

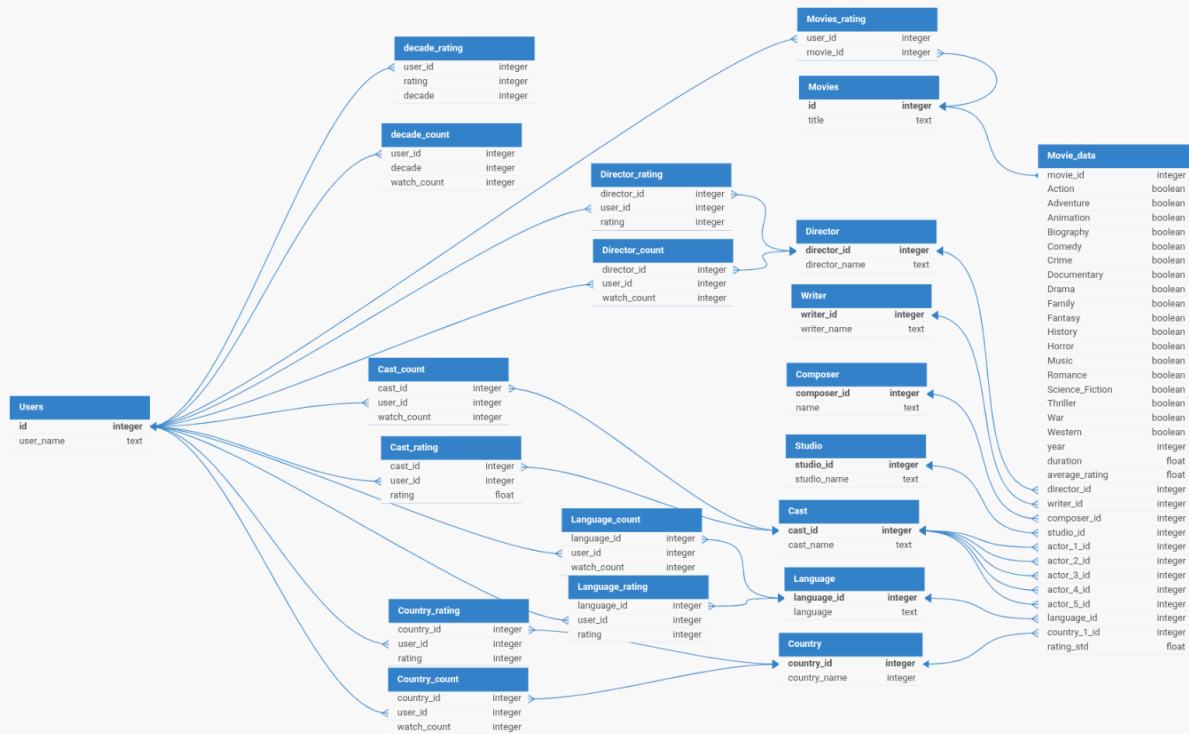
Phase 2

In The Name Of GOD

Datascience Prjct: Phase 2

Members: Sepanta Ghonoodi, Shahab Sherafat, Mahdi Yari

Database schema



Generated with [erd.dbdesigner.net](https://dbdesigner.net)

Link: <https://dbdesigner.page.link/QQwD3K8r53MfD5P17>

Queries On Database

1. A simple one:

```
SELECT movie_id, rating
FROM Movies_ratings;
```

This query gets the ratings of each movie from ratings table.

