

# IMDb Top 250 Movies Analysis

```
import pandas as pd
import matplotlib.pyplot as plt
import re

# Load Excel file
df = pd.read_excel("IMDb_Top_250.xlsx")
df.head()

# Function to extract numeric rating and vote count
def extract_rating_votes(rating_str):
    rating_str = rating_str.replace('\xa0', ' ') # Clean non-breaking space
    match = re.match(r"([\d.]+\s*([(\d.]+)([MK])\s*)", rating_str)
    if match:
        rating = float(match.group(1))
        count = float(match.group(2).replace(',', ''))
        multiplier = {'K': 1e3, 'M': 1e6}
        votes = count * multiplier[match.group(3)]
        return rating, int(votes)
    return None, None

df[['RatingValue', 'VoteCount']] = df['Rating'].apply(lambda x:
pd.Series(extract_rating_votes(x)))

# Convert Duration (e.g., "2h 22m") to minutes
def duration_to_minutes(duration_str):
    match = re.match(r'(\d+)h\s*(\d+)m', duration_str)
    if match:
        return int(match.group(1)) * 60 + int(match.group(2))
    return None

df['DurationMin'] = df['Duration'].apply(duration_to_minutes)
df.head()

plt.figure(figsize=(10, 6))
plt.scatter(df['Year'], df['RatingValue'], alpha=0.7, color='skyblue', edgecolor='gray')
plt.title("IMDb Ratings Over the Years")
plt.xlabel("Year")
plt.ylabel("Rating")
plt.grid(True)
plt.tight_layout()
plt.show()

top_voted = df.sort_values(by='VoteCount', ascending=False).head(20)

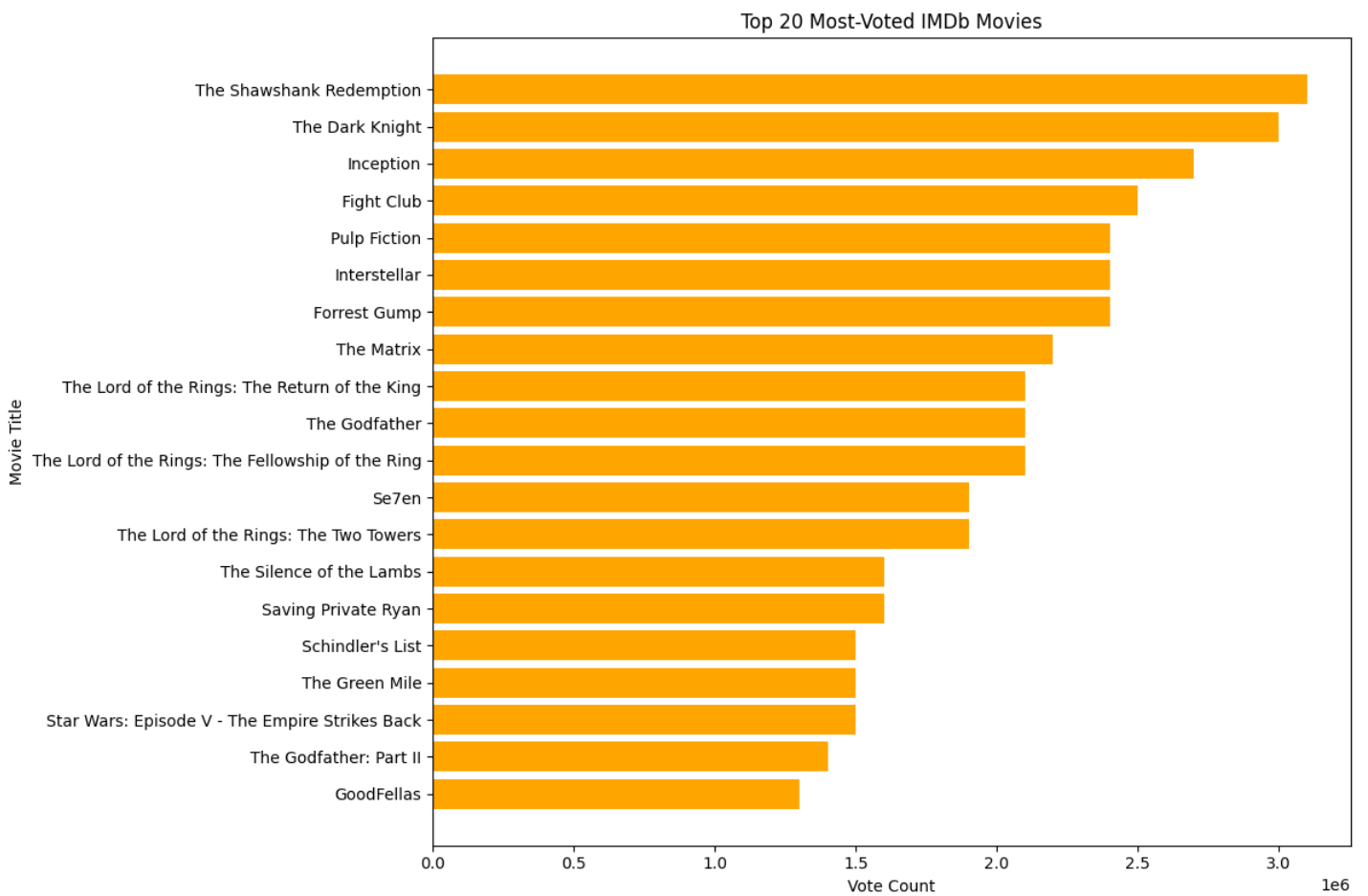
plt.figure(figsize=(12, 8))
plt.barh(top_voted['Title'][::-1], top_voted['VoteCount'][::-1], color='orange')
plt.title("Top 20 Most-Voted IMDb Movies")
plt.xlabel("Vote Count")
plt.ylabel("Movie Title")
plt.tight_layout()
```

```
plt.show()
```

