



DANMARKS TEKNISKE UNIVERSITET

01415

Computational Tools for Big Data

---

# Final Project

---

*Anonymous*

December 2016

# 1 Introduction

In this project, we will use the Project Gutenberg Corpus documents, which contains more than 150 GB(500 GB with audio and images) of uncleaned data. We will process the text documents of this dataset to build a system that is able to estimate the metadata (topic, popularity, author...) of a given document. Also, we will implement an Information Retrieval System, in which, given a document as a query, the system will output a list of the most similar documents, this way, a person can input its favourite novel, and find more books of the same kind.

## 1.1 Structure of the report

This report is separated in sections which are describing the different parts of the system. In each section we are representing our ideas but we are not including complete code samples in the sections, just a simplified pseudocode functions which are more accessible for a reader. If the reader is interested in complete source code he can find it in the appendices.

## 1.2 Development environment

We are working on a desktop computer which is now running as a server for the calculations. All work is done in a virtual container with access to 8GB RAM and quad core processor and the team members have access to the resources via ssh and rdp protocols.

# 2 Dataset

## 2.1 Data acquisition

We used for our project a dataset from Project Gutenberg. We downloaded the data by mirroring host server. [1] This can be achieved by running following command.

---

```
rsync -av --del ftp@ftp.ibiblio.org::gutenberg /gutenbergDataRsync
```

---

Size of the downloaded data is around 150 GB excluding an audio and images files. Because our aim for this project is to work with text files we are not using them in our analyse.

## 2.2 Data explanation

The downloaded data consist of ebooks files mostly in txt, pdf, html and epub file formats. The book files are mapped to unique *ID* which is connecting them with other information. For example the book *Frankenstein; Or, The Modern Prometheus* belongs to the id 84 and therefore there exists files 84.txt, 84.epub, 84.rdf and etc. We are using this ID to access multiple format for a single book and also for extracting meta-data information.

Metadata are stored in RDF format and contain most important facts about the book as author, title, downloads and topic, we are parsing those files to be able train our prediction.

---

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//
EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
<html lang="en" xml:lang="en" xmlns="http://www.w3.org/1999/xhtml">
<head>
<meta http-equiv="Content-Type" content="text/html; charset=UTF-8" />
<title>Zero the Slaver: A Romance of Equatorial Africa</title>
<meta content="text/html; charset=utf-8" http-equiv="Content-Type"/>
<meta content="Lawrence Fletcher" name="Author"/>
<meta content="Zero the Slaver" name="Description"/>
<style type="text/css">
```

---

Listing 1: Selection from a XHTML book file

---

```

<pgterms:downloads rdf:datatype="http://www.w3.org/2001/XMLSchema#integer">2611
</pgterms:downloads>
  <dcterms:hasFormat>
    <pgterms:file rdf:about="http://www.gutenberg.org/ebooks/10.txt.utf-8">
      <dcterms:format>
        <rdf:Description rdf:nodeID="Nbe07b875f95b483fbc4dead440533f41">
          <rdf:value rdf:datatype="http://purl.org/dc/terms/IMT">
            text/plain</rdf:value>
          <dcam:memberOf rdf:resource="http://purl.org/dc/terms/IMT"/>
        </rdf:Description>
      </dcterms:format>
    </pgterms:file>
  </dcterms:hasFormat>

```

---

Listing 2: Selection from a RDF metadata file

## 3 Data Preprocessing

### 3.1 Data standardization

As first step for data preparation we are converting different data formats to cleaned text. Even though the DDBB (corpus) of our system is based on the Gutenberg project database, we can process files from different kinds of sources, so that our document query can be a wide range of file type. We can process:

- Plain text files
- HTML documents stored in the computer DDBB
- URLs
- pdfs

Support for different format types can be easily added.

### 3.2 Text cleaning

Once the relevant plain text is obtained the NLTK 3.0 Python library is used to preprocess the data, performing a correction and normalization of the words and expressions which we will further be used for representing each document. As we will explain later, this preprocessing will be performed in parallel for each document for obvious reasons (the corpus does not fit into memory). This preprocessing stage consists of the following steps:

- Removing accents: In case there are special characters for some reason, we first remove them.
- Tokenization: It divides the whole plain text String into separate tokens, in this case String words.
- Lower Case: Since upper letters are not important for the purposes of the task, every letter is turned into lower case so that the IR system is not case sensitive.
- Removing Alpha-numeric words: It removes all the punctuation and words that are not a combination of letters and numbers.
- Removing stopwords: It removes common words such as articles or junctions that are used in language to glue subjects. These words are not very relevant to the IR system because they don't offer information about the context themselves. The NLTK corpus of English stopwords is used in order to identify them.
- Steaming: Normalize verbs, plurals and compositions so that they are treated as the same. The Python implementation of the Snowball steamers is used.

### 3.3 Metadata extraction

As we mentioned before the metadata are stored in a RDF format. We are loading the RDF file as XML document, defining namespace to be able to access elements and extracting the desired information from data. The implementation is described in following function.

---

```
import xml.etree.ElementTree as et #XML parser

def extract_metadata(filename):
    #load rdf file as tree
    tree = et.parse(filename)
    #namespace definition
    ns = {"dcterms": "http://purl.org/dc/terms/", "pgterms":
    ↪ "http://www.gutenberg.org/2009/pgterms/"}
    #extract data
    for entry in tree.findall("./dcterms:subject/rdf:Description/rdf:value", ns):
        subjects.append(entry.text)
    return {'Subject':subjects}

