# Introduction to Web Science

**Assignment 8**

Prof. Dr. Steffen Staab

staab@uni-koblenz.de

Korok Sengupta

koroksengupta@uni-koblenz.de

René Pickhardt

rpickhardt@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 : Yankee
Candidates :

1. Sabin Bhattarai - 216203590
   sbhattarai@uni-koblenz.de

2. Biplov K.C. - 216203865
   biplov@uni-koblenz.de

3. Syed Salman Ali. - 216203923
   salmanali@uni-koblenz.de

# 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 coeffcient 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 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`.

## Answer

Following is the result obtained from our implemented functions for 1.1.1 and 1.1.2 and 1.2. The code is presented in Answer section of 1.2.

```
Jaccard Similarity between two articles Germany and Europe is:
0.03504043126684636

Cosine Similarity between two articles Germany and Europe is:
5.42425451145e-05

Jacard with Graph Similarity between two articles Germany and Europe is:
0.27307692307692305
```

**Figure 1:** Obtained results from our implemented cosine, jaccards, jaccard with outlinks

## 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`.

---

**Answer 1.1.1** and **Answer 1.1.2** and **Answer 1.2** all uses the following code to initially have necessary tfidf calculated, necessary dictionary consisting of articles with respective word list, total word count, total document frequency for each word etc. calculated and saved in files.

Note: The cosine similarity funtion, Jaccard and Jaccard with graph simililarity functions are shown all together in 1.4 answer section.

