

TP2
IFT2015
Data structures
Fall 2022

In this TP, you will implement two different `Map` classes to serve basic natural language processing (NLP) functions. Particularly, you will constitute datasets of numerous documents. You will use the first `Map` to store the words of these documents (`WordMap`). The keys for the `WordMap` will be the words themselves, and the values will be references to the second type of `Maps` (`FileMaps`), which keep the word occurrences in the files. The keys for the `FileMaps` will be the file names, and the values will be a `List` of the positions of the associated words in the corresponding files.

Consider the following `dataset` which contains three documents (PS. a real `dataset` will contain many more files and longer texts). The files of a `dataset` are stored in a `directory`.

0000026.txt:

```
With stunning photos and stories, National Geographic Explorer  
Wade Davis celebrates the extraordinary diversity of the world's  
indigenous cultures, which are disappearing from the planet at  
an alarming rate.
```

0000048.txt:

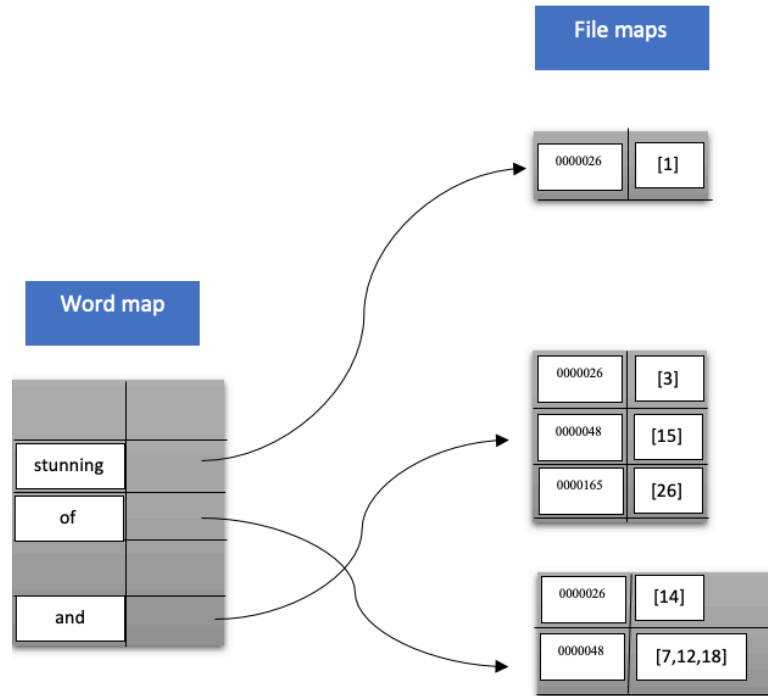
```
Photographer Phil Borges shows rarely-seen images of people from  
the mountains of Dharamsala, India, and the jungles of the  
Ecuadorean Amazon. In documenting these endangered cultures, he  
intends to help preserve them.
```

0000165.txt:

```
In this deceptively casual talk, Charles Leadbeater weaves a  
tight argument that innovation isn't just for professionals  
anymore. Passionate amateurs, using new tools, are creating  
products and paradigms that companies can't.
```

Figure 1 shows a fraction of the data structure for this example.

- Before building the data structure, you will need to preprocess the text by replacing all punctuation marks by spaces, replacing multiple spaces by a single one, and converting all characters to lowercase.
- To implement the `Maps`, you will override the functions in `WordMap.java` and `FileMap.java` provided as template classes, which both implement the `Map` interface. You are allowed to define any additional functions of your own.
- You can use the `hashCode()` method of the `Object` class, or override it.



- You must implement your load factor management so that if the insertion of a word causes the load factor to go above 0.75, then you should resize the capacity of the word Map by $2 \times \text{the previous size} + 1$.

Once your data structure is completed, you will implement two essential operations used in NLP:

1. suggesting the next most probable word like an auto-complete function which is related to the concept named **bi-grams** in NLP;
2. finding the most relevant document for a key word search which is related to the concept named **TFIDF** (term frequency-inverse document frequency) in NLP.

Bigrams

A bigram is a segment of text consisting of two consecutive words that occur in a provided text at least once. Bi-grams are very informative means to demonstrate the semantic associations between words. As an example, the bigrams in the following text: “*Photographer Phil Borges shows rarely-seen images of people from the mountains of Dharamsala, India*” are:

Photographer Phil

Phil Borges

Borges shows

shows rarely

rarely seen

seen images

...