Result:

	movie_id	rating
1237	186	3.0
1238	186	4.0
1239	186	4.0
1240	451	4.5
1241	451	2.5
1242	451	5.0
1243	451	5.0
1244	451	3.0
1245	451	5.0
1246	451	4.0
1247	451	3.0
1248	451	5.0

2. Selecting the movies with highest rates using filtering and sorting:

```
SELECT movie_id, AVG(rating) as avg_rating
FROM Movies_ratings
GROUP BY movie_id
HAVING AVG(rating) > 4.0
ORDER BY avg_rating DESC;
```

Phase 2

Result:

	movie_id	avg_rating
1	335	5.0
2	190	5.0
3	153	4.89285714285714
4	218	4.75
5	523	4.71875
6	724	4.70909090909091
7	163	4.69495091164095
8	649	4.67362924281984
9	303	4.62686567164179
10	598	4.62570888468809
11	517	4.6216814159292
12	429	4.61979166666667

3. Joining movies and ratings tables:

```
SELECT m.movie_id, m.movie_title, m.genre, AVG(r.rating) AS avg_rating
FROM Movies m
JOIN Movies_ratings r ON m.movie_id = r.movie_id
GROUP BY m.movie_id
ORDER BY avg_rating DESC;
```

Result:

	movie_id	title	avg_rating
1	335	Children of Heaven	5.0
2	190	12th Fail	5.0
3	153	The Grapes of Wrath	4.89285714285714
4	218	Das Boot	4.75
5	523	The Best Years of Our Lives	4.71875
6	724	The General	4.70909090909091
7	163	12 Angry Men	4.69495091164095
8	649	Parasite	4.67362924281984

Phase 2

4. Counting the number of movies each user's seen:

```
SELECT user_id, COUNT(DISTINCT movie_id) AS num_moviesRated
FROM Movies_ratings
GROUP BY user_id
ORDER BY num_moviesRated DESC;
```

Result:

	user_id	num_moviesRated
2	296	711
3	147	710
4	603	705
5	653	698
6	277	696
7	847	694
8	826	694
9	727	694
10	537	694
11	561	693
12	343	693
13	272	693

5. Selecting the movies based which were released after 2010:

```
SELECT movie_id, title, year
FROM Movies
WHERE year > 2010;
```

Result:

	movie_id	title	year
1	2	Ant-Man	2015
2	3	Jai Bhim	2021
3	4	The Nice Guys	2016
4	5	Aftersun	2022
5	6	The Cabin in the Woods	2011

Feature Engineering

For this part we added some new features which rationally will help the model recommend better. We tried to consider the logic behind what the model will do and give it some useful extra information about how each movie or each user are.

1. Favorite Genre For Each User:

According to the ratings we have we found the average rating of each genre given by each user and then we assigned the genre with the maximum rating for each user.

2. Ratings STD For Each Movie:

We know that the movies which have more varied ratings might help the model categorize and recommend users to each other better. So using the ratings we found the STD of the ratings for each movie and assigned it as a feature of movies.

3. Analyzing Each Person's Personality:

It is possible to find a person's characteristics by the movies they've seen! Using the big five model we tried to define a new one hot encoded feature for users which determines the five personality traits for each user according to the genres they like.



Phase 2

For each personality trait inside each user we calculated the average rating of the genres that trait might be related to. Then we assigned a threshold(3.5) for the ratings to know whether a person likes that genre or not. If a person likes all the genres of a trait, the one hot feature of that trait will be equal to 1.

We normalized the numerical features using sklearn and removed the null data. About removing unuseful features, because each feature is an independent feature from the others and we've concluded that there might be no irrelevant feature and no correlation between them. So we decided to keep all the features.

Phase 2

CI/CD Implementation

```
load_data.py script.sql -- SQLite.sql pivot_data.csv U res.txt U pipeline.yml X
.github > workflows > pipeline.yml
8 jobs:
9   pipeline-check:
11     steps:
12       - name: Checkout code
14
15       - name: Set up Docker Build
16         uses: docker/setup-buildx-action@v3 "buildx": Unknown word.
17
18       - name: Build Docker image
19         run: |
20           docker build -t phase_2_image .
21
22       - name: Run Docker container
23         run: |
24           docker run --rm phase_2_image
25
```

