

## Homework 2

### 1. Identifying and removing stop words

The word count mapper function and reducer function are created to get the word count output. The input files include three txt files: Halmet.txt, The\_Tragedie\_of\_Macbeth.txt, and The\_Tragedy\_of\_Romeo\_and\_Juliet.txt, and the output is in the form of "word count". Since in this step we just need to identify stop words, so we just print stop words out to the output.

```
mapper_stopwords.py
1  #!/usr/bin/env python
2  import os
3  import re
4  import sys
5
6  #works_path = r'works/'
7
8  # tokenize function
9  def tokenize(text):
10     min_length = 3 # typically, words with length smaller than 3 is not meaningful
11     words = map(lambda word: word.lower(), text.strip().split(' ')) # case insensitive
12     p = re.compile('[a-zA-Z]+$') # remove the mess up
13     filtered_words = list(filter(lambda token: p.match(token) and len(token)>=min_length, words))
14     return filtered_words
15
16 # get the word + count
17 for line in sys.stdin:
18     words = tokenize(line.decode('utf-8'))
19     if words:
20         for word in words:
21             print('%s\t%s' % (word,1))
22
```

```
reducer_stopwords.py
1  #!/usr/bin/env python
2  import sys
3  from os import system
4
5  previous = None
6  wc_list = {}
7  sum = 0
8  # input comes from STDIN
9  for line in sys.stdin:
10     # remove leading and trailing whitespace
11     line = line.strip()
12     key, value = line.split('\t',1)
13     # by key (here: word) before it is passed to the reducer
14     if key == previous:
15         sum += int(value)
16     else:
17         if previous:
18             wc_list[previous] = sum
19         sum = int(value)
20         previous = key
21 if previous == key: # the last word
22     wc_list[previous] = sum
23 sorted_wc_list = sorted(wc_list.items(),key=lambda (k,v):v)
24 # output stop words (count > 150)
25 for key,value in sorted_wc_list:
26     if value > 150:
27         print key + '\t' + str(value)
28
```

We wrote a tokizen function to get the cleaned words. What the function dose is keeping words that only contains characters, changing them to lowercase and removing those with less than three characters.

```
# tokenize function
def tokenize(text):
    min_length = 3 # typically, words with length smaller than 3 is not meaningful
    words = map(lambda word: word.lower(), text.strip().split(' ')) # case insensitive
    p = re.compile('[a-zA-Z]+$') # remove the mess up
    filtered_words = list(filter(lambda token: p.match(token) and len(token)>=min_length, words))
    return filtered_words
```

In order to identify stop words, we simply take the words with count bigger than 150 as stop words, and print them out to the output. In addition, we transfer the output from HDFS to the local file system and name it as stopwords.txt, which will be used in step 2.

```
# output stop words (count > 150)
for key,value in sorted_wc_list:
    if value > 150:
        print key + '\t' + str(value)
```

```
zhihua@namenode:~$ hadoop fs -cat /tmp/stopwords/part-00000
when 151
should 152
must 152
like 153
then 163
she 163
their 166
thee 168
let 174
was 174
hath 176
would 177
more 178
how 182
may 194
project 195
which 204
good 217
enter 220
they 234
him 243
her 245
from 277
shall 283
all 303
thy 305
are 318
our 324
have 326
will 374
what 420
thou 449
your 501
but 566
his 577
for 642
not 696
this 699
with 774
that 946
you 1024
and 2397
the 2907
```

```
# get the stopwords output to local file system and name it to stopwords.txt
hadoop fs -get /tmp/stopwords .
cp stopwords/part-00000 stopwords.txt
```

