

EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the [Privacy Policy](#).



MITx: 6.86x

Machine Learning with Python-From Linear Models to Deep Learning

[Help](#)[sandipan_dey](#)[Unit 1 Linear Classifiers and
Course > Generalizations \(2 weeks\)](#)[Project 1: Automatic Review
> Analyzer](#)

> 9. Feature Engineering

9. Feature Engineering

Frequently, the way the data is represented can have a significant impact on the performance of a machine learning method. Try to improve the performance of your best classifier by using different features. In this problem, we will practice two simple variants of the bag of words (BoW) representation.

Remove Stop Words

1/1 point (graded)

Try to implement stop words removal in your feature engineering code. Specifically, load the file **stopwords.txt**, remove the words in the file from your dictionary, and use features constructed from the new dictionary to train your model and make predictions.

Compare your result in the **testing** data on Pegasos algorithm using $T = 25$ and $L = 0.01$ when you remove the words in **stopwords.txt** from your dictionary.

Hint: Instead of replacing the feature matrix with zero columns on stop words, you can modify the `bag_of_words` function to prevent adding stopwords to the dictionary

Accuracy on the test set using the original dictionary: 0.8020

Accuracy on the test set using the dictionary with stop words removed:



You have used 1 of 20 attempts

✓ Correct (1/1 point)

Change Binary Features to Counts Features

1/1 point (graded)

Again, use the same learning algorithm and the same feature as the last problem. However, when you compute the feature vector of a word, use its count in each document rather than a binary indicator.

Hint: You are free to modify the `extract_bow_feature_vectors` function to compute counts features.

Accuracy on the test set using the dictionary with stop words removed and counts features:

0.7700



Try to compare your result to the last problem, and see the discussion in solution after answering the question.

Submit

You have used 1 of 20 attempts

✓ Correct (1/1 point)

Some additional features that you might want to explore are:

- Length of the text
- Occurrence of all-cap words (e.g. "AMAZING", "DON'T BUY THIS")
- Word embeddings

Besides adding new features, you can also change the original unigram feature set. For example,

- Threshold the number of times a word should appear in the dataset before adding them to the dictionary. For example, words that occur less than three times across the train dataset could be considered irrelevant and thus can be removed. This lets you reduce the number of columns that are prone to overfitting.

There are also many other things you could change when training your model. Try anything that can help you understand the sentiment of a review. It's worth looking through the dataset and coming up with some features that may help your model. Remember that not all features will actually help so you should experiment with some simpler ones before trying anything too complicated.

Discussion

Hide Discussion

Topic: Unit 1 Linear Classifiers and Generalizations (2 weeks):Project 1: Automatic Review Analyzer / 9. Feature Engineering

Add a Post

Show all posts	by recent activity
<div><div>💬</div><div>Change Binary Features to Counts Features</div><div>I replaced the binary representation in the table for the number of times this word is included in the review. However, as a result, the accuracy became lower. I obtained a...</div></div>	12
<div><div>?</div><div>Could we have more weeks released in advance?</div><div>I see the schedule, but the holiday season is coming.</div></div>	11
<div><div>?</div><div>The total words remain unchanged after loading the stopword data</div><div>Hi, for the first question. I loaded the stopwords by writing: <code>stopword_data = utils.load_data('stopwords.txt')</code> and modified the bag of words function to: <code>def bag_of_words...</code></div></div>	10
<div><div>?</div><div>Try to compare your result to the last problem, and see the discussion in solution after answering the question?</div><div>Where can I see this discussion? I've succesfully completed the problem, however I'd be happy to read the discussion. The problem is, I don't have a "Show Answer" butto...</div></div>	1
<div><div>?</div><div>When using Count features</div><div>Hi, I am using <code>count(word)</code> for every word from word list in every occurence of reviews in <code>extract_bow_feature_vectors</code> function. I am getting output which is almost 76% o...</div></div>	1
<div><div>?</div><div>Troubleshooting the code when answers marked wrong by grader for Stop Words</div><div>I have applied Stop words check criteria on the dictionary and verified that length of dictionary earlier and later got changed. On applying the Pegasos algorithm with its b...</div></div>	1
<div><div>?</div><div>Change binary to count</div><div>By count do you the number of times the word has appeared in the entire "reviews" list?</div></div>	3
<div><div>💬</div><div>Bag of words and stopwords versus negation</div><div>Consider these two possible reviews: - "I really enjoy this coffee. I would buy it again." - "I really don't enjoy this coffee. I would not buy it again." Using single words rips th...</div><div>👤Community TA</div></div>	1
<div><div>💬</div><div>There are periods in amongst the words.</div><div>Punctuation is leaking in. There are periods in the dictionary... Highest count of periods in one review in the training data? 580. Somebody wrote an essay...</div><div>👤Community TA</div></div>	1
<div><div>💬</div><div>It is great stuff!</div><div>thanks</div></div>	1
<div><div>💬</div><div>Unit 2 Release</div><div>Hi Admin, How about the unit 2? It's still not yet opened as the course schedule. Thanks</div></div>	5
<div><div>💬</div><div>parameter initialization</div><div>I am just wondering if there are some better ways for parameter initialization instead of zeros. More specifically for linear SVM, consider binary classification. Will it help if ...</div></div>	3

Learn About Verified Certificates