# 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: Oscar

# 1 Similarity - (40 Points)

This assignment will have one exercise which is dived 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 aply 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 coefficent of two word sets and return the value.

3. Compute the result for the articles `Germany` and `Europe`.

```
 1: #This is a special case where im calculating jacaard sim for
 2: #only Germany and Europe. I know it says in the assignmnt to
 3: #build the wordset for all documents, but i dont need it here.
 4: #So i haven't built it. However, I put it in 1.3 because it's
 5: #required.
 6:
 7: import pandas as pd
 8: import re
 9: import numpy as np
10: #Jaccard-Similarity on sets
11: store = pd.HDFStore('store.h5')
12: df1=store['df1']
13: df1['text']=df1['text'].str.lower()
14: df1.text.replace('', np.nan, inplace=True)
15: df1.dropna(subset=['text'], inplace=True)
16: df1.text=df1.text.apply(lambda x: re.findall(r'[0-9a-zA-Z]+', x))
17: df1.text=df1.text.apply(lambda x: ' '.join(map(str, x)))
```

```
18:
19: def calcJaccardSimilarity(wordset1, wordset2):
20:         #length of intersection set
21:         JK1= len(wordset1 & wordset2)
22:         #length of union set
23:         JK2 = len(wordset1 | wordset2)
24:         #Jaccard Formula
25:         JK = JK1/float(JK2)
26:         return JK
27:
28: #get id-s for articles named Germany and Europe
29: index1= df1[df1.name=="Germany"].index
30: index1=index1[0]
31:
32: index2= df1[df1.name=="Europe"].index
33: index2=index2[0]
34:
35: #building sets for each article
36: set1= set(df1.text[index1].split())
37: set2=set(df1.text[index2].split())
38:
39: print(calcJaccardSimilarity(set1,set2),"\n")
```

### 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 occuring 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: #This is also a special case solution designed for finding
 2: #cosine similarity between germany and europe where im
 3: #considering my text corpa as only the words involved,
 4: #which are words in the texts of germany and europe article.
 5: #I know it's not theoritcally correct as im supposed to cal-
 6: #culate it for all df1 but it works here because it's too
 7: #many zeros that we don't need. However, in 1.3 I've provided
 8: #a full solution covering all pd1 docs
 9: import pandas as pd
10: import numpy as np
11: import re
```

```
12: from collections import defaultdict
13: from math import log
14: from math import sqrt
15: store = pd.HDFStore('store.h5')
16: df1=store['df1']
17: df1['text']=df1['text'].str.lower()
18: df1.text.replace('', np.nan, inplace=True)
19: df1.dropna(subset=['text'], inplace=True)
20: df1.text=df1.text.apply(lambda x: re.findall(r'[0-9a-zA-Z]+', x))
21: df1.text=df1.text.apply(lambda x: ' '.join(map(str, x)))
22:
23:
24: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
25:         tfIdfDict1_mag=0
26:         tfIdfDict2_mag=0
27:         vec_mult=0
28:         for key, value in tfIdfDict1.items():
29:                 value2=tfIdfDict2[key]
30:                 type(value2)
31:                 tfIdfDict1_mag+=(value*value)
32:                 tfIdfDict2_mag+=(value2*value2)
33:
34:                 vec_mult+=(value*value2)
35:         tfIdfDict1_mag=sqrt(tfIdfDict1_mag)
36:         tfIdfDict2_mag=sqrt(tfIdfDict2_mag)
37:
38:         return vec_mult/float(tfIdfDict1_mag*tfIdfDict2_mag)
39:
40: #get id-s for articles named Germany and Europe
41: index1= df1[df1.name=="Germany"].index
42: index1=index1[0]
43:
44: index2= df1[df1.name=="Europe"].index
45: index2=index2[0]
46:
47: tf_idf1 = defaultdict(int)
48: tf_idf2 = defaultdict(int)
49:
50: docfrq=defaultdict(int)
51:
52: #preprocesing both tfidf docs data
53: for word in df1.text[index1].split():
54:         docfrq[word] = 0
55:         tf_idf1[word] =0
56:         tf_idf2[word] =0
57:
58:
59: for word in df1.text[index2].split():
60:         docfrq[word] = 0
```

```
61:             tf_idf1[word] =0
62:             tf_idf2[word] =0
63:
64:
65: #calculating tf and idf
66: for index, row in df1.iterrows():
67:         flag=False
68:         l=[]
69:         for word in row.text.split():
70:                 if(word in docfrq  and word not in l):
71:                         docfrq[word] +=1
72:                         l.append(word)
73:
74: N =len(df1)
75:
76: for word in df1.text[index1].split():
77:         tf_idf1[word] +=1
78:
79: for word in df1.text[index2].split():
80:         tf_idf2[word] +=1
81:
82: for key, value in docfrq.items():
83:         if(tf_idf1[key]!=0):
84:                 tf_idf1[key]=tf_idf1[key]*log(N/float(value))
85:         if(tf_idf2[key]!=0):
86:                 tf_idf2[key]=tf_idf2[key]*log(N/float(value))
87:
88: print(calculateCosineSimilarity(tf_idf1,tf_idf2))
```

## 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 aply 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`.

```
 1: #same as jacaard in 1.111 but for outlinks
 2:
 3: import pandas as pd
 4: import re
 5: import numpy as np
 6:
 7: #Jaccard-Similarity on sets
 8: store = pd.HDFStore('store.h5')
 9: df2=store['df2']
10:
```

```
11: def calcJaccardSimilarity(wordset1, wordset2):
12:        #length of intersection set
13:        JK1= len(wordset1 & wordset2)
14:        #length of union set
15:        JK2 = len(wordset1 | wordset2)
16:        #Jaccard Formula
17:        JK = JK1/float(JK2)
18:        return JK
19:
20: #get id-s for articles named Germany and Europe
21: index1= df2[df2.name=="Germany"].index
22: index1=index1[0]
23:
24: index2= df2[df2.name=="Europe"].index
25: index2=index2[0]
26:
27: #building sets for each article
28: set1= set(df2.out_links[index1])
29: set2=set(df2.out_links[index2])
30:
31: print(calcJaccardSimilarity(set1,set2),"\n")
```

## 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 answere 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 usefull for our task.

Answer:

1. Concerning the first question whether the most similar articles will be the same or not, the answer as i think is no because the similarity measure can sometime take core different approaches in calculating similarity. For example in the case of jacaard, the method only looks on the two sets at hand without a broader look to the whole corpus. However, cosine similarity takes that in mind with global weights in addition to the two docs being measured. This could leave to small similarities that when piled up can diverse the top ranked results of similarity measures.

2. Concerning question 2, i thought about it and wanted to get a new solution to it. Most 2 important factors in comparison are in my opinion similar docs and their ranking. The solution wasn't so new but i took it from other courses we took in ws. I decided to calculate levenstein distance between ranked results for each similarity approach. The inspiration behind this is looking a the lists of resulting documents as different words of same length and since we want to keep similar docs and their ranking in mind, words levenstein does exactly that. PS: I found that cosine similarity is closer to both jacaard text and graph more than they're closer to each other

```
 1: #This is the most elegant form for all other questions
 2: #in terms of modeling similarity as i do it here for
 3: #all documents given each similarity. In cosine sim
 4: #i followed a different and more realistic approach
 5: #by using only non zero entries in tfidf vectors thus
 6: #you will see that my calculateCosineSimilarity function
 7: #is different from here and for me it's the most generic
 8: #one that can apply on any case
 9: import pandas as pd
10: import re
11: import numpy as np
12: from math import log
13: from math import sqrt
14:
15:
16: #Jaccard-Similarity on sets
17: store = pd.HDFStore('store.h5')
18: df1=store['df1']
19: df1['text']=df1['text'].str.lower()
20: df1.text.replace('', np.nan, inplace=True)
21: df1.dropna(subset=['text'], inplace=True)
22: df1.text=df1.text.apply(lambda x: re.findall(r'[0-9a-zA-Z]+', x))
23: df1.text=df1.text.apply(lambda x: ' '.join(map(str, x)))
24: df2=store['df2']
25:
26:
27: #functions
28:
29: def calcJaccardSimilarity(wordset1, wordset2):
30:         #length of intersection set
```

```
31:        JK1= len(wordset1 & wordset2)
32:        #length of union set
33:        JK2 = len(wordset1 | wordset2)
34:        #Jaccard Formula
35:        JK = JK1/float(JK2)
36:        return JK
37:
38: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
39:        tfIdfDict1_mag=tfIdfDict1.values()
40:        tfIdfDict1_mag=[x**2 for x in tfIdfDict1_mag]
41:        tfIdfDict1_mag= sqrt(sum(tfIdfDict1_mag))
42:
43:        tfIdfDict2_mag=tfIdfDict2.values()
44:        tfIdfDict2_mag=[x**2 for x in tfIdfDict2_mag]
45:        tfIdfDict2_mag= sqrt(sum(tfIdfDict2_mag))
46:
47:        vec_mult=0
48:
49:        set1=set(tfIdfDict1.keys())
50:        set2=set(tfIdfDict2.keys())
51:        common_words=list(set1&set2)
52:
53:        for i in range(len(common_words)):
54:                vec_mult+=tfIdfDict1[common_words[i]]*tfIdfDict2[common_words[i]]
55:
56:        return vec_mult/float(tfIdfDict1_mag*tfIdfDict2_mag)
57:
58: def smart_push(l,item,n):
59:        if(len(l)<n):
60:                l.append(item)
61:        else:
62:                if(item[0]>l[0][0]):
63:                        l.remove(l[0])
64:                        l.append(item)
65:        l=sorted(l, key=lambda x: x[0])
66:        return l
67:
68: def levenshteinDistance(s1, s2):
69:     if len(s1) > len(s2):
70:         s1, s2 = s2, s1
71:
72:     distances = range(len(s1) + 1)
73:     for i2, c2 in enumerate(s2):
74:         distances_ = [i2+1]
75:         for i1, c1 in enumerate(s1):
76:             if(c1 == c2):
77:                 distances_.append(distances[i1])
78:             else:
79:                 distances_.append(1 + min((distances[i1], distances[i1 + 1], dist
```

```
 80:             distances = distances_
 81:         return distances[-1]
 82:
 83:
 84: target_doc_index= df1[df1.name=="Germany"].index
 85: target_doc_index=target_doc_index[0]
 86: target_doc_text = df1.text[target_doc_index]
 87: target_doc_outlinks = df2.out_links[target_doc_index]
 88:
 89: #cosine similarity pre processing
 90: N=len(df1)
 91:
 92: docfrq=dict()
 93: target_doc_tfidf=dict()
 94: docs_tfidf={}
 95:
 96: print('started preprocessing')
 97: #filled docfrqs and target_tfidt
 98: for index,row in df1.iterrows():
 99:         docs_tfidf[index]=dict()
100:         l=[]
101:         for word in row.text.split():
102:                 if(word in docfrq and word not in l):
103:                         docfrq[word]+=1
104:                 else:
105:                         if(word not in docfrq):
106:                                 docfrq[word]=1
107:                 l.append(word)
108:
109:
110:                 if(index==target_doc_index):
111:                         if(word in target_doc_tfidf):
112:                                 target_doc_tfidf[word]+=1
113:                         else:
114:                                 target_doc_tfidf[word]=1
115:
116:                 if(word in docs_tfidf[index]):
117:                         docs_tfidf[index][word]+=1
118:                 else:
119:                         docs_tfidf[index][word]=1
120:         if(index%10000==0):
121:                 print(index)
122:
123: for key,value in target_doc_tfidf.items():
124:         target_doc_tfidf[key]=value*log(N/float(docfrq[key]))
125:
126:
127: text_jacard=[]
128: graph_jacard=[]
```

9

```
129: cosine=[]
130: print('preprocessing ended ... calculating similarity')
131: for index,row in df1.iterrows():
132:         #calculating cosine
133:         for key,value in docs_tfidf[index].items():
134:                 docs_tfidf[index][key]=value*log(N/float(docfrq[key]))
135:         cosine=smart_push(cosine,(calculateCosineSimilarity(target_doc_tfidf,docs_
136:
137:         #calculating text jacaard
138:         target_text_set=set(target_doc_text.split())
139:         current_text_set=set(df1.text[index].split())
140:         text_jacard= smart_push(text_jacard,(calcJaccardSimilarity(target_text_set
141:         #calculating graph jacaard
142:         target_graph_set=set(target_doc_outlinks)
143:         current_graph_set=set(df2.out_links[index])
144:         graph_jacard=smart_push(graph_jacard,(calcJaccardSimilarity(target_graph_s
145:         if(index %10000==0):
146:                 print(index)
147:
148: #getting the docs ids ranked according to each sim measure
149: text_jacard = [int(i[1]) for i in text_jacard]
150: cosine = [int(i[1]) for i in cosine]
151: graph_jacard = [int(i[1]) for i in graph_jacard]
152:
153: print(text_jacard)
154: print(cosine)
155: print(graph_jacard)
156:
157: #grouping the output to calculate levenstein distance between each measure
158: #by using an function for levenstein distance i found online. I assign
159: #each doc with a letter and this way each ranked result will be like a word
160: #and at the end we calculate levenstein distance for all
161:
162: grouped=list(set(text_jacard)|set(cosine)|set(graph_jacard))
163: reference=dict()
164: for i in range(len(grouped)):
165:         reference[grouped[i]]=chr(97+i)
166:
167: text_jacard_string=''
168: cosine_string=''
169: graph_jacard_string=''
170: for i in range(len(text_jacard)):
171:         text_jacard_string+=reference[text_jacard[i]]
172:         cosine_string+=reference[cosine[i]]
173:         graph_jacard_string+=reference[graph_jacard[i]]
174:
175:
176: print(str(levenshteinDistance(text_jacard_string,cosine_string))+' levenstein dist
177: print(str(levenshteinDistance(text_jacard_string,graph_jacard))+' levenstein dista
```

```
178: print(str(levenshteinDistance(cosine_string,graph_jacard_string))+' levenstein dis
```

## 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?
   **Answer:** 1)Calculating cosine sim o(n) n=nmbr of docs (chosen docs) 2)Saving the result and sorting them o(n)log(n) considering that the sorter i used in python is efficient or else it would be o(n pow 2)

2. How much time would roughly be consumed to do all of these computations?
   **Answer:** our complexity = o(n) + o(n)log(n) our time= o(n)time + o(n)log(n)time assuming out time is 1 second to execute anything which is far not accurate for sure. PS: n is roughly 27000 so the time = 27000 + 27000*4.43=5.43*27000=146610 seconds=2443.5 minutes=40.725 hours=1.7 days

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?

**Answer:**I found out that selecting the longest 100 docs is more informative that just randomly selecting docs, especially in cosine simulation. Longer articles bigger chances for similarity values accumulation. It's rather more simulating to reality than random docs.

PS:I made two files, one for random 100 and one for longest articles.

Random

```
 1: #Solution for 100 randomly selected docs
 2: #This solution is also a special case solution
 3: #where im considering my docs corpa to be the
 4: #selected items + germany article
 5:
 6: import pandas as pd
 7: import re
 8: import numpy as np
 9: import copy
10: from math import log
```

```
11: from math import sqrt
12:
13: #Jaccard-Similarity on sets
14: store = pd.HDFStore('store.h5')
15: df1=store['df1']
16: df1['text']=df1['text'].str.lower()
17: df1.text.replace('', np.nan, inplace=True)
18: df1.dropna(subset=['text'], inplace=True)
19: df1.text=df1.text.apply(lambda x: re.findall(r'[0-9a-zA-Z]+', x))
20: df1.text=df1.text.apply(lambda x: ' '.join(map(str, x)))
21: df2=store['df2']
22:
23:
24: #functions
25:
26: def calcJaccardSimilarity(wordset1, wordset2):
27:         #length of intersection set
28:         JK1= len(wordset1 & wordset2)
29:         #length of union set
30:         JK2 = len(wordset1 | wordset2)
31:         #Jaccard Formula
32:         JK = JK1/float(JK2)
33:         return JK
34:
35: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
36:          tfIdfDict1_mag=0
37:          tfIdfDict2_mag=0
38:          vec_mult=0
39:          for key, value in tfIdfDict1.items():
40:                  value2=tfIdfDict2[key]
41:                  type(value2)
42:                  tfIdfDict1_mag+=(value*value)
43:                  tfIdfDict2_mag+=(value2*value2)
44:
45:                  vec_mult+=(value*value2)
46:          tfIdfDict1_mag=sqrt(tfIdfDict1_mag)
47:          tfIdfDict2_mag=sqrt(tfIdfDict2_mag)
48:          return vec_mult/float(tfIdfDict1_mag*tfIdfDict2_mag)
49:
50: def levenshteinDistance(s1, s2):
51:     if len(s1) > len(s2):
52:         s1, s2 = s2, s1
53:
54:     distances = range(len(s1) + 1)
55:     for i2, c2 in enumerate(s2):
56:         distances_ = [i2+1]
57:         for i1, c1 in enumerate(s1):
58:             if(c1 == c2):
59:                 distances_.append(distances[i1])
```

```
60:              else:
61:                  distances_.append(1 + min((distances[i1], distances[i1 + 1], dista
62:          distances = distances_
63:      return distances[-1]
64: random_samples = df1.sample(n=100)
65: target_doc_index= df1[df1.name=="Germany"].index
66: target_doc_index=target_doc_index[0]
67: target_doc_text = df1.text[target_doc_index]
68: target_doc_outlinks = df2.out_links[target_doc_index]
69:
70: #cosine similarity pre processing
71: N=len(df1)
72: docfrq=dict()
73: target_doc_tfidf=dict()
74: docs_tfidf={}
75: for index,row in random_samples.iterrows():
76:         docs_tfidf[index]={}
77:         for word in row.text.split():
78:                 docfrq[word]=0
79:                 target_doc_tfidf[word]=0
80:                 docs_tfidf[index][word]=0
81:
82: for word in target_doc_text.split():
83:         if(word in docfrq):
84:                 docfrq[word]+=1
85:         else:
86:                 docfrq[word]=1
87:         if(word in target_doc_tfidf):
88:                 target_doc_tfidf[word]+=1
89:         else:
90:                 target_doc_tfidf[word]=1
91:
92: for index,row in random_samples.iterrows():
93:         docs_tfidf[index]=copy.copy(docfrq)
94:
95: for index,row in random_samples.iterrows():
96:         for word in row.text.split():
97:                 if(index==target_doc_index):
98:                         target_doc_tfidf[word]+=1
99:                 else:
100:                        docs_tfidf[index][word]+=1
101:
102:
103: for index, row in df1.iterrows():
104:         flag=False
105:         l=[]
106:         for word in row.text.split():
107:                 if(word in docfrq  and word not in l):
108:                         docfrq[word] +=1
```

```
109:                            l.append(word)
110:
111: print("Data preprocessing is over")
112:
113: text_jacard=[]
114: graph_jacard=[]
115: cosine=[]
116:
117:
118: for index, row in random_samples.iterrows():
119:          #calculating text jacaard
120:          target_text_set=set(target_doc_text.split())
121:          current_text_set=set(df1.text[index].split())
122:          text_jacard.append((calcJaccardSimilarity(target_text_set,current_text_set
123:
124:          #calculating graph jacaard
125:          target_graph_set=set(target_doc_outlinks)
126:          current_graph_set=set(df2.out_links[index])
127:          graph_jacard.append((calcJaccardSimilarity(target_graph_set,current_graph_
128:
129:          #calculating cosine sim
130:          tf_idf1=target_doc_tfidf
131:          tf_idf2=docs_tfidf[index]
132:          for key, value in docfrq.items():
133:                  if(tf_idf1[key]!=0):
134:                          tf_idf1[key]=tf_idf1[key]*log(N/float(value))
135:                  if(tf_idf2[key]!=0):
136:                          tf_idf2[key]=tf_idf2[key]*log(N/float(value))
137:          cosine.append((calculateCosineSimilarity(tf_idf1,tf_idf2),index))
138:
139:
140: text_jacard.sort(key=lambda tup: tup[0])
141: graph_jacard.sort(key=lambda tup: tup[0])
142: cosine.sort(key=lambda tup: tup[0])
143:
144: #getting the docs ids ranked according to each sim measure
145: text_jacard = [int(i[1]) for i in text_jacard]
146: cosine = [int(i[1]) for i in cosine]
147: graph_jacard = [int(i[1]) for i in graph_jacard]
148:
149: print(text_jacard)
150: print(graph_jacard)
151: print(cosine)
152:
153:
154: #grouping the output to calculate levenstein distance between each measure
155: #by using an function for levenstein distance i found online. I assign
156: #each doc with a letter and this way each ranked result will be like a word
157: #and at the end we calculate levenstein distance for all
```

14

```
158:
159: grouped=list(set(text_jacard)|set(cosine)|set(graph_jacard))
160: reference=dict()
161: for i in range(len(grouped)):
162:         reference[grouped[i]]=chr(97+i)
163:
164: text_jacard_string=''
165: cosine_string=''
166: graph_jacard_string=''
167: for i in range(len(text_jacard)):
168:         text_jacard_string+=reference[text_jacard[i]]
169:         cosine_string+=reference[cosine[i]]
170:         graph_jacard_string+=reference[graph_jacard[i]]
171:
172:
173: print(str(levenshteinDistance(text_jacard_string,cosine_string))+' levenstein dist
174: print(str(levenshteinDistance(text_jacard_string,graph_jacard))+' levenstein dista
175: print(str(levenshteinDistance(cosine_string,graph_jacard_string))+' levenstein dis
```

Longest

```
 1: #Solution for longest 100 articles
 2: #This solution is also a special case solution
 3: #where im considering my docs corpa to be the
 4: #selected items + germany article
 5:
 6: import pandas as pd
 7: import re
 8: import numpy as np
 9: import copy
10: from math import log
11: from math import sqrt
12:
13: #Jaccard-Similarity on sets
14: store = pd.HDFStore('store.h5')
15: df1=store['df1']
16: df1['text']=df1['text'].str.lower()
17: df1.text.replace('', np.nan, inplace=True)
18: df1.dropna(subset=['text'], inplace=True)
19: df1.text=df1.text.apply(lambda x: re.findall(r'[0-9a-zA-Z]+', x))
20: df1.text=df1.text.apply(lambda x: ' '.join(map(str, x)))
21: df2=store['df2']
22:
23:
24: #functions
25:
26: def calcJaccardSimilarity(wordset1, wordset2):
27:         #length of intersection set
28:         JK1= len(wordset1 & wordset2)
29:         #length of union set
```

```
30:         JK2 = len(wordset1 | wordset2)
31:         #Jaccard Formula
32:         JK = JK1/float(JK2)
33:         return JK
34:
35: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
36:         tfIdfDict1_mag=0
37:         tfIdfDict2_mag=0
38:         vec_mult=0
39:         for key, value in tfIdfDict1.items():
40:                 value2=tfIdfDict2[key]
41:                 type(value2)
42:                 tfIdfDict1_mag+=(value*value)
43:                 tfIdfDict2_mag+=(value2*value2)
44:
45:                 vec_mult+=(value*value2)
46:         tfIdfDict1_mag=sqrt(tfIdfDict1_mag)
47:         tfIdfDict2_mag=sqrt(tfIdfDict2_mag)
48:         return vec_mult/float(tfIdfDict1_mag*tfIdfDict2_mag)
49:
50: def levenshteinDistance(s1, s2):
51:     if len(s1) > len(s2):
52:         s1, s2 = s2, s1
53:
54:     distances = range(len(s1) + 1)
55:     for i2, c2 in enumerate(s2):
56:         distances_ = [i2+1]
57:         for i1, c1 in enumerate(s1):
58:             if(c1 == c2):
59:                 distances_.append(distances[i1])
60:             else:
61:                 distances_.append(1 + min((distances[i1], distances[i1 + 1], dista
62:         distances = distances_
63:     return distances[-1]
64:
65: #I add a length column to the df1 to sort it
66: lengths=[]
67: for index,row in df1.iterrows():
68:         lengths.append(len(df1.text[index]))
69: df1['length']=lengths
70: df1=df1.sort_values(by=['length'], ascending=[False])
71: samples=df1.head(100)
72:
73:
74: target_doc_index= df1[df1.name=="Germany"].index
75: target_doc_index=target_doc_index[0]
76: target_doc_text = df1.text[target_doc_index]
77: target_doc_outlinks = df2.out_links[target_doc_index]
78:
```

```
79: #cosine similarity pre processing
80: N=len(df1)
81: docfrq=dict()
82: target_doc_tfidf=dict()
83: docs_tfidf={}
84: for index,row in samples.iterrows():
85:         docs_tfidf[index]={}
86:         for word in row.text.split():
87:                 docfrq[word]=0
88:                 target_doc_tfidf[word]=0
89:                 docs_tfidf[index][word]=0
90:
91: for word in target_doc_text.split():
92:         if(word in docfrq):
93:                 docfrq[word]+=1
94:         else:
95:                 docfrq[word]=1
96:         if(word in target_doc_tfidf):
97:                 target_doc_tfidf[word]+=1
98:         else:
99:                 target_doc_tfidf[word]=1
100:
101: for index,row in samples.iterrows():
102:         docs_tfidf[index]=copy.copy(docfrq)
103:
104: for index,row in samples.iterrows():
105:         for word in row.text.split():
106:                 if(index==target_doc_index):
107:                         target_doc_tfidf[word]+=1
108:                 else:
109:                         docs_tfidf[index][word]+=1
110:
111:
112: for index, row in df1.iterrows():
113:         flag=False
114:         l=[]
115:         for word in row.text.split():
116:                 if(word in docfrq  and word not in l):
117:                         docfrq[word] +=1
118:                         l.append(word)
119:
120: print("Data preprocessing is over")
121:
122: text_jacard=[]
123: graph_jacard=[]
124: cosine=[]
125:
126:
127: for index, row in samples.iterrows():
```

```
128:            #calculating text jacaard
129:            target_text_set=set(target_doc_text.split())
130:            current_text_set=set(df1.text[index].split())
131:            text_jacard.append((calcJaccardSimilarity(target_text_set,current_text_set
132:
133:            #calculating graph jacaard
134:            target_graph_set=set(target_doc_outlinks)
135:            current_graph_set=set(df2.out_links[index])
136:            graph_jacard.append((calcJaccardSimilarity(target_graph_set,current_graph_
137:
138:            #calculating cosine sim
139:            tf_idf1=target_doc_tfidf
140:            tf_idf2=docs_tfidf[index]
141:            for key, value in docfrq.items():
142:                    if(tf_idf1[key]!=0):
143:                            tf_idf1[key]=tf_idf1[key]*log(N/float(value))
144:                    if(tf_idf2[key]!=0):
145:                            tf_idf2[key]=tf_idf2[key]*log(N/float(value))
146:            cosine.append((calculateCosineSimilarity(tf_idf1,tf_idf2),index))
147:
148:
149: text_jacard.sort(key=lambda tup: tup[0])
150: graph_jacard.sort(key=lambda tup: tup[0])
151: cosine.sort(key=lambda tup: tup[0])
152:
153: #getting the docs ids ranked according to each sim measure
154: text_jacard = [int(i[1]) for i in text_jacard]
155: cosine = [int(i[1]) for i in cosine]
156: graph_jacard = [int(i[1]) for i in graph_jacard]
157:
158: print(text_jacard)
159: print(graph_jacard)
160: print(cosine)
161:
162:
163: #grouping the output to calculate levenstein distance between each measure
164: #by using an function for levenstein distance i found online. I assign
165: #each doc with a letter and this way each ranked result will be like a word
166: #and at the end we calculate levenstein distance for all
167:
168: grouped=list(set(text_jacard)|set(cosine)|set(graph_jacard))
169: reference=dict()
170: for i in range(len(grouped)):
171:         reference[grouped[i]]=chr(97+i)
172:
173: text_jacard_string=''
174: cosine_string=''
175: graph_jacard_string=''
176: for i in range(len(text_jacard)):
```

```
177:            text_jacard_string+=reference[text_jacard[i]]
178:            cosine_string+=reference[cosine[i]]
179:            graph_jacard_string+=reference[graph_jacard[i]]
180:
181:
182: print(str(levenshteinDistance(text_jacard_string,cosine_string))+' levenstein dis
183: print(str(levenshteinDistance(text_jacard_string,graph_jacard))+' levenstein dista
184: print(str(levenshteinDistance(cosine_string,graph_jacard_string))+' levenstein di
```

## 1.5 Hints:

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

   ```python
   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 preperation might allready take quite some runtime.

4. When computing the sparse tf-idf vectors you might allready want to store the eukleadan 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:

```python
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
#Iit will not change column itself, it will only show the result of aplication
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())
```

```python
#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 LaTeXengine to `LuaLaTeX`.