

# NWHoogle

## Workshop 2 Project

Jay Ning<sup>1</sup>   Wendy Wang<sup>2</sup>   Beatrice Hou<sup>3</sup>

UIC Data Science

2022

# Table of Contents

1. Data acquisition and processing
2. Indexing and Ranking
3. Server and UI
4. Demo

# Data acquisition and processing

## Crawl data

### Target parsing

<https://www.gutenberg.org/cache/epub/10000/pg10000.txt>

The Project Gutenberg eBook of The Magna Carta

This eBook is for the use of anyone anywhere in the United States and most other parts of the world at no cost and with almost no restrictions whatsoever. You may copy it, give it away or re-use it under the terms of the Project Gutenberg License included with this eBook or online at [www.gutenberg.org](http://www.gutenberg.org). If you are not located in the United States, you will have to check the laws of the country where you are located before using this eBook.

Title: The Magna Carta

Author: Anonymous

Release Date: October 15, 2003 [eBook #10000]  
[Most recently updated: December 20, 2021]

Language: English

Produced by: Michael Hart

\*\*\* START OF THE PROJECT GUTENBERG EBOOK THE MAGNA CARTA \*\*\*

The Magna Carta

Contents

The Text of Magna Carta  
Magna Carta 1215  
The text of THE MAGNA CARTA

A note from Michael Hart, preparer of the 0.1 version.

This file contains a number of versions of the Magna Carta, some of which were a little mangled in transit. I am sure our volunteers will find and correct errors I didn't catch, and that version 0.2 - 1.0 will have significant improvements, as well as at least one more version in Latin.

---

# Data acquisition and processing

## Crawl data

Use requests and define function to get book files

```
1 def storageToLocalFiles(storagePath, data):#store in the path
2     fhandle = open(storagePath,"w")
3     fhandle.write(data)
4     fhandle.close()

1 import requests
2 headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/67.0.3396.87 Safari/537.36'}
3 for i in range(10044,11001):#It should larger than 200MB
4     storagePath = 'books/'+str(i)+".txt"
5     link='https://www.gutenberg.org/cache/epub/'+str(i)+'pg'+str(i)+'.txt'
6     r=requests.get(link,headers=headers)
7     storageToLocalFiles(storagePath,r.content)
```

# Data acquisition and processing

## Data Cleaning

### Observed Data

Found that some books do not have TXT format, the content is HTML. It is characterized as beginning with `<!DOCTYPE html>`

```
1 <!DOCTYPE html>
2 <html class="client-nojs" lang="en" dir="ltr">
3 <head>
4   <meta charset="UTF-8"/>
5
6   <title>404 | Project Gutenberg</title>
7   <link rel="stylesheet" href="/gutenberg/style.css?v=1.1">
8   <link rel="stylesheet" href="/gutenberg/collapsible.css?v=1.1">
9   <link rel="stylesheet" href="/gutenberg/new_nav.css?v=1.321231">
10  <link rel="stylesheet" href="/gutenberg/pg-desktop-one.css">
11  <meta name="viewport" content="width=device-width, initial-scale=1">
12  <meta name="keywords" content="books, ebooks, free, kindle, android, iphone, ipad"/>
13  <meta name="google-site-verification" content="wucOEvsNj5kP3Ts_36OfP64laakK-1mVTg-ptrcG9io"/>
14  <meta name="alexaVerifyID" content="4wNaCljsE-A82vP_ih2H_UqXZVM"/>
15  <link rel="copyright" href="https://www.gnu.org/copyleft/fdl.html"/>
16  <link rel="shortcut icon" href="/gutenberg/favicon.ico?v=1.1"/>
17
18  <meta property="og:title" content="Project Gutenberg" />
19  <meta property="og:type" content="website" />
20  <meta property="og:url" content="https://www.gutenberg.org/" />
21  <meta property="og:description" content="Project Gutenberg is a library of free eBooks." />
22  <meta property="fb:admins" content="615269807" />
23  <meta property="fb:app_id" content="115319388529183" />
24  <meta property="og:site_name" content="Project Gutenberg" />
25  <meta property="og:image" content="https://www.gutenberg.org/gutenberg/pg-logo-144x144.png" />
26 </head>
27 <body>
28   <div class="container"><!-- start body --><nav>
29   <!--<div id="main_logo"> -->
30   <a id="main logo" href="/" class="no-hover">
```

