

Assignment 1

Information Retrieval and Web Search

Winter 2022

Total points: 100

Issued: 01/14/2022 Due: 01/27/2022

All the code has to be your own (exceptions to this rule are specifically noted below). The code must run on the CAEN environment without additional installation or additional files (except for the data files specified in the assignment).

You can discuss the assignment with others, but the code is to be written individually. You are to abide by the University of Michigan/Engineering honor code; violations will be reported to the Honor Council.

Please use Python 3 for the assignment. Whenever something is not specified in the assignment, that is a design choice you can make; if you make any assumptions, please include them in your write-up.

1. [45 points] Text preprocessing.

Write a Python program that preprocesses a collection of documents. You will test this program on the Cranfield dataset, which is a standard Information Retrieval text collection, consisting of 1400 documents from the aerodynamics field. The dataset `cranfield.zip` is available from the Files section on Canvas under Assignments/.

Programming guidelines:

Write a program called *preprocess.py* that preprocesses the collection. Assume that the documents are available in a folder whose name you will read from the command line. For testing purposes, use the *cranfieldDocs/* folder that you will obtain after unpacking the `cranfield.zip` archive.

Include the following functions in *preprocess.py*:

a. Function that removes the SGML tags.

Name: *removeSGML*;

Input: string;

Output: string

b. Function that tokenizes the text.

Name: *tokenizeText*; input: string; output: list (of tokens)

The tokenizer should separate the punctuation from the words, whenever the punctuation is not an integral part of the word.

For instance:

The current population of U.S.A. is 332,087,410 as of Friday, 01/22/2021, based on Worldometer elaboration of the latest United Nations' data.

should be tokenized as

The current population of U.S.A. is 332,087,410 as of Friday , 01/22/2021 , based on Worldometer elaboration of the latest United Nations ' data .

Your tokenizer should represent your best effort to correctly address the following cases, among others.:

- tokenization of . (do not tokenize acronyms, abbreviations, numbers)
- tokenization of ' (expand when needed, e.g., I'm -> I am; tokenize the possessive, e.g., Sunday's -> Sunday 's; etc.)
- tokenization of dates (keep dates together)
- tokenization of - (keep phrases separated by - together)
- tokenization of , (do not tokenize numbers)

Note that the use of regex is acceptable to the extent that your solution is not trivialized.

c. Function that removes the stopwords.

Name: *removeStopwords*; input: list (of tokens); output: list (of tokens)

Use the list of stopwords available under the Files section on Canvas, under Assignments/

d. Function that stems the words.

Name: *stemWords*;

Input: list (of tokens);

Output: list (of stemmed tokens)

You are allowed to use a publicly available implementation of the Porter stemmer.

E.g., <https://tartarus.org/martin/PorterStemmer/python.txt>

(others are acceptable as well)

The main program should perform the following sequence of steps:

- i. open the folder containing the data collection, provided as the first argument on the command line (e.g., *cranfieldDocs/*), and read one file at a time from this folder.
- ii. for each file, apply, in order: *removeSGML*, *tokenizeText*, *removeStopwords*, *stemWords*
- iii. in addition, write code to determine (this is after step ii above):
 - the total number of words in the collection (numbers should be counted as words)
 - vocabulary size (i.e., number of unique terms)

- most frequent 50 words in the collection, along with their frequencies (list in reverse order of their frequency, i.e., from most to least frequent)

The *preprocess.py* program should be run using a command like this:

```
% python preprocess.py cranfieldDocs/
```

It should produce a file called *preprocess.output* with the following content (words are to be listed in reverse order of their frequency):

```
Words [total-number-of-words]
Vocabulary [total-number-of-unique-words]
Top 50 words
Word1 [frequency-of-Word1]
Word2 [frequency-of-Word2]
...
Word50 [frequency-of-Word50]
```

Example of a preprocess.output file:

```
Words 125478
Vocabulary 339
Top 50 words
Word1 250
Word2 123
...
```

Write-up guidelines:

Create a text file called *preprocess.answers*, and include the following information:

- Total number of words in the Cranfield collection.
- Vocabulary size of the Cranfield collection.
- The minimum number of unique words in the Cranfield collection accounting for 25% of the total number of words in the collection?
Example: if the total number of words in the collection is 100,
and we have the following word-frequency pairs: airplane - 30 space - 10
clear - 8 cut - 7 etc. the answer to this question will be 1 (1 word accounts for 25%
of the total 100 words)
- Pick two subsets of the Cranfield dataset. Determine and report the number of words and the size of the vocabulary for the subsets you selected. Use this information to derive the K and beta parameters required by the application of the Heaps law. Use these values to predict what would be the vocabulary size if the corpus were to increase to 1,000,000 words, or to 100,000,000 words.

2. [55 points] Language identification.

Implement a language identifier, using the letter-based bigram probability model discussed in class, with add-one smoothing. The dataset to be used for this assignment is included in `languageIdentification.data.zip`, available from the Files section on Canvas, under Assignments/.

Programming guidelines:

Once you unpack the data archive, you will obtain a folder called *languageIdentification.data/*. Store this data folder in the same folder as your program. There are five files included in the *languageIdentification.data/* folder: three files with data to use for training (*English*, *French*, *Italian*, stored under a subfolder called *training/*), one file with data to use for test (*test*), and one solution file to use for evaluation (*solution*).

Write a Python program called *languageIdentification.py*. The program should include the following functions:

a. Function to train a bigram language model.

Name: *trainBigramLanguageModel*;

Input: string (training text in a given language);

Output: dictionary with character frequencies collected from the string; dictionary with character-bigram frequencies collected from the string

Given an input string, this function will calculate the frequencies for all the single characters and for all the bigram characters in the string.

b. Function to determine the language of a string.

Name: *identifyLanguage*

Input: string (text for which the language is to be identified); list of strings (each string corresponding to a language name); list of dictionaries with single character frequencies (each dictionary corresponding to the single character frequencies in a language); list of dictionaries with bigram character frequencies (each dictionary corresponding to the bigram character frequencies in a language);

Output: string (the name of the most likely language).

Note: in the input lists, elements at a given position K in the lists correspond to the same language L.

The main program should perform the following sequence of steps:

- i. Use the *trainBigramLanguageModel* function to build unigram and bigram dictionaries for each of the three language files provided as training.
- ii. Open the test file, provided as the first argument on the command line, and for each line in the test file, apply the *identifyLanguage* function.

The *languageIdentification.py* program should be run using a command like this:

```
% python languageIdentification.py languageIdentification.data/test
```

It should produce a file called *languageIdentification.output*, with the following content:

Line1 Language1

Line2 Language2

...

Line300 Language300

(where LineN represents the line number in the test file, and LanguageN is the language determined by the *identifyLanguage* function as the most likely one for that particular line in the test file)

Write-up guidelines:

Create a text file called *languageIdentification.answers*, and include the following information:

- Accuracy of your language identifier when comparing the output of your system with the solution provided? (in other words, what is the percentage of predictions that are correct)

General submission instructions:

- Include all the files for this assignment in a folder called *[your-username].Assignment1/* **Do not** include the data folders, i.e., *cranfieldDocs/*, *languageIdentification.data/*, *stopwords*. For instance, *mihalcea.Assignment1/* will contain *preprocess.py*, *file_with_Porter_stemmer_code*, *preprocess.output*, *preprocess.answers*, *languageIdentification.py*, *languageIdentification.output*, *languageIdentification.answers*.
- Archive the folder using zip and submit on Canvas by the due date.
- Include your name and username in each program and in all the .answers files
- Make sure all your programs run correctly on the CAEN machines.