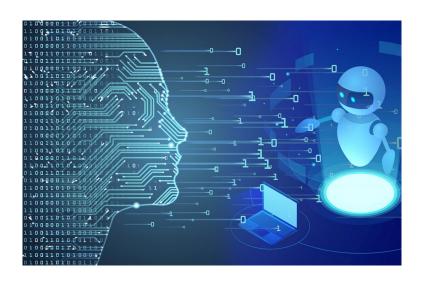
DURHAM UNIVERSITY

DEPARTMENT OF COMPUTER SCIENCE MASTER OF DATA SCIENCE



COMP42415 Text Mining and Language Analytics Workshops

Dr Stamos Katsigiannis

Preface

Welcome to the workshop manual for the COMP42415 Text Mining and Language Analytics course of Durham University. This manual has been meticulously crafted to serve as your comprehensive guide to the fascinating world of Natural Language Processing (NLP). As language plays a pivotal role in our daily lives, understanding and harnessing the power of natural language is increasingly crucial in various fields, ranging from artificial intelligence and data science to linguistics and computational linguistics.

This workshop is designed to provide you with a hands-on and in-depth exploration of NLP, equipping you with the knowledge and skills necessary to navigate the intricacies of natural language and to build practical applications. Whether you are a seasoned professional seeking to enhance your expertise or a beginner eager to delve into the realm of NLP, this manual is tailored to meet your needs.

Key Features of this Workshop Manual:

- 1. Comprehensive Coverage: This manual covers a wide array of topics, from the foundational concepts of NLP to advanced techniques and applications. Each section is structured to build upon the previous one, ensuring a logical progression of learning.
- 2. Hands-On Exercises: Learning by doing is at the core of this workshop. Throughout the manual, you will find hands-on exercises that allow you to apply the concepts you learn in real-world scenarios. These exercises are designed to reinforce your understanding and enhance your practical skills.
- 3. Real-World Applications: NLP has diverse applications, and this manual provides insights into how NLP techniques are used in real-world scenarios.
- 4. Do-It-Yourself: This manual is designed to support self-learning, providing detailed examples and explanations for the presented NLP techniques.

By the end of this workshop, you will have gained a solid understanding of NLP concepts, developed practical skills in implementing NLP solutions, and be well-prepared to tackle challenges in the ever-evolving landscape of NLP

We hope you find this workshop manual informative, engaging, and instrumental in your quest to master Natural Language Processing.

Best wishes for a rewarding learning experience!

Dr Stamos Katsigiannis

iv Preface

Required Python packages

Required python version: This workshop manual has been tested on Python 3.11.

The following Python packages are required:

- re
- \bullet nltk
- numpy
- \bullet scipy
- \bullet math
- \bullet scikit-learn
- \bullet matplotlib
- \bullet seaborn
- \bullet pandas
- \bullet tensorflow
- pickle

Contents

Pı	refac		iii
\mathbf{R}	equir	ed Python packages	v
1	Wo	kshop 1: Text query with Regular Expressions	1
_	1.1	Regular Expressions Definition	
	$1.1 \\ 1.2$	The "re" Python package	
	1.2	v 1	
		1.2.1 Example	
		1.2.2 RegEx functions in "re"	
		1.2.3 Regular expression metacharacters	
		1.2.4 Regular expression special sequences	
		1.2.5 Repeating regular expressions	
	1.3	Using Regular Expressions to match strings	
		1.3.1 Matching example	
		1.3.2 Matching to validate strings	5
		1.3.3 Credit card number validation	6
		1.3.4 Validation of United Kingdom's National Insurance numbers (NINO)	
		1.3.5 Validation of hexadecimal numbers	8
	1.4	Using Regular Expressions to search elements in files	ç
		1.4.1 Parsing XML files	ç
		1.4.2 Parsing HTML files	
		1.4.3 Parsing raw text	
	1.5	Using Regular Expressions to substitute strings	
	_	1.5.1 Substitution example	
		1.5.2 Email domain substitution	
	1.6	Exercises	
2	Wo	kshop 2: Text pre-processing	15
4	2.1	NLTK corpora	
	2.1	-	
		2.1.1 Import NLTK	
		2.1.2 Corpus file IDs	
		2.1.3 Corpus file categories	
		2.1.4 Corpus words	
		2.1.5 Selecting files from specific category	
		2.1.6 Corpus sentences	
		2.1.7 Accessing sentences from specific file	
		2.1.8 Number of documents in each category	
		2.1.9 Corpus raw text	17
2.2		Input text from text file	17
	2.3	Text pre-processing	18
		2.3.1 Sentence Tokenisation	18
		2.3.2 Word Tokenisation	18
		2.3.3 Lowercasing	19
		2.3.4 Stop words removal	20
		2.3.5 Punctuation removal	
		2.3.6 Stemming	
		2.3.7 Lemmatisation	
		2.3.8 Part of Speech (POS) tagging	
	2.4	Exercises	25

viii Contents

3	Wo	rkshop 3: Text representation	27
	3.1	Load corpus	27
	3.2	Vocabulary	
		3.2.1 Words in corpus	28
		3.2.2 Unique words in corpus	
		3.2.3 Vocabulary of multiple documents	29
	3.3	One-hot encoding	
		3.3.1 One-hot encoding of words in vocabulary	
		3.3.2 One-hot encoding of text	
	3.4	Term Frequency (TF) representation	32
		3.4.1 Compute TF of words in text	
		3.4.2 TF representation of documents	
	3.5	Term Frequency - Inverse Document Frequency (TF-IDF)	
		3.5.1 Document Frequency (DF)	
		3.5.2 Inverse Document Frequency (IDF)	
		3.5.3 Term Frequency - Inverse Document Frequency (TF-IDF)	
	3.6	Exercises	
\mathbf{A}	Apr	pendix: Test files	37
		movies.xml	37
	A.2	links.html	37
		emails.txt	
		cds.xml	
		alice.txt	
		dune tyt	42

Workshop 1: Text query with Regular Expressions

1.1 Regular Expressions Definition

A regular expression (shortened as regex or regexp) is a sequence of characters that define a search pattern. Usually such patterns are used by string-searching algorithms for "find" or "find and replace" operations on strings, or for input validation.

1.2 The "re" Python package

Python has a built-in package called "re", which can be used to work with Regular Expressions. Let's import this package and use a regular expression to check whether the sentence "I started studying this year at Durham University" ends with the word "University" or with the word "school".

1.2.1 Example

```
import re # Import the re package
txt = "I started studying this year at Durham University"
x1 = re.search("University$", txt) # Returns a Match object if there is a match anywhere in the
    string with the regex
x2 = re.search("school$", txt)
print("x1:",x1)
print("x2:",x2)
if(x1):
   print("The text ends with 'University'")
   print("The text does not end with 'University'")
   print("The text ends with 'school'")
   print("The text does not end with 'school'")
The output will look like:
x1: <_sre.SRE_Match object; span=(39, 49), match='University'>
x2: None
The text ends with 'University'
The text does not end with 'school'
```

1.2.2 RegEx functions in "re"

The re module offers a set of functions that allows us to search a string for a match:

Function	Description
findall(args)	Returns a list containing all matches
$\operatorname{search}(args)$	Returns a Match object if there is a match anywhere in the string
$\mathrm{match}(\mathit{args})$	Returns a Match object if there is a match starting at the beginning of the string
$\operatorname{split}(\mathit{args})$	Returns a list where the string has been split at each match
$\mathrm{sub}(\mathit{args})$	Replaces one or many matches with a string

1.2.3 Regular expression metacharacters

As you can see in our first example, we used the character "\$" in order to indicate that a text matching the regular expression should end with the string preceding the character "\$". For example, the regular expression "car\$" indicates that the text should and with the string "car". In this case, "\$" is considered as a metacharacter, i.e. a character with a special meaning. Below are the metacharacters supported by the "re" package:

Character	Description	Example
[]	A set of characters	"[a-f]"
\	Signals a special sequence (can also be used to escape special characters)	"\s"
	Any character (except newline character)	"unirsity"
^	Starts with	"^She"
\$	Ends with	"John\$"
*	Zero or more occurrences	"o*"
+	One or more occurrences	"l+"
?	Matches 0 or 1 repetitions of the preceding regex	ab? will match either "a" or "ab"
{}	Exactly the specified number of occurrences	"o{2}"
	Either or	"he she"
()	Capture and group	

The first metacharacters we'll look at are [and]. They're used for specifying a character class, which is a set of characters that you wish to match. Characters can be listed individually, or a range of characters can be indicated by giving two characters and separating them by a "-". For example, [abc] will match any of the characters a, b, or c; this is the same as [a-c], which uses a range to express the same set of characters. If you wanted to match only lowercase letters, your regex would be [a-z].

Metacharacters are not active inside classes. For example, [akm\$] will match any of the characters "a", "k", "m", or "\$". "\$" is usually a metacharacter, but inside a character class it is stripped of its special nature.

You can match the characters not listed within the class by complementing the set. This is indicated by including a "^" as the first character of the class. For example, [^5] will match any character except "5". If the caret appears elsewhere in a character class, it does not have special meaning. For example: [5^] will match either a "5" or a "^".

Perhaps the most important metacharacter is the backslash, \backslash . As in Python string literals, the backslash can be followed by various characters to signal various special sequences. It is also used to escape all the metacharacters so you can still match them in patterns. For example, if you need to match a [or \backslash , you can precede them with a backslash to remove their special meaning: \backslash [or \backslash].

Some of the special sequences beginning with "\" represent predefined sets of characters that are often useful, such as the set of digits, the set of letters, or the set of anything that isn't a white space.

1.2.4 Regular expression special sequences

Let's see some of the main regular expression special sequences. For a more detailed list, please refer to https://docs.python.org/3/library/re.html#re-syntax.

These sequences can be included inside a character class. For example, [\s,.] is a character class that will match any white space character, or "," or ".".

Sequence	Description
d	Matches any decimal digit. Equivalent to the class [0-9].
$\backslash D$	Matches any non-digit character. Equivalent to the class [^0-9].
$\setminus s$	Matches any white space character. Equivalent to the class [$\t \cdot \t $
$\backslash S$	Matches any non-white space character. Equivalent to the class $[\t^{r}]$.
$\setminus w$	Matches any alphanumeric character. Equivalent to the class [a-zA-Z0-9_].
$\setminus \mathbf{W}$	Matches any non-alphanumeric character. Equivalent to the class [^a-zA-Z0-9_].

1.2.5 Repeating regular expressions

Being able to match varying sets of characters is the first thing regular expressions can do that isn't already possible with the methods available on strings. However, if that was the only additional capability of regexes, they wouldn't be much of an advance. Another capability is that you can specify that portions of the regular expression must be repeated a certain number of times.

The first metacharacter for repeating things that we'll look at is "*". "*" does not match the literal character "*", but it specifies that the previous character can be matched zero or more times, instead of exactly once.

Example: "do*g" will match "dg" (zero "o" characters), "dog" (one "o" character), "doooog" (four "o" characters), and so on.

Repetitions such as "*", "+", and "?" are greedy. When repeating a regular expression, the matching engine will try to repeat it as many times as possible. If later portions of the pattern don't match, the matching engine will then back up and try again with fewer repetitions. If this behaviour is undesirable, you can add "?" after the qualifier ("*?", "+?", "??") to make it perform the match in non-greedy or minimal fashion, i.e. as few characters as possible will be matched.

Step-by-step example

Let's consider the expression a[bcd]*b. This matches the letter "a", zero or more letters from the class [bcd], and finally ends with a "b". Now imagine matching this regular expression against the string "abcbd".

Step	Matched	Explanation
1	a	The "a" in the regex matches.
2	abcbd	The engine matches "[bcd]*", going as far as it can, which is to the end of the string.
3	FAILED	The engine tries to match "b", but the current position is at the end of the string, so it fails.
4	abcb	Back up, so that "[bcd]*" matches one less character.
5	FAILED	Try "b" again, but the current position is at the last character, which is a "d".
6	abc	Back up again, so that "[bcd]*" is only matching "bc".
7	abcb	Try "b" again. This time the character at the current position is "b", so it succeeds.

The end of the regular expression has now been reached, and it has matched "abcb". This demonstrates how the matching engine goes as far as it can at first, and if no match is found it will then progressively back up and retry the rest of the regular expression again and again. It will back up until it has tried zero matches for "[bcd]*", and if that subsequently fails, the engine will conclude that the string does not match the regex at all.

Another repeating metacharacter is "+", which matches one or more times. Pay careful attention to the difference between "*" and "+". "*" matches zero or more times, so whatever's being repeated may not be present at all, while "+" requires at least one occurrence.

Example: "do+g" will match "dog" (one "o" character), "dooog" (three "o" characters), and so on, but will not match "dg" (zero "o" characters).

There are two more repeating qualifiers. The question mark character "?" matches either once or zero times. Think of it as marking something as being optional.

Example: "pre-?processing" matches either "preprocessing" or "pre-processing".

The most complicated repeated qualifier is " $\{m,n\}$ ", where m and n are decimal integers. This qualifier means there must be at least m repetitions, and at most n. For example, "a/ $\{1,3\}$ b" will match "a/b", "a//b", and "a///b", but it will not match "ab", which has no slashes, or "a///b", which has four slashes. You can omit either m or n. In this case, default values for m or n are used. Omitting m is interpreted as a lower limit of 0, while omitting n results in an upper bound of infinity.

Note: Some qualifiers are interchangeable. For example " $\{0,\}$ " is the same as "*", " $\{1,\}$ " is the same as "+", and " $\{0,1\}$ " is the same as "?". "*", "+", and "?" make the regular expression easier to read, so try to use them if possible.

1.3 Using Regular Expressions to match strings

1.3.1 Matching example

Let's use the text "I started studying this year at Durham University" again and find out whether the string "at" or the string "in" exists in the text.

```
txt = "I started studying this year at Durham University"

x = re.search("at|in", txt)
print(x.string) # Returns the string passed into the function
print(x.span()) # Returns a tuple containing the start, and end positions of the match
print(x.group()) # Returns the part of the string where there was a match

print(txt[x.span()[0]:x.span()[1]]) # Print the content of the string at the positions of the match
```

The output will look like:

```
I started studying this year at Durham University (15, 17) in in
```

As you can see, there was a match to our regex at the character with index 15 (counting starts from 0), ending at the character with index 17. Indeed, the string "in" was found within the word "studying".

However, if you read the input text, there should have been a second match for the word "at" but only the first match was returned. Note that if there is more than one match, only the first occurrence of the match will be returned by the search() function! We can use the findall() function to get a list of all matches in the order they are found.

```
x = re.findall("at|in", txt)
for match in x:
    print(match)
```

The output will look like:

```
in at
```

As expected, the findall() function returned two matches, "in" and "at".

Consider the string "stp stop stoop stoop

1.3.2 Matching to validate strings

```
text = list()
text.append("0123456789")
text.append("12345")
text.append("0000a00005")
text.append("+000001111")
text.append("00000011115")
text.append("2030405060")
regex = "[0-9]{10}"
result = list()
for t in text:
   x = re.match(regex, t)
   if(x != None):
       print(t,"->",x.group())
   else:
       print(t,"-> No match")
   result.append(x)
for i in range(len(text)):
   if(result[i]!=None and result[i].group()==text[i]):
       print(text[i],"is a valid identification number")
   else:
       print(text[i],"is NOT a valid identification number")
```

```
0123456789 -> 0123456789

12345 -> No match

0000a00005 -> No match

+000001111 -> No match

00000011115 -> 0000001111

2030405060 -> 2030405060

0123456789 is a valid identification number

12345 is NOT a valid identification number

0000a000005 is NOT a valid identification number
```

```
+000001111 is NOT a valid identification number 00000011115 is NOT a valid identification number 2030405060 is a valid identification number
```

Notice that string "00000011115" consists of 11 numerical digits, thus the regular expression matches the subset "0000001111". However, this is not a valid identification number according to the specification above. When validating input, remember to check whether the matched string is equal to the query string.

1.3.3 Credit card number validation

Let's try to validate whether the following strings are valid VISA or Mastercard credit card numbers: "10000000000000", "400000000000", "50000000000000", "50000000000000", "500000000000000", "40123456789". VISA credit card numbers should start with a 4 and have 13 or 16 digits. Mastercard credit card numbers start with a 5 and have 16 digits.

```
text = list()
text.append("100000000000") # 13 digits - Not valid
text.append("4000000000000") # 13 digits - Valid VISA
text.append("5000000000000") # 13 digits - Not valid
text.append("5000000000000000") # 16 digits - Valid Mastercard
text.append("50000a0000000c000") # Not valid, contains letters
text.append("40123456789") # 11 digits - Not valid
regex = "(5[0-9]{15})|(4([0-9]{12}|[0-9]{15}))" # Number 5 followed by 15 digits OR number 4 followed
    by either 12 or 15 digits
result = list()
for t in text:
   x = re.match(regex, t)
   if(x != None):
       print(t,"->",x.group())
       print(t,"-> No match")
   result.append(x)
print("")
for i in range(len(text)):
   if(result[i]!=None and result[i].group()==text[i]):
       print(text[i],"is a valid VISA or Mastercard number")
       print(text[i], "is NOT a valid VISA or Mastercard number")
```

The output will look like:

```
100000000000 -> No match
400000000000 -> 400000000000
5000000000000 -> No match
500000000000000 -> 50000000000000
50000a0000000000 -> No match
40123456789 -> No match

1000000000000 is NOT a valid VISA or Mastercard number
400000000000 is a valid VISA or Mastercard number
5000000000000 is NOT a valid VISA or Mastercard number
5000000000000 is NOT a valid VISA or Mastercard number
500000000000000 is NOT a valid VISA or Mastercard number
50000000000000000 is NOT a valid VISA or Mastercard number
40123456789 is NOT a valid VISA or Mastercard number
```

Let's analyse the regex that we used. Both VISA and Mastercard numbers start with a specific digit but Mastercard has exactly 16 digits in total, while VISA can have either 13 or 16 digits. Let's first create a regex for each case separately. For Mastercard, it should be "5[0-9]{15}", the digit 5 followed by exactly 15 digits (0-9), for a total of 16 digits. For VISA, it should be "4([0-9]{12}|[0-9]{15})", the digit 4 followed by either exactly 12 digits, for a total of 13 digits, or exactly 15 digits, for a total of 16 digits. Then, to include both the VISA and the Mastercard cases in our final regex, we can enclose each regex within parentheses and combine them with the OR ("|") operator.

Note that the expression [0-9] could be switched to [\d]:

```
regex = "(5[\d]{15})|(4([\d]{12}|[\d]{15}))" # Number 5 followed by 15 digits OR number 4 followed by
    either 12 or 15 digits

result = list()
for t in text:
    x = re.match(regex, t)
    if(x != None):
        print(t,"->",x.group())
    else:
        print(t,"-> No match")
    result.append(x)

print("")
for i in range(len(text)):
    if(result[i]!=None and result[i].group()==text[i]):
        print(text[i],"is a valid VISA or Mastercard number")
    else:
        print(text[i],"is NOT a valid VISA or Mastercard number")
```

The output will look like:

```
10000000000 -> No match
400000000000 -> 400000000000
500000000000 -> No match
500000000000000 -> 500000000000000
50000a0000000000 -> No match
40123456789 -> No match
1000000000000 is NOT a valid VISA or Mastercard number
400000000000 is a valid VISA or Mastercard number
5000000000000 is NOT a valid VISA or Mastercard number
5000000000000 is NOT a valid VISA or Mastercard number
500000000000000 is NOT a valid VISA or Mastercard number
5000000000000000 is NOT a valid VISA or Mastercard number
40123456789 is NOT a valid VISA or Mastercard number
```

1.3.4 Validation of United Kingdom's National Insurance numbers (NINO)

According to the rules for validating UK national insurance numbers¹, a NINO is made up of 2 letters, 6 numbers and a final letter, which is always A, B, C, or D. It looks something like this: QQ 12 34 56 A. The characters D, F, I, Q, U, and V are not used as either the first or second letter of a NINO prefix. The letter O is not used as the second letter of a prefix.

Let's create the required regex step-by-step and validate the following strings: "AA 123456 B", "AO 123456 B", "AO 123456 B", "AA 123456 B". We must also take white spaces into consideration. Let's consider the following two ways of writing a NINO: AA123456A and AA 12345 A.

- 1. The first letter should be any of A, B, C, E, G, H, J, K, L, M, N, O, P, R, S, T, W, X, Y: (A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y)
- 2. The second letter should be any of A, B, C, E, G, H, J, K, L, M, N, P, R, S, T, W, X, Y: (A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y) (A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y)
- 3. The third letter can optionally be a white space character: $(A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y) (A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y) [\slashed]$
- 4. Then, exactly 6 digits (0-9) are required: $(A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y) (A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y)[\backslash s]? \cite{Constraint} \cite$
- 5. The next letter can optionally be a white space character: $(A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y) (A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y) [\S]? [0-9] \{6\} [\S]?$

¹ https://www.gov.uk/hmrc-internal-manuals/national-insurance-manual/nim39110

6. The final letter must be one of A, B, C, or D: $(A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y)(A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y)[\s]?[0-9]\{6\}[\s]?(\mathbf{A}|\mathbf{B}|\mathbf{C}|\mathbf{D})$

```
text = list()
text.append("AA 123456 B") # Valid
text.append("AO 123456 B") # Not valid
text.append("AQ123456 B") # Not valid
text.append("QQ 123456 B") # Not valid
text.append("AA 123456 X") # Not valid
text.append("AA 12 34 56 B") # Not valid
text.append("AA123456B") # Valid
text.append("AA12345B") # Not valid
text.append("A 123456 B") # Not valid
text.append("A 123456 B") # Not valid
regex =
    "(A|B|C|E|G|H|J|K|L|M|N|O|P|R|S|T|W|X|Y)(A|B|C|E|G|H|J|K|L|M|N|P|R|S|T|W|X|Y)[\scalebox{0.5}{$\setminus$} ?[0-9]\{6\}[\scalebox{0.5}{$\setminus$} ?(A|B|C|D)"
result = list()
for t in text:
   x = re.match(regex, t)
   if(x != None):
       print(t,"->",x.group())
    else:
       print(t,"-> No match")
   result.append(x)
print("")
for i in range(len(text)):
    if(result[i]!=None and result[i].group()==text[i]):
       print("VALID\t",text[i])
    else:
       print("----\t",text[i])
The output will look like:
AA 123456 B -> AA 123456 B
AO 123456 B -> No match
AQ123456 B -> No match
QQ 123456 B -> No match
AA 123456 X -> No match
AA 12 34 56 B -> No match
AA123456B -> AA123456B
AA12345B -> No match
A 123456 B -> No match
A 123456 B -> No match
VALID AA 123456 B
---- AO 123456 B
---- AQ123456 B
---- QQ 123456 B
---- AA 123456 X
---- AA 12 34 56 B
VALID AA123456B
```

1.3.5 Validation of hexadecimal numbers

---- AA12345B ---- A 123456 B ---- A 123456 B

Let's use regular expressions to check whether a string corresponds to a hexadecimal number. Consider the strings "xAF1400BD", "1299ab32", "xFF00FF5R", "0xaa00bb". How can we check if these strings are representations of hexadecimal numbers? Remember that valid hexadecimal digits are [0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F,a,b,c,d,e,f] and that in computers, hexadecimal numbers may be denoted with an "x" or "0x" (either lowercase

or uppercase) in the beginning. For example, the hexadecimal number FFF, can also be written as fff, 0xfff, 0XFFF, xfff, XFFF and can also have mixed lowercase and uppercase letters.

```
text = list()
text.append("xAF1400BD") # Valid
text.append("1299ab32") # Valid
text.append("xFF00FF5R") # Not valid - Character R is not a valid headecimal digit
text.append("0xaa00bb") # Valid
text.append("xaa00bb4657AB000922334bce111A") # Valid
text.append("0xfff") # Valid
text.append("xFfF") # Valid
text.append("AAA") # Valid
text.append("ALA") # Not valid - Character L is not a valid headecimal digit
 \textbf{regex} = "(0x|0X|x|X)?[0-9a-fA-F] + " \# \texttt{One optional occurence of 0x, 0X, x, or X, followed by at least} 
    one digit from 0 to 9, or lowercase letter from a to f, or uppercase letter from A to F
result = list()
for t in text:
   x = re.match(regex, t)
   if(x != None):
       print(t,"->",x.group())
       print(t,"-> No match")
   result.append(x)
print("")
for i in range(len(text)):
   if(result[i]!=None and result[i].group()==text[i]):
       print("VALID HEX\t",text[i])
   else:
       print("-----\t",text[i])
The output will look like:
```

```
xAF1400BD -> xAF1400BD
1299ab32 -> 1299ab32
xFF00FF5R -> xFF00FF5
0xaa00bb -> 0xaa00bb
xaa00bb4657AB000922334bce111A -> xaa00bb4657AB000922334bce111A
Oxfff -> Oxfff
xFfF -> xFfF
AAA -> AAA
ALA -> A
VALID HEX xAF1400BD
VALID HEX 1299ab32
           xFF00FF5R
VALTD HEX
          0xaa00bb
VALID HEX
           xaa00bb4657AB000922334bce111A
VALID HEX
           Oxfff
VALID HEX
           xFfF
VALID HEX
           AAA
```

Note that in the case of the "ALA" string, the regex found the match "A", which is a valid hexadecimal number, but the full string "ALA" is not a valid hexadecimal number. When validating input, remember to check that the matched string is equal to the query string.

