

Movie recommendation system

Ahmed Kedić
Džanan Habibija
Anđela Jeftović

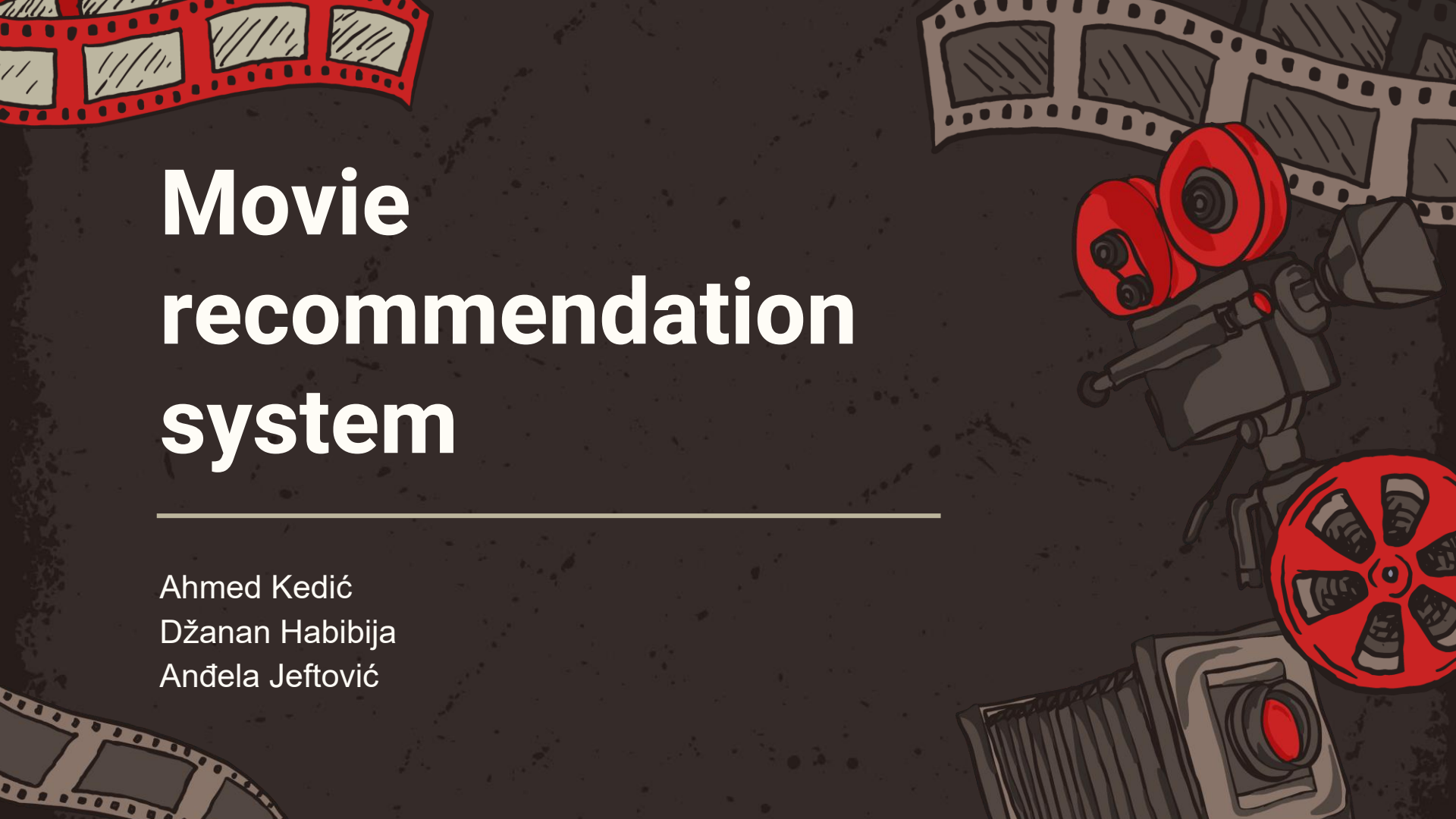


Table of contents

- 01 Introduction**
- 02 Literature Review**
- 03 Dataset**
- 04 Data Cleaning and Preprocessing**
- 05 Data Analysis**
- 06 Machine Learning Algorithms**
- 07 Conclusion**



The background is a dark, textured grey with a speckled pattern. In the corners, there are stylized, hand-drawn illustrations of film-related items. Top-left: a red film reel. Top-right: a red film strip with white frames, a red camera with two reels, and another red film reel. Bottom-left: a grey film strip. Bottom-right: a grey camera with two reels. The text '01' is in a light yellow, sans-serif font, and 'Introduction' is in a white, sans-serif font, both centered on the left side of the image.