## 2. Building the Inverted Index

The mapper\_index and reducer\_index files are generated to get the inverted index. The output of inverted index results includes word, wordcount, and location (document id and line number).

```
1 #!/usr/bin/env python
2 import os
3 import re
4 import sys
5
6
7 # tokenize function
8 def tokenize(text):
9     min_length = 3 # typically, words with length smaller than 3 is not meaningful
10    words = map(lambda word: word.lower(), text.strip().split(' ')) # case insensitive
11    p = re.compile('[a-zA-Z]+$') # remove the mess up
12    filtered_words = list(filter(lambda token: p.match(token) and len(token)>=min_length, words))
13    return filtered_words
14
15 # get the nonstopwords
16 word_list = []
17 with open("stopwords.txt", 'r') as f:
18     for line in f:
19         word, count = line.strip().split('\t')
20         word_list.append(word)
21
22 # get the word+line
23 fname = None
24
```

```
1 #!/usr/bin/env python
2 import sys
3
4 previous = None
5 wc = 0
6 count = 1
7 doclist = []
8 # input comes from STDIN
9 for line in sys.stdin:
10     # remove leading and trailing whitespace
11     line = line.strip()
12     key, doc, lc = line.split('\t')
13     # by key (here: word) before it is passed to the reducer
14     if key == previous:
15         count+=1
16         doclist.append((doc,lc))
17     else:
18         if previous:
19             print(previous + '\t' + str(count) + "\t" + str(doclist))
20             doclist = []
21         count=1
22         previous = key
23         doclist.append((doc,lc))
24 if previous == key:
25     print(previous + '\t' + str(count) + "\t" + str(doclist))
26
```

When we run the code, the final inverted index can be obtained, which won't contain the stop words (stopwords.txt) identified in Step 1. Then we transfer the output to local file system so that we can query the results locally.

```
# get the nonstopwords
word_list =[]
with open("stopwords.txt", 'r') as f:
    for line in f:
        word, count = line.strip().split('\t')
        word_list.append(word)

# get the index output to local file system
hadoop fs -get /tmp/index .
```

```
zhihua@namenode:~$ clear
zhihua@namenode:~$ hadoop fs -cat /tmp/index/part-00000 | head -10
abate 2      [('Halmet', '3701'), ('The Tragedy of Romeo and Juliet', '3450')]
abate       1      [('Halmet', '3704')]
abatement 1      [('The Tragedy of Romeo and Juliet', '1964')]
abbreviations 1      [('The Tragedie of Macbeth', '355')]
abed 1      [('The Tragedy of Romeo and Juliet', '2907')]
abhorred 3      [('The Tragedy of Romeo and Juliet', '4196'), ('The Tragedie of Macbeth', '3495'), ('The Tragedy of Romeo and Juliet', '3100')]
abide 4      [('The Tragedy of Romeo and Juliet', '4544'), ('Halmet', '4758'), ('The Tragedie of Macbeth', '2893')]
abiure 1      [('The Tragedie of Macbeth', '2893')]
able 3      [('The Tragedy of Romeo and Juliet', '4382'), ('The Tragedy of Romeo and Juliet', '206'), ('The Tragedie of Macbeth', '2890')]
cat: Unable to write to output stream.
zhihua@namenode:~$ hadoop fs -cat /tmp/index/part-00000 | head -20
abate 2      [('Halmet', '3701'), ('The Tragedy of Romeo and Juliet', '3450')]
abate       1      [('Halmet', '3704')]
abatement 1      [('The Tragedy of Romeo and Juliet', '1964')]
abbreviations 1      [('The Tragedie of Macbeth', '355')]
abed 1      [('The Tragedy of Romeo and Juliet', '2907')]
abhorred 3      [('The Tragedy of Romeo and Juliet', '4196'), ('The Tragedie of Macbeth', '3495'), ('The Tragedy of Romeo and Juliet', '3100')]
abide 4      [('The Tragedy of Romeo and Juliet', '4544'), ('Halmet', '4758'), ('The Tragedie of Macbeth', '2893')]
abiure 1      [('The Tragedie of Macbeth', '2893')]
able 3      [('The Tragedy of Romeo and Juliet', '4382'), ('The Tragedy of Romeo and Juliet', '206'), ('The Tragedie of Macbeth', '2890')]
aboue 2      [('The Tragedie of Macbeth', '3062'), ('The Tragedie of Macbeth', '2890')]
abound 1      [('The Tragedie of Macbeth', '2859')]
about 53      [('The Tragedie of Macbeth', '1204'), ('The Tragedie of Macbeth', '3214'), ('The Tragedie of Macbeth', '110'), ('The Tragedie of Macbeth', '1201'), ('The Tragedie of Macbeth', '2936'), ('The Tragedie of Macbeth', '2402'), ('The Tragedie of Macbeth', '59'), ('The Tragedie of Macbeth', '1337'), ('The Tragedie of Macbeth', '1973'), ('Halmet', '4186'), ('Halmet', '4966'), ('Halmet', '1659'), ('Halmet', '280'), ('Halmet', '1512'), ('Halmet', '4877'), ('Halmet', '2711'), ('Halmet', '1812'), ('Halmet', '3860'), ('Halmet', '3901'), ('Halmet', '1151'), ('Halmet', '1150'), ('Halmet', '1934'), ('Halmet', '5043'), ('The Tragedy of Romeo and Juliet', '2292'), ('The Tragedy of Romeo and Juliet', '4850'), ('The Tragedy of Romeo and Juliet', '4831'), ('The Tragedy of Romeo and Juliet', '4853'), ('The Tragedy of Romeo and Juliet', '4751'), ('The Tragedy of Romeo and Juliet', '574'), ('The Tragedy of Romeo and Juliet', '4796'), ('The Tragedy of Romeo and Juliet', '986'), ('The Tragedy of Romeo and Juliet', '1933'), ('The Tragedy of Romeo and Juliet', '2078'), ('The Tragedy of Romeo and Juliet', '3955'), ('The Tragedy of Romeo and Juliet', '3283'), ('The Tragedy of Romeo and Juliet', '1308'), ('The Tragedy of Romeo and Juliet', '2367'), ('The Tragedy of Romeo and Juliet', '865'), ('The Tragedie of Macbeth', '211'), ('Halmet', '1738')]
abridgement 1      [('Halmet', '1738')]
absence 2      [('The Tragedie of Macbeth', '2079'), ('The Tragedie of Macbeth', '1849')]
absent 2      [('The Tragedie of Macbeth', '3356'), ('Halmet', '4637')]
absolute 5      [('Halmet', '3958'), ('Halmet', '4291'), ('The Tragedie of Macbeth', '2330'), ('The Tragedie of Macbeth', '2330'), ('The Tragedie of Macbeth', '2330'), ('The Tragedie of Macbeth', '2330'), ('The Tragedie of Macbeth', '2330')]
```

