

CS 2230 Data Structures & Algorithms

Final Project: Word Similarity

60 points


Learning objectives


- Design, implement, and test an application based on a written specification
- Choose appropriate ADTs and efficient data structures for various tasks
- Use version control to collaborate on a coding project with another person

Unlike the previous programming assignments, we provide specification only at the level of the required behavior of the program, as opposed to at the level of the code or interfaces. It is your job to decide on the classes, interfaces, data structures, and algorithms to use to implement the program.

Submissions

- Nov 17, 11:59pm: Milestone 1 in GitHub (no slip days)
 - Finished Part 1
 - PROGRESS_REPORT_NOV17.txt
- Nov 29, 11:59pm: Milestone 2 in GitHub (no slip days)
 - Finished Part 3, and initial draft of Part 4's written answers
 - PROGRESS_REPORT_NOV29.txt
- Dec 6, 11:59pm: Final version in GitHub (up to 2 slip days if at least 1 partner has them)
 - Finished all Parts

Wherever you see a  that means to write something in your project report. Name the file `project.pdf` and include it in the base of your GitHub repository.

Wherever you see a  that means to write some more Java code. Make sure files you create go within the `edu.uiowa.cs.similarity` package.

See [Checklist.pdf](#) for all of these deliverables in a list.

Clone your repository

We've already created a repository for your team. You can find it at

<https://github.uiowa.edu/cs2230-fa17/project-team-NUM>

where NUM is an integer between 1 and 75, found from your Project Pair in ICON <https://uiowa.instructure.com/courses/64804/groups#tab-6093> .

The course staff members also have read/write access to your repository.
There is nothing in your repository to begin with.

First, have one partner do these steps

1. In NetBeans, clone <https://github.uiowa.edu/cs2230-assignments/semantic-similarity.git>
2. Team | Remote | Push...
3. Specify a Git Repository Location

Remote Name: mine

Repository URL: <https://github.uiowa.edu/cs2230-fa17/project-team-NUM>

- replace NUM with your team number

Username: your hawkid

Password: your hawkid password

Second, the other partner can get the code directly

In NetBeans, clone <https://github.uiowa.edu/cs2230-fa17/project-team-NUM>

- replace NUM with your team number

Introduction

In this project, you will write an application that understands the *semantic similarity*, or closeness of meaning, between two words. For example, the semantic similarity of "car" and "vehicle" is high, while the semantic similarity of "car" and "flower" is low. Semantic similarity has many uses, including finding synonyms, inferring sentiment, and inferring topic.

To compute the semantic similarity of two words, you will first compute each word's *semantic descriptor vector*. Then you will compute the closeness of two vectors using a similarity measure, such as cosine similarity or Euclidean distance.

Semantic descriptor vector

Given a text with n unique words denoted by (w_1, w_2, \dots, w_n) and a word w , let desc_w be the semantic descriptor vector of w computed using the text. desc_w is an n -dimensional vector. The i -th coordinate of desc_w is the number of sentences in which both w and w_i occur.

For example, suppose we are given the following text (the opening of Notes from the Underground by Fyodor Dostoyevsky, translated by Constance Garnett):

I am a sick man. I am a spiteful man. I am an unattractive man. I believe my liver is diseased. However, I know nothing at all about my disease, and do not know for certain what ails me.

For this text, the dimensions of the semantic descriptor vectors will be

[I, am, a, sick, man, spiteful, an, unattractive, believe, my, liver, is, diseased, However, know, nothing, at, all, about, disease, and, do, not, for, certain, what, ails, me]

That is, 1 dimension for each unique word of the text.

The word “man” appears in the first three sentences. Its semantic descriptor vector would be:

[I=3, am=3, a=2, sick=1, man=0, spiteful=1, an=1, unattractive=1, believe=0, my=0, liver=0, is=0, diseased=0, However=0, know=0, nothing=0, at=0, all=0, about=0, disease=0, and=0, do=0, not=0, for=0, certain=0, what=0, ails=0, me=0]

The word “liver” occurs in the fourth sentence, so its semantic descriptor vector is:

[I=1, am=0, a=0, sick=0, man=0, spiteful=0, an=0, unattractive=0, believe=1, my=1, liver=0, is=1, diseased=1, However=0, know=0, nothing=0, at=0, all=0, about=0, disease=0, and=0, do=0, not=0, for=0, certain=0, what=0, ails=0, me=0]

Cosine similarity

The cosine similarity between two vectors $u = (u_1, u_2, \dots, u_N)$ and $v = (v_1, v_2, \dots, v_N)$ is defined as:

$$\text{sim}(u, v) = \frac{u \cdot v}{|u||v|} = \frac{\sum_{i=1}^N u_i v_i}{\sqrt{(\sum_{i=1}^N u_i^2)(\sum_{i=1}^N v_i^2)}}$$

For example, the cosine similarity of “man” and “liver”, given the semantic descriptors above, is

$$\text{sim}(v_{\text{man}}, v_{\text{liver}}) = \frac{3}{\sqrt{26 * 5}} = 0.2631 \dots$$

Getting started: Run the Main program

We've given you a bare bones starter file Main.java, containing the main() method and some basic command-line options parsing. First, follow the directions in

Appendix A to eliminate the compiler error. After that, when you run the file, you'll see it print a usage message.

```
Missing required option: f
usage: Main -f <arg> [-h]
  -f, --file <arg>    input file to process
  -h                  print this help message
```

To provide the required -f option with an argument, set the command line arguments as follows.

Run | Set project configuration | Customize...

Fill in the following fields

Main Class: edu.uiowa.cs.similarity.Main

Arguments: -f "C:\Users\me\similarity\cleanup_test.txt"

Note that "C:\Users\me\similarity\cleanup_test.txt" must be replaced with the actual path to the cleanup_test.txt file your computer. If you get the path wrong, the program will print an error message.

Part 1: Processing and cleaning input (5 points)

To process a text file into useful semantic descriptor vectors, you need to 1) split the text into sentences (lists of words) and 2) clean up the words.

Reading text files into sentences



Write code to read the text files into sentences.

You should assume that the following punctuation always separates sentences: ".", "!", "?", and that is the only punctuation that separates sentences. You should assume that the only the following punctuation is present in the texts: [',', '--', ':', ';', '"', "'"]

We recommend that you use Scanner

(<https://docs.oracle.com/javase/8/docs/api/java/util/Scanner.html>) to read the file sentence by sentence. It allows you to specify a regular expression to use as a delimiter (see the Fish example in the documentation).

Cleaning up the words



Write code to clean up the words.

Small differences in words will cause there to be lots of "unique" words that we wouldn't really consider unique. First, **capitalization**: if we don't do any cleanup then "Man" and "man" will be considered different words. You should convert all words to lower case.

Second, **root words**: if we don't do any cleanup then "glass" and "glasses" will be considered different words. You should find the roots of all words by using "stemming". To do so, use the PorterStemmer (<https://opennlp.apache.org/docs/1.8.3/apidocs/opennlp-tools/opennlp/tools/stemmer/PorterStemmer.html>) from OpenNLP. See the appendix for setting up OpenNLP.

Another problem is that common words (such as "a", "the", "is") may not add much information to our vectors. We call these words **stop words**. We have provided a list of stop words in `stopword.txt`. You should remove any word that appears in file.

Testing

We've provided a file `cleanup_test.txt` that you can try your input processing on. Give the path of that file as the argument to the `-f` command line option. When we ran our input processing on this file, we got the following sentences.

Sentences

[[look, glum, night-cap], [feel, littl, breez], [ah], [whatev, mai, sai, good, aliv, dear, amédé]]

Num sentences

4

Your result may differ slightly but should certainly have words in the correct sentence, lower case (e.g. amédé not Amédé), no stop words (e.g., you), and stemmed words (e.g., littl not little).

Using a command line option to print the sentences

The upcoming Parts will require your program to print other information about the text. To control what information is printed, you will use command line options.



Add an option `-s` so that when it is included, the sentences are printed (e.g., as above) and the number of sentences. The particular formatting is **not important**, as long as we can tell what words are in each sentence and how many sentences there are.

So far, when you include the `-s` option, e.g.,

Arguments: `-f "cleanup_test.txt" -s`

the sentences should be printed. And, when you do not include the `-s` option, e.g.,

Arguments: `-f "cleanup_test.txt"`

nothing should be printed.

Part 2: build the semantic descriptor vectors (10 points)



Write more code that takes the sentences you generated from Part 1 and generate a semantic descriptor vector for each unique word.

You can use the provided file `vector_test.txt` that contains the text from the section "Semantic descriptor vector" to see if your program produces the same vectors as the example (except for *all* unique words, not just "man" and "liver").

Command line argument and output



Add another option `-v` to your program. When `-v` is included, your program should print every unique word along with its semantic descriptor vector. When `-v` is not included, it should not print the words and vectors.

You can pick the formatting of this output, as long as it is clear what the vector is for each word. **You may choose whether or not you want to print the entries in the vectors that are 0, depending on what is convenient.**

Tips

- Don't write all your code in the Main class! For example, we **highly recommend** that you define a class to represent a vector. You need to decide how to represent the vector in a data structure. To future proof your code, you might also consider a vector interface, so that you can try swapping in different implementations of vector.

Part 3: Top-J similar words (10 points)



Write code to compute the Top-J similar words to a given word.

Time to finally use the semantic descriptor vectors for something interesting! In this Part, you will implement a Top-J similar words query. That is, given a text, a query word Q, and an integer J, your program will print the J most similar words to Q.



To compute the similarity, you'll need to create a method that calculates the cosine similarity of two semantic descriptor vectors (see Cosine similarity). The mathematical calculation requires real numbers, so be sure that this method returns a double.

Command line argument and output



Your program needs another option -t that takes an argument of format Q, J. That option will cause the program to print the k most similar words to Q, along with their similarity score. Or, if there is no vector for Q (because Q did not appear in the text), then prints an error message "Cannot compute top-J similarity to Q".

Here are some examples to make sure you've done these steps correctly (again, exact formatting doesn't matter, as long as the information is clearly presented). Note that the order of equally similar words may be different. You should not return the query word itself (which always has cosine similarity of 1.0).

Arguments: -f "easy_sanity_test.txt" -t cat,5
[Pair{wolf,0.8}, Pair{tiger,0.8}, Pair{dog,0.8}, Pair{fox,0.8}, Pair{squirrel,0.8}]

Arguments: -f "easy_sanity_test.txt" -t cat,6
[Pair{wolf,0.8}, Pair{tiger,0.8}, Pair{dog,0.8}, Pair{fox,0.8}, Pair{squirrel,0.8}, Pair{ten,0.0}]

Arguments: -f "easy_sanity_test.txt" -t elephant,3
Cannot compute top-J similarity to elephant

Tips:

- You should stem Q before looking for its vector to make sure you don't accidentally reject the query.

Part 4: Try it on bigger data! (10 points)

Now is a good time to see how your program fares on a large text. Project Gutenberg (<http://www.gutenberg.org/>) has thousands of public domain books that you can download as plain text files. Try these (in order of size):

- <http://www.gutenberg.org/ebooks/41>
- <http://www.gutenberg.org/ebooks/7178>
- <http://www.gutenberg.org/ebooks/2600>

The time it takes to generate the semantic descriptors and run a Top-J query should be no more than ~2 minutes for the first two books. If it is taking significantly longer, see the following tips.



Try to improve your program's running time before continuing.

Tips for slow programs

- Most words do not appear with every other word. Therefore, for really large texts, most of the entries in the vectors will be 0 (this data is called "sparse"). You should consider that fact when designing your data structure for vectors.
- Does your algorithm for the Top-J query theoretically run in $O(N \log N)$ time where N is the number of unique words in the text?
- You can see the breakdown of where time is spent by using NetBeans' profiler. Use the **methods** option when you *configure session*. Under **main**, you'll find how much time was spent in each method.

main	8,206 ms (100%)	8,081 ms (100%)
edu.uiowa.cs.similarity.Main.main (String[])	8,032 ms (97.9%)	8,032 ms (99.4%)
edu.uiowa.cs.similarity.Main.getSentences (java.io.File)	7,919 ms (96.5%)	7,919 ms (98%)



Write about the following:

Let N =the number of unique words in the text file

S =the maximum number of unique words any word appears with

- a) What data structure did you finally use for vectors? What is the asymptotic memory usage of a vector? Of all of your vectors? Is this memory usage reasonable and why?
- b) What algorithm did you finally use for cosine similarity? What is its asymptotic running time? Is this running time reasonable and why?
- c) What algorithm did you finally use for the Top-J calculation? What is its asymptotic running time (might be in terms of J , too)? Is this running time reasonable and why?

- d) What improvements did you make from your *original* code to make it run faster? Give an example of your running time measurements before and after the changes. Describe the information that informed your choices (asymptotic running time analysis, asymptotic memory analysis, and/or profiling).

Part 5: More similarity measures (5 points)

There's many ways we can compare two vectors.



Implement two more similarity measures.

(negative) Euclidean distance between vectors

$$sim_{euc}(v_1, v_2) = -\|v_1 - v_2\|$$

Example test:

$$\begin{aligned}w_1 &= (1, 4, 1, 0, 0, 0) \\w_2 &= (3, 0, 0, 1, 1, 2) \\sim_{euc}(w_1, w_2) &= -5.1961524227\end{aligned}$$

(negative) Euclidean distance between normalized vectors

$$sim_{euc\ norm}(v_1, v_2) = -\left\| \frac{v_1}{|v_1|} - \frac{v_2}{|v_2|} \right\|$$

Example test:

$$\begin{aligned}w_1 &= (1, 4, 1, 0, 0, 0) \\w_2 &= (3, 0, 0, 1, 1, 2) \\sim_{euc\ norm}(w_1, w_2) &= -1.27861316602\end{aligned}$$

Testing



Create JUnit tests for these two similarity measures. You can use the above examples as test cases. Note that `assertEquals` for doubles takes a third argument, `epsilon`, since comparing for exact equality of doubles is fragile. You can use `epsilon=0.00001`.

