

Introduction to Web Science

Assignment 8

Prof. Dr. Steffen Staab

staab@uni-koblenz.de

René Pickhardt

rpickhardt@uni-koblenz.de

Korok Sengupta

koroksengupta@uni-koblenz.de

Olga Zagovora

zagovora@uni-koblenz.de

Institute of Web Science and Technologies
Department of Computer Science
University of Koblenz-Landau

Submission until: January 11, 2017, 10:00 a.m.

Tutorial on: January 13, 2017, 12:00 p.m.

Please look at all the lessons of part 2 in particular **Similarity of Text** and **graph based models**

For all the assignment questions that require you to write code, make sure to include the code in the answer sheet, along with a separate python file. Where screen shots are required, please add them in the answers directly and not as separate files.

Other than that this sheet is mainly designed to review and apply what you have learnt in part 2 it is a little bit larger but there is also more time over the x-mas break. In any case we wish you a mery x-mas and a happy new year.

Team Name: Tango

1. Mariya Chkalova
mchkalova@uni-koblenz.de
2. Arsenii Smyrnov
smyrnov@uni-koblenz.de
3. Simon Schauß
sschauss@uni-koblenz.de

4. Lukas Härtel
lukashaertel@uni-koblenz.de

1 Similarity - (40 Points)

This assignment will have one exercise which is divided into four subparts. The main idea is to study once again the web crawl of the Simple English Wikipedia. The goal is also to review and apply your knowledge from part 2 of this course.

We have constructed two data sets from it which are all the articles and the link graph extracted from Simple English Wikipedia. The extracted data sets are stored in the file <http://141.26.208.82/store.zip> which contains a pandas container and can be read with pandas in python. In subsection “1.5 Hints” you will find some sample python code that demonstrates how to easily access the data.

With this data set you will create three different models with different similarity measures and finally try to evaluate how similar these models are.

This assignment requires you to handle your data in efficient data structures otherwise you might discover runtime issues. So please read and understand the full assignment sheet with all the tasks that are required before you start implementing some of the tasks.

1.1 Similarity of Text documents (10 Points)

1.1.1 Jaccard - Similarity on sets

1. Build the word sets of each article for each article id.
2. Implement a function `calcJaccardSimilarity(wordset1, wordset2)` that can calculate the jaccard coefficient of two word sets and return the value.
3. Compute the result for the articles **Germany** and **Europe**.

1.1.2 TF-IDF with cosine similarity

1. Count the term frequency of each term for each article
2. Count the document frequencies of each term.
3. For each article id provide a dictionary of terms occurring in the article together with their tf-idf scores as the corresponding values.
4. Implement a function `calculateCosineSimilarity(tfIdfDict1, tfIdfDict2)` that computes the cosine similarity for two sparse tf-idf vectors and returns the value.
5. Compute the result for the articles **Germany** and **Europe**.

1.2 Similarity of Graphs (10 Points)

You can understand the similarity of two articles by comparing their sets of outlinks (and see how much they have in common). Feel free to reuse the `computeJaccardSimilarity` function from the first part of the exercise. This time do not apply it on the set of words within two articles but rather on the set of outlinks being used within two articles. Again compute the result for the articles **Germany** and **Europe**.

Answer on 1.1 and 1.2 We used the recompiled `store.h5` file with ASCII characters only.

