



School of Science and Engineering

**CSC538201 Artificial Intelligence for Digital**

**Movie-Roulette**

**Final Report**

Adnane Harachi

**Supervised by:**

Prof. Asmae Mourhir

## **Table of contents:**

|  |    |
|--|----|
| 1. Project Inception .....                                       | 3  |
| 1.1. Framing the Business Idea as an ML Problem .....            | 3  |
| 1.2. Feasibility analysis .....                                  | 5  |
| 2. ML Pipeline Development - From a Monolith to a Pipeline ..... | 9  |
| 2.1. Ensuring ML Pipeline Reproducibility .....                  | 9  |
| 2.2. Pipeline Components .....                                   | 16 |
| 3. Model Deployment .....  | 28 |
| 3.1 ML System Architecture .....                                 | 28 |
| 3.2 Application development .....                                | 29 |
| 3.3 Integration and Deployment .....                             | 32 |
| 3.4. Model Serving and online testing.....                       | 36 |
| 4. Monitoring and Continual Learning.....                        | 37 |
| 4.1. Resource Monitoring.....                                    | 37 |
| 4.2. Data distribution drift monitoring.....                     | 38 |
| 4.3. Continual Learning: CT/CD pipeline .....                    | 39 |
| 4.4. (3 pts) Pipeline orchestration .....                        | 40 |
| 5. Conclusion .....  | 40 |
| References .....   | 42 |

# 1. Project Inception

## 1.1. Framing the Business Idea as an ML Problem

- **Business case description:**

Movie Roulette is a movie recommendation system that uses machine learning algorithms to personalize movie recommendations based on users' tastes, watching history, and comments. The system provides unique suggestions to each user based on user interactions and movie features such as genre, actor, and narrative keywords. Furthermore, Movie Roulette uses sentiment analysis to assess user comments, which improves the accuracy and relevancy of the suggestions. Movie Roulette's user-friendly design and complex features create a dynamic platform for cinephiles to find, explore, and enjoy a curated selection of movies tailored to their individual interests.

- **Business value of using ML:**

**1) Enhanced User Experience:** Movie Roulette uses machine learning to monitor user behavior, preferences, and comments, resulting in a more engaging and personalized experience. This personalized strategy results in enhanced user happiness, retention rates, and client loyalty, which drives revenue growth.

**2) Improved recommendation Accuracy:** Machine learning algorithms evaluate massive quantities of data, such as movie content and user input, to generate extremely accurate and relevant suggestions. This accuracy in suggestions increases the likelihood that users will find material they appreciate, improving their entire experience and motivating them to spend more time on the site.

**3) Increased User Interaction:** Movie Roulette uses machine learning to adapt to user preferences and trends in real time, resulting in a dynamic and engaging browsing experience. This degree of response not only keeps users interested, but also encourages them to share comments and participate

in community conversations, resulting in a livelier and more active user base.

**4) Cost Efficiency:** Implementing machine learning can save money by automating the recommendation process and decreasing the need for manual curation. This efficiency enables the platform to grow more effectively, delivering a bigger audience without incurring additional operational expenditures.

**5) Competitive Advantage:** By utilizing superior machine learning algorithms, Movie Roulette gains a competitive advantage in the industry. The platform's ability to give superior tailored suggestions and a smooth user experience distinguishes it from competitors, drawing new users and gaining market share.

- **Data overview:**

The Movie Roulette system uses data from several sources to provide tailored movie suggestions. The Movie Database is based on multiple csv file that different operations like preprocessing and

validation was performed on.

A screenshot of a Jupyter Notebook interface titled "AJAX-Movie-Recommendation [Administrator]". The left sidebar shows a tree view of files and notebooks, including "clean\_data.py", "data\_testing.py", "Unit\_Integration\_Testing.py", "data.csv", and "evaluation.py". The main area has tabs for "clean\_data.py", "data\_testing.py", "Unit\_Integration\_Testing.py", "data.csv", and "evaluation.py". The "data.csv" tab is active, displaying a list of movie titles and their details. The output pane shows the raw CSV data:

```
1 director_name,actor_1_name,actor_2_name,actor_3_name,genres,movie_title
2 James Cameron,CCH Pounder,Joel David Moore,Wes Studi,Action Adventure Fantasy Sci-Fi,avatar
3 Gore Verbinski,Johnny Depp,Orlando Bloom,Jack Davenport,Action Adventure Fantasy,pirates of the caribbean: at world's end
4 Sam Mendes,Christopher Waltz,Rory Kinnear,Stephanie Sigman,Action Adventure Thriller,spectre
5 Christopher Nolan,Tom Hardy,Christian Bale,Joseph Gordon-Levitt,Action Thriller,the dark knight rises
6 Doug Walker,Doug Walker,Rob Walker,unknown,Documentary,star wars: episode vi - the force awakens
7 Andrew Stanton,Daryl Sabara,Samantha Morton,Polly Walker,Action Adventure Sci-Fi,john carter
8 Sam Raimi,J.K. Simmons,James Franco,Kirsten Dunst,Action Adventure Romance,spider-man 3
9 Nathan Greno,Brad Garrett,Donna Murphy,M.C. Gainey,Adventure Animation Comedy Family Fantasy Musical Romance,tangled
10 Joss Whedon,Chris Hemsworth,Robert Downey Jr.,Scarlett Johansson,Action Adventure Sci-Fi,avengers: age of ultron
11 David Yates,Alan Rickman,Daniel Radcliffe,Rupert Grint,Adventure Family Fantasy Mystery,harry potter and the half-blood prince
12 Zack Snyder,Henry Cavill,Lauren Cohan,Alan D. Purwin,Action Adventure Sci-Fi,batman v superman: dawn of justice
13 Bryan Singer,Kevin Spacey,Marlon Brando,Frank Langella,Action Adventure Sci-Fi,superman returns
14 Marc Forster,Giancarlo Giannini,Mathieu Amalric,Rory Kinnear,Action Adventure,quantum of solace
15 Gore Verbinski,Johnny Depp,Orlando Bloom,Jack Davenport,Action Adventure Fantasy,pirates of the caribbean: dead man's chest
16 Gore Verbinski,Johnny Depp,Ruth Wilson,Tom Wilkinson,Action Adventure Western,the lone ranger
17 Zack Snyder,Henry Cavill,Christopher Meloni,Harry Lennix,Action Adventure Fantasy Sci-Fi,man of steel
18 Andrew Adamson,Peter Dinklage,Pierfrancesco Favino,Damián Alcázar,Action Adventure Family Fantasy,the chronicles of narnia: the lion, the witch and the wardrobe
19 Joss Whedon,Chris Hemsworth,Robert Downey Jr.,Scarlett Johansson,Action Adventure Sci-Fi,the avengers
20 Rob Marshall,Johnny Depp,Sam Clafin,Stephen Graham,Action Adventure Fantasy,pirates of the caribbean: on stranger tides
21 Barry Sonnenfeld,Will Smith,Michael Stuhlbarg,Nicole Kidman,Action Adventure Comedy Family Fantasy Sci-Fi,men in black: international
22 Peter Jackson,Aidan Turner,Adam Brown,James Nesbitt,Adventure Fantasy,the hobbit: the battle of the five armies
23 Marc Webb,Emma Stone,Andrew Garfield,Chris Zylka,Action Adventure Fantasy,the amazing spider-man
24 Ridley Scott,Mark Addy,William Hurt,Scott Grimes,Action Adventure Drama History,robin hood
25 Peter Jackson,Aidan Turner,Adam Brown,James Nesbitt,Adventure Fantasy,the hobbit: the desolation of smaug
26 Chris Weitz,Christopher Lee,Eva Green,Kristin Scott Thomas,Adventure Family Fantasy,the golden compass
27 Peter Jackson,Naomi Watts,Thomas Kretschmann,Evan Parke,Action Adventure Drama Romance,king kong
28 James Cameron,Leonardo DiCaprio,Kate Winslet,Gloria Stuart,Drama Romance,titanic
29 Anthony Russo,Robert Downey Jr.,Scarlett Johansson,Chris Evans,Action Adventure Sci-Fi,captain america: civil war
30 Peter Berg,Liam Neeson,Alexander Skarsgård,Tadanobu Asano,Action Adventure Sci-Fi Thriller,battleship
```

In the snippet above we can clearly see the different csv files used for data.

Link: <https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/datasets>

- **Project archetype:**

Movie Roulette is a Content-Based Movie Recommendation System that uses machine learning algorithms to suggest movies based on users' tastes, watching history, and interactions. It primarily use content-based filtering algorithms, assessing movie features such as genre, actor, and narrative keywords to deliver individualized recommendations. The system's goal is to increase user engagement and pleasure by providing recommendations tailored to individual preferences and interests. While sentiment analysis is an option for future integration, Movie Roulette now relies on content-based filtering to provide suggestions. The algorithm is intended to grow over time, with the option of integrating sentiment analysis to improve suggestions depending on user feedback. Movie Roulette also intends to optimize user interactions and deliver a fluid browsing experience, focusing on usability and accessibility for cinephiles looking for handpicked movie selections.

## 1.2. Feasibility analysis

- **literature review**

- **Introduction:**

The incorporation of machine learning (ML) into recommendation systems has transformed tailored content distribution, greatly improving user experience and engagement across several platforms. This literature review looks at key and current research on collaborative filtering, content-based filtering, hybrid models, and the most recent developments in deep learning for recommendation systems.

- **Collaborative filtering**

Collaborative filtering (CF) is still one of the most extensively used approaches in recommendation systems. It works by analyzing trends in user behavior and predicting future preferences. Sarwar et al. (2001) developed user-based and item-based collaborative filtering algorithms and demonstrated their usefulness in suggesting things by identifying similarities between users or objects. These approaches

employ massive datasets to find connections and recommend things based on the interests of similar users.

Matrix factorization methods, such as Singular Value Decomposition (SVD), have improved CF. Koren, Bell, and Volinsky (2009) demonstrated the effectiveness of SVD in collecting latent characteristics in user-item interactions, resulting in more accurate recommendations. This method breaks down the user-item interaction matrix into lower-dimensional representations, revealing hidden linkages between users and things.

- **Content-Based Filtering**

Content-based filtering (CBF) suggests things similar to those a user has previously loved based on item characteristics. Pazzani and Billsus (2007) presented a thorough discussion of CBF approaches, highlighting the significance of item qualities like genre, director, and keywords in providing suggestions. These strategies create user profiles by studying the characteristics of things that users have rated highly and then utilizing these profiles to recommend comparable items.

Lops, Gemmis, and Semeraro (2011) emphasized the benefits of CBF for adjusting to new objects without needing considerable user engagement. CBF systems are especially successful in areas with extensive information, allowing for the suggestion of goods with precise attributes that match user preferences.

- **Hybrid Models**

Hybrid recommendation systems combine CF and CBF to take use of both methodologies' strengths while mitigating their respective shortcomings. Burke (2002) classified hybrid recommendation approaches, including weighted, switching, and mixed models. These systems improve robustness and accuracy by combining various sources of data.

Bell and Koren (2007) demonstrated the advantages of hybrid models in the Netflix Prize competition, where combining CF and CBF approaches resulted in much higher recommendation accuracy. Hybrid systems may overcome difficulties such as cold-start and data scarcity, resulting in more trustworthy and diversified suggestions.

- **Deep Learning Advances**

Recent advances in deep learning have expanded the capabilities of recommendation systems. He et al. (2017) presented Neural Collaborative Filtering (NCF), a deep learning-based method for modeling complicated user-item interactions using neural networks. NCF beats regular CF because it captures nonlinear linkages and higher-order interactions.

Covington, Adams, and Sargin (2016) examined the use of deep neural networks in YouTube's recommendation engine, emphasizing its scalability and efficacy while processing massive amounts of data. These models use several layers to learn abstract representations, resulting in accurate and tailored suggestions.

- **model choice/ specification of a baseline**

For the baseline model, I used a hybrid recommendation system that I took from huggingface that combines collaborative and content-based filtering. This strategy takes advantage of the qualities of both systems to give more accurate and tailored movie recommendations.

The Collaborative Filtering Component uses user interaction data (ratings, watch history) to forecast user preferences based on comparable users' activities. Matrix factorization techniques such as Singular Value Decomposition (SVD) will be used to identify latent components that capture user and object attributes.