#example call
id = '10015'
extract_metadata(id+'/pg'+id+'.rdf')
```

---

Listing 3: A simplified function to extract information from a RDF metadata file

## 4 Implementation

The following figure is showing how the implemented system works. As it can be observed, our system has an initial preprocessing step in which the documents of the system are treated in a parallel manner, using MRjob library. In this preprocessing, we extract the relevant information of the documents to have a standardized representation of them. When a query is given to the system, it is transformed to the same representation and compared in parallel with the documents in the dataset. We will go through each part separately in the subsections.

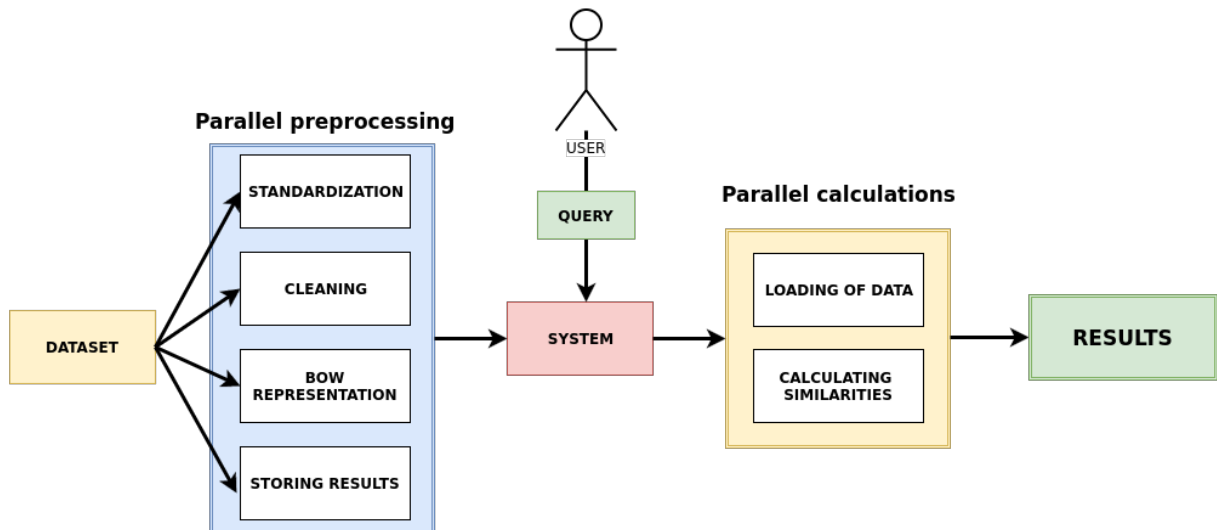


Figure 1: Diagram of the system

### 4.1 Document Representations

Every document is represented as the count of a set of normalized relevant words. This representation could be use for classification and regression techniques directly such as Naive Bayes or Random Forest in order to obtain patterns in documents. Furthermore we will transform each document into a TF-IDF

vector that contains, for every different word in the set, a floating number obtained as the product of two terms:

- $TF_{(t,d)}$ : Term Frequency (BoW). Frequency of term 't' in the document 'd'.
- $IDF_{(t,c)}$ : Inverse Document Frequency of term 't' in corpus 'c',

$$TF_{(t,d)} = \frac{\text{Occurrences of term t in document d}}{\text{Number of total terms in d}}$$

$$IDF_{(t,c)} = \log_2 \left( \frac{\text{Total number of documents in c}}{\text{Number of documents that contain t}} \right)$$

The resulting TF-IDF vector representing each document in the corpus has the form:

$$TF - IDF_{(t,d,c)} = TF_{(t,d)} \times IDF_{(t,c)}$$

The IDF term is a measure of how discriminative a word is in the whole corpus and the TF term is a measure of how representative a word is for a document.

After testing a few queries, we have found that the  $TF_{(t,d)}$  representation alone has a better retrieval of documents so it is the one we use in the final implementation.

The following figure shows the Bag of words for different files in a bar chart. We can appreciate that the most common enough could have enough representative power to characterize the documents to be able to draw similarities between them. Also, the next image, shows the Bag of Words of another document in a picture format to have a bigger overview of the weight of the words.

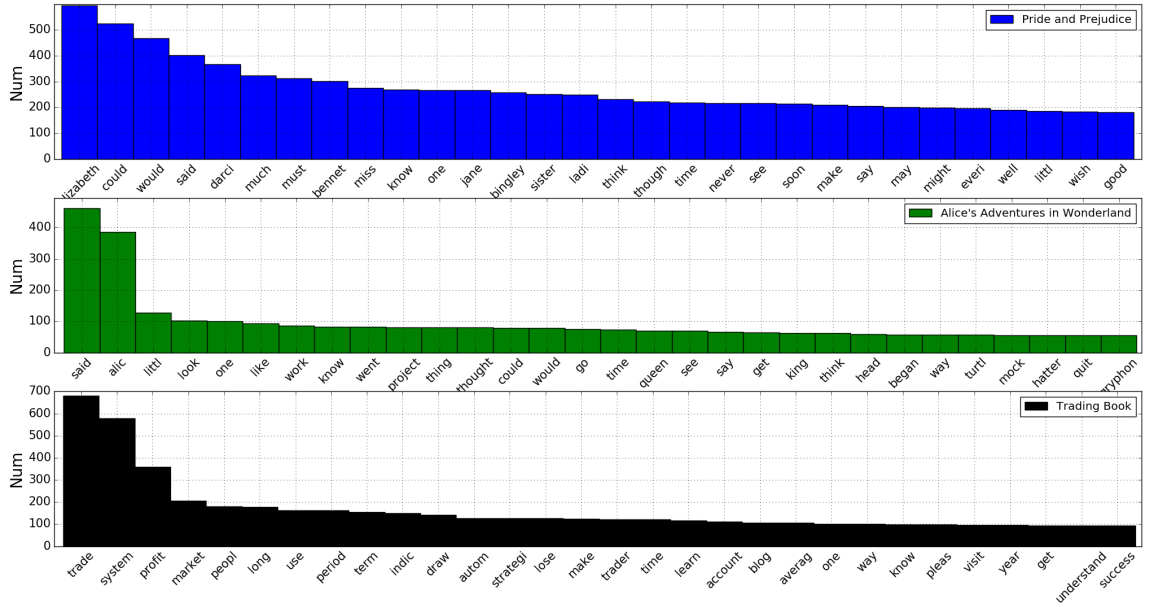


Figure 2: Bag of Words chart representation



the common final sparse representation. Once this is done, we have a common sparse representation of all the documents, we use the sparse matrix library of python in order to make operations faster.

These stages are not so trivial and we use the **pandas** library in order to deal with CSVs and merge the dictionaries in the right way. Once we obtain the common words between the query and the documents of the dataset, we can compute the euclidean distance measure between it and the rest of the documents to obtain the documents. Then we find the top 5 files with the highest similarity measure and check its resemblance to the query. We need to do this ourselves

## 4.4 Running time

Running our system on data-set of this size was challenging. We tried to implemented it as parallel as possible. Time for preprocessing the data was based with linear dependence on number of records loaded. For 1000 records preprocessing is taking around 1 hour. Because in our dataset we have around 50000 books in different formats, creating the BOW representation for whole dataset is slow process. We can increase speed by using a cluster or by upgrading our hardware. This should be reasonable simple because all code is prepared to be run distributively. When the representation is done calling a query can took from 3-15 seconds based on the number of record preprocessed.

## 5 Testing queries

We have performed several tries with different queries in order to tune the parameters and models (similarity functions, words per document...). The performance of the system has been tested by us, giving a query to the system and checking if the retrieved documents are related. Taking into account both precision and recall. Of course the methodology can be improved in a lot of points but we managed to create this full first implementation.

As an example query we use the document *Pride and Prejudice*, the 3 most similar documents retrieved by the system, that are not from the same author are:

- The Emancipated by George Gissing
- Denzil Quarrier by George Gissing
- Demos by George Gissing

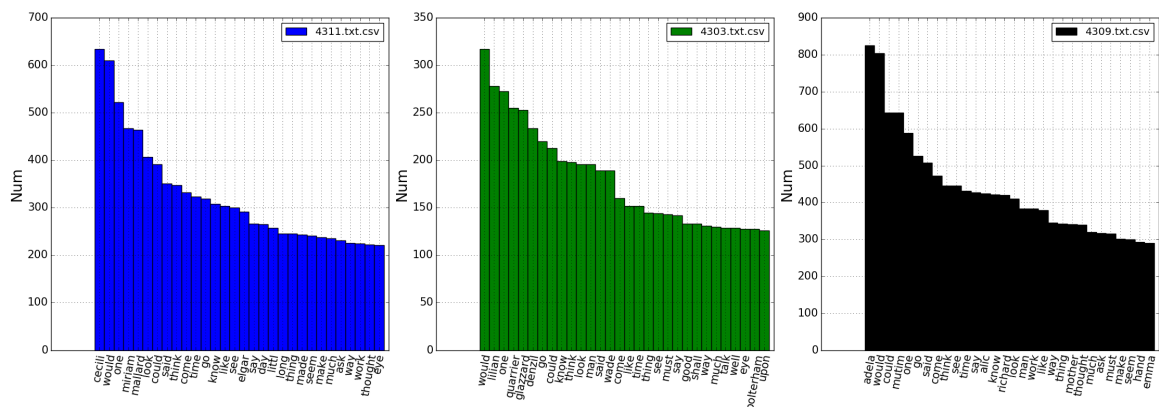


Figure 4: Bag of Words of the retrieved documents for *Pride and Prejudice*

All the retrived books are similar novels to the query one, but apparently the style of writting of the author has a lot of influence in this retrieval system. So it may be used for detecting plagiarism for example.

It seems like in order the improve the system, we should remove the common verbs and words that are non-informative. This is partially done by the IDF component, but in our dataset we found it more misleading than helping.

When queries not related to the dataset were given, such as a book about trading or a Python tutorial, the results of the retrieval system performed badly, because all the retrieved documents have very little to do with the query, but there is little we can do about it. Since the system is very scalable, there would

be no problem, adding a lot of more documents into it. For example, the retrieved documents for the trading book are:

- The Black Experience in America by Norman Coombs
- West Indian Fables by James Anthony Froude Explained by J. J. Thomas
- Political Ideals by Bertrand Russell

The first two documents are actually about slavery, but since the words "work", "black", "white", "time", "power" are very common, then the similarity is high (in trading the candlestick charts are common and they are in black and white, so it is a common word). The third document is actually about economy so it is a reasonable retrieval.

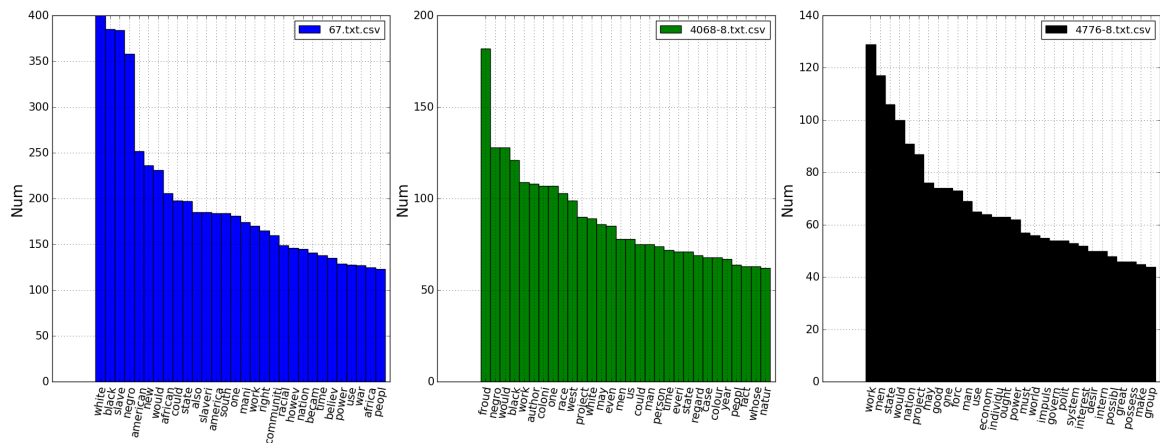


Figure 5: Bag of Words of the retrieved documents for a Trading Book

---

