# Introduction to Web Science

## Assignment 6

Prof. Dr. Steffen Staab        René Pickhardt

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

Korok Sengupta

koroksengupta@uni-koblenz.de

Institute of Web Science and Technologies
Department of Computer Science
University of Koblenz-Landau

Submission until:   December 6, 2016, 10:00 a.m.
Tutorial on:   December 9, 2016, 12:00 p.m.

Please look at the lessons 1) **Simple descriptive text models** & 2) **Advanced descriptive text 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.

Group name: uniform
Group members: Pradip Giri, Jalak Arvind Kumar Pansuriya, Madhu Rakhal Magar

# 1 Digging deeper into Norms (10 points)

You have been introduced to the concept of a norm and have seen that the uniform norm $||\cdot||_\infty$ fullfills all three axioms of a norm which are:

1. Positiv definite

2. Homogeneous

3. Triangle inequality

Recall that for a function $f : M \longrightarrow \mathbb{R}$ with $M$ being a finite set[1] we have defined the $L_1$-norm of $f$ as:

$$||f||_1 := \sum_{x \in M} |f(x)| \qquad (1)$$

In this exercise you should

1. calculate $||f - g||_1$ and $||f - g||_\infty$ for the functions $f$ and $g$ that are defined as

   - $f(0) = 2, f(1) = -4, f(2) = 8, f(3) = -4$ and
   - $g(0) = 5, f(1) = 1, g(2) = 7, g(3) = -3$

2. proof that all three axioms for norms hold for the $L_1$-norm.

## 1.1 Hints:

1. The proofs work in a very similar fashion to those from the uniform norm that was depicted in the videos.

2. You can expect that the proofs for each property also will be "three-liners".

3. Both parts of this exercise are meant to practice proper and clean mathematical notation as this is very helpfull when reading and understanding research papers. Discuss in your study group not only the logics of the calculation and the proof (before submission) but try to emphasize on the question whether your submission is able to communicate exactly what you are doing.

**Answer:**

1.1 We know that for a function $f : M \longrightarrow \mathbb{R}$ with $M$ being a finite set, we have defined the $L_1$-norm of $f$ as:

$$||f||_1 := \sum_{x \in M} |f(x)| \qquad (2)$$

---

[1]You could for example think of the function measuring the frequency of a word depening on its rank.

So, in our case:

$$||f - g||_1 := \sum_{x \in M} |f(x) - g(x)| \tag{3}$$

= |f(0) - g(0)| + |f(1) - g(1)| + |f(2) - g(2)| + |f(3) - g(3)|
= |2 - 5| + |-4 - 1| + |8 -7| + |-4 - -3|
= 3 + 5 + 1 + 1
= 10

Also,

$$||f - g||_\infty := \sum_{x \in M} max(|f(x) - g(x)| : 0...M) \tag{4}$$

= max(|f(0) - g(0)| , |f(1) - g(1)| , |f(2) - g(2)| , |f(3) - g(3)|)
= max(3, 5, 1, 1)
= 5

Hence, we got

$$||f - g||_1 := 10 \tag{5}$$

and

$$||f - g||_\infty := 5 \tag{6}$$

1.2 Here is the solution:

## 1.2

For a function $f : M \to R$ with $M$ being a finite set, we can defined $L_1$-norm of $f$ as:

$$||f||_1 := \sum_{x \in M} |f(x)| \quad\text{———}\quad \textcircled{1}$$

To prove all three axioms for norms holds for $L_1$-norm, Firstly take:

① Positive Definite $\left(||f||_1 = 0 \Rightarrow f = 0\right)$

    Proof:

$$||f||_1 = 0 \iff \sum_{x \in M} |f(x)| = 0$$

$$\Rightarrow |f(x)| = 0 \, \forall x$$

$$\Rightarrow f(x) = 0 \, \forall$$

$$\Rightarrow f = 0 \quad \text{proved//}$$

② Homogeneous $\left(||\alpha f||_1 = \alpha ||f||_1 \; ; \; \alpha \in R\right)$

· Proof:

$$||\alpha f||_1 = \sum_{x \in M} |\alpha f(x)|$$

$$= |\alpha| \sum_{x \in M} |f(x)|$$

$$= |\alpha| \, ||f||_1 \qquad \therefore \alpha \in R$$

proved

③ Triangle Inequality $\left( ||f+g||_1 \leq ||f||_1 + ||g||_1 \right)$

Proof:

$$||f+g||_1 = \sum_{x \in M} | f(x) + g(x) |_1$$

$$\leq \sum_{x \in M} |f(x)| + \sum_{x \in M} |g(x)|$$

$$\leq ||f||_1 + ||g||_1$$

$$= ||f||_1 + ||g||_1$$

Proved

## 2 Coming up with a research hypothesis (12 points)

You can find all the text of the articles from Simple English Wikipedia at `http://141.26.208.82/simple-20160801-1-article-per-line.zip` each line contains one single article.

