

Word Count Example in Spark

The following slides show how to implement the MapReduce Word count algorithm, both in **Java** and **Python**

To understand the code, refer to

- The standard Java and Python APIs and reference manuals.
- Introduction to Programming in Spark (*for this course*)
- Spark RDD Programming guide.
- Spark Java API.
- Spark Python API.

All links and the complete code (**WordCountExample.java** and **WordCountExample.py**) can be found in the course Moodle.

OUTLINE

- ① Functional programming
- ② 1-Round Word Count in Spark
- ③ 2-Round Word Count in Spark

Functional programming

Functional programming

One of the core ideas of functional programming is that **functions can be arguments to other functions**.

- In a MapReduce algorithm, the map and reduce functions used in each round can be regarded as "arguments" for the round.
- Spark's API relies heavily on passing functions. For example, many methods (e.g., RDD operations) require functions as arguments.

In Spark there are essentially **2 ways to pass functions as arguments** to a method.

- as **anonymous functions** or **lambdas** which are defined inline without a name, where the actual argument is expected.
- as **named functions** defined outside the method.

Functional programming: Java example

Let **values** be a variable of type **JavaRDD<Double>**

denotes a generic element of the RDD values

USE of ANONYMOUS FUNCTIONS

- Single statement:

```
JavaRDD<Double> squaredValues = values.map((x) -> x*x);
```

- Multiple statements:

```
double fixed = 1.5;
JavaRDD<Double> normValues = values.map((x) -> {
    double diff = fixed - x;
    return diff;
});
```

Obs: variable `fixed` is "captured" by the anonymous function

Functional programming: Java example

USE of NAMED FUNCTIONS

Create a separate class (in a separate file) such as:

```
public class myOps {  
    public static double mySquare(double x) {  
        return x * x;  
    }  
}
```

Then, in the class containing the main program write:

```
JavaRDD<Double> squaredValues = values.map(myOps::mySquare);
```

Functional programming: Python example

Let `values` be an RDD of double

USE of ANONYMOUS FUNCTIONS

```
squaredValues = values.map(lambda x: x * x)
```

Obs: Python does not allow multiple statements in anonymous functions.

USE of NAMED FUNCTIONS

```
def mySquare(x):  
    return x * x
```

```
squaredValues = values.map(mySquare)
```

1-Round Word Count in Spark

JAVA

Java: initialization

```
import ....  
  
public class WordCountExample{  
  
    public static void main(String[] args) throws IOException {  
  
        // SPARK SETUP  
        SparkConf conf = new SparkConf(true).setAppName("WordCount");  
        JavaSparkContext sc = new JavaSparkContext(conf);  
  
        // INPUT READING  
        int K = Integer.parseInt(args[0]);  
        JavaRDD<String> docs = sc.textFile(args[1]).repartition(K).cache();  
  
        Allows you to pass some configuration properties  
        as command-line arguments  
        CLI arguments  
        path to file
```

SC ≡ Spark context[†]

repartition(k) ≡ performs a randomized partitioning of the elements of the RDD

cache() ≡ upon occurrence of the first action tries to store all data in main memory

docs ≡ RDD of strings where each element is a line of the input file, that is one of the input documents.

Java: implementation of the round

```
JavaPairRDD<String, Long> wordCounts = docs  
    .flatMapToPair(f1) // <-- MAP PHASE  
    .groupByKey()      // <-- SHUFFLING+GROUPING  
    .mapValues(f2);    // <-- REDUCE PHASE
```

f1: transforms each document (a string of the RDD docs) into a set of (word, count) pairs

f2: for each (word, listOfCounts) transforms the list of counts into their sum

Java: implementation of the round

```
JavaPairRDD<String, Long> wordCounts = docs
    .flatMapToPair((document) -> {      // <-- MAP PHASE
        String[] tokens = document.split(" ");
        HashMap<String, Long> counts = new HashMap<>();
        ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();
        for (String token : tokens) {
            counts.put(token, 1L + counts.getOrDefault(token, 0L));
        }
        for (Map.Entry<String, Long> e : counts.entrySet()) {
            pairs.add(new Tuple2<>(e.getKey(), e.getValue()));
        }
        return pairs.iterator();
    })
    .groupByKey()
    .mapValues((it) -> { // <-- REDUCE PHASE (R1)
        long sum = 0;
        for (long c : it) sum += c;
        return sum;
    })
}
```

} function f₁

} function f₂

* flatMapToPair : takes an input
element into 0, 1 or more key value pairs.

