# Introduction to Web Science

## Assignment 7

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:  December 14, 2016, 10:00 a.m.
Tutorial on:  December 16, 2016, 12:00 p.m.

Please look at the lessons 1) **Similarity of Text** & 2) **Generative 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.

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

1

# 1 Modelling Text in a Vector Space and calculate similarity (10 points)

Given the following three documents:

$D_1$ = this is a text about web science

$D_2$ = web science is covering the analysis of text corpora

$D_3$ = scientific methods are used to analyze webpages

## 1.1 Get a feeling for similarity as a human

Without applying any modeling methods just focus on the semantics of each document and decide which two Documents should be most similar. Explain why you have this opinion in a short text using less than 500 characters.

## 1.2 Model the documents as vectors and use the cosine similarity

Now recall that we used vector spaces in the lecture in order to model the documents.

1. How many base vectors would be needed to model the documents of this corpus?

2. What does each dimension of the vector space stand for?

3. How many dimensions does the vector space have?

4. Create a table to map words of the documents to the base vectors.

5. Use the notation and formulas from the lecture to represent the documents as document vectors in the word vector space. You can use the term frequency of the words as coefficients. You can / should omit the inverse document frequency.

6. Calculate the cosine similarity between all three pairs of vectors.

7. According to the cosine similarity which 2 documents are most similar according to the constructed model.

## 1.3 Discussion

Do the results of the model match your expectations from the first subtask? If yes explain why the vector space matches the similarity given from the semantics of the documents. If no explain what the model lacks to take into consideration. Again 500 Words should be enough.

**Answers:**

**1.1**

Without applying any modeling methods we can say that two documents are similar if they contain some of the same term. The possible measures of similarity might be:
The length of the document, the number of terms in common, whether the terms are common or unusual and how many times the each term appears etc. While considering above facts we can say that the document D1 is much more similar to document D2. There is some common terms in document D1 and D2. The semantic of the D1 and D2 look same, both are somehow talking about web science.

**1.2**

1. In the given three documents there are 19 different words(terms) so, 19 base vectors would be needed to model the documents of this corpus.

2. Each dimensions of the vector space stands for term/token used to index a set of documents.

3. The vector space have 19 dimensions.

4. A table to map the words of the documents to the base vectors:

   In our documents $D_1$, $D_2$, and $D_3$ we have 19 different words.
   Let V = $<\vec{T_1}, \vec{T_2} .... , \vec{T_{19}} >$ be the vector space spanned by the words "this", "is", "a", "text"...... (all the different words in three documents.)

$$
\vec{T_1} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} ,
\vec{T_2} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} ,
\vec{T_3} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} ,......,
\vec{T_{19}} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}
$$

and so on...

Now, in every document there is no repetition of any words. That means the term frequency of each words in a document is $= 1$. The below table shows the Weighting by Term Frequency (tf):

| -     | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_5$ | $T_6$ | $T_7$ | $T_8$ | $T_9$ | $T_{10}$ | $T_{11}$ | $T_{12}$ | $T_{13}$ | $T_{14}$ | $T_{15}$ |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|----------|----------|
| $d_1$ | 1     | 1     | 1     | 1     | 1     | 1     | 1     | 0     | 0     | 0        | 0        | 0        | 0        | 0        | 0        |
| $d_2$ | 0     | 1     | 0     | 1     | 0     | 1     | 1     | 1     | 1     | 1        | 1        | 1        | 0        | 0        | 0        |
| $d_3$ | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0     | 0        | 0        | 0        | 1        | 1        | 1        |

| -     | $T_{16}$ | $T_{17}$ | $T_{18}$ | $T_{19}$ | length     |
|-------|----------|----------|----------|----------|------------|
| $d_1$ | 0        | 0        | 0        | 0        | $\sqrt{(7)}$ |
| $d_2$ | 0        | 0        | 0        | 0        | 3          |
| $d_3$ | 1        | 1        | 1        | 1        | $\sqrt{(7)}$ |

where T means Term in document vector space.

5. The document vector of a document $D_1$ is:

$$\vec{d_1} = \sum_{i=1}^{7} tf(w_i, D_1)w_i \tag{1}$$

$$= tf(T_1, D_1)\vec{T_1} + tf(T_2, D_1)\vec{T_2} + tf(T_3, D_1)\vec{T_3} + tf(T_4, D_1)\vec{T_4} + tf(T_5, D_1)\vec{T_5} + tf(T_6, D_1)\vec{T_6} + tf(T_7, D_1) \tag{2}$$

$$=$$

$$1 * \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + 1 * \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + ..... + 1 * \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \vec{d_1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Similarly, the document vector for $D_2$ and $D_3$ is given by:

$$\vec{d_2} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad and \quad \vec{d_3} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

6. The cosine similarity between two documents are given by:

$$cos(\theta) = \frac{d_1.d_2}{|d_1|.|d_2|} \quad (3)$$

|     | d1   | d2   | d3 |
|-----|------|------|----|
| d1  | 1    | 0.50 | 0  |
| d2  | 0.50 | 1    | 0  |
| d3  | 0    | 0    | 1  |

$cos(\theta) = 1$; documents exactly the same; $= 0$, totally different
From the above table we can see that the cosine similarity between d1 and d2 is $=$
0.50
The cosine similarity between d1 and d3 $= 0$.
The cosine similarity between d2 and d3 $= 0$

7. According to the cosine similarity document d1 document d2 are most similar.

**1.3**
The results of the model and our expectation from sub-task1.1 is matched. Each dimensions
of the vector space is represented by each term/token in the documents where term/token
refer to the unique word in the documents. As we know the cosine similarity of two
documents is give by $cos(\theta)$ as in equation1. If the number of similar words in two
document are higher then the value of $cos(\theta)$ is also high. Because of this our expectation
matched.