# Data acquisition and processing

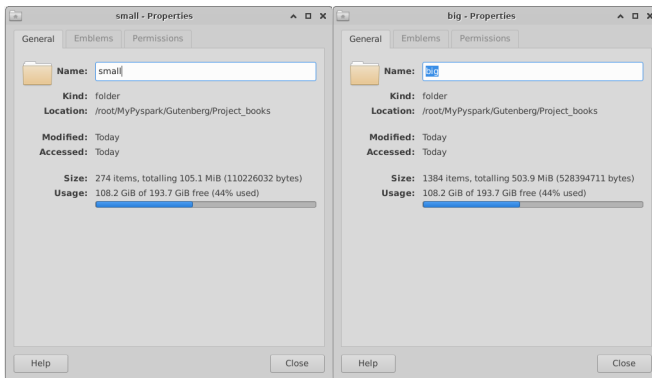
## Data Cleaning

Use `readline()` function to check first line. And use `os` module to delete wrong file

```
1 import os
2 path1 = "books"
3 files= os.listdir(path1)
4 d_list=[]
5 for file in files:
6     f = open("books/"+file,"r",encoding='utf-8')
7     if(f.readline()=="<!DOCTYPE html>\n"):
8         d_list.append(file)
9     # f.readline()
10 for file in d_list:
11     os.remove("books/"+file)
```

# Data acquisition and processing

## Result





# Indexing and Ranking

## Upload Files

### Upload Files

Firstly we divided the books obtained by the crawler into two parts: 100 megabytes and 500 megabytes, and upload them to 2 HDFS folders respectively (project small and project big) then check the file in the Hadoop browse directory.

<input type="checkbox"/>	<a href="#">drwxr-xr-x</a>	<a href="#">root</a>	<a href="#">supergroup</a>	0 B	May 22 11:11	<a href="#">0</a>	0 B	<a href="#">project_big</a>	
<input type="checkbox"/>	<a href="#">drwxr-xr-x</a>	<a href="#">root</a>	<a href="#">supergroup</a>	0 B	May 21 19:06	<a href="#">0</a>	0 B	<a href="#">project_small</a>	



# Indexing and Ranking

## Preprocessing

In the preprocessing part, we delete all punctuation, numbers and other symbols, only retain English characters. In addition, we quoted a text processing toolkit called nltk to delete the common stopwords, and restore the part of speech of English words with deformation.

```
import re
import math
def removeSomeTexts(text):
    return re.sub(r"[^a-z-'\s]",'', text)
stop_word = stopwords.words("english")
new_words=["aa","aaa","aaaahh"]
for i in new_words:
    stop_word.append(i)

def lemmatization(text):
    lemmatizer = WordNetLemmatizer()
    return lemmatizer.lemmatize(text)
def stem(text):
    lancaster_stemmer = LancasterStemmer()
    return lancaster_stemmer.stem(text)
```

# Indexing and Ranking

## Index

### Use the local session

```
1 from pyspark import SparkConf, SparkContext
2 import os, re
3 from pyspark.sql.session import SparkSession
4 from nltk.corpus import stopwords
5 from nltk.stem import WordNetLemmatizer
6 from nltk.stem.porter import PorterStemmer
7 from nltk.stem.lancaster import LancasterStemmer
8 import os, re
9 import time
10
11 spark = SparkSession.builder.master('local').appName('timetest').getOrCreate()
12
13 sc = spark.sparkContext
14 #sc = SparkContext.getOrCreate(SparkConf())
15 sc._conf.set('spark.executor.memory', '8g')
16 sc._conf.set('spark.driver.memory', '8g')
17 sc._conf.set('spark.driver.maxResultsSize', '0')
```