In this task we want you to be creative and do some research on this data set. The ultimate goal for this exercise is to practice the way of coming up with a research hypothesis and testable predictions.

In order to do this please **shortly**[2] answer the following questions:

1. What are some obervations about the data set that you can make? State at least three obervations.

2. Which of these observations make you curious and awaken your interest? Ask a question about why this pattern could occur.

3. Formulate up to three potentiel research hypothesis.

4. Take the most promesing hypothesis and develop testable predictions.

5. Explain how you would like to use the data set to test the prediction by means of descriptive statistics. Also explain how you would expect your outcome.

   (If you realize that the last two steps would not lead anywhere repeat with one of your other research hypothesis.)

### 2.1 Hints:

- The first question could already include some diagrams (from the lecture or ones that you did yourselves).

- In step 3 explain how each of your hypothesis is falsifiable.

- In the fifth step you could state something like: "We expect to see two diagrams. The first one has ... on the x-axis and ... on the y-axis. The image should look like a ... The second diagram ...". You could even draw a sketch of the diagram and explain how this would support or reject your testable hypothesis.

**Answer:**
2.1 Here is some observations that we made:

1. Each article in simple wiki is separated by newline.

2. Each article in simple wiki is from diverse field and describes about some particular unique topic.

---

[2]Depending on the question shortly could mean one or two sentences or up to a thousand characters. We don't want to give a harsh limit because we trust in you to be reasonable.

3. Each article in simple wiki is informative.
   (Assumption: An informative article informs the reader by explaining and giving details on a given topic. The informative article is not an argumentative article that tries to persuade the reader to one side or the other. It covers all the pertinent details: who, what, when, where and why.)

4. Simple wiki consist of certain number of foreign(loan) words.
   (Foreign(loan) words: is a word used directly from another language with little or no translation.)

2.2 "Simple wiki consist of certain number of foreign(loan) words." This observation make us curious.
Is this really true?
Each article tells us about one particular topic and gives the background of it.
The description of article provides specific information, data, facts, definition of terms, date, origin of words which are organized and presented within same context.
So there must be use of some foreign words in the article

2.3 Here are three potential research hypothesis:

1. The simple wiki consists of at least 10% of foreign words(or loan words).

2. Each article in simple wiki is informative.

3. Each article in simple wiki consists of at least 200 words.

Here, we can say that hypothesis is falsifiable if we got the counter example or counter facts for it.
For the first hypothesis, It will be falsifiable if we got less than 10% of foreign words no foreign words after doing the research.
For second hypothesis it will be falsifiable if there is no any kind of data, facts,date etc in article. Also, If it contains dummy text.
For third hypothesis we can say that it is falsifiable if any article in simple wiki is less than 200 words.

2.4 We found that "The simple wiki consists of at least 10% of foreign words(or loan words)." is the most promising hypothesis.

**Testable Predictions:**

1. Each article comes from diverse field and talks about unique particular topic.

2. The articles cover older information, usually from the beginning to a historic ancient

times.Which are from different parts of the world. Our assumption is "to describe the diversified topics simple wiki use some foreign words."

This means we have to find out the foreign words in English and compare these words with the words in the simple wiki.
-Example of counting the words.

2.5 Firstly, We are going to get the all foreign words in English from Wikipedia. After that we compare each words in an article with foreign words that we have. We will save number of foreign words in an article. After that we repeat the comparison process for all article.
We expect to see two diagrams. the first one is Cumulative Frequency Diagram(CDF) which has % of foreign words in y-axis and no.of article in x-axis. If this diagram shows the number of foreign words is greater than 10% then our hypothesis will be proved.
The second diagram will be a scatter plot where we plot the foreign words and total other words in the article.
We will also calculate the mean, median and percentage of foreign words in simple wiki.

