

```

1 !apt-get install openjdk-8-jdk-headless -qq > /dev/null
2 !wget -q https://apache.mirror.colo-serv.net/spark/spark-2.4.7/spark-2.4.7-bin-hadoop2.7.t
3 !tar xf spark-2.4.7-bin-hadoop2.7.tgz
4 !pip install -q findspark
5
6 import os
7 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
8 os.environ["SPARK_HOME"] = "/content/spark-2.4.7-bin-hadoop2.7"
9
10 import findspark
11 findspark.init("spark-2.4.7-bin-hadoop2.7")# SPARK_HOME
12
13 import pyspark
14 from pyspark.sql import *
15 from pyspark.sql.functions import *
16 from pyspark import SparkContext, SparkConf
17
18 sc = SparkContext.getOrCreate()
19 spark = SparkSession.builder.getOrCreate()

```

Create an RDD from a text file

Each line of the text file becomes an element of the RDD.

```

1 !wget http://www.gutenberg.org/files/2600/2600-0.txt -O war_and_peace.txt
2 textFile = sc.textFile('war_and_peace.txt')

--2021-03-17 01:59:58-- http://www.gutenberg.org/files/2600/2600-0.txt
Resolving www.gutenberg.org (www.gutenberg.org)... 152.19.134.47, 2610:28:3090:3000:0:b:
Connecting to www.gutenberg.org (www.gutenberg.org)|152.19.134.47|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 3359584 (3.2M) [text/plain]
Saving to: 'war_and_peace.txt'

war_and_peace.txt  100%[=====>]   3.20M   5.01MB/s   in 0.6s

2021-03-17 01:59:59 (5.01 MB/s) - 'war_and_peace.txt' saved [3359584/3359584]

```

```

1 #One common transformation is
2 #filtering data that matches a predicate.
3 #We can use this to create a new RDD
4 #holding just the strings that contain
5 #the word Anna.
6
7 # The filter() transformation returns a new RDD
8 # containing only the elements that satisfy a predicate.
9 # A predicate is a function that returns True or False

```

```

10 # given an element of the RDD.
11 # The following function "lambda x: "Anna" in x",
12 # given an element x of the RDD, a line in this case,
13 # returns condition '"Anna" in x', which can be True or False.
14 annaLines = textFile.filter(lambda x: "Anna" in x)
15
16 #One example of an action is first()
17 #which returns the first element in an RDD.
18 firstLine = annaLines.first()
19
20 print(firstLine)
21
22 #Another example of action is collecting
23 #all the elements of an RDD.
24 allAnnaLines = annaLines.collect()
25
26 print(allAnnaLines)

```

```

    It was in July, 1805, and the speaker was the well-known Anna Pávlovna
    ['It was in July, 1805, and the speaker was the well-known Anna Pávlovna', 'rank and imp

```

```

1 #RDD<String>
2 N = textFile.filter(lambda x: "Natasha" in x)
3 #RDD<String>
4 P = textFile.filter(lambda x: "Pierre" in x)
5 #RDD<String>
6 A = textFile.filter(lambda x: "Andrew" in x)
7 #RDD<String>
8 NaP = N.intersection(P)
9 #RDD<String>
10 NaA = N.intersection(A)
11 #RDD<String>
12 NaPoNaA = NaP.union(NaA)
13
14 #These RDD have only been defined; they are not computed yet.
15 #They will only be computed if some action is performed as follows.
16
17 print(len(NaP.collect()))
18 print(len(NaA.collect()))
19
20 print(NaPoNaA.collect())

```

```

0
0
[]

```

```

1 #map() takes in a function and applies it to each element in the RDD
2 #with the result of the function being the new value of each element
3 #in the resulting RDD.

```

```

4

```

```

5 rdd = sc.parallelize([1, 2, 3, 4]);
6 result = rdd.map(lambda x: x*x);
7 print(result.collect());

```

```
[1, 4, 9, 16]
```

```

1 #Sometimes we want to produce multiple output elements for each input element.
2 #The operation to do this is called flatMap().
3 #As with map(), the function we provide to flatMap() is called individually
4 #for each element in the input RDD.
5 #Instead of returning a single element, we return in this function an iterator
6 #with our return values.
7 #Rather than producing an RDD of iterators, flatMap() gives back an RDD
8 #of the elements from all of the iterators.

```

```
9
```

```

10 #A simple usage of flatMap() is splitting up an input string into words.
11 #From each line, we want to output multiple words.

```

```
12
```

```
13 words = textFile.flatMap(lambda x: x.split());
```

```
14
```

```
15 print(words.collect()[0:100])
```

```
16 print(words.count())
```

```
['The', 'Project', 'Gutenberg', 'EBook', 'of', 'War', 'and', 'Peace,', 'by', 'Leo', 'To]
566316
```

```

1 #Suppose we would like to transform our string RDD of words
2 #to an RDD of the word lengths so that we can compute different stats with ease.

```

```
3
```

```
4 wordLength = words.map(lambda x: len(x));
```

```
5
```

```
6 #Then, we can compute different stats on it. E.g.
```

```
7
```

```
8 wordAvgLength = wordLength.mean();
```

```
9
```

```
10 print(wordAvgLength)
```

```
11
```

```
12 #and quite a few others (min, max, stdev, histograms, etc).
```

```
13 print(wordLength.max())
```

```
4.669703840258791
```

```
36
```

```

1 #The most common action on basic RDDs you will likely use is reduce(),
2 #which takes a function that operates on two elements of the type in your RDD
3 #and returns a new element of the same type.

