

In [1]: `from pyspark.sql import functions as F`

Starting Spark application

ID	YARN Application ID	Kind	State	Spark UI	Driver log	Current session?
2	application_1635723418847_0003	pyspark	idle	Link	Link	✓

SparkSession available as 'spark'.

In [2]: `def countWords (fileName):
 textfile = sc.textFile(fileName)
 lines = textfile.flatMap(lambda line: line.split(" "))
 counts = lines.map (lambda word: (word, 1))
 aggregatedCounts = counts.reduceByKey (lambda a, b: a + b)
 return aggregatedCounts.top (200, key=lambda p : p[1])`

In [3]: `countWords("s3://chrisjermainebucket/text/Holmes.txt")`

```
[('the', 5404), ('', 3145), ('and', 2798), ('of', 2720), ('to', 2700),
('a', 2575), ('I', 2533), ('in', 1702), ('that', 1559), ('was', 1360),
('his', 1096), ('is', 1076), ('you', 1029), ('he', 1014), ('it', 976),
('my', 901), ('have', 893), ('with', 843), ('had', 806), ('as', 776), ('
which', 753), ('at', 739), ('for', 697), ('be', 612), ('not', 598), ('fr
om', 485), ('upon', 460), ('said', 448), ('but', 441), ('me', 414), ('we
', 413), ('this', 407), ('been', 385), ('very', 371), ('her', 367), ('yo
ur', 359), ('"I', 349), ('were', 336), ('on', 334), ('by', 334), ('an',
329), ('all', 321), ('so', 317), ('are', 316), ('would', 313), ('she', 3
05), ('It', 290), ('no', 286), ('one', 283), ('could', 280), ('has', 27
7), ('there', 275), ('The', 273), ('into', 272), ('out', 272), ('He', 26
4), ('what', 264), ('or', 260), ('Mr.', 259), ('when', 257), ('little',
257), ('him', 253), ('who', 253), ('will', 250), ('up', 250), ('some', 2
27), ('do', 217), ('should', 207), ('down', 204), ('may', 201), ('Holmes
', 197), ('our', 195), ('man', 193), ('if', 189), ('see', 184), ('am', 1
81), ('shall', 170), ('must', 168), ('can', 165), ('about', 163), ('over
', 161), ('than', 159), ('any', 157), ('only', 155), ('more', 151), ('ca
me', 142), ('other', 140), ('they', 140), ('before', 138), ('know', 13
7), ('You', 136), ('think', 132), ('two', 128), ('Holmes', 127), ('us',
126), ('did', 126), ('"It', 124), ('There', 121), ('might', 118), ('come
', 117), ('"You', 112), ('it.', 110), ('just', 110), ('such', 110), ('mu
ch', 107), ('back', 106), ('heard', 104), ('made', 102), ('time', 102),
('But', 100), ('where', 100), ('found', 100), ('"And', 99), ('how', 96),
('Sherlock', 96), ('now', 95), ('their', 95), ('it', 94), ('own', 94),
('never', 92), ('then', 92), ('like', 90), ('after', 90), ('however', 8
9), ('quite', 89), ('We', 89), ('most', 86), ('good', 85), ('through', 8
5), ('took', 84), ('tell', 84), ('them', 84), ('away', 84), ('She', 84),
('its', 84), ('saw', 84), ('And', 83), ('me,', 83), ('him.', 82), ('S
t.', 80), ('go', 80), ('Project', 80), ('way', 79), ('without', 79), ('f
ace', 79), ('nothing', 78), ('Holmes.', 78), ('Miss', 77), ('few', 77),
```

```
( 'make', 76), ( 'left', 76), ( 'matter', 75), ( 'every', 75), ( 'door', 75),
( 'small', 75), ( 'take', 74), ( 'last', 74), ( 'me.', 74), ( 'you,', 74), ( '
find', 74), ( 'until', 73), ( 'long', 73), ( 'young', 73), ( 'A', 73), ( '"Th
e', 73), ( 'say', 72), ( 'case', 72), ( 'As', 72), ( '"But', 72), ( 'he,', 6
9), ( 'these', 68), ( 'Then', 68), ( 'put', 67), ( '"Well,', 67), ( 'first',
67), ( 'then,', 66), ( 'once', 65), ( 'seemed', 65), ( 'round', 65), ( 'thoug
ht', 64), ( 'right', 64), ( 'even', 64), ( 'while', 64), ( 'him,', 64), ( 'we
nt', 63), ( 'If', 63), ( 'seen', 62), ( 'old', 62), ( 'ever', 61), ( 'three',
61), ( 'he.', 61), ( 'still', 61), ( 'himself', 61), ( 'hand', 61), ( 'those
', 60), ( 'rather', 59), ( 'though', 59), ( 'something', 59), ( '"Oh,', 58),
```

