

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.

Team Name: Bravo

Team Members: Shriharsh Ambhore, Kandhasamy Rajasekaran, Daniel Akbari

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.

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 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 cumulative 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.

Answer:

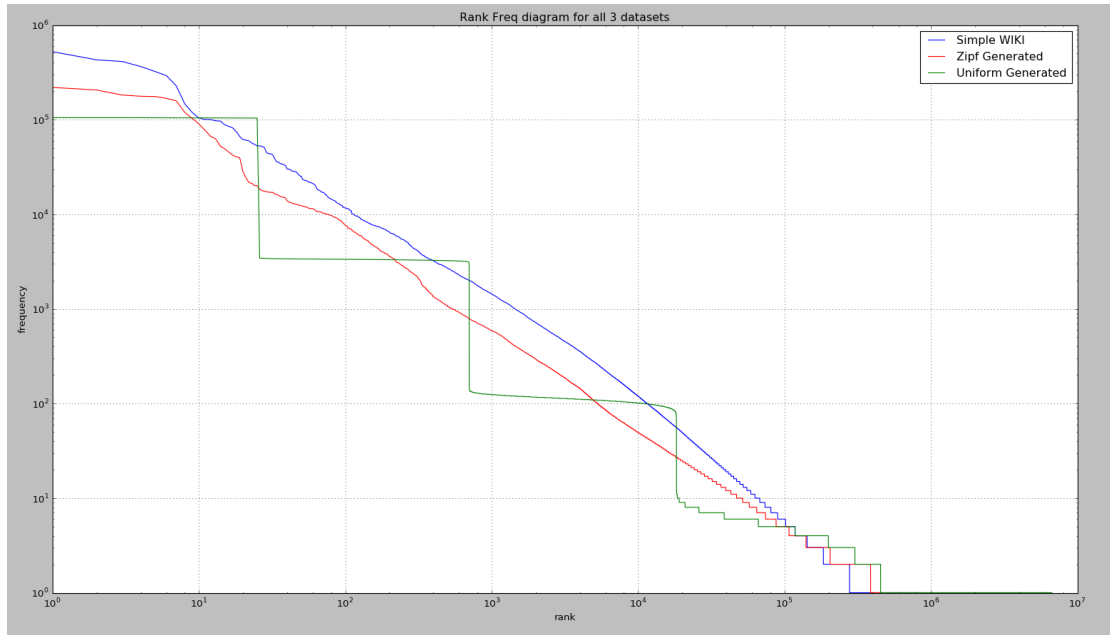


Figure 1: Rank Freq plot for all 3 data sets

Kolmogorov smirnov test value for Simple wiki and Zipf generative model is 0.4118572943437089
Kolmogorov smirnov test value for Simple wiki and Uniform generative model is 0.4362332736488038

Relatively, Zipf generative model is close to our Simple English wikipedia based on Kolmogorov Smirnov test.

On running the program multiple times, Zipf generative model is close to Simple English wikipedia than Uniform generative model.

2nd run: Kolmogorov smirnov test value for Simple wiki and Zipf generative model is 0.3924561212098345 Kolmogorov smirnov test value for Simple wiki and Uniform generative model is 0.5123234287653820

3rd run: Kolmogorov smirnov test value for Simple wiki and Zipf generative model is 0.3512987638253703 Kolmogorov smirnov test value for Simple wiki and Uniform generative model is 0.4569732539405873

