

## DSA210 Spring2025 Project Second Draft

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import ast
6 from scipy import stats
7 from itertools import combinations
8
9 def load_and_merge(movies_path: str, credits_path: str) -> pd.DataFrame: 1 usage
10     """Merging two dataset"""
11     movies = pd.read_csv(movies_path)
12     credits = pd.read_csv(credits_path)
13     df = pd.merge(movies, credits, left_on='id', right_on='movie_id', how='inner')
14     return df
15
```

```
16 def preprocess(df: pd.DataFrame) -> pd.DataFrame: 1 usage
17     """Data Process"""
18     df['release_date'] = pd.to_datetime(df['release_date'], errors='coerce')
19     df['year'] = df['release_date'].dt.year
20
21     df['genre_list'] = df['genres'].apply(lambda x: [d['name'] for d in ast.literal_eval(x)])
22     df['primary_genre'] = df['genre_list'].apply(lambda lst: lst[0] if lst else None)
23     df['genre_count'] = df['genre_list'].apply(len)
24
25     df['runtime_bucket'] = pd.cut(
26         df['runtime'],
27         bins=[0, 90, 120, 150, np.inf],
28         labels=['<=90', '91-120', '121-150', '150+']
29     )
30     return df
31
```

```
32 def print_all_pearson(df: pd.DataFrame, cols: list): 1 usage
33     """Print Pearson r and p-value for every unique pair in cols."""
34     print("\n== Pairwise Pearson Correlations with p-values ==")
35     for a, b in combinations(cols, r=2):
36         x, y = df[a], df[b]
37         mask = x.notna() & y.notna()
38         if mask.sum() < 2: continue
39         r, p = stats.pearsonr(x[mask], y[mask])
40         print(f"{a:12s} ↔ {b:12s} : r = {r:6.3f}, p = {p:.2e}")
41
```

```

=== Pairwise Pearson Correlations with p-values ===
budget      ⇔ revenue      : r = 0.731, p = 0.00e+00
budget      ⇔ popularity   : r = 0.505, p = 7.05e-310
budget      ⇔ vote_count   : r = 0.593, p = 0.00e+00
budget      ⇔ vote_average : r = 0.093, p = 9.95e-11
budget      ⇔ runtime      : r = 0.270, p = 6.91e-81
budget      ⇔ genre_count  : r = 0.269, p = 1.66e-80
revenue     ⇔ popularity   : r = 0.645, p = 0.00e+00
revenue     ⇔ vote_count   : r = 0.781, p = 0.00e+00
revenue     ⇔ vote_average : r = 0.197, p = 2.72e-43

revenue     ⇔ runtime      : r = 0.251, p = 6.25e-70
revenue     ⇔ genre_count  : r = 0.182, p = 4.05e-37
popularity  ⇔ vote_count   : r = 0.778, p = 0.00e+00
popularity  ⇔ vote_average : r = 0.274, p = 1.95e-83
popularity  ⇔ runtime      : r = 0.226, p = 2.09e-56
popularity  ⇔ genre_count  : r = 0.155, p = 3.48e-27
vote_count  ⇔ vote_average : r = 0.313, p = 1.19e-109
vote_count  ⇔ runtime      : r = 0.272, p = 3.64e-82
vote_count  ⇔ genre_count  : r = 0.154, p = 7.04e-27
vote_average ⇔ runtime      : r = 0.375, p = 3.29e-160
vote_average ⇔ genre_count : r = 0.086, p = 2.85e-09
runtime     ⇔ genre_count  : r = 0.098, p = 8.78e-12

```

```

42  def eda(df: pd.DataFrame): 1 usage
54      # 2. Descriptive stats
55      numeric = [
56          'budget', 'revenue', 'popularity',
57          'vote_count', 'vote_average', 'runtime', 'genre_count'
58      ]
59      print("\n=== Descriptive Statistics ===")
60      print(df[numeric].describe())
61

```