---

```
 1:
 2: #NOTE: This python file creates necessary files with calculated
 3: # values required for cosine and jackard and statistical measure
 4: # calculation. These files saved heryby will be loaded in next .py file.
 5:
 6:
 7: import pandas as pd
 8: import itertools
 9: from collections import Counter
10: import math
```

```
11: import pickle
12:
13: articles_with_totalwords = {}
14: articles_list_ofwords = {}
15: articles_list_ofoutlinks = {}
16: tf_idf_article_eachterm = {}
17:
18:
19: # tf_dict has id as article and value = dictionary of each term with associated
20: # term frequency in that document.
21: def tf_idf_cal(tf_dict , df_dict):
22:     global  tf_idf_article_eachterm
23:     tf_idf_article_eachterm = {}
24:     total_document_length = len(tf_dict.keys())
25:     for each_article in tf_dict:
26:         tf_idf = {}
27:         for each_term in tf_dict[each_article]:
28:             tf_idf[each_term]= tf_dict[each_article][each_term] * \
29:                         math.log(total_document_length / df_dict[each_term],10)
30:         tf_idf_article_eachterm[each_article] = tf_idf
31:
32:
33: # returns the dictionary of term frequency
34: def count_term_frequency_each_article(dict_df1):
35:     global articles_with_totalwords
36:     global articles_list_ofwords
37:     df1_withtermfrequency_eacharticle = {}
38:     for article_name in dict_df1.keys():
39:         list_words = dict_df1[article_name][0].split()
40:         articles_list_ofwords[article_name] = list_words
41:         articles_with_totalwords[article_name] = len(list_words)
42:         df1_withtermfrequency_eacharticle[article_name] = \
43:                 dict(Counter(dict_df1[article_name][0].split()).most_common())
44:     return df1_withtermfrequency_eacharticle
45:
46:
47: # updates the unique_term_dict value which contains the all unique terms
48: # by document frequency for given terms and dataframe
49: def count_document_frequency_for_eachterm(df1 ,unique_terms_dict ):
50:     document_frequency_for_eachterm_dic = {}
51:     copyDF = df1.copy()
52:     copyDF['wordset']= copyDF.text.map(lambda x: set(x.lower().split()))
53:     for eacharticleterms in copyDF['wordset']:
54:         for eachterm in eacharticleterms:
55:             if eachterm in document_frequency_for_eachterm_dic:
56:                 num=document_frequency_for_eachterm_dic[eachterm]+1
57:             else:
58:                 num=1
59:             document_frequency_for_eachterm_dic[eachterm] = num
```

```
60:     return document_frequency_for_eachterm_dic
61:
62: # Responsible for reading articles from file and evaluating tfIDF , creating
63: # dictionary of list of words per article and dictionry of article and len words
64: def read_articles_from_file_evaluate_tfidf(filename):
65:     global all_terms_basevector_dict
66:     global articles_list_ofoutlinks
67:     store = pd.HDFStore(filename)#read .h5 file
68:     df1=store['df1']
69:     df2=store['df2']
70:     df1["text"] = df1.text.str.lower()
71:
72:     # Dictionary of article names and its associated article text in list form
73:     dict_df1 = df1.set_index('name').T.to_dict('list')
74:     dict_df2 = df2.set_index('name').T.to_dict()
75:     articles_list_ofoutlinks = dict_df2
76:
77:     # get the term frequency for each article given the dictionary of datas
78:     articles_with_termfrequency_dict = count_term_frequency_each_article(dict_df1)
79:
80:     # now we get document frequency . first get all terms and put in the set
81:     all_terms = list(itertools.chain.from_iterable([list(\
82:                     term_frequency_dict.keys()) for term_frequency_dict \
83:                             in articles_with_termfrequency_dict.values()])))
84:     unique_terms_dict = {each_term : 0 for each_term in all_terms}
85:
86:     document_frequency_for_eachterm_diction = \
87:                     count_document_frequency_for_eachterm(df1, unique_terms_dict)
88:
89:     # Evaluate tfidf for each term in articles. Stored in dictionary
90:     tf_idf_cal(articles_with_termfrequency_dict , \
91:                                     document_frequency_for_eachterm_diction)
92:
93:
94: if __name__ == "__main__":
95:      # Here we read articles from file, create base vectors, calculate tfidf,
96:      # calculate document vectors and store all in global variable
97:      read_articles_from_file_evaluate_tfidf("store.h5")
98:
99:      # Save  ariicletotalwords, tfidfforeachterm , articles with outlinks
100:     # and articles with its words into a pickle file.
101:     try:
102:         pickle.dump( articles_with_totalwords, open( \
103:                                     "articles_with_totalwords.p","wb"))
104:         pickle.dump( tf_idf_article_eachterm, open( \
105:                                     "tf_idf_article_eachterm.p","wb"))
106:         pickle.dump( articles_list_ofwords, open( \
107:                                     "articles_list_ofwords.p","wb"))
108:         pickle.dump( articles_list_ofoutlinks, open( \
```

```
109:                          "articles_list_ofoutlinks.p","wb"))
110:      except Exception as e:
111:          print("ERROR: Failed to save to a file")
112:          import sys
113:          sys.exit(1)
```

## 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:**

After the research we decided to use Kendall's Tau as a statistical measure to compare two different ranking measure for the same object.
Kendall's Tau($\tau$) is a coefficient that represents the degree of concordance between two columns of ranked data.

$\tau = \frac{\text{number of concordant pairs} - \text{number of discordant pairs}}{\text{number of concordant pairs} + \text{number of discordant pairs}}$

Concordant pairs: the number of observed ranks below a particular rank which are larger than that particular rank.

Discordant pairs: the number of observed ranks below a particular rank which are smaller in value than that particular rank.

We would perform Kendall's tau and calculate statistical significance on three different combinations: Jaccard co-efficient with cosine similarity, Jaccard co-efficient with Graph with Jaccard co-efficient and Cosine similarity with Graph with Jaccard co-efficient. Below is an example where X could represent cosine similarity and Y could represent Jaccard's co-efficient. Here, X would be the baseline.

| X | Y | Concordant | Discordant |
|---|---|---|---|
| 1 | 2 | 10 | 1 |
| 2 | 1 | 10 | 0 |
| 3 | 4 | 8 | 1 |
| 4 | 3 | 8 | 0 |
| 5 | 6 | 6 | 1 |
| 6 | 5 | 6 | 0 |
| 7 | 8 | 4 | 1 |
| 8 | 7 | 4 | 0 |
| 9 | 10 | 2 | 1 |
| 10 | 9 | 2 | 0 |
| 11 | 12 | 0 | 1 |
| 12 | 11 | | |

Here,

$\tau = \frac{\text{sum of concordant} - \text{sum of discordant}}{\text{sum of concordant} + \text{sum of discordant}}$

or, $\tau = \frac{60-6}{60+6}$

or, $= 0.818$

So, for above given data Kendall's co-efficient is 0.818

Now,

For statistical significance (Z) we use:

$Z = \frac{3 * \tau * \sqrt{n(n-1)}}{\sqrt{2(2n+5)}}$

or, $Z = \frac{3 * 0.818 * \sqrt{12(12-1)}}{\sqrt{2(2*12+5)}}$

or, $Z = 3.7019$

For any value of 'Z' if 'Z' is greater than 1.96 the X and Y are statistically significant. After calculating statistical significance for 3 combinations of similarity measure (Jaccard co-efficient with cosine similarity, Jaccard co-efficient - Graph with Jaccard co-efficient and Cosine similarity with Jaccard co-efficient - Graph) on 100 articles- long or random we would plot a bar graph where each segment(Jaccard co-efficient with cosine similarity, Jaccard co-efficient with Graph with Jaccard co-efficient and Cosine similarity with Graph with Jaccard co-efficient) would consists of two bars; one representing number of values less than 1.96 and the other greater than equals to 1.96.

Once the graph is observed, depending upon number of count we could explain if certain article always be the same independent which model we use.

## 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?

---

**Answer 1.4**

Note it also has implementations for question part 1.1 and 1.2 in the following code snippet.

**1.**
The total number of similarity computations that would have to be done can be given as ????????

**2.**
The total time that would roughly be consumed to do all these computations would be

????????