```
get_similar('Pride and Prejudice')
#Results:
['The Emancipated by George Gissing',
'Denzil Quarrier by George Gissing',
'Demos by George Gissing']
#Running time 6.25s
```

---

Listing 4: Pseudofunction queries

## 5.1 Topic detection and downloads prediction

We also tried to develop a system to predict the numbers of downloads and topic from the BOW representation. About the downloads, we tried basic regression algorithms such as linear regression but the results were just random. This could be expected because the number of downloads we have is only for the last month and it is not something that intuitively seem related to the frequency of certain words.

About topic detection we ran into the problem that the documents were categorized in slightly different manner and also, they had multiple tags, so it made it hard for us to model it and create a good system. We leave this feature as a possible future line of development.

## 6 Conclusion

During the development of this project we have seen the special challenges that Big Data presents. We solved the computational limitations by paralleling the processing of the data using the library MRjob and we pretransformed and permanently stored the dataset in order to not have to recompute it every time there is a new query.

We believe that the overall performance of the system is acceptable, the system has many limitations but we were able to tackle the problems and come up with a feasible and time efficient solution. More work needs to be done regarding the representation and similarity measures of the system but this could be the starting point of a high quality information retrieval system.



## References

- [1] “Gutenberg:mirroring\_how-to,” 2013 (accessed November 25, 2016). [Online]. Available: [http://www.gutenberg.org/wiki/Gutenberg:Mirroring\\_How-To](http://www.gutenberg.org/wiki/Gutenberg:Mirroring_How-To)

# Appendices

## Code

### .1 main.py

