

CS159 - Assignment 5

Part 1: Friday November 4 at 6pm

Part 2: Friday, November 11 at 6pm

For our last assignment (yay!) we will be exploring different approaches for calculating word similarities based on distributional similarity. For example, given the word `dog` we might come up with the following similar words:

```
terrier  
breed  
dogs  
inu  
spitz
```

You may work with a partner on this assignment (I encourage you to try and work with a different partner) and you may use whatever programming language you like.

1 Part 1: Warm-up

To get started, we're going to do a similarity calculation "manually" to make sure that you understand the concepts and to give you some test data to check your answers against. You may submit this portion digitally or on paper.

1. Read through the rest of this handout and make sure you understand the problem description.
2. Calculate the similarity between the "words" `b` and `d` with TFIDF weighting using EUCLIDEAN distance for the following data:

```
a a b c c  
a a d c c  
a a b c c  
e e d g g  
e e f g g  
h h h h h
```

Follow the specifications below including length normalization, etc. where appropriate. Along with your answer, turn in intermediate calculations that will help me verify correctness, for example, the co-occurrence vectors for **b** and **d**, the IDF's for each of the words and the length normalized TFIDF co-occurrence vectors.

You should do this part “manually”, which doesn't necessarily mean on paper (e.g. you can use Excel, python or some other tool), but you should not be trying to write the whole program and should instead be doing the individual calculations.

2 Part 2: The program

2.1 Input

Like the last assignment, for this assignment you will be creating a script called `wordsim.sh` that outputs the most similar words to a set of input words based on a some of the different weighting schemes and similarity measures we've looked at in class. This script should take three arguments as parameters:

```
wordsim.sh <stoplist> <sentences> <inputfile>
```

- `<stoplist>`: a list of stop words, one per line, that should be ignored from the input.
- `<sentences>`: a list of sentences/text fragments, one per line, to be used for training the distributional similarity method.
- `<inputfile>`: a file consisting of lines of the following form:

```
<word>    <weighting>    <sim_measure>
```

where `<weighting>` is one of:

- TF: term frequency - use the number of times each word occurs in the word context.
- TFIDF: term frequency with inverse document frequency - use the term frequency times the inverse document frequency. When calculating the IDF, *treat each line as a separate document*.
- PMI: pointwise mutual information - use the pointwise mutual information weighting between the word and the feature. Calculate:

$$p(w, f) = \frac{\text{count}(w, f)}{\text{count}(\text{all_words})}$$

$$p(w) = \frac{\text{count}(w)}{\text{count}(\text{all_words})}$$

$$p(f) = \frac{\text{count}(f)}{\text{count}(\text{all_words})}$$

where $\text{count}(\text{all_words})$ is the total number of word occurrences in the corpus, $\text{count}(f)$ and $\text{count}(w)$ are both over the whole corpus (i.e. not only in word contexts) and $\text{count}(w, f)$ is the number of times word f was seen in the context of word w , counting duplicates (i.e. if f occurred twice in one context, you'd count it twice).

and `<sim_measure>` is one of:

- L1: L1 distance, normalized by the L2 (Euclidean) length of the vectors.
- EUCLIDEAN: Euclidean distance, normalized by the L2 (Euclidean) length of the vectors.
- COSINE: Cosine distance, normalized by the L2 (Euclidean) length of the vectors.

2.2 Output

For any run of the program, the first thing you should output are some statistics: the number of unique words (after stoplist removal, lowercasing, etc.), the number of word occurrences (again, after preprocessing) and the number of sentences/lines.

For each line in the `<input_file>` you will calculate the 10 most similar words based on the weighting and similarity measure specified (or less, if there are less than 10 words in the data set). Output these in order based on most similar (note that, one per line with the similarity score, tab separated. For example:

```
SIM: dog TF COSINE
terrier  0.6757599580387642
now      0.6732906004244482
kind     0.6564829111760163
called   0.6561324054538721
considered 0.6526006142705167
today    0.6505379660187838
type     0.6481430531459309
possible 0.6470210473696648
used     0.646499660989564
also     0.6403788033560914
```

2.3 Implementation details

- Split words based on whitespace.
- Lowercase all words
- Before calculating any statistics, contexts, etc. remove all stop words.

