# Advanced Database Management and Query Optimization for a Movie Dataset

*Phase 2*

Rohan Patel
ESDS
University At Buffalo
Buffalo, NY
rpatel38@buffalo.edu

Jay Thanki
ESDS
University At Buffalo
Buffalo, NY
jayyoges@buffalo.edu

*Abstract*— **This detailed report delves into the management and optimization of a complex movie database derived from various sources such as Movies.csv, Links.csv, Tags.csv, Ratings.csv, and ml-youtube.csv. The research highlights the importance of data integrity, efficient query execution, and database normalization. It provides a case study of the difficulties faced and creative solutions used while managing enormous datasets in a PostgreSQL environment.**

**Keywords— SQL, ER Model, Relational Model, MovieLens Data, Data Administration, Query Optimization, Deployment.**

## I. Handling Large Datasets (Task 5)

We had a number of noteworthy difficulties in handling this large dataset, especially in the data loading stage. The CSV files' various formats were the main source of the difficulty. In particular, several files applied the field terminators—double quotes—inconsistently in various rows and columns. To guarantee effective data integration, distinct data loading syntax had to be created for every file due to the formatting discrepancy.

The Ratings file, which contained more than two million entries, presented another well-known issue. Long loading times into the temporary table—a crucial stage in our data processing pipeline—were caused by the sheer volume of this data. This file's size not only made loading take longer, but it also made it take longer to change the data types in these columns before integrating them into the main table that is always open.

We used temporary tables, which run in memory and hence provide faster processing speeds than the main table, to speed up the data loading process. But indexing, which is usually a powerful strategy for handling enormous datasets, was considered impractical in this case. The choice was made based on the fact that these tables were only meant to be temporary; indexing them would have required an unnecessary use of resources and space because they were going to be deleted after the data was cleaned and loaded.

## II. SQL Query Testing and Execution (Task 6)

We have tested our database using a number of SQL queries in this project section. A broad range of actions, such as insertion, deletion, updating, and intricate selection procedures, were included in the design of these tests. An extensive summary of the queries utilized and their particular goals is given below:

### 1. Identifying Movies with Highest Average Ratings

- **Query:** The movies with the highest average ratings were found using a SELECT statement, with a minimum of 1000 ratings required for a movie to be considered. In order to organize the Movies and Ratings tables according to movie titles, an inner join was performed.

- **Purpose:** To demonstrate how aggregate functions with GROUP BY and HAVING clauses can yield useful information, like a list of the highest rated films in our dataset.



Fig. 1. Identifying Movies with Highest Average Ratings

### 2. Calculating Number of Ratings per Genre

- **Query:** To determine the total number of ratings and average rating per genre, a sophisticated SELECT statement linking the Genres, Ratings, and Genres_Master databases was required.

- **Purpose:** to illustrate the capacity to aggregate data across various dimensions—in this case, genres—and link multiple tables.



Fig. 2. Calculating Number of Ratings per Genre

### 3. Top Movies in Each Genre

- **Query:** This query, which required that the films have at least 1000 ratings, retrieved the top 3 films with the highest rankings in each category by using subqueries and window functions.

- **Purpose:** To demonstrate how to use window functions (RANK() OVER()) and subqueries, two advanced SQL features, for more in-depth data analysis.

Fig. 3 Top Movies in Each Genre

## 4. Average Movie Rating by Tag

- **Query:** The average rating for movies, sorted by tags, was determined by doing an INNER JOIN across the Tags, Ratings, and Tags_Master columns.

- **Purpose:** to provide an example of how to carry out intricate data aggregations across various categories—in this case, movie tags.



Fig. 4 Average Movie Rating by Tag

## 5. Identifying Most Popular Tags

- **Query:** To determine each tag's popularity, a straightforward yet powerful SELECT statement was utilized to count the occurrences of each one.

- **Purpose:** to emphasize how popular movie tags are determined by using the GROUP BY and ORDER BY clauses.