### 3. Query the Inverted Index

A query program on top of your full inverted file index that accepts a user-specified query (one or more words) and returns not only the document IDs but also the locations in the form of line numbers is created.

```
#!/usr/bin/env python

path = r'index/part-00000'

user_input = raw_input("Please input the words you want to query, seperate by space: ")

words = user_input.split(" ")
flag = 0 # check if the word can be found
with open(path, 'r') as f:
    for line in f:
        word, count, location = line.strip().split('\t')
        if word in words:
            print("The number of occurence of %s: %s" % (word, count))
            print("The location of occurence of %s: %s" % (word, location))
            flag = 1
if flag == 0:
    print("Sorry, the words you are querying cannot be found!")
```

```
zhihua@namenode:~$ ./query.py
Please input the words you want to query, seperate by space: hello
Sorry, the words you are querying cannot be found!
zhihua@namenode:~$ ./query.py
Please input the words you want to query, seperate by space: juliet
The number of occurence of juliet: 27
The location of occurence of juliet: [('The Tragedy of Romeo and Juliet', '2748'), ('The Tragedy of Romeo and Juliet', '46'), ('The Tragedy of Romeo and Juliet', '3529'), ('The Tragedy of Romeo and Juliet', '2476'), ('The Tragedy of Romeo and Juliet', '66'), ('The Tragedy of Romeo and Juliet', '1520'), ('The Tragedy of Romeo and Juliet', '4301'), ('The Tragedy of Romeo and Juliet', '4495'), ('The Tragedy of Romeo and Juliet', '2939'), ('The Tragedy of Romeo and Juliet', '1170'), ('The Tragedy of Romeo and Juliet', '1215'), ('The Tragedy of Romeo and Juliet', '4481'), ('The Tragedy of Romeo and Juliet', '1311'), ('The Tragedy of Romeo and Juliet', '3733'), ('The Tragedy of Romeo and Juliet', '4260'), ('The Tragedy of Romeo and Juliet', '2982'), ('The Tragedy of Romeo and Juliet', '3551'), ('The Tragedy of Romeo and Juliet', '4171'), ('The Tragedy of Romeo and Juliet', '4054'), ('The Tragedy of Romeo and Juliet', '19'), ('The Tragedy of Romeo and Juliet', '9'), ('The Tragedy of Romeo and Juliet', '2705'), ('The Tragedy of Romeo and Juliet', '2953'), ('The Tragedy of Romeo and Juliet', '1490'), ('The Tragedy of Romeo and Juliet', '4397'), ('The Tragedy of Romeo and Juliet', '1308'), ('The Tragedy of Romeo and Juliet', '2843')]
```