**Content-Based Filtering Component:** This component will assess movie properties (genres, actors, directors, and keywords) and propose films that are similar to those the user has previously appreciated. Feature vectors representing movie material will be created and used to assess similarity.

**Hybrid Integration:** The final suggestion will be a weighted blend of predictions from the collaborative and content-based components. The weights will be established using cross-validation to improve overall suggestion accuracy.

```

1 import streamlit as st
2 import pickle
3 import numpy as np
4 import pandas as pd
5
6 # Load the model from the pkl file
7
8 with open('deep_model_data.pkl', 'rb') as f:
9     df_filtered, smd, id_map, cosine_similarity, user_embeddings, movie_embeddings, user_id_mapping, movie_id_mapping = pickle.load(f)
10
11
12 def predict_similar_movies(userId, title, movie_count):
13     movieId = id_map.loc[title]['movieId'].astype('int')
14     # Map the userId and movieId to the corresponding integer values
15     user_input = np.array([user_id_mapping[userId]])
16     movie_input = np.array([movie_id_mapping[movieId]])
17
18     # Get the learned embeddings for the user and input movie
19     user_emb = user_embeddings[user_input][0]
20     movie_emb = movie_embeddings[movie_input][0]
21
22     # Compute cosine similarity between the input movie and all other movies
23     sim_scores = cosine_similarity([movie_emb], movie_embeddings)[0]
24
25     # Get the indices of the most similar movies
26     movie_indices = sim_scores.argsort()[-(movie_count+1):-1][::-1]
27
28     similar_movies = pd.DataFrame(columns=['title'])
29     for idx in movie_indices:
30         similar_movies = similar_movies.append(df_filtered[df_filtered['movieId'] == id_map.iloc[idx]['movieId']].iloc[0].to_frame().T, ignore_index=True)

```

Link: [https://github.com/adnanehar/Movie-Recommendation-Final/blob/main/deep\\_UI.py](https://github.com/adnanehar/Movie-Recommendation-Final/blob/main/deep_UI.py)

- Metrics for business goal evaluation

To assess the performance of the Movie Roulette system in meeting business objectives, we will use the following metrics:

### Recommendation Accuracy:

**Precision and recall:** Determine the fraction of relevant movies advised among all recommended movies and all relevant movies, respectively.

### User Engagement:

**The click-through rate (CTR)** is the proportion of recommended movies that users click on.

**Conversion Rate:** The proportion of recommended movies that consumers watch.

**User Retention Rate:** The percentage of users who return to the site over a given time period.

### User Satisfaction:

**The Net Promoter Score (NPS)** assesses consumers' likelihood of recommending the platform to

others.

**User comments Analysis:** A sentiment analysis of user evaluations and comments to determine overall satisfaction.

### **Business Impact:**

**income Growth:** Higher subscription rates and ad income due to better user engagement and retention.

**Market Share:** Increase in user base relative to competitors, reflecting the platform's ability to recruit and keep users.

## 2. ML Pipeline Development - From a Monolith to a Pipeline

### 2.1. Ensuring ML Pipeline Reproducibility

Reproducibility is essential in machine learning (ML) projects because it ensures that experiments can be reliably repeated, findings checked, and models kept across time. This section describes the measures used to assure repeatability in Movie Roulette's machine learning workflow.

- **Project structure definition and modularity**

A well-defined project structure encourages organization and modularity, which are critical for scalability and maintenance. The project is organized was organized using cookie cutter as follows:

**Src/:** The source code is structured into modules like data\_preprocessing, model\_training, model\_evaluation, and utils.

**data/:** Raw and processed data are stored here, along with subdirectories for training, validation, and test datasets.

**notebooks/:** Jupyter notebooks are used for exploratory data analysis and experimentation.

**models/:** Saved models and associated artifacts.

Utility scripts for data gathering and preparation. Configuration files for various settings and studies.

**logs:** Files for recording experiment runs and results.

**requirements.txt:** A list of dependencies needed to run the project.

Git is used for version control, which ensures that all code changes are logged, allowing for

collaborative development and quick rollbacks to earlier versions.

|   |              |            |
|---|--------------|------------|
|  .dvc              | final commit | 1 hour ago |
|  .github/workflows | final commit | 1 hour ago |
|  .zen              | final commit | 1 hour ago |
|  datasets        | final commit | 1 hour ago |
|  gx              | final commit | 1 hour ago |
|  model           | final commit | 1 hour ago |
|  pipelines       | final commit | 1 hour ago |
|  static          | final commit | 1 hour ago |
|  steps           | final commit | 1 hour ago |
|  templates       | final commit | 1 hour ago |
|  tests           | final commit | 1 hour ago |
|  .dvcignore      | final commit | 1 hour ago |
|  .gitattributes  | final commit | 1 hour ago |

## Advantages:

Branching and merging allow for feature development and integration in parallel.

History Tracking: Keeps a comprehensive record of modifications, allowing for reversal to earlier states.

Collaboration: Features like pull requests and code reviews enable collaborative development.

- **Code versioning**

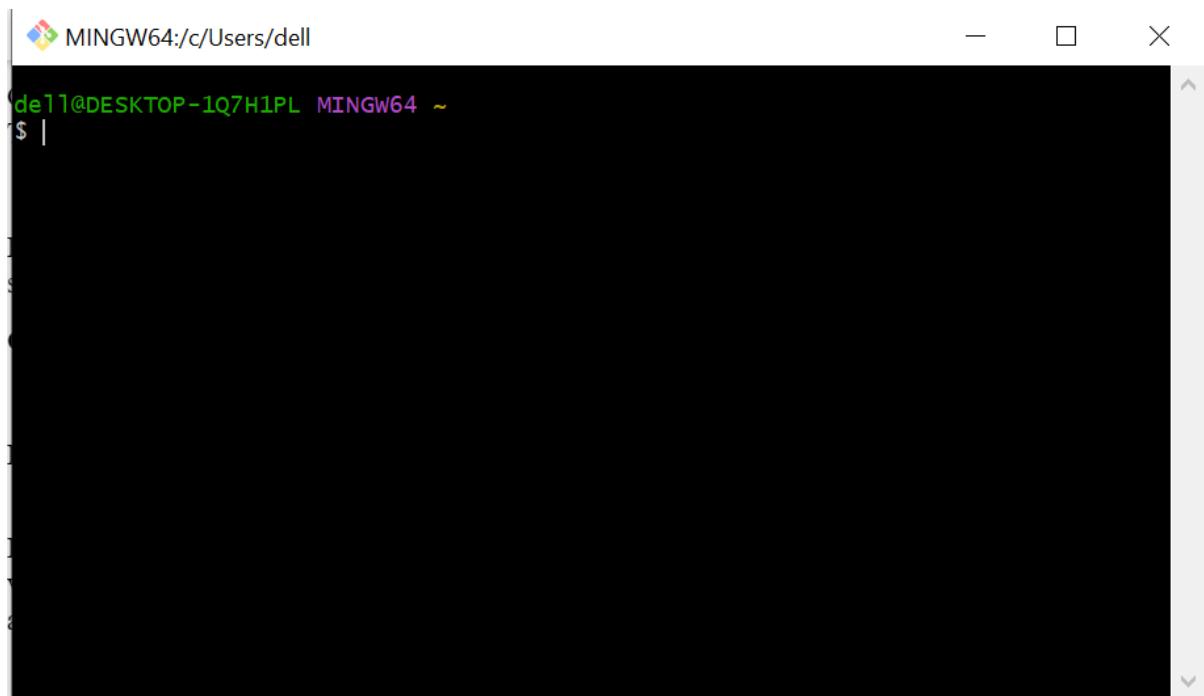
Git is used for extensive code versioning, which facilitates strong collaboration and change tracking. The key practices include:

**Branching Strategy:** Create feature branches for future innovations while keeping the main branch stable.

**Commit Messages:** Use traditional commit message rules to guarantee clarity and traceability.

**Pull Requests:** Creating a review mechanism for code changes to ensure code quality and uniformity.

In addition to Git, Data Version Control (DVC) is used to manage huge data files and dataset versions. DVC records data changes alongside code, allowing for seamless data versioning and assuring that the data used for training and assessment is consistent and repeatable.



The screenshot shows a terminal window titled "MINGW64:/c/Users/dell". The title bar includes standard window controls (minimize, maximize, close). The terminal itself is black with white text. It displays the command prompt "dell1@DESKTOP-1Q7H1PL MINGW64 ~" followed by a single dollar sign (\$) indicating where input can be entered. The rest of the screen is blank, showing only the scroll bars on the right side.

**Advantages:**

**Data management:** Handles huge datasets efficiently while avoiding bloating the Git repository.

Reproducibility ensures that the precise version of an experiment's data may be recovered.

**Integration with Git:** Works seamlessly with Git, providing version control of both code and data.

- **Data versioning**

DVC is essential for versioning datasets, as it ensures that each experiment can be tracked back to the exact data version utilized. This configuration includes:

**Data Storage:** Using DVC to handle data stored in distant storage systems (such as AWS S3 or Google Drive).

**Data Pipelines:** Creating data processing pipelines in DVC helps automate and replicate data preparation procedures.

**Version Tags:** Data versions used in separate studies are tagged for easier retrieval and consistency across pipeline stages.

Advantages:

**Storage Agnostic:** Works with a variety of storage backends (local and cloud).

**Pipeline Management:** Defines data processing pipelines to make data preparation repeatable.

**Data Lineage:** Maintains data lineage, indicating how data was generated and altered.

Link: <https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/.dvc>

- **(5 pts) Experiment tracking and model versioning**

MLflow is widely used for experiment tracking and model versioning, as it provides a centralized platform for managing and tracking experiments. This configuration includes:

**Experiment Tracking:**

MLFlow is a strong tool for experiment tracking and model versioning, allowing for full monitoring and comparison of various models and their performance. By tracking multiple metrics and

parameters for each experiment, MLFlow allows researchers and developers to follow the progress of their models and make better decisions.

In our project, we used MLFlow for experiment tracking and model versioning. MLFlow's metadata store automatically gathers and retains key information for both data and model samples, ensuring that our experiments are reproducible and transparent.

Our major aim with experiment monitoring is to determine which model performs best among the possibilities accessible to us. To do this, we created four critical assessment metrics: accuracy, precision, and recall.

**Logging:** Recording all essential information for each experiment, such as hyperparameters, dataset versions, code versions, and outcomes.

**Comparisons:** Using the MLflow UI, you can compare various runs and analyze performance across studies.

### **Model Versioning:**

**Model Registry:** Use MLflow's model registry to track different model versions.

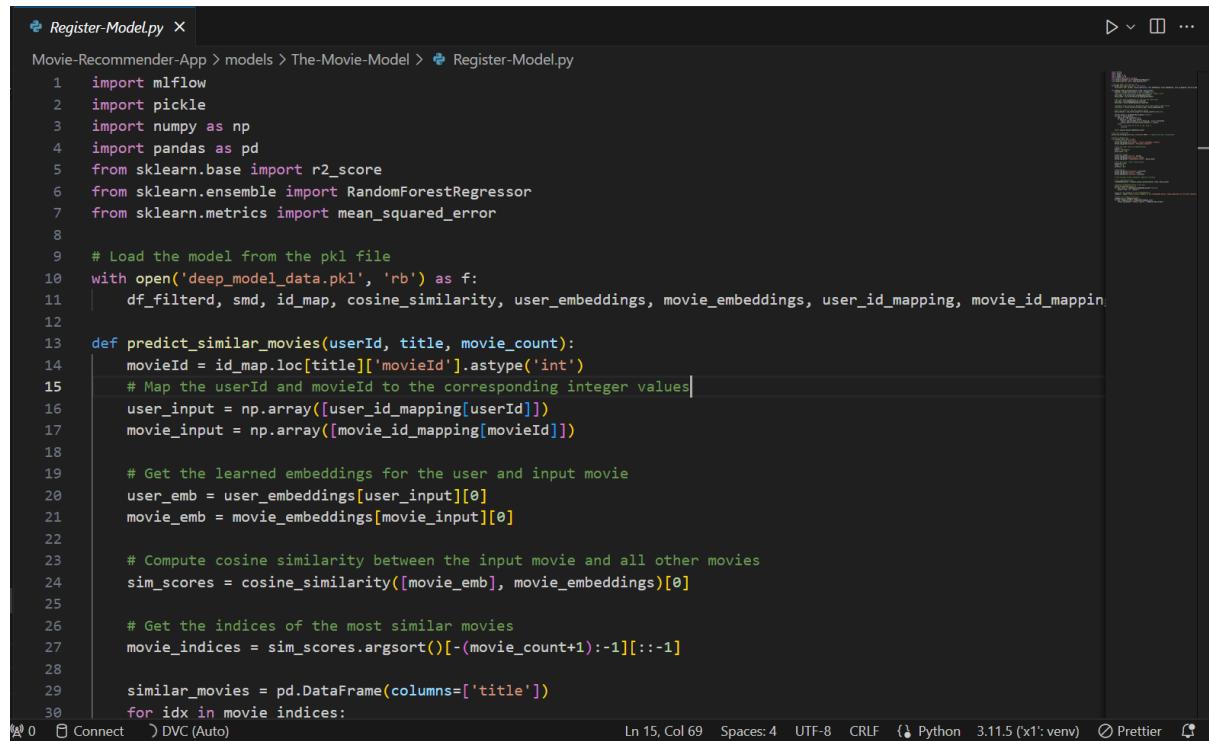
**Model Staging:** Designating models as "staging" or "production" to control lifecycle stages.