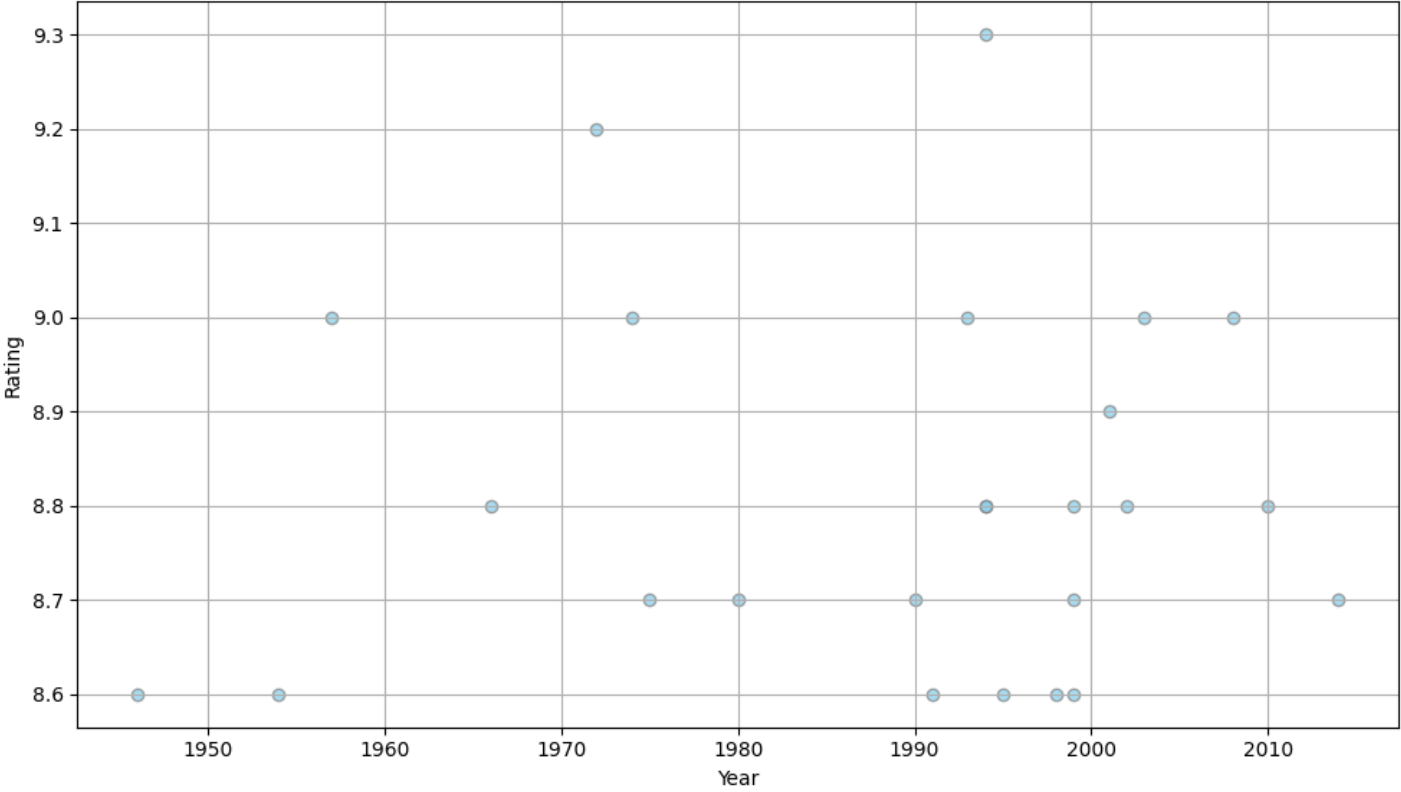
```
plt.figure(figsize=(10, 6))
plt.scatter(df['DurationMin'], df['RatingValue'], alpha=0.7, color='green')
plt.title("Duration vs IMDb Rating")
plt.xlabel("Duration (minutes)")
plt.ylabel("Rating")
plt.grid(True)
plt.tight_layout()
plt.show()
```

```
df['Decade'] = (df['Year'] // 10) * 10
decade_rating = df.groupby('Decade')['RatingValue'].mean().reset_index()
```