If instead you want to implement a 1-1
transformation ($1 \text{ element} \rightarrow 1 \text{ (key,value)pair}$)

then you can use mapToPair

The pairs returned by flatMapToPair are
returned as an iterator

* mapValues : it computes the following transformation of every element resulting from groupByKey()

$$(\text{word}, \underbrace{\text{list of counts}}_{f \text{ acts on}}) \rightarrow (\text{word}, \sum_{c \in \text{list of counts}} c)$$

Java: implementation of the round

Alternatively, the combination of `groupByKey` and `mapValues` can be replaced by `reduceByKey` as follows:

```
JavaPairRDD<String, Long> wordCounts = docs
    .flatMapToPair((document) -> {      // <-- MAP PHASE
        String[] tokens = document.split(" ");
        HashMap<String, Long> counts = new HashMap<>();
        ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();
        for (String token : tokens) {
            counts.put(token, 1L + counts.getOrDefault(token, 0L));
        }
        for (Map.Entry<String, Long> e : counts.entrySet()) {
            pairs.add(new Tuple2<>(e.getKey(), e.getValue()));
        }
        return pairs.iterator();
    })
    .reduceByKey((x, y) -> x+y);          // <-- REDUCE PHASE
```

↳ substitutes `groupByKey` and `mapValues`

Merge By Key

- * The argument is a function $D \times D \rightarrow D$
where D is the domain of the values of the pairs
The function must be COMMUTATIVE and ASSOCIATIVE
- * For each key k it combines the values of the pairs with key k as follows:

$$\begin{array}{c} (k, v_1) \\ (k, v_2) \\ (k, v_3) \\ (k, v_4) \end{array} \begin{array}{c} > \\ > \\ > \\ > \end{array} \begin{array}{c} (k, v_1 + v_2) \\ (k, v_3 + v_4) \end{array} \begin{array}{c} > \\ > \end{array} (k, \sum v_i)$$

The order in which the aggregation is executed is not specified but it EXPLOITS PARTITIONS

Java: implementation of the round

Also, one might define the function passed to flatMapToPair in a separate class MyMethods as follows:

```
public class MyMethods {  
    public static Iterator<Tuple2<String, Long>>  
        processDoc(String document) {  
        String[] tokens = document.split(" ");  
        HashMap<String, Long> counts = new HashMap<>();  
        ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();  
        for (String token : tokens) {  
            counts.put(token, 1L + counts.getOrDefault(token, 0L));  
        }  
        for (Map.Entry<String, Long> e : counts.entrySet()) {  
            pairs.add(new Tuple2<>(e.getKey(), e.getValue()));  
        }  
        return pairs.iterator();  
    }  
}
```

Java: implementation of the round

and use function `processDoc` in class `WordCountExample` as follows:

```
JavaPairRDD<String, Long> wordCounts = docs  
    .flatMapToPair(MyMethods::processDoc) // <-- MAP PHASE  
    .reduceByKey((x, y) -> x+y); // <-- REDUCE PHASE
```

PYTHON

Python: main function

```
def main():

# SPARK SETUP
conf = SparkConf(true).setAppName('WordCount').setMaster("local[*]")
sc = SparkContext(conf=conf)

# INPUT READING
K = sys.argv[1]
K = int(K)
data_path = sys.argv[2]
docs = sc.textFile(data_path).repartition(K).cache()

# COMPUTATION OF WORD COUNTS
wordCounts = word_count_1(docs)

if __name__ == "__main__":
    main()
```

⇒ docs is
an RDD of
strings

→ RDD of (word, count) pairs

`set Master("local [*]")` : configuration
for local mode . When you want
to execute the program on a cluster
you will write `yarn` instead of
`local [*]` (yarn = our cluster
manager)

`Word_Count_1` : function that
defines the transformations applied
to docs

Python: relevant functions

```
import ...  
  
def word_count_per_doc(document):  
    pairs_dict = {}  
    for word in document.split(' '):  
        if word not in pairs_dict.keys():  
            pairs_dict[word] = 1  
        else:  
            pairs_dict[word] += 1  
    return [(key, pairs_dict[key]) for key in pairs_dict.keys()]  
  
def word_count_1(docs):  
    word_count = (docs.flatMap(word_count_per_doc) # <-- MAP PHASE  
                  .groupByKey() # <-- SHUFFLE+GROUPING  
                  .mapValues(lambda vals: sum(vals))) # <-- REDUCE PHASE  
    return word_count
```

→ RDD of (word, count) pairs

Python: relevant functions

Alternatively:

```
import ...  
  
def word_count_per_doc(document):  
    pairs_dict = {}  
    for word in document.split(' '):  
        if word not in pairs_dict.keys():  
            pairs_dict[word] = 1  
        else:  
            pairs_dict[word] += 1  
    return [(key, pairs_dict[key]) for key in pairs_dict.keys()]  
  
def word_count_1(docs):  
    word_count = (docs.flatMap(word_count_per_doc) # <-- MAP PHASE  
                  .reduceByKey(lambda x, y: x + y)) # <-- REDUCE PHASE  
    return word_count
```