### **Advantages:**

**Centralized Metadata administration:** Stores all experiment data in one location for convenient access and administration.

**Search and Query:** Allows users to search and query experiments based on a variety of parameters.

**Visualization:** Provides a UI for displaying experiment runs, comparing metrics, and analyzing Performance.

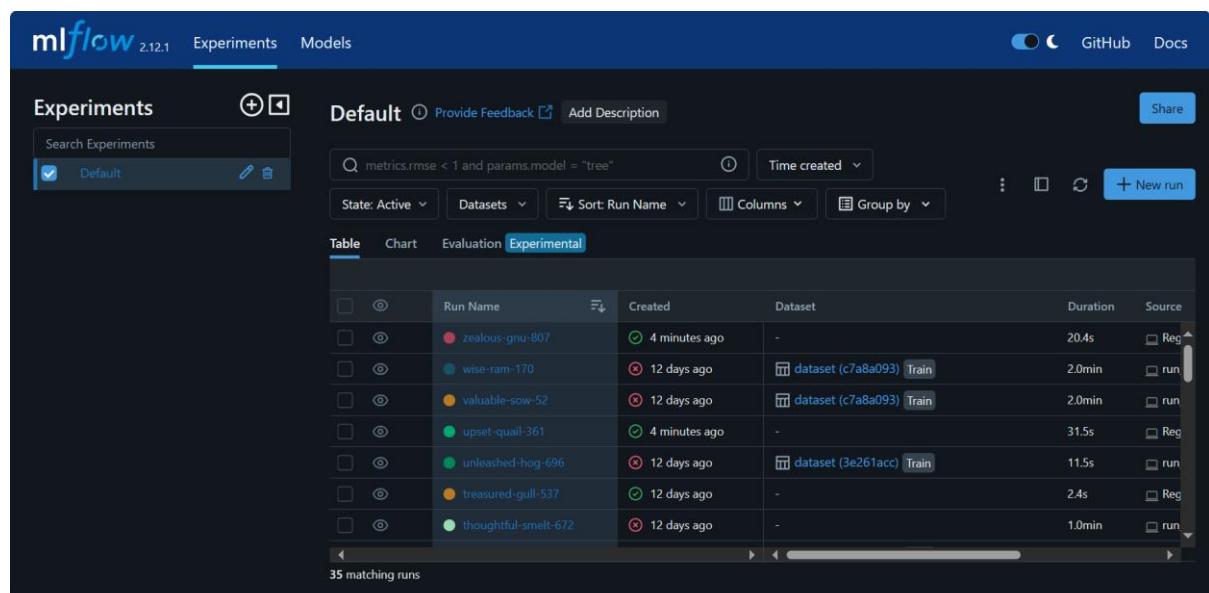


```

Register-Model.py X
Movie-Recommender-App > models > The-Movie-Model > Register-Model.py
1 import mlflow
2 import pickle
3 import numpy as np
4 import pandas as pd
5 from sklearn.base import r2_score
6 from sklearn.ensemble import RandomForestRegressor
7 from sklearn.metrics import mean_squared_error
8
9 # Load the model from the pkl file
10 with open('deep_model_data.pkl', 'rb') as f:
11     df_filtered, smd, id_map, cosine_similarity, user_embeddings, movie_embeddings, user_id_mapping, movie_id_mapping = pickle.load(f)
12
13 def predict_similar_movies(userId, title, movie_count):
14     movieId = id_map.loc[title]['movieId'].astype('int')
15     # Map the userId and movieId to the corresponding integer values
16     user_input = np.array([user_id_mapping[userId]])
17     movie_input = np.array([movie_id_mapping[movieId]])
18
19     # Get the learned embeddings for the user and input movie
20     user_emb = user_embeddings[user_input][0]
21     movie_emb = movie_embeddings[movie_input][0]
22
23     # Compute cosine similarity between the input movie and all other movies
24     sim_scores = cosine_similarity([movie_emb], movie_embeddings)[0]
25
26     # Get the indices of the most similar movies
27     movie_indices = sim_scores.argsort()[-(movie_count+1):-1][::-1]
28
29     similar_movies = pd.DataFrame(columns=['title'])
30     for idx in movie_indices:
31

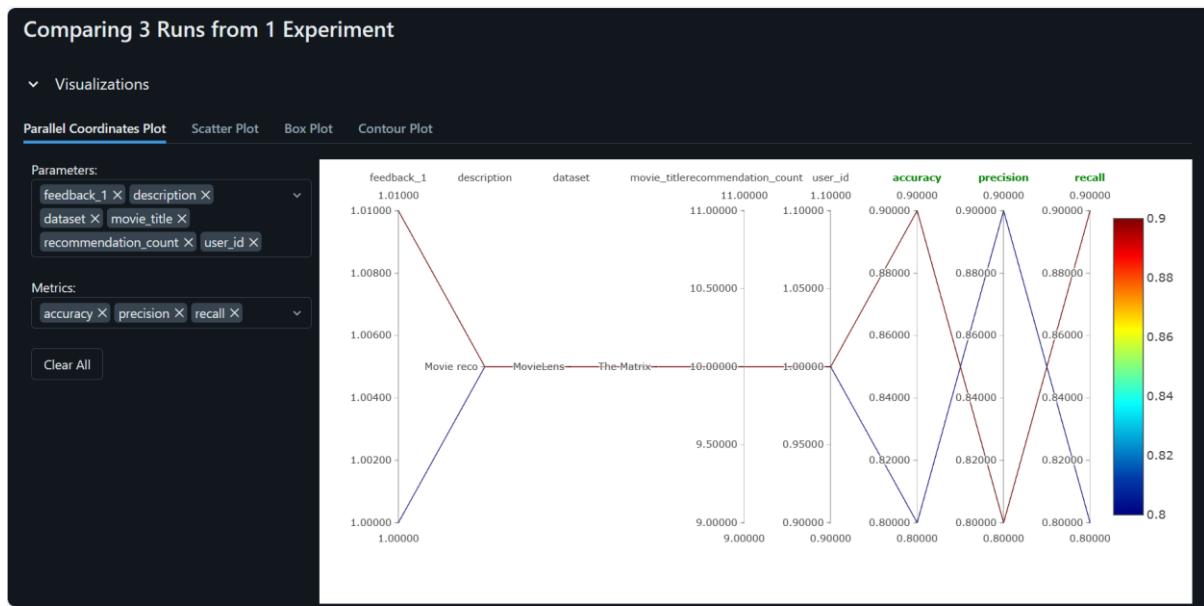
```

Ln 15, Col 69   Spaces: 4   UTF-8   CRLF   { Python 3.11.5 (x1: venv)   { Prettier



|                      | Run Name      | Created                  | Dataset | Duration | Source |
|----------------------|---------------|--------------------------|---------|----------|--------|
| zealous-gnu-807      | 4 minutes ago | -                        |         | 20.4s    | Reg    |
| wise-ram-170         | 12 days ago   | dataset (c7a8a093) Train |         | 2.0min   | run    |
| valuable-sow-52      | 12 days ago   | dataset (c7a8a093) Train |         | 2.0min   | run    |
| upset-quail-361      | 4 minutes ago | -                        |         | 31.5s    | Reg    |
| unleashed-hog-696    | 12 days ago   | dataset (3e261acc) Train |         | 11.5s    | run    |
| treasured-gull-537   | 12 days ago   | -                        |         | 2.4s     | Reg    |
| thoughtful-smelt-672 | 12 days ago   | -                        |         | 1.0min   | run    |

35 matching runs



- **Setting up a meta store for metadata**

The MLflow Tracking Server serves as a single store for all experiment metadata, ensuring that all details are saved and available.

**Advantages:**

**Centralized Metadata administration:** Stores all experiment data in one location for convenient access and administration.

**Search and Query:** Allows users to search and query experiments based on a variety of parameters.  
**Visualization:** Provides a UI for displaying experiment runs, comparing metrics, and analyzing performance.

- **Setting up the machine learning pipeline under an MLOps platform**

ZenML provides a powerful framework for coordinating the full machine learning workflow. ZenML simplifies each step, from data intake and preprocessing to model training, assessment, and deployment, resulting in a seamless integration and monitoring process. ZenML makes machine learning projects more productive by offering a standardized way to building and maintaining processes.

## 2.2. Pipeline Components

### 2.2.1 Setup of data pipeline within the larger ML pipeline/ MLOps Platform

The ML pipeline for Movie Roulette is built with a variety of components to provide quick data processing, model training, and assessment. This section describes how to build up the data pipeline inside the wider ML pipeline running on the MLOps platform, including data validation, preprocessing, feature engineering, model training, and offline assessment. Additionally, model behavioral tests are created to validate the models' dependability and usefulness.

#### **Data Pipeline Setup inside the Larger ML Pipeline/MLOps Platform.**

#### **Using zenml:**

##### **Advantages:**

Modularity: ZenML enables the design of modular and reusable pipeline components, which improves maintainability and scaling.

Versioning: Manages pipeline versions to provide repeatability and traceability.

Workflow Orchestration: Enables the orchestration of complicated data workflows, allowing for automation and monitoring.

Disadvantage: not enough documentation.

So here I created a pipeline where I have many steps, I ingested the data, clean it, evaluate it then make a pipeline called training pipeline to take all these steps and build a Zenml pipeline.

```
steps > clean_data.py > ...
14
15     @step
16     def clean_data(
17         data: pd.DataFrame,
18     ) -> Tuple[
19         Annotated[pd.DataFrame, "x_train"],
20         Annotated[pd.DataFrame, "x_test"],
21         Annotated[pd.Series, "y_train"],
22         Annotated[pd.Series, "y_test"],
23     ]:
24         """Data cleaning class which preprocesses the data and divides it into train and test data.
25
26         Args:
27             data: pd.DataFrame
28         """
29
30         try:
31             preprocess_strategy = DataPreprocessStrategy()
32             data_cleaning = DataCleaning(data, preprocess_strategy)
33             preprocessed_data = data_cleaning.handle_data()
34
35             divide_strategy = DataDivideStrategy()
36             x_train, x_test, y_train, y_test = divide_strategy.handle_data(preprocessed_data)
37
38             return x_train, x_test, y_train, y_test
39         except Exception as e:
40             logging.error(e)
41             raise e
```

Links: <https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/steps>

Screenshot of a code editor showing two files: `evaluation.py` and `ingest_data.py`.