1.4 Using Regular Expressions to search elements in files

1.4.1 Parsing XML files

Let's use regular expressions to parse an XML file. We will use the movies.xml file which contains the titles and release dates of 5 movies. We will use a regular expression to retrieve all the movie titles from the file. First,

copy the file "movies.xml" to your current working directory or retrieve the absolute path of the file. Then open the file, load its contents into a variable and close the file.

```
f = open("movies.xml", "r") # Opens the file for reading only ("r")
text = f.read() # Store the contents of the file in variable "text". read() returns all the contents
    of the file
f.close() # Close the file
print(text) # Print the contents of variable "text"
```

The output will look like:

```
<movies>
 <movie id="1">
   <title>And Now for Something Completely Different</title>
   <released>1971</released>
 </movie>
 <movie id="2">
   <title>Monty Python and the Holy Grail</title>
   <released>1974</released>
 </movie>
 <movie id="3">
   <title>Monty Python's Life of Brian</title>
   <released>1979</released>
 </movie>
 <movie id="4">
   <title>Monty Python Live at the Hollywood Bowl</title>
   <released>1982</released>
 </movie>
 <movie id="5">
   <title>Monty Python's The Meaning of Life</title>
   <released>1983</released>
 </movie>
</movies>
```

As you can see above, all movie titles in the movies.xml XML file are enclosed within the <title> and </title> tags. Let's use a regular expression to match all the strings that are enclosed within these tags.

```
['<title>And Now for Something Completely Different</title>', '<title>Monty Python and the Holy Grail</title>', "<title>Monty Python's Life of Brian</title>", '<title>Monty Python Live at the Hollywood Bowl</title>', "<title>Monty Python's The Meaning of Life</title>"]

Movie titles from movies.xml:
And Now for Something Completely Different
Monty Python and the Holy Grail
Monty Python's Life of Brian
Monty Python Live at the Hollywood Bowl
Monty Python's The Meaning of Life
```

1.4.2 Parsing HTML files

Let's read the links.html HTML file and use regular expressions to find all links within the file. Remember that in HTML, links are denoted using the $\langle a \rangle$ and $\langle a \rangle$ tags and the link url is provided using the "href" attribute within the $\langle a \rangle$ tag. For example a link for the main website of Durham University would be:

```
<a href="https://www.dur.ac.uk/">Durham University</a>
```

```
f = open("links.html", "r") # Opens the file for reading only ("r")
text = f.read() # Store the contents of the file in variable "text". read() returns all the contents
    of the file
f.close() # Close the file
print(text,"\n") # Print the contents of variable "text"

regex = '<a href=".*"'
x = re.findall(regex, text) # Find all matches of the regex

print(x,"\n")

print("Links from links.xml:")
for link in x:
    link = link.replace('<a href="',"") # Remove <a href=" by replacing it with empty string link = link.replace('"',"") # Remove " by replacing it with empty string print(link)</pre>
```

The output will look like:

```
<!doctype html>
<html lang="en-GB">
<head>
  <meta name="viewport" content="width=device-width, initial-scale=1">
  <title>Test page for Text Mining and Language Analytics</title>
</head>
<body>
  <h1>Test page for Text Mining and Language Analytics</h1>
  <a href="https://www.durham.ac.uk/departments/academic/computer-science/">Link 1</a><br/>br/>
  <a href="https://www.gov.uk/government/organisations/hm-revenue-customs">Link 2</a><br/>br/>
  <a href="https://www.gov.uk/">Link 3</a><br/>
  <a href="https://www.dur.ac.uk/">Link 4</a><br/>
   <a href="https://www.nhs.uk/">Link 5</a><br/>
</body>
</html>
['<a href="https://www.durham.ac.uk/departments/academic/computer-science/"', '<a
    href="https://www.gov.uk/government/organisations/hm-revenue-customs", '<a
    href="https://www.gov.uk/"', '<a href="https://www.dur.ac.uk/"', '<a href="https://www.nhs.uk/"']
Links from links.xml:
https://www.durham.ac.uk/departments/academic/computer-science/
https://www.gov.uk/government/organisations/hm-revenue-customs
https://www.gov.uk/
https://www.dur.ac.uk/
https://www.nhs.uk/
```

Note than we used single quotes to denote strings that included a double quote character.

As you can see above, we successfully retrieved the urls from links.html. Nevertheless, please note that using regular expressions is not the best approach for parsing HTML files due to the flexibility of HTML syntax. Solutions like XPath (https://developer.mozilla.org/en-US/docs/Web/API/XPathExpression) are more suitable for HTML parsing.

1.4.3 Parsing raw text

The file emails.txt contains a list of emails from various domains. Let's use regular expressions to find the emails from the durham.ac.uk domain.

```
f = open("emails.txt", "r") # Opens the file for reading only ("r")
text = f.read() # Store the contents of the file in variable "text". read() returns all the contents
    of the file
f.close() # Close the file
print(text,"\n") # Print the contents of variable "text"

regex = "[0-9a-zA-Z!#$%&'*+-/=?^_^{[]}~.]+@durham.ac.uk"

x = re.findall(regex, text) # Find all matches of the regex

print(x,"\n")

print("Emails from durham.ac.uk:")
for email in x:
    print(email)
```

The output will look like:

```
john+acme.co@hotmail.com
bob@gmail.com
tom@durham.ac.uk
jerry@durham.ac.uk
scrooge@durham.ac.uk
donald@yahoo.co.uk
huey@yahoo.co.uk
dewey@gmail.com
louie.duck@durham.ac.uk
gyro.gearloose@yahoo.co.uk
bart@yahoo.co.uk
homer@gmail.com
stan@hotmail.com
kyle-broflovski@durham.ac.uk
eric@yahoo.co.uk
kenny@gmail.com
butters@durham.ac.uk
wendy@hotmail.com
randy_marsh@durham.ac.uk
chef@gmail.com
['tom@durham.ac.uk', 'jerry@durham.ac.uk', 'scrooge@durham.ac.uk', 'louie.duck@durham.ac.uk',
    'kyle-broflovski@durham.ac.uk', 'butters@durham.ac.uk', 'randy_marsh@durham.ac.uk']
Emails from durham.ac.uk:
tom@durham.ac.uk
jerry@durham.ac.uk
scrooge@durham.ac.uk
louie.duck@durham.ac.uk
kyle-broflovski@durham.ac.uk
butters@durham.ac.uk
randy_marsh@durham.ac.uk
```

As you can see above, we retrieved all emails from the durham.ac.uk email.

1.5 Using Regular Expressions to substitute strings

1.5.1 Substitution example

Let's now convert any string that matches the "at|in" regex in the text "I started studying this year at Durham University" with the string "FOO". To achieve this, we are going to use the sub() function.

```
txt = "I started studying this year at Durham University"
x = re.sub("at|in","F00", txt)
print(x)
```

I started studyFOOg this year FOO Durham University

As expected, two matches of the regex were converted to "FOO". What if we wanted only the first match to be substituted with "FOO"? We can add an additional argument in the sub() function indicating the number of substitutions we would like to make.

```
x = re.sub("at|in", "F00", txt,1)
print(x)
```

The output will look like:

I started studyF00g this year at Durham University

1.5.2 Email domain substitution

Let's load again emails.txt and change the domain to "new.ac.uk" for all emails in the "durham.ac.uk", "gmail.com", and "yahoo.co.uk" domains.

```
f = open("emails.txt", "r") # Opens the file for reading only ("r")
text = f.read() # Store the contents of the file in variable "text". read() returns all the contents
    of the file
f.close() # Close the file
print("OLD EMAILS:")
print(text,"\n") # Print the contents of variable "text"

regex = "@durham.ac.uk|@gmail.com|@yahoo.co.uk"

print("NEW EMAILS:")
x = re.sub(regex,"@new.ac.uk", text)

print(x)
```

```
OLD EMAILS:
john+acme.co@hotmail.com
bob@gmail.com
tom@durham.ac.uk
jerry@durham.ac.uk
scrooge@durham.ac.uk
donald@yahoo.co.uk
huey@yahoo.co.uk
dewey@gmail.com
louie.duck@durham.ac.uk
gyro.gearloose@yahoo.co.uk
bart@yahoo.co.uk
homer@gmail.com
stan@hotmail.com
kyle-broflovski@durham.ac.uk
eric@yahoo.co.uk
kenny@gmail.com
butters@durham.ac.uk
wendy@hotmail.com
randy_marsh@durham.ac.uk
chef@gmail.com
NEW EMAILS:
john+acme.co@hotmail.com
```

bob@new.ac.uk tom@new.ac.uk jerry@new.ac.uk scrooge@new.ac.uk donald@new.ac.uk huey@new.ac.uk dewey@new.ac.uk louie.duck@new.ac.uk gyro.gearloose@new.ac.uk bart@new.ac.uk homer@new.ac.uk stan@hotmail.com kyle-broflovski@new.ac.uk eric@new.ac.uk kenny@new.ac.uk butters@new.ac.uk wendy@hotmail.com randy_marsh@new.ac.uk chef@new.ac.uk

1.6 Exercises

- Exercise 1.1 Validate the following identification numbers: "500011110000" (valid), "5000 1111 0000" (valid), "500001110000" (valid), "500001110000" (valid), "500001110000" (not valid), "500001110000" (not valid), "400001110000" (not valid), "5000011110000" (not valid). A valid number should be 12 digits long, the first digit should always be 5, the 5th digit should be 0 or 1, and the last digit cannot be 8 or 9. The numbers can also be grouped in groups of four digits with a white space between groups.
- Exercise 1.2 Redo the activity from Section 1.3.4 in order to also support the following writing of NINO strings: "AA 12 34 56 A"
- Exercise 1.3 Use regular expressions to print a list of all the album titles, a list of all the artists, and the average price of all CDs in the cds.xml file.
- Exercise 1.4 Create a regular expression that matches all the strings in the first column but none of those in the second column

affgfking	fgok
rafgkahe	a fgk
bafghk	affgm
baffgkit	afffhk
affgfking	fgok
rafgkahe	afg.K
bafghk	$\operatorname{aff}\operatorname{gm}$
baffgkit	afffhgk

Exercise 1.5 Use a regular expression to substitute the quantities of the bought items with "XX" in the following sentence: "Yesterday we bought 120 packs of A4 paper, 5 bottles of ink, 10 boxes of paperclips, 200 notebooks, and 5.35 litres of fuel. The total cost for order #1290 was £1000.43.". Note that the order number and cost must not be substituted. The resulting sentence must be: "Yesterday we bought XX packs of A4 paper, XX bottles of ink, XX boxes of paperclips, XX notebooks, and XX litres of fuel. The total cost for order #1290 was £1000.43."

Workshop 2: Text pre-processing

2.1 NLTK corpora

The NLTK python package comes pre-loaded with a lot of corpora that can be used to experiment with text pre-processing.

2.1.1 Import NLTK

Let's load one of the available corpora in NLTK, called movie_reviews. The movie_reviews corpus contains documents consisting of reviews about movies and annotated in relation to their sentiment.

```
import nltk # Import the NLTK library
nltk.download('movie_reviews') # Download movie_reviews from NLTK
from nltk.corpus import movie_reviews # Import the movie_reviews corpus from NLTK
```

2.1.2 Corpus file IDs

Let's check the file IDs in the movie_reviews corpus:

```
movie_reviews.fileids() # List file-ids in the corpus
```

The output will look like:

```
['neg/cv000_29416.txt', 'neg/cv001_19502.txt', 'neg/cv002_17424.txt', 'neg/cv003_12683.txt', 'neg/cv004_12641.txt', ...]
```

2.1.3 Corpus file categories

Let's examine what classes (categories) are included in the corpus:

```
movie_reviews.categories() # List categories in the corpus
```

The output will look like:

```
['neg', 'pos']
```

This means that the movie_reviews corpus contains documents that have been annotated as belonging to two distinct classes, more specifically the Negative class (neg) and the Positive class (pos), in relation to the sentiment of the respective movie review.

2.1.4 Corpus words

Let's see the list of words in the movie_review corpus and print their number:

```
print(movie_reviews.words())
length = len(movie_reviews.words())
print("Number of words in corpus: ", length)
```

```
['plot', ':', 'two', 'teen', 'couples', 'go', 'to', ...]

Number of words in corpus: 1583820
```

2.1.5 Selecting files from specific category

Use the following in order to view the file IDs for files belonging to a specific category. Note that you can provide a list with categories, e.g ['neg','pos','other']

```
movie_reviews.fileids(['neg']) # List file ids with 'neg' category
movie_reviews.fileids(['pos']) # List file ids with 'pos' category
```

The output will look like:

```
['neg/cv000_29416.txt', 'neg/cv001_19502.txt', 'neg/cv002_17424.txt', 'neg/cv003_12683.txt', 'neg/cv004_12641.txt', 'neg/cv005_29357.txt', 'neg/cv006_17022.txt', 'neg/cv007_4992.txt', ...]
['pos/cv000_29590.txt', 'pos/cv001_18431.txt', 'pos/cv002_15918.txt', 'pos/cv003_11664.txt', 'pos/cv004_11636.txt', 'pos/cv005_29443.txt', 'pos/cv006_15448.txt', 'pos/cv007_4968.txt', ...]
```

2.1.6 Corpus sentences

To access the tokenised sentences included in the corpus:

```
nltk.download('punkt') # Download the Punkt sentence tokenizer from NLTK

movie_reviews.sents()

The output will look like:
```

```
[['plot', ':', 'two', 'teen', 'couples', 'go', 'to', 'a', 'church', 'party', ',', 'drink', 'and', 'then', 'drive', '.'], ['they', 'get', 'into', 'an', 'accident', '.'], ...]
```

2.1.7 Accessing sentences from specific file

We can use the "fileids" argument in order to access the sentences from a specific file in the corpus:

```
movie_reviews.sents(fileids='pos/cv004_11636.txt')
```

The output will look like:

```
[['moviemaking', 'is', 'a', 'lot', 'like', 'being', 'the', 'general', 'manager', 'of', 'an', 'nfl', 'team', 'in', 'the', 'post', '-', 'salary', 'cap', 'era', '--', 'you', "'", 've', 'got', 'to', 'know', 'how', 'to', 'allocate', 'your', 'resources', '.'], ['every', 'dollar', 'spent', 'on', 'a', 'free', '-', 'agent', 'defensive', 'tackle', 'is', 'one', 'less', 'dollar', 'than', 'you', 'can', 'spend', 'on', 'linebackers', 'or', 'safeties', 'or', 'centers', '.'], ...]
```

2.1.8 Number of documents in each category

Let's see how many documents (files) does the movie_reviews corpus contain, and then how many documents of "neg" category and how many of "pos" category it contains:

```
documents = len(movie_reviews.fileids())
documents_neg = len(movie_reviews.fileids(['neg']))
documents_pos = len(movie_reviews.fileids(['pos']))
print("Number of documents: ", documents)
print("Number of documents in neg category: ", documents_neg)
print("Number of documents in pos category: ", documents_pos)
```

```
Number of documents in neg category: 1000 Number of documents in pos category: 1000
```

As you can see here, the corpus is perfectly balanced between negative and positive sentiment reviews, having 1000 documents in each category.

