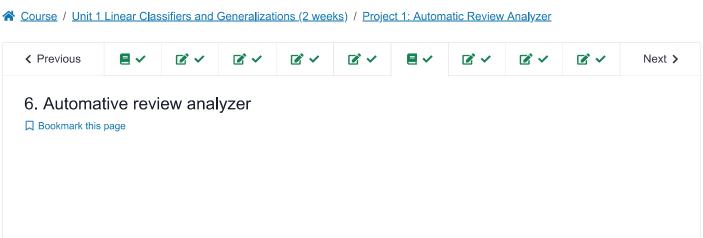
<u>Help</u> smitha_kannur -

Course <u>Progress</u> <u>Dates</u> **Discussion** Resources



Now that you have verified the correctness of your implementations, you are ready to tackle the main task of this project: building a classifier that labels reviews as positive or negative using text-based features and the linear classifiers that you implemented in the previous section!

The Data

The data consists of several reviews, each of which has been labeled with -1 or +1, corresponding to a negative or positive review, respectively. The original data has been split into four files:

- reviews_train.tsv (4000 examples)
- reviews_validation.tsv (500 examples)
- reviews_test.tsv (500 examples)

To get a feel for how the data looks, we suggest first opening the files with a text editor, spreadsheet program, or other scientific software package (like <u>pandas</u>).

Translating reviews to feature vectors

We will convert review texts into feature vectors using a **bag of words** approach. We start by compiling all the words that appear in a training set of reviews into a **dictionary**, thereby producing a list of d unique words.

We can then transform each of the reviews into a feature vector of length d by setting the $i^{\rm th}$ coordinate of the feature vector to 1 if the $i^{\rm th}$ word in the dictionary appears in the review, or 0 otherwise. For instance, consider two simple documents "Mary loves apples" and "Red apples". In this case, the dictionary is the set $\{\text{Mary}; \text{loves}; \text{apples}; \text{red}\}$, and the documents are represented as (1;1;1;0) and (0;0;1;1).

A bag of words model can be easily expanded to include phrases of length m. A **unigram** model is the case for which m=1. In the example, the unigram dictionary would be (Mary; loves; apples; red). In the **bigram** case, m=2, the dictionary is (Mary loves; loves apples; Red apples), and representations for each sample are (1;1;0), (0;0;1). In this section, you will only use the unigram word features. These functions are already implemented for you in the bag of words function.

In utils.py, we have supplied you with the load data function, which can be used to read the .tsv files and returns the labels and texts. We have also supplied you with the bag_of_words function in project1.py, which takes the raw data and returns dictionary of unigram words. The resulting dictionary is an input to extract_bow_feature_vectors which computes a feature matrix of ones and zeros that can be used as the input for the classification algorithms. Using the feature matrix and your implementation of learning algorithms from before, you will be able to compute θ and θ_0 .

Discussion

Hide Discussion

Topic: Unit 1 Linear Classifiers and Generalizations (2 weeks):Project 1: Automatic Review Analyzer / 6. Automative review analyzer

Add a Post

SI	how all posts 💙 by recent activity	ity 🗸
9	Opening the .tsv's in Pandas To open the reviews files in pandas, the following line would do the trick - train = pd.read csv('/reviews test.tsv',sep='\t',encoding='ISO	1
•	When is the recitation due? In the course overview I see "Unit 1 Linear Classifiers and Generalizations (2 weeks) - Recitation 1: Tuning the Regularization Hyperpara	3
7	linear separability of higher dimensional feature vectors. when forming feature vectors like using a bag of words approach as above. how can we be sure they are linearly separable or what if th	1
	- Encoding and positing insure	

Encoding and naming issue
The validation file name is "reviews_val.tsv", not "reviews_validation.tsv". Also, I think the encoding of the file is windows-1252.

NameError: name 'csv' is not defined

Previous

Next >

© All Rights Reserved



edX

About

Affiliates

edX for Business

Open edX

Careers

<u>News</u>

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>

/