→ nuplaus groupByKey + mapValues

* Java vs Python

flatRepToPair → flatMap

mapToPair → map

* word-count-per-doc : is applied
to each document to produce
(word, count) pairs (partial counts)

2-Round Word Count in Spark

2-Round Word Count in Spark

We will see 3 alternative implementations:

- ① Implementation based on Improved word count assigning random keys *explicitly*
- ② Implementation based on Improved word count assigning random keys *on the fly*
- ③ Implementation based on Spark partitions.

JAVA

Java: improved word count (explicit random keys)

```
JavaPairRDD<String, Long> wordCounts = docs → RDD of strings
    .flatMapToPair((document) -> { // <-- MAP PHASE (R1)
        String[] tokens = document.split(" ");
        HashMap<String, Long> counts = new HashMap<>();
        ArrayList<Tuple2<Integer, Tuple2<String, Long>>>
            pairs = new ArrayList<>();
        for (String token : tokens) {
            counts.put(token, 1L + counts.getOrDefault(token, 0L));
        }
        for (Map.Entry<String, Long> e : counts.entrySet()) {
            pairs.add(new Tuple2<>(randomGenerator.nextInt(K), ← random
                [ new Tuple2<>(e.getKey(), e.getValue()));           Key
            })
        }
        return pairs.iterator(); → (word, count)
    })
    .groupByKey()          // <-- SHUFFLE AND GROUPING → pair
```

grouping based on random keys

(Key, list of (word, count) pairs)

```
.flatMapToPair((element) -> { // <-- REDUCE PHASE (R1)
    HashMap<String, Long> counts = new HashMap<>();
    for (Tuple2<String, Long> c : element._2()) {
        counts.put(c._1(), c._2() + counts.getOrDefault(c._1(), 0L));
    }
    ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();
    for (Map.Entry<String, Long> e : counts.entrySet()) {
        pairs.add(new Tuple2<>(e.getKey(), e.getValue()));
    }
    return pairs.iterator();
})
.reduceByKey((x, y) -> x+y); // <-- REDUCE PHASE (R2)
```

First flatMapToPair $\xrightarrow{\text{removes key in } [\ell, k-1]}$

* Document $\rightarrow \{(key, (word, count))\}$

groupByKey : creates (key, list of (word, count) pairs)
pairs when keys are the removable keys

Second flatMapToPair

* removes the removable keys

* aggregates counts of each word limited
to the occurrences in the group

reduceByKey : compute final counts

Java: improved word count (random keys on-the-fly)

```
JavaPairRDD<String, Long> wordCounts = docs
    .flatMapToPair((document) -> { // <-- MAP PHASE (R1)
        String[] tokens = document.split(" ");
        HashMap<String, Long> counts = new HashMap<>();
        ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();
        for (String token : tokens) {
            counts.put(token, 1L + counts.getOrDefault(token, 0L));
        }
        for (Map.Entry<String, Long> e : counts.entrySet()) {
            pairs.add(new Tuple2<>(e.getKey(), e.getValue()));
        }
        return pairs.iterator();
    })
    .groupByKey() // <-- GROUP BY (R2)
    .mapPartitions((iter) -> {
        Random randomGenerator = new Random();
        return iter.map((key, value) -> {
            int randomKey = randomGenerator.nextInt(K);
            return new Tuple2<String, Long>(randomKey, value);
        });
    })
    .reduceByKey((key1, key2) -> key1 + key2); // <-- REDUCE PHASE (R3)
```

As in the round trip alg.

assignment of random keys on the fly

(key, list of (word, count) pairs)
Key is reusable

```
.flatMapToPair((element) -> { // <-- REDUCE PHASE (R1)
    //
    // AS BEFORE
    //
})
.reduceByKey((x, y) -> x+y); // <-- REDUCE PHASE (R2)
```

`groupBy(...)`

- * Assigns a random key to each (word, count) pair produced by the first flatMapToPair
- * groups data into (key, list of (word, count) pairs) pairs , when key is random.