`evaluation.py` (Top Editor)

```

steps > ingest_data.py > ingest_df
1 import logging
2 import pandas as pd
3 from zenml import step
4
5 class IngestData:
6
7     def __init__(self, data_path: str):
8         self.data_path = data_path
9
10    def get_data(self):
11        logging.info(f"Ingesting data from {self.data_path}")
12        # Read the CSV file
13        df = pd.read_csv(self.data_path)
14
15        # Convert columns to appropriate data types
16
17        df['title_year'] = df['title_year'].astype('string') # Assuming 'director_name' is a string
18        df['director_facebook_likes'] = df['director_facebook_likes'].astype('string') # Assuming 'actor_'
19
20        return df
21
22 @step
23 def ingest_df(data_path: str) -> pd.DataFrame:
24     """
25     Args:
26         data_path: str: Path to the data file.
27     Returns:
28         df: pd.DataFrame: DataFrame containing the ingested data.
29     """
30
31
32
33
34
35
36
37
38

```

`ingest_data.py` (Bottom Editor)

```

steps > evaluation.py > IngestData > get_data
5 class IngestData:
6
7     def get_data(self):
8
9         return df
10
11 @step
12 def ingest_df(data_path: str) -> pd.DataFrame:
13     """
14     Args:
15         data_path: str: Path to the data file.
16     Returns:
17         df: pd.DataFrame: DataFrame containing the ingested data.
18     """
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

```

The code in both files is identical, showing the definition of the `IngestData` class and the `ingest_df` step function.

```

steps > model_train.py > train_model
20
21 @step
22 def train_model(
23     x_train: Union[pd.DataFrame, None],
24     x_test: Union[pd.DataFrame, None],
25     y_train: Union[pd.Series, None],
26     y_test: Union[pd.Series, None],
27     config: ModelNameConfig,
28 ) -> RegressorMixin:
29     """
30     Args:
31         x_train: pd.DataFrame or None
32         x_test: pd.DataFrame or None
33         y_train: pd.Series or None
34         y_test: pd.Series or None
35     Returns:
36         model: RegressorMixin
37     """
38     try:
39         model = None
40         tuner = None
41
42         if config.model_name == "lightgbm":
43             mlflow.lightgbm.autolog()
44             model = LightGBMModel()
45         elif config.model_name == "randomforest":
46             mlflow.sklearn.autolog()
47             model = RandomForestModel()
48         elif config.model_name == "xgboost":
49             mlflow.xgboost.autolog()

```

```

pipelines > training_pipeline.py > train_pipeline
1  from numpy import str_
2  from zenml.config import DockerSettings
3  from zenml.integrations.constants import MLFLOW
4  from zenml import pipeline
5  from zenml.steps.base_step import BaseStep
6  from steps.clean_data import clean_data
7  from steps.evaluation import evaluation
8  from steps.ingest_data import ingest_df
9  from steps.model_train import train_model
10
11 docker_settings = DockerSettings(required_integrations=[MLFLOW])
12
13
14 @pipeline(enable_cache=True, settings={"docker": docker_settings})
15 def train_pipeline(data_path: str):
16     """
17     Args:
18         ingest_data: DataClass
19         clean_data: DataClass
20         model_train: DataClass
21         evaluation: DataClass
22     Returns:
23         mse: float
24         rmse: float
25     """
26     df = ingest_df(data_path)
27     x_train, x_test, y_train, y_test = clean_data(df)
28     model = train_model(x_train, x_test, y_train, y_test)
29     rmse, mse = evaluation(model, x_test, y_test)
30

```

Link: <https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/pipelines>

The only problem that this has is that it takes a long time to do the training which is not very optimal but I do believe that this problem can be solved by switching from a cpu to a gpu

- **Data Validation and Verification**

Great Expectation is an effective tool for evaluating and testing data in the machine learning pipeline. It allows users to establish data expectations, such as schema, distribution, and integrity, which are then automatically tested against incoming data. By introducing Great Expectations into the workflow, data scientists may improve data quality and consistency, resulting in more trustworthy and resilient machine learning models.

- **Preprocessing and Feature Engineering**

**Data Validation & Verification Tool:** Great Expectations, pandas, and numpy

**Why are there such high expectations?**

Great Expectations is an open-source package that includes comprehensive capabilities for data validation and profiling. It allows for the formulation, maintenance, and assessment of data expectations (assertions), assuring high-quality and consistent datasets for training and evaluation.

**Advantages:**

**Comprehensive Validation:** With Great Expectations, you can design thorough expectation suites that test data against a variety of criteria, including data types, ranges, missing values, and unique restrictions.

**Automated Documentation:** Creates concise, human-readable documentation based on expectations that can be shared with stakeholders to preserve data quality transparency.

**Integration Capabilities:** It seamlessly connects with a variety of data sources and pipeline architectures, making it adaptable to varied situations.

**Historical Traceability:** Keeps a record of data validation outcomes, which is essential for auditing and compliance.

Pandas and NumPy are two important libraries in the Python environment for data processing and numerical computing, respectively. They are critical in preprocessing and feature engineering, offering reliable and fast methods for dealing with big datasets and executing complicated transformations. Pandas provides strong data structures like as DataFrames, which facilitate data manipulation, cleaning, and transformation. It also has a variety of features for merging, joining, reshaping, and aggregating data. Pandas' straightforward syntax and comprehensive capability make data preparation chores simple, even for individuals with little programming knowledge.

NumPy, on the other hand, supports massive, multidimensional arrays and matrices, as well as a set of mathematical algorithms for rapidly manipulating these arrays. NumPy operations are highly optimized, and they can do large-scale numerical computations significantly quicker than standard Python. It also integrates smoothly with other scientific libraries in Python, making it an indispensable component of the data science toolset.

Pandas is used to manage missing values, duplicate records, and data discrepancies. It enables data type conversions and normalization. Pandas and NumPy are used to extract new features from existing data, such as categorical variable encoding, numerical feature scaling, and interaction term generation.

The screenshot shows a Jupyter Notebook interface with the following details:

- EXPLORER View:** Shows a tree structure of files and folders. The root folder is "AJAX-MOVIE-RECOMMENDATION". Inside it are ".vscode", ".zen", "datasets", "gx", "mlruns", "model", and "Movie-Recommender-App". "Movie-Recommender-App" contains "models", "The-Movie-Model", ".zen", "gx", "checkpoints", and "expectations". "expectations" contains "ge\_store\_backend\_id" and "movies metadata expectations".
- Code Editor:** Displays the code for "preprocessing 1.ipynb".

```
import numpy as np

data = pd.read_csv('movie_metadata.csv')
```
- Output:** Shows the result of the `data.head(10)` command, displaying a DataFrame with 8 rows and 9 columns. The columns are: color, director\_name, num\_critic\_for\_reviews, duration, director\_facebook\_likes, actor\_3\_facebook\_likes, actor\_2\_name, actor\_1\_name, and actor\_2\_facebook\_likes.

|   | color | director_name     | num_critic_for_reviews | duration | director_facebook_likes | actor_3_facebook_likes | actor_2_name     | actor_1_name |
|---|-------|-------------------|------------------------|----------|-------------------------|------------------------|------------------|--------------|
| 0 | Color | James Cameron     | 723.0                  | 178.0    | 0.0                     | 855.0                  | Joel David Moore |              |
| 1 | Color | Gore Verbinski    | 302.0                  | 169.0    | 563.0                   | 1000.0                 | Orlando Bloom    |              |
| 2 | Color | Sam Mendes        | 602.0                  | 148.0    | 0.0                     | 161.0                  | Rory Kinnear     |              |
| 3 | Color | Christopher Nolan | 813.0                  | 164.0    | 22000.0                 | 23000.0                | Christian Bale   |              |
| 4 | NaN   | Doug Walker       | Nan                    | Nan      | 131.0                   | Nan                    | Rob Walker       |              |
| 5 | Color | Andrew Stanton    | 462.0                  | 132.0    | 475.0                   | 530.0                  | Samantha Morton  |              |
| 6 | Color | Sam Raimi         | 392.0                  | 156.0    | 0.0                     | 4000.0                 | James Franco     |              |
| 7 | Color | Nathan Gruson     | 224.0                  | 100.0    | 15.0                    | 284.0                  | Donna            |              |

preprocessing 1.ipynb M mrs.ipynb M

Movie-Recommender-App > models > The-Movie-Model > mrs.ipynb > Adnane Harachi > 80561 > importing libraries

+ Code + Markdown | Run All Restart Clear All Outputs Variables Outline ... Python 3.11.5

## Convert pandas dataframe to Great Expectation

```
[ ] my_df = ge.from_pandas(df) Python
```

```
[11] type(my_df) Python
```

```
... great_expectations.dataset.pandas_dataset.PandasDataset
```

## GE Data Quality Tests

```
[38] # check number of rows in the dataset
my_df.expect_table_row_count_to_equal(1000) Python
```

```
... {
  "success": false,
  "result": {
    "observed_value": 5043
  },
  "meta": {}
} Cell 3 of 147
```

DVC (Global)

preprocessing 1.ipynb M mrs.ipynb M

Movie-Recommender-App > models > The-Movie-Model > mrs.ipynb > Adnane Harachi > 80561 > importing libraries

+ Code + Markdown | Run All Restart Clear All Outputs Variables Outline ... Python 3.11.5

## Test on columns

```
[13] my_df.expect_column_to_exist('movie_title') Python
```

```
... {
  "success": true,
  "result": {},
  "meta": {},
  "exception_info": {
    "raised_exception": false,
    "exception_traceback": null,
    "exception_message": null
  }
}
```

```
[14] my_df.expect_column_values_to_be_unique('color') Python
```

```
... {
  "success": false,
  "result": {
    "element_count": 5043,
    "missing_count": 19,
    "bad_records": 0
  }
} Cell 3 of 147
```

DVC (Global)

```
my_df.expect_column_values_to_not_be_null('movie_title')
[15] Python
```

```
... {
    "success": true,
    "result": {
        "element_count": 5043,
        "unexpected_count": 0,
        "unexpected_percent": 0.0,
        "unexpected_percent_total": 0.0,
        "partial_unexpected_list": []
    },
    "meta": {},
    "exception_info": {
        "raised_exception": false,
        "exception_traceback": null,
        "exception_message": null
    }
}
```

## Test values in a set (list)

```
df.genres.unique()
[16] DVC (Global) Spaces: 4 Cell 16 of 147 ✘ Prettier ⌂
```

```
Movie-Recommender-App > models > The-Movie-Model > mrs.ipynb > Adnane Harachi > 80561 > importing libraries
+ Code + Markdown | ▶ Run All ⚡ Restart ⏷ Clear All Outputs | 📁 Variables ⏷ Outline ...
```

```
Python 3.11.5
```

```
Save my test cases and re-use
```

```
my_df.save_expectation_suite('movies.expectations.json')
[32] Python
```

```
df2 = ge.read_csv('C:\Windows\System32\Movie-Roulette\movie_metadata.csv')
[33] Python
```

## Test with Config file

```
test_results = df2.validate(expectation_suite="movies.expectations.json")
[36] Python
```

## database preview

```
df.columns
[869] Python
```

```
... Index(['adult', 'belongs_to_collection', 'budget', 'genres', 'homepage', 'id',
       'original_language', 'original_title', 'overview', 'popularity',
       'production_companies', 'production_countries', 'release_date',
       'title', 'video', 'vote_average'],
      dtype='object')
```

```
DVC (Global) Spaces: 4 Cell 16 of 147 ⌂
```

Link: [https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/.ipynb\\_checkpoints](https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/.ipynb_checkpoints)

```

1  {
2    "data_asset_type": "Dataset",
3    "expectation_suite_name": "default",
4    "expectations": [
5      {
6        "expectation_type": "expect_column_to_exist",
7        "kwargs": {
8          "column": "movie_title"
9        },
10       "meta": {}
11     },
12     {
13       "expectation_type": "expect_column_values_to_not_be_null",
14       "kwargs": {
15         "column": "movie_title"
16       },
17       "meta": {}
18     },
19     {
20       "expectation_type": "expect_column_max_to_be_between",
21       "kwargs": {
22         "column": "imdb_score",
23         "max_value": 9.5,
24         "min_value": 1.5
25       },
26       "meta": {}
27     },
28     {
29       "expectation_type": "expect_column_mean_to_be_between",
30       "kwargs": {}
31     }
32   ]
33 }

```

Ln 1, Col 1 Spaces: 2 UTF-8 CRLF {} JSON ✓ Prettier

```

df.isnull().sum()

```

|     | color | director_name | num_critic_for_reviews | duration | director_facebook_likes | actor_3_facebook_likes | actor_2_name | actor_1_facebook_likes | gross | genres | actor_1_name | movie_title | num_voted_users | cast_total_facebook_likes | actor_3_name | facenumber_in_poster | plot_keywords | movie_imdb_link | num_user_for_reviews | language | Count |
|-----|-------|---------------|------------------------|----------|-------------------------|------------------------|--------------|------------------------|-------|--------|--------------|-------------|-----------------|---------------------------|--------------|----------------------|---------------|-----------------|----------------------|----------|-------|
| ... | 19    | 104           | 50                     | 15       | 104                     | 23                     | 13           | 7                      | 884   | 0      | 7            | 0           | 0               | 0                         | 23           | 13                   | 153           | 0               | 21                   | 14       | 5     |