01

Introduction

Introduction

- What is movie recommendation system?
- Collaborative filtering VS Content-based filtering
- How did we implement a movie recommendation system?



The background is a dark, textured grey. In the corners, there are stylized, hand-drawn illustrations of film-related items. Top-left: a red film reel. Top-right: a red film strip with two red reels and a black camera. Bottom-left: a grey film strip. Bottom-right: a grey film reel, a red film reel, and a black camera. The text is centered in the middle of the frame.

02

Literature Review



Social and economic landscape

Content-based filtering

- Recommend movies based on attributes
- Solves cold-start problem
- Personalized recommendations

MOVREC System

- Combines collaborative filtering + K-means clustering
- Uses IMDb dataset
- Reduces information overload


Hybrid approaches

- Collaborative, content-based, hybrid systems
- Cold-start, sparse ratings, large-scale data
- Suggests novel similarity metric

Hotel recommendation system

- Collaborative filtering, NLP, classification
- High precision and recall with large datasets
- ML applicability across domains

Systematic review

- Covers algorithms: filtering, clustering, metaheuristics
 - Datasets: MovieLens, IMDb, Netflix
 - Hybrid systems, future trends
- 

03 Dataset




Dataset

For our movie recommendation system we are using 9000+ Movies Dataset.

This dataset contains following columns:

- Release_Date
- Title
- Overview
- Popularity
- Vote_Count
- Vote_Average
- Original_Language
- Genre
- Poster_Url



The background is a dark, textured grey. In the top-left corner is a red film reel. In the top-right corner is a red movie camera with a red film strip passing through it. In the bottom-left corner is a grey film strip. In the bottom-right corner is a grey film reel. The text '04' is in a light yellow font, and 'Data Cleaning and Preprocessing' is in a white font, both centered on the left side of the image.

04

Data Cleaning and Preprocessing



Data cleaning and preprocessing

<code>.info()</code> <code>.head()</code> <code>.describe()</code>	Explored dataset
<code>.isnull().sum()</code> <code>.duplicated().any()</code>	Checked for missing values and duplicates
<code>_fillna()</code>	Replaced missing Vote_Count with column mean.
<code>.median()</code> <code>.mean()</code> <code>.mode()</code>	Analyzed distribution: mean, median, mode, quartiles
<code>.min()</code> <code>.max()</code>	Scaling: Min-Max scaling for numerical features



