

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: XXXX

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`.

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.

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: Answer to 1.1:

The Jaccard similarity of Germany and Europe is 0.04588607594936709

The Cosine similarity of Germany and Europe is 0.153587715144 **Answer to 1.2:**

The Jaccard similarity of the outlinks of Germany and Europe is 0.27307692307692305

Answer to 1.3:

The most similar articles to a certain article will be different, depending on which model is used. The reason for this is because the Jaccard Similarity only checks the occurrence of a term in two documents. The amount of occurrences, however, is not being considered. The Cosine Similarity on the other hand takes these occurrences into account and the similarity rating decreases, the more often the word occurs in a document. This would lead to different similarity rankings, because the Jaccard Similarity would rank documents with many filler terms higher than the Cosine Similarity, where the rating of those filler terms should be quite low due to the common usage.

One could measure the rankings by setting up a query to a certain subject, as for example a famous person in a certain field, and check if information about this person / subject is ranked higher than information not related to the person / subject. As an example, we take a famous musician and set up a query. The top k most similar articles should be about the person himself, his band, if any, other musicians who play the same instrument or other musicians who played with him. **Code:**

```
1: # assignment 8
2: # Andrea Mildes - mildes@uni-koblenz.de
3: # Sebastian Blei - sblei@uni-koblenz.de
4: # Johannes Kirchner - jkirchner@uni-koblenz.de
5: # Abdul Afghan - abdul.afghan@outlook.de
6:
7: import logging
8: import time
9: import pandas as pd
10: import re
11: import numpy as np
12: from numpy import zeros
13: import operator
14: np.set_printoptions(threshold=np.nan)
15:
16:
17: document_freq = {}
18: word_vector = {}
19: article_dic = {}
20: sparse_tfidf_vectors = {}
21: df1 = pd.DataFrame()
22: df2 = pd.DataFrame()
23:
```

```
24: # Configure logging and set start time
25: logging.basicConfig(filename='similarity.log', level=logging.DEBUG)
26: start_time = 0
27:
28:
29: # Read the given file into a string
30: def read_file(file):
31:     with open(file) as f:
32:         data = f.read().replace('\n', '')
33:     return data
34:
35:
36: # Write a file
37: def write_file(filename, generated_string):
38:     with open(filename, 'w') as f:
39:         f.write(generated_string)
40:
41:     t = str(round((time.time() - start_time), 2)).zfill(5)
42:     logging.info "[" + t + "]" +
43:         "Finished writing \"" + filename + "\" file. \n")
44:     return
45:
46:
47: def create_set(s):
48:     regex = re.compile("\w+")
49:     s_list = regex.findall(s)
50:     word_set = set()
51:
52:     for s in s_list:
53:         word_set.add(s.lower())
54:
55:     return word_set
56:
57:
58: def calc_term_freq(s):
59:     regex = re.compile("\w+")
60:     s_list = regex.findall(s)
61:     word_rank = {}
62:
63:     for s in s_list:
64:         s = s.lower()
65:         try:
66:             # Increase count by one
67:             word_rank[s] += 1
68:         except KeyError:
69:             # Create new dictionary entry for terms that
70:             # are not part of the dict yet
71:             word_rank[s] = 1
72:
```

```
73:     return word_rank
74:
75:
76: def calc_document_freq(s):
77:     # create set of the given string in order to get only one
78:     # occurrence of each word
79:
80:     s_set = create_set(s)
81:
82:     for s in s_set:
83:         try:
84:             # Increase count by one
85:             document_freq[s] += 1
86:         except KeyError:
87:             # Create new dictionary entry for chars that
88:             # are not part of the dict. yet
89:             document_freq[s] = 1
90:     return
91:
92:
93: def calcJaccardSimilarity(wordset1, wordset2):
94:     intersection = wordset1.intersection(wordset2)
95:     union = wordset1.union(wordset2)
96:
97:     jac = len(intersection) / len(union)
98:     return jac
99:
100:
101: # Compute the document frequency and store it in global
102: # dictionary document_freq
103: def log_doc_freq():
104:     t = str(round((time.time() - start_time), 2))
105:     logging.info "[" + t + "]" Compute document frequency ..."
106:
107:     df1.text.apply(calc_document_freq)
108:
109:     # sorted_x = sorted(document_freq.items(), key=operator.itemgetter(1))
110:     # for k, v in sorted_x:
111:     #     print(str(v) + " | " + k)
112:
113:     t = str(round((time.time() - start_time), 2))
114:     logging.info "[" + t + "]" Done! \n"
115:
116:     return
117:
118:
119: # tfidf(word, document) = tf(word, document) * log(|D| / df(word))
120: def calc_tfidf(term, term_freq):
121:     # 1: term frequency of the word in document
```

```
122:     # 2: Amount of Documents
123:     amount_of_documents = len(df1)
124:
125:     if len(document_freq) == 0:
126:         logging.error("document_freq dictionary is empty!")
127:         exit()
128:
129:     # 3: document frequency of the term
130:     document_freq_word = document_freq[term]
131:
132:     tfidf = term_freq * np.log(amount_of_documents / document_freq_word)
133:
134:     return tfidf
135:
136:
137: def create_dic_of_all_article_with_terms():
138:     t = str(round((time.time() - start_time), 2))
139:     logging.info "[" + t + "]" Compute tfidf dictionary ..."
140:
141:     dic = {}
142:
143:     # Iterate over df1
144:     for i in range(0, len(df1)):
145:         articles_dic = {}
146:         word_rank_dic = calc_term_freq(df1.loc[i].text)
147:
148:         # Iterate over the word rank dictionary of the article df1.loc[i]
149:         for k, v in word_rank_dic.items():
150:             # Store each word as key in the article_dic and the
151:             # tfidf of the word as value
152:             articles_dic[k] = calc_tfidf(k, v)
153:
154:         # Store the resulting article_dic as a value and the
155:         # article id as a key in the dic dictionary
156:         dic[i] = articles_dic
157:
158:     t = str(round((time.time() - start_time), 2))
159:     logging.info "[" + t + "]" Done! \n"
160:
161:     return dic
162:
163:
164: # Generate a dictionary with a unique vector for each word
165: def create_word_vectors():
166:     global word_vector
167:
168:     t = str(round((time.time() - start_time), 2))
169:     logging.info "[" + t + "]" Compute word vectors ..."
170:
```

```
171:     # Get the amount of all words in all documents
172:     l = len(document_freq)
173:     i = 0
174:
175:     for k, v in document_freq.items():
176:         v = zeros(l)
177:         v[i] = 1
178:         word_vector[k] = v
179:         i += 1
180:
181:     t = str(round((time.time() - start_time), 2))
182:     logging.info "[" + t + "]" Done! \n")
183:
184:     return
185:
186:
187: # Computes the sparse tfidf vector for an article and stores
188: # it in a global dictionary together with its euclidean
189: # length as a tuple
190: # We don't do this in the calculateCosineSimilarity method
191: # because we only want to compute the vector for every article once
192: # since this improves performance dramatically.
193: def compute_sparse_tfidf_vector(article_id):
194:     global sparse_tfidf_vectors
195:     a = article_dic[article_id]
196:
197:     # Calculate the vector for the given article
198:     article_vector = np.zeros(len(document_freq))
199:     for term, tfidf in a.items():
200:         # Get the vector corresponding to the current word
201:         vec = word_vector[term]
202:         # Multiply vector with tfidf
203:         article_vector += vec * tfidf
204:
205:     sparse_tfidf_vectors[article_id] = (article_vector,
206:                                         np.linalg.norm(article_vector))
207:
208:     return
209:
210:
211: # Iterates of a list of article ids and computes its sparse vectors
212: def compute_sparse_tfidf_vector_from_list(article_list):
213:     t = str(round((time.time() - start_time), 2))
214:     logging.info "[" + t + "]" Computing sparse vectors ..."
215:
216:     for i in range(0, len(article_list)):
217:         compute_sparse_tfidf_vector(article_list[i])
218:
219:     t = str(round((time.time() - start_time), 2))
```



```
220:     logging.info "[" + t + "]" + " Done!")
221:
222:     return
223:
224:
225: def calculateCosineSimilarity(tfIdfDict1, tfIdfDict2):
226:     vector1 = tfIdfDict1[0]
227:     vector2 = tfIdfDict2[0]
228:
229:     # Get the dot product of the vectors
230:     dot = vector1.dot(vector2)
231:
232:     # Get the length of both vectors
233:     length_vector1 = tfIdfDict1[1]
234:     length_vector2 = tfIdfDict2[1]
235:
236:     # Some articles are empty. If this is the case, mark it
237:     if length_vector1 == 0 or length_vector2 == 0:
238:         return -1
239:
240:     # Inverse cosine of the dot product divided by the product of
241:     # the length of both vectors
242:     cosine_sim = dot / (length_vector1*length_vector2)
243:
244:     return cosine_sim
245:
246:
247: # Calculate the length of all articles and choose the longest ones
248: def get_longest_articles():
249:     l = []
250:
251:     # Iterate over df1 and append a tuple to the list l
252:     # The tuple contains the index of the text combined with its length
253:     for i in range(0, len(df1)):
254:         l.append((i, len(df1.loc[i].text)))
255:
256:     # Sort the list by text-length descending in place
257:     l.sort(key=lambda tup: tup[1], reverse=True)
258:
259:     # Generate list with 100 entries containing only the article ids
260:     l2 = []
261:     for i in range(0, 100):
262:         l2.append(l[i][0])
263:
264:     return l2
265:
266:
267: # Select 100 random articles
268: def get_random_articles():
```

```
269:     l = []
270:
271:     for i in range(0, 100):
272:         random_index = np.random.choice(df1.index, replace=False)
273:         l.append(random_index)
274:
275:     return l
276:
277:
278: def compute_jaccard_similarity_of_all_articles(articles):
279:     t = str(round((time.time() - start_time), 2))
280:     logging.info "[" + t + "]" Compute jaccard similarities "
281:         "of all articles ...")
282:
283:     matrix = np.zeros(shape=(100, 100))
284:
285:     # Compute every possible combination of articles. We only need to
286:     # calculate the similarity once per pair and we do
287:     # not need to calculate the similarity of the article with itself.
288:     for i in range(0, len(articles)):
289:         i_article = create_set(df1.loc[articles[i]].text)
290:
291:         # Some articles are empty. If this is the case, skip it.
292:         if len(i_article) == 0: continue
293:
294:         for j in range(0, i):
295:             j_article = create_set(df1.loc[articles[j]].text)
296:
297:             # Some articles are empty. If this is the case, skip it.
298:             if len(j_article) == 0: continue
299:
300:             jaccard = calcJaccardSimilarity(i_article, j_article)
301:             matrix[i, j] = jaccard
302:
303:         # Logging
304:         if i % 10 == 0:
305:             t = str(round((time.time() - start_time), 2))
306:             logging.info "[" + t + "]" Finished " + str(i) +
307:                 "% of all articles")
308:
309:     t = str(round((time.time() - start_time), 2))
310:     logging.info "[" + t + "]" Done! \n")
311:
312:     return matrix
313:
314:
315: def compute_cosine_similarity_of_all_articles(articles):
316:     t = str(round((time.time() - start_time), 2))
317:     logging.info "[" + t + "]" Compute cosine similarities "
```

```
318:         "of all articles ...")
319:
320:     matrix = np.zeros(shape=(100, 100))
321:
322:     # First compute the sparse vectors for each article in the list
323:     compute_sparse_tfidf_vector_from_list(articles)
324:
325:     # Compute every possible combination of articles. We only need to
326:     # calculate the similarity once per pair and we do
327:     # not need to calculate the similarity of the article with itself.
328:     for i in range(0, len(articles)):
329:         i_vec = sparse_tfidf_vectors[articles[i]]
330:
331:         for j in range(0, i):
332:             j_vec = sparse_tfidf_vectors[articles[j]]
333:
334:             cos_sim = calculateCosineSimilarity(i_vec, j_vec)
335:             # Store the cosine similarity in a matrix
336:             matrix[i, j] = cos_sim
337:
338:     # Logging
339:     if i % 10 == 0:
340:         t = str(round((time.time() - start_time), 2))
341:         logging.info "[" + t + "]" + " Finished " + str(i) +
342:             "% of all articles"
343:
344:     t = str(round((time.time() - start_time), 2))
345:     logging.info "[" + t + "]" + " Done! \n"
346:
347:     return matrix
348:
349:
350: def compute_jaccard_similarity_for_outlinks_of_all_articles(articles):
351:     t = str(round((time.time() - start_time), 2))
352:     logging.info "[" + t + "]" + " Compute jaccard similarities for outlinks "
353:         "of all articles ..."
354:
355:     matrix = np.zeros(shape=(100, 100))
356:
357:     # Compute every possible combination of articles. We only need to
358:     # calculate the similarity once per pair and we do
359:     # not need to calculate the similarity of the article with itself.
360:     for i in range(0, len(articles)):
361:         i_article = set(df2.loc[articles[i]].out_links)
362:
363:         # Some articles are empty. If this is the case, skip it.
364:         if len(i_article) == 0: continue
365:
366:         for j in range(0, i):
```

```
367:         j_article = set(df2.loc[articles[j]].out_links)
368:
369:         # Some articles are empty. If this is the case, skip it.
370:         if len(j_article) == 0: continue
371:
372:         jaccard = calcJaccardSimilarity(i_article, j_article)
373:         matrix[i, j] = jaccard
374:
375:         # Logging
376:         if i % 10 == 0:
377:             t = str(round((time.time() - start_time), 2))
378:             logging.info "[" + t + "]" Finished " + str(i) +
379:                 "% of all articles"
380:
381:         t = str(round((time.time() - start_time), 2))
382:         logging.info "[" + t + "]" Done! \n"
383:
384:     return matrix
385:
386:
387: # Add all matrix values to a list and sort it,
388: # then return the highest 10 values
389: def find_best_matches(matrix):
390:     l = []
391:     for i in range(0, len(matrix)):
392:         for j in range(0, i):
393:             l.append((matrix[i, j], (i, j)))
394:
395:     l.sort(key=lambda tup: tup[0], reverse=True)
396:
397:     return l[:10]
398:
399:
400: def print_best_matches(matrix1):
401:     l1 = find_best_matches(matrix1)
402:
403:     for i in range(0, len(l1)):
404:         x = l1[i][1][0]
405:         y = l1[i][1][1]
406:         print("| %1.3f | %30s - %-30s |" %
407:             (l1[i][0],
408:              df1[df1.index == x].name.values[0],
409:              df1[df1.index == y].name.values[0]))
410:     return
411:
412:
413: def pretty_print_matrix(matrix):
414:     for i in range(0, 100):
415:         for j in range(0, 100):
```

```
416:         print("%.3f|" % (matrix[i, j]), end='')
417:     print()
418:
419:
420: def main():
421:     global start_time, df1, df2, article_dic
422:     # Set start time
423:     start_time = time.time()
424:
425:     # Reset log file
426:     with open('similarity.log', 'w'):
427:         pass
428:
429:     t = str(round((time.time() - start_time), 2))
430:     logging.info "[" + t + "]" --- Started --- \n")