# Indexing and Ranking

## Index

Use the local session to test the time when doing the indexing part with small document (105M)

```
1 #Use local session to test small document
2 start = time.time()
3 data = sc.wholeTextFiles("hdfs://10.20.0.151:9000/project_small")
4 numfiles = data.count()
5 wordcount = data.flatMap(lambda x: [(os.path.basename(x[0]), removeSomeTexts(i.lower()), 1) for i in re.split('\W+', x[1])]) \
6     .filter(lambda x: x[0][1] not in stop_word).filter(lambda x: x[0][1] != '').map(lambda x: ((x[0][0], lemmatization(x[0][1]), x[1])) \
7     .reduceByKey(lambda a, b: a+b).map(lambda x: (x[0][1], (x[0][0], x[1]))) \
8     .reduceByKey(lambda a, b: a+b).sortByKey()
9 #wordcount.take(10)
10 wordcount.saveAsTextFile('hdfs://10.20.0.151:9000/output/InvIndex_small')
11 end = time.time()
12 print(end-start, " seconds")
13
```

[Stage 3:]> (0 + 1) / 1]

168.81149792671204 seconds

168.81149792671204 seconds

# Indexing and Ranking

## Index

Use the local session to test the time when doing the indexing part with big document (503M)

```
1 #Use local session to test big document
2 start = time.time()
3 data = sc.wholeTextFiles('hdfs://10.20.0.151:9000/project_big')
4 numfiles = data.count()
5 wordcount = data.flatMap(lambda x: [((os.path.basename(x[0]), removeSomeTexts(i.lower())), 1) for i in re.split('\W+', x[1]))] \
6     .filter(lambda x: x[0][1] not in stop_word).filter(lambda x: x[0][1] != '').map(lambda x: ((x[0][0], lemmatization(x[0][1])), x[1])) \
7     .reduceByKey(lambda a, b: a+b).map(lambda x: (x[0][1], (x[0][0], x[1]))) \
8     .reduceByKey(lambda a, b: a+b).sortByKey()
9 #wordcount.take(10)
10 wordcount.saveAsTextFile('hdfs://10.20.0.151:9000/output/InvIndex_big')
11 end = time.time()
12 print(end-start, " seconds")

[Stage 7:>                                     (0 + 1) / 1]

814.2438313961029 seconds
```

814.2438313961029 seconds

# Indexing and Ranking

## Index

### Use the cluster session

```
from pyspark import SparkConf, SparkContext
import os, re
import nltk
from pyspark.sql.session import SparkSession
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.stem.porter import PorterStemmer
from nltk.stem.lancaster import LancasterStemmer
import os, re
import time

spark1 = SparkSession.builder.master('spark://DPWII-JayNing:7077').appName('MyInvIndex').getOrCreate()
sc1 = spark1.sparkContext
sc1._conf.set('spark.executor.memory', '8g')
sc1._conf.set('spark.driver.memory', '8g')
sc1._conf.set('spark.driver.maxResultSize', '0')
```

# Indexing and Ranking

## Index

Use the cluster session to test the time when doing the indexing part with small document(105M)

```
1 #Use cluster session to test small document
2 start = time.time()
3 data = sc1.wholeTextFiles('hdfs://10.20.0.151:9000/project_small')
4 numFiles = data.count()
5 wordcount = data.flatMap(lambda x: [(os.path.basename(x[0]), removeSomeTexts(i.lower())), 1] for i in re.split('\W+', x[1])) \
6     .filter(lambda x: x[0][1] not in stop_word).filter(lambda x: x[0][1] != '') \
7     .map(lambda x: ((x[0][0], lemmatization(x[0][1])), x[1])) \
8     .reduceByKey(lambda a, b: a+b).map(lambda x: (x[0][1], (x[0][0], x[1]))) \
9     .reduceByKey(lambda a, b: a+b).sortByKey()
10 #wordcount.take(10)
11 wordcount.saveAsTextFile('hdfs://10.20.0.151:9000/output/InvIndex_small_cluster')
12 end = time.time()
13 print(end-start, " seconds")

80.18091893196106 seconds
```

80.18091893196106 seconds

# Indexing and Ranking

## Index

Use the cluster session to test the time when doing the indexing part with big document(503M)

```
1 #Use local session to test big document
2 start = time.time()
3 data = sc1.wholeTextFiles('hdfs://10.20.0.151:9000/project_big')
4 numFiles = data.count()
5 wordcount = data.flatMap(lambda x: [(os.path.basename(x[0]), removeSomeTexts(i.lower()), 1) for i in re.split('\W+', x[1])]) \
6     .filter(lambda x: x[0][1] not in stop_word).filter(lambda x: x[0][1] != '').map(lambda x: ((x[0][0], lemmatization(x[0][1])), x[1])) \
7     .reduceByKey(lambda a, b: a+b).map(lambda x: (x[0][1], (x[0][0], x[1]))) \
8     .reduceByKey(lambda a, b: a+b).sortByKey()
9 #wordcount.take(10)
10 wordcount.saveAsTextFile('hdfs://10.20.0.151:9000/output/InvIndex_big_cluster1')
11 end = time.time()
12 print(end-start, " seconds")

[Stage 51:=====> (1 + 1) / 2]
439.3812403678894 seconds
```