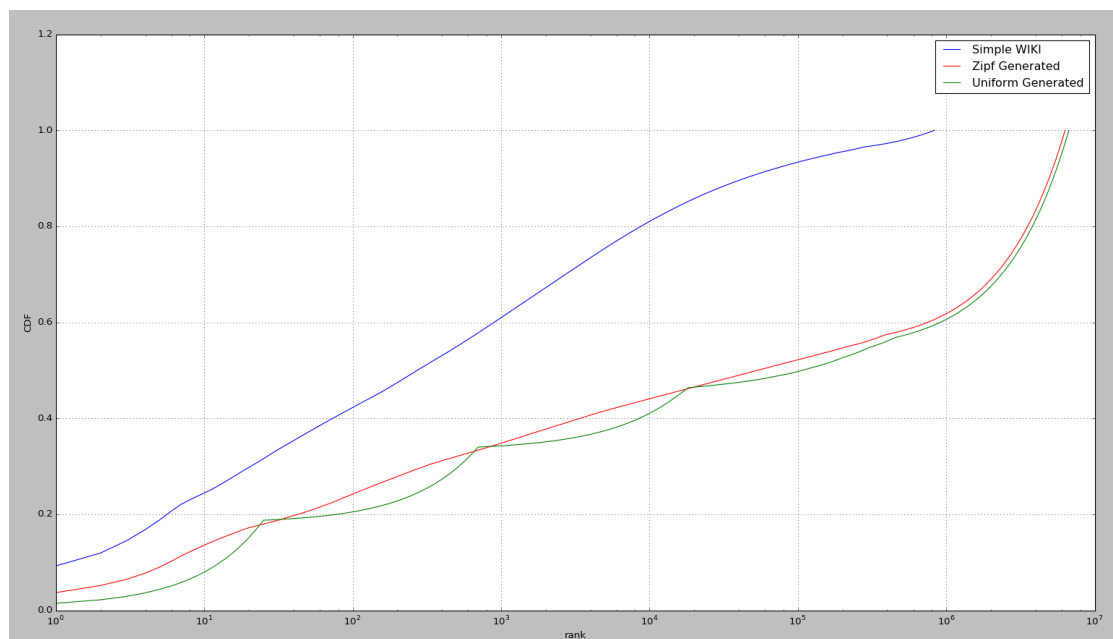
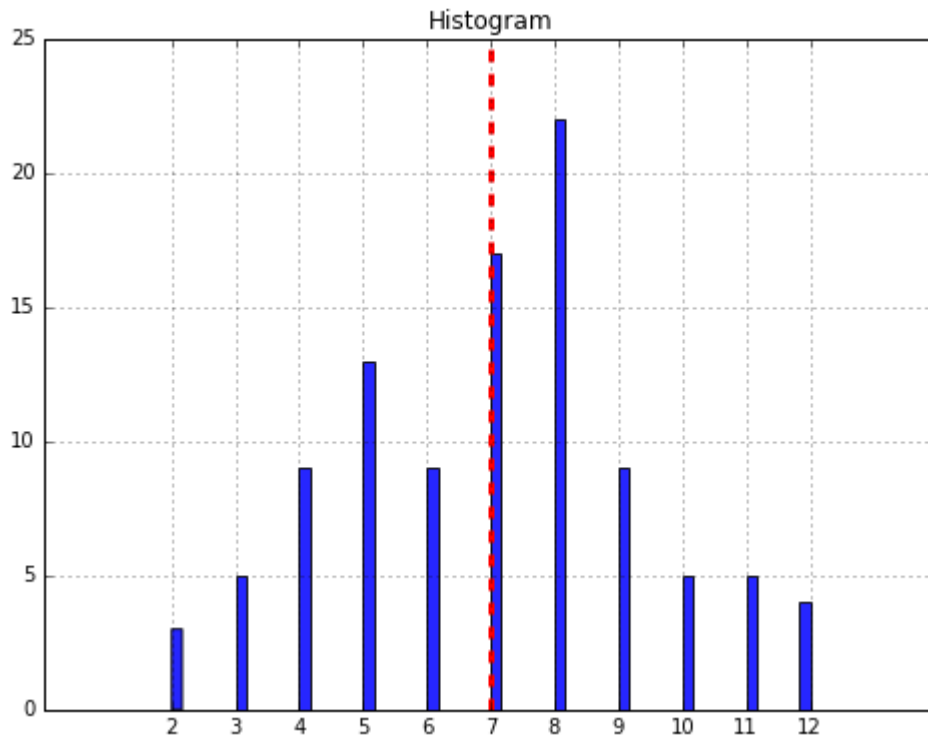


Figure 2: CDF plot for all 3 data sets

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.