```

=== Descriptive Statistics ===

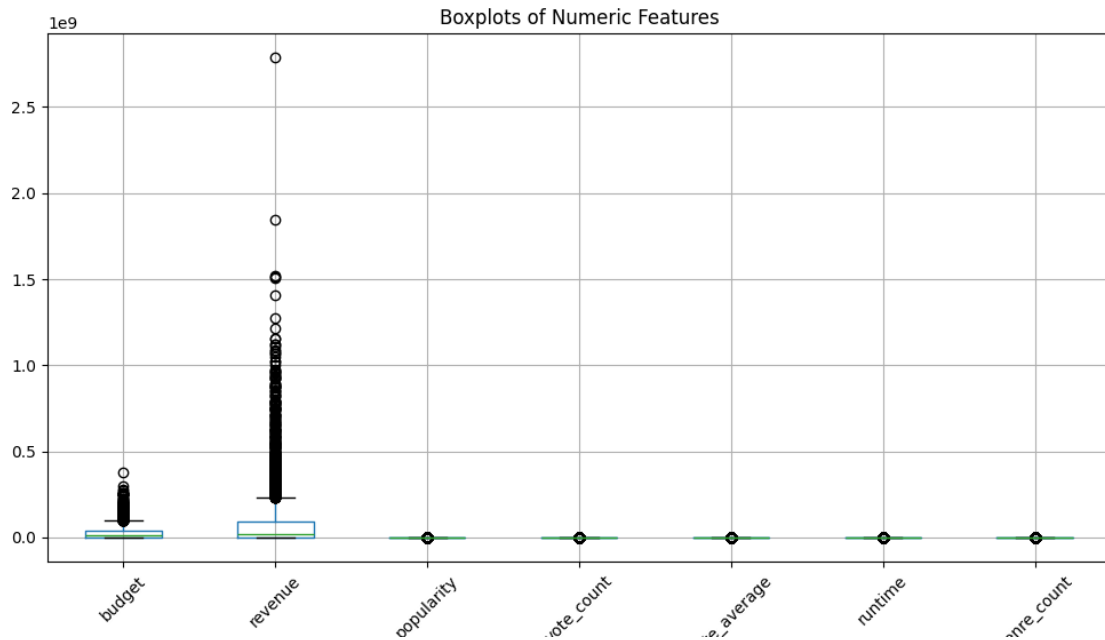
```

	budget	revenue	...	runtime	genre_count
count	4.803000e+03	4.803000e+03	...	4801.000000	4803.000000
mean	2.904504e+07	8.226064e+07	...	106.875859	2.531751
std	4.072239e+07	1.628571e+08	...	22.611935	1.120955
min	0.000000e+00	0.000000e+00	...	0.000000	0.000000
25%	7.900000e+05	0.000000e+00	...	94.000000	2.000000
50%	1.500000e+07	1.917000e+07	...	103.000000	2.000000
75%	4.000000e+07	9.291719e+07	...	118.000000	3.000000
max	3.800000e+08	2.787965e+09	...	338.000000	7.000000

```

42 def eda(df: pd.DataFrame): 1 usage
62     # 3. Boxplots
63     plt.figure(figsize=(12, 6))
64     df[numeric].boxplot()
65     plt.title("Boxplots of Numeric Features")
66     plt.xticks(rotation=45)
67     plt.show()
68

```

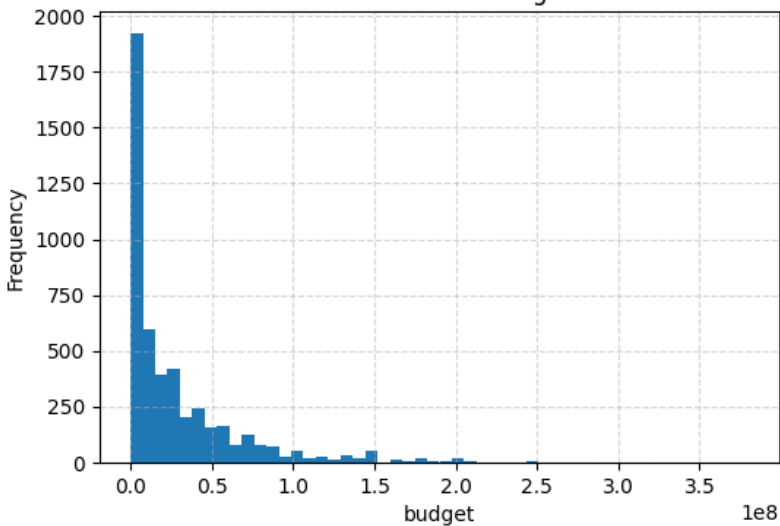