439.3812403678894 seconds

# Indexing and Ranking

## Index

Summary:

Session/FileSize	small	big
local	168 seconds	814 seconds
cluster	80 seconds	439 seconds



# Indexing and Ranking

## Index

In addition to the file index, we also set another index for the line of each word in each file

```
#add line index
data = sc.wholeTextFiles('hdfs://10.20.0.151:9000/project_small')
numFiles = data.count()
wordcount = data.flatMap(lambda x: [(os.path.basename(x[0]), i.lower()) for i in re.split('\n', x[1])])
wordcount1 = wordcount.zipWithIndex()
wordcount3 = wordcount1.flatMap(lambda x: [(x[0][0], x[1], removeSomeTexts(i.lower())) for i in re.split('\W+', x[0][1])]) \
    .filter(lambda x: x[1] != '') \
    .filter(lambda x: x[1] not in stop_word) \
    .map(lambda x: (x[0][0], lemmatization(x[1]), x[0][1])) \
    .map(lambda x: (x[0][0], x[0][1], x[1]))
wordcount3.take(15)
```

# Indexing and Ranking

## Rank

We first calculate the RDD which have the document id along with its document term frequency for each token

```
1 tf1 = wordcount.reduceByKey(lambda a,b:a+b).map(lambda x:(x[0][1],(x[0][0],x[1])))
2 #tf2 = tf1.map(lambda x:(x[0],(x[1][0],1+math.log10(x[1][1]))))
3 tf = tf1.map(lambda x:(x[0],(x[1][0], 1+math.log10(x[1][1]) )))
4 tf.take(10)
```

Another one is the RDD which has the IDF score

```
1 idf1 = tf1.map(lambda x:(x[0],1)).reduceByKey(lambda x,y:x+y)
2 idf = idf1.map(lambda x:(x[0],math.log10(numFiles / (x[1]))))
3 idf.take(10)
```

# Indexing and Ranking

## Rank

Perform an inner join to assign each token

```
1 tfidf1 = tf.join(idf)
2 tfidf = tfidf1.map(lambda x: (x[1][0][0], (x[0], x[1][0][1] * x[1][1]))).sortByKey()
3 index = tfidf.map(lambda x: (x[1][0], (x[0], x[1][1]))).filter(lambda x: x[0]!='')
4 index = index.sortByKey()
5 index.take(10)
```



Save lines and tfidf to MySQL

```
1 user = 'root'
2 pw = '12345678'
3 ## Database information
4 table_name = 'books'
5 url = 'jdbc:mysql://localhost:3306/project?user='+user+'&password='+pw
6 properties = { 'password': pw, 'user': user}
7 lines.write.jdbc(url=url, table=table_name, mode='append', properties=properties)
8
```

# Indexing and Ranking

## Table Join

Check if the tables with row indexes and TFIDF indexes have been imported into SQL workbench respectively.

1 • <b>SELECT * FROM</b> project.line_index;			1 • <b>SELECT * FROM</b> project.books		
t Grid  Filter Rows: <input type="text"/> Export			t Grid  Filter Rows: <input type="text"/>		
term	fileName	rowNum	term	fileName	TFIDF
project	10000.txt	0	abashed	10008.txt	0.9705248794454127
gutenberg	10000.txt	0	abashed	10017.txt	0.7459665670122424
ebook	10000.txt	0	abashed	10038.txt	0.7459665670122424
magna	10000.txt	0	abashed	10041.txt	0.7459665670122424
carta	10000.txt	0	abashed	10047.txt	0.7459665670122424
ebook	10000.txt	2	abashed	10048.txt	0.7459665670122424
use	10000.txt	2	abashed	10057.txt	0.7459665670122424
anyone	10000.txt	2	abashed	10062.txt	1.195083191878583
anywhere	10000.txt	2	abashed	10064.txt	1.4919331340244848
united	10000.txt	2	abashed	10068.txt	0.7459665670122424
state	10000.txt	2	abashed	10083.txt	0.9705248794454127
part	10000.txt	3	abashed	10094.txt	0.7459665670122424
world	10000.txt	3	abashed	10095.txt	1.195083191878583