431:
432:     store = pd.HDFStore('store2.h5')
433:     df1 = store['df1']
434:     df2 = store['df2']
435:
436:     # -----#
437:     # 1.1.1 Calculate the Jaccard coefficient for the #
438:     # articles "Germany" and "Europe" #
439:     # -----#
440:     word_set_germany = create_set(
441:         str(df1[df1.name == 'Germany'].text.values[0]))
442:     word_set_europe = create_set(
443:         str(df1[df1.name == 'Europe'].text.values[0]))
444:
445:     print("Jaccard Similarity of Germany and Europe: " +
446:           str(calcJaccardSimilarity(word_set_germany,
447:                                     word_set_europe)))
448:
449:     # -----#
450:     # 1.1.2 Calculate the cosine similarity for the #
451:     # article "Germany" and "Europe" #
452:     # -----#
453:     t = str(round((time.time() - start_time), 2))
454:     logging.info "[" + t + "]" --- Started calculations for 1.1.2 ---"
455:
456:     # Compute the document frequency and store it in
457:     # global dictionary document_freq
458:     log_doc_freq()
459:
460:     # Compute a dictionary with the article id as key and
461:     # a dictionary for each article containing the article terms
462:     # and the corresponding tfidf values as value
463:     article_dic = create_dic_of_all_article_with_terms()
464:
```

```
465:     # Compute a dictionary to match a unique vector to a specific word.
466:     # This is done before calculating the cosine similarity
467:     # so this only has to be executed once
468:     create_word_vectors()
469:
470:     # Get ID of the article "Germany" and "Europe"
471:     ger_id = df1[df1.name == "Germany"].index[0]
472:     eur_id = df1[df1.name == "Europe"].index[0]
473:
474:     # Create sparse vectors for each article
475:     compute_sparse_tfidf_vector_from_list([ger_id, eur_id])
476:
477:     # Calculate the cosine similarity of both articles
478:     cos_sim = calculateCosineSimilarity(sparse_tfidf_vectors[ger_id],
479:                                         sparse_tfidf_vectors[eur_id])
480:
481:     print("Cosine Similarity for Germany and Europe: " + str(cos_sim))
482:
483:     t = str(round((time.time() - start_time), 2))
484:     logging.info "[" + t + "]" --- Finished calculations for 1.1.2 --- \n")
485:
486:     # -----#
487:     # 1.2 Similarity of graphs #
488:     # -----#
489:     t = str(round((time.time() - start_time), 2))
490:     logging.info "[" + t + "]" --- Started calculations for 1.2 ---"
491:
492:     germany_outlinks = set(df2[df2.name == "Germany"].out_links.values[0])
493:     europe_outlinks = set(df2[df2.name == "Europe"].out_links.values[0])
494:     outlinks_jaccard_sim = calcJaccardSimilarity(germany_outlinks,
495:                                                  europe_outlinks)
496:
497:     print("Jaccard Similarity of outlinks of Germany and Europe: "
498:           + str(outlinks_jaccard_sim))
499:
500:     t = str(round((time.time() - start_time), 2))
501:     logging.info "[" + t + "]" --- Finished calculations for 1.2 --- \n")
502:
503:     # -----#
504:     # 1.4 Implement the measure #
505:     # -----#
506:     t = str(round((time.time() - start_time), 2))
507:     logging.info "[" + t + "]" --- Started calculations for 1.4 ---"
508:
509:     t = str(round((time.time() - start_time), 2))
510:     logging.info "[" + t + "]" --- Started calculations "
511:                  "for the 100 longest articles --- \n")
512:     # Find the 100 longest articles
513:     l_articles = get_longest_articles()
```

```
514:
515:     jaccard_matrix = \
516:         compute_jaccard_similarity_of_all_articles(l_articles)
517:     cosine_matrix = \
518:         compute_cosine_similarity_of_all_articles(l_articles)
519:     jaccard_outlinks_matrix = \
520:         compute_jaccard_similarity_for_outlinks_of_all_articles(l_articles)
521:
522:     print("Longest articles:")
523:     print_best_matches(jaccard_matrix)
524:     print("-")
525:     print_best_matches(cosine_matrix)
526:     print("-")
527:     print_best_matches(jaccard_outlinks_matrix)
528:
529:     t = str(round((time.time() - start_time), 2))
530:     logging.info "[" + t + "]" --- Started calculations "
531:                 "for 100 random articles --- \n")
532:
533:     # Find 100 random articles
534:     r_articles = get_random_articles()
535:     jaccard_matrix = \
536:         compute_jaccard_similarity_of_all_articles(r_articles)
537:     cosine_matrix = \
538:         compute_cosine_similarity_of_all_articles(r_articles)
539:     jaccard_outlinks_matrix = \
540:         compute_jaccard_similarity_for_outlinks_of_all_articles(r_articles)
541:
542:     print("Random articles: ")
543:     print_best_matches(jaccard_matrix)
544:     print("-")
545:     print_best_matches(cosine_matrix)
546:     print("-")
547:     print_best_matches(jaccard_outlinks_matrix)
548:
549:     t = str(round((time.time() - start_time), 2))
550:     logging.info "[" + t + "]" --- Finished calculations for 1.4 --- \n")
551:
552:     t = str(round((time.time() - start_time), 2))
553:     logging.info "[" + t + "]" --- Finished ---"
554:
555:     # Exit manually, because the log wont be complete otherwise
556:     exit(0)
557:
558: if __name__ == '__main__':
559:     main()
```