## Extra credit: Integrate queries with Spark

Several codes in the Spark framework are wrote to query the inverted index. The results are showed below.

```
In [1]: # read the data
path = r'Hw2/index/part-00000' #the part of index output file
data = sc.textFile(path)
```

```
In [14]: # query words you want

query_words = ['abroad', 'able', 'juliet']

filter_rdd = data.filter(lambda x: x.split("\t")[0] in query_words)
results = filter_rdd.collect()
if not results:
    print("Sorry, the words you are querying cannot be found!")
else:
    word_list = filter_rdd.map(lambda x: x.split("\t")[0]).collect()
    for word in query_words:
        if word not in word_list:
            print("The result of word {s} not found!" % word)
    for res in results:
        word, count, location = res.split("\t")
        print("The number of occurrence of {s}: %s" % (word, count))
        print("The location of occurrence of {s}: %s" % (word, location))
```

```
The result of word {abroad} not found!
The number of occurrence of {able}: 3
The location of occurrence of {able}: [('The_Tragedy_of_Romeo_and_Juliet', '4382'), ('The_Tragedy_of_Romeo_and_Juliet', '206'),
('Helmet', '4405')]
The number of occurrence of {juliet}: 27
The location of occurrence of {juliet}: [('The_Tragedy_of_Romeo_and_Juliet', '2748'), ('The_Tragedy_of_Romeo_and_Juliet', '46'),
('The_Tragedy_of_Romeo_and_Juliet', '3529'), ('The_Tragedy_of_Romeo_and_Juliet', '2476'), ('The_Tragedy_of_Romeo_and_Juliet',
'66'), ('The_Tragedy_of_Romeo_and_Juliet', '1520'), ('The_Tragedy_of_Romeo_and_Juliet', '4301'), ('The_Tragedy_of_Romeo_and_Juliet', '4495'), ('The_Tragedy_of_Romeo_and_Juliet', '2939'), ('The_Tragedy_of_Romeo_and_Juliet', '1170'), ('The_Tragedy_of_Romeo_and_Juliet', '1215'), ('The_Tragedy_of_Romeo_and_Juliet', '4481'), ('The_Tragedy_of_Romeo_and_Juliet', '1311'), ('The_Tragedy_of_Romeo_and_Juliet', '3733'), ('The_Tragedy_of_Romeo_and_Juliet', '4260'), ('The_Tragedy_of_Romeo_and_Juliet', '2982'), ('The_Tragedy_of_Romeo_and_Juliet', '3551'), ('The_Tragedy_of_Romeo_and_Juliet', '4171'), ('The_Tragedy_of_Romeo_and_Juliet', '4054'), ('The_Tragedy_of_Romeo_and_Juliet', '19'), ('The_Tragedy_of_Romeo_and_Juliet', '9'), ('The_Tragedy_of_Romeo_and_Juliet', '2705'), ('The_Tragedy_of_Romeo_and_Juliet', '2953'), ('The_Tragedy_of_Romeo_and_Juliet', '1490'), ('The_Tragedy_of_Romeo_and_Juliet', '4397'), ('The_Tragedy_of_Romeo_and_Juliet', '1308'), ('The_Tragedy_of_Romeo_and_Juliet', '2843')]
```