```
1 import numpy as np
2 import utilfunc as uf
3 import textprocesslib as tplib
4 import highfunc as hf
5 from subprocess import call # Execute commands
6 import pandas as pd
7
8 from scipy.sparse import csr_matrix
9 import pickle_lib as pkl
10
11 from graph_lib import gl
12 #rootFolder = "./GutenbergData"
13 rootFolder = "./GutenbergData/A"
14 outFolder = "./books"
15 csvFolder = "./csvs"
16
17 ## Read the first wrong library
18 read_preproc = 0
19
20 if (read_preproc == 1):
21     # If we need to read and preprocess the raw files
22     # and put then in the output folder
23
24     hf.read_and_preprocess(rootFolder, outFolder)
25
26 compute_BoW = 0
27 if (compute_BoW == 1):
28     ## Read the preprocessed data, preprocess the text and metadata
29     # Creating a bag of word for every file and pickleing it
30     # The files are already in the desired HTML format or any
31     # other that we modify later
32
33     use_MrJob = 1
34     if (use_MrJob == 0):
35         hf.compute_and_save_BoW(outFolder, csvFolder, MaxWords = 1000)
36         ## USE MRJOB INSTEAD !!!
37
38     elif (use_MrJob == 1):
39         allPaths = uf.get_allPaths(outFolder)
40         pd.DataFrame(allPaths).to_csv("./paths.txt")
41
42         call("python ./precompute_DDBB_MrJob.py ./paths.txt", shell=True)
43
44 load_BoWs = 0
45 if (load_BoWs == 1):
46     # Load all the BoWs and join them into a common format ?
47     All_bow, dict_files = hf.load_and_join_BoWs(csvFolder, MaxFiles = 500)
48
49 #####
50 ##### Query Formatting #####
51 #####
52 # Now we can enter a Query, then it will be transformed in the same way as the
53 # corpus and a similarity measure between it and all the documents will be done
54 # The system will output the N most similar elements.
```

```

55 import copy
56 query_f = 1
57 if(query_f == 1):
58
59     query_dirs = ["/books/0/Winchester", "Yorkshire Battles", "1342-0.txt", "11-0.txt"]
60     query_dir = query_dirs[2]
61     # BoW = hf.preprocess_file("https://docs.python.org/2/howto/urllib2.html",
62     #                           typefile = "URL",MaxWords = 1000)
63
64     # BoW = hf.preprocess_file("/free_ebook.pdf", typefile = "pdf",MaxWords = 1000)
65     # BoW = hf.preprocess_file("1342-0.txt", typefile = "text",MaxWords = 1000)
66     BoW = hf.preprocess_file("11-0.txt", typefile = "text",MaxWords = 1000)
67
68     BoW = BoW.set_index(['word']) # Set the index
69     BoW.index = BoW.index.str.encode('utf-8') # Changing
70
71     eu_dist = hf.get_doc_sims(All_bow, BoW)
72
73     best_common,best_common_i = uf.sort_and_get_order(eu_dist[:1].tolist(), reverse = False)
74     # print best_common_i[:5]
75     # print dict_files
76     # plt.close("all")
77     gl.set_subplots(1,3)
78     for i in range(3):
79         BoW_values = All_bow[str(best_common_i[i])].values
80         BoW_index = All_bow.index
81         Nwords = 30
82         b_v, b_vi = uf.sort_and_get_order(BoW_values, reverse = True)
83         gl.bar(BoW_index[b_vi[0:Nwords]].tolist(),
84               BoW_values[b_vi[0:Nwords]],
85               legend = [dict_files[best_common_i[i]].split("/")[-1]],
86               labels = ["", "", "Num"],
87               nf = 1)
88
89     ## Get the BoW of the most common:
90     plotting_thins = 0
91     if(plotting_thins == 1):
92         query_dirs = ["/books/0/Winchester", "Yorkshire Battles"]
93         query_dirs = ["1342-0.txt", "11-0.txt"]
94         labels = [ "Pride and Prejudice", "Alice's Adventures in Wonderland", "Trading Book"]
95
96     gl.set_subplots(3,1)
97     for i in range(3):
98
99
100         # HTML text
101         # BoW = hf.preprocess_file(query_dir, typefile = "text",MaxWords = 1000)
102         # BoW = hf.preprocess_file("https://docs.python.org/2/howto/urllib2.html",
103         #                           typefile = "URL",MaxWords = 1000)
104         if(i < 2):
105             query_dir = query_dirs[i]
106             BoW = hf.preprocess_file(query_dir, typefile = "text",MaxWords = 1000)
107         else:
108             BoW = hf.preprocess_file("/free_ebook.pdf", typefile = "pdf",MaxWords = 1000)
109
110         BoW = BoW.set_index(['word']) # Set the index
111         BoW.index = BoW.index.str.encode('utf-8') # Changing
112
113         # BoW = hf.load_Bow()
114
115         Nwords = 30

```

```

116     BoW = BoW.sort(['num'], ascending=[0])
117     print BoW.index[0:Nwords].shape
118     print BoW["num"].values[1:Nwords]
119
120
121     gl.bar(BoW.index[0:Nwords].tolist(), BoW["num"].values[0:Nwords],
122           legend = [labels[i]],
123           labels = ["", "", "Num"],
124           nf = 1)
125
126     caca = open("./BoW.txt", "w+")
127     for i in range(len(BoW.index.tolist())):
128         for j in range(BoW.iloc[i]["num"]):
129             caca.write(BoW.iloc[i].name + " ")
130     caca.close()
131

```

## .2 precompute/<sub>D</sub>DBB<sub>MrJob</sub>.py

```

1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3  import utilfunc as uf
4  import highfunc as hf
5  from subprocess import call
6  import os
7
8
9  ## This file is meant to be called from the main program
10 # in order to fucking preprocess all the DDBB obtaining their
11 ## BoW and storing it into different files.
12
13 rootFolder = "./books"
14
15 class MRWordVC(MRJob):
16
17     def mapper_VC(self, _, line):
18         # This mapper obtains all the paths of all the files in the DDBB
19         # and then fucking yielding a task for all of them.
20         print os.getcwd()
21         outFolder = "./csvs"
22         # outFolder = "/home/montoya/Desktop/DTU Lec/1st Semester/5. Computational Tools for Big Data"
23
24         filedir = line
25         BoW = hf.preprocess_file(filedir, typefile = "HTML", MaxWords = 1000)
26         file_name = filedir.split("/")[-1]
27         filedir = filedir.split(file_name)[0]

```

## .3 textprocesslib.py

