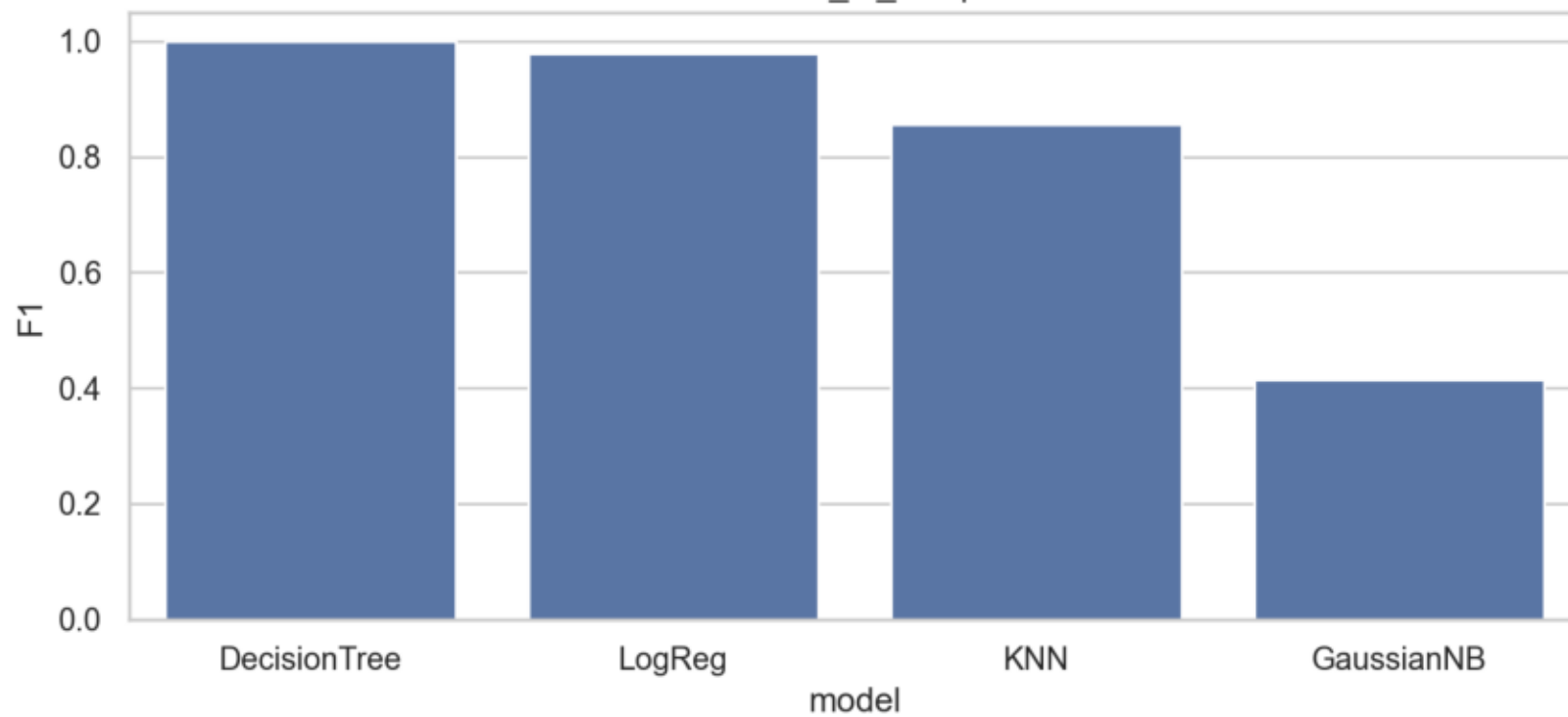






bar_classification_f1_comparison

classification_f1_comparison



scatter_predicted_vs_actual_linear_regression