# Indexing and Ranking

## Table Join

To speed up the query for corresponding terms, we use natural joins to automatically match the records of columns with the same name in both tables and merge the two tables.

```
CREATE TABLE tfitfidfd AS SELECT * FROM project.books NATURAL JOIN project.line_index;
```







Filter Rows: <input type="text"/>				
Edit:    Export/Import:   Wrap Cell Content: <input type="checkbox"/> Fetch rows:				
term	fileName	TFIDF	rowNum	
abbot	10000.txt	2.094431030211474	59	
abbot	10000.txt	2.094431030211474	180	
abbey	10000.txt	1.564524419258904	344	
abbot	10000.txt	2.094431030211474	346	
abolished	10000.txt	0.993710655438553	356	
abandon	10000.txt	0.538086675634899	385	
abbey	10000.txt	1.564524419258904	392	
abandon	10000.txt	0.538086675634899	394	
abandon	10000.txt	0.538086675634899	426	
abbot	10000.txt	2.094431030211474	540	
able	10000.txt	0.0822909851809811	630	
abbot	10000.txt	2.094431030211474	665	
abbey	10000.txt	1.564524419258904	835	
abolished	10000.txt	0.993710655438553	848	
abbey	10000.txt	1.564524419258904	886	
abbot	10000.txt	2.094431030211474	1034	
able	10000.txt	0.0822909851809811	1123	

# Indexing and Ranking

## Table Join

Add a self-incrementing ID to the new table and import the table into the booktfidf table created in Django

```
1 INSERT INTO project.booktfidf SELECT * FROM project.tfidf;  
SELECT * FROM project.booktfidf;
```

t Grid  Filter Rows: <input type="text"/> Edit:    Export/Import:   W				
id	term	fileName	TFIDF	rowNum
1	abbot	10000.txt	2.094431030211474	59
2	abbot	10000.txt	2.094431030211474	180
3	abbey	10000.txt	1.564524419258904	344
4	abbot	10000.txt	2.094431030211474	346
5	abolished	10000.txt	0.993710655438553	356
6	abandon	10000.txt	0.538086675634899	385
7	abbey	10000.txt	1.564524419258904	392
8	abandon	10000.txt	0.538086675634899	394
9	abandon	10000.txt	0.538086675634899	426
10	abbot	10000.txt	2.094431030211474	540
11	abbot	10000.txt	2.094431030211474	540

# Indexing and Ranking

## Search speed up

In Django sever we define the model with unique\_index:

```
class BookTFIDF(models.Model):
    term = models.CharField(max_length = 40)
    fileName = models.CharField(max_length = 10)
    TFIDF = models.FloatField()
    rowNum = models.IntegerField()
    class Meta:
        unique_together = ["term", "fileName"]
        managed = True
        db_table = 'booktfidf'
        constraints = [models.UniqueConstraint(fields=["term", "fileName"], name = 'unique_token')]
```

Which make a index in MySQL session to speed up search:

Indexes in Table				Index Details	
Visible	Key	Type	Uniq Columns	Key Name:	unique_token
<input checked="" type="checkbox"/>	PRIMARY	BTREE	YES id	Index Type:	BTREE
<input checked="" type="checkbox"/>	unique_token	BTREE	YES term, fileName	Allows NULL:	
				Cardinality:	205530
				Comment:	
				User Comment:	
Columns in table				Packed:	
				Unique: YES	
Column	Type	Nulla	Indexes		
id	bigint	NO	PRIMARY		
term	varchar(40)	NO	unique_token		
fileName	varchar(10)	NO	unique_token		
TFIDF	double	NO			

# Server and UI

Search page

NWHoogle



friendly|



Search NWHoogle



## index Function

```
61 def index(request):
62     context={'results':{}, 'count':0, 'search_term':'empty'}
63     # get keyword of searching.
64     if request.method == "POST":
65         # just submit the result to the result url.
66         return sql_result(request)
67     else:
68         #If it is get, just stay in this page.
69         return render(request, 'minisearch/header.html', context)
```