Data cleaning and preprocessing – Scaling

```
import pandas as pd
from sklearn.preprocessing import minmax_scale
import seaborn as sns
import matplotlib.pyplot as plt
from scipy import stats

data['Vote_Count'] = pd.to_numeric(data['Vote_Count'], errors='coerce')
data['Vote_Average'] = pd.to_numeric(data['Vote_Average'], errors='coerce')

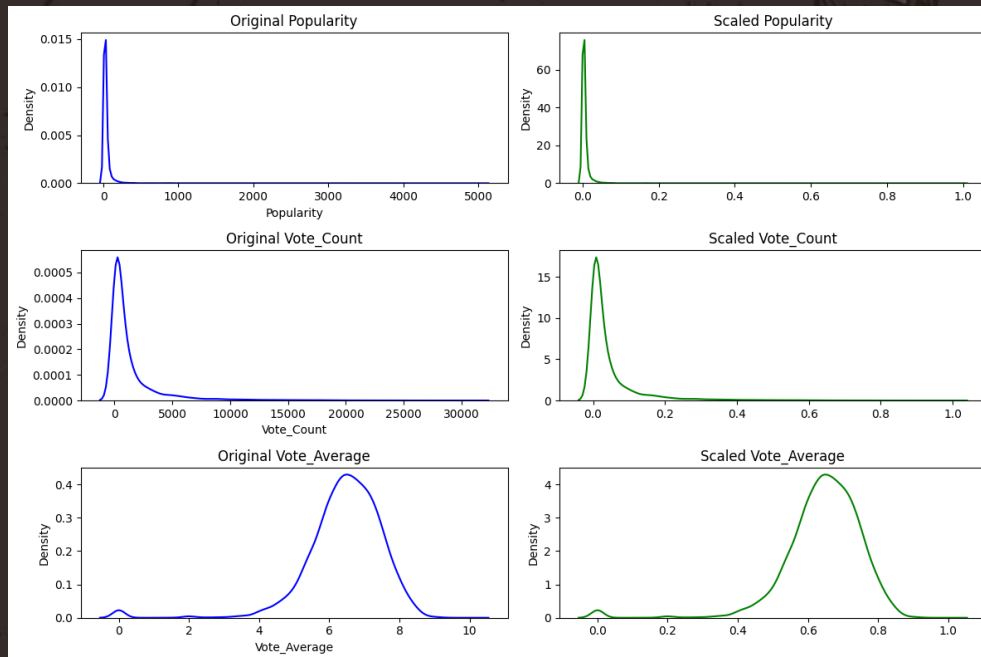
columns_to_scale = ['Popularity', 'Vote_Count', 'Vote_Average']

scaled_data = minmax_scale(data[columns_to_scale], axis=0)

fig, ax = plt.subplots(len(columns_to_scale), 2, figsize=(12, 8))

for i, col in enumerate(columns_to_scale):
    sns.kdeplot(data[col], ax=ax[i][0], color='blue')
    ax[i][0].set_title(f"Original {col}")
    sns.kdeplot(scaled_data[:, i], ax=ax[i][1], color='green')
    ax[i][1].set_title(f"Scaled {col}")

plt.tight_layout()
plt.show()
```





Data cleaning and preprocessing – Normalization

```
from scipy import stats
```

```
positive_popularity = data['Popularity'][data['Popularity'] > 0]
```

```
normalized_popularity, _ = stats.boxcox(positive_popularity)
```

```
fig, ax = plt.subplots(1, 2, figsize=(12, 6))
```

```
sns.kdeplot(positive_popularity, ax=ax[0], color='blue')
```

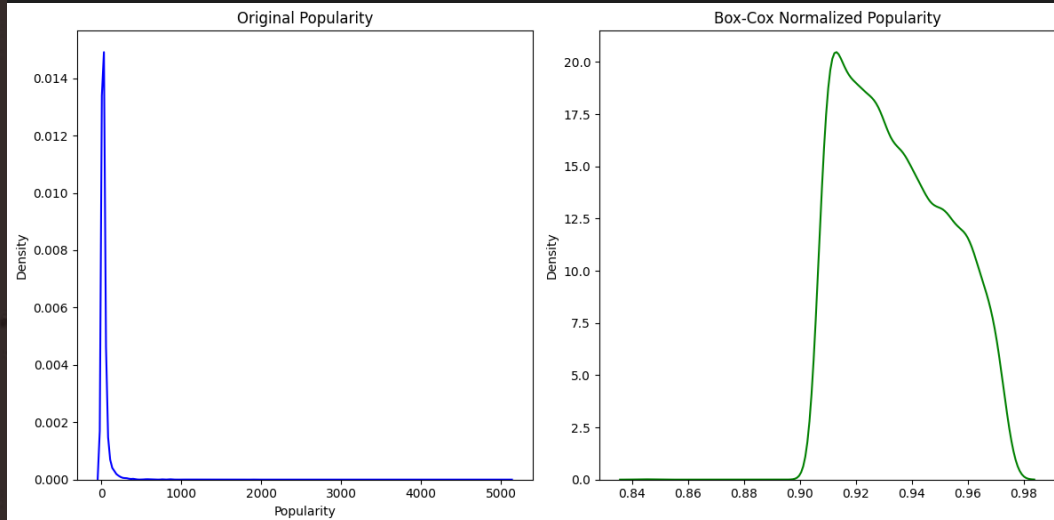
```
ax[0].set_title("Original Popularity")
```

```
sns.kdeplot(normalized_popularity, ax=ax[1], color='green')
```

```
ax[1].set_title("Box-Cox Normalized Popularity")
```

```
plt.tight_layout()
```

```
plt.show()
```



A stylized illustration of a film reel and a camera, rendered in a comic book style with bold lines and a limited color palette of red, black, and white. The reel is positioned in the upper left corner, and the camera is in the upper right corner. The background is a dark, textured grey with scattered white specks, resembling a film negative or a starry sky.