Listing 1 Tasks 1.1 and 1.2

```
1:
2: # coding: utf-8
3:
4: import pandas as pd
5: import string
6: import math
7:
8: store = pd.HDFStore('store2.h5')
9: df1=store['df1']
10: df2=store['df2']
11: word_dict = dict()
12: tf_dict = dict()
13: df_dict = dict()
14: tfidf_dict = dict()
15: euclead_dict = dict()
16: ARTICLES = len(df1)
17:
18: def create_set(text):
19:     words = text.lower().split()
20:     words_pure = set(w.strip(string.punctuation) for w in words)
21:     return words_pure
22:
23: def create_dict(index):
24:     try:
25:         word_dict[index] = create_set(df1.text.loc[index])
26:     except:
27:         print (article_name)
28:         print (df1[df1.name==article_name].text)
29:
30:
31: # 1.1.1. 2) Implementation of the function calcJaccardSimilarity
32:
33: def calcJaccardSimilarity(set1, set2):
34:     intersection = set1.intersection(set2)
35:     union = set1.union(set2)
```

```
36:     return len(intersection)/len(union)
37:
38: def getId(article):
39:     idx = df1[df1.name==article].index.tolist()[0]
40:     return idx
41:
42: def calc_tf_df(index):
43:     a_tf_dict = dict()
44:     text = df1.text.loc[index]
45:     a_unique_words = word_dict[index]
46:     for term in a_unique_words:
47:         tf = text.count(str(term))
48:         a_tf_dict[term] = tf
49:         if term in df_dict:
50:             df_dict[term] += 1
51:         else:
52:             df_dict[term] = 1
53:     tf_dict[index] = a_tf_dict
54:
55: def calc_tfidf(index):
56:     a_tfidf = dict()
57:     a_unique_words = word_dict[index]
58:     sum_tfidf = 0
59:     for term in a_unique_words:
60:         a_tf = tf_dict[index]
61:         tfidf = a_tf[term] * math.log(ARTICLES/df_dict[term])
62:         sum_tfidf = sum_tfidf + tfidf * tfidf
63:         a_tfidf[term] = tfidf
64:     tfidf_dict[index] = a_tfidf
65:     euclead_dict[index] = math.sqrt(sum_tfidf)
66:
67: def calc_scalar(tfIdfDict1, tfIdfDict2):
68:     wordset1 = set(tfIdfDict1.keys())
69:     wordset2 = set(tfIdfDict2.keys())
70:     common_words = wordset1.intersection(wordset2)
71:     sum_products = 0
72:     for term in common_words:
73:         try:
74:             sum_products = sum_products + tfIdfDict1[term] *
75:                 ↪ tfIdfDict2[term]
76:         except:
77:             print (term)
78:     return sum_products
79:
80: def calc_Euclead(tfIdfDict):
81:     coords = list(tfIdfDict.values())
82:     sum_coords = 0
83:     for value in coords:
84:         sum_coords = sum_coords + value * value
```

```
84:     return math.sqrt(sum_coords)
85:
86:
87: # 1.1.2 4) Implementation of the fuction calculateCosineSimilarity
88:
89: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
90:     return calc_scalar(tfIdfDict1, tfIdfDict2) / (calc_Euclead(
91:         ↪ tfIdfDict1) * calc_Euclead(tfIdfDict2))
92:
93: if __name__ == '__main__':
94:
95:     # 1.1.1. 1) create dictionary that contains set of unique
96:     ↪ words (value) for each article id (key)
97:
98:     for index in df1.index.tolist():
99:         create_dict(index)
100:
101:     # 1.1.1. 3) Jaccard Similarity for articles Germany, Europe
102:
103:     Jaccard_sim = calcJaccardSimilarity(word_dict[getId('Germany')]
104:         ↪ ], word_dict[getId('Europe')])
105:     print ('Jaccard Similarity coefficient: ', Jaccard_sim)
106:
107:     # 1.1.2. 1), 2) populate dictionaries tf_dict and df_dict:
108:     # tf_dict has article id as a key and a dictionary (term (key)
109:     ↪ - term frequency (value)) as a value
110:     # df_dict contains unique terms in all articles as keys and a
111:     ↪ number of document in which this term is occured as
112:     ↪ value
113:
114:     for article_index in df1.index.tolist():
115:         calc_tf_df(article_index)
116:
117:     # 1.1.2. 3) populate dictionary tfidf_dict that has article id
118:     ↪ as a key and a dictionary (term (key) - tfidf (value))
119:     ↪ as a value
120:
121:     for article_index in df1.index.tolist():
122:         calc_tfidf(article_index)
123:
124:     # 1.1.2. 5) Cosine Similarity for articles Germany, Europe
125:
126:     cosine_sim = calculateCosineSimilarity(tfidf_dict[getId('
127:         ↪ Germany')], tfidf_dict[getId('Europe')])
```

```
C:\Users\321\Documents\Uni\Uni\webse intro\store2>python "Assignment8_tasks_1-1_and_1-2.py"
Jaccard Similarity coefficient: 0.046031746031746035
Cosine Similarity coefficient: 0.0047778531080394345
Jaccard Similarity coefficient for outlinks: 0.27307692307692305
```