answer : 7

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

answer : 0.78

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

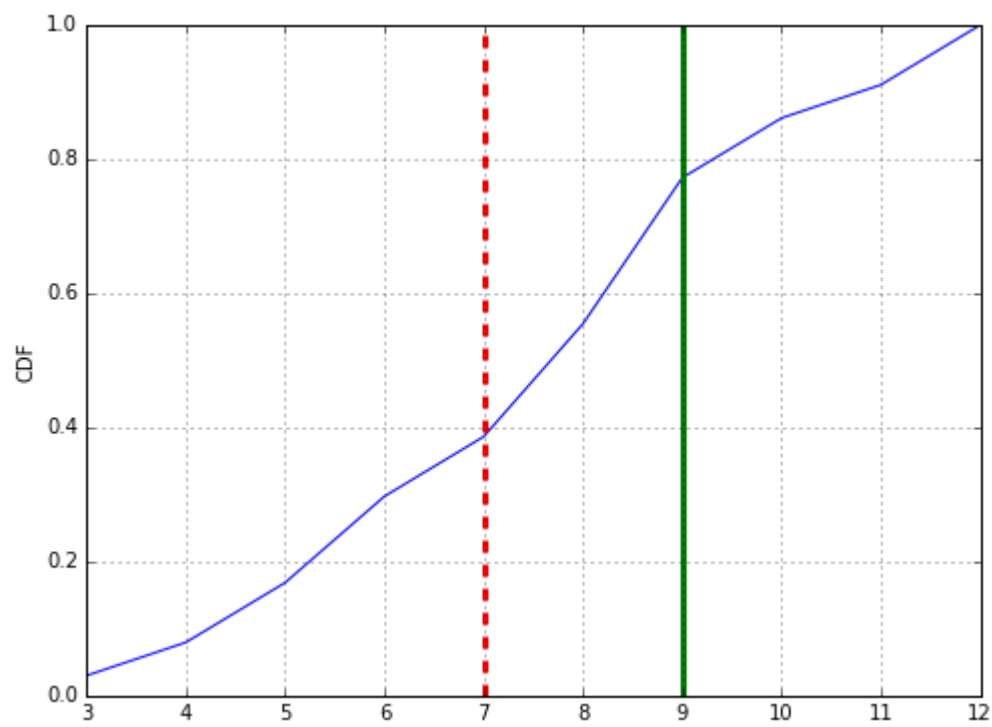


Figure 3: CDF plot for 100 Toss

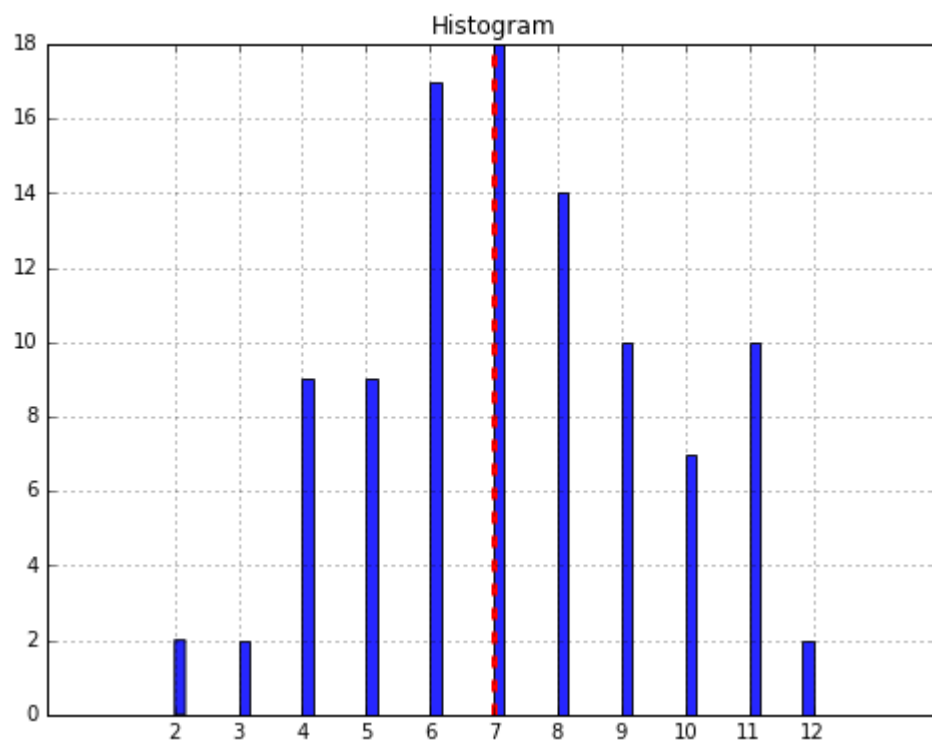


Figure 4: Hist plot for 100 Toss second simulation

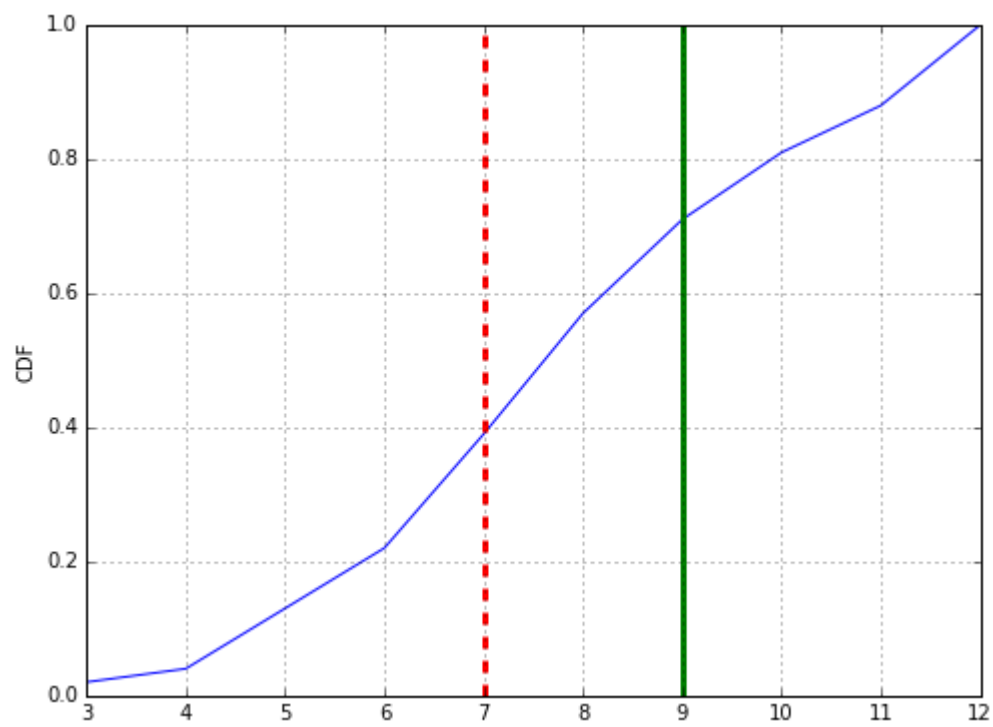


Figure 5: CDF plot for 100 Toss second simulation

Answer: 0.1

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

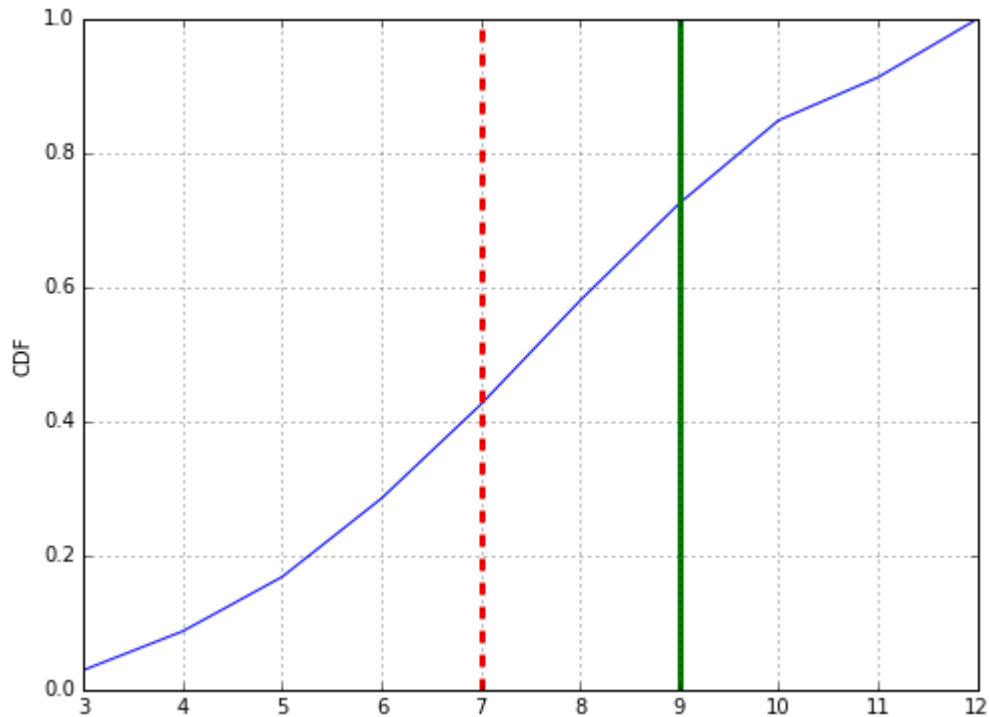


Figure 6: CDF plot for 1000 Toss first simulation

Answer: 0.035

6. What conclusion can you draw from increasing the number of steps in the simulation?
Answer: The maximum point wise distance in 100 toss is greater than 1000 toss. So we can conclude that when we increase the tosses the result would get closer together therefore the distance would get smaller. And also the distribution in histogram is more likely in 1000 rather than 100. Median for both cases is 7.

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.