# 2 Building generative models and compare them to the observed data (10 points)

This week we provide you with two probability distributions for characters and spaces which can be found next to the exercise sheet on the WeST website. Also last week we provided you with a dump of Simple English Wikipedia which should be reused this week.

## 2.1 build a generator

Count the characters and spaces in the Simple English Wikipedia dump. Let the combined number be $n$. Use the sampling method from the lecture to sample $n$ characters (which could be letters or a space) from each distribution. Store the result for the generated text for each distribution in a file.

## 2.2 Plot the word rank frequency diagram and CDF

Count the resulting words from the provided data set and from the generated text for each of the probability distributions. Create a word rank frequency diagram which contains all 3 data sets. Also create a CDF plot that contains all three data sets.

## 2.3 Which generator is closer to the original data?

Let us assume you would want to creat a test corpus for some experiments. That test corpus has to have a similar word rank frequency diagram as the original data set. Which of the two generators would you use? You should perform the Kolmogorov Smirnov test as discussed in the lecture by calculating the maximum pointwise distance of the CDFs.

**How do your results change when you generate the two text corpora for a second or third time? What will be the values of the Kolmogorov Smirnov test in these cases?**

**Answer 2.1**

**Listing 1** helper.py

```
1: """
2:   Helper Methods
3: """
4: # pylint: disable-msg=C0103
5: from threading import Thread
6: import random
7: import bisect
8: from collections import Counter
```

```
 9: from functools import reduce
10: from copy import deepcopy
11: import numpy as np
12: import matplotlib.pyplot as plt
13:
14: THREADS = 30
15:
16: def remove_special_char(regex, string):
17:     """removes any character than a-z and 0-9"""
18:     return regex.sub('', string.lower())
19:
20:
21: def read_file_content(file_path):
22:     """ open given file, read and return the content
23:     None will be returned if error while opening the file
24:     """
25:     content = None
26:     with open(file_path, 'r+') as f:
27:         content = f.read()
28:     return content
29:
30:
31: def chunks(lis, num):
32:     """ returns the chunks """
33:     n = max(1, num)
34:     return (lis[i:i+n] for i in range(0, len(lis), n))
35:
36:
37: def save_file(file_name, content):
38:     """ writes to the file"""
39:     with open(file_name, 'w+') as f:
40:         f.write(content)
41:     print("{} written sucessfully".format(file_name))
42:
43:
44: def cal_probability(list_of_touple):
45:     """ calculates the probability each key from given list"""
46:     total = reduce(lambda x, y: x + y[1], list_of_touple, 0)
47:     return[{x[0]: x[1]/total} for x in list_of_touple]
48:
49:
50: def cal_cum_frequency(input_list_of_dict):
51:     """ calculates the cumulative frequency"""
52:     list_of_dict = deepcopy(input_list_of_dict)
53:     for index, obj in enumerate(list_of_dict, start=1):
54:         req_dicts = list_of_dict[:index]
55:         total = 0
56:         for req_dict in req_dicts:
57:             for key in req_dict.keys():
```

```
58:                    if key != 'cum_freq':
59:                        total += req_dict[key]
60:            obj['cum_freq'] = total
61:        return list_of_dict
62:
63:
64: def _generate_text(char_count, keys, keys_cdfs, result, index):
65:     """ generate text using thread"""
66:     characters = []
67:     for x in range(0, char_count // THREADS):
68:         randomValue = random.random()
69:         idx = bisect.bisect(keys_cdfs, randomValue)
70:         characters.append(keys[idx])
71:     text = "".join(characters)
72:     result[index] = text
73:     return text
74:
75: def gen_random_text(char_count, keys, keys_cdfs):
76:     """ generate random text"""
77:     threads = [None] * 1000
78:     results = {}
79:     i = 0
80:     for i in range(0, len(threads)):
81:         threads[i] = Thread(target=_generate_text,\
82:                                     args=(char_count, keys, keys_cdfs, results, i
83:         threads[i].start()
84:
85:     for thread in threads:
86:         thread.join()
87:
88:     return ''.join(list(results.values()))
89:
90:
91: def get_word_with_probability(text):
92:     """retun words with its probability"""
93:     freq_each_words = Counter(text.split()).most_common()
94:     sum_of_total_char_occ = sum([frequency for (key, frequency) in freq_each_words
95:     prob_each_word = {}
96:     for (key, frequency) in freq_each_words:
97:         prob_each_word[key] = frequency / sum_of_total_char_occ
98:     return prob_each_word
99:
100:
101: def draw_plot(listElements_S, listElements_Z, listElements_U, xLabel, yLabel):
102:     """ draws the plot"""
103:     x_S = [x for x in range(1, len(listElements_S)+1)]
104:     y_S = np.array(listElements_S)
105:     x_Z = [x for x in range(1, len(listElements_Z)+1)]
106:     y_Z = np.array(listElements_Z)
```

```
107:        x_U = [x for x in range(1, len(listElements_U)+1)]
108:        y_U = np.array(listElements_U)
109:        plt.figure(figsize=(12, 9))
110:        plt.plot(x_S, y_S, 'r', label="Simple English")
111:        plt.plot(x_Z, y_Z, 'b', label="Zips Distribution Words")
112:        plt.plot(x_U, y_U, 'g', label="Uniform Distribution Words")
113:        plt.xlabel(xLabel)
114:        plt.ylabel(yLabel)
115:        plt.yscale('log')
116:        plt.xscale('log')
117:        plt.legend(loc='upper right')
118:        plt.grid()
119:        plt.show()
120:
121:
122: def draw_cdf_plot(listElements_S, listElements_Z, listElements_U, xLabel, yLabel)
123:        """draws the cdf plot"""
124:        x_S = [x for x in range(1, len(listElements_S) + 1)]
125:        y_S = np.array(listElements_S)
126:        x_Z = [x for x in range(1, len(listElements_Z) + 1)]
127:        y_Z = np.array(listElements_Z)
128:        x_U = [x for x in range(1, len(listElements_U) + 1)]
129:        y_U = np.array(listElements_U)
130:        plt.figure(figsize=(12, 9))
131:        plt.plot(x_S, y_S, 'r', label="Simple English")
132:        plt.plot(x_Z, y_Z, 'b', label="Zips Distribution Words")
133:        plt.plot(x_U, y_U, 'g', label="Uniform Distribution Words")
134:        plt.xlabel(xLabel)
135:        plt.ylabel(yLabel)
136:        plt.xscale('log')
137:        plt.yscale('log')
138:        plt.ylim(0, 1.2)
139:        plt.legend(loc='bottom right')
140:        plt.grid()
141:        plt.show()
```