Task2

In [4]:

```
import re
import numpy as np

# load up all of the 19997 documents in the corpus
corpus = sc.textFile ("s3://chrisjermainebucket/comp330_A6/20_news_same_

# each entry in validLines will be a line from the text file (that has a
validLines = corpus.filter(lambda x : 'id' in x)

# now we transform it into a bunch of (docID, text) pairs
keyAndText = validLines.map(lambda x : (x[x.index('id="')+ 4 : x.index(

# now we split the text in each (docID, text) pair into a list of words
# after this, we have a data set with (docID, ["word1", "word2", "word3"
# we have a bit of fancy regular expression stuff here to make sure that
# die on some of the documents
regex = re.compile('[^a-zA-Z]')
keyAndListOfWords = keyAndText.map(lambda x : (str(x[0]), regex.sub(' ',

# now get the top 20,000 words... first change (docID, ["word1", "word2"
# to ("word1", 1) ("word2", 1)...
allWords = keyAndListOfWords.flatMap(lambda x: ((j, 1) for j in x[1]))

# now, count all of the words, giving us ("word1", 1433), ("word2", 3423
allCounts = allWords.reduceByKey(lambda a, b: a + b)

# and get the top 20,000 words in a local array
# each entry is a ("word1", count) pair
topWords = allCounts.top (20000, lambda x : x[1])
```

In [5]:

```
# And we'll create a RDD that has a bunch of (word, dictNum) pairs
# start by creating an RDD that has the number 0 thru 20000
# 20000 is the number of words that will be in our dictionary
twentyK = sc.parallelize(range(20000))

# now, we transform (0), (1), (2), ... to ("mostcommonword", 1) ("nextmo
# the number will be the spot in the dictionary used to tell us where th
# HINT: make use of topWords in the lambda that you supply
dictionary = twentyK.map(lambda x: (topWords[x][0], x))

# finally, print out some of the dictionary, just for debugging
dictionary.top(10)
# print(dictionary)
```

```
[('zz', 6488), ('zyxel', 13879), ('zyeh', 18361), ('zy', 8935), ('zx', 4
113), ('zw', 9729), ('zvm', 19090), ('zv', 3579), ('zurich', 15636), ('z
uma', 3634)]
```

Assignment 4

Task 1:

- First, get an RDD encoding your dictionary, where the RDD has a bunch of (word, posInDictionary) pairs.
- Next, create a second RDD that effectively has a bunch of (word, docID) pairs, where the word occurs in the given document (you get this just like the code from lab, where you flatMap the document corpus)

In [39...

```

import re
import numpy as np

# load up all of the 19997 documents in the corpus
corpus = sc.textFile ("s3://chrisjermainebucket/comp330_A6/20_news_same_

# each entry in validLines will be a line from the text file (that has a
validLines = corpus.filter(lambda x : 'id' in x)

# now we transform it into a bunch of (docID, text) pairs
keyAndText = validLines.map(lambda x : (x[x.index('id="') + 4 : x.index(

# now we split the text in each (docID, text) pair into a list of words
# after this, we have a data set with (docID, ["word1", "word2", "word3"
# we have a bit of fancy regular expression stuff here to make sure that
# die on some of the documents
regex = re.compile('[^a-zA-Z]')
keyAndListOfWords = keyAndText.map(lambda x : (str(x[0]), regex.sub(' ',

# now get the top 20,000 words... first change (docID, ["word1", "word2"
# to ("word1", 1) ("word2", 1)...
word_docID_unmapped = keyAndListOfWords.flatMap(lambda x: ((j, x[0]) for
word_docID_unmapped_local = word_docID_unmapped.collect()