```

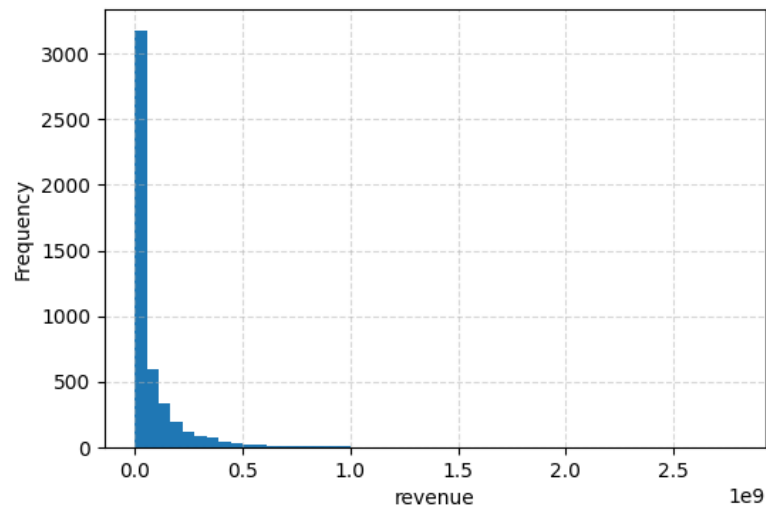
42 def eda(df: pd.DataFrame): 1 usage
69     # 4. Histograms
70     for col in numeric:
71         plt.figure(figsize=(6, 4))
72         df[col].dropna().hist(bins=50)
73         plt.title(f'Distribution of {col}')
74         plt.xlabel(col)
75         plt.ylabel('Frequency')
76         plt.grid(visible=True, linestyle='--', alpha=0.5)
77         plt.show()
78

```

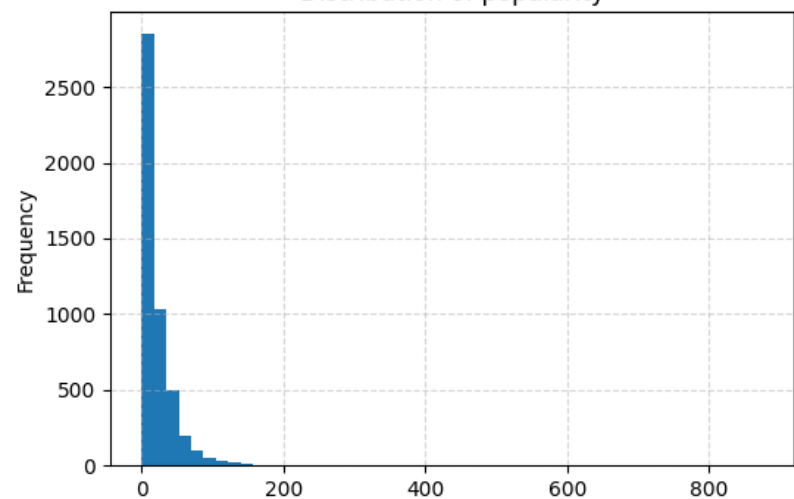
Distribution of budget