```

1  import unicodedata
2  from nltk.tokenize import word_tokenize, wordpunct_tokenize, sent_tokenize
3  from nltk.stem import SnowballStemmer
4  from nltk.corpus import stopwords
5
6  import numpy as np
7  import pandas as pd
8  ## TEXT PREPROCESSING FUNCTIONS
9
10 def doc_preprocess(document, mode = 0):
11     # Preprocesses a document: Tokenization, lemmatization...
12     if (mode == 1): print document

```

```

13     document = strip_accents(document)
14     if (mode == 1):print document
15     document = doc_tokeniz(document,mode)
16     if (mode == 1):print document
17     document = doc_lowercase(document,mode)
18     if (mode == 1):print document
19     document = doc_rem_punctuation(document,mode)
20     if (mode == 1):print document
21     document = doc_rem_stopwords(document,mode)
22     if (mode == 1):print document
23     document = doc_stem(document,mode)
24     if (mode == 1):print document
25     return document
26
27 def strip_accents(s):
28     # s = s.decode('utf-8').encode('ascii', 'replace')
29     # s = s.encode('utf-8', 'replace')
30
31     return ''.join(c for c in unicodedata.normalize('NFD', s)
32                    if unicodedata.category(c) != 'Mn')
33
34 def doc_tokeniz(document, mode):
35     tokens = word_tokenize(document)
36     return tokens
37
38 def doc_lowercase (document, mode):
39     low_text = [w.lower() for w in document]
40     return low_text
41
42 def doc_rem_stopwords(document, mode):
43     stopwords_en = stopwords.words('english')
44     clean_text = [word for word in document if not word in stopwords_en]
45     return clean_text
46
47 def doc_stem(document, mode):
48     stemmer = SnowballStemmer('english')
49     stemmed_text = [stemmer.stem(word)for word in document]
50     return stemmed_text
51
52 def doc_rem_punctuation(document, mode):
53     clean_text = [w for w in document if w.isalnum()]
54     return clean_text
55
56 #####
57 ##### Bag of Words #####
58
59 from sklearn.feature_extraction.text import CountVectorizer
60 from sklearn.feature_extraction.text import TfidfVectorizer
61 # Initialize the "CountVectorizer" object, which is scikit-learn's
62 # bag of words tool.
63
64
65 def obtain_BoW(document, MaxWords = 1000):
66     # Obtains the BoW and transforms it to a pandas Dataframe
67     # Document is a list of words
68     vectorizer = CountVectorizer(analyzer = "word", \
69                                tokenizer = None, \
70                                preprocessor = None, \
71                                stop_words = None, \
72                                max_features = MaxWords)
73

```

```

74     BoWdoc = vectorizer.fit_transform(document)
75     #     print BoWdoc
76
77     # BoW = Pairs that tell you for each word,
78     # the index file it belongs to
79     # ## Create the dictionary with the words
80     #     vectorizer.fit(document)
81     #     # Transform the file to Sparse form matrix
82     #     BoWdoc = vectorizer.transform(document)
83
84     # Get the vocabulary
85     vocab = vectorizer.get_feature_names()
86     #     print vocab
87     # Transform it to an array form
88     BoWdoc_array = BoWdoc.toarray()
89     # Sum up the counts of each vocabulary word
90     count = np.sum(BoWdoc_array, axis=0)
91     #     print count
92     BoW = dict()
93     BoW["word"] = vocab
94     BoW["num"] = count
95
96     datFr = pd.DataFrame(BoW)
97
98     return datFr
99
100 def obtain_Tfidf(document):
101     # Document is a list of words
102     vectorizer = TfidfVectorizer(analyzer = "word", \
103                                 tokenizer = None, \
104                                 preprocessor = None, \
105                                 stop_words = None, \
106                                 max_features = 50, \
107                                 use_idf= False)
108
109     ## Create the dictionary with the words
110     vectorizer.fit(document)
111
112     # Get the vocabulary
113     vocab = vectorizer.get_feature_names()
114     print vocab
115
116     # Transform the file to Sparse form matrix
117     BoWdoc = vectorizer.transform(document)
118     # Pairs that tell you for each word, the index file it belongs to
119
120     # Transform it to an array form
121     BoWdoc_array = BoWdoc.toarray()
122     # Sum up the counts of each vocabulary word
123     dist = np.sum(BoWdoc_array, axis=0)
124
125     print BoWdoc
126     print dist

```

#### .4 utilfunc.py

```

1 import os
2 import magic
3 import shutil
4 from bs4 import BeautifulSoup # For HTML web treatment.
5 import json

```

```

6 import numpy as np
7 #####
8 ##### Files library #####
9 #####
10
11 def create_folder_if_needed (folder):
12     # This function creates a path if it does not exist
13     if not os.path.exists(folder):
14         os.makedirs(folder)
15
16 def get_allPaths(rootFolder, fullpath = "yes"):
17     ## This function finds all the files in a folder
18     ## and its subfolders
19
20     allPaths = []
21
22     for dirName, subdirList, fileList in os.walk(rootFolder): # FOR EVERY DOCUMENT
23         #         print "dirName"
24         for fname in fileList:
25             # Read the file
26             path = dirName + '/' + fname;
27             if (fullpath == "yes"):
28                 allPaths.append(os.path.abspath(path))
29             else:
30                 allPaths.append(path)
31
32     return allPaths
33
34 def type_file(filedir):
35     mime = magic.Magic()
36     filetype = mime.id_filename(filedir)
37     #     filetype = mime.id_filename(filedir, mime=True)
38
39     # This will be of the kind "image/jpeg" so "type/format"
40     filetype = filetype.split(",")[0]
41     return filetype
42
43 def copy_file(file_source, file_destination, new_name = ""):
44     # Copies a file into a new destination.
45     # If a name is given, it changes its name
46
47     file_name = ""
48     file_path = ""
49
50     file_name = file_source.split("/")[-1]
51     file_path = file_source.split("/")[0]
52
53     if (len(new_name) == 0): # No new name specified
54         file_name = file_source.split("/")[-1]
55     else:
56         file_name = new_name
57
58     create_folder_if_needed(file_destination)
59
60     shutil.copy2(file_source, file_destination + "/" + file_name)
61
62 def loadJsonFromFile(filename):
63     try:
64         with open(filename) as data_file:
65             jsonData = json.load(data_file)
66             # return created json object

```