**Listing 2** assignment7.py

```
 1: """
 2:  Generative Model
 3: """
 4: # pylint: disable-msg=C0103
 5: import re
 6: import threading
 7: from timeit import default_timer as timer
 8: from pathlib import Path
 9: from collections import Counter
10: from operator import itemgetter
11: import numpy as np
12: from helper import read_file_content
```
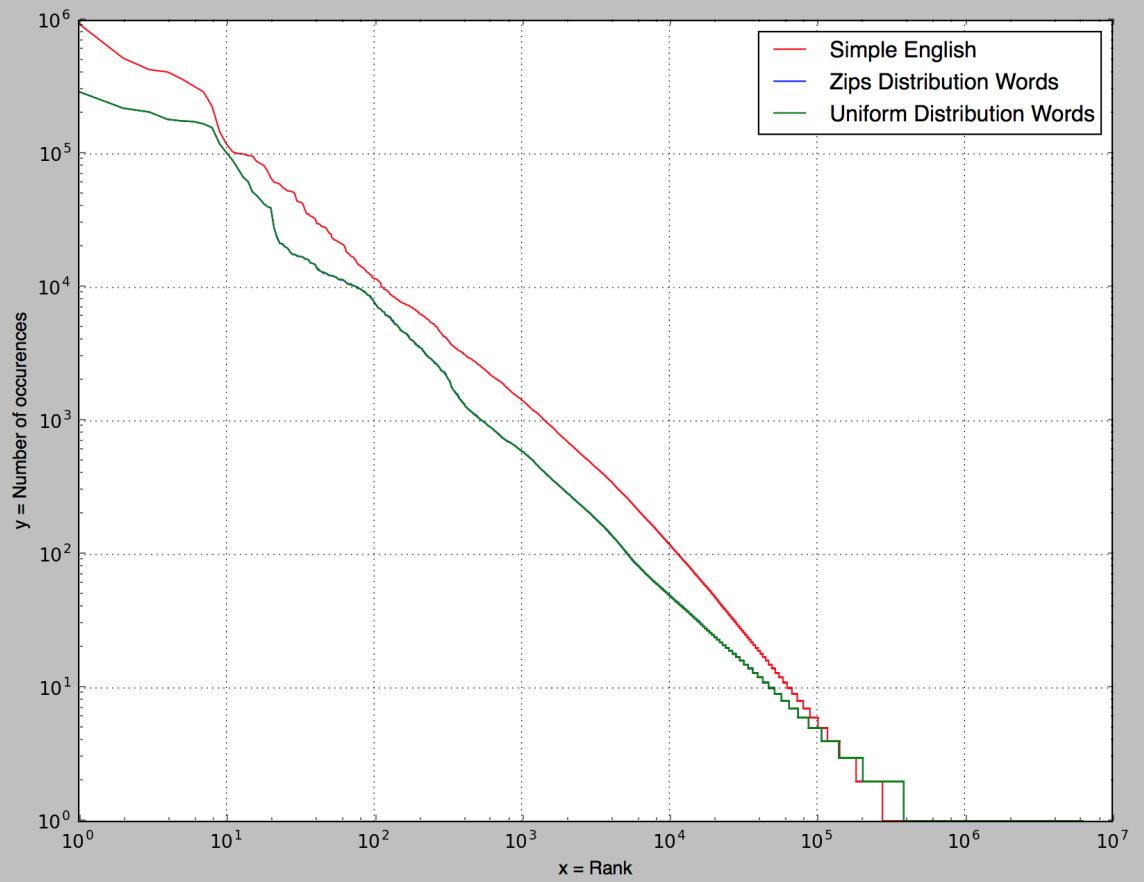
```
13: from helper import chunks
14: from helper import remove_special_char
15: from helper import save_file
16: from helper import cal_probability
17: from helper import cal_cum_frequency
18: from helper import draw_plot
19: from helper import gen_random_text
20: from helper import get_word_with_probability
21: from helper import draw_cdf_plot
22:
23:
24: counter = Counter()
25: total_character_in_english_text = 0
26: probabilities_file = 'probabilities.py.txt'
27: regex = re.compile('[^a-zA-Z0-9 ]')
28:
29: def cdf_calculaton(prob_dict):
30:     """ calculates cdf"""
31:     arr_for_cdf = list(prob_dict.values())
32:     a = np.array(arr_for_cdf)
33:     cdfEvalDict = np.cumsum(a)
34:     return (list(prob_dict.keys()), cdfEvalDict)
35:
36: def count_alphanumberic(string):
37:     """ count frequency of each alphabet"""
38:     global counter
39:     clean_string = remove_special_char(regex, string)
40:     ctr = Counter(clean_string)
41:     counter = counter + ctr
42:
43: def perform_kolmogro_test(cdf_main , cdf_compare):
44:     """ perform kolmogro test"""
45:     max_point_wise_distance = max([abs(cdf_S - cdf_Zipf) \
46:                                 for (cdf_S, cdf_Zipf) in \
47:                                     zip(cdf_main, cdf_compare)])
48:     return max_point_wise_distance
49:
50: def process_articles(articles):
51:     """ access all the article from articles"""
52:     for article in articles:
53:         count_alphanumberic(article)
54:
55:
56: def use_multithread(articles):
57:     """" using multi threading"""
58:     group_articles = chunks(articles, int(len(articles) / 20))
59:     thread_list = []
60:     for art in group_articles:
61:         t = threading.Thread(target=process_articles, args=(art, ))
```

```python
62:            thread_list.append(t)
63:            t.start()
64:
65:        for thread in thread_list:
66:            thread.join()
67:
68:
69:    def without_multi(articles):
70:        """ processing in single thread"""
71:        for article in articles:
72:            count_alphanumberic(article)
73:
74:
75:    def process_output_file(file_path):
76:        """ read already generated output and process it"""
77:        global counter
78:        output = None
79:        with file_path.open() as f:
80:            output = f.read()
81:        counter = eval(output)
82:
83:
84:    def load_file_and_process(file_path):
85:        """ load the simple article file and process it"""
86:        global counter
