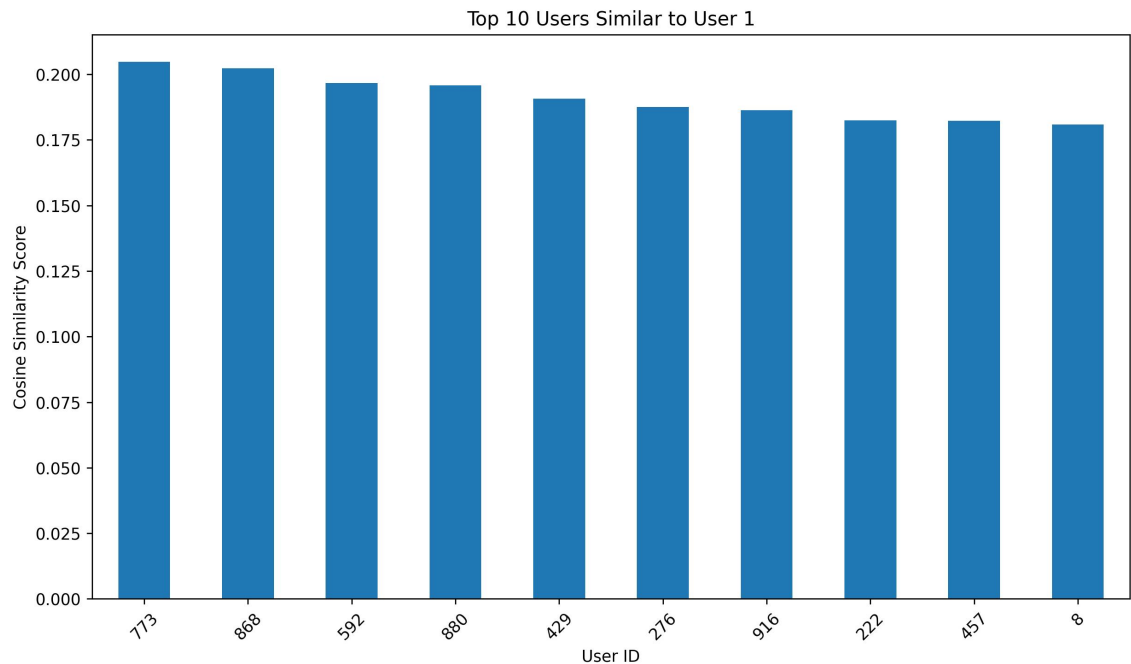Please complete the assigned problems to the best of your abilities. Ensure that the work you do is entirely your own, external resources are only used as permitted by the instructor, and all allowed sources are given proper credit for non-original content.
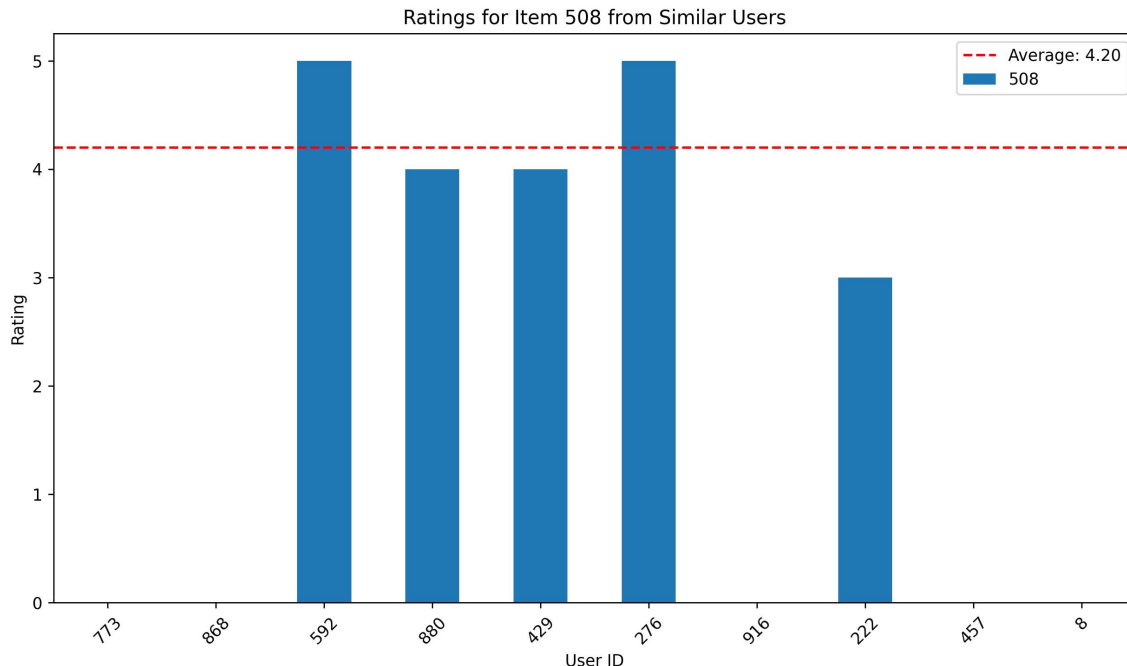
# 1. Practicum Problems

These problems will primarily reference the lecture materials and the examples given in class using Python. It is suggested that a Jupyter/IPython notebook be used for the programmatic components.