```

- Now, a lot of those words in the documents won't actually appear in the dictionary. But if you join the two RDDs, you'll have a bunch of (word, (docID, posInDictionary)) pairs, where the given document has the given word at the given position in the dictionary.

In [40...

```
dictionary.top(10)
```

```

[('zz', 6488), ('zyxel', 13879), ('zyeh', 18361), ('zy', 8935), ('zx', 4
113), ('zw', 9729), ('zvm', 19090), ('zv', 3579), ('zurich', 15636), ('z
uma', 3634)]

```

In [41...

```
word_docID_unmapped.top(10)
```

```

[('zzzzzzt', '20_newsgroups/rec.sport.baseball/104569'), ('zzzzzz', '20_
newsgroups/rec.sport.hockey/53841'), ('zzzzzz', '20_newsgroups/rec.spor
t.baseball/105004'), ('zzzzzz', '20_newsgroups/rec.sport.baseball/105002
'), ('zzzzzz', '20_newsgroups/rec.sport.baseball/104795'), ('zzzzzz', '2
0_newsgroups/rec.sport.baseball/104540'), ('zzzzzz', '20_newsgroups/rec.
motorcycles/105113'), ('zzzzzz', '20_newsgroups/rec.motorcycles/104730
'), ('zzzz', '20_newsgroups/comp.sys.ibm.pc.hardware/60262'), ('zzzz', '
20_newsgroups/comp.sys.ibm.pc.hardware/60262')]

```

In [42...

```
word_docID_unmapped = word_docID_unmapped.join(dictionary)
```

In [43...

```
word_docID_unmapped.top(10)
```

```
[('zz', ('20_newsgroups/talk.politics.guns/54380', 6488)), ('zz', ('20_newsgroups/talk.politics.guns/54380', 6488)), ('zz', ('20_newsgroups/talk.politics.guns/54380', 6488)), ('zz', ('20_newsgroups/sci.med/59185', 6488)), ('zz', ('20_newsgroups/sci.crypt/15545', 6488)), ('zz', ('20_newsgroups/sci.crypt/15545', 6488)), ('zz', ('20_newsgroups/rec.sport.baseball/105004', 6488)), ('zz', ('20_newsgroups/rec.sport.baseball/105002', 6488)), ('zz', ('20_newsgroups/rec.sport.baseball/104795', 6488)), ('zz', ('20_newsgroups/rec.sport.baseball/104540', 6488))]
```

- Next, process this RDD (using an appropriate Spark operation) so that you get a bunch of (docid, (listOfAllDictionaryPos)) pairs. Not surprisingly, listOfAllDictionaryPos lists all of the posInDictionary values found for that document.

In [44...

```
word_docID_mapped = word_docID_unmapped.map(lambda x: (x[1][0], x[1][1]))
word_docID_mapped.top(1)
```

```
[('20_newsgroups/talk.religion.misc/84570', 19904)]
```

In [45...

```
word_docID_mapped = word_docID_mapped.groupByKey()
word_docID_mapped.top(1)
```

```
[('20_newsgroups/talk.religion.misc/84570', <pyspark.resultiterable.ResultIterable object at 0x7f5c51209150>)]
```

In [46...

```
word_docID_mapped = word_docID_mapped.map(lambda x : (x[0], list(x[1])))
```

- Then finally, you will write a map () that will take that RDD and convert into the listOfAllDictionaryPos values to a NumPy array.