# 3 Statistical Validity (8 points)

In the above question, you were asked to formulate your hypothesis. In this one, you should follow your own defined roadmap from task 2 validate (or reject) your hypothesis.

**Answer:**
**Observation:**
After reading first few articles we observe that some of the word are from medieval period and early ages with historical background. And found some of the word used by Roman period which suggest us that there must be foreign English word which are borrowed from other languages. We observed there is at least 10% words which are loaned to English language from other language.

**Thinking of Interesting Questions:**
We are interested in finding how other language help(influence) in development of English language? Are they directly used in English Language or with some modifications?.What is the percentage of foreign words in simple wiki?

**Hypothesis Formulation:**
The simple wiki consists of at least 10% of foreign words(or loan words).

**Development of Testable Predictions:**
If our hypothesis is correct, We will get at least 10% of loanwords in simple wiki.

**Data Collection:**
First we collect all the required loaned word from `https://en.wikipedia.org/wiki/Lists_of_English_words_by_country_or_language_of_origin` Wikipedia and clean it and make it unique.

**Test Prediction:**
We load all the articles and select only article having more than 15 words because we think an article must have minimum 15 words. And we check each and every word of the given article against our set of English loaned word. During comparison we trim the word so that space will not effect our comparison.

**Test Result:**
In the process we use 12513 unique loan(foreign) words. After calculation we found that the mean of the loaned word in the given articles is 18.439. Which means every article have 18 loaned word or There is 18.43% loanwords in the given simple wiki file. The median of the calculation is 8.0, which means 50% of the article have less or equal to 8 loaned word and 50 % of the article have more than 8 loan words.

**Listing 1** download.py

```
 1: # pylint: disable-msg=C0103
 2: """ Downloads all the list form given url"""
 3: from urllib.request import urlopen
 4: from urllib.error import HTTPError
 5: from urllib.parse import urlparse
 6: from bs4 import BeautifulSoup
 7:
 8: soup = None
 9: restricted_names = ["place names", "references",
10:     "see also",
11:     "Common in South African English",
12:     "External links"]
13:
14: def save_content(file_name, content):
15:     """ save the file with given name and givn content """
16:     with open(file_name, 'w+') as f:
17:         f.write(content)
18:     print('{} written sucessfully'.format(file_name))
19:
20:
21: def download_url(input_url):
22:     """ returns html content if request is successful and None for other
23:     responses"""
24:     try:
25:         html = urlopen(input_url)
26:     except HTTPError:
27:         return None
28:     return html
29:
30:
31: def find_inside_given_id(soup, selector):
32:     """ find inside the given class"""
33:     uls = soup.select(selector)
34:     borrowed_word = set()
35:     for ul in uls:
36:         lis = ul.select("li")
37:         print(lis)
38:         for li in lis:
39:             a = li.select("a")[0]
40:             borrowed_word.add(a.text)
41:
42:     return borrowed_word
43:
44:
45: def find_next_to_selector(soup, selector):
46:     """ find next to the given class name """
47:     spans = soup.select(selector)
48:     borrowed_word = set()
```

```
49:        for span in spans:
50:            if span.string.lower() not in restricted_names:
51:                uls = span.parent.findNext("ul")
52:                lis = uls.select("li")
53:                for li in lis:
54:                    found_workd = li.find("a").text
55:                    borrowed_word.add(found_workd)
56:
57:                dls = span.parent.findNext("dl")
58:                if dls is not None:
59:                    dts = dls.select("dt")
60:                    for dt in dts:
61:                        atags = dt.select("a")
62:                        for a in atags:
63:                            borrowed_word.add(a.text)
64:        return borrowed_word
65:
66:
67: def process_html(input_html):
68:     """ process the downloaded html"""
69:     soup = BeautifulSoup(input_html, 'html.parser')
70:     spans = soup.select('span.mw-headline')
71:     length = len(spans)
72:     print(length)
73:     if length > 2:
74:         return find_next_to_selector(soup, 'span.mw-headline')
75:     else:
76:         return find_inside_given_id(soup, '#mw-content-text')
77:
78:
79: def extract_file_name(input_url):
80:     """returns the required file name extracting from url"""
81:     path = urlparse(input_url).path
82:     return path.split('/')[-1].strip() + '.txt'
83:
84:
85: def main(input_url):
86:     """ Entry of the programme """
87:     global restricted_names
88:     restricted_names = list(map(lambda x: x.lower().strip(), restricted_names))
89:     html = download_url(input_url)
90:     if html is not None:
91:         print("Processing {}".format(input_url))
92:         html_content = html.read().decode()
93:         words = list(process_html(html_content))
94:         file_name = extract_file_name(input_url)
95:         save_content(file_name, str(words))
96:
97:
```

```
 98: if __name__ == "__main__":
 99:     url = 'https://en.wikipedia.org/wiki/List_of_English_words_of_French_origin_(8
100:     main(url)
```