Java: use of Spark partitions

```
JavaPairRDD<String, Long> wordCounts = docs
    .flatMapToPair((document) -> {      // <-- MAP PHASE (R1)
        //
        // SAME AS IN THE 1-ROUND ALGORITHM
        //
    })
    .mapPartitionsToPair((element) -> {      // <-- REDUCE PHASE (R1)
        HashMap<String, Long> counts = new HashMap<>();
        while (element.hasNext()){
            Tuple2<String, Long> tuple = element.next();
            counts.put(tuple._1(), tuple._2() +
                counts.getOrDefault(tuple._1(), 0L));
        }
        ArrayList<Tuple2<String, Long>> pairs = new ArrayList<>();
        for (Map.Entry<String, Long> e : counts.entrySet()) {
            pairs.add(new Tuple2<>(e.getKey(), e.getValue()));
        }
        return pairs.iterator();
    })
}

iterator of (word, count) pairs from one partition
```

(word, count)

```
.groupByKey()          // <-- SHUFFLE+GROUPING
.mapValues((it) -> { // <-- REDUCE PHASE (R2)
    long sum = 0;
    for (long c : it) sum += c;
    return sum;
});
// reduceByKey CAN BE USED IN PLACE OF groupByKey AND mapValues
```

`mapPartitionsToPair` exploits the Spark partitions



- * processes each partition "element", that is a list of (word, count) pairs provided as an iterator
- * creates (word, count) pairs when counts are aggregated at partition level

PYTHON

Python: relevant functions

```
def word_count_per_doc(document, K=-1):
    pairs_dict = {}
    for word in document.split(' '):
        if word not in pairs_dict.keys():
            pairs_dict[word] = 1
        else:
            pairs_dict[word] += 1
        if K == -1:
            return [(key, pairs_dict[key]) for key in pairs_dict.keys()]
    else:
        return [(rand.randint(0,K-1),(key, pairs_dict[key]))  

                for key in pairs_dict.keys()]
```

$K = -1 \quad \text{Doc} \rightarrow \{(word, count)\}$

$K \neq -1 \quad \text{Doc} \rightarrow \{(key, (word, count))\}$

random in $[0, K-1]$

Python: relevant functions

```
def gather_pairs(pairs):
    pairs_dict = {}
    for p in pairs[1]:
        word, occurrences = p[0], p[1]
        if word not in pairs_dict.keys():
            pairs_dict[word] = occurrences
        else:
            pairs_dict[word] += occurrences
    return [(key, pairs_dict[key]) for key in pairs_dict.keys()]
```

(key, list of (word, count) pairs)

aggregates counts for the same word

Python: improved word count (explicit random keys)

RDD & strings element of docs

```
def word_count_2(docs, K):
    word_count = (docs.flatMap(lambda x: word_count_per_doc(x, K))
                  # <-- MAP PHASE (R1)
                  .groupByKey()                      # <-- SHUFFLE+GROUPING
                  .flatMap(gather_pairs)            # <-- REDUCE PHASE (R1)
                  .reduceByKey(lambda x, y: x + y)) # <-- REDUCE PHASE (R2)
    return word_count
```

After first flatMap : RDD whose elements are pairs $(\text{key}, (\text{word}, \text{count}))$ $\text{key} \in [0, k-1]$
random

After groupByKey : RDD whose elements are pairs $(\text{key}, \text{list of } (\text{word}, \text{count}) \text{ pairs})$

After second flatMap : RDD whose elements are pairs $(\text{word}, \text{count})$ (at most k pairs for every word)

After reduceByKey : RDD whose elements are pairs $(\text{word}, \text{count})$ 1 per word

Python: improved word count (random keys on-the-fly)

(word, count)

parameter K = -1

```
def word_count_3(docs, K):
    word_count = (docs.flatMap(word_count_per_doc) # <-- MAP PHASE (R1)
                  .groupBy(lambda x: (rand.randint(0,K-1))) # <-- SHUFFLE+GROUPING
                  .flatMap(gather_pairs)                      # <-- REDUCE PHASE (R1)
                  .reduceByKey(lambda x, y: x + y))          # <-- REDUCE PHASE (R2)
    return word_count
```

{
as before

Python: use of Spark partitions

```
def gather_pairs_partitions(pairs):  
    pairs_dict = {}  
    for p in pairs:  
        word, occurrences = p[0], p[1]  
        if word not in pairs_dict.keys():  
            pairs_dict[word] = occurrences  
        else:  
            pairs_dict[word] += occurrences  
    return [(key, pairs_dict[key]) for key in pairs_dict.keys()]
```

```
def word_count_with_partition(docs):  
    word_count = (docs.flatMap(word_count_per_doc) # <-- MAP PHASE (R1)  
                  .mapPartitions(gather_pairs_partitions) # <-- REDUCE PHASE (R1)  
                  .groupByKey() # <-- SHUFFLE+GROUPING  
                  .mapValues(lambda vals: sum(vals))) # <-- REDUCE PHASE (R2)  
    return word_count
```

or alternatively, reduceByKey

$K = -1$