Command line arguments



To tell your program which measure to use, add another command line option `-m` that can take the arguments `cosine`, `euc`, or `eucnorm`. If the `-m` option is not included, the measure should default to cosine similarity.

Examples:

Arguments: -f "easy_sanity_test.txt" -t cat,5
Uses cosine similarity

Arguments: -f "easy_sanity_test.txt" -t cat,5 -m euc
Uses negative Euclidean distance

Arguments: -f "easy_sanity_test.txt" -t cat,5 -m eucnorm
Uses negative Euclidean distance of norms



Good design: Code that calls the similarity measures (e.g., the Top-J query code) should be reusable! Restructure your program, if needed, to promote reuse (hint: interface or generic type?). The litmus test is to ask yourself: "if I now add a 4th similarity measure, do I need to change the Top-J code? If the answer is "yes", then you are not done with this step.



Run an interesting Top-J query on one of the larger texts. Run it for all three similarity measures and report your findings. Discuss what is the same and different about the results for the 3 similarity measures.

Part 6: Clustering similar words (10 points)

Next, we want to find *groups* of words that are similar to each other. The problem of **clustering** is to partition n points (in our case, a point is a semantic descriptor vector) into clusters of similar points. In k-means clustering, given a desired number of clusters, k , we find k **cluster means** (a center point) and assign each of the n points to the nearest mean.

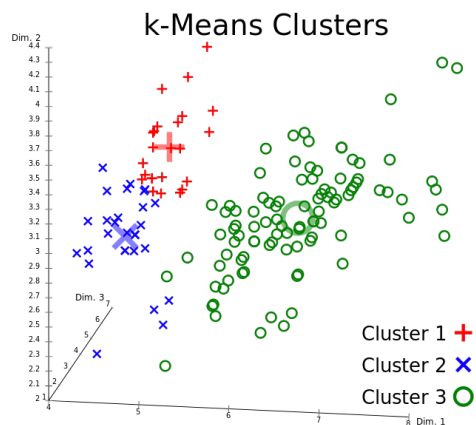


Figure 1: Illustration of k-means clusters where $k=3$ (meaning 3 clusters) and the number of dimensions is 3.

In general k-means clustering tries to figure out what the optimal choice of the cluster means is ("optimal" here means that the distance of each point from the mean is minimized). The **k-means algorithm** approximates a good set of means using *iterative refinement*. We start with

random means, and each iteration moves the means to new places. We can stop when we have an iteration where no points switch between clusters or after a fixed number of iterations.

Here is the pseudocode, where it will always run for a fixed number of iterations `iters`.

```
means = generate k random points
for i in 0 to iters
    clusters = [[]] // k of these
    // assign points to clusters
    for p in points
        find m in means s.t. euclidean_distance(m, p) is minimized
        add p to the cluster for m
    // calculate new means
    for i in 0 to k
        means[i] = centroid of the points in clusters[i]
```

Definition: the centroid of a set of r points w_1, w_2, \dots, w_r is $\frac{w_1 + w_2 + \dots + w_r}{r}$



Implement k-means clustering using this algorithm. To make sure the initial means are reasonable, initialize the k random points by randomly sampling k words without replacement from the database. **NOTE:** k-means only needs to be coded to use Euclidean distance, not the other similarity measures.

Command line options



Add a new option `-k` that takes an argument of the form `k, iters` to tell the program to run k-means clustering using k clusters and running for `iters` iterations. When finished, it should print the words in each cluster. We recommend running the following examples.

Arguments: `-f "easy_sanity_test.txt" -k 3,2`

Run k-means clustering with $k=3$ and 2 iterations. Your results will vary, but you'll probably see that the categories are still mixed up. We need to run for more iterations.

Cluster 0

nine,seven,two,three,eight,four,ten,five,on,

Cluster 1

six,

Cluster 2

banana,parslei,wolf,basil,appl,tomato,tiger,fox,squirrel,squash,cat,dog,pumpkin,

Arguments: `-f "easy_sanity_test.txt" -k 3,100`

Run k-means clustering with $k=3$ and 100 iterations. Did it find the appropriate clusters?