Fig. 5 Top Movies in Each Genre

## 6. & 7. Movies Filtered by Specific Genres and Tags

- **Queries:** These SELECT lines demonstrate the database's capacity to manage simple data retrieval based on predetermined criteria by filtering movies by particular genres and tags.

Fig. 6 Movies Filtered by Specific Genre



Fig. 7 Movies Filtered by Specific Tag

**8. Updating Missing Links in MovieLinks**

- **Query:** A UPDATE statement was used to add missing YouTube, TMDb, and IMDb links in the MovieLinks table using matching data from the Movies table.

- **Purpose:** to demonstrate the use of UPDATE procedures to synchronize data across linked tables.



Fig. 8 Updating Missing Links in MovieLinks

**9. Deletion of Redundant Data**

- **Query:** The process entailed eliminating films and categories that weren't mentioned in the Ratings or Tags tables.

- **Purpose:** to guarantee that there are no orphaned records in the database and to preserve data integrity.



Fig. 9 Deletion of Redundant Data

**10. Inserting a New Movie and Associated Links**

- **Query:** We added 'Dune: Part Two' as a new movie and its associated links to the database by using the DO command.

- **Purpose:** to demonstrate the insertion procedure and the connections between various tables in our database.



Fig. 10 Insert New Movie

Fig. 11 Insert Links of the Movie

## 11. Updating and Deleting Movie Data

- **Query:** These queries concerned the removal of movie links and individual movies, as well as the upgrading of particular movie details (such as a new trailer link).

- **Purpose:** To illustrate the dynamic nature of the database in handling updates and deletions, reflecting real-world changes.



Fig. 12 Update Movies. Year



Fig. 13 Updating New Trailer Link



Fig. 14 Deleting MovieLinks of Movie



Fig. 15 Deleting Specific Movie from Movies

We ran each of these queries on our database and took screen grabs of the responses. These screenshots demonstrate the efficiency of the queries and the resilience of our database to handle different kinds of data operations. The variety of these queries, which range from straightforward data retrievals to intricate updates and aggregations, highlights the adaptability and effectiveness of our database system in handling and evaluating huge datasets.

## III. Query Execution Analysis (Task 7)

We performed a critical study of query execution throughout this project phase in order to pinpoint performance bottlenecks and apply fixes. We examined three troublesome queries using PostgreSQL's EXPLAIN tool, calculating their costs and coming up with ways to make them better.

## 1. Average Rating of Movies by Tag

- **Initial Query and Issue:** The purpose of the query was to find the average movie rating by tag, but because of the quantity of the dataset, it was originally taking an unfeasible 6-7 minutes. The query was executing extremely inefficient sequential scans across three tables, as the EXPLAIN investigation showed.

Fig. 16. Explain tool result on query



Fig. 17. Before optimization

- **Optimization Strategy:** Using the current composite primary key index on UserId and MovieID, we changed the join condition to include R.UserID = T.UserID in order to increase efficiency. The query execution time was slashed to a matter of seconds as a result.



Fig. 18. Post optimization

- **Alternative Approach:** Another potential solution could be to create a new table containing aggregated data for each movie. While this would significantly enhance query efficiency, it would involve a trade-off with the granularity of user data.

## 2. Movies with Most Number of Ratings

- **Initial Query and Issue:** The purpose of this query was to retrieve the highest-rated movies, and it was initially executing in roughly 5.5 seconds per execution. The main reason for the delay was found to be the lack of a specific index on the Ratings table's MovieID column.



Fig. 19. Explain tool result on query



Fig. 20. Before Optimization

- **Optimization Strategy:** We created a non-clustered index on the MovieID column of the Ratings table. This optimization led to a 40% performance improvement, reducing the query execution time to approximately 3.5 seconds.

```
1  CREATE INDEX idx_ratings_movieID ON Ratings (MovieID);
2  SELECT M.Title, COUNT(*) Ratings
3  FROM Ratings R
4  INNER JOIN Movies M ON M.MovieID = R.MovieID
5  GROUP BY M.Title
6  ORDER BY COUNT(*) DESC
7  --DROP INDEX IF EXISTS idx_ratings_movieID;
```

| | title<br>character varying (1000) | ratings<br>bigint |
|---|---|---|
| 1 | Pulp Fiction (1994) | 67310 |
| 2 | Forrest Gump (1994) | 66172 |
| 3 | Shawshank Redemption, The (1994) | 63366 |
| 4 | Silence of the Lambs, The (1991) | 63299 |
| 5 | Jurassic Park (1993) | 59715 |
| 6 | Star Wars: Episode IV - A New Hope (1977) | 54502 |
| 7 | Braveheart (1995) | 53769 |
| 8 | Terminator 2: Judgment Day (1991) | 52244 |
| 9 | Matrix, The (1999) | 51334 |
| 10 | Schindler's List (1993) | 50054 |
| 11 | Toy Story (1995) | 49695 |
| 12 | Fugitive, The (1993) | 49581 |
| 13 | Apollo 13 (1995) | 47777 |
| 14 | Independence Day (a.k.a. ID4) (1996) | 47048 |
| 15 | Usual Suspects, The (1995) | 47006 |

Total rows: 1000 of 26729    Query complete 00:00:03.163

Fig. 21. Post optimization

## 3. Fetching Movies Year-wise in Range

**Initial Query and Issue:** Using the LIKE operator on the movie title, the query sought to filter movies inside a specified range of years. This method had a substantial negative impact on performance because it required a full table scan.

Fig. 22. . Explain tool result on query

```
1  SELECT M.Title, GM.Genre, CAST(AVG(Rating) AS NUMERIC(10,2)) Rating, COUNT(*) Ratings
2  FROM Movies M
3  INNER JOIN Genres G ON G.MovieID = M.MovieID
4  INNER JOIN Genres_Master GM ON GM.ID = G.GenreID
5  INNER JOIN Ratings R ON R.MovieID = M.MovieID
6  WHERE (
7          Title LIKE '%2010%'
8      OR Title LIKE '%2011%'
9      OR Title LIKE '%2012%'
10     OR Title LIKE '%2013%'
11     OR Title LIKE '%2014%'
12     OR Title LIKE '%2015%'
13         )
14 GROUP BY M.Title, GM.Genre
```

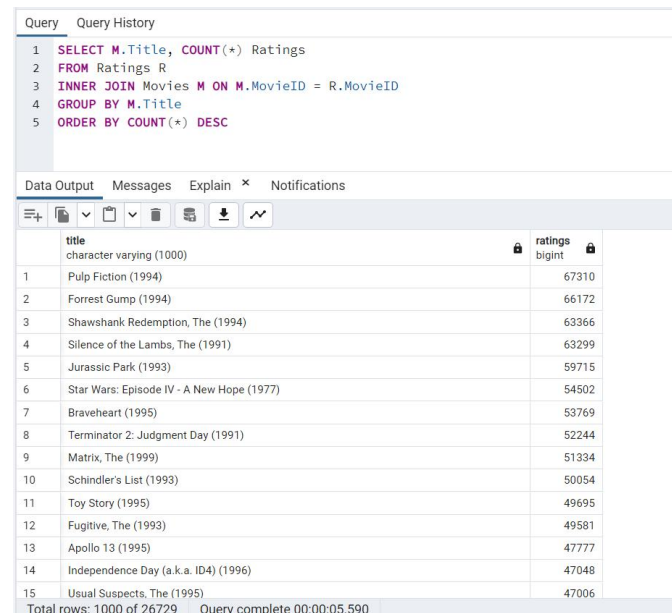| | title<br>character varying (1000) | genre<br>character varying (100) | rating<br>numeric (10,2) | ratings<br>bigint |
|---|---|---|---|---|
| 1 | '71 (2014) | Action | 3.66 | 35 |
| 2 | '71 (2014) | Drama | 3.66 | 35 |
| 3 | '71 (2014) | Thriller | 3.66 | 35 |
| 4 | '71 (2014) | War | 3.66 | 35 |
| 5 | #chicagoGirl: The Social Network Takes on a Dictator (2013) | Documentary | 3.67 | 3 |
| 6 | $ellebrity (Sellebrity) (2012) | Documentary | 2.00 | 2 |
| 7 | (A)sexual (2011) | Comedy | 3.33 | 3 |
| 8 | (A)sexual (2011) | Documentary | 3.33 | 3 |
| 9 | (A)sexual (2011) | Drama | 3.33 | 3 |

Total rows: 1000 of 9007    Query complete 00:00:02.766

Fig. 23. Before Optimization

**Optimization Strategy:** We introduced a new Year column in the Movies table, separating the year from the title. This allowed for direct querying on the Year column, thereby avoiding full table scans and significantly improving the query's efficiency.

```
1  SELECT M.Title, GM.Genre, CAST(AVG(Rating) AS NUMERIC(10,2)) Rating, COUNT(*) Ratings
2  FROM Movies M
3  INNER JOIN Genres G ON G.MovieID = M.MovieID
4  INNER JOIN Genres_Master GM ON GM.ID = G.GenreID
5  INNER JOIN Ratings R ON R.MovieID = M.MovieID
6  WHERE M.Year BETWEEN 2010 AND 2015
7  GROUP BY M.Title, GM.Genre
```

| | title<br>character varying (1000) | genre<br>character varying (100) | rating<br>numeric (10,2) | ratings<br>bigint |
|---|---|---|---|---|
| 1 | '71 (2014) | Action | 3.66 | 35 |
| 2 | '71 (2014) | Drama | 3.66 | 35 |
| 3 | '71 (2014) | Thriller | 3.66 | 35 |
| 4 | '71 (2014) | War | 3.66 | 35 |
| 5 | #chicagoGirl: The Social Network Takes on a Dic... | Documentary | 3.67 | 3 |
| 6 | $ellebrity (Sellebrity) (2012) | Documentary | 2.00 | 2 |
| 7 | (A)sexual (2011) | Comedy | 3.33 | 3 |
| 8 | (A)sexual (2011) | Documentary | 3.33 | 3 |
| 9 | (A)sexual (2011) | Drama | 3.33 | 3 |
| 10 | ??3D (2012) | Horror | 1.50 | 1 |
| 11 | [REC] 4: Apocalypse (2014) | Horror | 2.75 | 6 |
| 12 | [REC] 4: Apocalypse (2014) | Thriller | 2.75 | 6 |
| 13 | [REC]¡ 3 Gènesis (2012) | Horror | 2.67 | 46 |
| 14 | [REC]¡ 3 Gènesis (2012) | Thriller | 2.67 | 46 |
| 15 | +1 (2013) | Sci-Fi | 3.00 | 13 |

Total rows: 1000 of 9000    Query complete 00:00:02.736

Fig. 24. Post Optimization

These optimizations highlight the importance of understanding the underlying database structure and query execution plan. By analyzing and modifying our queries and database schema, we achieved significant performance improvements, ensuring faster and more efficient data retrieval. This exercise underscores the critical role of database tuning in managing large datasets and provides valuable insights into practical database optimization techniques.

## IV. Database Deployment (Bonus Task)

The task involved creating an interactive web application using Streamlit, a popular Python library for building web applications, to visualize and display query results from a PostgreSQL database. The primary focus was on a movie database, aptly named "Movie Mania", which allows users to search for movies based on different criteria such as title, year, rating, and genre.

Fig. 24. Home Page

### Technical Implementation

1. **Streamlit Integration**: Streamlit was chosen for its simplicity and effectiveness in creating data-driven web applications. The application's layout includes a sidebar for navigation and different sections for user interaction.
2. **Database Connection**: The application connects to a PostgreSQL database using SQLAlchemy, a Python SQL toolkit, ensuring a secure and efficient link for querying the database.