Dharamsala, India

Bigrams can be used to suggest the next most probable word that can appear after a typed word in a search engine, and to improve the prediction of an auto-completion system.

Consider the following seven texts:

1. Thank you so much for your help.
2. I **really** appreciate your help.
3. I **really** like what you said.
4. I apologize for coming over unannounced like this.
5. Sorry, do you know what time it is?
6. I'm **really** sorry for not inviting you.
7. I **really** like your watch.

Suppose after training, our model learns the occurrences of every two-words to determine the most probable one, W_2 , after another, W_1 . For example, from the sentences 2, 3, 6, and 7, after the word "**really**" the words "appreciate", "like", or "sorry" occur.

Let's calculate the probability of observing the second word, W_2 , occurring after word W_1 . We apply the relation of the conditional probability:

$$P(W_2 | W_1) = C(W_1, W_2) / C(W_1)$$

The probability of word W_2 given the preceding word W_1 , $P(W_2 | W_1)$, equals the count of their bigram, or the co-occurrence of the two words, $C(W_1, W_2)$ divided by the count of W_1 , $C(W_1)$.

From our example, let's find the word that has the highest probability to appear after the word "**really**". We need $C(\text{"really"}) = 4$, and $C(\text{"really"}, \text{"appreciate"}) = 1$, $C(\text{"really"}, \text{"like"}) = 2$, and $C(\text{"really"}, \text{"sorry"}) = 1$. The most probable word to come after "**really**" is "like", with $P(\text{"like"} | \text{"really"}) = 0.5$, $C(\text{"really"}, \text{"like"}) / C(\text{"really"}) = 2/4 = 0.5$. The probabilities $P(\text{"appreciate"} | \text{"really"}) = P(\text{"sorry"} | \text{"really"}) = 0.25$.

To the query "*the most probable bigram of really*" your program, given the above dataset, should suggest the next most probable word that can occur after the word "**really**", or the most probable bigram of "**really**", which is the word "like". *If there are two words or more with the same probability, return the smallest word based on the lexicographic order.*

TFIDF (Term frequency-inverse document frequency)

This is a score that reflects the importance of a word in a document. In NLP, a word is effective for a file to be categorized if it occurs frequently in that file, but has very few

occurrences in other texts in the dataset. To calculate TFIDF, we first calculate the word (term) frequency:

$$TF(w) = \text{count}(w) / \text{totalW}$$

where $\text{count}(w)$ is the number of times w appears in a document, and totalW is the total number of words in the document (length in terms of words). Then the inverse document frequency is calculated to weigh down the words that are frequent also in other files while scale up the rare ones:

$$IDF(w) = \ln(\text{totalD} / \text{count}(d, w))$$

where totalD is the total number of documents considered, and $\text{count}(d, w)$ is the number of documents with w in it. And, finally, the TFIDF is calculated by:

$$TFIDF(w) = TF(w) \times IDF(w)$$

Imagine a search engine that uses TFIDF to rank documents based on a provided word. You will be given a word and you should propose the most relevant document in the dataset that ranks first based on it. PS. When a word is not present in any document in the dataset, $TF(w) = 0$ and $IDF(w) = \text{infinity}$, then consider

$TFIDF(w) = 0$. *If there are two or more documents with the same $TFIDF(w)$, return the one with the smallest document name based on the lexicographic order.*

Consider the following dataset of four documents:

902.txt:

This article is about the astronomical object. For other uses, see Planet (disambiguation). A planet is a large, rounded astronomical body that is neither a star nor its remnant. The best available theory of planet formation is the nebular hypothesis, which posits that an interstellar cloud collapses out of a nebula to create a young protostar orbited by a protoplanetary disk. Planets grow in this disk by the gradual accumulation of material driven by gravity, a process called accretion.

900.txt:

The discovery of other solar system wanderers rivaling Pluto in size suddenly had scientists asking what wasn't a planet. They put their heads together in 2006 and came up with three conditions for planethood: A planet must orbit the sun, be large enough so that its own gravity molds it into a spherical shape, and it must have an orbit free of other small objects. Unfortunately for Pluto, our one time ninth planet failed to meet the third condition.

901.txt:

The largest known small bodies, in the conventional sense, are several icy Kuiper belt objects found orbiting the Sun beyond the orbit of Neptune. Ceres which is the largest main belt asteroid and is now considered a dwarf planet is roughly 950 km (590 miles) in diameter.

