

# SNLP 2016

## Exercise 2

- 1) (1 point) Use set theory and the axioms defining a probability function to show that:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

- 2) (2 points) Consider the following joint probability distribution:

$x$	0	0	1	1
$y$	0	1	0	1
$p(X = x, Y = y)$	0.32	0.08	0.48	0.12

Compute the marginal distributions:  $P(X)$ ,  $P(Y)$  and the conditional distributions  $P(X|Y)$  and  $P(Y|X)$ . Are these random variables independent?

- 3) (2 points) Suppose the vocabulary has  $30k$  words. Let's say the average sentence length in your corpus is 5 words. How many parameters do you need to store and estimate in order to describe the probability distribution:  $P(w_1, w_2 \dots w_5)$ ?
- 4) (5 points) In this exercise you will compare the probability distributions  $P(W_i|W_{i-1} = \text{"of"})$  and  $P(W_i|W_{i-1} = \text{"the"})$ . The distribution of the words given the previous word is "of" or "the" respectively.
- Download the Brown corpus from the web [http://www.nltk.org/nltk\\_data/](http://www.nltk.org/nltk_data/) or through the python NLTK toolkit.
  - Tokenize and lowercase each token.
  - Estimate the conditional probability distributions  $P(W_i|W_{i-1} = \text{"of"})$  and  $P(W_i|W_{i-1} = \text{"the"})$  with maximum likelihood estimation.
  - The expected value of the function  $f(x) = -\log_2 P(x)$  is called the entropy of the probability distribution  $P(X)$ . It shows how unexpected an event is on average.  
 $E[-\log_2 P(X)] = \sum_{i=1}^N P(x_i) * -\log_2 P(x_i)$
  - Plot the frequency distribution (unnormalized frequency counts) or the probability distribution for the 50 most frequent tokens for both distributions. Based on the plots which distribution do you expect to have a higher entropy ?
  - Compute the entropy for both distributions and verify your guess.

## 1 Submission Instructions: Read carefully

- You can form groups of maximum 3 people.
- Submit only 1 archive file in the ZIP format with name containing the MN of all the team members, e.g.:

Exercise\_01\_MatriculationNumber1\_MatriculationNumber2\_MatriculationNumber3.zip

- Provide in the archive:

- your code, accompanied with sufficient comments,
- a PDF report with answers, solutions, plots and brief instructions on executing your code,
- a README file with the group member names, matriculation numbers and emails,
- Data necessary to reproduce your results <sup>1</sup>
- The subject of your submission mail must contain the string [SNLP] (including the braces) and explicitly denoting that it is an exercise submission, e.g:

[SNLP] Exercise Submission 01

- Depending on your tutorial group, send your assignment to the corresponding tutor:
  - Sedigheh Eslami: *eslami@mpi-inf.mpg.de*
  - Naszdi Kata: *b.naszadi@gmail.com*
  - Stephanie Lund: *stflund@gmail.com*

## 2 General Information

- In your mails to us regarding the tutorial please add the tag [SNLP] in the subject accompanied by an appropriate subject briefly describing the contents.
- Feel free to use any programming language of your liking. However we strongly advise in favor of Python, due to the abundance of available tools (also note that Python3 comes with an excellent native support of UTF8 strings).
- Avoid using libraries that solve what we ask you to do (unless otherwise noted).
- Avoid building complex systems. The exercises are simple enough.
- Do not include any executable files in your submission, as this may cause the e-mail server to reject it.
- **In case of copying, all the participants (including the original solution) will get 0 points for the whole assignment. Note: it is rather easy to identify a copied solution. Plagiarism is also not tolerated.**
- **Missing the deadline even for a few minutes, will result in 50% point reduction. Submission past the next tutorial, is not corrected, as the solutions will already be discussed.**
- **Please submit in your solutions necessary to support your claims. Failure to do so, might results in reduction of points in the relevant questions.**
- Each assignment has 10 points and perhaps some bonus points (usually 2 or 3). In order to qualify for the exams, you need to have 2/3 of the total points. For example, in case there are 12 assignments, you need to collect at least 80 out of the 120 points to be eligible for the exams. A person that gets 10 plus 2 bonus points in every exercise, needs to deliver only 7 assignments in order to be eligible for the exams, since  $7 \cdot 12 = 84$ .
- Attending the tutorial gives 2 points increase for the corresponding assignment.
- Exercise points (including any bonuses) guarantee only the admittance to the exam, however have no further effect on the final exam grade.

---

<sup>1</sup>If you feel that these files are beyond reasonable size for an email submission and also reasonably convenient, please provide a means for us to access them online