

CSE 255, Winter 2015: Assignment 2

Instructions

In this assignment you will build a **recommender system** to make predictions related to reviews of Electronics products on *Amazon*.

Please submit your solution electronically to Dongcai Shen (doshen@cs.ucsd.edu) on or before March 10 (the day after the final lecture of the quarter). If you did *not* participate on Kaggle please attach your prediction files for each of the three tasks in the specified format. Your grades will be determined by your performance on the predictive tasks as well as a brief written report about the approaches you took.

This assignment should be completed **individually**.

To begin, download the files for this assignment from:

<http://jmcauley.ucsd.edu/cse255/data/assignment2.tar.gz>

Files

train.json.gz 1,000,000 reviews to be used for training. It is not necessary to use *all* reviews for training, for example if doing so proves too computationally intensive. While these files are one-json-per-line (much as we have seen so far in class), you may find it useful to represent them more concisely in order to produce a more efficient solution. The fields in this file are:

itemID The ID of the item. This is a hashed product identifier from Amazon.

reviewerID The ID of the reviewer. This is a hashed user identifier from Amazon.

helpful Helpfulness votes for the review. This has two subfields, 'nHelpful' and 'outOf'. The latter is the total number of votes this review received, the former is the number of those that considered the review to be helpful.

reviewText The text of the review. It should be possible to successfully complete this assignment *without* making use of the review data, though an effective solution to the 'helpfulness prediction' task will presumably make use of it.

summary Summary of the review.

unixReviewTime Time of the review in seconds since 1970.

reviewTime Plain-text representation of the review time.

category Category labels of the product being reviewed.

pairs_Rating.txt Pairs (userIDs and itemIDs) on which you are to predict ratings (see the tasks below).

pairs_Purchase.txt Pairs on which you are to predict whether a user purchased an item or not.

pairs_Helpful.txt Pairs on which you are to predict helpfulness votes. A third column in this file is the total number of votes, from which you should predict how many were helpful.

helpful.json.gz The review data associated with the helpfulness prediction *test* set. The 'nHelpful' field has been removed from this data, since that is the value you need to predict above. This data will only be of use for the helpfulness prediction task.

baselines.py A simple baseline for each task, described below.

Please do not try to crawl these products from Amazon, or to reverse-engineer the hashing function I used to anonymize the data. I assure you that doing so will not be easier than successfully completing the assignment.

Tasks

You are expected to complete the following three tasks:

Rating prediction Predict people's star ratings as accurately as possible, for those (user,item) pairs in 'pairs_Rating.txt'. Accuracy will be measured in terms of the (*root*) *mean-squared error* (RMSE).

Purchase prediction Predict given a (user,item) pair from ‘pairs_Purchase.txt’ whether the user purchased the item (really, whether it was one of the items they reviewed). Accuracy will be measured in terms of the *categorization accuracy* (1 minus the Hamming loss). The test set has been constructed such that exactly 50% of the pairs correspond to purchased items and the other 50% do not.

Helpfulness prediction Predict whether a user’s review of an item will be considered helpful. The file ‘pairs_Helpful.txt’ contains (user,item) pairs, with a third column containing the number of votes the user’s review of the item received, you must predict how many of them were helpful. Accuracy will be measured in terms of the total *absolute error*, i.e., you are penalized one according to the difference $|n_{\text{Helpful}} - \text{prediction}|$, where ‘nHelpful’ is the number of helpful votes the review actually received, and ‘prediction’ is your prediction of this quantity.

These three error measures are described on *Kaggle*:

RMSE <https://www.kaggle.com/wiki/RootMeanSquaredError>

Classification accuracy <https://www.kaggle.com/wiki/HammingLoss>

Absolute error <https://www.kaggle.com/wiki/AbsoluteError>

A competition page will be set up on Kaggle to keep track of your results compared to those of other members of the class. The leaderboard will show your results on *half of* the test data, but your ultimate score will depend on your predictions across the *whole* dataset.

Grading and Evaluation

You will be graded on the following aspects. Each of the three tasks is worth 10% of your grade.

- Your written report. This should describe the approaches you took to each of the three tasks. To obtain good performance, you should not need to invent new approaches (though you are more than welcome to!) but rather you will be graded based on your decision to apply reasonable approaches to each of the given tasks (3/30% for each task).
- Your ability to obtain a solution which outperforms the baselines on *the unseen portion of* the test data (3/30% for each task). Obtaining full marks requires a solution which is substantially better (i.e., at least several percent) than baseline performance.
- Your ranking for each of the three tasks compared to other students in the class (2/30% for each task).
- Obtain a solution which outperforms the baselines on *the seen portion of* the test data (i.e., the leaderboard). This is a consolation prize in case you overfit to the leaderboard. (2/30% for each task).

Baselines

Simple baselines have been provided for each of the three tasks. These are included in ‘baselines.py’ among the files above. These three baselines operate as follows:

Rating prediction Return the global average rating, or the user’s average if we have seen them before in the training data.

Purchase prediction Find the most popular products that account for 50% of purchases in the training data. Return ‘1’ whenever such a product is seen at test time, ‘0’ otherwise.

Helpfulness prediction Multiply the number of votes by the global average helpfulness rate, or the user’s rate if we saw this user in the training data.

Running ‘baselines.py’ produces three files containing predicted outputs. Your submission files should have the same format.

Kaggle

We have set up a Kaggle page to help you evaluate your solution (please contact Dongcai if you didn't receive the invitation). The Kaggle pages for each of the three tasks are:

<https://inclass.kaggle.com/c/cse-255-assignment-2-task-1-rating-prediction/>

<https://inclass.kaggle.com/c/cse-255-assignment-2-task-2-purchase-prediction/>

<https://inclass.kaggle.com/c/cse-255-assignment-2-task-3-helpfulness-prediction/>