

EDAN20 Language Technology -Lab 1.

Nils Romanus

September 2022

1 Creating an Indexer

Using the regex library coupled with baseline python text operations one may create an indexer for novels written by Selma Lagerlöf.

1.1 Creating a Tokenizer

```
1 # Write your regex here
2 regex = r'\p{L}+'
3 # Write your code here
4 def tokenize(text: str) -> list:
5     return re.findall( r'\p{L}+', text)
```

Listing 1: Tokenizer

By using regex the regex pattern in code snippet 1 coupled with the *findall* call one may create a tokenizer that returns a list of all the words in a text.

1.2 Extracting indices

```
1 # Write your code here
2 from typing import Dict, List
3 def text_to_idx(words: List[str], text: str) -> Dict[str, List[int]]:
4     # get unique words
5     words = set(words)
6     idxs = {}
7
8     for word in words:
9         word_delimited = r'\b' + word + r'\b'
10        for match in re.finditer(word_delimited, text):
11            if word not in idxs.keys():
12                idxs[word] = [match.start()]
13            else:
14                idxs[word] += [match.start()]
15
16    return idxs
```

Listing 2: Indexer

By using the tokenizer from code snippet 1 one may create the indexer displayed in code snippet 2. This is done by using the *finditer* regex function with the regex pattern for every word in the list of word extracted from the tokenizer. For every word its indices are generated by using the *start* function on each match. The results are subsequently stored in a map structure where the words constitutes the keys and the list of indices the corresponding values. This function can subsequently be applied to a tokenized version of an entire file, thus generating an index for the entire novel found in that file.

1.3 Creating a Master Index

```

1 # write your code here
2 def get_master_index(file_names):
3     # find full vocabulary across documents
4     file_tokens = {} # map file name to tokens
5     master_tokens = set()
6     for file_name in file_names:
7         with open("Selma/" + file_name, "r") as f:
8             tokens = set(tokenize(f.read().lower()))
9
10            master_tokens = master_tokens.union(tokens)
11
12            # save tokens for later
13            file_tokens[file_name] = tokens
14
15    # create index for each document
16    file_to_index = {} # {file_name: {token: idx}}
17    for file_name in file_names:
18        with open("Selma/" + file_name, 'r') as f:
19            tokens = file_tokens[file_name]
20            file_to_index[file_name] = text_to_idx(tokens, f.read()
21            .lower())
22
23    # create master index for each word
24    _master_index = {
25        token: {file_name: idx[token]
26                for file_name, idx in file_to_index.items()
27                if token in idx}
28        for token in master_tokens
29    }
30    return _master_index

```

Listing 3: Master Indexer

By using the indexer on each file in the set of files in the *Selma* corpus one may create a master index covering the entire corpus. As code snippet 3 shows this is done by iterating through all files, generating all tokens for each file and calling the indexer function displayed in code snippet 2. The results are subsequently stored in a nested map structure with the token as primary with its value being a map with the file names as keys and list of indices as values.

1.4 Generating Concordances

```

1 # Write your code here
2 def concordance(word, master_index, window) -> None:
3     assert word in master_index
4
5     for file_name, indices in master_index[word].items():
6         print(file_name)
7
8         with open("Selma/" + file_name, 'r') as f:
9             text = f.read().lower().replace("\n", " ")
10            for start in indices:
11                print("\t", text[start - window:start + window])

```

Listing 4: Concordances

The master index can subsequently be used to generate all concordances of a word in the entire corpus. This is done by using the index of the word in the master index map and printing the slice of the text where the word occurs including the *window* number of characters before and after.

1.5 Representing Documents with *tf-idf*

```
1 # Write your code here
2 from collections import Counter
3 import numpy as np
4 import itertools
5
6 def get_tfidf(file_names):
7     # create index for each document
8     file_to_count = {}
9     all_tokens = set()
10    for file_name in file_names:
11        with open("Selma/" + file_name, 'r') as f:
12            counter = Counter(tokenize(f.read().lower()))
13
14            file_to_count[file_name] = counter # {file_name: {token
15            : count}}
16            all_tokens = all_tokens.union(counter.keys())
17
18    # Create tf
19    tf = {file_name: {} for file_name in file_names} # {file_name:
20    {token: tf}}
21    for file_name, token_to_count in file_to_count.items():
22        for token, count in token_to_count.items():
23            tf[file_name].update({token : count/sum(list(
24            token_to_count.values()))})
25
26    # Compute idf
27    idf = {} # {token: idf}
28    N = len(file_names)
29    for token in all_tokens:
30        idf[token] = np.log10(N / len([_
31        for file_name in file_names
32        if token in file_to_count[
33        file_name]]))
34
35    # combine into tfidf
36    tfidf = {file_name : {} for file_name in file_names}
37    for all_token in all_tokens:
38        for file_name, token_to_tf in tf.items():
39            value = token_to_tf[all_token]*idf[all_token] if
40            all_token in token_to_tf else 0.0
41
42            tfidf[file_name].update({all_token: value})
43
44    return tfidf
```