87:        content = read_file_content(file_path)
88:        start = timer()
89:        articles = content.split("\n")
90:        use_multithread(articles)
91:        end = timer()
92:        print("Total time taken:", end - start)
93:        # sorting the counter
94:        counter = sorted(counter.items(), key=itemgetter(0))
95:        save_file('output.txt', str(counter))
96:
97:
98:    def process_given_info_file(file_path):
99:        """ process the given file which contains probability information"""
100:        def get_text(text, start, end):
101:            """ returns slice of text from stat to end:"""
102:            return text[text.find(start): text.find(end) + 1]
103:
104:        def process_list(given_list):
105:            """ returns dict"""
106:            return list([{x[0]: x[1]} for x in given_list])
107:        file_content = read_file_content(file_path)
108:        text_arr = file_content.split("\n")
109:        texts = list(filter(lambda x: len(x.strip()) > 0, text_arr))
110:        result = list(map(lambda x: get_text(x, '{', '}'), texts))
```
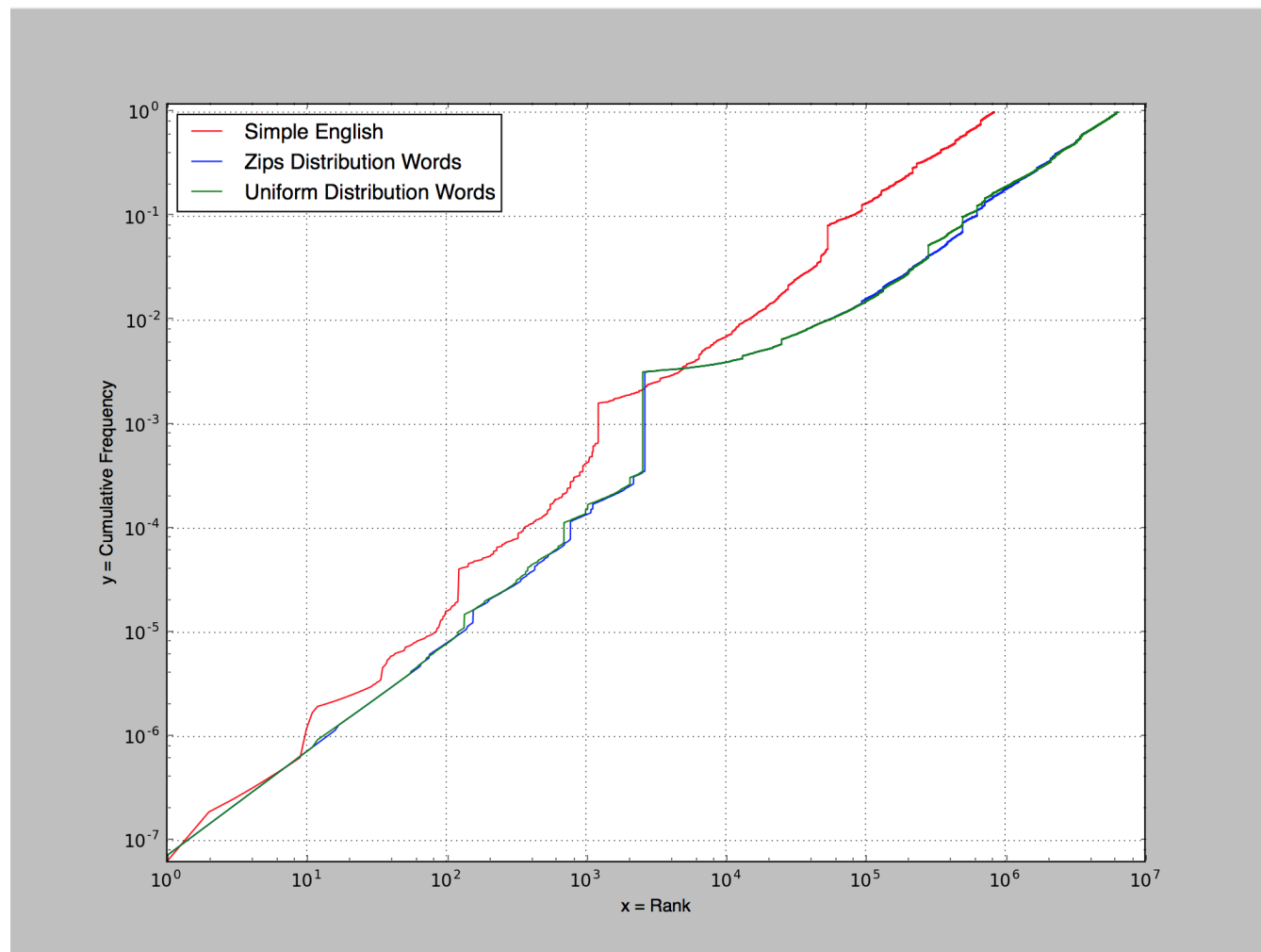
```
111:        final = []
112:        for res in result:
113:            jt = eval(res)
114:            mp = sorted(jt.items(), key=itemgetter(0))
115:            final.append(mp)
116:        return list(map(lambda x: process_list(x), final))
117:
118:
119:
120: def get_keys_and_cdfs(list_of_dict):
121:        """ returns the keys and its cumulative frequency"""
122:        cums = []
123:        keys = []
124:        for cur_dict in list_of_dict:
125:            for key in cur_dict.keys():
126:                if key == 'cum_freq':
127:                    cums.append(cur_dict[key])
128:                else:
129:                    keys.append(key)
130:        return [keys, cums]
131:
132:
133: def main(file_path):
134:        """ entry point of the application"""
135:        global total_character_in_english_text
136:        simple_english_text = read_file_content(file_path)
137:        total_character_in_english_text = len(simple_english_text)
138:        prob_each_word_dict = get_word_with_probability(simple_english_text)
139:        _, assoc_cdf_s = cdf_calculaton(prob_each_word_dict)
140:
141:        # for Zipf and uniform distrubation
142:        prob_dist_given = process_given_info_file(probabilities_file)
143:        zipf_cum = cal_cum_frequency(prob_dist_given[0])
144:        unif_cum = cal_cum_frequency(prob_dist_given[0])
145:        all_keys_zipf, asso_zipf_cum = get_keys_and_cdfs(zipf_cum)
146:        all_keys_unif, asso_unif_cum = get_keys_and_cdfs(unif_cum)
147:
148:        text_zipf = gen_random_text(total_character_in_english_text, all_keys_zipf, as
149:        save_file('text_zipf.txt', text_zipf)
150:
151:        text_unif = gen_random_text(total_character_in_english_text, all_keys_unif, as
