

# Worksheet 7: Document Analysis

## Statistical Analysis and Document Mining

MSIAM 1<sup>st</sup> year / ENSIMAG 2<sup>nd</sup> year

### Question 1

You are willing to analyse a document containing 1,000,000 words.

1. Let  $k = 10$ ,  $\beta = 0.5$ . Following the Heaps' law, what is the number of distinct words in the document?
2. It was found that the 7% of all words are "the" article. Following the Zipf's law, estimate the value of  $\lambda$ . What is the frequency of the second most frequent word? The third most frequent one?

### Question 2

We have the following pre-processed training documents:

- $d_1 : \{ \text{"cat"}, \text{"sit"}, \text{"mat"}, \text{"cat"}, \text{"jump"}, \text{"bed"}, \text{"cat"}, \text{"good"}, \text{"sit"} \}$
- $d_2 : \{ \text{"cat"}, \text{"dog"}, \text{"jump"}, \text{"sit"}, \text{"dog"}, \text{"cat"}, \text{"animal"} \}$
- $d_3 : \{ \text{"table"}, \text{"sit"}, \text{"write"}, \text{"think"}, \text{"good"}, \text{"book"}, \text{"dog"}, \text{"sit"} \}$

1. Can you guess the context of each document?
2. What is the tf-idf weights for the following words: "cat", "dog", "sit", "animal", "book"?
3. Comment the question above.

### Question 3

1. Represent in the compressed sparse row (CSR) format the following matrix:

$$\mathbf{X} = \begin{pmatrix} 0 & 1 & 0 & 0 & 4 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 2 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 2 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 4 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 4 & 3 & 0 & 1 & 1 \end{pmatrix}.$$

2. Write a pseudo-code to multiply a CSR sparse matrix  $A$  by a (dense) vector  $b$ .