2.1.9 Corpus raw text

To access the raw text of a specific file, use the ".raw()" function:

```
rawtext = movie_reviews.raw('neg/cv002_17424.txt').strip()[:500] # strip() removes blank spaces in
    the beginning and end of a string. [:500] is used to only retrieve the first 500 characters of
    the string
print(rawtext)
```

The output will look like:

```
it is movies like these that make a jaded movie viewer thankful for the invention of the timex indiglo watch .

based on the late 1960's television show by the same name , the mod squad tells the tale of three reformed criminals under the employ of the police to go undercover .

however , things go wrong as evidence gets stolen and they are immediately under suspicion .

of course , the ads make it seem like so much more .

quick cuts , cool music , claire dane's nice hair and cute outfits , car
```

2.2 Input text from text file

NLTK corpora are very useful for testing NLP and text processing algorithms. But what if we wanted to load text from a text file? Let's try to load the text from the "alice.txt" file. First, copy the file "alice.txt" to your current working directory or retrieve the absolute path of the file. Then open the file, load its contents into a variable and close the file.

```
f = open("alice.txt", "r") # Opens the file for reading only ("r")
text = f.read() # Store the contents of the file in variable "text". read() returns all the contents
    of the file
f.close() # Close the file
print(text) # Print the contents of variable "text"
```

- Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book," thought Alice "without pictures or conversations?"
- So she was considering in her own mind (as well as she could, for the hot day made her feel very sleepy and stupid), whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a White Rabbit with pink eyes ran close by her.
- There was nothing so very remarkable in that; nor did Alice think it so very much out of the way to hear the Rabbit say to itself, "Oh dear! Oh dear! I shall be late!" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually took a watch out of its waistcoat-pocket, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and burning with curiosity, she ran across the field after it, and fortunately was just in time to see it pop down a large rabbit-hole under the hedge.

2.3 Text pre-processing

In the text output above, you can see that the text contains various sentences, a mix of lowercase and uppercase letters, there are some parentheses, as well as some punctuation marks.

2.3.1 Sentence Tokenisation

Let's tokenise the text, i.e. chop the text into pieces. In this case, the text is in the English language and punctuation marks are used to separate sentences from each other and a blank space is used to separate words from each other. For splitting a string into sentences, NLTK has the sent_tokenize() default tokeniser function.

Let's divide the text from alice.txt into sentences:

```
from nltk import sent_tokenize # Import the sent_tokenize function from NLTK sent_tokenize(text) # Tokenise "text" into sentences and print the output
```

The output will look like:

['Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book," thought Alice "without pictures or conversations?"', 'So she was considering in her own mind (as well as she could, for the hot day made her feel very sleepy and stupid), whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a White Rabbit with pink eyes ran close by her.', 'There was nothing so very remarkable in that; nor did Alice think it so very much out of the way to hear the Rabbit say to itself, "Oh dear!', 'Oh dear!', 'I shall be late!"', '(when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually took a watch out of its waistcoat-pocket, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and burning with curiosity, she ran across the field after it, and fortunately was just in time to see it pop down a large rabbit-hole under the hedge.']

2.3.2 Word Tokenisation

Now let's divide the text from alice.txt into words using the word_tokenize() NLTK function:

```
from nltk import word_tokenize # Import the word_tokenize function from NLTK
word_tokenize(text) # Tokenise "text" into words and print the output
```

```
['Alice', 'was', 'beginning', 'to', 'get', 'very', 'tired', 'of', 'sitting', 'by', 'her', 'sister',
    'on', 'the', 'bank', ',', 'and', 'of', 'having', 'nothing', 'to', 'do', ':', 'once', 'or',
    'twice', 'she', 'had', 'peeped', 'into', 'the', 'book', 'her', 'sister', 'was', 'reading', ',',
    'but', 'it', 'had', 'no', 'pictures', 'or', 'conversations', 'in', 'it', ',', '
    'what', 'is', 'the', 'use', 'of', 'a', 'book', ',', "''", 'thought', 'Alice', '``', 'without',
    'pictures', 'or', 'conversations', '?', "''", 'So', 'she', 'was', 'considering', 'in', 'her',
    own', 'mind', '(', 'as', 'well', 'as', 'she', 'could', ',', 'for', 'the', 'hot', 'day', 'made',
    'her', 'feel', 'very', 'sleepy', 'and', 'stupid', ')', ',', 'whether', 'the', 'pleasure', 'of',
    'making', 'a', 'daisy-chain', 'would', 'be', 'worth', 'the', 'trouble', 'of', 'getting', 'up',
    'and', 'picking', 'the', 'daisies', ',', 'when', 'suddenly', 'a', 'White', 'Rabbit', 'with',
    'pink', 'eyes', 'ran', 'close', 'by', 'her', '.', 'There', 'was', 'nothing', 'so', 'very',
    'remarkable', 'in', 'that', ';', 'nor', 'did', 'Alice', 'think', 'it', 'so', 'very', 'much',
    'out', 'of', 'the', 'way', 'to', 'hear', 'the', 'Rabbit', 'say', 'to', 'itself', ',', '``', 'Oh',
    'dear', '!', 'Oh', 'dear', '!', 'I', 'shall', 'be', 'late', '!', "''", '(', 'when', 'she',
    'thought', 'it', 'over', 'afterwards', ',', 'it', 'occurred', 'to', 'her', 'that', 'she',
    'ought', 'to', 'have', 'wondered', 'at', 'this', ',', 'but', 'at', 'the', 'time', 'it', 'all',
    'seemed', 'quite', 'natural', ')', ';', 'but', 'when', 'the', 'Rabbit', 'actually', 'took', 'a',
    'watch', 'out', 'of', 'its', 'waistcoat-pocket', ',', 'and', 'looked', 'at', 'it', ',', 'and',
    'then', 'hurried', 'on', ',', 'Alice', 'started', 'to', 'her', 'feet', ',', 'for', 'it',
    'flashed', 'across', 'her', 'mind', 'that', 'she', 'had', 'never', 'before', 'seen', 'a',
    'rabbit', 'with', 'either', 'a', 'waistcoat-pocket', ',', 'or', 'a', 'watch', 'to', 'take',
```

```
'out', 'of', 'it', ',', 'and', 'burning', 'with', 'curiosity', ',', 'she', 'ran', 'across',
'the', 'field', 'after', 'it', ',', 'and', 'fortunately', 'was', 'just', 'in', 'time', 'to',
'see', 'it', 'pop', 'down', 'a', 'large', 'rabbit-hole', 'under', 'the', 'hedge', '.']
```

What if we wanted to tokenise each individual sentence of the text one by one?

```
for sentence in sent_tokenize(text): # Iterate the sentences in "text"
    print(word_tokenize(sentence)) # Print the word tokenisation results. Don't forget the four
    spaces indentation that Python expects!
```

The output will look like:

```
['Alice', 'was', 'beginning', 'to', 'get', 'very', 'tired', 'of', 'sitting', 'by', 'her', 'sister',
         'on', 'the', 'bank', ',', 'and', 'of', 'having', 'nothing', 'to', 'do', ':', 'once', 'or',
         'twice', 'she', 'had', 'peeped', 'into', 'the', 'book', 'her', 'sister', 'was', 'reading', ',',
         'but', 'it', 'had', 'no', 'pictures', 'or', 'conversations', 'in', 'it', ',', '``', 'and',
         'what', 'is', 'the', 'use', 'of', 'a', 'book', ',', "''", 'thought', 'Alice', '``', 'without',
         'pictures', 'or', 'conversations', '?', "''"]
['So', 'she', 'was', 'considering', 'in', 'her', 'own', 'mind', '(', 'as', 'well', 'as', 'she',
         "could", ",", "for", "the", "hot", "day", "made", "her", "feel", "very", "sleepy", "and", "like the state of the state o
         'stupid', ')', ',', 'whether', 'the', 'pleasure', 'of', 'making', 'a', 'daisy-chain', 'would',
         'be', 'worth', 'the', 'trouble', 'of', 'getting', 'up', 'and', 'picking', 'the', 'daisies', ',',
         'when', 'suddenly', 'a', 'White', 'Rabbit', 'with', 'pink', 'eyes', 'ran', 'close', 'by', 'her',
['There', 'was', 'nothing', 'so', 'very', 'remarkable', 'in', 'that', ';', 'nor', 'did', 'Alice', 'think', 'it', 'so', 'very', 'much', 'out', 'of', 'the', 'way', 'to', 'hear', 'the', 'Rabbit', 'say', 'to', 'itself', ',', '``', 'Oh', 'dear', '!']
['Oh', 'dear', '!']
['I', 'shall', 'be', 'late', '!', "''"]
['(', 'when', 'she', 'thought', 'it', 'over', 'afterwards', ',', 'it', 'occurred', 'to', 'her',
         'that', 'she', 'ought', 'to', 'have', 'wondered', 'at', 'this', ',', 'but', 'at', 'the', 'time',
         'it', 'all', 'seemed', 'quite', 'natural', ')', ';', 'but', 'when', 'the', 'Rabbit', 'actually',
         'took', 'a', 'watch', 'out', 'of', 'its', 'waistcoat-pocket', ',', 'and', 'looked', 'at', 'it',
         ',', 'and', 'then', 'hurried', 'on', ',', 'Alice', 'started', 'to', 'her', 'feet', ',', 'for',
         'it', 'flashed', 'across', 'her', 'mind', 'that', 'she', 'had', 'never', 'before', 'seen', 'a',
         'rabbit', 'with', 'either', 'a', 'waistcoat-pocket', ',', 'or', 'a', 'watch', 'to', 'take',
         'out', 'of', 'it', ',', 'and', 'burning', 'with', 'curiosity', ',', 'she', 'ran', 'across', 'the', 'field', 'after', 'it', ',', 'and', 'fortunately', 'was', 'just', 'in', 'time', 'to',
         'see', 'it', 'pop', 'down', 'a', 'large', 'rabbit-hole', 'under', 'the', 'hedge', '.']
```

2.3.3 Lowercasing

As you can see in the list of tokens above, some words have uppercase letters. Let's convert all characters in our tokens to lowercase. We can use the lower() function to convert all characters in a string to lowercase. For example:

```
test_string = "uNIverSIty"
test_string_lowercase = test_string.lower()
print(test_string_lowercase)
```

The output will look like:

```
university
```

Let's now convert all words in the text from alice.txt to lowercase, iterating through each sentence and word, and saving the lowercase words in a list:

```
words_lowercase = []
for sentence in sent_tokenize(text): # Iterate the sentences in "text"
    for word in word_tokenize(sentence): # Iterate the words in "sentence"
        words_lowercase.append(word.lower()) # Convert word to lowercase and add it to the
        "words_lowercase" list
print(words_lowercase) # Print the list of lowercase words
```

```
['alice', 'was', 'beginning', 'to', 'get', 'very', 'tired', 'of', 'sitting', 'by', 'her', 'sister',
     'on', 'the', 'bank', ',', 'and', 'of', 'having', 'nothing', 'to', 'do', ':', 'once', 'or', 'twice', 'she', 'had', 'peeped', 'into', 'the', 'book', 'her', 'sister', 'was', 'reading', ',', 'but', 'it', 'had', 'no', 'pictures', 'or', 'conversations', 'in', 'it', ',', '``', 'and',
     'what', 'is', 'the', 'use', 'of', 'a', 'book', ',', "''", 'thought', 'alice', '``', 'without',
     'pictures', 'or', 'conversations', '?', "''", 'so', 'she', 'was', 'considering', 'in', 'her',
     'own', 'mind', '(', 'as', 'well', 'as', 'she', 'could', ',', 'for', 'the', 'hot', 'day', 'made', 'her', 'feel', 'very', 'sleepy', 'and', 'stupid', ')', ',', 'whether', 'the', 'pleasure', 'of',
     'making', 'a', 'daisy-chain', 'would', 'be', 'worth', 'the', 'trouble', 'of', 'getting', 'up',
     'and', 'picking', 'the', 'daisies', ',', 'when', 'suddenly', 'a', 'white', 'rabbit', 'with',
     'pink', 'eyes', 'ran', 'close', 'by', 'her', '.', 'there', 'was', 'nothing', 'so', 'very',
     'remarkable', 'in', 'that', ';', 'nor', 'did', 'alice', 'think', 'it', 'so', 'very', 'much',
     'out', 'of', 'the', 'way', 'to', 'hear', 'the', 'rabbit', 'say', 'to', 'itself', ',', '``', 'oh',
     'dear', '!', 'oh', 'dear', '!', 'i', 'shall', 'be', 'late', '!', "''", '(', 'when', 'she',
     'thought', 'it', 'over', 'afterwards', ',', 'it', 'occurred', 'to', 'her', 'that', 'she',
     'ought', 'to', 'have', 'wondered', 'at', 'this', ',', 'but', 'at', 'the', 'time', 'it', 'all',
     'seemed', 'quite', 'natural', ')', ';', 'but', 'when', 'the', 'rabbit', 'actually', 'took', 'a',
     'watch', 'out', 'of', 'its', 'waistcoat-pocket', ',', 'and', 'looked', 'at', 'it', ',', 'and',
     'then', 'hurried', 'on', ',', 'alice', 'started', 'to', 'her', 'feet', ',', 'for', 'it',
     'flashed', 'across', 'her', 'mind', 'that', 'she', 'had', 'never', 'before', 'seen', 'a',
     'rabbit', 'with', 'either', 'a', 'waistcoat-pocket', ',', 'or', 'a', 'watch', 'to', 'take',
     'out', 'of', 'it', ',', 'and', 'burning', 'with', 'curiosity', ',', 'she', 'ran', 'across',
     'the', 'field', 'after', 'it', ',', 'and', 'fortunately', 'was', 'just', 'in', 'time', 'to', 'see', 'it', 'pop', 'down', 'a', 'large', 'rabbit-hole', 'under', 'the', 'hedge', '.']
```

As you can see, there are no uppercase characters left in the words list.

2.3.4 Stop words removal

Stop words are common words in any language that don't hold semantic meaning of their own. In many NLP applications, it is useful to remove stop words from a text. To achieve this, we use stop words lists that have been compiled for each language. We can access the list of English stop words from the NLTK library using the following:

```
nltk.download('stopwords') # Download the stopwords lists from NLTK

from nltk.corpus import stopwords # Import the stop words lists from NLTK

stopwords_english = stopwords.words('english') # Load the stop words list for English in variable
    "stopwords_english"

print(stopwords_english) # Print the "stopwords_english" list
```

The output will look like:

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll",
    "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she',
    "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their',
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these',
    'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
    'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
    'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
    'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in',
    'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'here', 'there', 'when',
    'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some',
    'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can',
    'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've',
    'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn',
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn't", 'mightn't",
    'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't"]
```

As you can see, the list of words from alice.txt contains a number of stop words. Let's use NLTK's English stop words list to remove them:

```
words_lowercase_nostopwords = [] # Create empty list for remaining words
words_removed = [] # Create empty list for removed words
for word in words_lowercase: # Iterate lowercase words' list
    if word not in stopwords_english: # Check if word is in the stop words list
        words_lowercase_nostopwords.append(word)
    else:
        words_removed.append(word)
print(words_lowercase_nostopwords) # Print list of remaining words
print(words_removed) # Print list of removed words
```

```
['alice', 'beginning', 'get', 'tired', 'sitting', 'sister', 'bank', ',', 'nothing', ':', 'twice',
     'peeped', 'book', 'sister', 'reading', ',', 'pictures', 'conversations', ',', '``', 'use',
     'book', ',', "''", 'thought', 'alice', '``', 'without', 'pictures', 'conversations', '?', "''", 'considering', 'mind', '(', 'well', 'could', ',', 'hot', 'day', 'made', 'feel', 'sleepy',
     'stupid', ')', ',', 'whether', 'pleasure', 'making', 'daisy-chain', 'would', 'worth', 'trouble',
     'getting', 'picking', 'daisies', ',', 'suddenly', 'white', 'rabbit', 'pink', 'eyes', 'ran',
    'close', '.', 'nothing', 'remarkable', ';', 'alice', 'think', 'much', 'way', 'hear', 'rabbit', 'say', ',', '``', 'oh', 'dear', '!', 'oh', 'dear', '!', 'shall', 'late', '!', "''", 'thought', 'afterwards', ',', 'occurred', 'ought', 'wondered', ',', 'time', 'seemed', 'quite',
     'natural', ')', ';', 'rabbit', 'actually', 'took', 'watch', 'waistcoat-pocket', ',', 'looked', ',', 'hurried', ',', 'alice', 'started', 'feet', ',', 'flashed', 'across', 'mind', 'never',
     'seen', 'rabbit', 'either', 'waistcoat-pocket', ',', 'watch', 'take', ',', 'burning',
     'curiosity', ',', 'ran', 'across', 'field', ',', 'fortunately', 'time', 'see', 'pop', 'large',
     'rabbit-hole', 'hedge', '.']
['was', 'to', 'very', 'of', 'by', 'her', 'on', 'the', 'and', 'of', 'having', 'to', 'do', 'once',
     'or', 'she', 'had', 'into', 'the', 'her', 'was', 'but', 'it', 'had', 'no', 'or', 'in', 'it',
     'and', 'what', 'is', 'the', 'of', 'a', 'or', 'so', 'she', 'was', 'in', 'her', 'own', 'as', 'as',
     'she', 'for', 'the', 'her', 'very', 'and', 'the', 'of', 'a', 'be', 'the', 'of', 'up', 'and',
     'the', 'when', 'a', 'with', 'by', 'her', 'there', 'was', 'so', 'very', 'in', 'that', 'nor',
     'did', 'it', 'so', 'very', 'out', 'of', 'the', 'to', 'the', 'to', 'itself', 'i', 'be', 'when',
     'she', 'it', 'over', 'it', 'to', 'her', 'that', 'she', 'to', 'have', 'at', 'this', 'but', 'at',
     'the', 'it', 'all', 'but', 'when', 'the', 'a', 'out', 'of', 'its', 'and', 'at', 'it', 'and',
     'then', 'on', 'to', 'her', 'for', 'it', 'her', 'that', 'she', 'had', 'before', 'a', 'with', 'a',
     'or', 'a', 'to', 'out', 'of', 'it', 'and', 'with', 'she', 'the', 'after', 'it', 'and', 'was',
     'just', 'in', 'to', 'it', 'down', 'a', 'under', 'the']
```

As you can see in the list of the remaining words and the list of the removed words, we removed multiple words from the alice.txt text, including "was", "to", "she", "it", "for", etc.

2.3.5 Punctuation removal

If you look at the list of the remaining words, you can see that some of its contents are not actual words, but they are punctuation marks instead. Let's remove these punctuation marks from the words list:

```
from string import punctuation # Import the punctuation marks string
print(punctuation)
print("Variable type: ", type(punctuation))
```

The output will look like:

```
!"#$%&'()*+,-./:;<=>?@[\]^_`{|}~
Variable type: <class 'str'>
```

As you can see, the variable "punctuation" contains all the potential punctuation marks. However, the variable is a string. To help us iterate and compare with the list of words, we will first convert the string "punctuation" to a list with one character (punctuation mark) per element.

```
punctuation_list = list(punctuation) # Convert punctuation to a list
print(punctuation_list)

words_lowercase_nostopwords_no_punctuation = []
for word in words_lowercase_nostopwords:
```

```
if word not in punctuation_list:
    words_lowercase_nostopwords_no_punctuation.append(word)
print(words_lowercase_nostopwords_no_punctuation)
```

```
['!', '"', '#', '$', '%', '&', """, '(', ')', '*', '+', ',', '-', '.', '/', ':', ';', '<', '=', '>', '?', '@', '[', '\\', ']', '-', '-', ':', '\{', '|', '\}', '\']

['alice', 'beginning', 'get', 'tired', 'sitting', 'sister', 'bank', 'nothing', 'twice', 'peeped', 'book', 'sister', 'reading', 'pictures', 'conversations', ''`', 'use', 'book', "''", 'thought', 'alice', '``', 'without', 'pictures', 'conversations', "''", 'considering', 'mind', 'well', 'could', 'hot', 'day', 'made', 'feel', 'sleepy', 'stupid', 'whether', 'pleasure', 'making', 'daisy-chain', 'would', 'worth', 'trouble', 'getting', 'picking', 'daisies', 'suddenly', 'white', 'rabbit', 'pink', 'eyes', 'ran', 'close', 'nothing', 'remarkable', 'alice', 'think', 'much', 'way', 'hear', 'rabbit', 'say', '``', 'oh', 'dear', 'oh', 'dear', 'shall', 'late', "''", 'thought', 'afterwards', 'occurred', 'ought', 'wondered', 'time', 'seemed', 'quite', 'natural', 'rabbit', 'actually', 'took', 'watch', 'waistcoat-pocket', 'looked', 'hurried', 'alice', 'started', 'feet', 'flashed', 'across', 'mind', 'never', 'seen', 'rabbit', 'either', 'waistcoat-pocket', 'watch', 'take', 'burning', 'curiosity', 'ran', 'across', 'field', 'fortunately', 'time', 'see', 'pop', 'large', 'rabbit-hole', 'hedge']
```

As you can see, the majority of the punctuation marks where removed. But what about the remaining "''" and "''"? The punctuation marks list that we used did not contain double quotes and as a consequence, "'"" and "''" were treated as valid words. We can address this issue by adding the respective double quote characters in the list of punctuation marks.

Note that there are alternative ways for removing punctuation marks or special characters from text, like for example the use of regular expressions.

2.3.6 Stemming

Consider the words "walk", "walks", "walking", and "walked". It is evident that all these words are different forms of the word "walk". We can use stemming to reduce each word to its respective stem , i.e. the core meaning-bearing unit of a word. NLTK comes with some common stemming algorithms. Let's use the Porter Stemming algorithm to reduce these four words to their stems:

```
from nltk.stem import PorterStemmer # Import the Porter stemmer from NLTK

porter = PorterStemmer() # Create a Porter stemmer object
for word in ['walk','walks','walking','walked']:
    print(word,"->",porter.stem(word))
```

The output will look like:

```
walk -> walk
walks -> walk
walking -> walk
walked -> walk
```

Let's now stem the words from the alice.txt text using the Porter stemmer:

```
words_stemmed = []
for word in words_lowercase_nostopwords_no_punctuation:
    words_stemmed.append(porter.stem(word))
print(words_stemmed)
```

```
['alic', 'begin', 'get', 'tire', 'sit', 'sister', 'bank', 'noth', 'twice', 'peep', 'book', 'sister', 'read', 'pictur', 'convers', 'use', 'book', 'thought', 'alic', 'without', 'pictur', 'convers', 'consid', 'mind', 'well', 'could', 'hot', 'day', 'made', 'feel', 'sleepi', 'stupid', 'whether', 'pleasur', 'make', 'daisy-chain', 'would', 'worth', 'troubl', 'get', 'pick', 'daisi', 'suddenli', 'white', 'rabbit', 'pink', 'eye', 'ran', 'close', 'noth', 'remark', 'alic', 'think', 'much', 'way', 'hear', 'rabbit', 'say', 'oh', 'dear', 'oh', 'dear', 'shall', 'late', 'thought',
```

```
'afterward', 'occur', 'ought', 'wonder', 'time', 'seem', 'quit', 'natur', 'rabbit', 'actual', 'took', 'watch', 'waistcoat-pocket', 'look', 'hurri', 'alic', 'start', 'feet', 'flash', 'across', 'mind', 'never', 'seen', 'rabbit', 'either', 'waistcoat-pocket', 'watch', 'take', 'burn', 'curios', 'ran', 'across', 'field', 'fortun', 'time', 'see', 'pop', 'larg', 'rabbit-hol', 'hedg']
```

As you can see, the Porter stemming algorithm reduced the words from alice.txt to their stems. However, it is evident that a lot of the stems do not correspond to real words from the English language.

2.3.7 Lemmatisation

Lemmatisation can address this issue by reducing each word to the respective dictionary headword. Let's lemmatise the words "walk", "walks", "walking", and "walked" using NLTK's WordNetLemmatizer:

```
nltk.download('wordnet') # Download the WordNetLemmatizer package
from nltk.stem import WordNetLemmatizer # Import the WordNetLemmatizer
wnl = WordNetLemmatizer() # Create a WordNetLemmatizer object
for word in ['walk','walks','walking','walked']:
    print(word,"->",wnl.lemmatize(word))
```

The output will look like:

```
walk -> walk
walks -> walk
walking -> walking
walked -> walked
```

As you can see, the words "walk" and "walks" were converted to the lemma "walk", but "walking" and "walked" were not changed. The reason for this is that the WordNetLemmatizer considers by default all inputs as nouns, thus "walks" is considered as the plural form of "walk" and is converted to "walk", but "walking" and "walked" are valid lemmas and remain unchanged. To lemmatise all words to their base verb form, we must indicate that we are inputting verb forms:

```
for word in ['walk','walks','walking','walked']:
    print(word,"->",wnl.lemmatize(word,pos='v'))
```

The output will look like:

```
walk -> walk
walks -> walk
walking -> walk
walked -> walk
```

Using the "pos" argument, the input words were handled as verb forms from the WordNetLemmatizer and returned the base form "walk" for all words. The options for "pos" are noun (n), verb (v), adverb (r), adjective (a). More information about the WordNetLemmatizer here: https://www.nltk.org/api/nltk.stem.wordnet.html

Let's lemmatise the text from alice.txt, treating the input as nouns and then as verbs:

```
lemmas_noun = []
lemmas_verb = []
for word in words_lowercase_nostopwords_no_punctuation:
    lemmas_noun.append(wnl.lemmatize(word,pos='n'))
    lemmas_verb.append(wnl.lemmatize(word,pos='v'))
print(lemmas_noun)
print(lemmas_verb)
```

```
['alice', 'beginning', 'get', 'tired', 'sitting', 'sister', 'bank', 'nothing', 'twice', 'peeped',
    'book', 'sister', 'reading', 'picture', 'conversation', 'use', 'book', 'thought', 'alice',
    'without', 'picture', 'conversation', 'considering', 'mind', 'well', 'could', 'hot', 'day',
    'made', 'feel', 'sleepy', 'stupid', 'whether', 'pleasure', 'making', 'daisy-chain', 'would',
    'worth', 'trouble', 'getting', 'picking', 'daisy', 'suddenly', 'white', 'rabbit', 'pink', 'eye',
```