152:        save_file('text_unif.txt', text_unif)
153:
154:        prob_for_each_ord_dict_zipf = get_word_with_probability(text_zipf)
155:        _, ass_cdf_new_zipf = cdf_calculaton(prob_for_each_ord_dict_zipf)
156:
157:        prob_for_each_word_dict_unif = get_word_with_probability(text_unif)
158:        _, asso_cdf_new_unif = cdf_calculaton(prob_for_each_word_dict_unif)
159:        #-----------------------------------------------------------------------
```

```
160:        # Plots for Rank Frequency STARTS
161:        #----------------------------------------------------------------------
162:        #Note: _S is for simple English
163:        freq_each_word_s = Counter(simple_english_text.split()).most_common()
164:        list_for_rank_freq_diag_s = [freq for (word, freq) in freq_each_word_s]
165:
166:        #Note _Z is for Zipf
167:        freq_each_word_z = Counter(text_zipf.split()).most_common()
168:        list_for_rank_freq_diag_z = [freq for (word, freq) in freq_each_word_z]
169:
170:        #Note_U is for Uniform
171:        freq_each_word_u = Counter(text_unif.split()).most_common()
172:        list_for_rank_freq_diag_u = [freq for (word, freq) in freq_each_word_u]
173:
174:        draw_plot(list_for_rank_freq_diag_s, list_for_rank_freq_diag_z,\
175:                list_for_rank_freq_diag_u, "x = Rank", "y = Number of occurences")
176:        #----------------------------------------------------------------------
177:        # Plots for Rank Frequency ENDS
178:        #----------------------------------------------------------------------
179:
180:        #----------------------------------------------------------------------
181:        # Plots for Rank and CDF STARTS
182:        #----------------------------------------------------------------------
183:
184:        draw_cdf_plot(assoc_cdf_s, ass_cdf_new_zipf, asso_cdf_new_unif,\
185:         "x = Rank", "y = Cumulative Frequency")
186:
187:        #----------------------------------------------------------------------
188:        # Plots for Rank and CDF ENDS
189:        #----------------------------------------------------------------------
190:        # Now we perform Kolmogorov Smirnov test  by calculating the maximum
191:        # pointwise distance of CDFs
192:        max_point_wised_s_zipf = perform_kolmogro_test(assoc_cdf_s, ass_cdf_new_zipf)
193:        max_point_wised_s_unif = perform_kolmogro_test(assoc_cdf_s, asso_cdf_new_unif)
194:        print("Obtained Maximum Pointwise Distance between simple English \
195:                                        and Zipf is : ", max_point_wised_s_zipf)
196:        print("Obtained Maximum Pointwise Distance between simple English \
197:                                        and Unif is : ", max_point_wised_s_unif)
198:
199: if __name__ == "__main__":
200:        file = 'simple-20160801-1-article-per-line'
201:        main(file)
```