```

67         return jsonData
68
69     except IOError:
70         print "Error: File does not appear to exist."
71         return 0
72
73     #####
74     ##### HTML treating #####
75     #####
76
77     ### AUTHOR AND TITLE
78     #<div class="dochead">
79     #<h2 class="author">Lawrence Fletcher</h2>
80     #<h2 class="title">"Zero the Slaver"</h2>
81     #<hr/>
82     #</div>
83
84     ### CHAPTER ARE LIKE
85     #<div class="bodytext">
86     #<a href="" name="chap02"></a>
87     #<h3>Chapter Two.</h3>
88     #<h4>A Night of Horror.</h4>
89
90
91     def check_cover(filedir):
92         # This function checks if the file has the words cover
93         # If it does, it is just the data of it,
94         # I couldnt find a proper way to separate them really
95         # They do not follow an easy format
96
97         filename = filedir.split("/")[-1]
98
99         if (filename.find("(cover)") != -1):
100             return 1
101         else:
102             return 0
103
104     def process_HTML_doc(filedir):
105         ## Read the file with HTML and close it
106         fd = open(filedir, 'r')
107         doc_HTML = fd.read()
108         fd.close()
109         # Use BeautifulSoup to process the HTML, get tittle, plain text...
110         soup = BeautifulSoup(doc_HTML) # Transform plain text HTML into soup structure
111
112         ## Estructure of the HTMLS !!
113
114         ## First check that it is a file of a document
115         # <meta content="Zero the Slaver" name="Description"/>
116         # <meta name="description" content="Project Gutenberg Ebooks." />
117         # content_type = soup.find("meta", {"name": "Description"})
118         # if(content_type == None):
119         #     content_type = soup.find("meta", {"name": "description"})
120         #
121         # if(content_type == None):
122         #     return None
123         # content_type = content_type["content"]
124         # if (content_type == "Project Gutenberg Ebooks."):
125         #     return None
126
127         useful_text= ""

```



```

128     ##### FIND ALL CHAPTER #####
129
130     useful = soup.findAll("p") # Directori (Dessert > Bannaan )
131     for elem in useful:
132         useful_text += " " + elem.text
133     #         keywords += " " + elem.text
134
135     #     useful = soup.findAll("div", {"class": "bodytext"}) # Directori (Dessert > Bannaan )
136     #     for elem in useful:
137     #         useful_text += " " + elem.text
138     ##         keywords += " " + elem.text
139
140
141     return useful_text
142
143 def sort_and_get_order (x, reverse = True ):
144     # Sorts x in increasing order and also returns the ordered index
145     x = np.array(x)
146     x = x.flatten() # Just in case we are given a matrix vector.
147     order = range(len(x))
148
149     if (reverse == True):
150         x = -x
151
152     x_ordered, order = zip(*sorted(zip(x, order)))
153
154     if (reverse == True):
155         x_ordered = -np.array(x_ordered)
156
157     return np.array(x_ordered), np.array(order)
158

```

## .5 highfunc.py

```

1  import os
2  import magic
3  import shutil
4  from bs4 import BeautifulSoup # For HTML web treatment.
5  import json
6  import utilfunc as uf
7  import textprocesslib as tplib
8  import pickle_lib as pkl
9  import pandas as pd
10 import urllib2
11 from pyPdf import PdfFileWriter, PdfFileReader
12 from textract import process
13 def getURLContent(path = 'http://python.org/'):
14     response = urllib2.urlopen(path)
15     html = response.read()
16     soup = BeautifulSoup(html, 'html.parser')
17     return soup.get_text()
18
19 def getPDFContent(path = "../.../P3.pdf"):
20     # content = ""
21     # Load PDF into pyPDF
22     # input_file = file(path, "rb")
23     # pdf = PdfFileReader(input_file)
24     # Iterate pages
25     # for i in range(0, pdf.getNumPages()):
26     #     # Extract text from page and add to content
27     #     content += pdf.getPage(i).extractText() + "\n"

```

```

28 #         print pdf.getPage(i).extractText()
29 #     # Collapse whitespace
30 #     content = " ".join(content.replace(u"\xa0", " ").strip().split())
31 #     input_file.close()
32 ##     print content
33     content = process(path)
34     return content
35
36 def read_and_preprocess(rootFolder, outFolder, MaxFiles = 1000):
37     # This function reads the documents and preprocess them
38     # It only selects the valid files and extracts the
39     # important metainformation and text.
40     # It also renames them and puts them in other folders.
41
42     # Max    MaxFiles files per folder
43     i = 0
44     aux_f = 0
45
46     ## Get all the paths of the folder
47     allPaths = uf.get_allPaths(rootFolder)
48
49     # For every path TODO in map-reduce
50     for filedir in allPaths:
51         filetype = uf.type_file(filedir)
52
53         print filedir
54         print filetype
55         # For all files, we check the extension
56         if (filetype == "HTML document"):
57             # Copy the file
58             coverf = uf.check_cover(filedir)
59             if (coverf == 1):
60                 # Copy the file in the new destination
61                 uf.copy_file(filedir, outFolder + "_cover" + "/" + str(aux_f))
62             else:
63                 uf.copy_file(filedir, outFolder + "/" + str(aux_f))
64
65             i = i + 1
66
67             if (i >= MaxFiles):
68                 i = 0
69                 aux_f += 1
70
71 def preprocess_file(filedir, typefile = "text", MaxWords = 1000):
72     # This function will preprocess one query, obtaining its BoW
73     # text:      If the file is a text file in the system.
74     # HTML:      If the files an HTML document in your computer
75     # URL:       If the file is an URL
76     # pdf:       If the file is a pdf
77
78     ##### READ THE PROPER TEXT OF THE DOCUMENT #####
79     if (typefile == "text"):
80         caca = open(filedir)
81         document = caca.read()
82         document = document.decode('utf-8', errors='replace').encode('utf-8')
83         document = unicode(document, "utf-8")
84         caca.close()
85
86     elif (typefile == "HTML"):
87         document = uf.process_HTML_doc(filedir)
88         print filedir

```