The following shows the code alongside the plot which implements the above mentioned experiment for statistical measure.

```
1:
2: import pickle
3: import numpy as np
4: import random
5: from collections import Counter
6: import operator
7: import math
8: import pandas as pd
9: import matplotlib.pyplot as plt
10:
11: from_file_articles_with_totalwords = {}
12: from_file_articles_list_ofwords = {}
13: from_file_tf_idf_article_eachterm = {}
14: from_file_articles_list_ofoutlinks = {}
15:
16: def calculateCosineSimilarity(article1, article2):
17:     try:
18:         tf_idf_article1 = from_file_tf_idf_article_eachterm[article1]
19:         tf_idf_article2 = from_file_tf_idf_article_eachterm[article2]
20:
21:         scalar_product_article1_article2 = 0
22:         for each_term_article1 in tf_idf_article1:
23:             if each_term_article1 in tf_idf_article2.keys():
24:                 scalar_product_article1_article2 = \
25:                             scalar_product_article1_article2 + \
26:                             (tf_idf_article1[each_term_article1] * \
27:                                 tf_idf_article2[each_term_article1])
28:
29:         euc_dist_tfIdfDict1 = np.sum(np.square(list(tf_idf_article1.values())))
30:         euc_dist_tfIdfDict2 = np.sum(np.square(list(tf_idf_article2.values())))
31:
32:         cosineSimilarity = scalar_product_article1_article2 / \
33:                             (euc_dist_tfIdfDict1 * euc_dist_tfIdfDict2)
34:
35:     except Exception as e:
36:         print("ERROR: Failed to get article for calculating Cosine")
37:         import sys
38:         sys.exit(1)
39:     return cosineSimilarity
40:
41: def calcJaccardSimilarity(article1, article2):
42:     wordset1 = from_file_articles_list_ofwords[article1]
43:     wordset2 = from_file_articles_list_ofwords[article2]
```

```
44:        wordset1=set(wordset1)
45:        wordset2=set(wordset2)
46:        inter=wordset1.intersection(wordset2)
47:        union=wordset1.union(wordset2)
48:        if (len(union) > 0 ):
49:            jc=(len(inter)/len(union))
50:        else:
51:            jc = 1
52:        return jc
53:
54: def calcJaccardGraphSimilarity(article1, article2):
55:        wordset1 = from_file_articles_list_ofoutlinks[article1]["out_links"]
56:        wordset2 = from_file_articles_list_ofoutlinks[article2]["out_links"]
57:        wordset1=set(wordset1)
58:        wordset2=set(wordset2)
59:        inter=wordset1.intersection(wordset2)
60:        union=wordset1.union(wordset2)
61:        if (len(union) > 0 ):
62:            jc=(len(inter)/len(union))
63:        else:
64:            jc = 1
65:        return jc
66:
67:
68:
69: def initialise_loading():
70:        global from_file_articles_with_totalwords
71:        global from_file_articles_list_ofwords
72:        global from_file_tf_idf_article_eachterm
73:        global from_file_articles_list_ofoutlinks
74:
75:        # Lets load all necessary already calculated values
76:        try:
77:
78:            from_file_articles_with_totalwords = pickle.load(open\
79:                                            ( "articles_with_totalwords.p","rb"))
80:            from_file_tf_idf_article_eachterm = pickle.load(open\
81:                                            ( "tf_idf_article_eachterm.p","rb"))
82:            from_file_articles_list_ofwords = pickle.load(open\
83:                                            ( "articles_list_ofwords.p","rb"))
84:            from_file_articles_list_ofoutlinks = pickle.load(open\
85:                                            ( "articles_list_ofoutlinks.p","rb"))
86:        except Exception as e:
87:            print("ERROR: Failed to load from file")
88:            import sys
89:            sys.exit(1)
90:
91:
92:
```

```
 93:
 94:
 95: #-------------------------------------------------------------------------------
 96: # Begins Experiment for 1.3 and 1.4
 97: #-------------------------------------------------------------------------------
 98:
 99: def longest_article(how_many):
100:     article_in_order = Counter(from_file_articles_with_totalwords).most_common()
101:     if (len(article_in_order) <= how_many):
102:         return [article for (article , numberofwords) in article_in_order]
103:     allarticles = [article for (article , numberofwords) in article_in_order]
104:     return allarticles[:how_many]
105:
106: def random_articles(how_many):
107:     if (len(from_file_articles_with_totalwords) <= how_many):
108:         return random.sample(from_file_articles_with_totalwords.keys(),len(\
109:                                     from_file_articles_with_totalwords))
110:     return random.sample(from_file_articles_with_totalwords.keys(),how_many)
111:
112: def giverankto_tuplelist(tuple_list):
113:     i = 1
114:     ranking_tuple = {}
115:     for (key , val) in tuple_list:
116:         ranking_tuple[key] = i
117:         i = i + 1
118:     return ranking_tuple
119:
120:
121: def calculate_statistical_significance(ranked_column):
122:     C = [sum(i > val for i in ranked_column[idx + 1:]) for idx, val in \
123:                                     enumerate(ranked_column)]
124:     try:
125:         C.pop()
126:     except:
127:         print("No elements to get the Statistic measure\n")
128:
129:
130:     D = [sum(i < val for i in ranked_column[idx + 1:]) for idx, val in \
131:                                     enumerate(ranked_column)]
132:     try:
133:         D.pop()
134:     except:
135:         print("No elements to get the Statistic measure\n")
136:
137:     kendalls_tau = (sum(C) - sum(D)) / (sum(C) + sum(D))
138:
139:     n = len(ranked_column)
140:     statistical_significance_Z = (3 * kendalls_tau * math.sqrt(n * (n - 1))) / \
141:                                     math.sqrt(2 * (2*n+5))