05

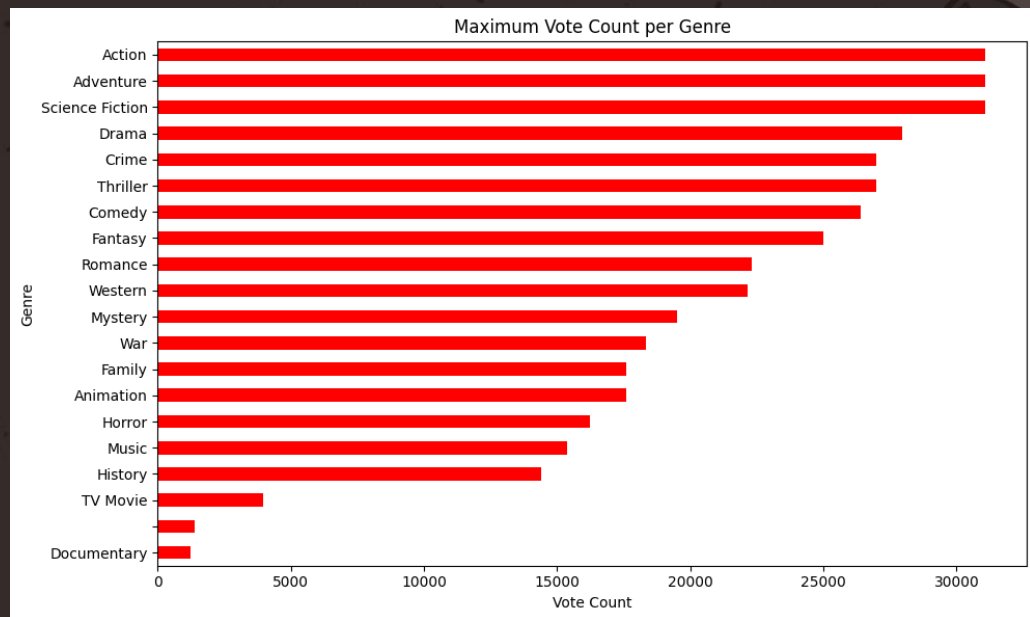
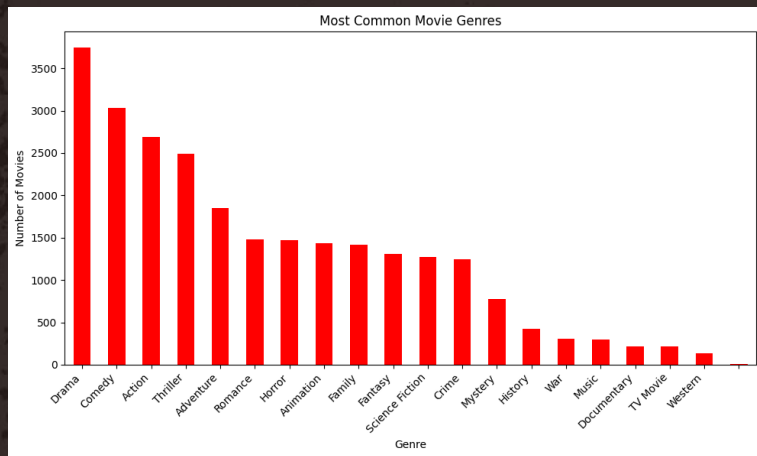
Data Analysis