Console Output:

```

/Library/Frameworks/Python.framework/Versions/3.5/bin/python3.5 /Users/johanneskirchne
Jaccard Similarity of Germany and Europe: 0.04588607594936709
Cosine Similarity for Germany and Europe: 0.153587715144
Jaccard Similarity of outlinks of Germany and Europe: 0.27307692307692305
Longest articles:
| 0.388 | Seychelles - Democratic_Republic_of_the_Congo_6e9c |
| 0.355 | 1944 - Solid |
| 0.322 | Namibia - Plant |
| 0.295 | Democratic_Republic_of_the_Congo_6e9c - Solid |
| 0.290 | Plant - European_Union_e368 |
| 0.290 | Latvia - Red |
| 0.290 | Seychelles - Solid |
| 0.288 | 1944 - Democratic_Republic_of_the_Congo_6e9c |
| 0.284 | Angola - Language |
| 0.276 | Hard_and_soft_drugs - Black |
-
| 0.817 | 1944 - Solid |
| 0.736 | Namibia - Plant |
| 0.734 | Namibia - European_Union_e368 |
| 0.706 | Angola - Language |
| 0.682 | Plant - European_Union_e368 |
| 0.622 | Mozambique - Wikipedia |
| 0.614 | Seychelles - Democratic_Republic_of_the_Congo_6e9c |
| 0.552 | Tolerance - World |
| 0.544 | Hard_and_soft_drugs - Black |
| 0.500 | 1944 - Seychelles |
-
| 0.355 | Angola - Language |
| 0.350 | Hard_and_soft_drugs - Black |
| 0.304 | Namibia - Plant |
| 0.280 | Hard_and_soft_drugs - 2002 |
| 0.275 | Latvia - Red |
| 0.266 | 2002 - Black |
| 0.261 | Plant - European_Union_e368 |
| 0.219 | Tolerance - Official_language |
| 0.214 | Namibia - European_Union_e368 |
| 0.206 | Seychelles - Solid |

