

Programming Project Collaborative Filtering
Algorithms for Data Science



08/10/2025 – 01/11/2024

Herrá Alpuente, Luis Miguel:

luis-miguel.herra-alpuente@universite-paris-saclay.fr

Teacher:

Henri Paris, Pierre:

pierre-henri@universite-paris-saclay.fr

Programming Project: Collaborative Filtering

Algorithms for Data Science

Student: Herrá Alpuente, Luis Miguel

Teacher: Henri Paris, Pierre

Date: 08/10/2025 – 01/11/2025

Table of Contents

1. Problem Analysis (p. 3)

First Idea (p. 4)

Why collaborative filtering (p. 4)

Item-Item and User-User approach (p. 4)

User-User problems (p. 5)

2. Implementation (p. 5)

Data Structures Decision (p. 6)

Cosine Similarity (p. 7)

Compute movies similarities (p. 7)

Predict Rating (p. 8)

Recommendation function by predictions (p. 9)

Storing every user recommendation (p. 9)

Code usage (p. 11)

3. Complexity Analysis (p. 12)

Time Complexity

Space Complexity

4. Improving This (p. 13)

Results non parallel (p. 13)

Results parallel (p. 14)

Why Not Theoretical 8x Speedup with 8 Threads? (p. 15)

5. Experimental Evaluation (p.15)

Optimal K size Analysis (p.16)

Execution times analysis (p.17)

Code Implementation Analysis (p.17)

Problem Analysis:

The objective of this programming assignment is to be able to implement a system for recommending items (e.g., movies), using collaborative filtering methods.

This project must be completed *individually* and without the use of LLMs such as ChatGPT.

1 Preliminaries

The input will be a dataset that contains user-movie ratings:

```
user_id, movie_id, rating, timestamp
```

The dataset is available at https://phparis.net/uploads/m2_ds_algods_ratings.zip. In the following, we assume that user and movie IDs are 1-indexed, i.e., the second line represents the first rating.

Your task is to implement a command-line program, *in Python*, (**not a Notebook**), which takes the following two parameters:

1. the name of the file containing the user-movie ratings
2. the similarity threshold (optional, to filter out weak recommendations)

The command line syntax is thus the following:

```
./<program_name> <ratings_file> <similarity_threshold>
```

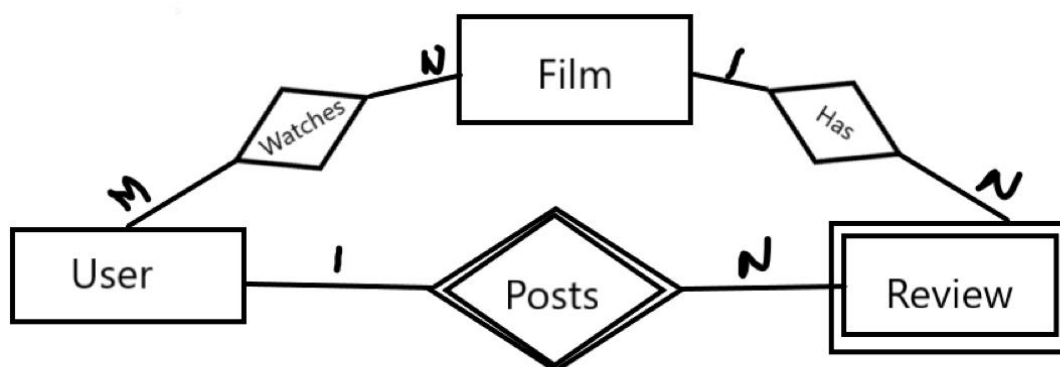
The console output will be the pairs of user IDs and recommended movie IDs, followed by their predicted rating:

```
./<program_name> ratings.csv 0.05
```

will return the following output:

```
0 12 4.5  
1 34 3.9
```

We should implement a regular python file with a collaborative filtering recommendation system for a films streaming service. The retrieved data is stored on https://phparis.net/uploads/m2_ds_algods_ratings.zip with the following structure: user_id, movie_id, rating, timestamp. We can assume this information is obtained from their databases. Seeing the relationship of the data, I think this would be the entity relationship diagram:



Has relationship could be also seen as a weak relationship (review's point of view).

First Idea:

At first I did not read the collaborative filter method approach, so purely based on the course content I thought this implementation.

My idea was to apply association rule using the Apriori algorithm. The goal was to generate association rules of size $k = 4$ that combine the movies a user has already reviewed with any others in the dataset. From the resulting association rules, I would then analyze the individual items; movies that the user has not yet seen. The intuition was that if a movie frequently appears in association rules together with those already reviewed by the user, it is likely to be similar and therefore a good candidate for recommendation. The goal was having a relatively fast online algorithm.

Why collaborative filtering:

Once I realized that the suggested approach was collaborative filtering, I began to wonder why would it be better than mine. The Apriori algorithm might be faster, but it has several disadvantages. Cold start problem for new users, since they have no reviewed films. It also doesn't account for partial or medium similarities; an association either exists or it doesn't. Moreover, it ignores individual user preferences, often resulting in recommendations dominated by the most popular movies, which leads to a less personalized user experience overall. Since collaborative filtering methods have higher complexity, I realized companies such as Netflix, Prime Video probably perform them offline. After doing some research I confirmed my hypothesis, they precompute all the recommendations. That is the reason the recommendations appear almost instantly.

Item-Item and User-User approach:

The first decision we need to make before starting to code is which data structures to use for storing the CSV information. This choice is strongly influenced by whether we adopt an item-item or user-user approach. But what exactly are the differences between these two options?

Item-item collaborative filtering focuses on finding movies with rating patterns similar to those that the user has already reviewed. It then predicts how the user would rate these similar movies, recommending those with the highest predicted scores.

User-user collaborative filtering, on the other hand, compares one user to others with similar rating behaviors. The idea is that if users A and B both liked Movie 1, and A also liked Movie 2, then B is likely to enjoy Movie 2 as well. Although the user-user approach might seem less complex, it presents several downsides.

User-User problems:

The **number of users** is usually much larger than the number of movies, which increases computational cost. Moreover, every time a **user adds a new rating, the entire similarity structure needs to be updated**, whereas in item-item filtering only “movie-movie relations” are affected. User-user filtering also suffers from the cold-start problem for new users. Therefore, its **main advantage** would only appear in early production stages, when the number of movies significantly exceeds the number of users, and even then, only if supported by evidence that it performs better in that specific context. If I had to choose an approach for my starting streaming service, I would implement item-item collaborative filtering due to its scalability. While we might face higher computational costs initially, the system won't require changes as we grow. Additionally, since this process runs offline, we can handle longer processing times during the early stages. Research shows that most companies in the industry follow this approach.

Implementation:

```
def load_ratings(filename):  
    """  
    OBJ: Load ratings from CSV file  
    str -> dict, dict  
    """  
    users = {}  
    movies = {}  
  
    with open(filename, 'r') as f:  
        next(f) # Skip header  
        for line in f:  
            line = line.strip()  
            if not line:  
                continue  
  
            parts = line.split(',')  
            user_id = int(parts[0])  
            movie_id = int(parts[1])  
            rating = float(parts[2])  
            # timestamp = int(parts[3]) # Not needed now, for UI perhaps  
            if user_id not in users:  
                users[user_id] = {}  
            if movie_id not in movies:  
                movies[movie_id] = {}  
            users[user_id][movie_id] = rating  
            movies[movie_id][user_id] = rating  
  
    return users, movies
```

Data Structures Decision:

I chose a **dual dictionary approach** with `users[user_id][movie_id] = rating` and `movies[movie_id][user_id] = rating` because it provides $O(1)$ lookup time for both user-centric and item-centric queries. This structure gives efficient similarity computation in item-item collaborative filtering, as we frequently need to access all users who rated a specific movie and all movies rated by a specific user (for predictions). The two-index structures make both the offline similarity computation and online prediction phases highly efficient, while the space complexity is similar to the other options.

An alternative approach I considered was storing **statistics** per movie (median of average ratings, similar to the Flajolet-Martin algorithm) instead of individual user-rating pairs. While this would dramatically reduce memory usage, storing one value per movie instead of thousands of user ratings, it fundamentally breaks collaborative filtering. We need individual ratings to compute meaningful similarities between items based on user preference patterns. Aggregate statistics would only tell us which movies are generally good, not which movies appeal to similar user tastes. Moreover, we will need to recompute after each review the associated rating.

I also evaluated a **class-based object-oriented approach** with `User` (id, list of reviews), `Film` (average rating), and `Review` (user, film, score) classes. While this structure is more intuitive, legible (perhaps better for avoiding race conditions), it performs poorly for our access patterns. Retrieving all ratings for a specific movie would require iterating through all users review lists or maintaining separate index structures anyway. The dictionary approach provides the same functionality with better performance and simpler implementation, making the memory overhead of storing ratings twice a worthwhile trade-off.

Final Decision: Dual dictionary approach.

```
def cosine_similarity(movie1_ratings, movie2_ratings):
    common_users = set(movie1_ratings.keys()) & set(movie2_ratings.keys())
    if len(common_users) == 0:
        return 0.0
    dot_product = 0.0
    magnitude1 = 0.0
    magnitude2 = 0.0
    for user in common_users:
        r1 = movie1_ratings[user]; r2 = movie2_ratings[user]
        dot_product += r1 * r2
        magnitude1 += r1**2; magnitude2 += r2**2
    magnitude1 = math.sqrt(magnitude1)
    magnitude2 = math.sqrt(magnitude2)
    if magnitude1 == 0 or magnitude2 == 0: return 0.0

    return dot_product / (magnitude1 * magnitude2)
```

Cosine Similarity:

The `cosine_similarity` function computes the similarity between two movies based on their rating vectors. It first identifies common users who rated both movies using set intersection, then calculates the cosine of the angle between the two rating vectors. The numerator is the dot product of the ratings (sum of $r_1 * r_2$), while the denominator is the product of their magnitudes (Euclidean norms). This produces a similarity score between -1 and 1, where values closer to 1 indicate highly similar movies. The function returns 0.0 if there are no common users or if either movie has zero magnitude, preventing division by zero and handling edge cases where meaningful comparison isn't possible.

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

I also considered using a **simple rating difference approach** where movies would be similar if their ratings differed by less than a threshold (e.g., $|\text{rating}_1 - \text{rating}_2| < 0.5$). However, this method fails to account for different user rating behaviors; some users are tough critics (3-4 stars) while others are nicer (4-5 stars) for the same quality. **Cosine similarity** is better because it focuses on the **pattern of who rates what**, not the absolute values. It measures whether the same types of users like both movies regardless of their rating scale, which is exactly what we need for collaborative filtering. In other words, we take into consideration each person reviewing behaviours for recommending the best experiences, unlike simple rating.

Compute movies similarities:

The `compute_movie_similarities` function performs the item-item collaborative filtering by computing pairwise similarities between all movies. It iterates through all movie pairs exactly once using nested loops where the inner loop starts at $i+1$, resulting in $\mathbf{C(k,2) = k(k-1)/2}$ comparisons for k movies; this avoids redundant calculations since similarity is symmetric. For each pair, it computes the cosine similarity and stores it (both $\text{movie}_1 \rightarrow \text{movie}_2$ and $\text{movie}_2 \rightarrow \text{movie}_1$) if the similarity exceeds the threshold, allowing $O(1)$ lookup later. Finally, it sorts each movie's similarity list in descending order to enable efficient top- k . Since this runs **offline**, we can tolerate the $O(k^2)$ computational complexity, as it only needs to be executed periodically (e.g., nightly per country) rather than in real-time during user requests.

```

def compute_movie_similarities(movies, similarity_threshold=0.0):
    """
    OBJ: Compute similarities between all movie pairs
    dict, float -> dict
    """
    movie_ids = list(movies.keys())
    similarities = {}

    print(f"Computing similarities for {len(movie_ids)} movies...")

    for i in range(len(movie_ids)):
        movie1 = movie_ids[i]
        if i % 100 == 0:
            print(f"Processed {i}/{len(movie_ids)} movies")

        if movie1 not in similarities:
            similarities[movie1] = []

        for movie2 in movie_ids[i+1:]:
            sim = cosine_similarity(movies[movie1], movies[movie2])

            if sim >= similarity_threshold:
                similarities[movie1].append((movie2, sim))

                if movie2 not in similarities:
                    similarities[movie2] = []
                    similarities[movie2].append((movie1, sim))

        # Sort each movie's similar movies by similarity (descending)
        for movie_id in similarities:
            similarities[movie_id].sort(key=lambda x: x[1], reverse=True)

    print("Similarity computation complete!")
    return similarities

```

Predict Rating:

The `predict_rating` function estimates what rating a user would give to an unseen movie using item-item collaborative filtering. It retrieves the *k* most similar movies to the target movie (already precomputed and sorted), then filters for only those the user has actually rated. The prediction is calculated as a weighted average: each similar movie's rating is weighted by its similarity score to the target movie, then divided by the sum of all similarity weights. This means movies that are more similar to the target have greater influence on the predicted rating. If no similar movies have been rated by the user or the movie has no computed similarities, the function returns `None`, indicating we cannot make a reliable prediction for this user-movie.


```
def predict_rating(user_id, movie_id, users, similarities, k=10):
    """
    OBJ: Predict rating for a user-movie pair using Item-Item
    collaborative filtering
    int, int, dict, dict, int -> float or None
    """
    if movie_id not in similarities:
        return None
    user_ratings = users[user_id] #Movies that were rated by the user
    # Get k most similar movies that the user has rated
    similar_movies = similarities[movie_id][:k]

    weighted_sum = 0.0
    similarity_sum = 0.0

    for movie, sim in similar_movies:
        if movie in user_ratings:
            weighted_sum += sim * user_ratings[movie]
            similarity_sum += abs(sim)
    if similarity_sum == 0: return None
    return weighted_sum / similarity_sum
```

Recommendation function by predictions:

The `recommend_for_user` function generates personalized movie recommendations for a specific user. We analyze every movie that the user has not rated: `all_movies - user_ratings`. Then, it calls `predict_rating` for each unrated movie to estimate how much the user would like it based on their rating history and movie similarities. Finally, it sorts all predictions by predicted rating in descending order and returns the top-N highest-rated movies as recommendations. This gives us the `n_top` movies most likely liked movies by the user. We will call this function for every user in the CSV to have our precomputed recommendations.

Storing every user recommendation:

`Generate_recommendations` method produces recommendations for all users in the dataset by iterating through each user and calling `recommend_for_user` with top n=1 to get only their single best recommendation. It collects each user's top-predicted movie along with the predicted rating into a list of tuples (`user_id`, `movie_id`, `pred_rating`). This is part of the **offline processing** approach for item-item collaborative filtering; we precompute recommendations for all users periodically rather than in real-time, which aligns with the scalability advantage of item-item filtering I mentioned earlier. Since this runs offline, we can tolerate the computational cost of iterating through all users and predicting ratings for their unrated movies. Progress is printed every 100 users to monitor the processing.

```

def recommend_for_user(user_id, users, movies, similarities, k=10,
top_n=5):
    """
    OBJ: Generate top-N recommendations for a single user
    int, dict, dict, dict, int, int -> list
    """
    if user_id not in users:
        print(f"User {user_id} not found in dataset")
        return []

    user_ratings = users[user_id]
    all_movie_ids = set(movies.keys())

    # Find movies the user hasn't rated
    unrated_movies = all_movie_ids - set(user_ratings.keys())

    # Predict ratings for unrated movies
    predictions = []
    for movie_id in unrated_movies:
        pred = predict_rating(user_id, movie_id, users, similarities, k)
        if pred is not None:
            predictions.append((movie_id, pred))
    # Sort by predicted rating and take top N
    predictions.sort(key=lambda x: x[1], reverse=True)
    return predictions[:top_n]

```

```

def generate_recommendations(users, movies, similarities,
similarity_threshold, k=10):
    """
    OBJ: Generate recommendations for all users
    dict, dict, dict, float, int -> list
    """
    recommendations = []
    print(f"Generating recommendations for {len(users)} users...")
    user_items = list(users.items())
    for i in range(len(user_items)):
        user_id, user_ratings = user_items[i]
        if i % 100 == 0:
            print(f"Processed {i}/{len(users)} users")
        # Use recommend_for_user but only take the best recommendation
        user_recs = recommend_for_user(user_id, users, movies,
similarities, k, top_n=1)
        if len(user_recs) != 0:
            movie_id, pred_rating = user_recs[0]
            recommendations.append((user_id, movie_id, pred_rating))
    print("Recommendations complete!")
    return recommendations

```

Code usage:

The main execution block parses command-line arguments for the ratings file, optional similarity threshold, and optional specific user ID. Loading the data into the dual dictionary structure, performs the **offline similarity computation** for all movie pairs, then generates recommendations. If a specific user ID is provided, it shows top-5 recommendations for that user only; otherwise, it displays the top 1 recommendation per user. Works as a instant called main class.

```
if len(sys.argv) < 2:
    print("Usage: ./collab_filter.py <ratings_file>
[similarity_threshold] [user_id]*")
    print("*If user_id is provided, shows recommendations only for that
user")
    sys.exit(1)

ratings_file = sys.argv[1]
similarity_threshold = float(sys.argv[2]) if len(sys.argv) > 2 else 0.0
single_user_id = int(sys.argv[3]) if len(sys.argv) > 3 else None

# Load data
print("Loading ratings...")
users, movies = load_ratings(ratings_file)
print(f"Loaded {len(users)} users and {len(movies)} movies")

# Compute similarities
similarities = compute_movie_similarities(movies, similarity_threshold)

# Generate recommendations
if single_user_id is not None:
    # Single user mode
    print(f"\nTop recommendations for user {single_user_id}:")
    recommendations = recommend_for_user(single_user_id, users, movies,
similarities)
    for movie_id, rating in recommendations:
        print(f"{single_user_id} {movie_id} {rating:.1f}")
else:
    # ALL users mode
    recommendations = generate_recommendations(users, movies,
similarities, similarity_threshold)
    for user_id, movie_id, rating in recommendations:
        print(f"{user_id} {movie_id} {rating:.1f}")
```

Complexity Analysis:

Time Complexity:

load_ratings: $O(r)$ where r = number of ratings

cosine_similarity: $O(u)$ where u = number of common users between two movies

compute_movie_similarities: $O(m^2 \times u_{avg})$: m = mvs, u_{avg} = average users/ movie

predict_rating: $O(k)$ where k = number of neighbors considered

recommend_for_user: $O(m \times k)$ for predicting all unrated movies

generate_recommendations: $O(n \times m \times k)$ where n = number of users

Total Time: $O(m^2 \times u_{avg} + n \times m \times k)$

Dominated by offline similarity computation: $O(m^2 \times u_{avg})$ and recommendations: $O(n \times m \times k)$.

Space Complexity:

users dictionary: $O(r)$ storing all ratings

movies dictionary: $O(r)$ storing all ratings (duplicate)

similarities dictionary: $O(m^2 \times s)$ where s = average similar movies per movie

Temporary prediction lists: $O(m)$ per user

Total Space: $O(r + m^2 \times s)$ - dominated by the dual rating storage $O(r)$ and precomputed similarity matrix $O(m^2 \times s)$

Conclusions:

The $O(m^2 \times u_{avg})$ time complexity for similarity computation and $O(m^2 \times s)$ space complexity for storing the similarity matrix are precisely why this approach must run **offline**. Computing all pairwise movie similarities is quadratic in the number of movies, which would be incredibly expensive for real-time requests, for example, with 10,000 movies, we'd need 50 million comparisons. However, since movies are added "infrequently" compared to user interactions, we can precompute these similarities tactically and store them in memory. Once precomputed, the online prediction phase is very efficient at $O(k)$.

Improving this:

As I was working on this system, I realized it was perfect for **parallelization** because computing similarities between movies is completely independent; each movie pair can be compared without depending on other calculations. We can distribute the workload across $n-1$ cores (reserving one for system overhead) and process multiple comparisons simultaneously, dramatically reducing precomputation time. For example, with 8 cores we could speed up the process almost 7-fold, transforming a task that would take hours into one completed in minutes. This makes it viable to update similarities more frequently and keep recommendations current without compromising system performance. I decided to implement this new idea in Java (as it is where I have worked in this area) and migrate it to Python by using the multiprocessing lib and a Pool with as many worker threads as possible in parallel. The results were sensational.

Results non parallel:

Loading ratings...

Loaded 6040 users and 3675 movies in 0.86 seconds

Computing similarities for 3675 movies...

Processed 0/3675 movies

Processed 3600/3675 movies

Similarity computation complete!

