## Introduction to Web Science

### Assignment 7 with solutions

Prof. Dr. Steffen Staab René Pickhardt

staab@uni-koblenz.de rpickhardt@uni-koblenz.de

Korok Sengupta Olga Zagovora

koroksengupta@uni-koblenz.de 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.

Team Name: XXXX



## 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.



#### 1.4 Solution

#### 1.1 Example solution:

From my point of view, sentence  $D_1$  and  $D_2$  are the most similar, because they both are about "web science". However, sentence  $D_3$  should be also similar to sentence  $D_1$  and  $D_2$ , since they all are about "science" and "web".

- 1.2
  - 1. 19
  - 2. each unique word in the dataset
  - 3. 19

	word	[optional column]	base vector
4.	this	$w_1$	$\vec{e_1} = (1, 0, 0, 0, 0, 0, \dots, 0)^T$
	is	$w_2$	$\vec{e_2} = (0, 1, 0, 0, 0, 0, \dots, 0)^T$
	a	$w_3$	$\vec{e_3} = (0, 0, 1, 0, 0, 0, \dots, 0)^T$
	text	$w_4$	$\vec{e_4} = (0, 0, 0, 1, 0, 0, \dots, 0)^T$
	about	$w_5$	$\vec{e_5} = (0, 0, 0, 0, 1, 0, \dots, 0)^T$
	web	$w_6$	$\vec{e_6} = (0, 0, 0, 0, 0, 1, \dots, 0)^T$
	science	$w_7$	$\vec{e_7} = (0, 0, 0, 0, 0, 0, 1, 0, \dots, 0)^T$
	covering	$w_8$	$\vec{e_8} = (0, 0, 0, 0, 0, 0, 0, 1, \dots, 0)^T$
	:	:	
	webpages	$w_6$	$\vec{e_{19}} = (0, 0, 0, 0, 0, 0, \dots, 0, 1)^T$

5.

$$D_{1} \longrightarrow \vec{d_{1}} = \sum_{i=1}^{19} tf(w_{i}, D_{1})\vec{e_{i}} = 1 \cdot \vec{e_{1}} + 1 \cdot \vec{e_{2}} + 1 \cdot \vec{e_{3}} + 1 \cdot \vec{e_{4}} + 1 \cdot \vec{e_{5}} + 1 \cdot \vec{e_{6}} + 1 \cdot \vec{e_{7}} + 0 \cdot \vec{e_{8}} + \dots = \begin{pmatrix} 1\\1\\1\\1\\1\\0\\0\\\vdots\\0 \end{pmatrix}$$

$$(1)$$

$$D_2 \longrightarrow \vec{d_2} = \sum_{i=1}^{19} t f(w_i, D_2) \vec{e_i} = 0 \cdot \vec{e_1} + 1 \cdot \vec{e_2} + 0 \cdot \vec{e_3} + 1 \cdot \vec{e_4} + 0 \cdot \vec{e_5} + 1 \cdot \vec{e_6} + \dots + 1 \cdot \vec{e_{12}} + 0 \cdot \vec{e_{13}} + \dots = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$D_{3} \longrightarrow \vec{d_{3}} = \sum_{i=1}^{19} t f(w_{i}, D_{3}) \vec{e_{i}} = 0 \cdot \vec{e_{1}} + 0 \cdot \vec{e_{2}} + \dots + 0 \cdot \vec{e_{12}} + 1 \cdot \vec{e_{13}} + \dots + 1 \cdot \vec{e_{19}} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$(3)$$

6.

$$cos(\Theta) = \frac{<\vec{a}, \vec{b}>}{\|\vec{a}\| \|\vec{b}\|} = sim(\vec{a}, \vec{b}) \tag{4}$$

$$sim(D_1, D_2) = \frac{\langle \vec{d_1}, \vec{d_2} \rangle}{\|\vec{d_1}\| \|\vec{d_2}\|} = \frac{\sum_{i=1}^{19} tf(w_i, D_1) tf(w_i, D_2)}{\sqrt{\sum_{i=1}^{19} tf(w_i, D_1)^2} \sqrt{\sum_{i=1}^{19} tf(w_i, D_2)^2}}$$
(5)

$$sim(D_1,D_3) = \frac{<\vec{d_1},\vec{d_3}>}{\|\vec{d_1}\|\|\vec{d_3}\|} = \frac{\sum_{i=1}^{19} tf(w_i,D_1)tf(w_i,D_3)}{\sqrt{\sum_{i=1}^{19} tf(w_i,D_1)^2}\sqrt{\sum_{i=1}^{19} tf(\ w_i,D_3)^2}} \qquad (6)$$



7. Documnets  $D_1$  and  $D_2$  are the most similar.

#### 1.3 Example solution:

The results of the model partly match my expectations. On the one hand, documents  $D_1$  and  $D_2$  are the most similar as I expected. On the other hand, (surprisingly) both  $D_1$  and  $D_2$  have similarity of 0 with the document  $D_3$ . From my point of view, the model lack of features which could measure similarity between semanticly close words.



## 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 i $\|$ n 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 create 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?

#### 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?



### **Important Notes**

#### **Submission**

- Solutions have to be checked into the github repository. Use the directory name groupname/assignment7 with solutions/in your group's repository.
- The name of the group and the names of all participating students must be listed on each submission.
- Solution format: all solutions as one PDF document. Programming code has to be submitted as Python code to the github repository. Upload all .py files of your program! Use UTF-8 as the file encoding. Other encodings will not be taken into account!
- Check that your code compiles without errors.
- Make sure your code is formatted to be easy to read.
  - Make sure you code has consistent indentation.
  - Make sure you comment and document your code adequately in English.
  - Choose consistent and intuitive names for your identifiers.
- Do not use any accents, spaces or special characters in your filenames.

#### **Acknowledgment**

This latex template was created by Lukas Schmelzeisen for the tutorials of "Web Information Retrieval".

### **LATEX**

Currently the code can only be build using LuaLaTeX, so make sure you have that installed. If on Overleaf, there's an error, go to settings and change the LaTeX engine to LuaLaTeX.