Listing 5: tf-idf

In order to compare documents one may use the *tf-idf*. This metric can be generated by using the function displayed in code snippet 5. The functions work

by firstly creating a nested map with file name as the primary key and a value consisting of a mapping between the tokens in the text and their occurrence counts. The relative term frequency (tf) can subsequently be computed and stored in a similar nested map structure by using token counts. Similarly, the inverse document frequency (idf) can be computed by iterating through all tokens in the full corpus and checking which documents they occur in. The tf and idf can be subsequently combined to form the tf-idf which is stored in a nested map structure using the file name as the primary key and token as secondary key.

1.6 Comparing Documents

```

1 # Write your code here
2 from scipy.spatial.distance import cosine
3 import pandas as pd
4
5 def cosine_similarity(document1, document2, tfidf=tfidf, file_names
6 =get_files('Selma', 'txt')): # e.g. troll.txt
7     if not tfidf:
8         tfidf = get_tfidf(file_names)
9
10    return np.dot(list(tfidf[document1].values()), list(tfidf[
11 document2].values())) \
        / np.linalg.norm(list(tfidf[document1].values()))
        \
        / np.linalg.norm(list(tfidf[document2].values()))

```

Listing 6: Cosine Similarity

The tf-idf metric is useful for comparing documents. By using the normalized dot product (cosine similarity) of two documents tf-idf one may compute a similarity score between a pair of documents as shown in code snippet 6. One may perform this operation for all documents in the corpus and thereby generate the similarity matrix displayed in table ??.

	troll.txt	kejsaren.txt	marbacka.txt	herrgard.txt	nils.txt	osynliga.txt	jerusalem.txt	bannlyst.txt	gosta.txt
troll.txt	1.000000	0.181284	0.147154	0.004074	0.188475	0.192597	0.007057	0.088621	0.195738
kejsaren.txt	0.181284	1.000000	0.071121	0.000740	0.049663	0.051108	0.001834	0.024009	0.048018
marbacka.txt	0.147154	0.071121	1.000000	0.003614	0.084741	0.093168	0.004873	0.036810	0.080168
herrgard.txt	0.004074	0.000740	0.003614	1.000000	0.005068	0.004826	0.370689	0.000949	0.003110
nils.txt	0.188475	0.049663	0.084741	0.005068	1.000000	0.110574	0.004539	0.050982	0.104826
osynliga.txt	0.192597	0.051108	0.093168	0.004826	0.110574	1.000000	0.028300	0.052055	0.124755
jerusalem.txt	0.007057	0.001834	0.004873	0.370689	0.004539	0.028300	1.000000	0.006460	0.004321
bannlyst.txt	0.088621	0.024009	0.036810	0.000949	0.050982	0.052055	0.006460	1.000000	0.049037
gosta.txt	0.195738	0.048018	0.080168	0.003110	0.104826	0.124755	0.004321	0.049037	1.000000

Table 1: Similarity Matrix

2 Google Indexing

Jeff Dean's presentation on Google indexing techniques covers Google's indexing techniques and the infrastructure around them. Slide 14 titled *Research Project*

1997 displays an indexing architecture that is similar to the one used in the lab. In the indexing script, one does not really have a doc server however the *Selma* directory has the same function. The `open(Selma/ + file_name, r)` constitutes a query to our so-called 'docs server' and the only index shard is our master index map. Comparing this to slide 37 *Dealing with growth* one may note that our program does not make use of a cache server functionality for common queries. We also do not use any encodings for the document, instead they are 'encoded' with simply their file name, e.g 'troll.txt'. On the other hand the encoding technique is, in a way, similar to the one displayed on slide 45. In 1997 Google also used a map structure where 'docid' corresponds to a list of hit mappings. Similar to in our program the hits contain the index, however, it also contains other information such as font type or size. This information is absent in our index mapping. Also, we do not store the number of occurrences of a certain hit, this is instead accessed through computing the length of the list containing the indices of the hit. We also do not make use of index shard replicas as shown in the slide. This works for the simple application of indexing a few of Selma Lagerlöfs novels however as Dean notes in his presentation it also comes with some problems. One benefit of our program is that it is easy to set up and keep track of. On the other hand, our program is computationally expensive and vulnerable to machine failure.