In [53...

```
import numpy as np

def list_to_np(ls):
    arr = np.zeros(20000)
    for i in ls:
        arr[i] += 1
    return arr

res = word_docID_mapped.map(lambda x: ((x[0]), list_to_np(x[1])))
```

In [66...

```
# [tup[1] for tup in res if tupp[0] == "20_newsgroups/comp.graphics/37261"]
import numpy as np

result1 = np.array(res.lookup("20_newsgroups/comp.graphics/37261"))
result1[result1.nonzero()]
```

```
array([ 8.,  2.,  6.,  3., 12.,  4.,  3.,  6.,  2.,  1.,  1.,  5.,  2.,
        2.,  2.,  3.,  1.,  1.,  1.,  1.,  3.,  1.,  1.,  2.,  3.,  4.,
        1.,  1.,  1.,  1.,  1.,  3.,  1.,  1.,  1.,  2.,  1.,  1.,  1.,
        2.,  1.,  1.,  2.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  2.,
        1.,  1.,  2.,  2.,  1.,  2.,  1.,  1.,  1.,  3.,  4.,  1.,  1.,
        1.,  1.,  2.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  2.,  1.,
        1.,  1.,  1.,  5.,  2.,  2.,  1.,  1.,  5.,  1.,  4.,  1.,  1.,
        1.,  2.,  1.,  2.,  1., 11.,  1.,  1.,  1.,  1.,  2.,  2.,  2.,
        5.,  1.,  2.,  1.,  1.,  1.,  1.,  2.,  1.,  2.,  2.,  4.,  1.,
        1.,  1.,  5.,  1.,  1.,  1.,  1.,  2.,  4.,  1.,  1.,  1.,  3.,
        1.,  1.,  1.,  1.,  3.,  2.,  2.,  1.,  1.,  6.,  1.,  6.,  1.,
        1.,  3.,  1.,  1.,  2.,  1.,  1.,  1.,  1.,  2.,  7.,  1.,  1.,
        1.,  1.,  1.]])
```

