

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.

Group name: uniform

Group members: Pradip Giri, Jalak Arvind Kumar Pansuriya, Madhu Rakhal Magar

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 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.

Listing 1 Helper file

```
1: """ Helper file for assignment 8 """
2: import re
3: import math
4: from collections import Counter
5: import numpy as np
6:
7: def split_with_space(text):
8:     """ Split given text with spaces """
9:     text = re.sub(r'[^A-Za-z0-9\s]+' , '', text.lower())
10:    return re.split(r'\s+', text)
11:
12:
13: def word_set_from_text(text):
14:     """ Builds set of words from the text """
15:     # split the text using space
16:     words = split_with_space(text)
17:     tkon = set()
18:     for word in words:
19:         tkon.add(word)
20:     return tkon
21:
22:
23: def calc_jaccard_similarity(wordset1, wordset2):
24:     """Calculates jaccard similarity from given sets"""
25:     intersection = wordset1.intersection(wordset2)
26:     union = wordset1.union(wordset2)
27:     return len(intersection) / len(union)
```

```
28:
29:
30: def cal_term_freq(text_document):
31:     """ Calculate frequency of each word in given document """
32:     words = split_with_space(text_document)
33:     gen_freq = dict(Counter(words))
34:     return gen_freq
35:
36:
37: def tfidf(term_freq, total_doc, match_doc):
38:     """ Calculates term frequency inverse document frequency """
39:     return term_freq * math.log((total_doc/match_doc), 10)
40:
41:
42: def calculate_cosine_similarity(tfid_dict_1, tfid_dict_2):
43:     """ Calculates cosines similarity """
44:     product = 0
45:     for term in tfid_dict_1:
46:         if term in tfid_dict_2.keys():
47:             product = product + (tfid_dict_1[term] * tfid_dict_2[term])
48:     ecu_dic_tf1 = np.sum(np.square(list(tfid_dict_1.values())))
49:     ecu_dic_tf2 = np.sum(np.square(list(tfid_dict_2.values())))
50:     return product / (ecu_dic_tf1 * ecu_dic_tf2)
```

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**.

Answers

1. Build the word set

```
1: def gen_word_set(data_frame):
2:     """ Generate word set for the given data frame """
3:     results = []
4:     for _, row in data_frame.iterrows():
5:         results.append(word_set_from_text(row.text))
6:     data_frame['word_set'] = results
```

2. Implement function calcjaccardsimilarity

```
1: def calc_jaccard_similarity(wordset1, wordset2):  
2:     """Calculates jaccard similarity from given sets"""  
3:     intersection = wordset1.intersection(wordset2)  
4:     union = wordset1.union(wordset2)  
5:     return len(intersection) / len(union)
```

3. Compute the result

```
1: ger_wordset = df1[df1.name == "German"].iloc[0].word_set  
2: eu_wordset = df1[df1.name == "Europe"].iloc[0].word_set  
3: jacc_cof_eu_germany = calc_jaccard_similarity(eu_wordset, ger_wordset)
```

Jaccard Coefficient in Document 0.11428571428571428

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**.

Answers

1. Count the term frequency

```
1: def count_freq_each_term_each_article(data_frame):  
2:     """ Calculates the frequency of each term in each article """  
3:     results = []  
4:     for _, row in data_frame.iterrows():  
5:         results.append(cal_term_freq(row.text))  
6:     data_frame['term_freq'] = results  
7:  
8: def document_frequency_of_all_term(input_terms, documents):  
9:     """ Calculate the document frequency of each term in given documents """  
10:    results = dict()  
11:    for term in input_terms:
```

```
12:         print("Processing ", term)
13:         result = cal_doc_freq_of_term(term, documents)
14:         results[term] = result[term]
15:     return results
```

2. Count the document frequencies

```
1: def document_frequency_of_all_term(input_terms, documents):
2:     """ Calculate the document frequency of each term in given documents """
3:     results = dict()
4:     for term in input_terms:
5:         print("Processing ", term)
6:         result = cal_doc_freq_of_term(term, documents)
7:         results[term] = result[term]
8:     return results
```

3. For each article id provide a dictionary of terms occurring in the article together with their tf-idf scores as the corresponding values.

```
1: def cal_tfidf_each_term(data_frame, document_term_freq):
2:     """ Calculates the tfidf of each term """
3:     results = []
4:     doc_length = len(data_frame.term_freq)
5:     for _, row in data_frame.iterrows():
6:         term_set = row.term_freq
7:         res = dict()
8:         for term in term_set:
9:             res[term] = tfidf(term_set[term], doc_length, document_term_freq[term])
10:        results.append(res)
11:    data_frame['tf_idf'] = results
```