```

```
142:
143:     #Note Any Z value greater than 1.96 is going to statistically significant.
144:     return statistical_significance_Z
145:
146:
147: # Cosine vs Jacakard
148: def get_statistical_significance_cosine_jacakard(firstArticle , \
149:                                         remainingArticles_list):
150:     col1_from_cosine_firstcompared_toothers = {remainingArticles_list[i]: \
151:         calculateCosineSimilarity(firstArticle , remainingArticles_list[i])\
152:                                         for i in range (0,100)}
153:     sorted_cosine = sorted(col1_from_cosine_firstcompared_toothers.items(), \
154:                         key=operator.itemgetter(1) , reverse=True)
155:     ranked_cosine_dict = giverankto_tuplelist(sorted_cosine)
156:
157:
158:     col2_from_jackard_firstcompared_toothers = {remainingArticles_list[i]: \
159:         calcJaccardSimilarity(firstArticle , remainingArticles_list[i]) \
160:                                         for i in range (0,100)}
161:     sorted_jacard = sorted(col2_from_jackard_firstcompared_toothers.items(), \
162:                         key=operator.itemgetter(1) , reverse=True)
163:     ranked_jacard_column = [ranked_cosine_dict[term] for (term , value) \
164:                                         in sorted_jacard]
165:     return calculate_statistical_significance(ranked_jacard_column)
166:
167: #Cosine vs Jackard with outlinks - Graphs
168: def get_statistical_significance_cosine_jacakardGraph(firstArticle , \
169:                                         remainingArticles_list):
170:     col1_from_cosine_firstcompared_toothers = {remainingArticles_list[i]: \
171:        calculateCosineSimilarity(firstArticle , remainingArticles_list[i]) \
172:                                         for i in range (0,100)}
173:     sorted_cosine = sorted(col1_from_cosine_firstcompared_toothers.items(), \
174:                             key=operator.itemgetter(1) , reverse=True)
175:     ranked_cosine_dict = giverankto_tuplelist(sorted_cosine)
176:
177:
178:     col2_from_jackardGraph_firstcompared_toothers = {remainingArticles_list[i]:\
179:        calcJaccardGraphSimilarity(firstArticle , remainingArticles_list[i]) \
180:                                         for i in range (0,100)}
181:     sorted_jacard = sorted(col2_from_jackardGraph_firstcompared_toothers.items(),\
182:                             key=operator.itemgetter(1) , reverse=True)
183:     ranked_jacard_column = [ranked_cosine_dict[term] for (term , value) in \
184:                                         sorted_jacard]
185:     return calculate_statistical_significance(ranked_jacard_column)
186:
187:
188: # Jackard vs Jackard with outlinks - Graphs
189: def get_statistical_significance_jacakard_jacakardGraph(firstArticle , \
190:                                         remainingArticles_list):
```

```
191:        col1_from_cosine_firstcompared_toothers = {remainingArticles_list[i]: \
192:                calcJaccardSimilarity(firstArticle , remainingArticles_list[i])\
193:                                                for i in range (0,100)}
194:        sorted_cosine = sorted(col1_from_cosine_firstcompared_toothers.items(), \
195:                                    key=operator.itemgetter(1) , reverse=True)
196:        ranked_cosine_dict = giverankto_tuplelist(sorted_cosine)
197:
198:
199:        col2_from_jackardGraph_firstcompared_toothers = {remainingArticles_list[i]:\
200:              calcJaccardGraphSimilarity(firstArticle , remainingArticles_list[i])\
201:                                                for i in range (0,100)}
202:        sorted_jacard = sorted(col2_from_jackardGraph_firstcompared_toothers.items(),\
203:                                    key=operator.itemgetter(1) , reverse=True)
204:        ranked_jacard_column = [ranked_cosine_dict[term] for (term , value) in \
205:                                                sorted_jacard]
206:        return calculate_statistical_significance(ranked_jacard_column)
207:
208:
209: # We will evaluate each pair and do the plot by measuring whish was more
210: # significant compared to each other via frequency plot
211: def performStatistical_significance_and_plot(type_of_article_random_or_longest):
212:        cosine_jackard_dict = {}
213:        cosine_jackard_graph_dict = {}
214:        jackard_jackard_graph_dict = {}
215:        for article in type_of_article_random_or_longest:
216:            temp = type_of_article_random_or_longest[:]
217:            temp.remove(article)
218:            cosine_jackard_dict[article] =  \
219:                    get_statistical_significance_cosine_jacakard(article , temp)
220:            cosine_jackard_graph_dict[article] =  \
221:                get_statistical_significance_cosine_jacakardGraph(article , temp)
222:            jackard_jackard_graph_dict[article] =  \
223:                get_statistical_significance_jacakard_jacakardGraph(article , temp)
224:
225:
226:        # ------------------------------------------------------------------------
227:        # Counting the total significance for each pair and plotting
228:        Cosine_JackardDict = (list(cosine_jackard_dict.values()))
229:        Cosine_JackardGraphDict = (list(cosine_jackard_graph_dict.values()))
230:        Jackard_JackardGraphDict = (list(jackard_jackard_graph_dict.values()))
231:
232:        df = pd.DataFrame(
233:            {'Cosine_Jackard': Cosine_JackardDict,
234:             'Cosine_JackardGraph': Cosine_JackardGraphDict,
235:             'Jackard_JackardGraph': Jackard_JackardGraphDict
236:            })
237:        df_seperated = pd.DataFrame({'Insignificant' : df[df <= 1.96].count(),
238:            'Significant' : df[df > 1.96].count()})
239:        df_seperated.plot(kind='bar')
```