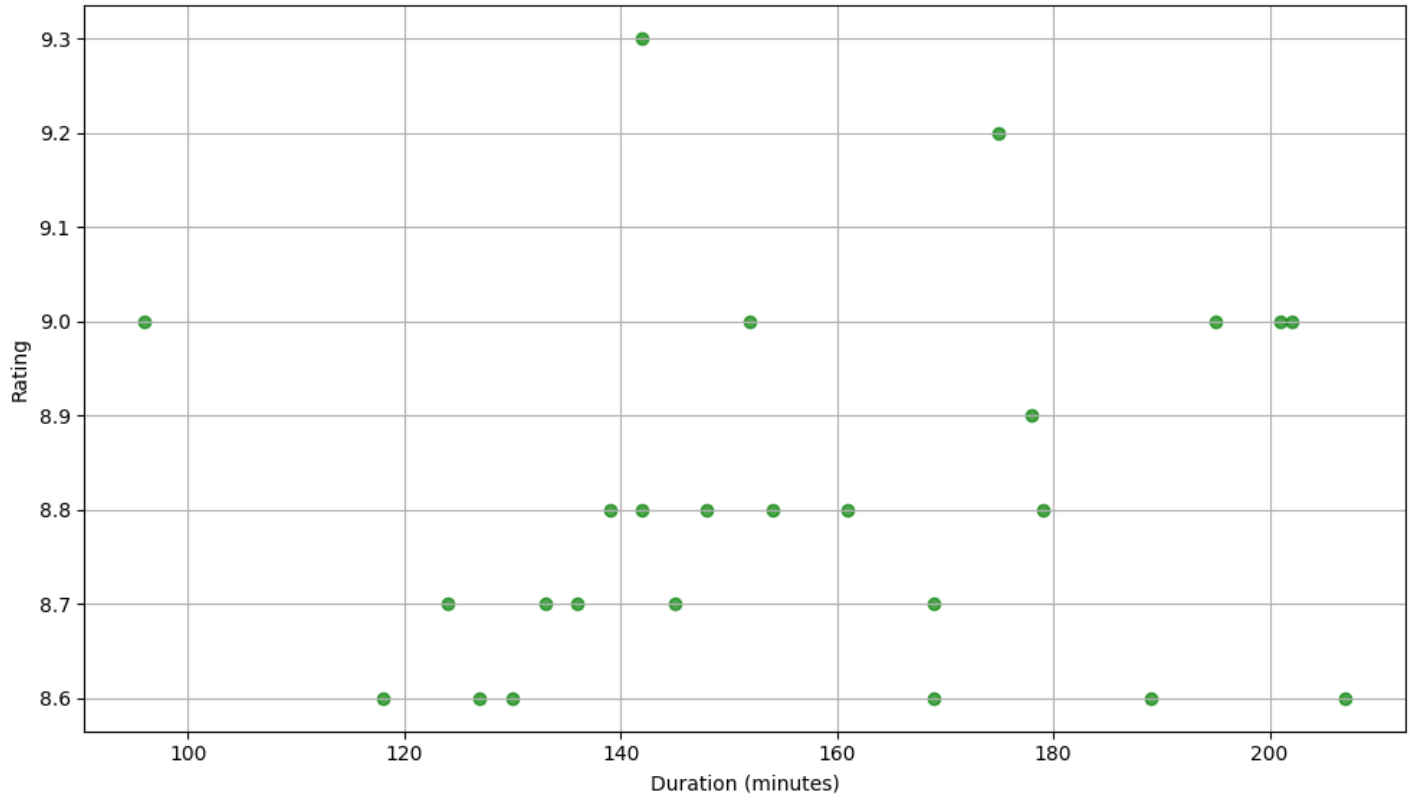
```
plt.figure(figsize=(10, 6))
plt.plot(decade_rating['Decade'], decade_rating['RatingValue'], marker='o',
         linestyle='-', color='purple')
plt.title("Average IMDb Rating Per Decade")
plt.xlabel("Decade")
plt.ylabel("Average Rating")
plt.grid(True)
plt.tight_layout()
plt.show()
```



IMDb Ratings Over the Years



Duration vs IMDb Rating



Average IMDb Rating Per Decade