3. **User Interface Design**: The web application features an engaging user interface, with a banner and various input options like text input, number input, and sliders. This design enhances user experience, making it intuitive to search for movies.
4. **Dynamic Querying**: The application dynamically constructs SQL queries based on user input. These queries fetch movie data like title, genres, trailer links, average ratings, and number of ratings. This approach allows for real-time, customized data retrieval.
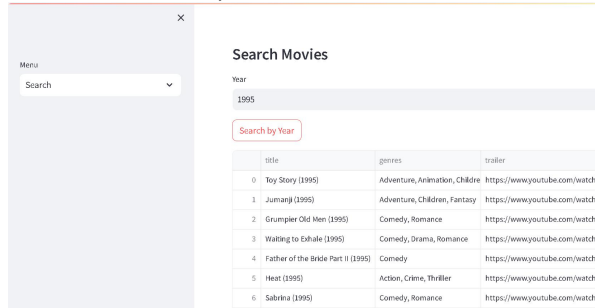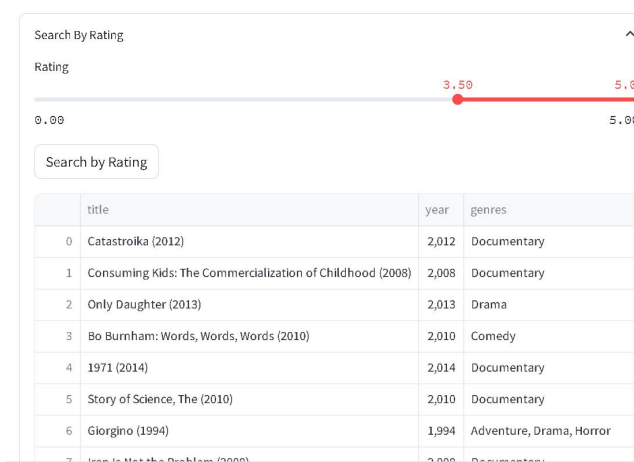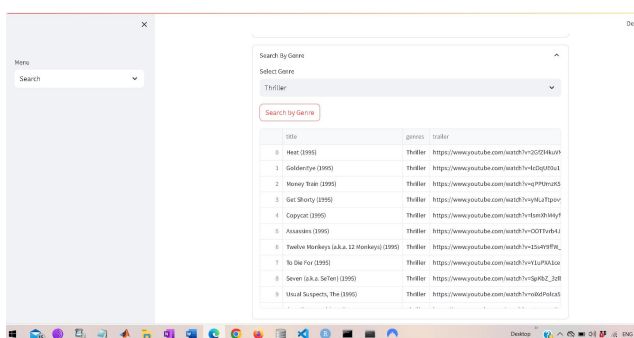


Fig. 25.Search by year



Fig. 26.Search by rating



Fig. 27.Search by genre



Fig. 28.Search by title

5. **Data Visualization**: Results from queries are displayed in a tabular format, providing a clear and organized presentation of data. Additional features like movie trailers are embedded directly into the application for an enriched user experience.
6. **Responsive Search Features**: The application offers various search functionalities:
   o **Search by Title**: Allows users to enter a movie name and retrieve relevant data.
   o **Search by Year**: Users can select a year to find movies released in that year.
   o **Search by Rating**: A slider to choose a rating range, returning movies within that rating spectrum.
   o **Search by Genre**: Users can select from a list of genres to find movies in the selected category.
7. **Error Handling and User Feedback**: The application provides user feedback and warnings, such as when no movies are found or a movie name is not entered.
8. **Accessibility and Performance**: Streamlit's lightweight framework ensures the application is accessible and performs efficiently, even with large datasets.

REFERENCES

[1] https://grouplens.org/datasets/movielens/20m/
[2] https://grouplens.org/datasets/movielens/20m-youtube/