github.com/sepanta-ghonoodi/letterbox-recommendation-system/actions

sepanta-ghonoodi / letterbox-recommendation-system

Code Issues Pull requests Actions Projects Security Insights Settings

Actions New workflow

All workflows

pipeline

Management

- Caches
- Attestations
- Runners
- Usage metrics
- Performance metrics

All workflows

Showing runs from all workflows

16 workflow runs

	Event	Status	Branch	Actor
adding data file	pipeline #16: Commit 0026b11 pushed by sepanta-ghonoodi	main	1 minute ago	...
adding feature engineering	pipeline #15: Commit bbd81f8 pushed by sepanta-ghonoodi	main	9 minutes ago	...
updating pipeline	pipeline #14: Commit 4130118 pushed by sepanta-ghonoodi	main	4 hours ago	...
fixing dockerfile	pipeline #13: Commit d4849dd pushed by sepanta-ghonoodi	main	2 days ago	...
changing git actions for docker	pipeline #12: Commit 0a9eb63 pushed by sepanta-ghonoodi	main	2 days ago	...

Phase 2

MLOps

Docker file:

```
Dockerfile X pipeline.py preprocess.ipynb db_export.ipynb feature_engineering.ipynb
Dockerfile > ...
You, 2 days ago | 1 author (You)
1 FROM python:3.8-slim
2 WORKDIR /app
3 COPY requirements.txt .
4
5 RUN pip install --upgrade pip && pip install -r requirements.txt
6 COPY phase2/ . You, 2 days ago * fixing dockerfile ...
7 CMD ["python3", "pipeline.py"]
8
```

Docker build:

```
~/D/IS4/DSC/letterbox-project (main*) > docker build --network=host -t phase_2_image . wren@wren
2025/05/05 23:40:02 In: [!string()]
2025/05/05 23:40:02 Parsed entitlements: []
[+] Building 2.6s (10/10) FINISHED
=> [internal] load build definition from Dockerfile
=> == transferring dockerfile: 289B 0.1s
=> [internal] load metadata for docker.io/library/python:3.8-slim 0.4s
=> [internal] load .dockerignore 1.7s
=> == transferring context: 2B 0.1s
=> [1/5] FROM docker.io/library/python:3.8-slim@sha256:1d52838af602b4b5a031beb13a0e4a0873280665ea7be7769ce2382f29c5a613f 0.4s
=> [internal] load build context 0.1s
=> == transferring context: 1.03kB 0.8s
=> CACHED [2/5] WORKDIR /app 0.8s
=> CACHED [3/5] COPY requirements.txt 0.8s
=> CACHED [4/5] RUN pip install --upgrade pip && pip install -r requirements.txt 0.8s
=> [5/5] COPY phase2/ . 0.2s
=> exporting to image 0.2s
=> == exporting layers 0.1s
=> == writing image sha256:6d29d96cdd6042d50137b88350e65210a9820e477f104be4edf14ede304c8fd1 0.0s
=> == naming to docker.io/library/phase_2_image 0.0s
```

Docker run:

```
~/D/IS4/DSC/letterbox-project (main*) > docker run --rm -v $(pwd)/data:/app/data phase_2_image wren@wren
write Cast_df to file
write Director_df to file
write Language_df to file
write Country_df to file
write decade_rating_df to file
write decade_count_df to file
write User_df to file
write Director_rating_df to file
write Director_count_df to file
write Language_rating_df to file
write Languages_count_df to file
write Countries_count_df to file
write Countries_rating_df to file
write Cast_rating_df to file
write Cast_count_df to file
write new_data_df to file
write Movie Cast_df to file
write sqlite sequence df to file
write Movie Genres df to file
write Writer_df to file
write Composer_df to file
write Studio_df to file
write Movies_df to file
write Movie Composer_df to file
write Movie Director_df to file
write Movie Writer_df to file
write Movie Country_df to file
write Movie Language_df to file
write Movie Studio_df to file
write ratings_df to file
write Genre_df to file
write User Favorite Genre_df to file
write Final Movie Data df to file
write Movies ratings df to file
Executing: 100% [ ] 6/6 [00:02:00:00, 2.85cell/s]
Executing: 100% [ ] 11/11 [00:01:00:00, 6.78cell/s]
```