903.txt:

The collapse of International Coffee

Organization, ICO, talks on export quotas yesterday removes the immediate need to reinstate U.S. legislation allowing the customs service to monitor coffee imports, analysts here said. The Reagan administration proposed in trade legislation offered Congress last month that authority to monitor coffee imports be resumed. That authority lapsed in September 1986. A bill also was introduced by Rep. Frank Guarini (DN.J.) However, the failure of the ICO talks in London to reach agreement on export quotas means the U.S. legislation is not immediately needed, one analyst said. Earlier supporters of the coffee bill hoped it could be passed by Congress quickly. "You're going to have a hard time convincing Congress (now) this is an urgent issue," the coffee analyst said.

Imagine a browser asked to find the most relevant document given word "**planet**": search planet. You will have to compute the TFIDF of the word "**planet**" in each document of the dataset:

$\text{TFIDF}(\text{"planet"})$ in document 902 = $\text{TF}(\text{"planet"}) \times \text{IDF}(\text{"planet"}) = 0.0109$

$\text{TF}(\text{"planet"}) = 3/79$

$\text{IDF}(\text{"planet"}) = \ln(4/3)$

$\text{TFIDF}(\text{"planet"})$ in document 900 = $\text{TF}(\text{"planet"}) \times \text{IDF}(\text{"planet"}) = 0.0109$

$\text{TF}(\text{"planet"}) = 3/79$

$\text{IDF}(\text{"planet"}) = \ln(4/3)$

$\text{TFIDF}(\text{"planet"})$ in document 901 = $\text{TF}(\text{"planet"}) \times \text{IDF}(\text{"planet"}) = 0.0061$

$\text{TF}(\text{"planet"}) = 1/47$

$\text{IDF}(\text{"planet"}) = \ln(4/3)$

$\text{TFIDF}(\text{"planet"})$ in document 903 = $\text{TF}(\text{"planet"}) \times \text{IDF}(\text{"planet"}) = 0.0$

$\text{TF}(\text{"planet"}) = 0$

$\text{IDF}(\text{"planet"}) = \ln(4/3)$

Therefore, the most relevant document about the word "planet" is document 900 with the optimal TFIDF score of 0.0109. This score is the same for document 902, but in the lexicographic order 900 comes before 902.

What your program should do

Your program will read an input file that consists of multiple queries (see below). You will write your answers, one by line on `stdout`. There are two types of queries. One is for suggesting the next most probable word that appear after a given word:

the most probable bigram of <word>

The second is for retrieving the most relevant document of a given word:

search <word>

Input (two file names)

1. The directory name of the dataset, which contains the documents. This directory will be stored at the level of your `Java` project;
2. The query file name, where each line is asking for one of the two query types above. Note that the query file will also be stored at the level of your `Java` project.

Output (full example)

You should perform the queries and write the answers, one by line, on `stdout`. Given the directory `dataset2` and query file `query.txt`:

```
search planet
```

```
the most probable bigram of new
```

```
search species
```

and the command line: `java Main dataset2 query.txt`, your program will output on `stdout`:

```
900.txt
```

```
new york
```

```
63.txt
```

Java code and submission

- You may form teams of two or less programmers
- You must follow the rules of OOP; use `CamelCase` for your identifiers; and comment your code
- All `Java` code files and input files must be in the same directory: no packages. The class that contains the main function must be called `Main.java`
- For submission on `StudiUM`, archive your directory (only the names `TP2.tar.gz` or `TP2.zip` will be allowed)
- You have until Sunday December 9 at 23h59 to submit your solution on `StudiUM`
- A penalty of 20% per day starting on December 10 at 00h00 will be applied if you submit after the submission deadline.

- You must follow the exact output format of one answer per line on the `stdout`, corresponding to each line in the query file.

Evaluation

- 10% if your code compile properly without any argument
- 10% if your code executes and find the correct results on the provided solved examples
- 40% if your code executes and find the correct results on unseen examples
- 20% if your code is efficient (how this will be done is yet to be determined)
- 20% if your code follows OOP principles, is lisible, and contains appropriate comments.

NB. Don't forget to include your names and matricule numbers in all your `Java` classes.

Questions. If you have questions, use TP2's forum on StudiUM, or contact one of the TAs or the Prof.

HAVE FUN & GOOD LUCK!