Data analysis overview



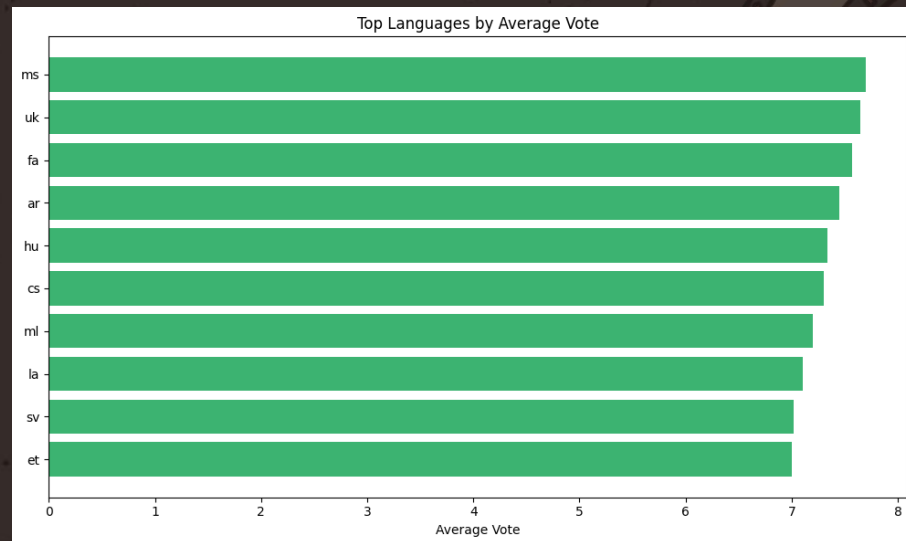
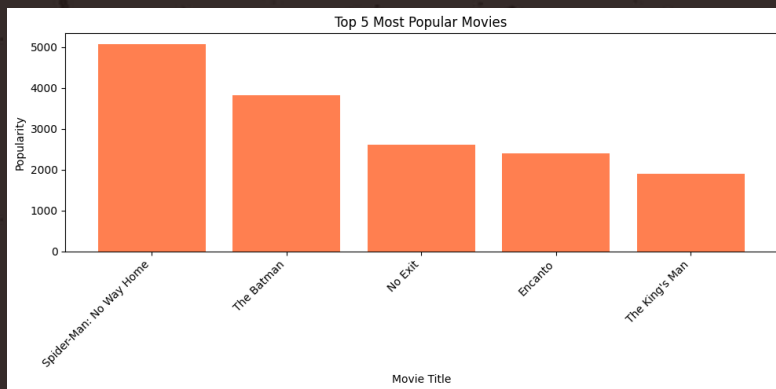
Genre analysis

- Genre popularity analysis
- Maximum vote count per genre
- Biggest vote average per genre
- Most common genres
- Top rated movies per genre



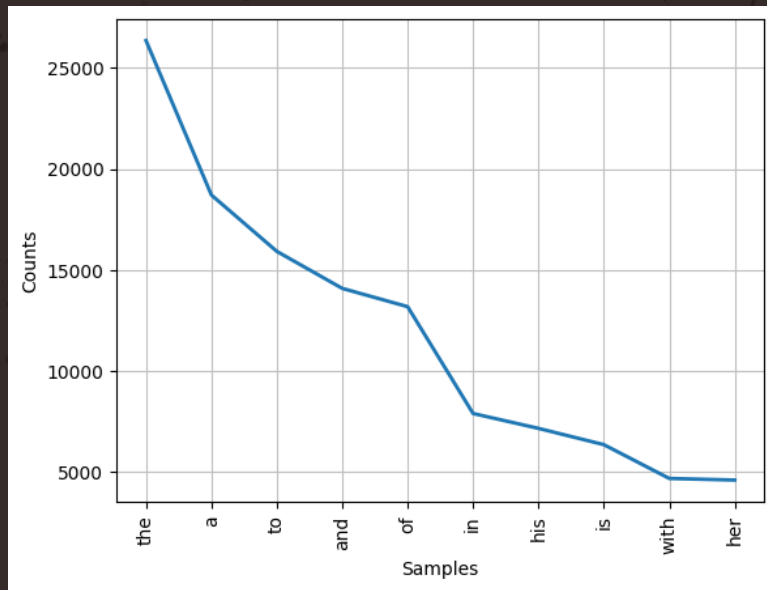
Popularity, top movies and language insight

- Popularity per year
- Top 5 most popular movies per language
- The most popular movie for each language
- Top largest movies with largest popularity
- Top rated movies per genre
- Top languages by average vote



