

# CSE 158, Winter 2017: Assignment 1

## Instructions

In this assignment you will build **recommender systems** to make predictions related to reviews of Clothing, Shoes, and Jewelry on *Amazon*.

Solutions will be graded on Kaggle (see below), with the competition closing at **5pm, Monday February 27** (note that the time reported on the competition webpage is in UTC!).

You will also be graded on a brief report, to be submitted electronically on gradescope by the following day. Your grades will be determined by your performance on the predictive tasks as well as your 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/data/assignment1.tar.gz>

## Files

**train.json.gz** 200,000 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.

**price** Price of the item.

**reviewHash** Hash of the review (essentially a unique identifier for 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\_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.

**pairs\_Category.txt** Pairs (userID and reviewHash) on which you are to predict the category of an item.

**pairs\_Rating.txt** Pairs (userIDs and itemIDs) on which you are to predict ratings. **Not relevant for CSE158**

**test\_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.

**test\_Category.json.gz** The review data associated with the category prediction *test* set. Again, the field that you are trying to predict has been removed.

**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 tasks:

**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  $|\text{nHelpful} - \text{prediction}|$ , where ‘nHelpful’ is the number of of helpful votes the review actually received, and ‘prediction’ is your prediction of this quantity.

**Category prediction** Predict the category of an item from a review and product metadata. Five categories are used for this task, namely Women’s, Men’s, Girls’, Boys’, and Baby clothing. Performance will be measured in terms of the fraction of correct classifications.

These error measures are described on *Kaggle*:

**Mean Absolute error** <https://www.kaggle.com/wiki/MeanAbsoluteError>

**Classification accuracy** Is equivalent to 1 minus the Hamming Loss: <https://www.kaggle.com/wiki/HammingLoss>

A competition page has been 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

This assignment is worth 25% of your grade. You will be graded on the following aspects. Each of the two tasks is worth 10 marks (i.e., 10% of your grade), plus 5 marks for the written report.

- Your ability to obtain a solution which outperforms the baselines on *the unseen portion of* the test data (5 marks 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 tasks compared to other students in the class (3 marks 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 mark for each task).

Finally, your written report should describe the approaches you took to each of the 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 (5 marks total).

## Baselines

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

**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.

**Category prediction** Look for a few likely words that may appear in reviews of each category (e.g. if the word ‘baby’ appears, classify as Baby clothing).

Running ‘baselines.py’ produces 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. You should be able to access the competition via your UCSD address. The Kaggle pages for each of the tasks are:

<https://inclass.kaggle.com/c/cse158-258-helpfulness-prediction>

<https://inclass.kaggle.com/c/cse158-categorization>

There is also a rating prediction task being used for the graduate class (CSE258), which you're more than welcome participate in if you're interested: <https://inclass.kaggle.com/c/cse258-rating-prediction>