Cell 32 of 147

```

df['genres'] = df['genres'].fillna('[]').apply(literal_eval).apply(lambda x: [i['name'] for i in x] if isinstance(x, list) else [])
df['production_companies'] = df['production_companies'].fillna('[]').apply(literal_eval).apply(lambda x: [i['name'] for i in x] if isinstance(x, list) else [])
df['production_countries'] = df['production_countries'].fillna('[]').apply(literal_eval).apply(lambda x: [i['name'] for i in x] if isinstance(x, list) else [])
df['spoken_languages'] = df['spoken_languages'].fillna('[]').apply(literal_eval).apply(lambda x: [i['name'] for i in x] if isinstance(x, list) else [])
df['year'] = pd.to_datetime(df['release_date'], errors='coerce').apply(lambda x: str(x).split('-')[0] if x != np.nan else np.nan)
df.head().transpose()

```

The Logic Behind  
IMDB's weighted rating formula:

$$WR = \frac{(v \div (v+m)) \times R + (m \div (v+m)) \times C}{C}$$

R = average for the movie (mean)  
v = number of votes for the movie  
m = minimum votes required to be listed in the Top 250 MOVIES  
C = the mean vote across the whole report

```

vote_counts = df[df['vote_count'].notnull()]['vote_count'].astype('int')
vote_averages = df[df['vote_average'].notnull()]['vote_average'].astype('int')
C = vote_averages.mean()
print('The Mean Value of the voting averages= ', C)
m = vote_counts.quantile(0.96)
print('The Minimum Vote count for a movie to consider= ', m)

```

[877] ... The Mean Value of the voting averages= 5.244896612406511  
The Minimum Vote count for a movie to consider= 576.6399999999994

```

qualified = df[(df['vote_count'] >= m) & (df['vote_count'].notnull()) & (df['vote_average'].notnull())][['title', 'vote_count', 'vote_average']]
qualified['vote_count'] = qualified['vote_count'].astype('int')
qualified['vote_average'] = qualified['vote_average'].astype('int')
print('The structure of the qualified database is= ', qualified.shape)

```

[878] ... The structure of the qualified database is= (1819, 6)

| Therefore, a movie has to have at least 576.63 votes on TMDB. We also see that the average rating for a movie on TMDB is 5.244 on a scale of 10. Only 1899 Movies qualify to be on our chart.

fe

The rest is in the notebook provided with the code to not make the report very long.

## 2.2.2 (5 pts) Integration of model training and offline evaluation into the ML pipeline / MLOps Platform

Model training and offline assessment are successfully integrated into the ML pipeline utilizing ZenML's model training and evaluation components. ZenML manages the whole training process, ensuring that models are trained effectively and that performance is measured using complete offline metrics. One of the primary benefits of utilizing ZenML for this purpose is its strong experiment tracking capabilities. ZenML painstakingly monitors model training runs and evaluation results, providing precise information about the performance of various model setups. This tracking is extremely useful for understanding how various model parameters and configurations affect overall performance, allowing data scientists to make educated judgments regarding model tweaks and enhancements. Furthermore, ZenML's scalability is a valuable advantage, allowing the model training and assessment procedures to be scaled up as necessary. This scalability enables large-scale experimentation, which is required for the development and refinement of high-quality machine learning models. By allowing for prolonged experimentation, ZenML assures that the ML pipeline can manage larger volumes of data and more complicated model architectures, resulting in more accurate and efficient models. Overall, using ZenML to integrate model training and offline assessment into the ML pipeline improves the machine learning workflow's resilience and efficacy.

### **2.2.3 Development of model behavioral tests**

Pytest is used for behavioral testing, which ensures the data pipeline's reliability and correctness.

**Data Tests:** These check the ability to read labels and tweets from files, assuring smooth functioning.

**ML Infrastructure Testing:** External dependencies, such as Cassandra database connections, are simulated to ensure that the data pipeline integrates seamlessly.

**Unit Tests:** Critical components such as preprocessing, assessment, and model inference are thoroughly tested to ensure correctness.

**Monitoring Tests:** Continuous evaluations are performed to detect data drifts in production, preventing possible model deterioration."

```

1 import pandas as pd
2 import pytest
3
4 def test_data(csv_file):
5     # Load the CSV file into a DataFrame
6     df = pd.read_csv(csv_file)
7
8     # Check for missing values
9     missing_values = df.isnull().sum()
10    if missing_values.any():
11        print("Missing values found:")
12        print(missing_values)
13    else:
14        print("No missing values found.")
15
16    # Check for duplicate rows
17    duplicate_rows = df.duplicated().sum()
18    if duplicate_rows:
19        print("Duplicate rows found:", duplicate_rows)
20    else:
21        print("No duplicate rows found.")
22
23    # Check data types
24    print("Data types:")
25    print(df.dtypes)
26
27    # Check for unique values in categorical columns
28    categorical_columns = df.select_dtypes(include=['object']).columns
29    for column in categorical_columns:
30        unique_values = df[column].unique()

```

```

tests > Unit_Integration_Testing.py
1 from matplotlib import _preprocess_data
2 import pandas as pd
3 import pytest
4 import logging
5 from unittest.mock import MagicMock, patch
6
7 from steps.clean_data import clean_data
8
9
10 sample_csv_path = r'C:\Users\dell\cleaned_data.csv'
11
12 @pytest.fixture
13 def test_successful_data_load():
14     # Arrange
15     test_csv_path = r'C:\Users\dell\Documents\cleaned_data.csv'
16
17     # Act
18     result = clean_data(test_csv_path)
19
20     # Assert
21     assert isinstance(result, pd.DataFrame)
22
23 def test_file_not_found():
24     # Arrange
25     non_existent_file_path = r'C:\Users\dell\Documents\cleaned_data.csv'
26
27     # Act and Assert
28     with pytest.raises(FileNotFoundError):
29         clean_data(non_existent_file_path)

```

Ln 26, Col 56 Spaces: 4 UTF-8 CRLF ⚡ Python 3.11.5 ('x1: venv') ⚡ Prettier ⚡

Link: <https://github.com/adnanehar/Movie-Recommendation-Final/tree/main/tests>

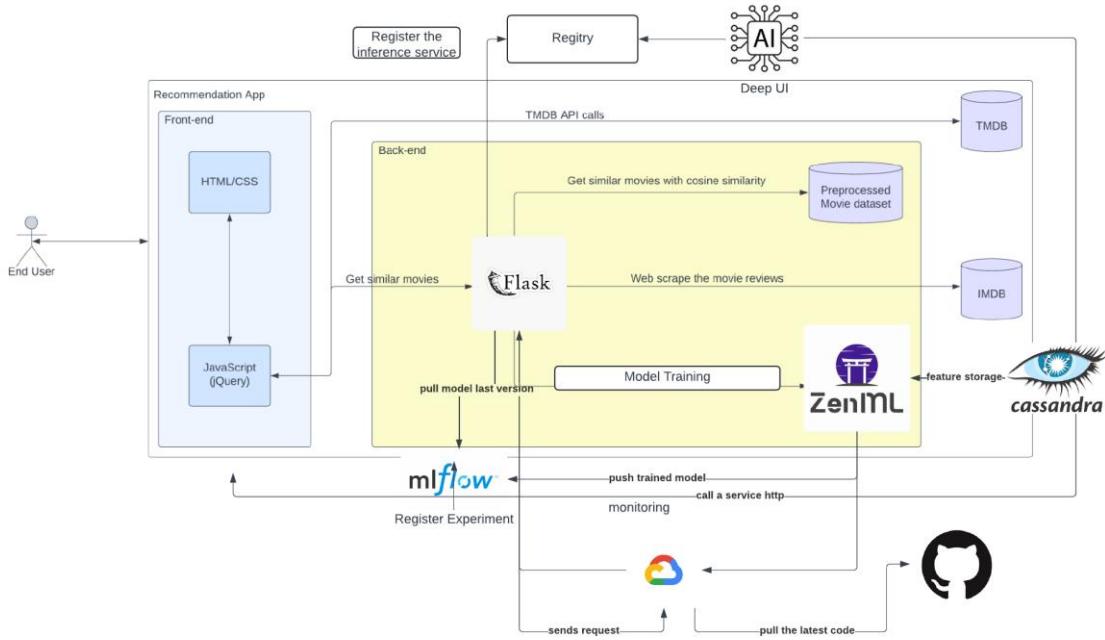
```
tests > 🐍 Infrastructure_test.py
 1 import pytest
 2 import os
 3 import shutil
 4 import tempfile
 5 import subprocess
 6 import sys
 7
 8 def test_python_version():
 9     assert sys.version_info >= (3, 6), "Python version must be 3.6 or higher"
10
11 def test_required_packages():
12     required_packages = ['numpy', 'pandas', 'pytest', 'scikit-learn', 'tensorflow']
13     missing_packages = [pkg for pkg in required_packages if pkg not in sys.modules]
14     assert not missing_packages, f"Required packages not installed: {missing_packages}"
15
16 def test_environment_variables():
17     required_env_vars = ['API_KEY', 'DATABASE_URL']
18     missing_env_vars = [var for var in required_env_vars if var not in os.environ]
19     assert not missing_env_vars, f"Required environment variables missing: {missing_env_vars}"
20
21 def test_directory_structure():
22     assert os.path.exists('data'), "Data directory missing"
23     assert os.path.exists('models'), "Models directory missing"
24     assert os.path.exists('logs'), "Logs directory missing"
25
26 def test_temp_directory_permission():
27     temp_dir = tempfile.mkdtemp()
28     assert os.access(temp_dir, os.W_OK), "No write permission in temp directory"
29     shutil.rmtree(temp_dir)
30
```

Connect ⌘ 1 file to analyze ⌂ DVC (Auto) Ln 1, Col 1 (13 selected) Spaces: 4 UTF-8 CRLF ⚙ Python 3.11.5 ('x1': venv) ⚙ Prettier