- Before calculating any statistics remove all words that are not exclusively letters (i.e. no numbers, punctuation, etc.).
- We will use a context of 2 words on each side of a word. If a word occurs multiple times in a context, count it multiple times. If you are at a line boundary (beginning or end) only count to the boundary. For example, if we had the input sentence:

`I like to eat bananas too .`

and “to” was in our stoplist, then we would preprocess and get:

`i like eat bananas too`

and if we then wanted the context for `bananas` we would get the words `{like eat too}`.

- For PMI the log should be base 10.
- When looking for the 10 most similar words, only consider words that have occurred 3 times or more.
- Make sure to follow the input and output specifications exactly. Look at the test input and output examples in the data directory and make sure yours match formatting, etc.
- My version takes less than a minute to load the data and calculate the test example. As with all of the assignments, do think about efficient data structures, etc. when you’re thinking about the implementation.

2.4 Hints

- You’ve walked through one example by hand, but I’d encourage you to look at some of the other ones manually as well.
- Be careful about converting between `ints` and `doubles`.
- For each of the similarity/distance measures make sure you think about whether smaller or larger means “more similar”.
- There is a strong relationship between `EUCLIDEAN` and `COSINE`.

2.5 Data

I have provided you with some data to get you started at:

<http://www.cs.middlebury.edu/~dkauchak/classes/cs457/assignments/assignment5/assign5-starter.tar.gz>

which includes the following files in the `data` directory:

- `stoplist`: a stoplist

- `sentences`: the data
- `test`: an example input file
- `test.out`: the output for the test input file

Because of rounding errors, etc. your numbers may be just slightly different, but your formatting and list should be the same.

2.6 Writeup

Once you have everything working, play with some different examples to see how well each of the approaches work. Include in your writeup the output from your system on two or more different words over a range of weightings and similarity measures. Your experiments should highlight strengths and weaknesses in the different approaches and allow you to provide some concrete analysis of the approaches.

In addition to the results, provide a short (i.e. a paragraph or two) analyzing your results and explaining them. Include all of this as a file called either `writeup.txt` or `writeup.pdf`.)

Your writeup will be graded based on the quality of your analysis as well as the quality of your chosen examples.

2.7 Extra credit

For extra credit you may implement additional weighting schemes and/or additional similarity measures. If you do this, include examination of these in your writeup and also include the extra parameters in the README. Points will be awarded based on the difficulty of the approach as well as its effectiveness.

2.8 When you're done

When you're all done, follow the directions on the course web page for submitting your work in the dropbox. Make sure that your code compiles, that your files are named as specified. If you get an error when submitting, try changing the name of the folder to include a version number and resubmit.

Include all of your code and your writeup. Your script `wordsim.sh` should be in the base directory and should work from there.

If you worked with a partner, put both people's last names on the submitted directory, but only submit one copy.

If the naming of the files isn't obvious about where various things are being done, please include a README file explaining your organization.

Commenting and code style

Your code should be commented appropriately (though you don't need to go overboard). The most important things:

- Your name (or names) and the assignment number should be at the top of each file you modify.
- Each class and method should have appropriate JavaDoc comments (doc strings if you do it in Python).
- If anything is complicated, put a short note in there to help me out if there are any issues.

This is a non-trivial assignment and it can get complicated, which makes code style and comments very important so that I can understand what you did. For this reason, you will lose points for poorly commented or poorly organized code.

Grading

Part	points
General	20
TF	5
TFIDF	5
PMI	5
L1	5
Euclidean	5
Cosine	5
Efficiency	5
style/commenting	5
extra credit	5
total	60 + 5 extra