**Answer 2.2**

Word frequencies depending on word

Cumulative word probability depending on word rank

**Answer 2.3**
Obtained Maximum Pointwise Distance between simple English and Zipf is : 0.860301461905
Obtained Maximum Pointwise Distance between simple English and Unif is : 0.861857581856

Although our model is far off, from Figure 2 and its resulting Kolmogorov Smirnov test We would choose data set created by Zipf distribution because it has lower distance measure with our original data set.This suggests Zipf distribution model reflects more to our descriptive model than data set created from uniform distribution

Two text corpora for second time
Obtained Maximum Pointwise Distance between simple English and Zipf is : 0.863258314669
Obtained Maximum Pointwise Distance between simple English and Unif is : 0.852805206287

Two text corpora for third time
Obtained Maximum Pointwise Distance between simple English and Zipf is : 0.857980183468
Obtained Maximum Pointwise Distance between simple English and Unif is : 0.862024209875

Our choice of a model earlier can be called o and be said inconsistent as we now observe that during the repetition of generating two text corpora our Pointwise distance di ered in each iteration.

## 2.4 Hints:

1. Build the cummulative distribution function for the text corpus and the two generated corpora

2. Calculate the maximum pointwise distance on the resulting CDFs

3. You can use `Collections.Counter`, `matplotlib` and `numpy`. You shouldn't need other libs.

# 3 Understanding of the cumulative distribution function (10 points)

Write a fair 6-side die rolling simulator. A fair die is one for which each face appears with equal likelihood. Roll two dice simultaneously n (=100) times and record the sum of both dice each time.

1. Plot a readable histogram with frequencies of dice sum outcomes from the simulation.

2. Calculate and plot cumulative distribution function.

3. Answer the following questions using CDF plot:

   What is the median sum of two dice sides? Mark the point on the plot.

   What is the probability of dice sum to be equal or less than 9? Mark the point on the plot.

4. Repeat the simulation a second time and compute the maximum point-wise distance of both CDFs.

5. Now repeat the simulation (2 times) with n=1000 and compute the maximum point-wise distance of both CDFs.

6. What conclusion can you draw from increasing the number of steps in the simulation?

## 3.1 Hints

1. You can use function from the lecture to calculate rank and normalized cumulative sum for CDF.

2. Do not forget to give proper names of CDF plot axes or maybe even change the ticks values of x-axis.

## 3.2 Only for nerds and board students (0 Points)

Assuming 20 groups of students. What is the likelihood that at least two groups come up with the same histograms in the case for n (=100)?