Figure 1: Similarity coefficients for articles Germany, Europe

```
124:     print ('Cosine Similarity coefficient: ', cosine_sim)
125:
126:
127:     # 1.2. Cosine Similarity for out_links of the articles Germany
    ↪     , Europe
128:
129:     Jaccard_sim_graph = calcJaccardSimilarity(set(df2.loc[getId('
    ↪     Germany')].out_links), set(df2.loc[getId('Europe')].
    ↪     out_links))
130:     print ('Jaccard Similarity coefficient for outlinks: ',
    ↪     Jaccard_sim_graph)
```

1.3 How similar have our similarities been? (10 Points)

Having implemented these three models and similarity measures (text with Jaccard, text with cosine, graph with Jaccard) our goal is to understand and quantify what is going on if they are used in the wild. Therefore in this and the next subtask we want to try to give an answer to the following questions.

- Will the most similar articles to a certain article always be the same independent which model we use?
- How similar are these measures to each other? How can you statistically compare them?

Assume you could use the similarity measure to compute the top k most similar articles for each article in the document collection. We want to analyze how different the rankings for these various models are.

Do some research to find a statistical measure (either from the lectures of part 2 or by doing a web search and coming up with something that we haven't discussed yet) that could be used best to compare various rankings for the same object.

Explain in a short text which measure you would use in such an experiment and why you think it is useful for our task.

Answer The most similar articles will probably not be the same, as the *Jaccard Coefficient* is only taking into account that articles share the same text. They do not perform any clean up on the data, therefore, the range of similarities is influenced by stopwords a

lot. Another issue is that count is not taken into the measure, as it is done in *TF-IDF*. The *Jaccard Coefficient on graph* similarity is highly semantic, as it spans across explicitly given links, strongly connected to the represented knowledge; it might be close to the *TF-IDF* approach but that's just a feeling.

Regarding the second question, we want to find out how well the results of the similarity functions fit. A statistical method to find this is the *Residual Sum of Squares* or *RSS*. It could be calculated pairwise on the similarity measures S_1 and S_2 , where the function will be $f(a, b)$ the similarity of articles a, b as rated by S_1 and the data-points $x_{a,b}$ the similarity as rated by S_2 . The sum would iterate over all $a \in D, b \in D$ where D is the set of all articles.

1.4 Implement the measure and do the experiment (10 Points)

After you came up with a measure you will most likely run into another problem when you plan to do the experiment.

Since runtime is an issue we cannot compute the similarity for all pairs of articles. Tell us:

1. How many similarity computations would have to be done if you wished to do so?
2. How much time would roughly be consumed to do all of these computations?

A better strategy might be to select a couple of articles for which you could compute your measure. One strategy would be to select the 100 longest articles. Another strategy might be to randomly select 100 articles from our corpus.

Computer your three similarity measures and evaluate them for these two strategies of selecting test data. Present your results. Will the results depend on the method for selecting articles? What are your findings?