## 3. Model Deployment

### 3.1 ML System Architecture

- (5 pts) Drawing with architecture highlights

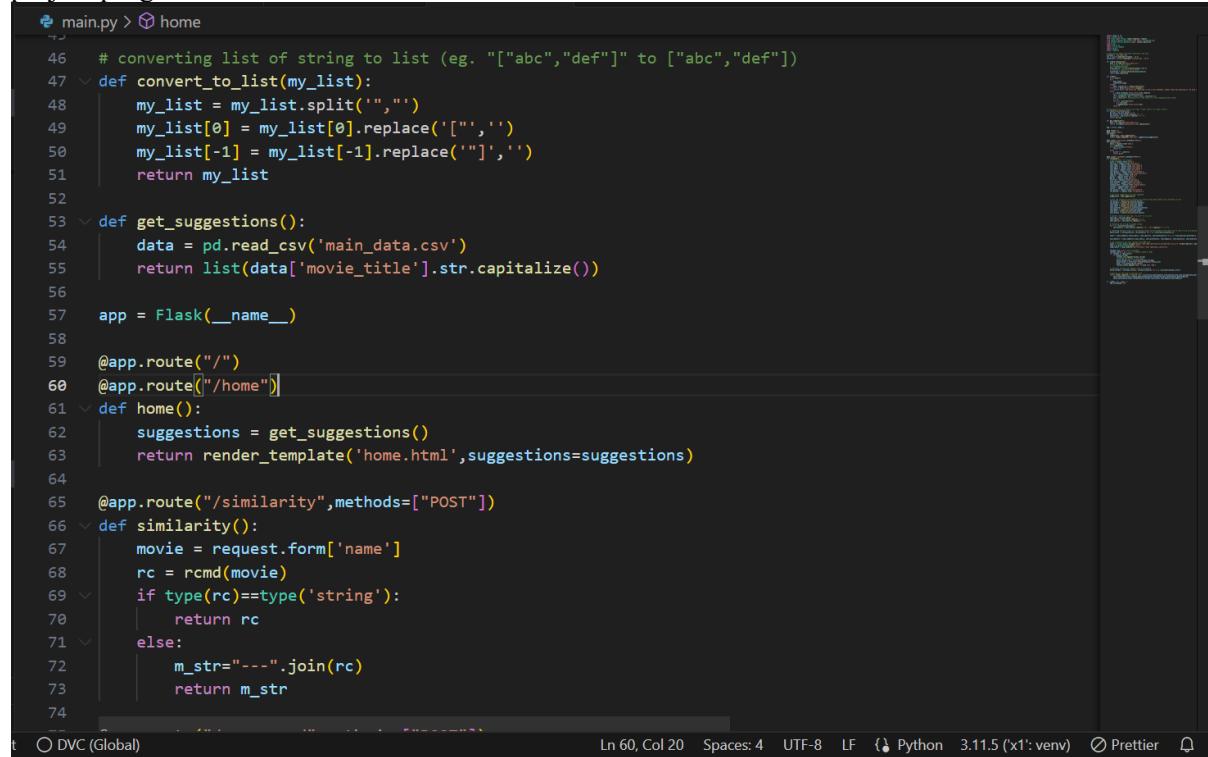


## 3.2 Application development

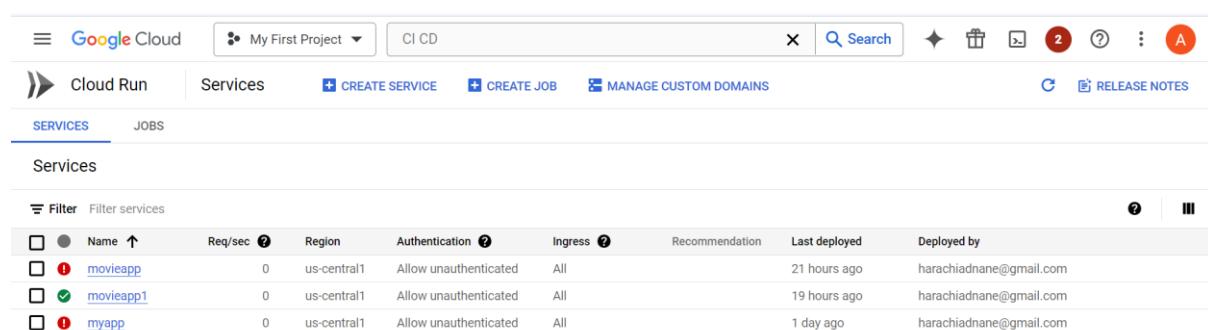
- (5 pts) Model service development

Flask is the core framework for building the model service, offering a fast and lightweight method for delivering machine learning models. Its simplicity and adaptability make it an ideal candidate for creating the backend infrastructure that handles model inference requests. Using Flask's lightweight nature, developers may quickly create microservices suited to certain model functionality. In addition, I used Docker containers and Docker images to package the program. I then deployed it using Google Cloud Platform. Flask serves as the foundation for building the model service, offering a quick and resource-efficient way to deploy machine learning models. Its inherent simplicity and versatility make it an ideal candidate for developing the backend infrastructure that handles model inference requests. With Flask's lightweight architecture, developers can quickly construct microservices suited to certain model functionality. This simplified method not only speeds up development, but also improves maintainability and extensibility as the system grows."

project progresses.



```
main.py > home
46     # converting list of string to list (eg. ["abc","def"] to ["abc","def"])
47     def convert_to_list(my_list):
48         my_list = my_list.split(",") 
49         my_list[0] = my_list[0].replace('[','')
50         my_list[-1] = my_list[-1].replace(']', '')
51         return my_list
52
53     def getSuggestions():
54         data = pd.read_csv('main_data.csv')
55         return list(data['movie_title'].str.capitalize())
56
57     app = Flask(__name__)
58
59     @app.route("/")
60     @app.route("/home")
61     def home():
62         suggestions = getSuggestions()
63         return render_template('home.html', suggestions=suggestions)
64
65     @app.route("/similarity", methods=["POST"])
66     def similarity():
67         movie = request.form['name']
68         rc = rcmd(movie)
69         if type(rc)==type('string'):
70             return rc
71         else:
72             m_str="---".join(rc)
73             return m_str
74
```



DVC (Global) Ln 60, Col 20 Spaces: 4 UTF-8 LF Python 3.11.5 ('x1': venv) Prettier

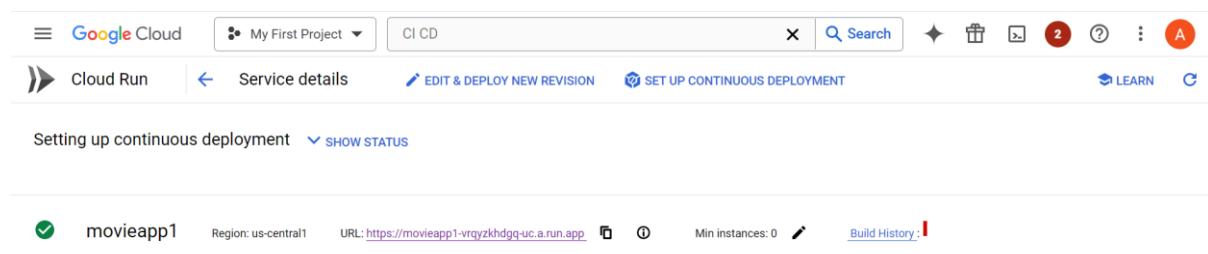
Google Cloud My First Project CI CD Search

Cloud Run Services + CREATE SERVICE + CREATE JOB MANAGE CUSTOM DOMAINS RELEASE NOTES

SERVICES JOBS

Services

| Name      | Req/sec | Region      | Authentication        | Ingress | Recommendation | Last deployed | Deployed by             |
|-----------|---------|-------------|-----------------------|---------|----------------|---------------|-------------------------|
| movieapp  | 0       | us-central1 | Allow unauthenticated | All     |                | 21 hours ago  | harachiadnane@gmail.com |
| movieapp1 | 0       | us-central1 | Allow unauthenticated | All     |                | 19 hours ago  | harachiadnane@gmail.com |
| myapp     | 0       | us-central1 | Allow unauthenticated | All     |                | 1 day ago     | harachiadnane@gmail.com |



Google Cloud My First Project CI CD

Cloud Run Service details EDIT & DEPLOY NEW REVISION SET UP CONTINUOUS DEPLOYMENT LEARN

Setting up continuous deployment SHOW STATUS

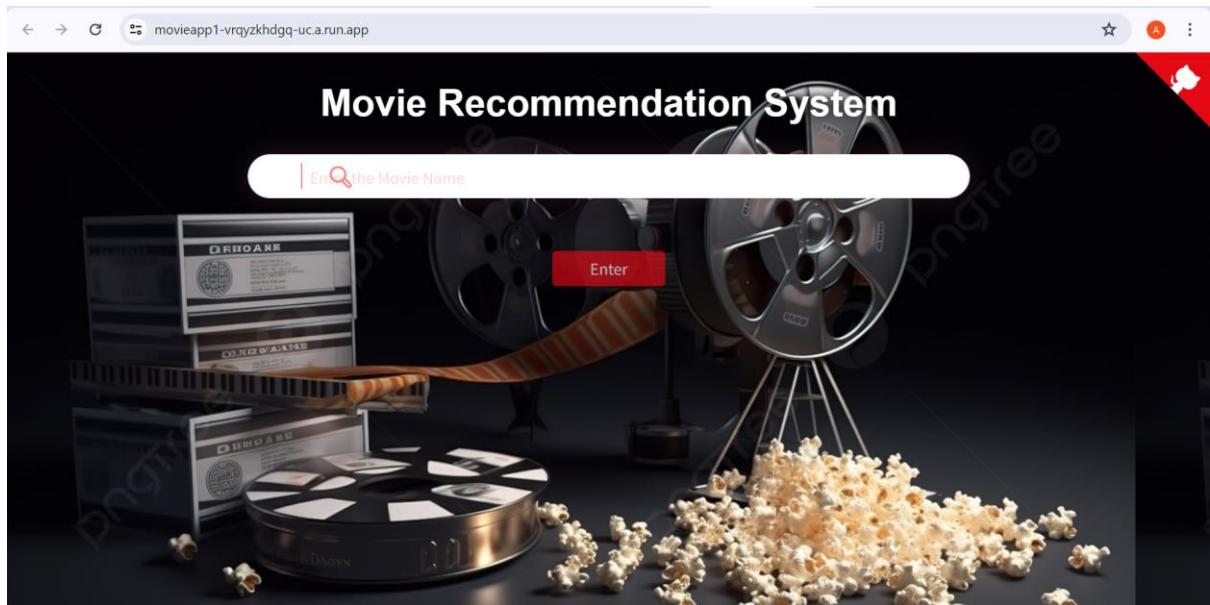
movieapp1 Region: us-central1 URL: https://movieapp1-vrqyzkhdgq-uc.a.run.app Min instances: 0 Build History

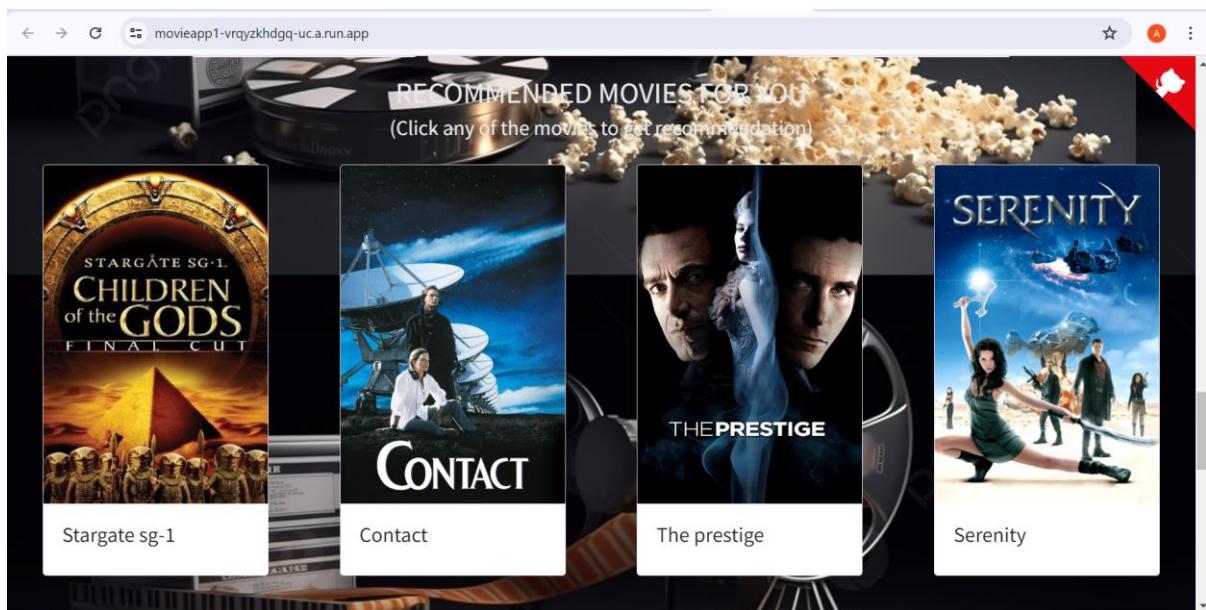
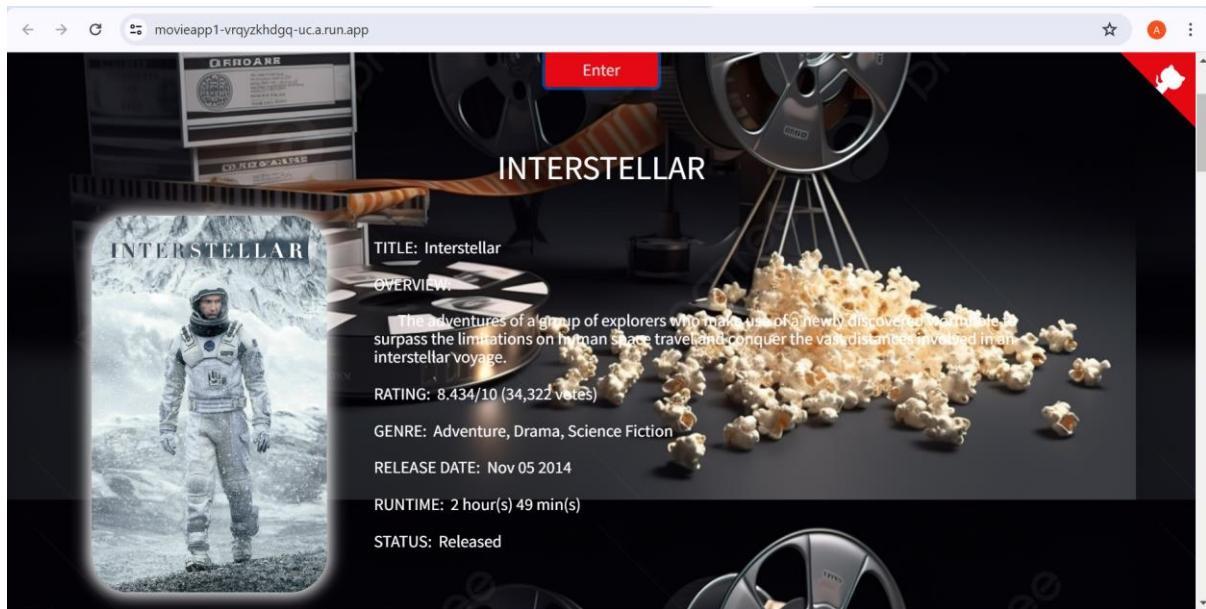
- **(5 pts) Front-end client development**