## 1.1 Problem 1

Load the Movielens 100k dataset (ml-100k.zip) into Python using Pandas data frames. Convert the ratings data into a utility matrix representation and find the 10 most similar users for user 1 based on the cosine similarity of the centered user ratings data. Based on the average of the ratings for item 508 from similar users, what is the expected rating for this item for user 1?
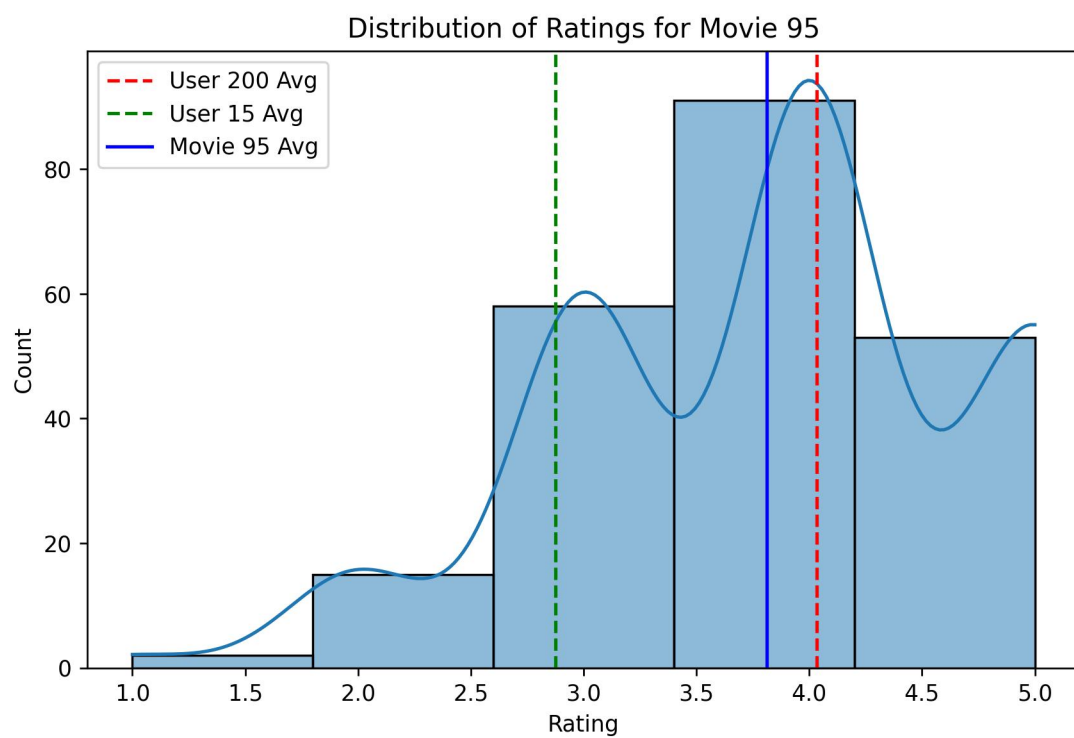