```

89         print len(document)
90         ## If the HTML file is not a book
91         if (type(document) == type(None)):
92             return pd.DataFrame([]) # Exit this loop execution
93         if(len(document) == 0):
94             return pd.DataFrame([])
95
96     elif (typefile == "URL"):
97         document = getURLContent(filedir)
98     #     document = unicode(document, "utf-8")
99
100    elif (typefile == "pdf"):
101        document = getPDFContent(filedir)
102        document = unicode(document, "utf-8")
103    #     print document
104    ##### Preprocess document and obtain BoW #####
105    processed = tplib.doc_preprocess(document)
106    BoW = tplib.obtain_BoW(processed, MaxWords)
107    return BoW
108
109 def compute_and_save_BoW(rootFolder, outFolder, MaxFiles = 1000, MaxWords = 1000):
110     ## Get all the paths of the folder
111     allPaths = uf.get_allPaths(rootFolder)
112     # For every path TODO in map-reduce
113     for filedir in allPaths:
114         BoW = preprocess_file(filedir, typefile = "text", MaxWords = 1000)
115     #     print ' '.join(processed)
116     #     BoW = tplib.obtain_BoW([' '.join(processed[1:1000]),
117     #         ' '.join(processed[1000:2000])], MaxWords)
118     #     BoW = tplib.obtain_Tfidf([processed[:1000],processed[1000:2000]])
119     ## Now we pickle the processed data into the machine disk
120
121     file_name = filedir.split("/")[-1]
122     uf.create_folder_if_needed(outFolder)
123     #     print BoW
124     #     pickle.store_pickle(outFolder + "/" + file_name + ".pkl",BoW,1)
125     BoW.to_csv(outFolder + "/" + file_name + ".csv",
126         encoding = "utf-8",
127         index = False) # Do not write the index number
128
129 def load_BoW(filedir):
130     BoW = pd.DataFrame.from_csv(filedir, sep=',',
131         index_col=1)
132     return BoW
133
134 def load_and_join_BoWs(rootFolder, MaxFiles = 100):
135     ## Get all the paths of the folder
136     allPaths = uf.get_allPaths(rootFolder)
137
138     # For every path TODO in map-reduce
139     Total_BoW = pd.DataFrame()
140
141     index_doc = 0
142
143     # Dictionary between the file and its number
144     dict_files = {}
145     for filedir in allPaths:
146         # Load the BoW
147         BoW = pd.DataFrame.from_csv(filedir, sep=',',
148             index_col=1)
149         # filedir.split("/")[-1]

```

```

150     BoW = BoW.rename(columns = {'num':str(index_doc)})
151     dict_files[index_doc] = filedir
152
153     # Concatenate them !
154     Total_BoW = pd.concat([Total_BoW, BoW], axis=1)
155     # print BoW
156     # print Total_BoW
157     index_doc += 1
158     # print dict_files
159     if (index_doc >= MaxFiles): # If we are done loading
160         # Fill the NaNs with 0
161         Total_BoW = Total_BoW.fillna(0)
162     # print Total_BoW.shape
163     return Total_BoW, dict_files # Return the total BoW
164     Total_BoW = Total_BoW.fillna(0)
165
166     return Total_BoW, dict_files
167
168 from scipy.sparse import csr_matrix
169 import numpy as np
170 def get_doc_sims(All_bow, BoW):
171     # This function gives back the similarity measures of the documents,
172     # compared to the database.
173     # print BoW
174     joined_inner = pd.concat([All_bow, BoW], axis=1, join='inner')
175     # print joined_inner.shape
176     A = csr_matrix(joined_inner)
177
178     v = np.array(joined_inner["num"].values)
179     # print v.shape
180     v = v/float(v.sum())
181     # print v.shape
182     A = A.T
183
184     # print A.sum(axis = 1).shape
185     A = A/A.sum(axis = 1) # Get the frequencies
186
187     A = A - v
188     A = np.multiply(A,A)
189     # print v.shape
190     eu_dist = A.dot(np.ones((v.size,1)))
191
192     return eu_dist

```

## .6 pickle\_lib.py

```

1  import pickle
2  import gc
3  import os
4
5  # Library for loading and storing big amounts of data into
6  # different files because pickle takes a lot of RAM otherwise.
7  # If the number of partitions = 1, then it just loads like a
8  # regular pickle file.
9  # It uses gc also to remove garbage variables.
10
11 def store_pickle (filename, li, partitions = 1, verbose = 1):
12     gc.collect()
13     splitted = filename.split(".")
14     if (len(splitted) == 1): # If there was no extension
15         fname = filename

```

```

16         fext = ""
17     else:
18         fname = '.'.join(splitted[:-1])    # Name of the file
19         fext = "." + splitted[-1]         # Extension of the file
20
21     # li: List of variables to save.
22     # It saves the variables of the list in "partitions" files.
23     # This function stores the list li into a number of files equal to "partitions" in pickle form
24     # If "partitions" = 1 then it is a regular load and store
25     num = int(len(li)/partitions);
26
27     if (partitions == 1): # Only 1 partition
28         if (verbose == 1):
29             print "Creating file: " + fname + fext
30             with open(fname + fext, 'wb') as f:
31                 pickle.dump(li, f)
32     else: # Several partitions
33         for i in range(partitions - 1):
34             if (verbose == 1):
35                 print "Creating file: " + fname + str(i) + fext
36                 with open(fname + str(i)+ fext, 'wb') as f:
37                     pickle.dump(li[i*num:(i+1)*num], f)
38                     # We dump only a subset of the list
39             # Last partition to create
40             if (verbose == 1):
41                 print "Creating file: " + fname + str(partitions - 1) + fext
42
43             with open(filename + str(partitions - 1)+ fext, 'wb') as f:
44                 pickle.dump(li[num*(partitions - 1):], f)
45                 # We dump the last subset.
46     gc.collect()
47
48 def load_pickle (filename, partitions = 1, verbose = 0):
49
50     gc.collect()
51     total_list = []
52     splitted = filename.split(".")
53     if (len(splitted) == 1): # If there was no extension
54         fname = filename
55         fext = ""
56     else:
57         fname = '.'.join(splitted[:-1])    # Name of the file
58         fext = "." + splitted[-1]         # Extension of the file
59
60     if (partitions == 1): # Only 1 partition
61         if (verbose == 1):
62             print "Loading file: " + fname + fext
63
64         if (os.path.exists(fname + fext) == True): # Check if file exists !!
65
66             with open(fname + fext, 'rb') as f:
67                 total_list = pickle.load(f) # We read the pickle file
68         else:
69             print "File does not exist: " + fname + fext
70             return []
71     else: # Several partitions
72         for i in range(partitions):
73             if (verbose == 1):
74                 print "Loading file: " + fname + str(i)+ fext
75
76             if (os.path.exists(fname + str(i)+ fext) == True): # Check if file exists !!

```