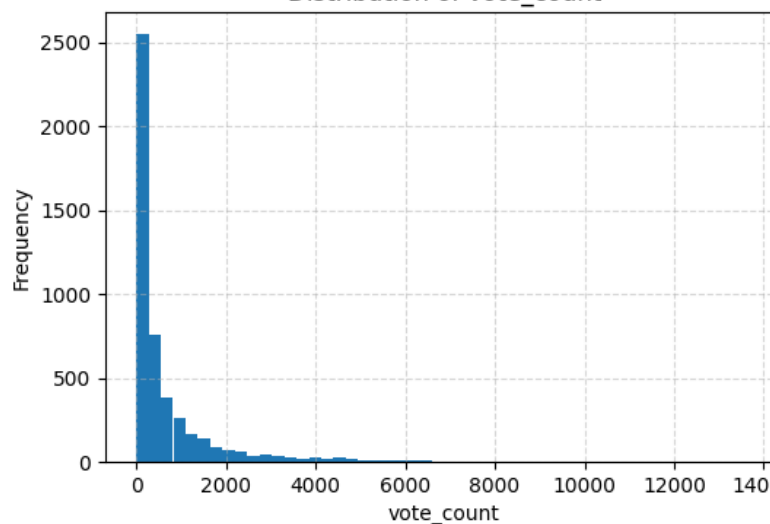
Distribution of revenue



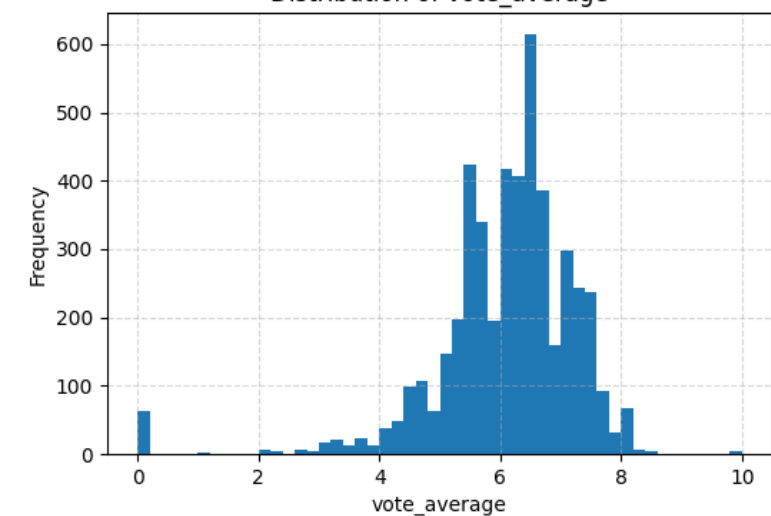
Distribution of popularity



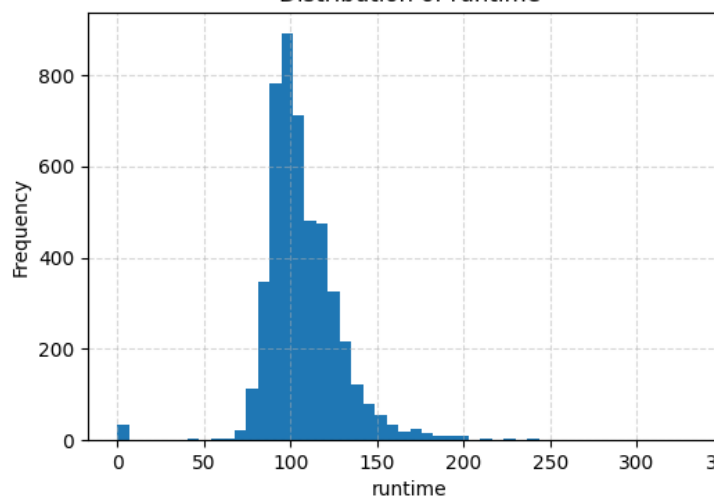
Distribution of vote\_count



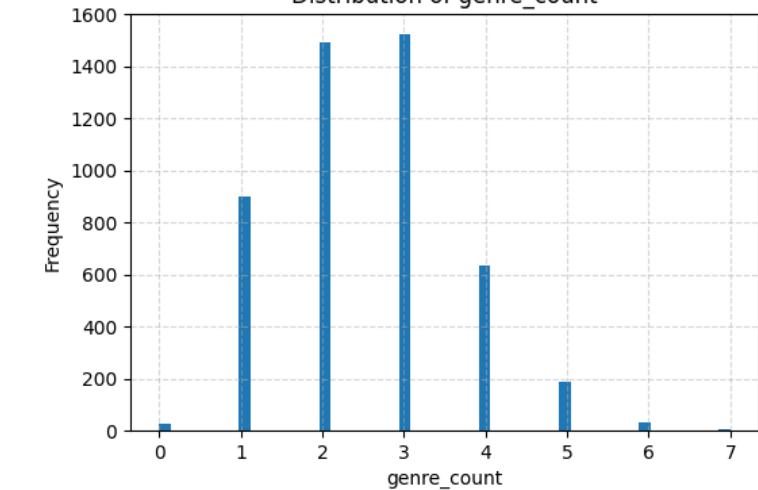
Distribution of vote\_average



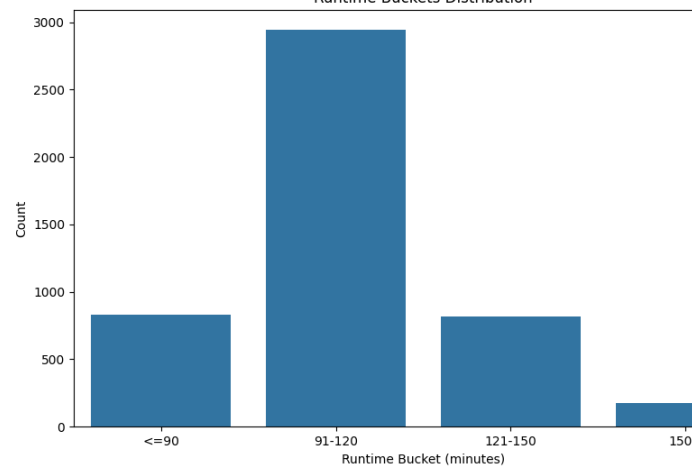
Distribution of runtime



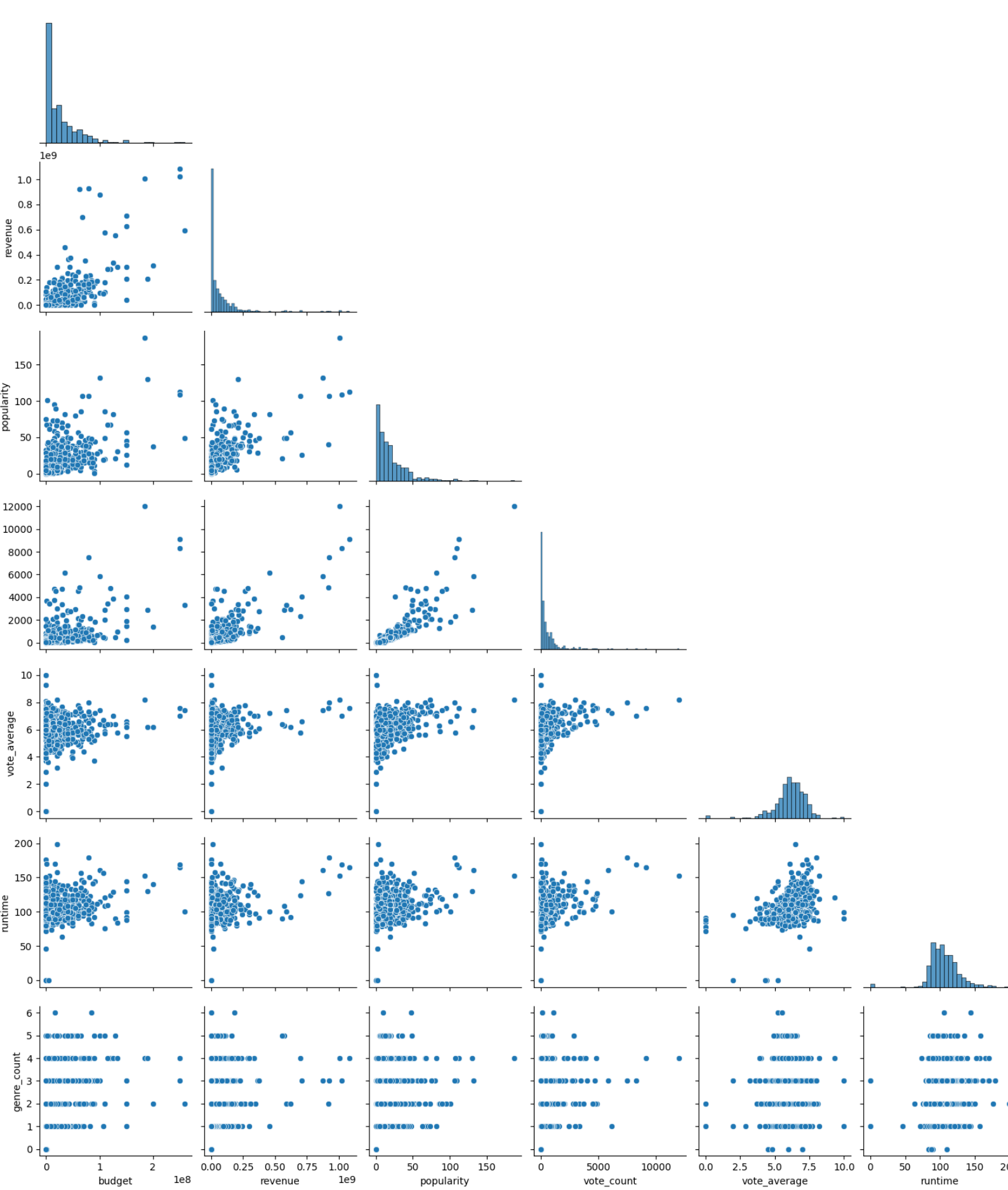
Distribution of genre\_count



Runtime Buckets Distribution