13

```
240:        plt.ylabel("Frequency")
241:        plt.show()
242:
243:
244: if __name__ == "__main__":
245:
246:        initialise_loading()
247:        print("\n Jaccard Similarity between two articles Germany and Europe is: \n"\
248:                                  , calcJaccardSimilarity("Germany" , "Europe" ))
249:        print("\n Cosine Similarity between two articles Germany and Europe is: \n"\
250:                                  , calculateCosineSimilarity("Germany" , "Europe" ))
251:        print("\n Jacard with Graph Similarity between two articles Germany and " + \
252:              "Europe is: \n" , calcJaccardGraphSimilarity("Germany" , "Europe"))
253:
254:        #-------------------------------------------------------------------------
255:        # Answer 1.4 Performing Experiment and testing with plots
256:        #-------------------------------------------------------------------------
257:
258:        # Here I will have to get hundred Longest and random Article comparisons
259:        hundredlongestArticle = longest_article(101)
260:        randomarticles = random_articles(101)
261:
262:        performStatistical_significance_and_plot(hundredlongestArticle)
263:        performStatistical_significance_and_plot(randomarticles)
```
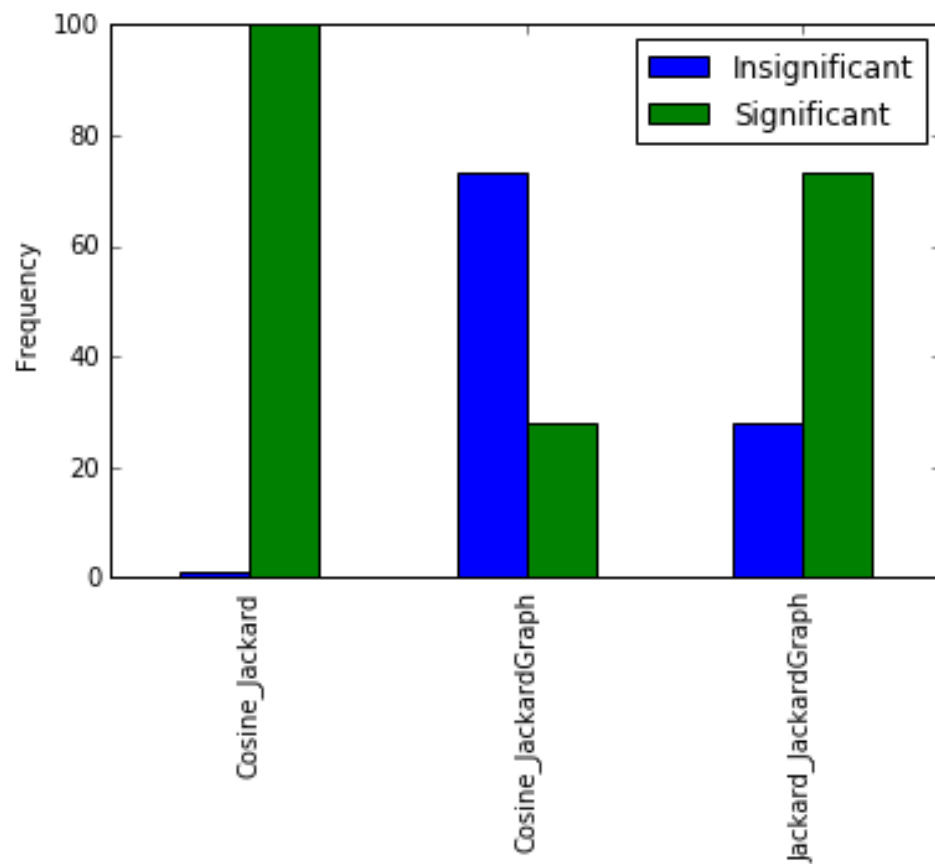
## Strategy 1 : 100 Longest Articles