Flask is an important tool for front-end client development since it allows you to create dynamic and adaptable web interfaces. Using Flask's features, developers can quickly design and build interactive

user interfaces, improving the overall user experience. Furthermore, REST requests improve communication between the front-end and backend, allowing for smooth integration and data sharing. This strategy has various advantages. For starters, separating front-end and backend components makes development and maintenance easier, allowing teams to collaborate freely and quickly. Furthermore, RESTful APIs encourage flexibility and scalability, allowing for smooth integration and development of several components.

Developers may create visually appealing and user-friendly interfaces using the front-end client's HTTP requests and CSS style, increasing user engagement and satisfaction. In summary, using Flask for front-end development in conjunction with RESTful APIs produces a versatile, scalable, and responsive web interface that caters to both developers and users. Furthermore, when installed on Google Cloud Platform, a website address is supplied post-deployment, allowing quick access to the deployed application.





To access the application: <https://movieapp1-vrqyzkhdgq-uc.a.run.app/>

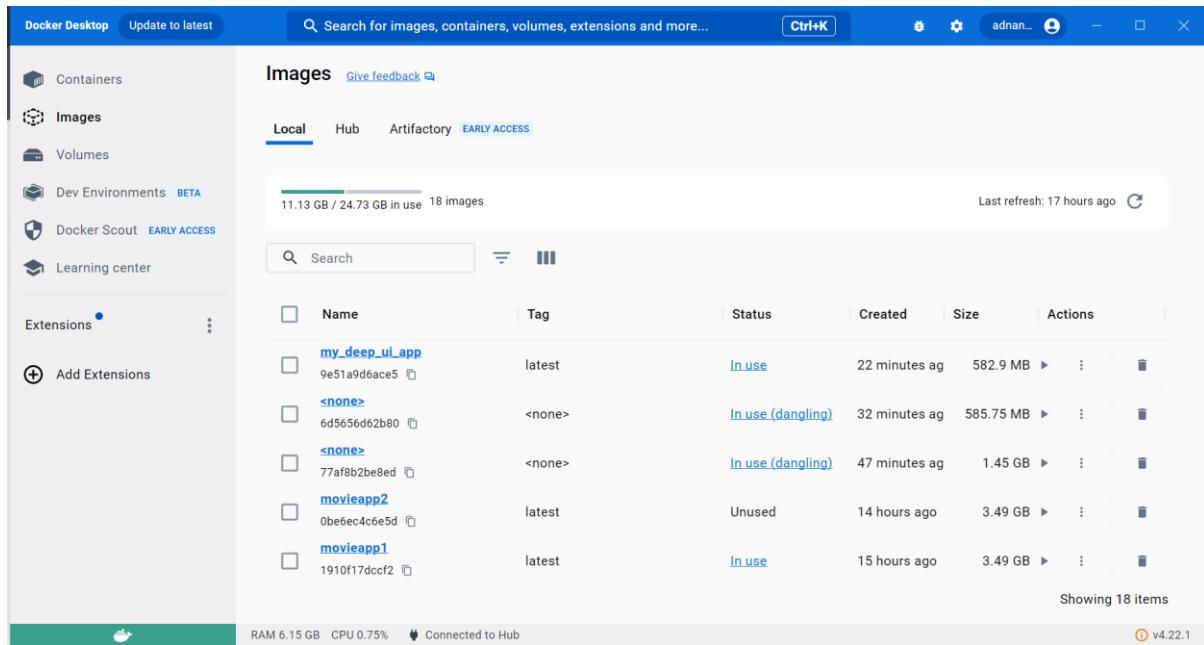
### 3.3 Integration and Deployment

- Packaging and containerization

Docker is a valuable tool for packaging and containerization since it encapsulates a program and its dependencies in portable containers. These containers create a standardized environment that ensures consistency between the development and production environments. Using Docker's containerization features, the packed application is deployed as a Docker container, allowing for smooth deployment and portability across several platforms. Furthermore, Docker Hub acts as a centralized repository for

storing and maintaining Docker images, simplifying the distribution and deployment process. This strategy has several advantages. First, Docker containers are portable, allowing them to run on any platform that supports Docker. This ensures consistent behavior across environments, simplifies the deployment process and reduces compatibility difficulties. Containers also enable isolation, effectively separating processes and resources to avoid conflicts between applications and dependencies. This separation improves security and stability, lowering the chance of system failure and providing consistent performance in a variety of computing scenarios. Overall, Docker's containerization and packaging features, together with Docker Hub's repository management, enable rapid and consistent application deployment while encouraging portability and isolation for increased dependability and scalability.

```
👉 dockerfile > ...
1  # Use the official Python image from the Docker Hub
2  FROM python:3.8-slim
3
4  # Set the working directory in the container
5  WORKDIR /app
6
7  # Copy the current directory contents into the container at /app
8  COPY . /app
9
10 # Install any needed packages specified in requirements.txt
11 RUN pip install --no-cache-dir -r requirements.txt
12
13 # Make port 5000 available to the world outside this container
14 EXPOSE 5000
15
16 # Define environment variable
17 ENV NAME World
18
19 # Run app.py when the container launches
20 CMD ["python", "main.py"]
21
```



- **Integration with a CI/CD Pipeline**

With Google Cloud Run, users can seamlessly integrate their GitHub repositories with the deployed container. This configuration ensures that whenever new updates are pushed to the repository, the deployed container automatically redeploys with the latest version.

Google Cloud Build interface showing the creation of a trigger. The left sidebar includes options for Dashboard, History, Repositories, Triggers (selected), and Settings. The main area shows the 'Create trigger' form with the following settings:

- Source** section:
  - Repository generation: 1st gen (selected)
  - Repository: adnanehar/Movie-Recommendation-Final (GitHub App)
  - Branch: ^main\$
  - Trigger only for a branch that matches the given regular expression [Learn more](#)
  - Invert Regex
  - Matches the branch: main
- Configuration** section:
  - Type:

Screenshot of the Google Cloud Platform Cloud Build interface showing the 'Create trigger' configuration page.

**Left Sidebar:**

- Dashboard
- History
- Repositories
- Triggers** (selected)
- Settings
- Release Notes

**Right Panel - Configuration:**

Type: Autodetected (radio button selected). A cloudbuild.yaml or Dockerfile will be detected in the repository.

Cloud Build configuration file (yaml or json) (radio button)

Dockerfile (radio button)

Buildpacks (radio button)

Location: Repository (radio button selected). adnanehar/Movie-Recommendation-Final (GitHub App).

Inline (radio button). Write inline YAML.

**Advanced:**

Approval: Require approval before build executes (checkbox)

Build logs: Send build logs to GitHub (checkbox)

Screenshot of the Google Cloud Platform Cloud Build interface showing the 'Triggers' list page.

**Left Sidebar:**

- Dashboard
- History
- Repositories
- Triggers** (selected)
- Settings
- Release Notes

**Top Bar:**

- + CREATE TRIGGER
- CONNECT REPOSITORY
- MANAGE REPOSITORIES
- LEARN

**Region:** europe-west9

**Filter:** Enter property name or value

| Name   | Description | Repository                           | Event           |
|--|-------------|--------------------------------------|-----------------|
| movie  | -           | adnanehar/My_Movie_Recommender       | Push to br: RUN |
| movie20  | -           | adnanehar/Movie-Recommendation-Final | Push to br: RUN |
| rmpgab-movieapp1-us-central1-adnanehar-My-Movie-Recommenderdrk | Build a...  | adnanehar/My_Movie_Recommender       | Push to br: RUN |

**Bottom Panel - Trigger Details:**

**Trigger:** rmpgab-movieapp1-us-central1-adnanehar-My-Movie-Recommenderdrk

**Source:** adnanehar/My\_Movie\_Recommender

**Branch:** main

**Commit:** 854a077

**Build Summary:** 3 Steps

- 0: Build: build --no-cache -t \${\_AR\_H...}
- 1: Push: push \${\_AR\_HOSTNAME}/\${\_S...
- 2: Deploy: gcloud run services update...

**Environment Variables:**

|                        |   |
|------------------------|---|
| Branch                 | main  |
| Commit                 | 854a077   |
| Image                  | \$_AR_HOSTNAME/\${PROJECT_ID}/cloud-run-source-deploy/\${REPO_NAME}/\${SERVICE_NAME}:\${COMMIT_SHA}   |
| Service Account        | 970408374802-compute@developer.gserviceaccount.com  |
| User substitutions     | <ul style="list-style-type: none"> <li>_SERVICE_NAME: movieapp1</li> <li>_DEPLOY_REGION: us-central1</li> <li>_trigger_id: 37a35ea2-8f8b-4ba-8de9-9957fb364ad4</li> <li>_PLATFORM: managed</li> <li>_AR_HOSTNAME: us-central1-docker.pkg.dev</li> </ul> |
| Built-in substitutions | <ul style="list-style-type: none"> <li>BRANCH_NAME: main</li> <li>BUILD_ID: 36ccb00e-4f6b-4231-b94a-514309867b5a</li> <li>COMMIT_SHA: 854a077cc8fd6a317341bef2e758b7f874ae422c</li> </ul>   |

- **Hosting the application**

The application is hosted on Google Cloud Run

The screenshot shows the Google Cloud Run interface. At the top, there's a navigation bar with 'Google Cloud', 'My First Project', 'CI CD', and various icons. Below that, a main header says 'Cloud Run' and 'Service details'. There are buttons for 'EDIT & DEPLOY NEW REVISION' and 'SET UP CONTINUOUS DEPLOYMENT'. On the right, there are 'LEARN' and 'GCP' links. The main content area shows a service named 'movieapp1' with a green checkmark. It lists 'Region: us-central1' and 'URL: https://movieapp1-vrqyzkhdgq-uc.a.run.app'. It also shows 'Min instances: 0' and a 'Build History' link.

### 3.4. Model Serving and online testing

- Model serving runtime

The model is served through google cloud with exposing the endpoints using flask.

- Serving mode (batch, on demand to a human, on demand to a machine)

**Serving Mode: On Demand to a Human**

The key interaction in Movie Roulette is for users to search for movies based on their tastes and then receive personalized suggestions. Given this user behavior, offering predictions on demand to human users is well aligned with the application's functionality and user experience objectives. Here's why.

**Real-time Interaction:** Movie suggestions are a dynamic and personalized process that is significantly affected by user feedback and preferences. Serving predictions on-demand enables users to receive real-time suggestions as they engage with the application, improving their entire experience.

**User interaction:** Movie Roulette promotes active user interaction and discovery by providing predictions on demand. Users may search for movies depending on their current mood, tastes, or particular criteria, resulting in more meaningful interactions and longer time spent on the site.

**Tailored Recommendations:** Every user has different movie choices, and their search queries reflect their distinct likes and interests. Serving predictions on demand allows Movie Roulette to provide personalized suggestions in response to each user's search, assuring relevance and customisation.