```
# 5. Pairplot (scatter-matrix) of numeric features (sampled)
sample_df = df[numeric].dropna().sample(min(len(df), 500))
g = sns.pairplot(sample_df, corner=True)
g.fig.suptitle('Pairwise Scatterplots (sampled)', y=1.02)
plt.show()
```



```

88 # 6. Correlation matrix & heatmap
89 corr = df[numeric].corr()
90 print("\n=== Correlation Matrix ===")
91 print(corr)
92 plt.figure(figsize=(8, 6))
93 sns.heatmap(corr, annot=True, cmap="coolwarm", fmt=".2f")
94 plt.title("Correlation Heatmap of Numeric Features")
95 plt.show()

```

```

=== Correlation Matrix ===

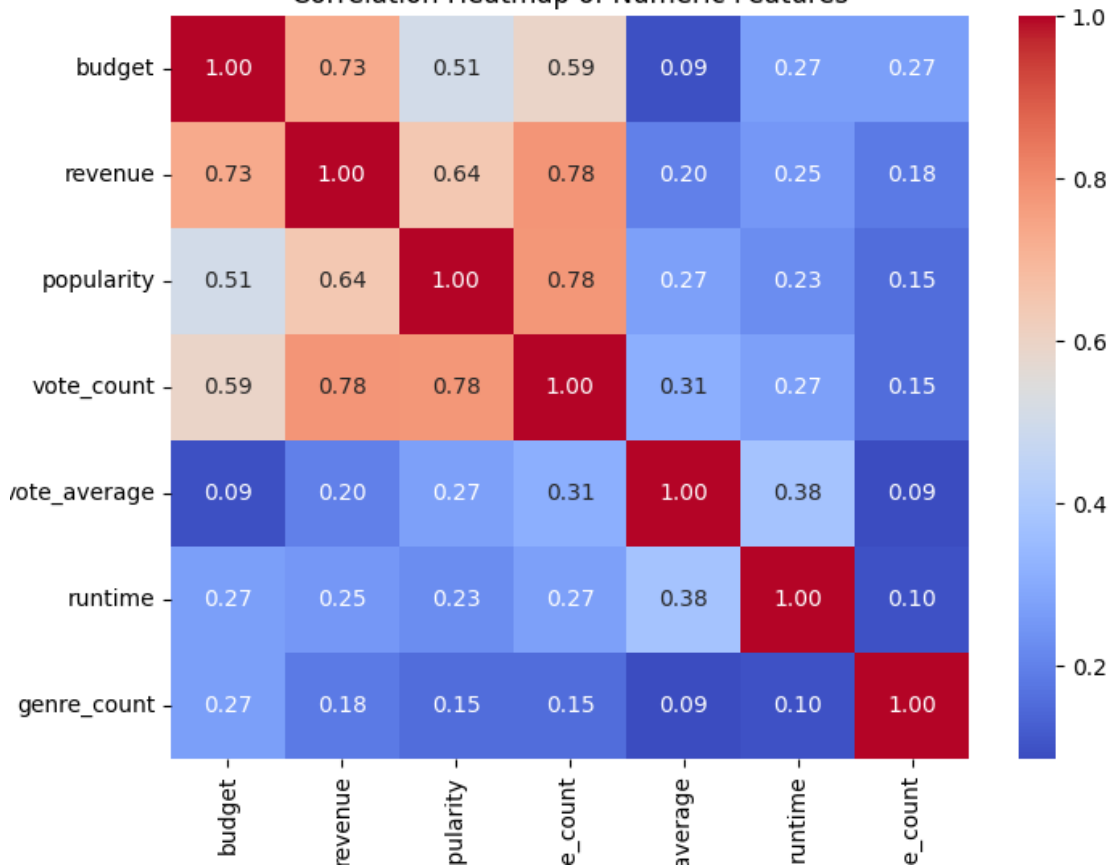
```

```

      budget  revenue  ...  runtime  genre_count
budget    1.000000  0.730823  ...  0.269851    0.269170
revenue    0.730823  1.000000  ...  0.251093    0.182185
popularity  0.505414  0.644724  ...  0.225502    0.154918
vote_count  0.593180  0.781487  ...  0.271944    0.154000
vote_average 0.093146  0.197150  ...  0.375046    0.085577