4. Implement Cosine Similarity functions taking two term frequency dictionary

```
1: def calculate_cosine_similarity(tfid_dict_1, tfid_dict_2):
2:     """ Calculates cosines similarity """
3:     product = 0
4:     for term in tfid_dict_1:
5:         if term in tfid_dict_2.keys():
6:             product = product + (tfid_dict_1[term] * tfid_dict_2[term])
7:     ecu_dic_tf1 = np.sum(np.square(list(tfid_dict_1.values())))
8:     ecu_dic_tf2 = np.sum(np.square(list(tfid_dict_2.values())))
9:     return product / (ecu_dic_tf1 * ecu_dic_tf2)
```

5. Compute the result for the articles Germany and Europe.

```
1: cal_tfidf_each_term(df1, term_document_freq)
2: ger_tf_idf = df1[df1.name == "German"].iloc[0].tf_idf
3: eu_tf_idf = df1[df1.name == "Europe"].iloc[0].tf_idf
4: val = calculate_cosine_similarity(ger_tf_idf, eu_tf_idf)
```

Cosine 0.000245755091545

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

```
1: ger_outlink = df2[df2.name == "German"].iloc[0].out_links
2: eu_outlink = df2[df2.name == "Europe"].iloc[0].out_links
3: jacc_cof_eu_germany_links = calc_jaccard_similarity(set(ger_outlink), set(eu_outlink))
4: print("Jaccard Coefficient in Links", jacc_cof_eu_germany_links)
```

Jaccard Coefficient in Links 0.0202020202020204

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.

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.