```
'ran', 'close', 'nothing', 'remarkable', 'alice', 'think', 'much', 'way', 'hear', 'rabbit',
    'say', 'oh', 'dear', 'oh', 'dear', 'shall', 'late', 'thought', 'afterwards', 'occurred', 'ought',
    'wondered', 'time', 'seemed', 'quite', 'natural', 'rabbit', 'actually', 'took', 'watch',
    'waistcoat-pocket', 'looked', 'hurried', 'alice', 'started', 'foot', 'flashed', 'across', 'mind',
    'never', 'seen', 'rabbit', 'either', 'waistcoat-pocket', 'watch', 'take', 'burning', 'curiosity',
    'ran', 'across', 'field', 'fortunately', 'time', 'see', 'pop', 'large', 'rabbit-hole', 'hedge']
['alice', 'begin', 'get', 'tire', 'sit', 'sister', 'bank', 'nothing', 'twice', 'peep', 'book',
    'sister', 'read', 'picture', 'conversations', 'use', 'book', 'think', 'alice', 'without',
    'picture', 'conversations', 'consider', 'mind', 'well', 'could', 'hot', 'day', 'make', 'feel',
    'sleepy', 'stupid', 'whether', 'pleasure', 'make', 'daisy-chain', 'would', 'worth', 'trouble',
    'get', 'pick', 'daisies', 'suddenly', 'white', 'rabbit', 'pink', 'eye', 'run', 'close',
    'nothing', 'remarkable', 'alice', 'think', 'much', 'way', 'hear', 'rabbit', 'say', 'oh', 'dear',
    'oh', 'dear', 'shall', 'late', 'think', 'afterwards', 'occur', 'ought', 'wonder', 'time', 'seem',
    'quite', 'natural', 'rabbit', 'actually', 'take', 'watch', 'waistcoat-pocket', 'look', 'hurry',
    'alice', 'start', 'feet', 'flash', 'across', 'mind', 'never', 'see', 'rabbit', 'either',
    'waistcoat-pocket', 'watch', 'take', 'burn', 'curiosity', 'run', 'across', 'field',
    'fortunately', 'time', 'see', 'pop', 'large', 'rabbit-hole', 'hedge']
```

However, this approach is not practical. Ideally, we would like to know what part of speech each word is and use the lemmatiser accordingly.

2.3.8 Part of Speech (POS) tagging

Part of Speech (POS) tagging is used to detect which part of speech each word in a sentence refers to. Let's use NLTK's POS tagging algorithm to assign POS tags to each of the words in the sentence "I had been a student here for a long time.":

```
nltk.download('averaged_perceptron_tagger')
from nltk import pos_tag

pos_tagged_sentence = pos_tag(word_tokenize('I had been a student here for a long time'))
print(pos_tagged_sentence)
```

The output will look like:

```
[('I', 'PRP'), ('had', 'VBD'), ('been', 'VBN'), ('a', 'DT'), ('student', 'NN'), ('here', 'RB'), ('for', 'IN'), ('a', 'DT'), ('long', 'JJ'), ('time', 'NN')]
```

As you can see, each word in the sentence has been annotated with a POS tags. These POS tags refer to the Penn Treebank POS tags (https://catalog.ldc.upenn.edu/docs/LDC95T7/c193.html). However, the WordNetLemmatizer expects different tag names. To address this issue, we have to first convert the Penn Treebank POS tags to the format expected by the WordNetLemmatizer.

```
def penn_to_wordnet(penn_pos_tag):
    tag_dictionary = {'NN':'n', 'JJ':'a','VB':'v', 'RB':'r'}
    try:
        # If the first two characters of the Penn Treebank POS tag are in the ``tag_dictionary''
        return tag_dictionary[penn_pos_tag[:2]]
    except:
        return 'n' # Default to Noun if no mapping avalable.

lemmas = []
for word, tag in pos_tagged_sentence:
    lemmas.append(wnl.lemmatize(word.lower(), pos=penn_to_wordnet(tag)))

print('I have been a student here for a long time')
print(lemmas)
```

```
I had been a student here for a long time
['i', 'have', 'be', 'a', 'student', 'here', 'for', 'a', 'long', 'time']
```

As you can see, the sentence was properly lemmatised using POS tagging to inform the lemmatiser about the part of speech that each word refers to.

Note: Please note that the conversion from Penn Treebank POS tags to WordNet format used here is a simplification. Full conversion tables should be used for better results. Also, the code above uses exception handling via the "try" and "except" statements. You can read more about exception handling in Python from here: https://docs.python.org/3/tutorial/errors.html

Let's now lemmatise the alice.txt text:

```
lemmas_alice = []
for sent in sent_tokenize(text): # Tokenise text into sentences
   pos_tagged_sentence_alice = pos_tag(word_tokenize(sent)) # Get POS tags for each sentence
   for word, tag in pos_tagged_sentence_alice: # Iterate though POS tagged words
        if word.lower() not in punctuation_list: # Ignore words that are punctuation marks
        lemmas_alice.append(wnl.lemmatize(word.lower(), pos=penn_to_wordnet(tag))) # Lemmatise word
print(lemmas_alice)
```

The output will look like:

```
['alice', 'be', 'begin', 'to', 'get', 'very', 'tired', 'of', 'sit', 'by', 'her', 'sister', 'on', 'the', 'bank', 'and', 'of', 'have', 'nothing', 'to', 'do', 'once', 'or', 'twice', 'she', 'have', 'peep', 'into', 'the', 'book', 'her', 'sister', 'be', 'read', 'but', 'it', 'have', 'no', 'picture', 'or', 'conversation', 'in', 'it', 'r', 'and', 'what', 'be', 'the', 'use', 'of', 'a', 'book', "''", 'think', 'alice', '``', 'without', 'picture', 'or', 'conversation', "''", 'so', 'she', 'be', 'consider', 'in', 'her', 'own', 'mind', 'as', 'well', 'a', 'she', 'could', 'for', 'the', 'hot', 'day', 'make', 'her', 'feel', 'very', 'sleepy', 'and', 'stupid', 'whether', 'the', 'pleasure', 'of', 'make', 'a', 'daisy-chain', 'would', 'be', 'worth', 'the', 'trouble', 'of', 'get', 'up', 'and', 'pick', 'the', 'daisy', 'when', 'suddenly', 'a', 'white', 'rabbit', 'with', 'pink', 'eye', 'run', 'close', 'by', 'her', 'there', 'be', 'nothing', 'so', 'very', 'remarkable', 'in', 'that', 'nor', 'do', 'alice', 'think', 'it', 'so', 'very', 'much', 'out', 'of', 'the', 'way', 'to', 'hear', 'the', 'rabbit', 'say', 'to', 'itself', '``', 'oh', 'dear', 'oh', 'dear', 'i', 'shall', 'be', 'late', "''", 'when', 'she', 'think', 'it', 'over', 'afterwards', 'it', 'occur', 'to', 'her', 'that', 'she', 'ought', 'to', 'have', 'wonder', 'at', 'this', 'but', 'at', 'the', 'time', 'it', 'all', 'seem', 'quite', 'natural', 'but', 'when', 'the', 'rabbit', 'actually', 'take', 'a', 'watch', 'out', 'of', 'it', 'waistcoat-pocket', 'and', 'look', 'at', 'it', 'and', 'then', 'that', 'she', 'nave', 'never', 'before', 'see', 'a', 'rabbit', 'with', 'cither', 'a', 'waistcoat-pocket', 'an, 'rabbit', 'with', 'either', 'a', 'waistcoat-pocket', 'or', 'a', 'watch', 'to', 'take', 'out', 'of', 'it', 'and', 'burn', 'with', 'curiosity', 'she', 'run', 'across', 'the', 'field', 'after', 'it', 'and', 'burn', 'with', 'curiosity', 'she', 'run', 'across', 'the', 'field', 'after', 'it', 'and', 'burn', 'with', 'curiosity', 'she', 'run', 'across', 'the', 'field', 'after', 'it', 'and', 'burn', 'with',
```

2.4 Exercises

- Exercise 2.1 Section 2.3.5: Create a new punctuation marks list to address the issue of the remaining "'' and "'' .".
- Exercise 2.2 Section 2.3.5: Remove stop words and punctuation marks without iterating twice through all words.
- Exercise 2.3 Load the text from dune.txt. Compute the number of words in the text, not including punctuation marks.
- Exercise 2.4 Lemmatise the text from dune.txt and print a list of all the lemmas. Remember to convert all words to lowercase and to remove punctuation marks.
- Exercise 2.5 Create a list of the unique lemmas in dune.txt, count their number and print the list and the number of lemmas.
- Exercise 2.6 Create a custom function to divide English text into sentences without using NLTK or other tokenisers. Consider that a sentence ends when one of the following characters occurs: ".", "?", "!". Also remember to take into consideration the line change character "\n". Test your function on the text from dune.txt.

Workshop 3: Text representation

In this lab we will work on different ways to represent text.

3.1 Load corpus

We will use the Gutenberg corpus from NLTK. Let's first load the Gutenberg corpus:

```
import nltk # Import the NLTK library
nltk.download('gutenberg') # Download the gutenberg corpus
from nltk.corpus import gutenberg # Import the gutenberg corpus from NLTK
gutenberg.fileids() # List file-ids in the corpus
```

The output will look like:

```
['austen-emma.txt',
 'austen-persuasion.txt',
 'austen-sense.txt',
 'bible-kjv.txt',
 'blake-poems.txt'
 'bryant-stories.txt',
 'burgess-busterbrown.txt',
 'carroll-alice.txt',
 'chesterton-ball.txt'
 'chesterton-brown.txt',
 'chesterton-thursday.txt',
 'edgeworth-parents.txt',
 'melville-moby_dick.txt',
 'milton-paradise.txt',
 'shakespeare-caesar.txt',
 'shakespeare-hamlet.txt',
 'shakespeare-macbeth.txt',
 'whitman-leaves.txt']
```

Let's compute some statistics for each document in the Gutenberg corpus, like the number of words, the number of sentences, and the number of characters:

```
print("Chars\tWords\tSents\tFile")
for fileid in gutenberg.fileids(): # Iterate through files in corpus
    num_chars = len(gutenberg.raw(fileid))
    num_words = len(gutenberg.words(fileid))
    num_sents = len(gutenberg.sents(fileid))
    \label{eq:print("%7.0f\t%7.0f\t%s" % (num_sents,num_words,num_chars,fileid))} print("%7.0f\t%7.0f\t%s" % (num_sents,num_words,num_chars,fileid))
```

```
Words Sents File
Chars
887071 192427
                7752 austen-emma.txt
466292
       98171
              3747 austen-persuasion.txt
673022 141576 4999 austen-sense.txt
4332554 1010654 30103 bible-kjv.txt
 38153 8354 438 blake-poems.txt
```

```
249439
         55563
                  2863 bryant-stories.txt
 84663
         18963
                  1054 burgess-busterbrown.txt
144395
         34110
                  1703 carroll-alice.txt
457450
         96996
                  4779 chesterton-ball.txt
406629
         86063
                  3806 chesterton-brown.txt
320525
        69213
                  3742 chesterton-thursday.txt
935158
        210663
                 10230 edgeworth-parents.txt
1242990
        260819
                 10059 melville-moby_dick.txt
468220
         96825
                  1851 milton-paradise.txt
         25833
112310
                  2163 shakespeare-caesar.txt
162881
         37360
                  3106 shakespeare-hamlet.txt
100351
         23140
                  1907 shakespeare-macbeth.txt
711215 154883
                  4250 whitman-leaves.txt
```

3.2 Vocabulary

As we can see above, each document consists of a few thousand words. But these words are not unique. Natural language consists of words that convey meaning and are re-used and combined to form different sentences. The set of unique words that are used in each document constitutes its vocabulary.

3.2.1 Words in corpus

Consider the text "The new table is red. The blue table is broken." Let's compute its vocabulary. First we should tokenise the text into words, convert them all to lowercase and remove punctuation marks:

```
from nltk import word_tokenize # Import the word_tokenize function from NLTK

from string import punctuation

punctuation_list = list(punctuation) # Convert string with punctuation marks to list

text = "The new table is red. The blue table is broken."

text_tokens_processed = []

text_tokens = word_tokenize(text) # Tokenise the text into words

for token in text_tokens: # Iterate through the available tokens
    if token not in punctuation_list: # Omit tokens that are punctuation marks
        text_tokens_processed.append(token.lower()) # Add lowercase version of token to list

print("List of processed words:",text_tokens_processed)

The output will look like:

List of processed words: ['the', 'new', 'table', 'is', 'red', 'the', 'blue', 'table', 'is', 'broken']
```

3.2.2 Unique words in corpus

As you can see in the list of words, the words "the", "table" and "is" appear two times each in the text. Let's now compute the vocabulary of this text, i.e. the list of unique words used in this text. To achieve this, we will use Python's **set** type. A set is similar to a list but allows only unique elements.