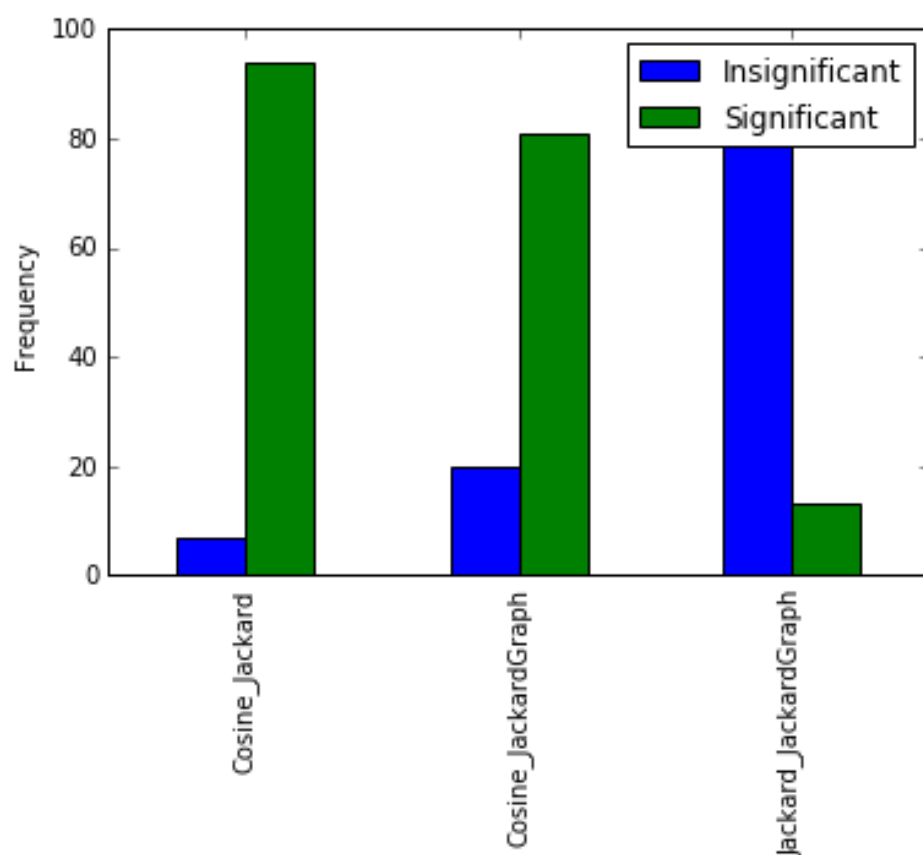**Figure 2:** Obtained significance evaluated through frequency measure.

**Strategy 2 : 100 Random Articles**



**Figure 3:** Obtained significance evaluated through frequency measure.

## 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"
```

```python
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())

#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`.