**Answers:**

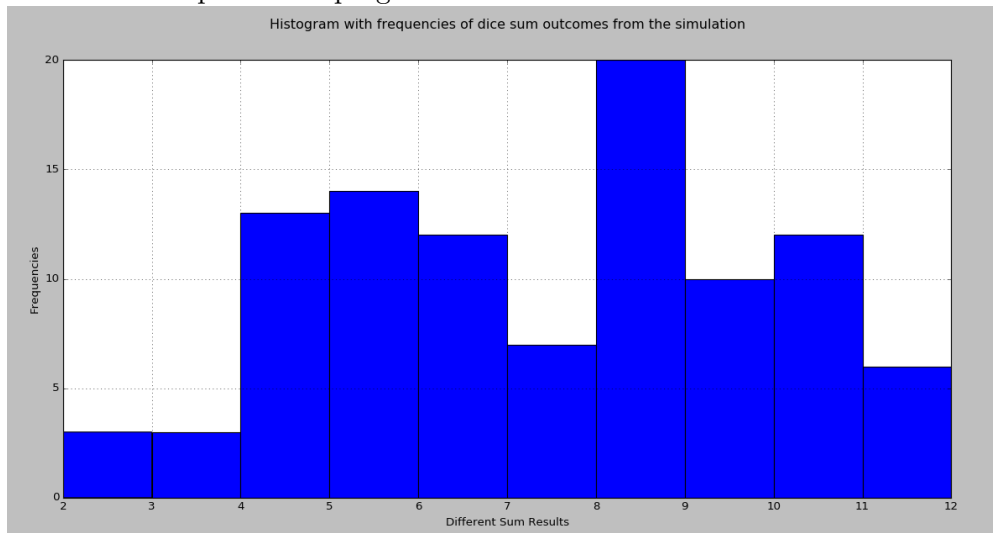The required python program for above question is given below:

```
1: [language=Python, frame=single, caption=uniform_assignment7_Q3.py]
2: import random
3: import numpy as np
4: import matplotlib.pyplot as plt
```
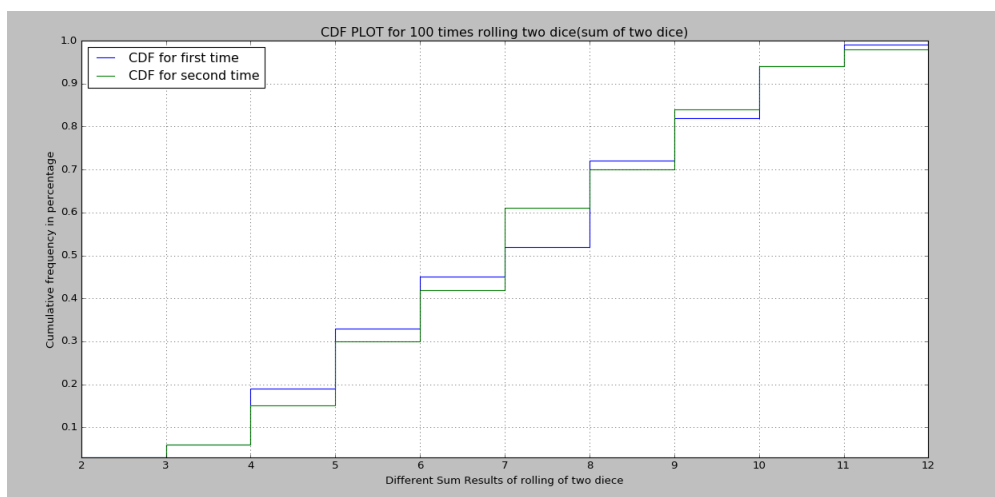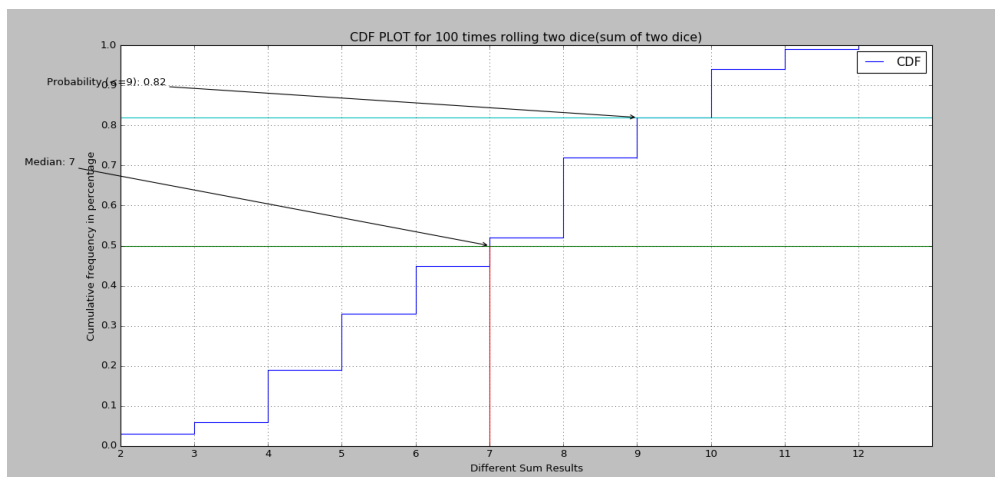
```
 5: import collections
 6:
 7: #function to rolling two diece simultaneously and returns the sum of results
 8: def rollDiceGetSum():
 9:     return random.randint(1, 6) + random.randint(1, 6)
10:
11: #functions to calculate CDF
12: def calculateCDF(result):
13:     frequencyNum = collections.Counter(result).most_common()
14:     total = sum([frequency for (key, frequency) in frequencyNum])
15:     probabilityForEachNumDict = {}
16:     for (key, frequency) in frequencyNum:
17:         probabilityForEachNumDict[key] = frequency / \
18:                                          total
19:
20:     arrayForCDFCalc = list(probabilityForEachNumDict.values())
21:     a = np.array(arrayForCDFCalc)  # Gets us CDF
22:     cdfEvalDict = np.cumsum(a)
23:     return (list(probabilityForEachNumDict.keys()), cdfEvalDict)
24:
25: # function to draw the histogram of sum of results
26: def drawHistogram(sumResult):
27:     plt.hist(sumResult, bins=[2,3,4,5,6,7,8,9,10,11,12])
28:     plt.title('Histogram with frequencies of dice sum outcomes from the simulation
29:     plt.ylabel('Frequencies')
30:     plt.xlabel('Different Sum Results')
31:     plt.xticks(np.arange(2, 13, 1))
32:     plt.grid('on')
33:     plt.show()
34:
35: #functions to draw CDF and shows the median and prob. of sum less or equal to 9
36: def drawSingleCDF(sumResult, median, probLessNnine):
37:     (xval, cdfval) = calculateCDF(sumResult)
38:
39:     plt.title("CDF PLOT for 100 times rolling two dice(sum of two dice)")
40:     plt.xlabel("Different Sum Results")
41:     plt.ylabel("Cumulative frequency in percentage")
42:     plt.grid('on')
43:     plt.xticks(np.arange(2, 13, 1))
44:     plt.yticks(np.arange(0, 1.1, 0.1))
45:     plt.plot(xval, cdfval, drawstyle='steps-post', label='CDF')
46:
47:     plt.plot((2, 13), (.5, .5))
48:     plt.plot((median, median), (0, .5))
49:     plt.annotate('Median: %d' % median, (median, .5), xytext=(0.7, 0.7), arrowpro
50:
51:     plt.plot((2, 13), (probLessNnine, probLessNnine))
52:     plt.annotate('Probability (<=9): %s' % probLessNnine, (9, probLessNnine), xyt
53:                  arrowprops=dict(arrowstyle='->'))
```
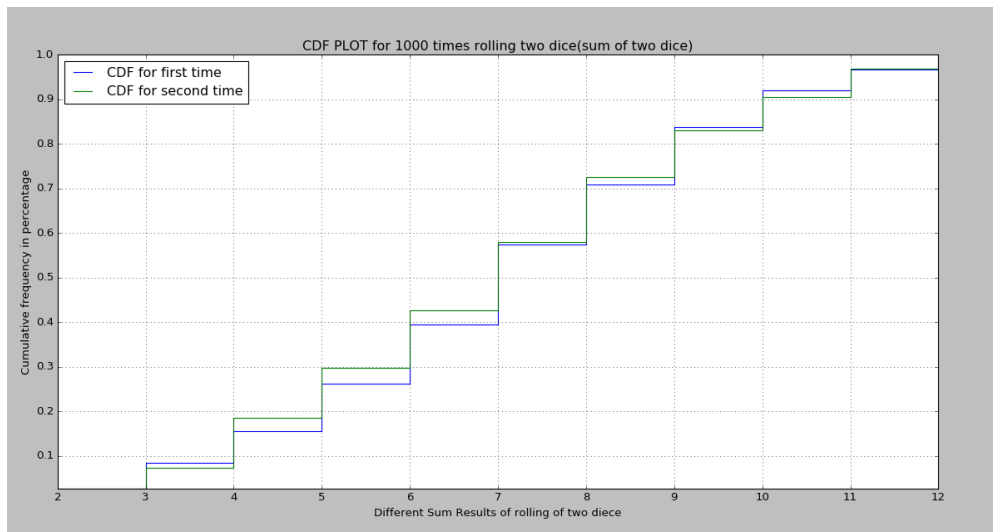
```
54:
55:     plt.legend(loc=0)
56:     plt.show()
57:
58: #functions to draw two CDF according to random number n
59: def drawDoubleCDF(result1, result2, n):
60:     (xval1, cdfval1) = calculateCDF(result1)
61:     (xval2, cdfval2) = calculateCDF(result2)
62:     maxDist = getMaxDistance(cdfval1, cdfval2)
63:
64:     print(" The maximum pointwise distance between two CDFs for n= " + str(n)+ ",
65:
66:     plt.title("CDF PLOT for " + str(n) + " times rolling two dice(sum of two dice)
67:     plt.xlabel("Different Sum Results of rolling of two diece")
68:     plt.ylabel("Cumulative frequency in percentage")
69:     plt.xticks(np.arange(2, 13, 1))
70:     plt.yticks(np.arange(0, 1.1, 0.1))
71:
72:     plt.plot(xval1, cdfval1, drawstyle='steps-post', label='CDF for first time')
73:     plt.plot(xval2, cdfval2, drawstyle='steps-post', label='CDF for second time')
74:
75:     plt.grid('on')
76:     plt.legend(loc=0)
77:     plt.show()
78:
79: #functions to roll two dice uniformly upto n times.
80: def rollDice(n):
81:     sumResult = []
82:     for _ in range(n):
83:         result = rollDiceGetSum()
84:         sumResult.append(result)
85:
86:     sumResult = np.sort(sumResult)
87:     return sumResult
88:
89: #functions to calculate median from the give array
90: def getMedian(sumResult):
91:     a = np.array(sumResult)
92:     print("Median val: ", np.median(a))
93:     return np.median(a)
94:
95: #functions to calculate probability of given value in a array
96: def getProbByVal(arr, val):
97:     gen = [x for x in arr if x<= val]
98:     return (len(gen)/len(arr))
99:
100: #functions to calculate maximum pointwise distance
101: def getMaxDistance(result1, result2):
102:     diffResult = [abs(x - y) for x, y in zip(result1, result2)]
```

```
103:      return max(diffResult)
104:
105: #main function to do random sampling according to question
106: def main():
107:      sumResult = rollDice(100)
108:
109:      # Draw histogram
110:      drawHistogram(sumResult)
111:
112:      median = getMedian(sumResult)
113:      val = getProbByVal(sumResult, 9)
114:      # Draw CDF
115:      drawSingleCDF(sumResult, median, val)
116:
117:      sumResult1 = rollDice(100)
118:
119:      # Draw CDF for n = 100
120:      drawDoubleCDF(sumResult, sumResult1, 100)
121:
122:      # Draw CDF for n= 1000
123:      result1 = rollDice(1000)
124:      result2 = rollDice(1000)
125:      drawDoubleCDF(result1, result2, 1000)
126:
127:
128:
129: if __name__ == '__main__':
130:      main()
```

Here is the output of the program:

CDF PLOT for 100 times rolling two dice(sum of two dice)



CDF PLOT for 100 times rolling two dice(sum of two dice)

```
Median val:  7.0
 The maximum pointwise distance between two CDFs for n= 100, is =0.09
 The maximum pointwise distance between two CDFs for n= 1000, is =0.035
```

3.6

While increasing the number of steps in the simulation we got smaller value of maximum point-wise distance. This means if the number of steps or no. of experiment is increased then we will get the more similar CDF.

## Important Notes

### Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment7/` 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`.