

Worth: 5%

One type of question encountered in the TOEFL, SAT, and GRE tests is the “Synonym Question”, where students are asked to pick a synonym of a word out of a list of alternatives. For example:

1. vexed (Answer: (a) annoyed)
 (a) annoyed
 (b) amused
 (c) frightened
 (d) excited

For this assignment, you will build an intelligent system that can learn to answer questions like this one. In order to do that, the system will approximate the *semantic similarity* of any pair of words. The semantic similarity between two words is the measure of the closeness of their meanings. For example, the semantic similarity between “car” and “vehicle” is high, while that between “car” and “flower” is low.

In order to answer the TOEFL question, you will compute the semantic similarity between the word you are given and all the possible answers, and pick the answer with the highest semantic similarity to the given word. More precisely, given a word w and a list of potential synonyms s_1, s_2, s_3, s_4 , we compute the similarities of $(w, s_1), (w, s_2), (w, s_3), (w, s_4)$ and choose the word whose similarity to w is the highest.

We will measure the semantic similarity of pairs of words by first computing a *semantic descriptor vector* of each of the words, and then taking the similarity measure to be the *cosine similarity* between the two vectors.

Given a text with n 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 efficiency’s sake, we will store the semantic descriptor vector as a dictionary, not storing the zeros that correspond to words which don’t co-occur with w . 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.

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

`{"i": 3, "am": 3, "a": 2, "sick": 1, "spiteful": 1, "an": 1, "unattractive": 1}`

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

`{"i": 1, "believe": 1, "my": 1, "is": 1, "diseased": 1}`

We store all words in all-lowercase, since we don’t consider, for example, “Man” and “man” to be different words. We do, however, consider, *e.g.*, “believe” and “believes”, or “am” and “is” to be different words. We discard all punctuation.

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:

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

We cannot apply the formula directly to our semantic descriptors since we do not store the entries which are equal to zero. However, we can still compute the cosine similarity between vectors by only considering the positive entries.

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

$$\frac{3 \cdot 1 \text{ (for the word "i")}}{\sqrt{(3^2 + 3^2 + 2^2 + 1^2 + 1^2 + 1^2 + 1^2)(1^2 + 1^2 + 1^2 + 1^2 + 1^2)}} = 3/\sqrt{130} = 0.2631 \dots$$

Part 1.

Subpart (a) `cosine_similarity(vec1, vec2)`

This function returns the cosine similarity between the sparse vectors `vec1` and `vec2`, stored as dictionaries. For example,

```
cosine_similarity({"a": 1, "b": 2, "c": 3}, {"b": 4, "c": 5, "d": 6})
```

should return approximately 0.70.

Subpart (b) `build_semantic_descriptors(sentences)`

This function takes in a list `sentences` which contains lists of strings (words) representing sentences, and returns a dictionary `d` such that for every word `w` that appears in at least one of the sentences, `d[w]` is itself a dictionary which represents the semantic descriptor of `w` (note: the variable names here are arbitrary). For example, if `sentences` represents the opening of *Notes from the Underground* above:

```
[['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']],
```

part of the dictionary returned would be:

```
{ 'man': {'i': 3, 'am': 3, 'a': 2, 'sick': 1, 'spiteful': 1, 'an': 1,
          'unattractive': 1},
  'liver': {'i': 1, 'believe': 1, 'my': 1, 'is': 1, 'diseased': 1},
  ... }
```

with as many keys as there are distinct words in the passage.

Subpart (c) `build_semantic_descriptors_from_files(filenamees)`

This function takes a list of `filenames` of strings, which contains the names of files (the first one can be opened using `open(filenamees[0], "r", encoding="utf-8")`), and returns the a dictionary of the semantic descriptors of all the words in the files `filenamees`, with the files treated as a single text.

You should assume that the following punctuation always separates sentences: `(".", "!", "?")`, and that is the only punctuation that separates sentences. You should also assume that that is the only punctuation that separates sentences. You should assume that the only the following punctuation is present in the texts:

```
[',', ' ', '-', '--', ':', ';', '"', '"']
```

Subpart (d) `most_similar_word(word, choices, semantic_descriptors, similarity_fn)`

This function takes in a string `word`, a list of strings `choices`, and a dictionary `semantic_descriptors` which is built according to the requirements for `build_semantic_descriptors`, and returns the element of `choices` which has the largest semantic similarity to `word`, with the semantic similarity computed using

the data in `semantic_descriptors` and the similarity function `similarity_fn`. The similarity function is a function which takes in two sparse vectors stored as dictionaries and returns a `float`. An example of such a function is `cosine_similarity`. If the semantic similarity between two words cannot be computed, it is considered to be -1 . In case of a tie between several elements in `choices`, the one with the smallest index in `choices` should be returned (*e.g.*, if there is a tie between `choices[5]` and `choices[7]`, `choices[5]` is returned).

Subpart (e) `run_similarity_test(filename, semantic_descriptors, similarity_fn)`

This function takes in a string `filename` which is the name of a file in the same format as `test.txt`, and returns the percentage (i.e., `float` between 0.0 and 100.0) of questions on which `most_similar_word()` guesses the answer correctly using the semantic descriptors stored in `semantic_descriptors`, using the similarity function `similarity_fn`.

The format of `test.txt` is as follows. On each line, we are given a word (all-lowercase), the correct answer, and the choices. For example, the line:

```
feline cat dog cat horse
```

represents the question:

```
feline:
(a) dog
(b) cat
(c) horse
```

and indicates that the correct answer is “cat”.

Part 2.

Subpart (a) Experimenting on a large corpus of text

Download the novels *Swann's Way* by Marcel Proust, and *War and Peace* by Leo Tolstoy from Project Gutenberg, and use them to build a semantic descriptors dictionary. Report how well the program performs (using the cosine similarity similarity function) on the questions in `test.txt`, using those two novels at the same time. Note: the program may take several minutes to run (or more, if the implementation is inefficient). The novels are available at the following URLs:

<http://www.gutenberg.org/cache/epub/7178/pg7178.txt>

<http://www.gutenberg.org/cache/epub/2600/pg2600.txt>

If your program takes too long to run, report the results on a shorter text.

Subpart (b) Experimenting with alternative similarity measures

Two alternative similarity measures between two vectors are: the negative distance in Euclidean space between the vectors, and the negative distances in Euclidean space between the normalized versions of the vectors:

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

$$sim_{euc_{norm}}(v_1, v_2) = -||v_1/|v_1| - v_2/|v_2|||$$

Report the performance of the algorithm using the text of the two novels on `test.txt` using these two similarity measures.

Subpart (c) Experimenting with smaller corpora of text: performance

Plot the performance on `test.txt` on 10%, 20%, ..., 90%, 100% of the text of the two novels, using the cosine similarity measure, as well as the two other possible similarity measures.