Listing 2 Tasks 1.4

```
1: from collections import Counter
2: from functools import partial, reduce
3: from itertools import combinations
4: import pandas as pd
5: import math
6: import re
7: import time
8:
9:
10: def compose(*functions):
11:     return lambda args: reduce(lambda acc, f: f(acc), reversed(
12:         ↪ functions), args)
13:
14: def split(text):
```



```
15:     # split text by words
16:     return re.split(r'\W+', text)
17:
18:
19: def calcJaccardSimilarity(wordset1, wordset2):
20:     intersection = len(wordset1 & wordset2)
21:     union = len(wordset1 | wordset2)
22:     # return 0 if set union 0
23:     return 0 if union is 0 else intersection / union
24:
25:
26: def calcTermFrequency(text):
27:     return Counter(split(text))
28:
29:
30: def calcTfIdf(d_term_frequencies, doc_count, wordset,
    ↪ term_frequencies):
31:     return {term: term_frequencies[term] * math.log(doc_count /
    ↪ d_term_frequencies[term]) for term in wordset}
32:
33:
34: def flatten(l):
35:     # http://stackoverflow.com/questions/952914/making-a-flat-list
    ↪ -out-of-list-of-lists-in-python
36:     return [item for sublist in l for item in sublist]
37:
38:
39: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
40:     # gather all terms
41:     terms = set(tfIdfDict1.keys()) | set(tfIdfDict2.keys())
42:     # calculate sparse matrix multiplication
43:     numerator = sum(tfIdfDict1.get(term, 0) * tfIdfDict2.get(term,
    ↪ 0) for term in terms)
44:     # calculate vector length
45:     denominator = math.sqrt(sum([v**2 for v in tfIdfDict1.values()
    ↪ ])) * math.sqrt(
46:         sum([v**2 for v in tfIdfDict2.values()])))
47:     return numerator / denominator
48:
49:
50: with pd.HDFStore('store2.h5', 'r', encoding='utf-8') as store:
51:     # define sample length
52:     n = 100
53:
54:     df2 = store['df2']
55:     # track rows by name
56:     df2 = df2.set_index('name')
57:     # map list of out links to set
58:     df2['out_links_set'] = df2['out_links'].map(set)
```

```
59:
60:     df1 = store['df1']
61:     # count articles
62:     df1_article_count = len(df1)
63:     # map text to list of words and then to word set
64:     df1['word_set'] = df1['text'].map(compose(set, split))
65:     # calculate overall term frequencies
66:     document_term_frequencies = Counter(flatten(df1['word_set']))
67:     # map text to term frequency
68:     df1['term_frequency'] = df1['text'].map(calcTermFrequency)
69:     # map wordset with term frequency to tf idf
70:     df1['tf_idf'] = list(map(lambda t: calcTfIdf(
71:         ↪ document_term_frequencies, df1_article_count, *t),
72:             zip(df1['word_set'], df1['
73:                 ↪ term_frequency'])))
74:
75:     # map text to text length
76:     df1['text_length'] = df1['text'].map(len)
77:     # sort by text length descending and track rows by name
78:     df1 = df1.sort_values('text_length', ascending=False).
79:         ↪ set_index('name')
80:
81:     # initialize counter for computations
82:     df1_count = 0
83:     # calculate combinations to calculate
84:     df1_combination_count = n * (n - 1) / 2
85:     # n first index
86:     df1_sets_with_index = df1.index.values[:n + 1]
87:     # initialize similarity dict
88:     similarities = {}
89:     # initialize calculation start time
90:     start_time = time.time()
91:     # iterate over pairwise combinations of word sets with index
92:     for left_index, right_index in combinations(
93:         ↪ df1_sets_with_index, 2):
94:         # calculate jaccard similarity for text
95:         text_jaccard = calcJaccardSimilarity(df1['word_set'].loc[
96:             ↪ left_index],
97:
98:
99:             df1['word_set'].loc[
100:                 ↪ right_index])
101:
102:         # calculate cosine similarity for terms
103:         cosine_similarity = calculateCosineSimilarity(df1['tf_idf'
104:             ↪ ].loc[left_index],
105:
106:
107:             df1['tf_idf'
108:                 ↪ ].loc[
109:                     ↪ right_index
110:                     ↪ ])
111:
112:         # calculate jaccard similarity for out links
113:         link_jaccard = calcJaccardSimilarity(df2['out_links_set'].
114:             ↪ loc[left_index],
115:
116:
117:             df2['out_links_set'].
118:                 ↪ loc[right_index])
```

```

                                                    ↪ loc[right_index
                                                    ↪ ])
97:         # store similarity pairs in similarity dictionary
98:         similarities.update({(left_index, right_index): [
           ↪ text_jaccard, link_jaccard, cosine_similarity]})
99:         similarities.update({(right_index, left_index): [
           ↪ text_jaccard, link_jaccard, cosine_similarity]})
100:        # increase computation counter
101:        df1_count += 1
102:        # print for feedback for every 100 computations
103:        if df1_count % 100 is 0:
104:            print(
105:                'calculated %i from %i (%f %%)' % (
106:                    df1_count, df1_combination_count, 100 *
           ↪ df1_count / df1_combination_count))
107:        print("first %s pairwise similarity calculations: %s seconds"
           ↪ % (n, time.time() - start_time))

