

Final

DSE 220: Machine Learning

Due Date: 11 June 2017

1 Instructions

The deadline for the kaggle competition is 11 June 2017 (11:59 PM UTC/ 4:59 PM PST). The report for the final should be submitted on Gradescope and the code should be submitted on github by 12 June 2017 (11:59 PM PST). The code will only be evaluated from github so make sure you have correct code on github. To secure full marks both the report and the code should be in sync and logically correct. Please only submit relevant and legible code. **Please complete the final individually.**

2 Overview

For the take-home final of DSE 220: Machine Learning, we combine all the concepts we have learned so far and apply them to a real world problem. The problem, in very broad terms, is to predict how useful an Amazon product review is going to be. The problem is open-ended and you can use any method you like.

Amazon is a leading e-commerce website and uses Machine Learning applications extensively in their products. For example, they recommend products based on your browsing history and previous purchases so as to maximize their profit (by selling you another product!). For your final, you are given Amazon data. Given a review, the item for which the review is and the user who has written the review, you need to predict how helpful the review will be.

Solutions will be graded on Kaggle. Please follow the link -

<https://inclass.kaggle.com/c/dse-220-final>

to view the webpage and signup using your UCSD email id (@ucsd.edu and not @eng.ucsd.edu). **Note that the time reported on Kaggle are in UTC and not PST.**

You will also be graded on a brief report, to be submitted electronically on Gradescope. Your grades will be determined by your performance on the predictive task as well as your written report about the approaches you took.

3 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, you may find it useful to represent them more concisely in order to produce a more efficient solution. The fields in this file are:

1. **itemID** The ID of the item. This is a hashed product identifier from Amazon.
2. **reviewerID** The ID of the reviewer. This is a hashed user identifier from Amazon.
3. **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.
4. **reviewText** The text of the review. It should be possible to successfully complete this task without making use of the review data, though an effective solution to the helpfulness prediction task will presumably make use of it.
5. **summary** Summary of the review.
6. **price** Price of the item.
7. **reviewHash** Hash of the review (essentially a unique identifier for the review).
8. **unixReviewTime** Time of the review, i.e, time elapsed in seconds since 1970. For example - 1400544000
9. **reviewTime** Plain-text representation of the review time in human readable format. For example - '05 20, 2014'
10. **category** Category labels of the product being reviewed.
11. **rating** Rating given by the reviewer.

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.

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.

Baseline.ipynb A simple baseline for each task, described later.

Please do not try to crawl these products from Amazon/any third-party website, to reverse-engineer the hashing function used to anonymize the data, or to copy

someone else's code off the internet. We assure you that doing so will not be easier than successfully completing the final. Since we are asking for your code, we will know whether or not you parsed the data, and we will check for the originality of code. Any unfair means used could end up earning you a zero on your final exam.

4 Task

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 mean absolute error, i.e., you are penalized one according to the **absolute difference of nHelpful and prediction**, where 'nHelpful' is the number of of helpful votes the review actually received, and 'prediction' is your prediction of this quantity.

The error measure is described on Kaggle:

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

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.**

5 Grading and Evaluation

You will be graded on the following aspects.

1. Your ability to obtain a solution which outperforms the baselines on the total test data (including the unseen portion). Obtaining full marks requires a solution which is substantially better than baseline performance. Following are the marks you'll get for beating every benchmark.

Baseline: 0.261

Benchmark I: ≤ 0.190 5 marks

Benchmark II: ≤ 0.185 10 marks

Benchmark III: ≤ 0.179 20 marks

Benchmark IV: ≤ 0.176 25 marks

Benchmark V: ≤ 0.173 30 marks

2. Your ranking compared to other students in the class. The top one-third students will receive the highest marks, followed by middle third and then the final third.

Top third of class: 20 marks

Middle third of class: 10 marks

Bottom third of class: 5 marks

3. Obtain a solution which outperforms the baselines on the seen portion of the test data (i.e., the public leaderboard) to obtain 20 marks. This is a consolation prize in case you overfit to the leaderboard.
4. The report accounts for 30 marks. It should describe the approaches you took to perform the task. Make sure that the methods you describe in the report include all the aspects of your final model including pre-processing, feature engineering etc. The aim is that anyone with your report should be able to recreate your results. Even if your model doesn't perform well, you can obtain marks in this section for the comprehensiveness of your analysis.

You can obtain a maximum of 100 marks in this project, which will be scaled down to 40% of the total course assessment. To obtain good performance, you should not need to invent new approaches (though you are more than welcome to!).

6 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(@ucsd.edu, not @eng.ucsd.edu).

<https://inclass.kaggle.com/c/dse-220-final>

You can submit only ten submissions per day to Kaggle. This is to ensure that you don't learn from the test data. Please ensure that you use a good validation set for measuring the performance of your model. **Tuning your models based on the leaderboard position can lead to overfitting and you might end up losing all the marks based on the private leaderboard.**

You need to select two top submissions on which you want us to evaluate you at the end of the competition. If top two submissions are not selected, we pick the top two submissions according to the public leaderboard.

7 Baselines

Simple baseline solution has been provided for the task. These are included in 'Baseline.ipynb' among the files above. The baseline operates as follows:

Multiply the number of votes by the global average helpfulness ratio (nHelpful/outOf), or the user's average helpfulness ratio if we saw the same user in the training data.

Running the given iPython notebook produces files containing predicted outputs. Your submission files should have the same format. Hints will be provided to you on Piazza by the TAs from time to time depending upon the class performance. Please monitor Piazza posts for updates.