**Listing 2** assignment6.py

```
 1: """
 2: How to use it
 3: first run the script
 4: python3 assignment6.py
 5: if you have already run script before you can only plot by using
 6: python3 assignment6.py plot
 7: """
 8: import os
 9: import sys
10: import glob
11: import numpy as np
12: import matplotlib.pyplot as plt
13:
14:
15: path = os.getcwd() + '/' + 'foreign-words'
16: article_file = os.getcwd() + '/' + 'simple-20160801-1-article-per-line'
17: all_foreign_words = set([])
18:
19: # this is list wich will contain number of foreign words
20: artile_infos_with_foreign_words = []
21:
22: def open_file(file_name):
23:     file_content = None
24:     with open(file_name, 'r+') as f:
25:         file_content = f.read()
26:     return file_content
27:
28: def process_foreign_text_file(file_name):
29:     file_content = open_file(file_name)
30:     list_of_words = eval(file_content)
31:     for word in list_of_words:
32:         all_foreign_words.add(word.strip())
33:
34: def process_articles(articles):
35:     all_articles = articles.split('\n')
36:     for id, article in enumerate(all_articles):
37:         words_arr = article.split(" ")
38:         #article having less than 15 words are discarded
39:         if len(words_arr) > 15:
40:             found_world = 0
41:             for word in words_arr:
42:                 if word.strip() in all_foreign_words:
43:                     found_world += 1
44:             artile_infos_with_foreign_words.append({"id": id, "loan_word": found_w
```

```
45:
46: def mean(arr):
47:     return np.mean(arr)
48:
49: def median(arr):
50:     return np.median(arr)
51:
52:
53: def plot_graph(data):
54:     num_bins = 100
55:     counts, bin_edges = np.histogram(data, bins=num_bins, normed=True)
56:     cdf = np.cumsum(counts)
57:     plt.plot(bin_edges[1:], cdf)
58:     plt.xlabel('Article')
59:     plt.ylabel('Loaned World %')
60:     plt.show()
61:
62:
63: def main():
64:     all_txt_files = glob.glob(path + '/*.txt')
65:     for text_file in all_txt_files:
66:         process_foreign_text_file(text_file)
67:     artile_content = open_file(article_file)
68:     process_articles(artile_content)
69:     with open("result.txt", "w+") as f:
70:         sorted_list = sorted(artile_infos_with_foreign_words, key=lambda info: in
71:         f.write(str(sorted_list))
72:     arr = list(map(lambda x: (x['total_word'], x['loan_word']), artile_infos_with
73:     arr = list(map(lambda x: x[1], arr))
74:     print("median", median(arr))
75:     print("mean", mean(arr))
76:     print("unique loaned word used ", len(all_foreign_words))
77:     process_diagram()
78:
79:
80: def draw_scatter_plot(arr):
81:     total_word = [x[0] for x in arr]
82:     foreign_word = [x[1] for x in arr]
83:     N = len(total_word)
84:     colors = np.random.rand(N)
85:     plt.scatter(total_word, foreign_word, s=15, c=colors, alpha=1)
86:     plt.ylabel('Loaned Word')
87:     plt.xlabel('Total Words')
88:     plt.xlim(0, 10000)
89:     plt.ylim(0, 1500)
90:     plt.show()
91:
92: def draw_histogram(data):
93:     plt.hist(data, bins= np.arange(0, 1300, 5), normed=True)
```

```
 94:        plt.ylabel('Loaned Word Frequency')
 95:        plt.xlabel('Number of Loaned Word')
 96:        plt.show()
 97:
 98: def process_diagram():
 99:        raw_data = None
100:        with open("result.txt", "r+") as f:
101:            raw_data = eval(f.read())
102:        if raw_data:
103:            arr = list(map(lambda x: (x['total_word'], x['loan_word']), raw_data))
104:            plot_graph(arr)
105:            draw_scatter_plot(arr)
106:            mod_arr = list(map(lambda x: x[1], arr))
107:            draw_histogram(mod_arr)
108:
109: if __name__ == "__main__":
110:        plot = None
111:        try:
112:            plot = sys.argv[1]
113:        except:
114:            pass
115:        if plot == "plot":
116:            process_diagram()
117:        else:
118:            main()
```