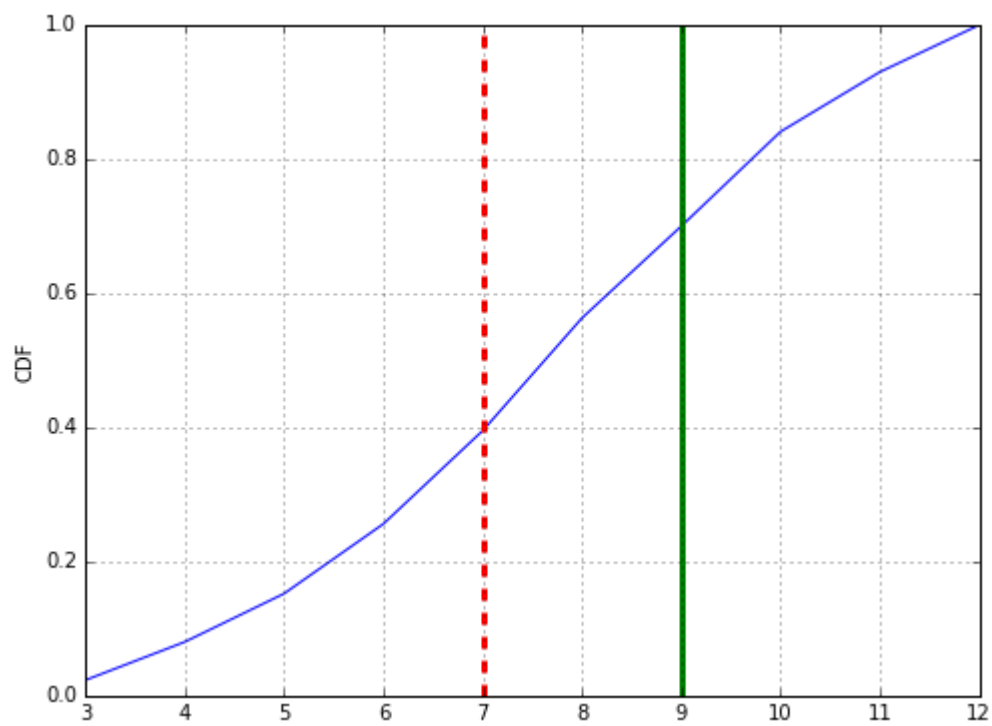


Figure 7: CDF plot for 1000 Toss second simulation

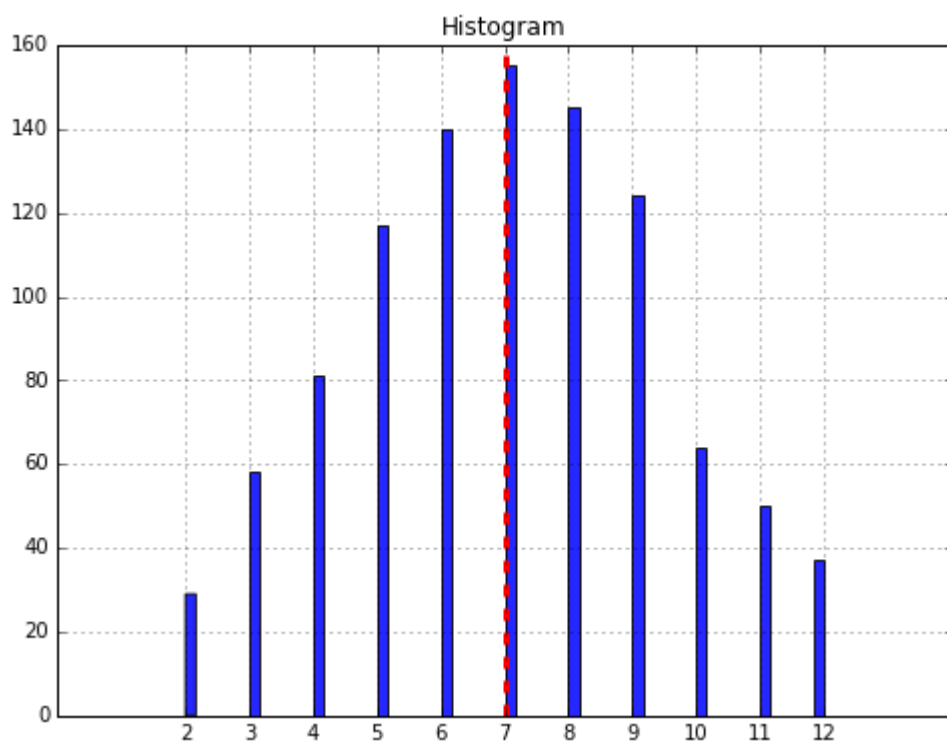


Figure 8: Hist plot for 1000 Toss first simulation

