

CSE 5243: Homework #3

Deadline: 11:59PM on 10/23/2018.
No late submissions will be accepted.

Instructions.

You will have around 3 weeks to work on this programming assignment. We currently use a 100-point scale for this homework, but it will take 10% of your final grade.

What you should turn in:

Please turn your code and all documentations/reports to Carmen. **Please try to be considerate of the TA and make sure he can run your code with minimal effort. :-)**

Questions?

Please create a post on Carmen or Piazza discussion areas to get timely help from other students, the TA, and the instructor. Students who actively answer others' questions will get extra credit! :-) Besides, everyone can benefit from checking what have been asked previously. Please try to avoid directly sending emails to the instructor/TA. Thanks!

1 Task: Sentiment Classification (100 points)

This assignment is the second part of a longer-term project. The objective is to give you the experience of building classifiers for your preprocessed real data in HW2.

1.1 Review HW2.

What you have done for HW2:

- **Data:** Sentiment Labelled Sentences Data Set, which contains sentences labelled with positive or negative sentiment. It can be downloaded here <http://archive.ics.uci.edu/ml/datasets/Sentiment+Labelled+Sentences>. Read their readme.txt file for detailed information. There are three subsets respectively from IMDB, Amazon and Yelp. Please merge them as a single dataset, which should contain 3,000 sentences in total.

- **Data Format:** Each data file is .txt where each row has two columns: sentence body and sentence label. For example, one sample sentence is “Very little music or anything to speak of. 0”, where the first column “Very little music or anything to speak of.” is the content of a sentence while the second column “0” is its sentiment label (1 means positive; 0 means negative).
- **Results from HW2:** In HW2, you should have constructed a feature vector for each sentence in the data set. Let us assume you use the frequency of unique words in the sentence body to construct a feature vector. If there are totally M sentences and N unique words in the dataset, you should have constructed a $M \times N$ matrix D , where $D_{i,j}$ means the count of word j in sentence i .

Please note that your feature vector for each sentence should only be based on the sentence content, and does not contain the sentiment label.

1.2 Task Description for HW3.

1.2.1 Overall goal.

The goal of HW3 is to build and evaluate classifiers to predict the sentiment of a given sentence. The true label of each sentence is binary (i.e., positive vs negative) and given in the dataset.

1.2.2 Dataset split.

To build a classifier, you need a training, validation, and testing set. Now that you have 3,000 sentences in total, you can randomly sample $p\%$ of them as training, $q\%$ as validation, $r\%$ as testing, where $p+q+r=100$. For example, sample 1/3 as training, 1/3 as validation, 1/3 as testing; or, 60% as training, 20% as testing, 20% as validation, or other percentage you prefer. Please detail how you split your data in your report.

1.2.3 Feature selection.

You will also play around with methods to select “important” features. Let us assume you start with a feature vector that comprises or keeps track of 2048 words. What you will need to do is to identify and leverage a measure of importance to pare down the feature vector size down to, say, 256 or 512 words, and compare and contrast the performance of **the original feature set vs pruned feature set** on two classification algorithms.

You have the freedom to define your own measure of importance. You may consider one of the following two intuitive choices:

- (1) One intuitive measure is the overall frequency of a word in the entire dataset. You can keep the top- K most frequent words (e.g., top 1000 most frequent words in the entire dataset) and prune less frequent ones. If you choose

to use this measure, you may first remove stop words like “the” and “a” from your dataset, as they are not informative but very frequent.

(2) TF-IDF. TF-IDF is the product of two statistics, term frequency and inverse document frequency. The tf-idf value increases proportionally to the number of times a word appears in the document, but is often offset by the frequency of the word in the entire corpus, which helps to adjust for the fact that some words appear more frequently in general. For more information, please refer to <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>.

For word t that occurs in sentence d , you can calculate its tf-idf in the simplest way in the previous link:

- $tf(t, d) = f_{t,d}$, where $f_{t,d}$ is simply the frequency of word t in document d .
- $idf(t, D) = \log(\frac{N}{|d \in D : t \in d|})$, is a measure of whether word t is common or rare across all sentences D . Note in your experiment, you can use your training sentence set as D . It can be obtained by dividing the total number of sentences (i.e., N) by the number of sentences containing the term (i.e., $|d \in D : t \in d|$), and then taking the logarithm of that quotient.
- Finally, the tf-idf score of word t in sentence d , is $tf(t, d) * idf(t, D)$.

Note that you can decide what importance measures to use, how many features should be pruned, etc. Please detail what you used in your report.

1.2.4 What classification algorithms to use?

You are expected to test out **at least two** independent classifiers (K-nearest neighbors, decision trees and Naïve Bayes are possible choices but you may pick others like simple neural networks we talked about in class). You can choose to implement your own version of these classifiers or download and use free software/packages.

Our *strong recommendation* is that you implement at least one classifier by yourself (e.g., K-NN). Keep in mind that if you use existing software you may have to transform your feature vector dataset to match their input specifications. Your report should describe any further data transformations you may have had to work with these software packages.

1.2.5 What results to expect?

To sum up, assuming you pick two classifiers (CL-1; CL-2) and you have the original feature vector from the last assignment (FV-1) (see notes below on things you may need to correct) and a pared down feature vector where you have selected some features (FV-2) as discussed in Section 1.2.3. You will need to run 4 sets of experiments (**each classifier on each type of feature vector** – think of it as a 2X2 experiment configuration matrix).

For each set of experiments (configuration) you will need to report:

- the overall scalability of the resulting approach (time to build classifier model – this is also known as the offline efficiency cost)
- time to classify a new tuple (also known as the online efficiency cost)
- the accuracy of the classifiers respectively on training, validation and testing sets.
- **comments on the advantages (or disadvantages, depending on what you find) using the feature selection methodology you adopted.** Examples could include – feature selection helps improve accuracy and efficiency costs, or feature selection is ineffective. Back up your observations with numbers/graphs/charts.

Notes: Please note that if the TA provides feedback on your preprocessing step and obtained features in HW2, please make sure the feedback is addressed in this assignment.

1.3 What You Should Turn in.

Use the submit command to submit all source files (such as README, source, data, report etc.), as you did for HW2:

- All source code. Please make it easy to run your code.
- A detailed README file that contains all the information about the directory, e.g., how to run your program, how to interpret the resulting output, etc. You can decide what your code outputs to show the TA. Examples would be (1) listing a few sentences, their predicted labels and their true labels; (2) a summary of your classifiers' results, such as their accuracy for each set of experiments.
- **A detailed report, listing all underlying assumptions, and explanations for performance obtained is expected. Prior to your submitting this report, if the TA provides feedback on the previous programming assignment in HW2, please make sure the feedback is addressed in this assignment's report.**

Detailing what you did is very important even if it did not work. Please describe any difficulties you may have encountered and any assumptions you are making.

1.4 Additional Notes

(1) Programming Requirement. You have complete choice on the programming environment as long as it will execute on the STD LINUX environment. Note – you are welcome to implement standalone code (Python suggested) for this

project and/or leverage free software or packages such as NLTK (Natural Language Toolkit) as long as you can concisely and precisely explain how to run the program in your report or README file.

(2) Submit on Carmen (and see instructions in HW2, if needed).

1.5 Acknowledgement

The dataset was originally used in the following paper:

[1] “From Group to Individual Labels using Deep Features”, Kotzias et al., SIGKDD 2015.