```



```
Random articles:
| 0.778 | International_System_of_Units_e6f4 - Holy_Roman_Empire_f314 |
| 0.737 | Moldova - Democratic_Republic_of_the_Congo_6e9c |
| 0.636 | Assen - Holy_Roman_Empire_f314 |
| 0.636 | Assen - International_System_of_Units_e6f4 |
| 0.625 | Utrecht - U.S._customary_units_c1a4 |
| 0.238 | Angola - Holy_Roman_Empire_f314 |
| 0.238 | Angola - International_System_of_Units_e6f4 |
| 0.237 | 1944 - Official_language |
| 0.222 | Assen - Togo |
| 0.211 | Angola - Togo |
-
| 0.372 | Utrecht - U.S._customary_units_c1a4 |
| 0.358 | 1944 - Official_language |
| 0.309 | List_of_Dutch_people_c3ba - Seychelles |
| 0.288 | Milk - Language_families_and_languages |
| 0.280 | Moldova - Democratic_Republic_of_the_Congo_6e9c |
| 0.267 | Assen - International_System_of_Units_e6f4 |
| 0.248 | Assen - Holy_Roman_Empire_f314 |
| 0.241 | 1944 - Belgium |
| 0.229 | International_System_of_Units_e6f4 - Holy_Roman_Empire_f314 |
| 0.184 | Conservatism - Telephone |
-
| 0.667 | Utrecht - U.S._customary_units_c1a4 |
| 0.562 | City - National_anthem |
| 0.364 | International_System_of_Units_e6f4 - Holy_Roman_Empire_f314 |
| 0.333 | Curve - Angola |
| 0.316 | List_of_Dutch_people_c3ba - Seychelles |
| 0.222 | Assen - International_System_of_Units_e6f4 |
| 0.200 | Assen - Holy_Roman_Empire_f314 |
| 0.150 | Same-sex_marriage - Somaliland |
| 0.143 | Hard_and_soft_drugs - National_anthem |
| 0.143 | Hard_and_soft_drugs - City |
Closing remaining open files:store2.h5...done
Process finished with exit code 0
```

Log Output:

```
1 INFO:root:[0.0] --- Started ---
2
3 INFO:root:[0.5] --- Started calculations for 1.1.2 ---
4 INFO:root:[0.5] Compute document frequency ...
5 INFO:root:[3.36] Done!
6
7 INFO:root:[3.36] Compute tfidf dictionary ...
8 INFO:root:[17.4] Done!
9
10 INFO:root:[17.4] Compute word vectors ...
11 INFO:root:[18.04] Done!
12
13 INFO:root:[18.05] Computing sparse vectors ...
14 INFO:root:[18.35] Done!
15 INFO:root:[18.35] --- Finished calculations for 1.1.2 ---
16
17 INFO:root:[18.35] --- Started calculations for 1.2 ---
18 INFO:root:[18.36] --- Finished calculations for 1.2 ---
19
20 INFO:root:[18.36] --- Started calculations for 1.4 ---
21 INFO:root:[18.36] --- Started calculations for the 100 longest articles ---
22
23 INFO:root:[21.74] Compute jaccard similarities of all articles ...
24 INFO:root:[21.74] Finished 0% of all articles
25 INFO:root:[21.97] Finished 10% of all articles
26 INFO:root:[22.52] Finished 20% of all articles
27 INFO:root:[23.3] Finished 30% of all articles
28 INFO:root:[24.25] Finished 40% of all articles
29 INFO:root:[25.42] Finished 50% of all articles
30 INFO:root:[26.76] Finished 60% of all articles
31 INFO:root:[28.24] Finished 70% of all articles
32 INFO:root:[29.91] Finished 80% of all articles
33 INFO:root:[31.7] Finished 90% of all articles
34 INFO:root:[33.39] Done!
35
36 INFO:root:[33.39] Compute cosine similarities of all articles ...
37 INFO:root:[33.39] Computing sparse vectors ...
38 INFO:root:[62.95] Done!
39 INFO:root:[62.95] Finished 0% of all articles
40 INFO:root:[62.97] Finished 10% of all articles
41 INFO:root:[62.98] Finished 20% of all articles
42 INFO:root:[63.0] Finished 30% of all articles
43 INFO:root:[63.02] Finished 40% of all articles
44 INFO:root:[63.05] Finished 50% of all articles
45 INFO:root:[63.09] Finished 60% of all articles
46 INFO:root:[63.14] Finished 70% of all articles
47 INFO:root:[63.19] Finished 80% of all articles
48 INFO:root:[63.25] Finished 90% of all articles
49 INFO:root:[63.33] Done!
```

```
50
51 INFO:root:[63.33] Compute jaccard similarities for outlinks of all articles ...
52 INFO:root:[63.34] Finished 0% of all articles
53 INFO:root:[63.35] Finished 10% of all articles
54 INFO:root:[63.38] Finished 20% of all articles
55 INFO:root:[63.42] Finished 30% of all articles
56 INFO:root:[63.48] Finished 40% of all articles
57 INFO:root:[63.55] Finished 50% of all articles
58 INFO:root:[63.65] Finished 60% of all articles
59 INFO:root:[63.75] Finished 70% of all articles
60 INFO:root:[63.87] Finished 80% of all articles
61 INFO:root:[64.0] Finished 90% of all articles
62 INFO:root:[64.13] Done!
63
64 INFO:root:[64.17] --- Started calculations for 100 random articles ---
65
66 INFO:root:[64.22] Compute jaccard similarities of all articles ...
67 INFO:root:[64.22] Finished 0% of all articles
68 INFO:root:[64.24] Finished 10% of all articles
69 INFO:root:[64.27] Finished 20% of all articles
70 INFO:root:[64.33] Finished 30% of all articles
71 INFO:root:[64.42] Finished 40% of all articles
72 INFO:root:[64.54] Finished 50% of all articles
73 INFO:root:[64.66] Finished 60% of all articles
74 INFO:root:[64.83] Finished 70% of all articles
75 INFO:root:[65.02] Finished 80% of all articles
76 INFO:root:[65.4] Done!
77
78 INFO:root:[65.4] Compute cosine similarities of all articles ...
79 INFO:root:[65.4] Computing sparse vectors ...
80 INFO:root:[67.45] Done!
81 INFO:root:[67.45] Finished 0% of all articles
82 INFO:root:[67.45] Finished 10% of all articles
83 INFO:root:[67.46] Finished 20% of all articles
84 INFO:root:[67.48] Finished 30% of all articles
85 INFO:root:[67.51] Finished 40% of all articles
86 INFO:root:[67.54] Finished 50% of all articles
87 INFO:root:[67.58] Finished 60% of all articles
88 INFO:root:[67.64] Finished 70% of all articles
89 INFO:root:[67.69] Finished 80% of all articles
90 INFO:root:[67.75] Finished 90% of all articles
91 INFO:root:[67.83] Done!
92
93 INFO:root:[67.83] Compute jaccard similarities for outlinks of all articles ...
94 INFO:root:[67.83] Finished 0% of all articles
95 INFO:root:[67.84] Finished 10% of all articles
96 INFO:root:[67.86] Finished 20% of all articles
97 INFO:root:[67.9] Finished 30% of all articles
98 INFO:root:[67.95] Finished 40% of all articles
99 INFO:root:[68.01] Finished 50% of all articles
100 INFO:root:[68.11] Finished 60% of all articles
101 INFO:root:[68.21] Finished 70% of all articles
102 INFO:root:[68.31] Finished 80% of all articles
103 INFO:root:[68.44] Finished 90% of all articles
104 INFO:root:[68.56] Done!
105
106 INFO:root:[68.6] --- Finished calculations for 1.4 ---
107
108 INFO:root:[68.6] --- Finished ---
```

1.5 Hints:

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

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**.