Compute 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 assignment8.py

```
1: # pylint: disable-msg=C0103
2: """ Assignment Number 8 """
3: from pathlib import Path
4: import pandas as pd
5: from helper import word_set_from_text
6: from helper import calc_jaccard_similarity
7: from helper import cal_term_freq
8: from helper import tfidf
9: from helper import calculate_cosine_similarity
10:
11:
12: gen_set_of_docs = None
13:
14: def gen_word_set(data_frame):
15:     """ Generate word set for the given data frame """
16:     results = []
17:     for _, row in data_frame.iterrows():
18:         results.append(word_set_from_text(row.text))
19:     data_frame['word_set'] = results
20:
21:
22: def cal_doc_freq_of_term(term, list_of_docs):
23:     """ Calculates the document frequency of given term """
24:     # create set from the doc
25:     # this code seems too slow
26:     if gen_set_of_docs is None:
27:         gen_set_of_docs = [word_set_from_text(doc) for doc in list_of_docs]
28:     result = {term: 0}
29:     for gen_set in gen_set_of_docs:
```

```
30:         if term in gen_set:
31:             result[term] += 1
32:     return result
33:
34:
35: # For question 1.1.2
36: def count_freq_each_term_each_article(data_frame):
37:     """ Calculates the frequency of each term in each article """
38:     results = []
39:     for _, row in data_frame.iterrows():
40:         results.append(cal_term_freq(row.text))
41:     data_frame['term_freq'] = results
42:
43:
44: def document_frequency_of_all_term(input_terms, documents):
45:     """ Calculate the document frequency of each term in given documents """
46:     results = dict()
47:     for term in input_terms:
48:         print("Processing ", term)
49:         result = cal_doc_freq_of_term(term, documents)
50:         results[term] = result[term]
51:     return results
52:
53:
54: def cal_tfidf_each_term(data_frame, document_term_freq):
55:     """ Calculates the tfidf of each term """
56:     results = []
57:     doc_length = len(data_frame.term_freq)
58:     for _, row in data_frame.iterrows():
59:         term_set = row.term_freq
60:         res = dict()
61:         for term in term_set:
62:             res[term] = tfidf(term_set[term], doc_length, document_term_freq[term])
63:         results.append(res)
64:     data_frame['tf_idf'] = results
65:
66:
67: def compare_cosine_for_random_100_column(data_frame, cosine_similarity):
68:     """ calculate and returns """
69:     # lets select 100 ids
70:     sample_data = data_frame.sample(100)
71:
72:     # calculating cosines
73:     first_row_ids = sample_data.iloc[0].tf_idf
74:
75:     results = []
76:     for _, row in sample_data.iterrows():
77:         res = cosine_similarity(first_row_ids, row.tf_idf)
78:         results.append(res)
```



```
79:     sample_data['cosine'] = results
80:     sample_data.sort_values(['cosine'], ascending=[False], inplace=True)
81:     print(sample_data.head())
82:
83:
84:
85: def main():
86:     """ Entry point of the program """
87:     print("Processing Starts")
88:     store = pd.HDFStore('store2.h5')
89:     df1 = store['df1']
90:     df2 = store['df2']
91:     gen_word_set(df1)
92:     # calculating term frequency
93:     count_freq_each_term_each_article(df1)
94:
95:     # Question 1.1.1
96:     ger_wordset = df1[df1.name == "German"].iloc[0].word_set
97:     eu_wordset = df1[df1.name == "Europe"].iloc[0].word_set
98:     jacc_cof_eu_germany = calc_jaccard_similarity(eu_wordset, ger_wordset)
99:     print("Jaccard Coefficient in Document", jacc_cof_eu_germany)
100:
101:     # Question 1.2
102:     # 1.2 Similarity of Graphs
103:     ger_outlink = df2[df2.name == "German"].iloc[0].out_links
104:     eu_outlink = df2[df2.name == "Europe"].iloc[0].out_links
105:     jacc_cof_eu_germany_links = calc_jaccard_similarity(set(ger_outlink), set(eu_outlink))
106:     print("Jaccard Coefficient in Links", jacc_cof_eu_germany_links)
107:
108:
109:     all_terms = set()
110:     for word_set in df1['word_set'].values:
111:         all_terms |= word_set
112:
113:     all_docs = []
114:     for doc in df1['text']:
115:         all_docs.append(doc)
116:
117:     doc_freq = Path('./doc_freq.txt')
118:     term_document_freq = None
119:
120:     if doc_freq.is_file():
121:         with open('doc_freq.txt', 'r+') as f:
122:             term_document_freq = eval(f.read())
123:     else:
124:         term_document_freq = document_frequency_of_all_term(all_terms, all_docs)
125:         with open('doc_freq.txt', 'w+') as f:
126:             f.write(str(term_document_freq))
127:
```

```
128:     # Calculate ifidf scores
129:     cal_tfidf_each_term(df1, term_document_freq)
130:     ger_tf_idf = df1[df1.name == "German"].iloc[0].tf_idf
131:     eu_tf_idf = df1[df1.name == "Europe"].iloc[0].tf_idf
132:
133:     print("Cosine", calculate_cosine_similarity(ger_tf_idf, eu_tf_idf))
134:     compare_cosine_for_random_100_column(df1, calculate_cosine_similarity)
135:
136:
137:
138: if __name__ == "__main__":
139:     main()
```

```
name \
9932      Humanism
17694  Johann_Strauss_I_db43
13930      Ballpoint_pen
14780  Dog_Day_Afternoon_afb2
4659      Capetian_dynasty

text \
9932  **Humanism** is a series of philosophies that ...
17694  **Johann Strauss I** (born Vienna, March 14 18...
13930  A **ballpoint pen** (also called a **biro** \-...
14780  _**Dog Day Afternoon**_ is a 1975 movie. It wa...
4659  The **Capetian Dynasty** includes any of the d...

word_set \
9932  {one, and, do, people, all, if, series, someth...
17694  {court, usually, and, was, while, jnr, radetzk...
13930  {them, was, entirely, good, pronounced, people...
14780  {life, dog, and, was, by, pierson, chris, day,...
4659  {years, and, both, duke, carlos, spain, luxemb...

term_freq \
9932  {'that': 1, 'one': 1, 'affirm': 1, 'to': 1, 'a...
17694  {'usually': 2, 'was': 3, 'march': 1, 'monarchy...
13930  {'soon': 1, 'usually': 1, 'them': 1, 'ones': 1...
14780  {'life': 1, 'dog': 1, 'and': 2, 'was': 1, 'chr...
4659  {'years': 1, 'and': 2, 'both': 1, 'duke': 1, '...

tf_idf      cosine
9932  {'that': 0.42279518204517086, 'one': 0.6485780... 0.011414
17694  {'court': 1.869911404050953, 'usually': 2.0426... 0.000005
13930  {'soon': 1.6841730472709278, 'usually': 1.0213... 0.000019
14780  {'is': 0.08613876745201976, 'john': 1.54165822... 0.000134
4659  {'1589': 3.536195326674055, 'that': 0.42279518... 0.000023
```

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 useful 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".

LaTeX

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 **LaTeX**engine to **LuaLaTeX**.