Track the convergence



The k-means algorithm's iterative refinement attempts to minimize the average distance from each cluster mean to the points in the cluster. Enhance your k-means algorithm so that at the end of each iteration, it calculates the average Euclidean distance each point to its cluster's mean. Print one average for each iteration, along with the iteration number (1,2,...,[iters]).

Run the following example again.

Arguments: `-f "easy_sanity_test.txt" -k 3,100`



Using the output, create a scatterplot (using any program you like, e.g., Excel) showing the average distance vs. iteration. After how many iterations did the clustering seem to converge?



Run k-means clustering on one of the larger texts. Run it once for at least 2 different values of k . For each value of k , create a scatterplot as above. Pick the number of iterations to be big enough that the algorithm converges.

Part 7: Cluster representatives (10 points)

It can be hard to read the final output of k-means clustering when the number of points and/or clusters is large. To make the output more interpretable, use your Top-J code to pick J representative words for each cluster. The J representative words should be the words closest to the mean for that cluster. The idea is that J should be small, perhaps 2 to 4.



You may need to modify your Top-J code a bit so that it can be called directly using a vector (the mean) as the query instead of a word.

Command line options



Add a new option `-j` that takes an argument of the form `k, iters, J` to tell the program to run k-means clustering using k clusters and running for `iters` iterations and print the final clusters using only the Top- J representatives for each cluster.

A simple test:

Arguments: -f "easy_sanity_test.txt" -j 3,100,1

Should print a number in one of the clusters, an animal in another cluster, and a food in the other cluster.



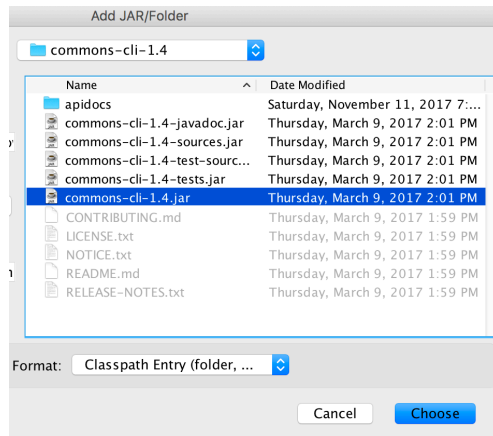
Run the k-means with representative words on one of the larger texts. Pick J between 2-4. Pick a reasonable k. Include a copy of your results with any interesting observations.

Appendix A

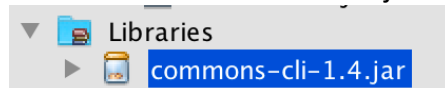
Add Apache CLI to your project

You will use Apache CLI to process command line options. Here are the steps to "install" it.

1. Download Apache CLI http://commons.apache.org/proper/commons-cli/download_cli.cgi and unzip it anywhere.
2. In your NetBeans project, right-click Libraries | Add JAR/folder
3. Navigate to where you unzipped the file. Choose commons-cli-*.jar



4. You should now see that jar file under Libraries



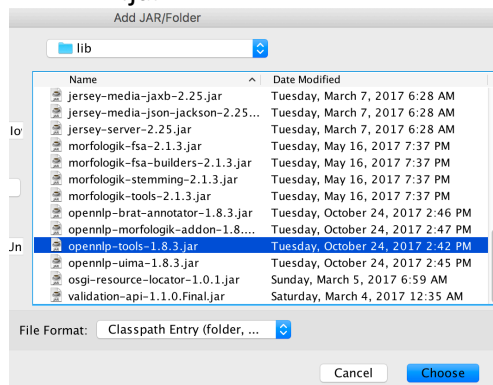
5. The `import org.apache.commons.cli.*` statement in `Main.java` should no longer have errors.

Appendix B

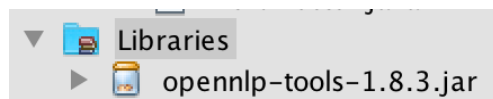
Add OpenNLP to your project

You will use OpenNLP Tools to stem (find roots of) words. Here are the steps to "install" it.

1. Download OpenNLP <https://opennlp.apache.org/download.html> and unzip it anywhere.
2. In your NetBeans project, right-click Libraries | Add JAR/folder
3. Navigate to where you unzipped the file, into the `lib/` folder, and choose `opennlp-tools-*.jar`



1. You should now see that jar file under Libraries



2. You can now import the stemmer library into your Java file without any errors.

```
import opennlp.tools.stemmer.*;
```

Acknowledgements

The concept for this assignment and some of the examples are drawn from ____

Icon files:

- JAVA File by AomAm from the Noun Project
- Pencil by Adrien Coquet from the Noun Project