```

77
78         with open(fname + str(i)+ fext, 'rb') as f:
79             part = pickle.load(f)      # We read the pickle file
80             total_list.extend(part)
81
82     else:
83         print "File does not exist: " + fname + str(i)+ fext
84         return []
85
86     gc.collect()
87     return total_list
88
89     #n = 3
90     #lista = [10, 23, 43, 65, 34, 98, 90, 84, 98]
91     #store_pickle("lista",lista,n)
92     #lista2 = load_pickle("lista",n)

```

## .7 simlib.py

```

1  ##### SIMILARITY LIBRARY #####
2
3  import numpy as np
4  from utilfunc import *
5  from textprocesslib import *
6  # Define similarity values between the query and the corpus.
7  def get_similarities(query,database):
8      query_words = doc_preprocess(query,1) # Preprocess query
9      query_words_tfidf = tf_idf_doc(query_words, database) # Get tfidf of query
10
11     for i in range (len(query_words_tfidf[0])):
12         print str(query_words_tfidf[1][i]) + "\t\t" + query_words_tfidf[0][i]
13     print " "
14     N_docs_database = len(database)
15
16     eu_dists = np.zeros([N_docs_database,1])
17     cos_sims = np.zeros([N_docs_database,1])
18     n_commun = np.zeros([N_docs_database,1])
19     tfidf_sum = np.zeros([N_docs_database,1])
20
21     for i in range (N_docs_database):
22         # Common words
23         (vq ,vd) = get_common_tfidf_vectors(query_words_tfidf,database[i])
24         # Non-common words
25         (vnq ,vnd) = get_noncommon_tfidf_vectors(query_words_tfidf,database[i])
26
27         # If no words in commun
28         if(vq.size == 0):
29             eu_dists[i] = -1
30             cos_sims[i] = -1
31             n_commun[i] = vq.size
32         else:
33             eu_dists[i] = get_euc_dist(vq,vd,vnq,vnd)
34             cos_sims[i] = get_cos_sim(vq,vd)
35             n_commun[i] = vq.size
36             tfidf_sum[i] = np.sum(vd)
37
38     similarities = (eu_dists,cos_sims, n_commun, tfidf_sum)
39     # Similatiries obtained are:
40     # - The euclidean distance
41     # - The cosine similarity
42     # - The number of words in commun.

```

```

43     #         (Important coz, they can have very high values in the eu y cos
44     #         but just because they have very few words in commun)
45     # - tfidf_sum: Sum of the tdidf of the values
46
47     return similarities
48
49 def rank_documents(Similarity, N_top, type_simi):
50     # RANK WITH EVERY INDEPENDENT SIMILARITY MEASURE.
51     n_doc = len(Similarity)
52     # 1) Euclidean distance: The closer, the better, so the smaller the value, the better
53     if (type_simi == "euclidean"):
54         # For this similarity, the vectors that did not have words in common have
55         # similarity -1, so we cannot order them directly like this.
56         for i in range(len(Similarity)):
57             if (Similarity[i] == -1):
58                 Similarity[i] = 10000000000;
59         indexes = np.array(range(0,n_doc))
60         together = zip(Similarity, indexes) # Zip both arrays together
61         sorted_together = sorted(together) # Sorted Increasing order
62         ordered_indexes = [x[1] for x in sorted_together]
63         # Since ordered_indexes contains the indexes of the best documents
64         # (obtained from decreasing similarity)
65
66     # 2) Cosine distance: The bigger the value, the better
67     if (type_simi == "cosine"):
68         # Wrong vectors have similarity -1 so no problem, the bigger the better.
69         indexes = np.array(range(0,n_doc))
70         together = zip(Similarity, indexes) # Zip both arrays together
71         sorted_together = sorted(together, reverse=True) # Sorted Decreasing order
72         ordered_indexes = [x[1] for x in sorted_together]
73         #print ordered_indexes
74
75     # 3) Use the number of words in common:
76     if (type_simi == "common"):
77         indexes = np.array(range(0,n_doc))
78         together = zip(Similarity, indexes) # Zip both arrays together
79         sorted_together = sorted(together, reverse=True) # Sorted Decreasing order
80         ordered_indexes = [x[1] for x in sorted_together]
81
82     # ) Use the sum of tfidf values:
83     if (type_simi == "tdidf_sum"):
84         indexes = np.array(range(0,n_doc))
85         together = zip(Similarity, indexes) # Zip both arrays together
86         sorted_together = sorted(together, reverse=True) # Sorted Decreasing order
87         ordered_indexes = [x[1] for x in sorted_together]
88
89     top_indexes = ordered_indexes[0:N_top]
90     return top_indexes
91
92 def get_combined_rank(Ranks, N_top, type_comb):
93     # Ok... we have several similarity measures, now we have to output the ranking
94     # based on those similarities. How do we combine them to get the best ranking ?
95     n_rank, n_doc = np.shape(Ranks)
96
97     eu_w = 0.5; # Weight of the euclidean distance importance
98     cos_w = 0.5; # Weight of the cosine distance importance
99     n_w = 1;    # Weight of the number of common words importance
100
101     total_similarity = np.zeros([n_doc,1])
102
103     for i in range(n_doc):

```

```

104         total_similarity [Ranks[i]] += eu_w * Ranks[0][i]
105         total_similarity [Ranks[i]] += cos_w * Ranks[1][i]
106         total_similarity [Ranks[i]] -= n_w * Ranks[2][i]
107
108         # Get the top indexes rank from the combined similarity methods:
109         indexes = np.array(range(0,n_doc))
110         together = zip(total_similarity, indexes) # Zip both arrays together
111         sorted_together = sorted(together) # Sorted Increasing order
112         ordered_indexes = [x[1] for x in sorted_together]
113         top_indexes = ordered_indexes[0:N_top]
114         return top_indexes
115
116     def get_euc_dist(vq,vd,vnq,vnd):
117         distance = np.sqrt(((vq - vd)*(vq - vd)).sum(axis=0))
118
119
120         return distance
121
122     def get_cos_sim(vq,vd):
123         modules = (np.sqrt((vq*vq).sum(axis=0)) * np.sqrt((vd*vd).sum(axis=0)))
124
125         if (modules == 0):
126             print "Error, cosine similarity. One vector is all 0s."
127             return -1
128         cos_sim = np.dot(vq.transpose(),vd)
129         cos_sim = cos_sim / modules
130
131         return cos_sim

```