```

```
4
```

```
5 #A simple example of such a function is +, which we can use to sum our RDD.
```

```
6 #With reduce(), we can easily sum the elements of our RDD.
```

```

6 #with reduce(), we can easily sum the elements of our RDD,
7 #count the number of elements, and perform other types of aggregations.
8
9 rdd = sc.parallelize([1, 2, 3, 4]);
10 sum = rdd.reduce(lambda x,y: x+y);
11 print(sum)

```

10

```

1 #reduce() requires that the return type of our result be the same type as that
2 #of the elements in the RDD we are operating over.
3 #This works well for operations like sum,
4 #but sometimes we want to return a different type.
5
6 #For example, when computing a running average,
7 #we need to keep track of both the count so far and the number of elements,
8 #which requires us to return a pair.
9 #We could work around this by first using map() where we transform every element
10 #into the element and the number 1, which is the type we want to return,
11 #so that the reduce() function can work on pairs.
12
13 rdd = sc.parallelize([1, 2, 3, 4])
14 sumcnt = rdd.map(lambda x: (x,1) ).reduce(lambda t,r: (t[0]+r[0], t[1]+r[1]) )
15
16 avg = sumcnt[0] / sumcnt[1]
17 print(avg)

```

2.5

```

1 #The aggregate() action frees us from the constraint of having the return
2 #be the same type as the RDD we are working on.
3 #With aggregate(), we supply:
4 #(1) An initial “zero” value of the type we want to return.
5 #(2) A function to combine the elements from our RDD with the “accumulator”.
6 #(3) A second function to “merge” two accumulators,
7 #   given that each machine accumulates its own results locally.
8
9 #We can use aggregate() to compute the average of an RDD,
10 #avoiding a map() before the reduce().
11
12 rdd = sc.parallelize([1, 2, 3, 4])
13 sumcnt = rdd.aggregate((0, 0),
14                         lambda acc, value: (acc[0] + value, acc[1] + 1),
15                         lambda acc1, acc2: (acc1[0] + acc2[0], acc1[1] + acc2[1]) )
16 avg = sumcnt[0] / sumcnt[1]
17 print(avg)

```

2.5

RDDs of key/value pairs

Spark provides operations on RDDs containing key/value pairs. These RDDs are called pair RDDs. Pair RDDs allow you to act on each key in parallel. For example, pair RDDs have a `reduceByKey()` method (analogous to `reduce` for regular RDDs) that can aggregate data separately for each key. We can create pair RDDs from existing RDDs. E.g.

```

1 import re
2 words = textFile.flatMap(lambda x: re.findall('\w+', x));
3
4 lw = words.map( lambda x: (len(x), x) );
5
6 # This creates an RDD of length-word pairs.
7 # What can we do with it?
8 # We can find for example the number of words for each length.
9
10 r = lw.countByKey();
11 print(r)
12
13 # Or, we can collect all the words of length >= 16.
14
15 longwordsRDD = lw.groupByKey().filter(lambda x: x[0] >= 16)
16
17 print(longwordsRDD.collect())
18
19 #What we get back is an object which allows iterating over the results.
20 #Turn the results of groupByKey into a list by calling list() on the values, e.g.
21
22 print(longwordsRDD.map(lambda x : (x[0], list(x[1]))).collect())

defaultdict(<class 'int'>, {3: 143657, 7: 42864, 9: 17486, 5: 58450, 2: 95460, 4: 98716,
[(16, <pyspark.resultiterable.ResultIterable object at 0x7f1ab6a402d0>), (18, <pyspark.r
[(16, ['enthusiastically', 'circumstantially', 'incomprehensible', 'misunderstanding',

```

Word count

```

1 textFile = sc.textFile('war_and_peace.txt')
2
3 word_counts = textFile.flatMap(lambda x: x.split()) \
4                          .map(lambda word: (word,1)) \
5                          .reduceByKey(lambda a,b: a+b)
6
7 print(word_counts.collect())
8
9 # Those familiar with the combiner concept from MapReduce should note that
10 # calling reduceByKey() will automatically perform combining locally
11 # on each machine before computing global totals for each key.
12 # The user does not need to specify a combiner.

```

```
[('The', 2550), ('Project', 78), ('EBook', 2), ('of', 14857), ('Peace,', 2), ('Leo', 4),
```

Word count with stopwords removed

```
1 !wget "https://gist.githubusercontent.com/sebleier/554280/raw/7e0e4a1ce04c2bb7bd41089c9821
2
3 textFile = sc.textFile('war_and_peace.txt')
4 stopwords = sc.textFile('stopwords.txt')
5
6 word_counts = textFile.flatMap(lambda x: x.split()) \
7                          .map(lambda word: (word.lower(),1)) \
8                          .subtractByKey(stopwords.map(lambda word: (word, 1))) \
9                          .reduceByKey(lambda a,b: a+b)
10
11 print(word_counts.collect())
```

```
--2021-03-17 02:00:14-- https://gist.githubusercontent.com/sebleier/554280/raw/7e0e4a1c
Resolving gist.githubusercontent.com (gist.githubusercontent.com)... 185.199.108.133, 18
Connecting to gist.githubusercontent.com (gist.githubusercontent.com)|185.199.108.133|:4
HTTP request sent, awaiting response... 200 OK
Length: 622 [text/plain]
Saving to: 'stopwords.txt'
```

```
stopwords.txt      100%[=====>]      622  --.-KB/s    in 0s
```

```
2021-03-17 02:00:14 (30.3 MB/s) - 'stopwords.txt' saved [622/622]
```

```
[('gutenberg', 24), ('leo', 4), ('tolstoy', 3), ('whatsoever.', 2), ('may', 260), ('it,
```