runtime    0.269851  0.251093  ...  1.000000    0.098290
genre_count 0.269170  0.182185  ...  0.098290    1.000000

```

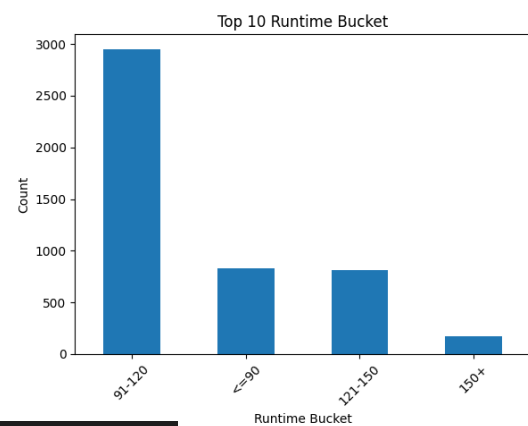
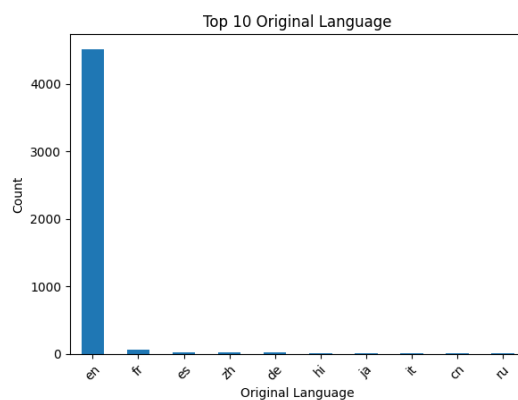
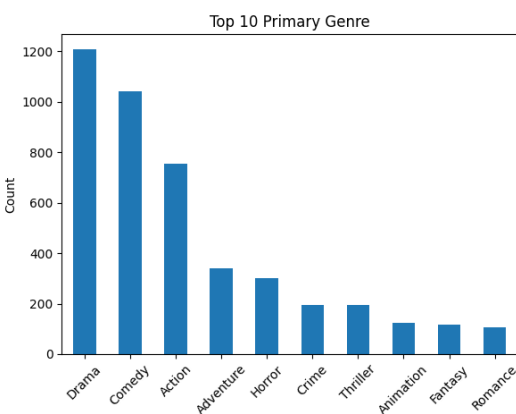
Correlation Heatmap of Numeric Features



```

100     # 7. Categorical distributions
101     cats = ['primary_genre', 'original_language', 'runtime_bucket']
102     fig, axes = plt.subplots(nrows=1, len(cats), figsize=(18, 5))
103     for ax, var in zip(axes, cats):
104         counts = df[var].value_counts().nlargest(10)
105         counts.plot(kind='bar', ax=ax)
106         ax.set_title(f'Top 10 {var.replace(__old: "_", __new: " ").title()}')
107         ax.set_xlabel(var.replace(__old: "_", __new: " ").title())
108         ax.set_ylabel("Count")
109         ax.tick_params(axis='x', rotation=45)
110     plt.tight_layout()
111     plt.show()
112

```

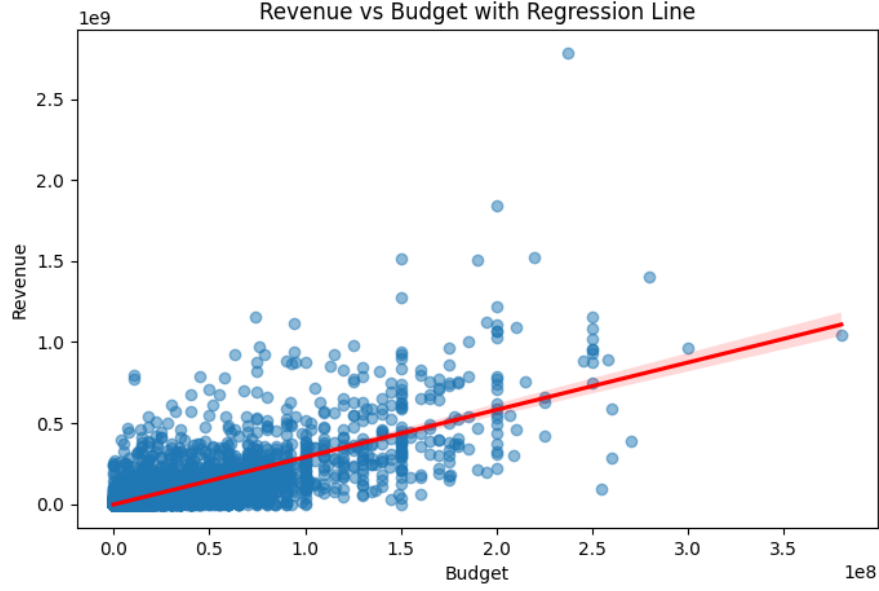


