

## Lab Exercise 1

Lab Date: 22nd of January 2021

## 1. Guidelines

Deadline for submitting your solution: **23:59 4th of February 2021.**

Submission is a zipped folder with `{your name}` containing:

- Source files (Python): format name as `Lab{LabNumber}-{Exercise Number}.py`.
- A summary/report file formatted as a pdf which explains and presents the results with respect to the input values.

## 2. Exercises

**Exercise 1:** With the following list of words

`['she', 'sells', 'sea', 'shells', 'by', 'the', 'sea', 'shore']` perform the following tasks:

- Print all words beginning with *sh*
- Print all words longer than four characters

**Exercise 2:** Read in the texts of the State of the Union addresses, using the *state\_union* corpus reader.

- Count occurrences of *men*, *women*, and *people* in each document.
- Explain what has happened to the usage of these words over time using graphs.

**Exercise 3:** *Pig Latin* is a simple transformation of English text. Each word of the text is converted as follows: move any consonant (or consonant cluster) that appears at the start of the word to the end, then append *ay*, e.g. *string* → *ingstray*, *idle* → *idleay*

[http://en.wikipedia.org/wiki/Pig\\_Latin](http://en.wikipedia.org/wiki/Pig_Latin)

- Write a function to convert a word to Pig Latin.
- Write code that converts text, instead of individual words.
- Explain how you would decode Pig Latin, either with words or code.

**Exercise 4:** Reddit is a social media where people can share posts. Using the old web-page `old.reddit.com` as the source, build a web scrapper and present the following information about the posts.<sup>1</sup>

- Amount of upvotes
- Which subreddit the post was made in
- The title of the post
- The time when the post was made

**Exercise 5:** Write a program to build your own custom *Corpus* (to a folder on your local computer) from tweets (e.g., using the results of 10 search queries or tweets from 10 accounts) on Twitter.<sup>2</sup>

- Write code to tokenize the tweets and remove stopwords.
- Print out the 10 most common words in the corpus.
- Print out the 10 most common words from each user/search.
- Print out the 10 most used hashtags in the corpus.

<sup>1</sup>If the page is down or you have issues with the scrapping of the webpage, you can find a dump here.

<sup>2</sup>If you have issues getting tweets from Twitter, you can pick terms from this set.