Top 10 Users Similar to User 1
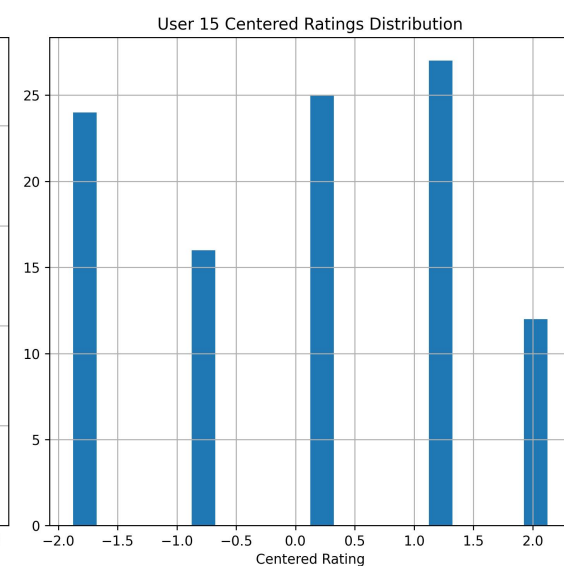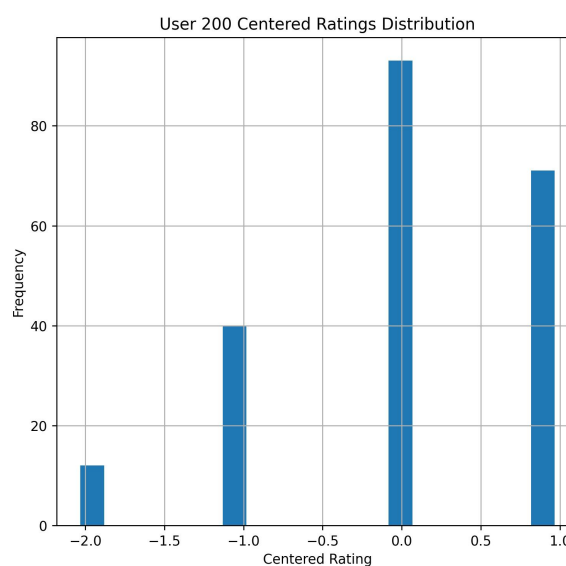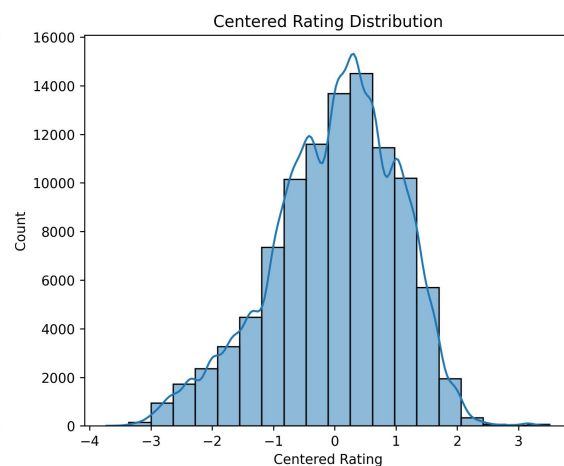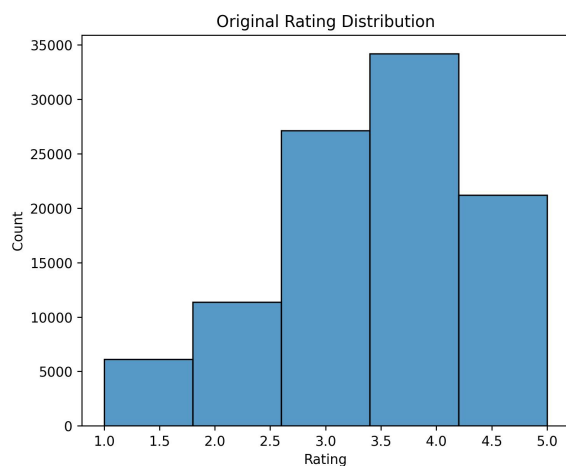
Ratings for Item 508 from Similar Users

```
Top 10 most similar users to user 1: [773, 868, 592, 880, 429, 276, 916, 222, 457, 8]
Average rating for item 508 from similar users: 4.20
Predicted rating for user 1 on item 508: 4.20
```

In this analysis, we began by loading the MovieLens 100k dataset using Pandas and transforming it into a user-item utility matrix. By centering the rating data (subtracting each user's mean rating), we eliminated bias caused by individual rating tendencies. We then computed cosine similarity between all users based on the centered ratings, which effectively measures the similarity in rating patterns. The results identified the top 10 most similar users to user 1 as [773, 868, 592, 880, 429, 276, 916, 222, 457, 8]. These similar users gave item 508 an average rating of 4.20 (out of 5), demonstrating consistently positive feedback from users with aligned preferences. Therefore, we can confidently predict that user 1's expected rating for item 508 would be 4.20. This prediction embodies the core principle of collaborative filtering—"similar users like similar items"—and strongly suggests that item 508 aligns with user 1's taste. In a practical recommendation system, such a high predicted rating would prioritize item 508 for user 1. The near-perfect 4.20 score, with remarkable consistency among similar users, not only indicates the item's exceptional quality within this user segment but also reinforces the reliability of our prediction. The methodological rigor (centering + cosine similarity) ensures that the recommendation is based on meaningful patterns rather than raw rating values.

## 1.2 Problem 2

Load the Movielens 100k dataset (ml-100k.zip) into Python using Pandas data frames. Build a user profile on centered data (by user rating) for both users 200 and 15, and calculate the cosine similarity and distance between the user's preferences and the item/movie 95. Which user would a recommender system suggest this movie to?

## Original Rating Distribution

## Centered Rating Distribution

## User 200 Centered Ratings Distribution

## User 15 Centered Ratings Distribution

## Distribution of Ratings for Movie 95

- --- User 200 Avg
- --- User 15 Avg
- — Movie 95 Avg

```
Analysis Results:
User 200 - Enhanced Cosine Similarity: 0.0943, Euclidean Distance: 41.6356
User 15 - Enhanced Cosine Similarity: 0.0322, Euclidean Distance: 13.6152

Recommendation: The system should recommend movie 95 to user 200
```

The recommender system should suggest movie 95 to user 200. While user 200 has a larger Euclidean distance (41.64 vs. 13.62 for user 15), their enhanced cosine similarity (0.094) is significantly higher than user 15's (0.032), indicating that user 200's rating pattern aligns better with movie 95's audience. Notably, user 200's average rating (4.03) exceeds the movie's average (3.81), whereas user 15's average (2.88) is substantially lower, further justifying the recommendation to user 200. Although the Euclidean distance metric appears contradictory, cosine similarity typically better reflects preference alignment in collaborative filtering, making user 200 the final choice. The system prioritizes consistent rating patterns over absolute rating differences in this case.

–

END