```
vocabulary = set() # Create an empty set
for word in text_tokens_processed: # Iterate through available words
    vocabulary.add(word) # Add word to set

print("Vocabulary:",vocabulary)
print("Vocabulary size:",len(vocabulary))

vocabulary2 = set(text_tokens_processed)
print("\nVocabulary2:",vocabulary2)
print("Vocabulary2 size:",len(vocabulary2))
```

```
Vocabulary: {'the', 'red', 'broken', 'blue', 'table', 'is', 'new'}
Vocabulary size: 7

Vocabulary2: {'the', 'red', 'broken', 'blue', 'table', 'is', 'new'}
Vocabulary2 size: 7
```

As you can see, the vocabulary used by the text "The new table is red. The blue table is broken." consists of the seven following words: is, table, blue, red, broken, new, the. Note that you don't have to add the contents of a list in a set one-by-one, as shown for variable "vocabulary2".

3.2.3 Vocabulary of multiple documents

Let's now compute the vocabulary for each document in the Gutenberg corpus:

```
for fileid in gutenberg.fileids(): # Iterate through documents in corpus
  vocabulary_of_document = set() # Create empty set
  for word in gutenberg.words(fileid): # Iterate through words in document
      if word not in punctuation_list: # Omit tokens that are punctuation marks
           vocabulary_of_document.add(word.lower())
  print("%6.0f\t%s" % (len(vocabulary_of_document),fileid))
```

The output will look like:

```
7328 austen-emma.txt
5820
      austen-persuasion.txt
6388 austen-sense.txt
12755 bible-kjv.txt
1521 blake-poems.txt
3925 bryant-stories.txt
1547 burgess-busterbrown.txt
2622 carroll-alice.txt
8313 chesterton-ball.txt
7780 chesterton-brown.txt
6335 chesterton-thursday.txt
8432 edgeworth-parents.txt
17215 melville-moby_dick.txt
9007
      milton-paradise.txt
3019
      shakespeare-caesar.txt
4703
      shakespeare-hamlet.txt
3451
      shakespeare-macbeth.txt
12437
      whitman-leaves.txt
```

As you can see, we computed the vocabulary for each document in the Gutenberg corpus and printed its size. However, vocabularies from different documents are expected to have similar words in them since all texts are in the same language. Let's compute the vocabulary that covers all documents in the dataset:

```
vocabulary_of_corpus = set() # Create empty set
for fileid in gutenberg.fileids(): # Iterate through documents in corpus
    for word in gutenberg.words(fileid): # Iterate through words in document
        if word not in punctuation_list: # Omit tokens that are punctuation marks
            vocabulary_of_corpus.add(word.lower())

print("Vocabulary of Gutenberg corpus:",len(vocabulary_of_corpus),"words")
```

The output will look like:

```
Vocabulary of Gutenberg corpus: 42314 words
```

As you can see, the vocabulary that covers all documents in the Gutenberg corpus is much larger than individual document vocabularies but significantly smaller that the sum of all the individual vocabularies.

One-hot encoding 3.3

One-hot encoding of words in vocabulary 3.3.1

Consider again the text "The new table is red. The blue table is broken." We have already computed its

```
vocabulary and would like to compute the One-Hot representation of each word in the vocabulary.
from numpy import array # Import array type from numpy
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
vocabulary = ['is', 'table', 'blue', 'red', 'broken', 'new', 'the']
data = array(vocabulary) # Convert to array because it is required by the LabelEncoder() object
print(data,"\n")
# Integer encoding - Assigns a unique index to each unique word
label_encoder = LabelEncoder()
integer_encoded = label_encoder.fit_transform(data)
print(integer_encoded,"\n")
# One-Hot encoding - Assigns a One-Hot binary representation to each word
onehot_encoder = OneHotEncoder(sparse_output=False)
integer_encoded = integer_encoded.reshape(len(integer_encoded), 1)
onehot_encoded = onehot_encoder.fit_transform(integer_encoded)
print(onehot_encoded,"\n")
for i in range(len(data)):
   print(onehot_encoded[i],"->",data[i])
The output will look like:
['is' 'table' 'blue' 'red' 'broken' 'new' 'the']
[2 5 0 4 1 3 6]
[[0. 0. 1. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 1. 0.]
 [1. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 1. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0.]
 [0. 0. 0. 1. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 1.]]
[0. \ 0. \ 1. \ 0. \ 0. \ 0. \ 0.] \rightarrow is
[0. 0. 0. 0. 0. 1. 0.] -> table
[1. 0. 0. 0. 0. 0. 0.] -> blue
[0. 0. 0. 0. 1. 0. 0.] \rightarrow red
[0. 1. 0. 0. 0. 0. 0.] -> broken
[0. 0. 0. 1. 0. 0. 0.] -> new
[0. 0. 0. 0. 0. 0. 1.] \rightarrow the
Let's create a dictionary to easily encode our text and one-hot encode the word "red":
dictionary = {}
for i in range(len(data)):
   dictionary[data[i]] = onehot_encoded[i]
print(dictionary)
print("\nred =",dictionary['red'])
The output will look like:
```

{'is': array([0., 0., 1., 0., 0., 0., 0.]), 'table': array([0., 0., 0., 0., 0., 1., 0.]), 'blue': array([1., 0., 0., 0., 0., 0., 0.]), 'red': array([0., 0., 0., 0., 1., 0., 0.]), 'broken': array([0., 1., 0., 0., 0., 0., 0.]), 'new': array([0., 0., 0., 1., 0., 0., 0.]), 'the':

```
array([0., 0., 0., 0., 0., 0., 1.])}
red = [0. 0. 0. 0. 1. 0. 0.]
```

Consider the one-hot encoded word (0, 0, 0, 0, 1, 0, 0). How can we convert it back to its respective real word? Let's create a function to do this:

```
def get_label_from_dictionary(dictionary,value):
    for word, one_hot in dictionary.items(): # Iterate all pairs (word, one-hot representation) in
    the dictionary
        if (one_hot == value).all(): # Compare equality between numpy arrays
            return word

print(get_label_from_dictionary(dictionary,[0,0,0,0,1,0,0]))
```

The output will look like:

red

3.3.2 One-hot encoding of text

Consider again the text "The new table is red. The blue table is broken." We have already computed its vocabulary and the one-hot representation of each word in the vocabulary. How can we one-hot encode the whole text? To do so, we have to perform a logical OR operation between the one-hot vectors of its constituent words.

```
from numpy import logical_or # Import the element-wise logical OR function from numpy
from numpy import zeros # Import the zeros function from numpy

text = "The new table is red. The blue table is broken."
print("List of processed words:",text_tokens_processed)

result = zeros(len(dictionary)) # Start from zero-valued vector - Convert to numpy array
for word in text_tokens_processed: # Iterate words in text
    print(result.astype(int), "OR", dictionary[word],"= ",end='')
    result = logical_or(result,dictionary[word]) # Compute the element-wise logical or between the
    partial result and the one-hot representation of the word
    print(result.astype(int))

print("\nOne-Hot encoded text:",result.astype(int))
```

The output will look like:

```
List of processed words: ['the', 'new', 'table', 'is', 'red', 'the', 'blue', 'table', 'is', 'broken']

[0 0 0 0 0 0 0] OR [0. 0. 0. 0. 0. 0. 1.] = [0 0 0 0 0 0 0 1]

[0 0 0 0 0 0 1] OR [0. 0. 0. 1. 0. 0. 0.] = [0 0 0 1 0 0 1]

[0 0 0 1 0 0 1] OR [0. 0. 0. 0. 1. 0.] = [0 0 0 1 0 1 1]

[0 0 0 1 0 1 1] OR [0. 0. 1. 0. 0. 0.] = [0 0 1 1 0 1 1]

[0 0 1 1 0 1 1] OR [0. 0. 0. 0. 0. 1. 0. 0.] = [0 0 1 1 1 1 1]

[0 0 1 1 1 1 1] OR [0. 0. 0. 0. 0. 0. 1.] = [0 0 1 1 1 1 1]

[0 0 1 1 1 1 1] OR [1. 0. 0. 0. 0. 0.] = [1 0 1 1 1 1 1]

[1 0 1 1 1 1 1] OR [0. 0. 1. 0. 0. 0.] = [1 0 1 1 1 1 1]

[1 0 1 1 1 1 1] OR [0. 0. 1. 0. 0. 0. 0.] = [1 1 1 1 1 1]

One-Hot encoded text: [1 1 1 1 1 1]
```

Let's create a function for one-hot encoding text and do the same for the text "the broken table":

```
def get_text_one_hot_encoding(words_list, dictionary):
    result = zeros(len(dictionary)) # Start from zero-valued vector - Convert to numpy array
    for word in words_list: # Iterate words in text
        result = logical_or(result,dictionary[word]) # Compute the element-wise logical or between the
    partial result and the one-hot representation of the word
```

```
return result.astype(int)
words_list = ['the','broken','table']
print("\nOne-Hot encoded text:",get_text_one_hot_encoding(words_list,dictionary))
The output will look like:
One-Hot encoded text: [0 1 0 0 0 1 1]
```

3.4 Term Frequency (TF) representation

3.4.1 Compute TF of words in text

Let's now compute the term frequency of each word in the text "The new table is red. The blue table is broken."

```
from nltk import FreqDist # Import the FreqDist function from NLTK

text = "The new table is red. The blue table is broken."

text_tokens_processed = ['the', 'new', 'table', 'is', 'red', 'the', 'blue', 'table', 'is', 'broken']

vocabulary = {'new', 'broken', 'the', 'blue', 'table', 'red', 'is'}

tf = FreqDist(text_tokens_processed) # Compute term frequency of words

print(tf,"\n")

vocabulary = sorted(vocabulary) # Sort alphabetiacally for better presentation

for word in vocabulary:
    print("%5.0f %s" % (tf[word],word))
```

The output will look like:

```
<FreqDist with 7 samples and 10 outcomes>

1 blue
1 broken
2 is
1 new
1 red
2 table
2 the
```

3.4.2 TF representation of documents

Let's now compute the TF representation of the text "The new table is red. The blue table is broken." and the text "The new table is broken":

```
text_tf = []
for word in vocabulary:
    text_tf.append(tf[word])

print(text_tf,"->",text)

text2 = "The new table is broken"
text2_tokens_processed = ['the','new','table','is','broken']
tf2 = FreqDist(text2_tokens_processed) # Compute term frequency of words

text2_tf = []
for word in vocabulary:
    text2_tf.append(tf2[word])

print(text2_tf,"->",text2)
```

The output will look like:

```
[1, 1, 2, 1, 1, 2, 2] -> The new table is red. The blue table is broken.
[0, 1, 1, 0, 1, 1] -> The new table is broken
```

3.5 Term Frequency - Inverse Document Frequency (TF-IDF)

Consider a corpus consisting of the three following documents:

- 1. "The new table is red. The blue table was broken."
- 2. "The new movie that we watched yesterday was terrible."
- 3. "We raised the red and blue flag yesterday."

3.5.1 Document Frequency (DF)

Let's compute the Document Frequency (DF) of each word in the above corpus. Remember that the DF of a word is equal to the number of documents in a corpus that the word appears in.

```
# Create a list of the lists of lowercase words without punctuation marks for each document
texts_words_processed = []
texts_words_processed.append(['the','new','table','is','red','the','blue','table','was','broken'])
texts_words_processed.append(['the','new','movie','that','we','watched','yesterday','was','terrible'])
texts_words_processed.append(['we','raised','the','red','and','blue','flag','yesterday'])
print(texts_words_processed)
# Create the vocabulary
vocabulary_texts = set()
for doc in texts_words_processed:
   for word in doc:
       vocabulary_texts.add(word)
vocabulary_texts = sorted(vocabulary_texts) # Sort vocabulary alphabetically for better presentation
print("\nVocabulary:",vocabulary_texts)
DF = dict() # Create an empty dictionary
for word in vocabulary_texts: # Iterate through words in vocabulary
   cnt = 0
   for doc in texts_words_processed: # Iterate through documents
       if word in doc:
          cnt += 1 # cnt += 1 is equal to cnt = cnt + 1
   DF[word] = cnt
print("\nDocument frequencies:",DF)
The output will look like:
[['the', 'new', 'table', 'is', 'red', 'the', 'blue', 'table', 'was', 'broken'], ['the', 'new',
    'movie', 'that', 'we', 'watched', 'yesterday', 'was', 'terrible'], ['we', 'raised', 'the', 'red',
    'and', 'blue', 'flag', 'yesterday']]
Vocabulary: ['and', 'blue', 'broken', 'flag', 'is', 'movie', 'new', 'raised', 'red', 'table',
    'terrible', 'that', 'the', 'was', 'watched', 'we', 'yesterday']
Document frequencies: {'and': 1, 'blue': 2, 'broken': 1, 'flag': 1, 'is': 1, 'movie': 1, 'new': 2,
    'raised': 1, 'red': 2, 'table': 1, 'terrible': 1, 'that': 1, 'the': 3, 'was': 2, 'watched': 1,
    'we': 2, 'yesterday': 2}
```

3.5.2 Inverse Document Frequency (IDF)

The Inverse Document Frequency (IDF) of a word is the logarithmically scaled inverse fraction of the documents that contain the word (obtained by dividing the total number of documents by the number of documents

containing the term, and then taking the logarithm of that quotient). IDF is defined as:

$$IDF(t,D) = \log\left(\frac{N}{DF(t,D)}\right)$$
 (3.1)

where t is a word (term), D the corpus, and N the number of documents in the corpus.

Let's compute the IDF for the examined corpus:

The output will look like:

As you can see above, the more documents that a word appeared in, the lower the IDF for the word. Note that IDF is 0 for the word "the" that appeared in all documents (as a consequence of $N = DF \Rightarrow IDF = log1 = 0$)

3.5.3 Term Frequency - Inverse Document Frequency (TF-IDF)

Term Frequency - Inverse Document Frequency (TF-IDF) is defined as the product of TF and IDF:

$$TFIDF(t, d, D) = TF(t, d) \cdot IDF(t, D)$$
(3.2)

where d is a document of corpus D.

Let's compute the TF-IDF for each word in each document of the examined corpus