## sql\_result Function

```
70 def sql_result(request):
71     # split the result, and find the preview in hdfs.
72     #time counter and variable initialization.
73     start_time = time.time()
74     results = [] # store the filename.
75     search_term = '' # store the term that user search.
76     context = dict() # used to render html.
77     # read in the index and filename related to input.
78     if request.method == "POST":
79         search_term=request.POST['keyword']
80         # call function in the search_call.py at current path.
81         results=sql_search(query=search_term)
82         results_list = list(results[1])
83         context={'results':results_list,'count':results[0],'search_term':search_term}
84         # acquire filename, TFIDF, rowNum and preview of text.
85         context['fileName'] = [list(results_list[i])[0] for i in range(results[0])]
86         context['TFIDF'] = [list(results_list[i])[1] for i in range(results[0])]
87         # context['rowNum'] = [list(results_list[i])[2] for i in range(results[0])]
88         # get the english meaning.
89         context['title1'] = "English Explantion"
90         context['explanation'] = get_meaning(search_term)
91         # get chinese translation.
92         context['title2'] = "Synonyms"
93         context['synonyms']= get_synonyms(search_term)
94
95         # get final result of the html page
96         context['final'] = ""
97         count = 1
98         fileName = list(context['fileName'])
```

# Server and UI

## Link to UI

```
102     for i in range(results[0]):
103         # get final results, with a special form.
104         # we can preview 3 lines besides the article.
105         context['final'] += "<table><th><tr>"
106         content,bookName = get_finalRes(fileName[i])
107         weblink = "<a href=http://localhost:8000/static/tfidf-index/"+fileName[i]+">"
108         context['final'] += weblink
109         context['final'] += "<h5>"
110         context['final'] += str(count)
111         context['final'] += ". "
112         context['final'] += bookName
113         context['final'] += "</h5></a></tr>"
114         context['final'] += content
115         context['final'] += "</h6></th></table>"
116         context['final'] += '<br>'
117         count+=1
118
119     # use amount to store the num of true term in the result.
120     # maybe the num of result < 20.
121     context['amount'] = count - 1
122     # end of time calculation, add time into context dictionary.
123     end_time = time.time()
124     context['time'] = str(round(end_time-start_time,2))
125     return render(request, 'minisearch/result.html', context)
126
```

## sql\_search Function

```
43 # here we use mysql statement to select the index.  
44 def sql_search(query):  
45     with connection.cursor() as cursor:  
46         cursor.execute("SELECT * FROM (SELECT filename, TFIDF FROM booktfidf WHERE term = '" + str(query) + "'" + ") AS a ORDER BY TFIDF DESC LIMIT 20")  
47         row = cursor.fetchall()  
48         num = len(row)  
49     return num, row
```

Here, we connect to MySQL database to fetch index, including file name, value of TFIDF which correspond to search term, and we set only take 20 of result which TFIDF is large. Pass number of result and data in database to sql\_result function.

## get\_finalRes Function

```
205 def get_finalRes(filename):
206     bookName=''
207     fileFetcher = Process_Data_Hdfs(filename)
208     content = ""
209     local_path = "/root/MyPyspark/Gutenberg/Django/dpwII_project/static/tfidf-index/"
210     if(os.path.exists(local_path + filename) == False):
211         fileFetcher.get_from_hdfs(local_path)
212
213     # after download the file, we open it to fetch info.
214     with open(local_path + filename) as f:
215         for lines in f.readlines():
216             # get the title of book.
217             if 'Title: ' in lines:
218                 bookName=lines[7:]
219             if 'Author: ' in lines:
220                 print(lines)
221                 content += '<tr><h6>'
222                 content += lines
223                 content += '</h6></tr>'
224                 break
225     return content,bookName
```

L.

