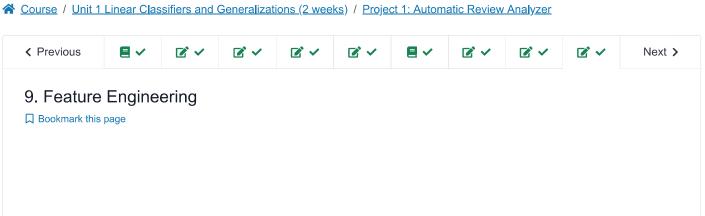
Course <u>Progress</u> <u>Dates</u>

Discussion

Resources



Project due Oct 7, 2020 05:29 IST Completed

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 <code>bag_of_words</code> 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:



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:



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

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

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

Add a Post

Show all posts ✓ by rec	cent activity 🗸
Remove stop words - how to debug? Am I missing something on removing the stop words? When I removed the stop words from the dictionary, my training accuracy we	19 e
I'm still not entirely sure why using count instead of binary representation decreased the accuracy? Even though we removed articles, verbs and pronouns from our bag of words model, why does count decrease accuracy? Is it specentary.	1 ci
? Help with additional Python resources Could someone please recommend some additional sources to acquire further Python skills that are needed for projects like this or	3 <u>In</u>
Poll: Average number of hours to complete Project 1 Out of curiosity, how many hours did it take you to complete Project 1? In my case it took me ~15 hours.	19
? Question about Recitation The value I2 in the parameter **penalty='l2'**, where does it come from?	2
all methods are passing, but returning incorrect validation error All my methods are passing and correct, but the validation error is incorrect. I've also cleared global variables and reran my tests, but the validation error is incorrect.	3 <u>u.</u>
? Dear STAFF For those of us, who didn't manage to complete the section 9. Feature Engineering, would the respective code be available to us?	1
Project 1 is very well engineered Project 1 is very well engineered: Robust code, well organized, very didactic. A model of how questions should be stated. Congratule	20 <u>la</u>
Change Binary Features to Counts Features What is meaning of document in this "However, when you compute the feature vector of a word, use its count in each document ra	8 at
? [To Staff] Problem 8 and 9 I have managed to complete problem 8 and 9 by making slight adjustments to existing codes (e.g. utils.py, project1.py). Are we expense.	2 <u>e</u>
Mathematics for Machine Learning Book Hello For those interested, this books it's a great resource for maths behind machine learning and it has an excellent chapter about Community TA	<u>t</u> 7
Tip: Change Binary to Count Features Using the Counter class from collections package made the task really simple for me. Just initialized the Counter class with the word	4 dl
? Trouble with Removing stop words extract bow features function	5 🔻

Previous	Next >
----------	--------



edX

<u>About</u>

Affiliates

edX for Business

Open edX

Careers

News

Legal

Terms of Service & Honor Code

Privacy Policy

Accessibility Policy

Trademark Policy

<u>Sitemap</u>

Connect

<u>Blog</u>

Contact Us

Help Center

Media Kit

Donate















© 2020 edX Inc. All rights reserved. 深圳市恒宇博科技有限公司 <u>粤ICP备17044299号-2</u>