```

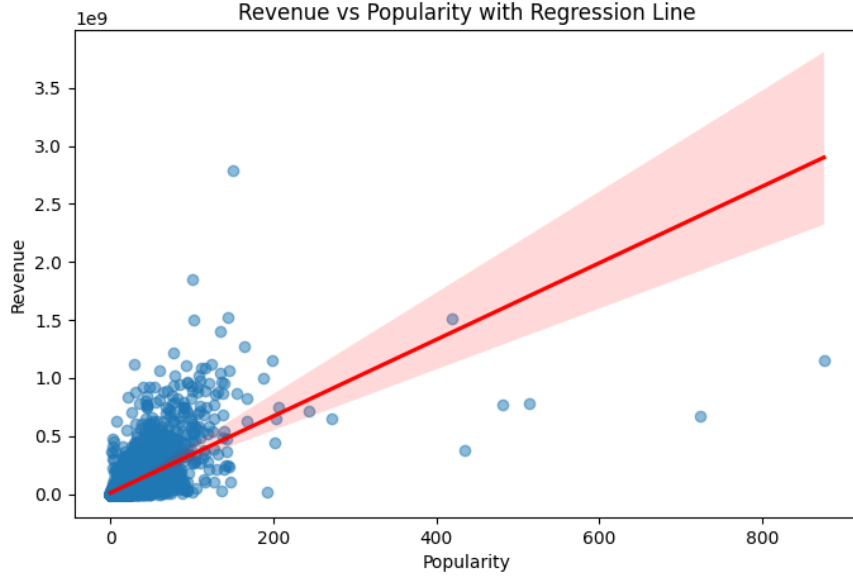
45     def eda(df: pd.DataFrame): 1 usage
113     # 8. Scatter plots with regression lines
114     pairs = [
115         ('budget', 'revenue'),
116         ('popularity', 'revenue'),
117         ('vote_average', 'revenue')
118     ]
119     for x, y in pairs:
120         plt.figure(figsize=(8, 5))
121         sns.regplot(
122             x=x, y=y, data=df,
123             scatter_kws={'alpha':0.5},
124             line_kws={'color':'red'}
125         )
126         plt.title(f"{y.title()} vs {x.title()} with Regression Line")
127         plt.xlabel(x.replace(__old: '_', __new: ' ').title())
128         plt.ylabel(y.replace(__old: '_', __new: ' ').title())
129         plt.show()
130

```

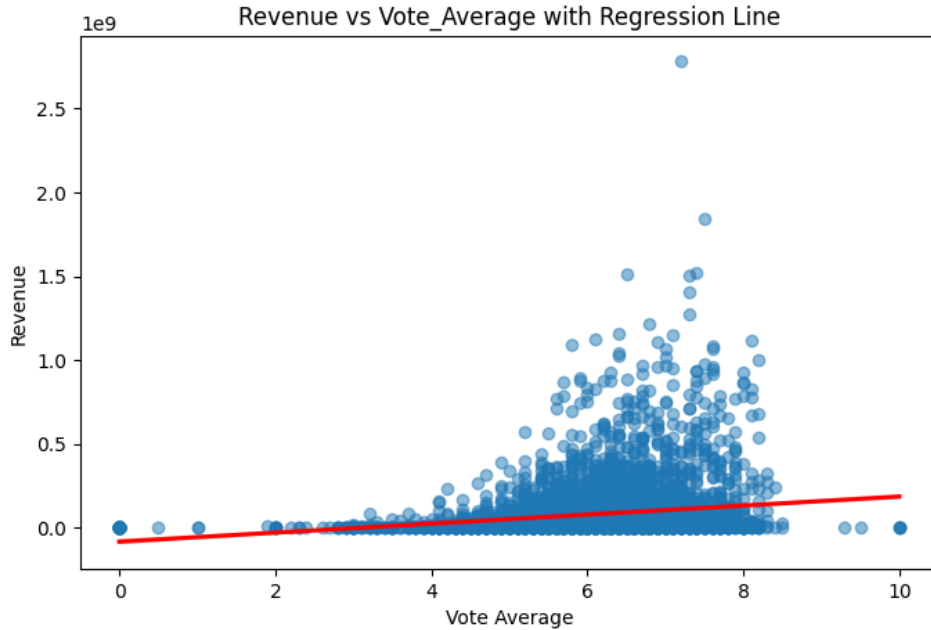
Revenue vs Budget with Regression Line



Revenue vs Popularity with Regression Line



Revenue vs Vote\_Average with Regression Line