Keywords analysis

- Most Common Words in Movie Overviews





06

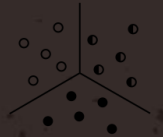
Machine Learning Algorithms



Machine learning algorithms



- K-Nearest Neighbors (KNN)



- KMeans Clustering



- Random Forest Classifier



K-Nearest Neighbors (KNN)

- Based on movie **overview text** using **TF-IDF + Cosine Similarity**
- Returns top 5 similar movies for any input movie
- Simple and effective for small/medium datasets

```
def recommend_movies(title, num_recommendations=5):  
    title = title.lower()  
    if title not in indices:  
        return f"Movie '{title}' not found in dataset."  
  
    idx = indices[title]  
    distances, indices_rec = nn_model.kneighbors(tfidf_matrix[idx], n_neighbors=num_recommendations + 1)  
    recommended_indices = indices_rec[0][1:] # Exclude the input movie itself  
  
    return data[['Title', 'Overview']].iloc[recommended_indices]
```

```
recommend_movies("Spider-man")
```

	Title	Overview
7949	Beyond the Ultimate Spin: The Making of 'Spide...	Documentary on the making of 'Spider-Man.'
201	Spider-Man 3	The seemingly invincible Spider-Man goes up ag...
1500	Spider-Man	When an extortionist threatens to force a mult...
168	Spider-Man: Homecoming	Following the events of Captain America: Civil...
191	Spider-Man: Into the Spider-Verse	Miles Morales is juggling his life between bei...

KNN Evaluation

```
def evaluate_knn_recommender(title, num_recommendations=5):
    title = title.lower()
    if title not in indices:
        return f"Movie '{title}' not found in dataset."

    idx = indices[title]
    distances, indices_rec = nn_model.kneighbors(tfidf_matrix[idx], n_neighbors=num_recommendations + 1)
    recommended_indices = indices_rec[0][1:]

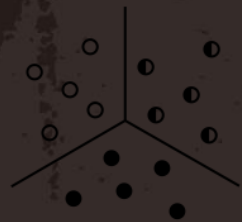
    # Similarity scores (1 - distance)
    similarities = 1 - distances[0][1:]
    avg_similarity = np.mean(similarities)

    print(f"\nKNN Recommender Evaluation for '{title}':")
    print(f"Average Cosine Similarity to Recommended Movies: {avg_similarity:.3f}")
    return data[['Title', 'Overview']].iloc[recommended_indices]

evaluate_knn_recommender("Parasyte: Part 1")
```

KNN Recommender Evaluation for 'parasyte: part 1':
Average Cosine Similarity to Recommended Movies: 0.334

	Title	Overview
1746	Parasyte: Part 2	Alien pods come to Earth and, naturally, start...
9731	Crayon Shin-chan: Fierceness That Invites Stor...	The Noharas get abducted by aliens who claim t...
7158	Alien Nation	A few years from now, Earth will have the firs...
3649	Attraction	After an alien ship crash lands in a Russian c...
9627	Crayon Shin-chan: Action Mask vs. Leotard Devil	Everyone's favorite TV superhero Action Mask s...



KMeans Clustering

- Used MultiLabelBinarizer to convert genres to binary matrix
- Grouped movies into **5 clusters** based on genre similarity
- Recommends movies from the same cluster

```
score = silhouette_score(genre_matrix, data['Genre_Cluster'])  
print(f"Silhouette Score: {score:.3f}")
```

Silhouette Score: 0.217

```
def recommend_by_genre_cluster(title, n=5):  
    title_lower = title.lower()  
    matches = data.index[data['Title'].str.lower() == title_lower]  
    if len(matches) == 0:  
        return f"Movie '{title}' not found."  
  
    idx = matches[0]  
    cluster_id = data.at[idx, 'Genre_Cluster']  
  
    peers = data[data['Genre_Cluster'] == cluster_id].drop(idx)  
    if peers.empty:  
        return f"No other movies in cluster {cluster_id}."  
  
    sample = peers.sample(min(n, len(peers)), random_state=42)  
    return sample[['Title', 'Genre_list']]
```

```
recommend_by_genre_cluster("The Batman", n=5)
```

	Title	Genre_list
1885	Friday the 13th Part 2	[Horror, Thriller]
7559	Sexy Beast	[Crime, Drama, Thriller]
2514	Run and Gun	[Action, Thriller]
4421	The Silence	[Horror, Drama, Thriller, Fantasy]
6791	Psycho II	[Horror, Mystery, Thriller]