Similarity computation took 149.23 seconds

Generating recommendations...

Top recommendations for user 5412:

5412 33 5.0

5412 37 5.0

5412 55 5.0

5412 78 5.0

5412 83 5.0

Recommendation generation took 0.01 seconds

=====

TOTAL EXECUTION TIME: 150.11 seconds

=====

Results parallel:

Loading ratings...

Loaded 6040 users and 3675 movies in 0.81 seconds

Computing similarities...

Computing similarities for 3675 movies using 8 threads...

Similarity computation complete!

Similarity computation took 57.86 seconds

Generating recommendations...

Top recommendations for user 5412:

5412 33 5.0

5412 37 5.0

5412 55 5.0

5412 78 5.0

5412 83 5.0

Recommendation generation took 0.01 seconds

=====

TOTAL EXECUTION TIME: 58.69 seconds

=====

DIFFERENCE EXECUTION TIME: 92.25 seconds

TIME REDUCTION PERCENTAGE: 61.45% faster

SPEEDUP: 2.56x

=====

Why Not Theoretical 8x Speedup with 8 Threads?

The parallel version achieves only 2.58x speedup instead of 8x due to unbalanced workload distribution; movie 0 compares with 3,674 others while movie 3,600 compares with only 75, leaving some threads idle while others work. Additional overhead comes from process communication, data serialization between threads, and sequential result merging. Finally, Data Loading and recommendation cannot be parallelized. Although it can be improved, it displays the beneficial impact of parallelism in the precomputing phase. Using AI and the better Load Balancing and Cache Ideas I have managed to improve it to 32 seconds execution time and on rerun 1 second (it is true that no changes were done to the file), this code is in `AIEnhancedParallelRecommendation.py`.

Experimental Evaluation:

[Code on the Experimental.py on the root folder.](#)

=====

EXPERIMENTAL EVALUATION SUMMARY

=====

Optimal k value: 5

RMSE at k=5: 1.5055

MAE at k=5: 1.1305

Accuracy Metrics for Different k Values:

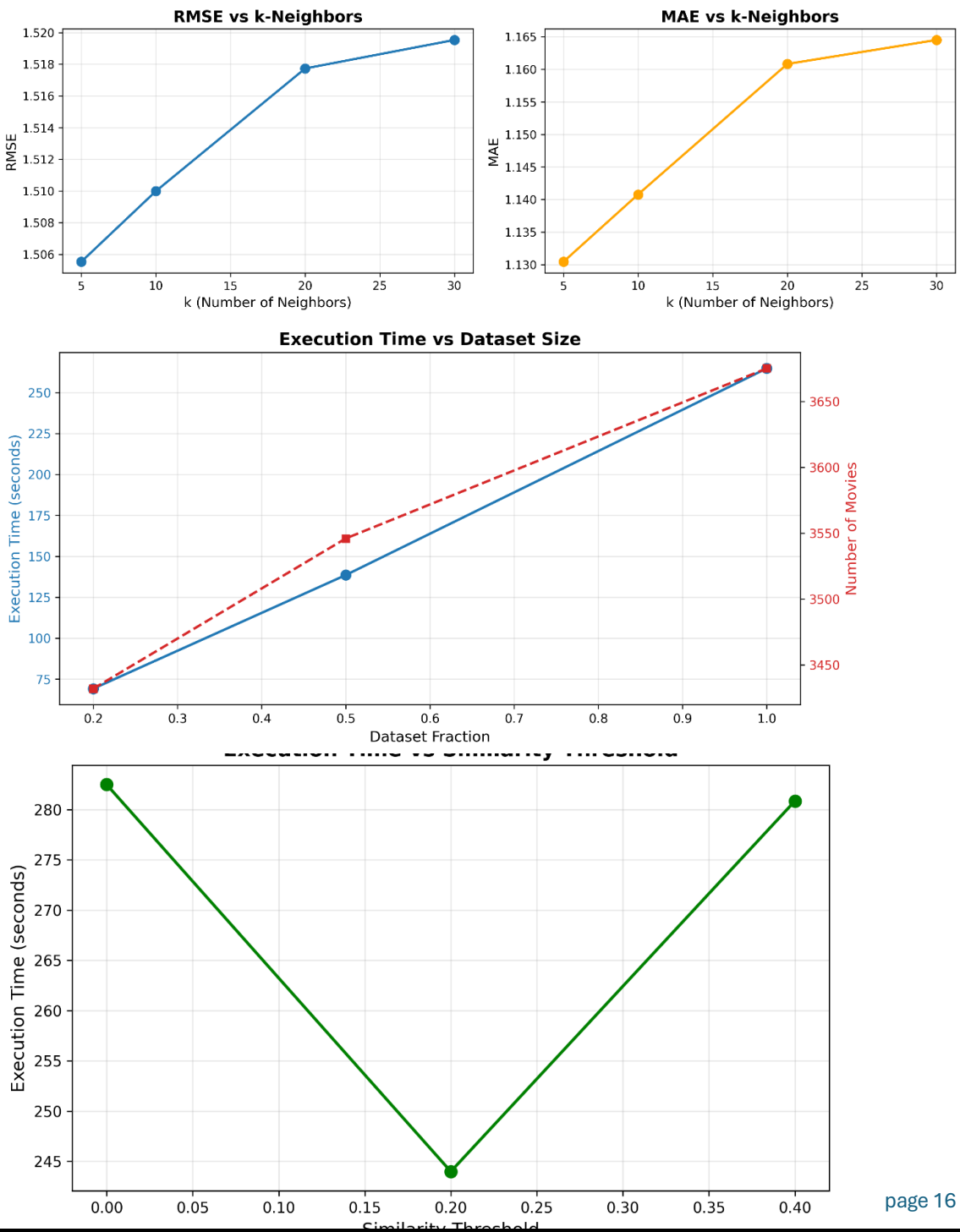
k	RMSE	MAE

5	1.5055	1.1305
10	1.5100	1.1408
20	1.5177	1.1608
30	1.5195	1.1645

=====

K Result Analysis:

The experimental evaluation reveals that **k=5 neighbors provides the optimal performance** with the lowest error metrics (RMSE=1.5055, MAE=1.1305). As k increases from 5 to 30, both RMSE and MAE gradually deteriorate, indicating that including more neighbors introduces noise rather than improving prediction accuracy. This degradation occurs because higher k values incorporate less similar movies into the weighted average, diluting the influence of truly relevant recommendations. The relatively small difference between k=5 and k=30 (RMSE increases only approximately 1.3%) demonstrates the algorithm's robustness.



Analysis of Execution Time Results

The **Execution Time vs Dataset Size** plot demonstrates the expected **quadratic scaling behavior** ($O(m^2)$) we talked about in the time complexity section, of the similarity computation phase, where execution time grows from 70 seconds at 20% of the data to aprox. 265 seconds at 100%, closely tracking the near-linear increase in movie count (3,440 to 3,675 movies). This confirms our complexity analysis and validates why offline precomputation is essential for production systems.

The **Execution Time vs Similarity Threshold** plot reveals an interesting U-shaped curve, with optimal performance at threshold=0.2 (~244 seconds). At threshold=0.0, the algorithm must store and sort all movie pairs regardless of similarity, increasing memory overhead and processing time. At high thresholds (0.4), fewer pairs pass the filter but the overhead of checking each pair remains constant, while subsequent sorting operations become more expensive due to data structure management. This suggests that a moderate threshold around 0.2 provides the best balance.

Code Implementation: simulate scenarios to validate our shared conclusions

1. **Data Splitting:** Splits user ratings into training (80%) and test (20%) sets to evaluate prediction accuracy
2. **Three Key Experiments:**
 - **Accuracy vs k-neighbors:** Tests how the number of similar movies (k) affects prediction accuracy (RMSE/MAE)
 - **Time vs Similarity Threshold:** Measures how filtering out weak similarities affects computation time
 - **Time vs Dataset Size:** Evaluates scalability by testing execution time on different fractions of the data (20%, 50%, 100%)
3. **Metrics Computed:**
 - **RMSE** (Root Mean Squared Error): Average prediction error
 - **MAE** (Mean Absolute Error): Average absolute prediction error
 - **Execution Time:** How long similarity computation takes
4. **Output:**
 - Generates 3 PNG plots visualizing the results (as suggested)
 - Prints a summary table showing optimal k value and accuracy metrics

Our objective: Ensuring all of our conclusions where correct.