```

130
131 def hypothesis_tests(df: pd.DataFrame): 1 usage
132     """Run hypothesis tests with clear null/alternative statements."""
133     clean = df[['budget', 'revenue']].dropna()
134     r, p = stats.pearsonr(clean['budget'], clean['revenue'])
135     print("\n--- Hypothesis Test 1: Budget vs Revenue ---")
136     print("H0:  $\rho = 0$  (no correlation)\nHa:  $\rho \neq 0$  (non-zero correlation)")
137     print(f"Pearson r = {r:.3f}, p-value = {p:.3e}")
138     print("→", "Reject H0" if p < 0.05 else "Fail to reject H0")
139
140     median_rating = df['vote_average'].median()
141     grp_high = df[df['vote_average'] > median_rating]['revenue'].dropna()
142     grp_low = df[df['vote_average'] <= median_rating]['revenue'].dropna()
143     t_stat, p_val = stats.ttest_ind(grp_high, grp_low, equal_var=False)
144     print("\n--- Hypothesis Test 2: High vs Low Rating Revenue ---")
145     print("H0:  $\mu_{\text{high}} = \mu_{\text{low}}$ \nHa:  $\mu_{\text{high}} \neq \mu_{\text{low}}$ ")
146     print(f"t-statistic = {t_stat:.3f}, p-value = {p_val:.3e}")
147     print("→", "Reject H0" if p_val < 0.05 else "Fail to reject H0")
148

```

```

--- Hypothesis Test 1: Budget vs Revenue ---
H0:  $\rho = 0$  (no correlation)
Ha:  $\rho \neq 0$  (non-zero correlation)
Pearson r = 0.731, p-value = 0.000e+00
→ Reject H0

```

```

--- Hypothesis Test 2: High vs Low Rating Revenue ---
H0:  $\mu_{\text{high}} = \mu_{\text{low}}$ 
Ha:  $\mu_{\text{high}} \neq \mu_{\text{low}}$ 
t-statistic = 10.824, p-value = 6.837e-27
→ Reject H0

```

## Results from the Exploratory Data Analysis and Tests

### Predictor | r with Revenue | Interpretations

**1) Vote Count** |  $\approx 0.78$  | Very strong positive association. Films that more people vote on tend to earn much more—both a signal and an outcome of popularity.

**2) Budget** |  $\approx 0.73$  | Strong positive link. Bigger investment generally yields bigger gross.

**3) Popularity** |  $\approx 0.64$  | Moderate-strong. TMDB's composite popularity score tracks revenue well.

**4) Vote Average** |  $\approx 0.20$  | Weak positive. Higher average ratings alone aren't a great predictor of gross.

**5) Runtime** |  $\approx 0.10$ – $0.20$  (approximately) | Slight positive; standard feature lengths (90–150 min) do tend to earn more.

**6) Genre Count** |  $\approx 0.05$  | Virtually no linear trend between “how many genres” a movie is tagged with and its revenue.

# Hypothesis Tests

## 1. Budget ↔ Revenue

- **H<sub>0</sub>:** No correlation ( $\rho = 0$ )
- **Result:**  $r \approx 0.73$ ,  $p \ll 0.001 \rightarrow$  **Reject H<sub>0</sub>**
- **Conclusion:** There is a highly significant, strong linear relationship between budget and revenue.

## 2. High vs. Low Rating Revenue

- Splitting our sample at the median for IMDb-style rating.
- **H<sub>0</sub>:** Mean revenue of “high-rating” equals “low-rating” group
- **Result:**  $t \approx 10.8$ ,  $p \approx 6.8 \times 10^{-27} \rightarrow$  **Reject H<sub>0</sub>**
- **Conclusion:** Even though the  $r$  is small ( $\approx 0.20$ ), movies with above-median ratings earn significantly more on average than lower-rated ones.

- Also further hypothesis tests can be concluded for seeing relationships between other entities too.

## Overall Conclusions From The Analysis

**Vote\_count** (how many users rate the film) and **popularity** are the strongest correlates of box-office revenue. They combine both pre- and post-release results .

**Budget** remains a powerful predictor: larger production and marketing budgets yield higher revenues.

**Ratings** (vote\_average) play a mediocre role. While they don't linearly track revenue as tightly, the t-test shows that better-rated movies still deliver a clear boost in mean earnings and revenue.

**Runtime** and **genre count** have minimal linear effects.