**Instant Feedback Loop:** Serving predictions on-demand establishes an instant feedback loop between the user's search query and the recommendations made. Users get instant answers based on

their input, allowing them to swiftly adjust their search criteria or explore alternate options, increasing their overall happiness with the program.

**Scalability and Efficiency:** Unlike batch processing, which requires vast amounts of data to be processed at preset intervals, Movie Roulette's on-demand predictions allow it to scale dynamically in response to customer demand. This provides effective resource use and responsiveness, especially during peak demand times.

**Adaptability to User purpose:** Users' search queries frequently express implicit purpose, such as studying a certain genre, discovering new releases, or locating films comparable to their favorites. Serving predictions on-demand allows Movie Roulette to respond to changing user intentions in real time, offering appropriate suggestions that match user expectations.

- Online testing:

A/B testing are an important approach for comparing different model iterations in order to improve performance. Through iterative analysis of test results, the model is constantly refined, bringing it closer to set performance criteria. Google Cloud Run emerges as a critical enabler in this process, providing a streamlined approach for running A/B tests through an easy dashboard interface. This integrated feature enables teams to quickly manage and monitor various models, providing real-time insights and informed decision-making to support continuous improvement programs.

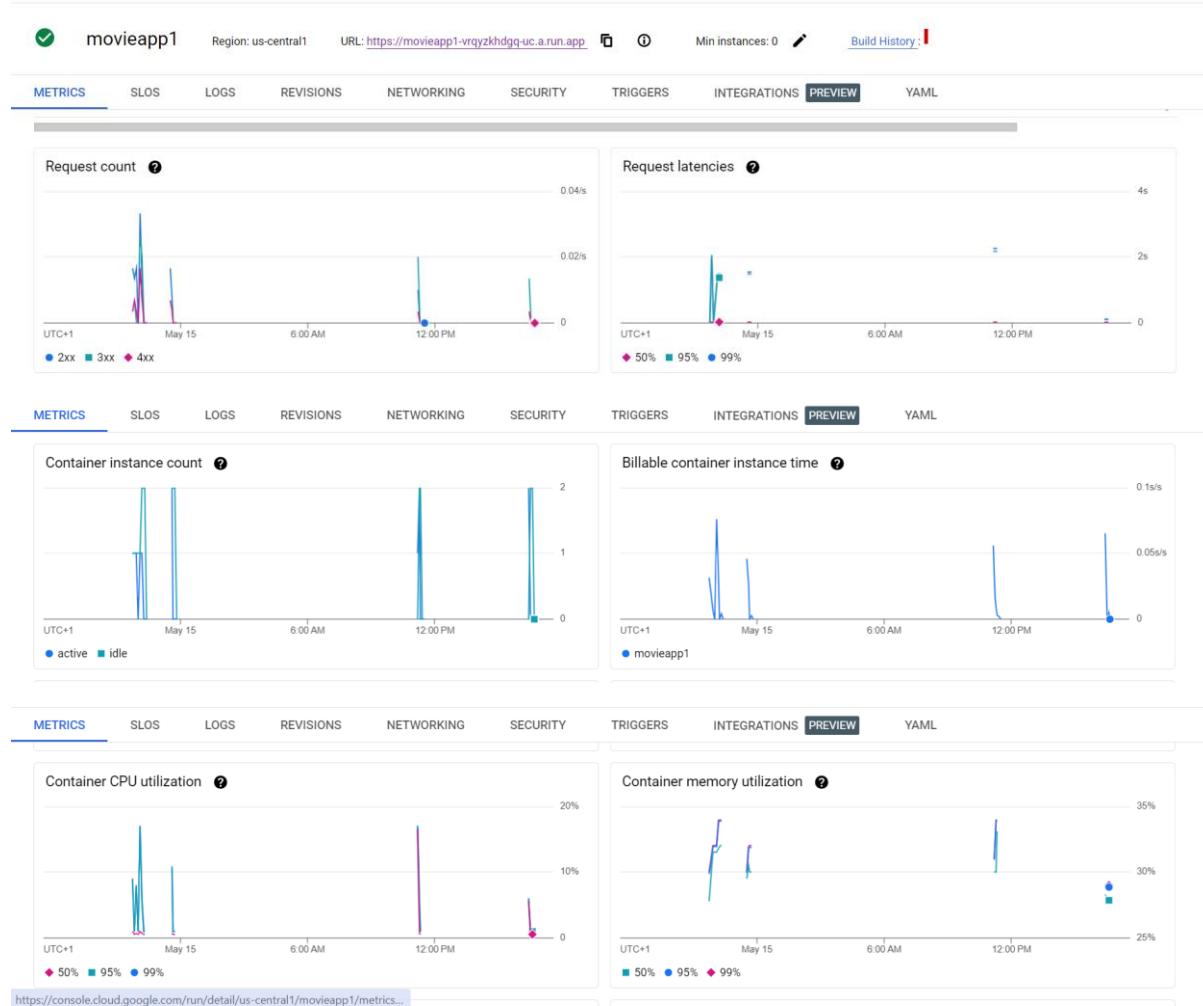
## 4. Monitoring and Continual Learning

Monitoring and continuous learning are key components of the Movie Roulette machine learning system, as they ensure the model's accuracy, efficiency, and adaptability to changing data and user preferences. Here, we look at the specific tools and approaches used, emphasizing their benefits and the reasoning behind their selection.

### 4.1. Resource Monitoring

Resource monitoring is critical to ensure the best performance and dependability of any system. Teams may maintain smooth operations by regularly monitoring critical indicators such as CPU utilization, memory consumption, and network activity. I created extensive monitoring solutions targeted to the unique requirements of the deployed application by leveraging Google Cloud's strong monitoring technologies, such as Stackdriver Monitoring and Logging. This provided real-time

visibility into the health and performance of multiple components, allowing for rapid interventions and ongoing optimization efforts to maintain service dependability and customer happiness.

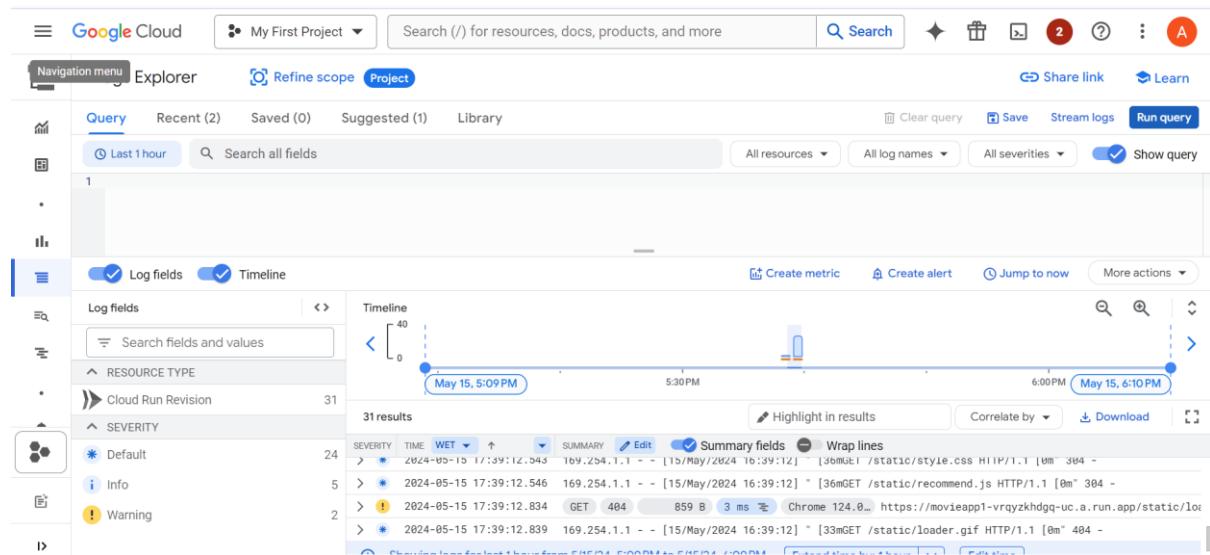


## 4.2. Data distribution drift monitoring

Data distribution shift refers to changes in the statistical features of input data that occur over time and have a substantial influence on machine learning model performance. Monitoring for such adjustments is critical for ensuring the correctness and dependability of deployed models.

To solve this, I used Google Cloud's Logging and Monitoring services to set up continuous monitoring of data dissemination. Using these services, I established automated methods to identify shifts in data distribution using statistical approaches such as the Kolmogorov-Smirnov (KS) test. This test analyzes the cumulative distribution functions of two datasets, allowing for the discovery of substantial variations in distributions.

When a data distribution shift is observed, Google Cloud Logging creates extensive logs that describe the nature and scope of the move. These logs give vital insights into the changes occurring within the input data, assisting in identifying possible issues influencing model performance.



#### 4.3. Continual Learning: CT/CD pipeline

For Continuous Training, I used Google Cloud's robust data processing and storage capabilities to rapidly ingest and preprocess enormous amounts of data. I built scalable pipelines to handle a variety of data sources and formats using data warehousing technologies such as BigQuery and Cloud Dataflow. This guaranteed that the model training datasets were constantly up to date and represented the most recent information.

In case of any data drifts, it will be displayed here:

#### 4.4. (3 pts) Pipeline orchestration

In our project, pipeline orchestration on Google Cloud is seamlessly achieved through Cloud Composer, a fully managed workflow orchestration service built on Apache Airflow. We define our workflows as directed acyclic graphs (DAGs), encapsulating the sequence of tasks required for data processing, model training, and deployment. These DAGs are deployed to our Cloud Composer environment, where we schedule their execution at predetermined intervals or trigger them manually as needed. Monitoring and management of these workflows are facilitated through the Airflow UI, providing real-time visibility into task statuses and logs. Leveraging integrations with other Google Cloud services, such as BigQuery for data processing and AI Platform for model training, our pipeline orchestrates the entire machine learning lifecycle with efficiency and reliability.

## 5. Conclusion

To summarize, the creation and implementation of the Movie Roulette machine learning system has been a journey filled with important accomplishments, useful lessons learned, and promising future paths. Movie Roulette has evolved into a complex platform that provides individualized movie suggestions to consumers all over the world thanks to painstaking planning, clever use of cutting-edge technology, and an unwavering dedication to quality.

- **Summary of Achievements**

Throughout the project, I gained a varied variety of skills and expertise that dramatically improved our machine learning pipeline:

## **Technology Acquired:**

- ZenML, a powerful MLOps framework that provides complete control over machine learning pipelines.
- Developed skills in Flask backend development, allowing for the implementation of fast and scalable model services.
- Acquired experience in Docker and DockerHub for application containerization and simplified deployment.
- Gained experience with CI/CD pipelines and google cloud.
- Designed and built a scalable ML pipeline architecture to provide dependability and flexibility.
- Tools such as Prometheus and Grafana have been integrated for real-time resource monitoring and performance improvement.
- Continuous Learning and Adaptation.
- CT/CD pipelines were implemented to allow for continuous learning, guaranteeing that the model can adapt to changing data distributions.
- Established monitoring tools for model performance and data distribution drift, allowing for proactive reactions to aberrations.

- **Lessons Learned**

- Lessons learned and project insights:
- Recognized the vital role of automation in expediting development processes, while highlighting the importance of rigorous planning and excellent mistake management.
- Recognized the problems of A/B testing, highlighting the significance of thorough testing and validation of automation scripts.
- Highlighted the importance of preserving data quality in machine learning systems, which necessitates continual monitoring and validation techniques.

- Scalability issues were addressed, highlighting the significance of designing systems that can handle increasing user demand and data volumes.
  - To improve user pleasure, we prioritized user-centric design concepts, concentrating on intuitive interfaces and tailored suggestions.
- **Future Directions**
- Plan to investigate sophisticated experimental approaches such as multi-armed bandit testing to improve model performance and user engagement.
  - Consider moving the application to cloud infrastructure for better scalability, dependability, and accessibility.
  - Natural language processing (NLP) techniques will be integrated to analyze user reviews and comments, improve recommendation accuracy, and increase user happiness.
  - Improve collaborative filtering algorithms to include user preferences and behavior, resulting in more accurate and personalized suggestions.
  - Consider creating a thriving community around Movie Roulette by encouraging user participation and cooperation to promote ongoing improvement and innovation.

## References