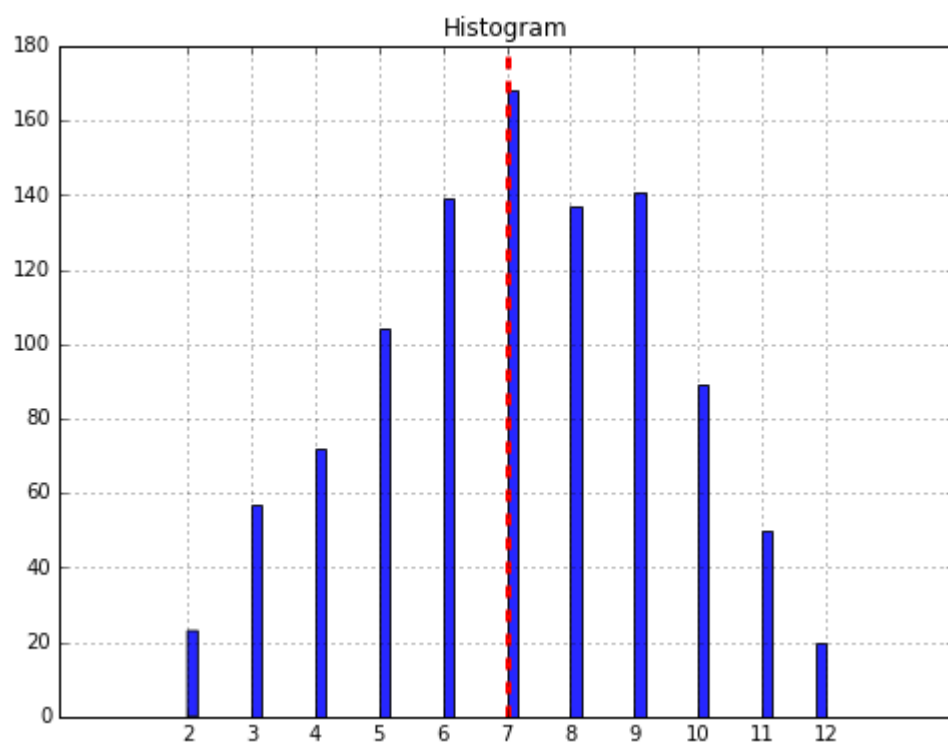


Figure 9: Hist plot for 1000 Toss second simulation

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 $n (=100)$?