hello

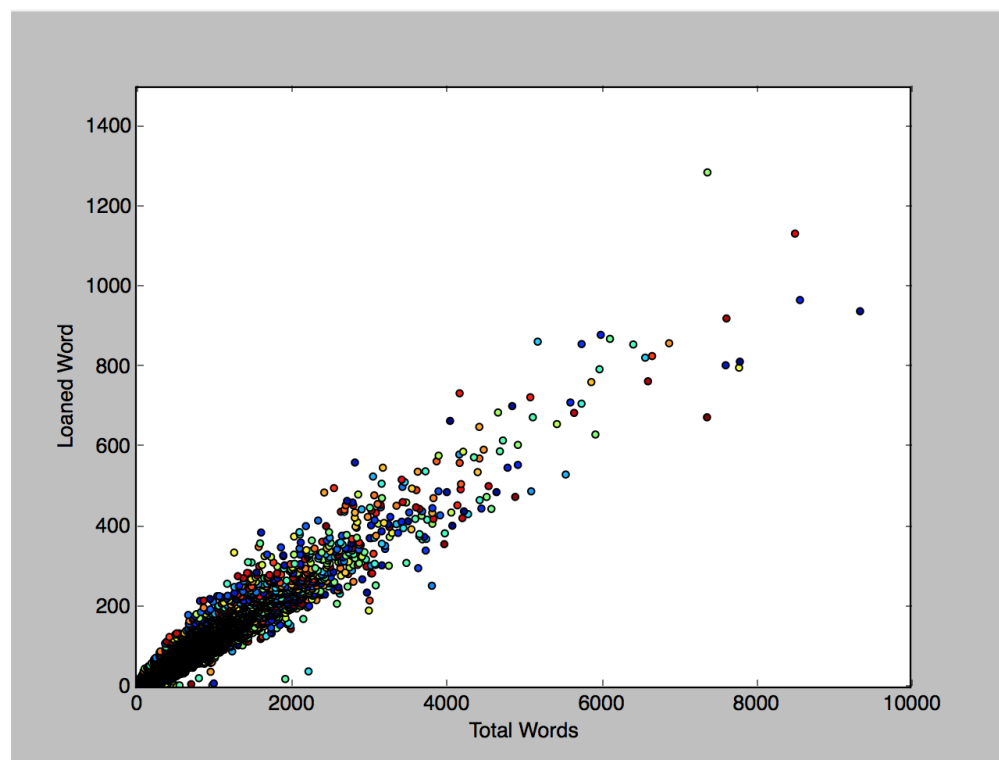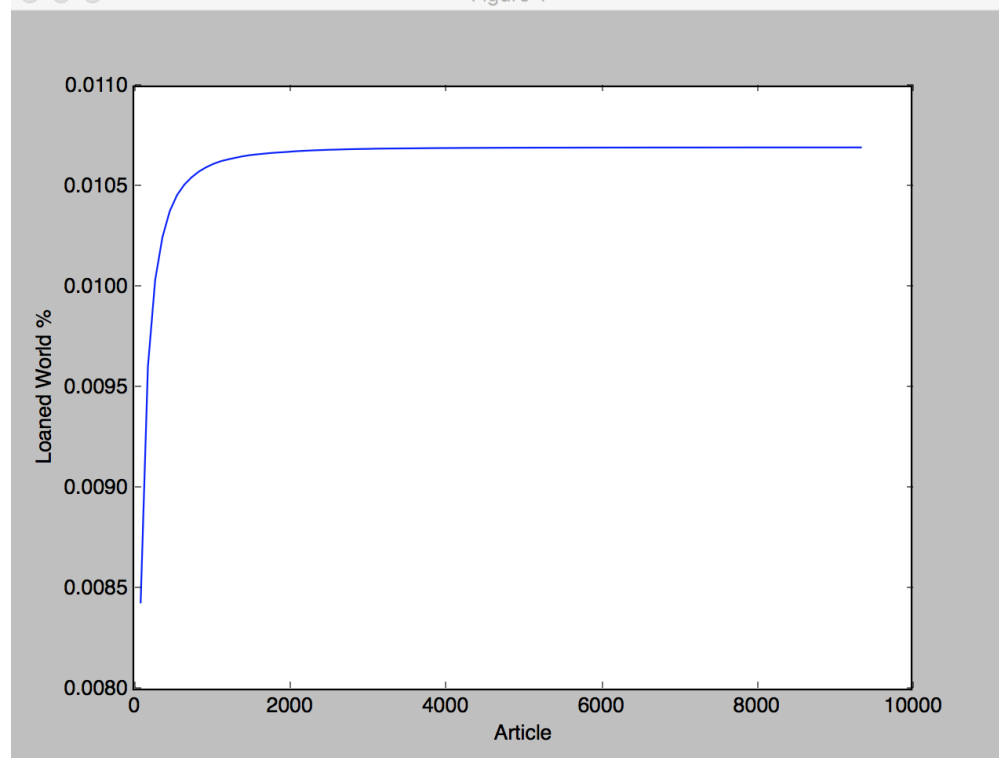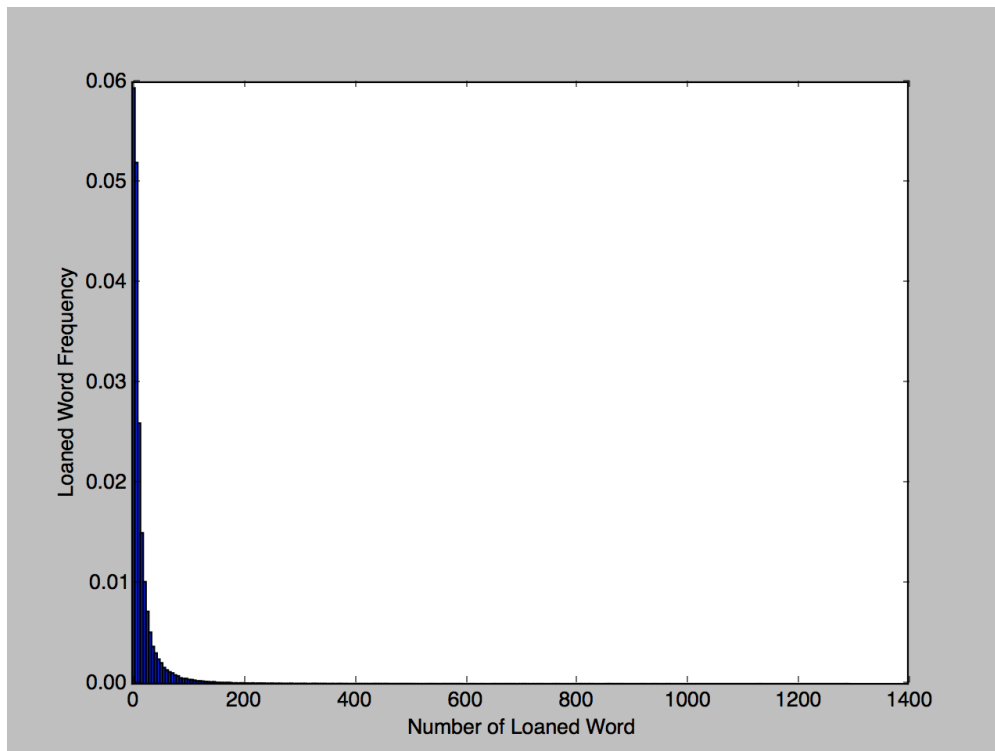**Conclusion**

From calculated mean, median and the graph(CDF, histogram and scatter plot) we found that simple wiki consist of at least 10% of the loaned English word.

**Problem faced**

1. Data collection
   During collection the list of loaned word. We crawl the wikidedia page but the pages have different HTML structure which take a lot of modification in the downloading script for specific page.

2. Length of Paragraph
   We assume the separation of the articles are by
   n and while processing the articles some are only blank lines and some of the articles have only few words like 3 and 7, 9 etc so we remove the blank lines and only takes the articles having more than 15 words.

## 3.1 Hints:

- In case feel uncomfortable to test one of the predictions from task 2 you can "steal" one of the many hypothesis (and with them imlicitly associated testable predictions)

or diagrams depicted from the lecture and reproduce it. However in that case you cannot expect to get the total amount of points for task 3.

# Important Notes

## Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment6/` 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 LaTeXengine to `LuaLaTeX`.