## Process\_Data\_Hdfs Class

```
17 # define a class used to fetch files from hdfs.
18 class Process_Data_Hdfs():
19     # initialize an object.
20     def __init__(self, filename):
21         self.client = Client('http://localhost:9870')
22         self.filename = "/project_big/" + filename
23
24     # download file from hdfs.
25     def get_from_hdfs(self, local_path):
26         self.client.download(self.filename, local_path, overwrite=False)
27         return
28
29     # return file under directory on hdfs.
30     def list(self):
31         return self.client.list(self.filename, status=False)
32
```

## get\_meaning & get\_synonyms Function

```
270 #get the english meanings of the keywords.
271 def get_meaning(post_data):
272     # which means the input is a phrase,
273     # simply show a warning.
274     if(len(post_data.split())>1):
275         return "Sorry, A phrase is not supported to get the English explanation for now!"
276     # otherwise, create a pydictionary object.
277     dictionary = PyDictionary()
278     dict = dictionary.meaning(post_data)
279     try:
280         # try to get the lenth of the dictionary.
281         length = len(dict[list(dict.keys())[0]])
282     except:
283         # which means there is no explanation in the pydictionary.
284         return "The keyword you input has no explanation in the dictionary!"
285     # if try successful, then get the meaning of keywords.
286     value = list(dict.values())[0]
287     string = ""
288     count = 1
289     for i in range(length):
290         if count > 5:
291             break
292         # output all meanings of keywords.
293         string += str(i+1)
294         string += ". "
295         string += value[i]
296         # a new line.
297         string += "<br>"
298         count += 1
```

# Server and UI

Link to UI

```
301 #get the synonyms of the keywords.
302 def get_synonyms(post_data):
303     synonyms = []
304     try:
305         for syn in wordnet.synsets(post_data):
306             for l in syn.lemmas():
307                 synonyms.append(l.name())
308     except:
309         return "The keyword you input has no synonyms in the dictionary!"
310     if(len(synonyms)==0):
311         return "The keyword you input has no synonyms in the dictionary!"
312     str1 = ""
313     count = 1
314     synonyms = list(set(synonyms))
315     for i in synonyms:
316         if count > 5:
317             break
318         str1 += (str(count)+". "+i+"<br>")
319         count += 1
320     return str1
```

Call methods in PyDictionary package, to get further information and explanation of search term. We only select top 5 item in package. And pass them in specific form combining with HTML.



## search result



# NWHoogle Books Searcher

### About Us

NING Jiayi	2030026110
WANG Ruoxin	2030026150
HOU Ailing	2030026052

### Features

#### English Explanation

1. (usually followed by 'to')
2. have the skills and qualifications to do things well
3. having inherent physical or mental ability or capacity
4. having a strong healthy body

#### Synonyms

1. capable
2. able
3. able-bodied

### Supports

Time Cost: **1.35s**  
We found **20 results** for your search "able"

1. [The Works of Samuel Johnson, Vol. 11.](#)  
Author: Samuel Johnson
2. [The Works of Samuel Johnson, Vol. 10.](#)  
Author: Samuel Johnson
3. [The Works of Samuel Johnson, Vol. 6](#)  
Author: Samuel Johnson
4. [Fruitfulness \(Fécondité\)](#)  
Author: Émile Zola
5. [Maezli](#)  
Author: Johanna Spyri
6. [The Minute Boys of the Mohawk Valley](#)  
Author: James Otis
7. [The Book of Household Management](#)  
Author: Mrs. Isabella Beeton
8. [Gutta-Percha Willie](#)

## Preview Book

The Project Gutenberg EBook of The Poetical Works of William Wordsworth  
Edited by William Knight

This eBook is for the use of anyone anywhere at no cost and with  
almost no restrictions whatsoever. You may copy it, give it away or  
re-use it under the terms of the Project Gutenberg License included  
with this eBook or online at [www.gutenberg.org](http://www.gutenberg.org)

Title: The Poetical Works of William Wordsworth  
Volume 1 of 8

Author: (Edited by William Knight)

Release Date: November 23, 2003 [EBook #10219]

Language: English

\*\*\* START OF THIS PROJECT GUTENBERG EBOOK POETRY OF WORDSWORTH \*\*\*

Produced by Jonathan Ingram, Clytie Siddall and the Online Distributed  
Proofreading Team.

## Meaning and Synonyms

### Features

#### English Explantion

1. (usually followed by 'to')
2. have the skills and qualifications to do things well
3. having inherent physical or mental ability or capacity
4. having a strong healthy body

#### Synonyms

1. capable
2. able
3. able-bodied

# DEMO !