```

1. $|computations| = 3 \cdot \frac{|articles| \cdot (|articles| - 1)}{2} = 3 \cdot \frac{27493 \cdot (27493 - 1)}{2} = 1133756334$
2. $\frac{|computations|}{|computations_{first100}|} \cdot time_{first100} = \frac{3.377918778}{3.4950} \cdot 3.2seconds = 2d19h$

1.5 Hints:

1. In order to access the data in python, you can use the following piece of code:

```

import pandas as pd
store = pd.HDFStore('store.h5')
df1=store['df1']
df2=store['df2']

```

2. Variables df1 and df2 are pandas DataFrames which is tabular data structure. df1 consists of article's texts, df2 represents links from Simple English Wikipedia articles. Variables have the following columns:
 - "name" is a name of Simple English Wikipedia article,
 - "text" is a full text of the article "name",
 - "out_links" is a list of article names where the article "name" links to.
3. In general you might want to store the counted results in a file before you do the similarity computations and all the research for the third and fourth subtask. Doing all this counting and preparation might already take quite some runtime.
4. When computing the sparse tf-idf vectors you might already want to store the euclidean length of the vectors. otherwise you might discover runtime issues when computing the length again for each similarity computation.

5. Finding the top similar articles for a given article id requires you to compute the similarity of the given article with comparison to all the other known articles and extract the top 5 similarities. Bare in mind that these are quite a lot of similarity computations! You can expect a runtime to find the top similar articles with respect to one of the methods to be up to 10 seconds. If it takes significant longer then you probably have not used the best data structures handle your data.
6. **Even though many third party libraries exist to do this task with even less computational effort those libraries must not be used.**
7. You can find more information about basic usage of pandas DataFrame in [pandas documentation](#).
8. Here are some usefull examples of operations with DataFrame:

```
import pandas as pd

store = pd.HDFStore('store.h5')#read .h5 file
df1=store['df1']
df2=store['df2']
print df1['name'] # select column "name"
print df1.name # select column "name"
print df1.loc[9] #select row with id equals 9
print df1[5:10] #select rows from 6th to 9th (first row is 0)
print df2.loc[0].out_links #select outlinks of article with id=0

#show all columns where column "name" equals "Germany"
print df2[df2.name=="Germany"]

#show column out_links for rows where name is from list ["Germany","Austria"]
print df2[df2.name.isin(["Germany","Austria"])]out_links

#show all columns where column "text" contains word "good"
print df1[df1.text.str.contains("good")]

#add word "city" to the beginning of each text value
#(IT IS ONLY SHOWS RESULT OF OPERATION, see explanation below!)
print df1.text.apply(lambda x: "city "+x)

#make all text lower case and split text by spaces
df1[["text"]]=df1.text.str.lower().str.split()

def do_sth(x):
    #here is your function
    #
    #
```

```
    return x

#apply do_sth function to text column
#It will not change column itself, it will only show the result of application
print df1.text.apply(do_sth())

#you always have to assign result to , e.g., column,
#in order it affects your data.
#Some functions indeed can change the DataFrame by
#applying them with argument inplace=True
df1[["text"]]=df1.text.apply(do_sth())

#delete column "text"
df1.drop('text', axis=1, inplace=True)
```

Important Notes

Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment8/` in your group's repository.
- The name of the group and the names of all participating students must be listed on each submission.
- Solution format: all solutions as *one* PDF document. Programming code has to be submitted as Python code to the github repository. Upload *all* `.py` files of your program! Use **UTF-8** as the file encoding. *Other encodings will not be taken into account!*
- Check that your code compiles without errors.
- Make sure your code is formatted to be easy to read.
 - Make sure you code has consistent **indentation**.
 - Make sure you comment and document your code adequately in English.
 - Choose consistent and intuitive names for your identifiers.
- Do *not* use any accents, spaces or special characters in your filenames.

Acknowledgment

This latex template was created by Lukas Schmelzeisen for the tutorials of "Web Information Retrieval".

LA_TE_X

Currently the code can only be build using **LuaLaTeX**, so make sure you have that installed. If on Overleaf, there's an error, go to settings and change the **L**A_TE_Xengine to **LuaLaTeX**.