In [69...

```
import numpy as np
result2 = np.array(res.lookup("20_newsgroups/talk.politics.mideast/75944"))
result2[result2.nonzero()]
```

```
array([135.,  37.,  71.,  28.,  49.,  19.,  46.,  16.,  13.,  22.,  9.,
        22.,  11.,  7.,  7.,  6.,  4.,  6.,  12.,  11.,  10.,  3.,
        10.,  4.,  2.,  21.,  5.,  4.,  2.,  2.,  1.,  1.,  1.,
        5.,  1.,  23.,  5.,  2.,  1.,  6.,  8.,  4.,  7.,  3.,
        3.,  2.,  1.,  1.,  6.,  4.,  4.,  7.,  1.,  8.,  7.,
        13.,  4.,  4.,  10.,  3.,  3.,  2.,  2.,  3.,  7.,  4.,
        1.,  2.,  4.,  8.,  4.,  7.,  2.,  1.,  1.,  2.,  1.,
        2.,  2.,  1.,  5.,  3.,  3.,  3.,  1.,  1.,  1.,  2.,
        1.,  4.,  3.,  1.,  3.,  3.,  4.,  7.,  1.,  2.,  1.,
        3.,  2.,  1.,  4.,  6.,  3., 11.,  1.,  6.,  3.,  1.,
        3.,  1.,  2.,  1.,  1.,  1.,  3.,  3.,  2.,  5.,  2.,
        2.,  2.,  2.,  1.,  1.,  1.,  1.,  1.,  3.,  1.,  1.,
        1.,  1.,  3.,  3.,  4.,  1.,  1.,  5.,  1.,  1.,  2.,
        6.,  2.,  2.,  1.,  1.,  1.,  1.,  1.,  1.,  3.,  5.,
        1.,  1.,  1.,  1.,  2.,  1.,  1.,  1.,  6.,  1.,  1.,
        2.,  1.,  3.,  2.,  1.,  1.,  3.,  2.,  2.,  3.,  8.,
        1.,  1.,  1.,  1.,  2.,  3.,  1.,  2.,  6.,  1.,  1.,
        2.,  1., 13.,  4.,  1.,  1.,  1.,  1.,  1.,  3.,  1.,
        1.,  2.,  2.,  1.,  1.,  1.,  2.,  1.,  1.,  1.,  1.,
        3.,  1.,  2.,  4., 26.,  1.,  1.,  3.,  2.,  2.,  1.,
        3.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  2.,  6.,
        1.,  1.,  1.,  6.,  1.,  1.,  1.,  4.,  2.,  1.,  1.,
        3.,  4.,  4.,  4.,  3.,  1.,  3.,  1.,  1.,  1.,  1.,
        1.,  1.,  4.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,  1.,
        1.,  1.,  1.,  1.,  2.,  1.,  1.,  2.,  1.,  9.,  1.,
        1.,  2.,  1.,  1.,  1.,  6.,  1.,  1.,  1.,  1.,  1.,
        3.,  1.,  2.,  1.,  1.,  1.,  2.,  1.,  2.,  9.,  1.,
        1.,  1.,  2.,  1.,  1.,  2.,  2.,  2.,  1.,  1.,  2.,
```

```

4., 1., 1., 1., 2., 1., 1., 1., 1., 11., 2.,
1., 1., 1., 1., 1., 1., 1., 2., 1., 1., 1.,
1., 1., 1., 2., 9., 1., 2., 7., 3., 3., 2.,
1., 1., 2., 1., 1., 1., 2., 2., 1., 1., 3.,
1., 8., 1., 1., 1., 1., 1., 1., 1., 8., 1.,
2., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
1., 10., 1., 1., 1., 1., 2., 2., 3., 4., 1.,
1., 1., 1., 1., 1., 1., 3., 1., 1., 1., 1.,
1., 1., 8., 1., 1., 1., 1., 1., 1., 1., 6.,
1., 1., 1., 1., 1., 1., 1., 1., 1., 2., 1.,
1., 1., 1., 1., 11., 2., 1., 1., 1., 1., 1.,
2., 1., 1., 3., 3., 1., 1., 1., 1., 1., 1.,
1., 1., 1., 1., 1., 2., 1., 3., 1., 1., 1.,
1., 1., 1., 1., 1., 1., 1., 1., 2., 1., 4.,
2., 3., 1., 1., 1., 4., 1., 1., 4., 1., 1.,
1., 2., 4., 1., 1., 1., 1., 2., 3., 1., 1.,
1., 2., 1., 2., 1., 1., 1., 1., 1., 1., 1.,
1., 1., 1., 2., 1., 1., 1., 1., 1., 2., 1.,
2., 1., 2., 2., 1., 1., 1., 1., 1., 4., 1.,
1., 1., 1., 8., 1., 2., 1., 1., 1., 2., 3.,
1., 1., 1., 1., 1., 3., 1., 1., 1., 1., 1.,
4., 1., 1., 1., 3., 2., 1., 1., 1., 1., 2.,
1., 1., 1., 1., 2., 1., 1., 1., 1., 1., 1.,
1., 1., 1., 2., 1., 1., 1., 1., 1., 1.,
1., 1., 1., 1., 1., 2., 1., 1., 1., 1.,
4 3 2 1 2 1 1 2 1

```

In [70...

```

import numpy as np
result3 = np.array(res.lookup("20_newsgroups/sci.med/58763"))
result3[result3.nonzero()]

```

```

array([4., 4., 3., 2., 1., 1., 4., 3., 1., 2., 1., 5., 1., 2., 1., 1.,
1.,
2., 1., 1., 1., 1., 1., 1., 2., 1., 1., 1., 1., 1., 1., 1.,
1.,
1., 1., 1., 1., 1., 1., 1., 2., 2., 1., 1., 1., 1., 1., 1.,
1.,
1., 1., 1., 1., 2., 1., 1., 1., 1., 2., 1., 1., 1., 2., 1., 5.,
1.,
1., 1., 1., 1., 1., 1., 2., 1., 2., 1., 1., 1., 1., 1., 1.,
1.,
1., 3., 1., 1.])

```

In []: