

Data Analytics 1

Assignment 3

Association Rule Mining

Release: 9 September 2024
Deadline: 27 September 2024 (11:55 pm)

The objective of the assignment is to build an association rule-based movie recommender system.

Dataset Description

Provided is a dataset comprising 100,836 ratings and 3,683 tag applications across 9,742 movies. This dataset captures user ratings, with a 5-star rating system, from MovieLens, a movie recommendation service. For additional information, please consult the README.

Data Files

- **ratings.csv**: Contains user ratings. Each line includes userId, movie, rating, and timestamp.
- **movies.csv**: Contains movie information. Each line includes movieId, title, and genre.

Assignment Tasks

Data Preprocessing

1. Form the transactional data set, which consists of entries of the form $\langle \text{user id}, \{\text{movies rated above } 2\} \rangle$. Consider only those users who have rated more than 10 movies.
2. Divide the data set into 80% training set and 20% test set. Remove 20% of movies watched from each user and create a test set using the removed movies.

Association Rule Mining

1. From the training set, extract the set of all association rules of form $X \rightarrow Y$, where X contains a single movie and Y contains the set of movies from the training set by employing the apriori or FPGrowth approach. Set some minsup and minconf (e.g., 50 and 0.1 respectively).
2. **Recommendation**: Generate two sets of lists:
 - The initial list includes the top 100 association rules, arranged in order of their support.
 - The second list comprises the top 100 rules, prioritizing them according to confidence.

Identify the rules that appear in both lists, and then arrange these shared rules based on their confidence score.

3. For each user in the test set, select association rules of the form $X \rightarrow Y$, where X is the movie in the training set. Compute the average precision and average recall by varying the number of rules from 1 to 10 and plot the graphs.
4. Take a sample example of users and their movie ratings from the test set and display precision and recall graphs.
5. Include the plots in the Report (md or pdf) along with your justification of selection of your algorithm and also briefly explain how you built the entire recommendation system in the report.

Submission Format

- Mention any inference from plots or results in notebook Markdown itself.
- Submit a zip folder named `<assignment3_teamId>` containing the following files (you can include any other implementation files as well):
 - `<TeamId>.report` (.pdf or .md)
 - `<TeamId>.recommender` (.py or etc)
 - `<TeamId>_top100RulesByConf` (.txt)
 - `<TeamId>_top100RulesBySup` (.txt)