```
from nltk import FreqDist # Import the FreqDist function from NLTK
TF = []
for doc in texts_words_processed: # Iterate though documents
   TF.append(FreqDist(doc)) # Compute word frequency
print(TF,"\n")
TFIDF = []
for tf_doc in TF:
   tfidf_doc = dict()
   for word in vocabulary_texts: # Iterate through words in vocabulary
       tfidf_doc[word] = tf_doc[word] * IDF[word] # Compute TF-IDF - tf_doc is of type FreqDist and
    will return 0 for words that don't exist
   TFIDF.append(tfidf_doc)
cnt = 0
for tfidf_doc in TFIDF:
   print("Text",cnt,"TF-IDF:",tfidf_doc,"\n")
   cnt += 1
```

The output will look like:

```
[FreqDist({'the': 2, 'table': 2, 'new': 1, 'is': 1, 'red': 1, 'blue': 1, 'was': 1, 'broken': 1}),
    FreqDist({'the': 1, 'new': 1, 'movie': 1, 'that': 1, 'we': 1, 'watched': 1, 'yesterday': 1,
    'was': 1, 'terrible': 1}), FreqDist({'we': 1, 'raised': 1, 'the': 1, 'red': 1, 'and': 1, 'blue':
    1, 'flag': 1, 'yesterday': 1})]
Text 0 TF-IDF: {'and': 0.0, 'blue': 0.4054651081081644, 'broken': 1.0986122886681098, 'flag': 0.0,
    'is': 1.0986122886681098, 'movie': 0.0, 'new': 0.4054651081081644, 'raised': 0.0, 'red':
    0.4054651081081644, 'table': 2.1972245773362196, 'terrible': 0.0, 'that': 0.0, 'the': 0.0, 'was':
    0.4054651081081644, 'watched': 0.0, 'we': 0.0, 'yesterday': 0.0}
Text 1 TF-IDF: {'and': 0.0, 'blue': 0.0, 'broken': 0.0, 'flag': 0.0, 'is': 0.0, 'movie':
    1.0986122886681098, 'new': 0.4054651081081644, 'raised': 0.0, 'red': 0.0, 'table': 0.0,
    'terrible': 1.0986122886681098, 'that': 1.0986122886681098, 'the': 0.0, 'was':
    0.4054651081081644, 'watched': 1.0986122886681098, 'we': 0.4054651081081644, 'yesterday':
    0.4054651081081644}
Text 2 TF-IDF: {'and': 1.0986122886681098, 'blue': 0.4054651081081644, 'broken': 0.0, 'flag':
    1.0986122886681098, 'is': 0.0, 'movie': 0.0, 'new': 0.0, 'raised': 1.0986122886681098, 'red':
    0.4054651081081644, 'table': 0.0, 'terrible': 0.0, 'that': 0.0, 'the': 0.0, 'was': 0.0,
    'watched': 0.0, 'we': 0.4054651081081644, 'yesterday': 0.4054651081081644}
```

3.6 Exercises

- Exercise 3.1 Create a dictionary with the one-hot encoding of the vocabulary of the carroll-alice.txt file from the Gutenberg corpus.
- Exercise 3.2 Based on the one-hot encoding of the vocabulary of the carroll-alice.txt from the Gutenberg corpus, one-hot encode the sentences "this is an old house", "this is a new house", and "he left his house" and compute their cosine distance.
- Exercise 3.3 Using the vocabulary of the carroll-alice.txt document from the Gutenberg corpus, compute the TF-IDF representations and the respective cosine and euclidean distances of the documents in a corpus containing alice.txt and dune.txt.
- Exercise 3.4 Use alice.txt to create a vocabulary and add to it the "unknown" word "<UNK>". Use this vocabulary to create the one-hot, the TF, the log normalised TF, and the TF-IDF representations of the alice.txt and the dune.txt documents.

Appendix: Test files

A.1 movies.xml

```
<movies>
 <movie id="1">
   <title>And Now for Something Completely Different</title>
   <released>1971</released>
 </movie>
 <movie id="2">
   <title>Monty Python and the Holy Grail</title>
   <released>1974</released>
 </movie>
 <movie id="3">
   <title>Monty Python's Life of Brian</title>
   <released>1979</released>
 </movie>
   <title>Monty Python Live at the Hollywood Bowl</title>
   <released>1982</released>
 </movie>
 <movie id="5">
   <title>Monty Python's The Meaning of Life</title>
   <released>1983</released>
 </movie>
</movies>
```

A.2 links.html

38 Appendix: Test files

A.3 emails.txt

john+acme.co@hotmail.com bob@gmail.com tom@durham.ac.uk jerry@durham.ac.uk scrooge@durham.ac.uk donald@yahoo.co.uk huey@yahoo.co.uk dewey@gmail.com louie.duck@durham.ac.uk ${\it gyro.gear loose@yahoo.co.uk}$ bart@yahoo.co.uk homer@gmail.com stan@hotmail.com kyle-broflovski@durham.ac.uk eric@vahoo.co.uk kenny@gmail.com butters@durham.ac.uk wendy@hotmail.com $randy_marsh@durham.ac.uk$ chef@gmail.com

A.4 cds.xml

```
<?xml version="1.0" encoding="IS08859-1" ?>
<CATALOG>
<CD>
<TITLE>Empire Burlesque</TITLE>
<ARTIST>Bob Dylan</ARTIST>
<COUNTRY>USA</COUNTRY>
<COMPANY>Columbia</COMPANY>
<PRICE>10.90</PRICE>
<YEAR>1985</YEAR>
</CD>
<CD>
<TITLE>Hide your heart</TITLE>
<ARTIST>Bonnie Tylor</ARTIST>
<COUNTRY>UK</COUNTRY>
<COMPANY>CBS Records</COMPANY>
<PRICE>9.90</PRICE>
<YEAR>1988</YEAR>
</CD>
<CD>
<TITLE>Greatest Hits</TITLE>
<ARTIST>Dolly Parton</ARTIST>
<COUNTRY>USA</COUNTRY>
<COMPANY>RCA</COMPANY>
<PRICE>9.90</PRICE>
<YEAR>1982</YEAR>
</CD>
<CD>
<TITLE>Still got the blues</TITLE>
<ARTIST>Gary More</ARTIST>
<COUNTRY>UK</COUNTRY>
<COMPANY>Virgin redords</COMPANY>
<PRICE>10.20</PRICE>
<YEAR>1990</YEAR>
</CD>
<CD>
<TITLE>Eros</TITLE>
<ARTIST>Eros Ramazzotti</ARTIST>
```

```
<COUNTRY>EU</COUNTRY>
<COMPANY>BMG</COMPANY>
<PRICE>9.90</PRICE>
<YEAR>1997</YEAR>
</CD>
<CD>
<TITLE>One night only</TITLE>
<ARTIST>Bee Gees
<COUNTRY>UK</COUNTRY>
<COMPANY>Polydor</COMPANY>
<PRICE>10.90</PRICE>
<YEAR>1998</YEAR>
</CD>
<CD>
<TITLE>Sylvias Mother</TITLE>
<ARTIST>Dr.Hook</ARTIST>
<COUNTRY>UK</COUNTRY>
<COMPANY>CBS</COMPANY>
<PRICE>8.10</PRICE>
<YEAR>1973</YEAR>
</CD>
<CD>
<TITLE>Maggie May</TITLE>
<ARTIST>Rod Stewart</ARTIST>
<COUNTRY>UK</COUNTRY>
<COMPANY>Pickwick</COMPANY>
<PRICE>8.50</PRICE>
<YEAR>1990</YEAR>
</CD>
<CD>
<TITLE>Romanza</TITLE>
<ARTIST>Andrea Bocelli</ARTIST>
<COUNTRY>EU</COUNTRY>
<COMPANY>Polydor</COMPANY>
<PRICE>10.80</PRICE>
<YEAR>1996</YEAR>
</CD>
<CD>
<TITLE>When a man loves a woman</TITLE>
<ARTIST>Percy Sledge</ARTIST>
<COUNTRY>USA</COUNTRY>
<COMPANY>Atlantic</COMPANY>
<PRICE>8.70</PRICE>
<YEAR>1987</YEAR>
</CD>
<CD>
<TITLE>Black angel</TITLE>
<ARTIST>Savage Rose</ARTIST>
<COUNTRY>EU</COUNTRY>
<COMPANY>Mega</COMPANY>
<PRICE>10.90</PRICE>
<YEAR>1995</YEAR>
</CD>
<CD>
<TITLE>1999 Grammy Nominees</TITLE>
<ARTIST>Many</ARTIST>
<COUNTRY>USA</COUNTRY>
<COMPANY>Grammy</COMPANY>
<PRICE>10.20</PRICE>
<YEAR>1999</YEAR>
</CD>
<CD>
<TITLE>For the good times</TITLE>
<ARTIST>Kenny Rogers</ARTIST>
<COUNTRY>UK</COUNTRY>
<COMPANY>Mucik Master</COMPANY>
```

<PRICE>8.70</PRICE> <YEAR>1995</YEAR> </CD> <CD> <TITLE>Big Willie style</TITLE> <ARTIST>Will Smith</ARTIST> <COUNTRY>USA</COUNTRY> <COMPANY>Columbia</COMPANY> <PRICE>9.90</PRICE> <YEAR>1997</YEAR> </CD> <CD> <TITLE>Tupelo Honey</TITLE> <ARTIST>Van Morrison</ARTIST> <COUNTRY>UK</COUNTRY> <COMPANY>Polydor</COMPANY> <PRICE>8.20</PRICE> <YEAR>1971</YEAR> </CD> <CD> <TITLE>Soulsville</TITLE> <ARTIST>Jorn Hoel</ARTIST> <COUNTRY>Norway</COUNTRY> <COMPANY>WEA</COMPANY> <PRICE>7.90</PRICE> <YEAR>1996</YEAR> </CD> <CD> <TITLE>The very best of</TITLE> <ARTIST>Cat Stevens <COUNTRY>UK</COUNTRY> <COMPANY>Island</COMPANY> <PRICE>8.90</PRICE> <YEAR>1990</YEAR> </CD> <CD> <TITLE>Stop</TITLE> <ARTIST>Sam Brown</ARTIST> <COUNTRY>UK</COUNTRY> <COMPANY>A and M</COMPANY> <PRICE>8.90</PRICE> <YEAR>1988</YEAR> </CD> <CD> <TITLE>Bridge of Spies</TITLE> <ARTIST>T`Pau</ARTIST> <COUNTRY>UK</COUNTRY> <COMPANY>Siren</COMPANY> <PRICE>7.90</PRICE> <YEAR>1987</YEAR> </CD> <CD> <TITLE>Private Dancer</TITLE> <ARTIST>Tina Turner</ARTIST> <COUNTRY>UK</COUNTRY> <COMPANY>Capitol</COMPANY> <PRICE>8.90</PRICE> <YEAR>1983</YEAR> </CD> <CD> <TITLE>Midt om natten</TITLE> <ARTIST>Kim Larsen</ARTIST> <COUNTRY>EU</COUNTRY> <COMPANY>Medley</COMPANY> <PRICE>7.80</PRICE>

<YEAR>1983</YEAR>

</CD> <CD> <TITLE>Pavarotti Gala Concert</TITLE> <ARTIST>Luciano Pavarotti</ARTIST> <COUNTRY>UK</COUNTRY> <COMPANY>DECCA</COMPANY> <PRICE>9.90</PRICE> <YEAR>1991</YEAR> </CD> <CD> <TITLE>The dock of the bay</TITLE> <ARTIST>Otis Redding</ARTIST> <COUNTRY>USA</COUNTRY> <COMPANY>Atlantic</COMPANY> <PRICE>7.90</PRICE> <YEAR>1987</YEAR> </CD> <CD> <TITLE>Picture book</TITLE> <ARTIST>Simply Red</ARTIST> <COUNTRY>EU</COUNTRY> <COMPANY>Elektra</COMPANY> <PRICE>7.20</PRICE> <YEAR>1985</YEAR> </CD> <CD> <TITLE>Red</TITLE> <ARTIST>The Communards <COUNTRY>UK</COUNTRY> <COMPANY>London</COMPANY> <PRICE>7.80</PRICE> <YEAR>1987</YEAR> </CD> <C:D> <TITLE>Unchain my heart</TITLE> <ARTIST>Joe Cocker</ARTIST> <COUNTRY>USA</COUNTRY> <COMPANY>EMI</COMPANY> <PRICE>8.20</PRICE> <YEAR>1987</YEAR> </CD> </CATALOG>

A.5 alice.txt

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book," thought Alice "without pictures or conversations?"

So she was considering in her own mind (as well as she could, for the hot day made her feel very sleepy and stupid), whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a White Rabbit with pink eyes ran close by her.

There was nothing so very remarkable in that; nor did Alice think it so very much out of the way to hear the Rabbit say to itself, "Oh dear! I shall be late!" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually took a watch out of its waistcoat-pocket, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and burning with curiosity, she ran across the field after it, and fortunately was just in time to see it pop down a large rabbit-hole under the hedge.

42 Appendix: Test files

A.6 dune.txt

In the week before their departure to Arrakis, when all the final scurrying about had reached a nearly unbearable frenzy, an old crone came to visit the mother of the boy, Paul.

It was a warm night at Castle Caladan, and the ancient pile of stone that had served the Atreides family as home for twenty-six generations bore that cooled-sweat feeling it acquired before a change in the weather.

The old woman was let in by the side door down the vaulted passage by Paul's room and she was allowed a moment to peer in at him where he lay in his bed.

By the half-light of a suspensor lamp, dimmed and hanging near the floor, the awakened boy could see a bulky female shape at his door, standing one step ahead of his mother. The old woman was a witch shadow - hair like matted spiderwebs, hooded 'round darkness of features, eyes like glittering jewels.