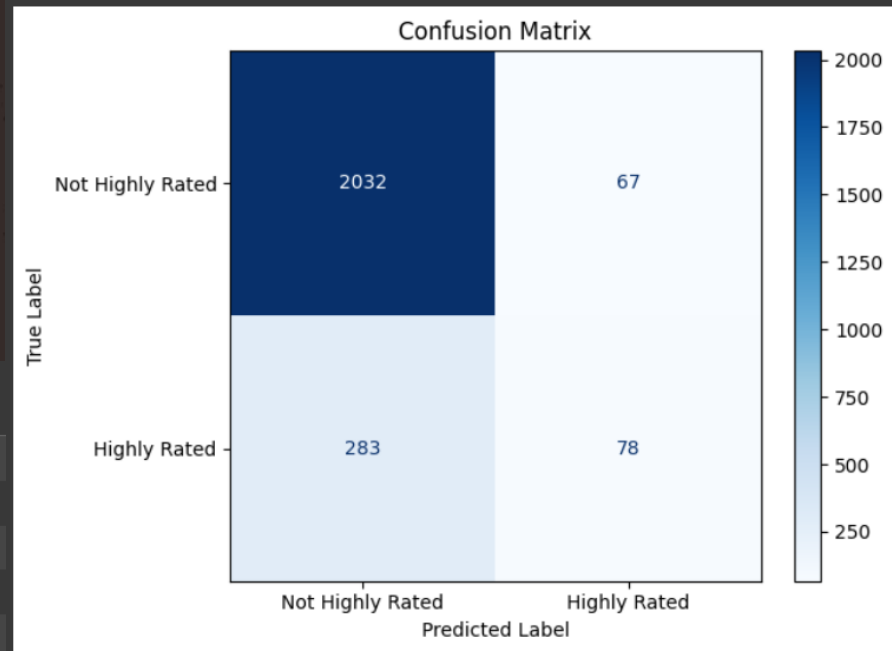
Random Forest Classifier

- Supervised learning to predict if a movie is **Highly Rated** (≥ 7.5)
- Used features: Popularity, Vote Count, Genre
- Trained on 75% of data, tested on 25%

Recommendations for 'Batman':

	Title	Overview
260	Batman Returns	While Batman deals with a deformed man calling...
1639	The Batman vs. Dracula	Gotham City is terrorized not only by recent e...
2162	Batman Beyond: Return of the Joker	The Joker is back with a vengeance, and Gotham...
221	The Dark Knight	Batman raises the stakes in his war on crime. ...
2487	Batman Unmasked: The Psychology of the Dark Kn...	Delve into the world of Batman and the vigilan...

	precision	recall	f1-score	support
0	0.88	0.97	0.92	2099
1	0.54	0.22	0.31	361
accuracy			0.86	2460
macro avg	0.71	0.59	0.61	2460
weighted avg	0.83	0.86	0.83	2460



ML Algorithms Comparison



	KNN	KMeans	Random Foresr
Goal	Recommended similar movies based on Overviews.	Clear understanding with well-developed ideas	Predict whether a movie will be highly rated
Type	Unsupervised	Unsupervised	Supervised
Personalization	Yes	Yes	No
Output	Top N movies based on similarities	Top N movies in same cluster	Wheather movie is highly rated or not
Performance	Good	Clusters are overlapping	Good prediction
Strengths	Fine-grained, story-based similarity. Great for personalized recommendations	Captures genre-based groupings. Simple, fast	Predictive power Uses multiple features -Interpretable output
Weaknesses	Ignores genre or popularity. Doesn't explain why results are similar	Clusters may be ambiguous due to genre overlap. Less personalized	Needs labeled data. Not used for content recommendation





07

Conclusion



Conclusion

Each model had its own purpose and strengths:

- **KNN** gave accurate and personalized results based on movie overviews.
- **KMeans** grouped movies effectively by genre, though with some overlap.
- **Random Forest** achieved good predictive performance using multiple features.

Our system demonstrates that combining multiple ML models creates a more accurate and informative movie recommendation and analysis tool.





Thank you!