Solutions for Question 1

Our Modeling choice to solve Question 1 was to ignore non significant words from the documents D1,D2,D3 (we have discussed this approach in the flipped classroom session).

D1= text about web science D2= web science covering analysis text corpora

D3=scientific methods used analyze web pages

ANS 1.1 : By analyzing the semantics of Documents D1,D2,D3 we can say that documents D1 and D2 are said to be similar. This is because if we look at the semantics of Document D1 which is about web science field. Similarly document D2 also refers to web science. However document D3 any no possible correlation between D1 and D2. Hence it can be said that semantically Document D1 and D2 are most similar.

Ans 1.2.1 : 12 base vectors would be needed to model the document

Ans 1.2.2 : Each dimension in the vectors space stands for a unique word observed in the document corpus.

Ans 1.2.3 : the Vector space will have 12 Dimensions

Ans 1.2.4 :

	About	Analysis	Analyze	Corpora	Covering	Methods	Science	Scientific	Text	Used	Web	webpages
1	1	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	1	0	0	0	0	0
0	0	0	0	0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	1

Figure 10: map of Word to base vector

	D1	D2	D3
About	1		
Analysis		1	
Analyze			1
Corpora		1	
Covering		1	
Methods			1
Science	1	1	
Scientific			1
Text	1	1	
Used			1
Web	1	1	
webpages			1

Figure 11: term frequency of words

Ans 1.2.5 and 1.2.6 please refer Figure 3-7

Ans 1.2.7 : According to the cosine similarity calculated , document D1,D2 are more similar

Ans 1.3 : Yes,the results of the model matches as per the results of the sub task 1.2. The vector space model matches the similarity given from the semantics of the documents as vector space model represents documents in a collection as a point in space and points that are semantically similar are close to each other.By calculating the similarity (in our case cosine similarity) between two documents we can conclude which documents are similar and which are not. In our model document D1,D2 were similar as the cosine similarity between them was 0.61 as compared to (D2,D3) and (D1,D3) were 0,0 respectively.

Handwritten mathematical derivation on lined paper:

Top left: $d_2 =$ (vertical vector of 0s)

Top right: $d_2 =$ (vertical vector of 0s)

Center: $d_2 =$ (sum of six terms, each with a coefficient and a vector)

Bottom left: $d_2 =$ (vertical vector of 0s)

Figure 14: 3

Ans 1.2.6

cosine similarity between $(\vec{d}_1, \vec{d}_2, \vec{d}_3)$

$\cos \theta = \frac{\langle \vec{d}_1, \vec{d}_2 \rangle}{\ \vec{d}_1\ _2 * \ \vec{d}_2\ _2}$ $= \frac{3}{\sqrt{4} * \sqrt{6}}$ $= \frac{3}{4.8989}$ $\cos \vec{d}_1, \vec{d}_2 = 0.61$	$\cos \theta = \frac{\langle \vec{d}_2, \vec{d}_3 \rangle}{\ \vec{d}_2\ _2 * \ \vec{d}_3\ _2}$ $\cos \theta = \frac{0}{\sqrt{6} * \sqrt{5}}$ $= 0$ $\cos \vec{d}_2, \vec{d}_3 = 0$
--	--

Figure 15: 4

Handwritten work on a spiral notebook showing cosine similarity calculations for vectors \vec{d}_1 , \vec{d}_2 , and \vec{d}_3 .

The calculation for $\cos \theta$ between \vec{d}_1 and \vec{d}_3 is shown as:

$$\cos \theta = \frac{\langle \vec{d}_1, \vec{d}_3 \rangle}{\|\vec{d}_1\|_2 * \|\vec{d}_3\|_2}$$
$$= \frac{0}{\sqrt{4} * \sqrt{5}}$$

Below this, it is noted that $\vec{d}_1, \vec{d}_3 = 0$.

Then, the cosine similarity between the vectors is calculated:

$$(\vec{d}_1, \vec{d}_2) = 0.61$$
$$(\vec{d}_2, \vec{d}_3) = 0$$
$$(\vec{d}_1, \vec{d}_3) = 0$$

Finally, a conclusion is drawn based on the calculations:

Ans 1.2.7
According to cosine similarity calculated document D_1 & D_2 are more similar.

Figure 16: 5

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".

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^AT_EX engine to LuaLaTeX.