1. **Apache Spark Professional Training with Hands On Lab Sessions**
2. **Oreilly Databricks Apache Spark Developer Certification Simulator**

# SPARK PAIR RDD FUNCTIONS

By www.HadoopExam.com

**Note: These instructions should be used with the HadoopExam Apache Spark: Professional Trainings. Where it is executed and you can do hands on with trainer.**

**Cloudera CCA175 (Hadoop and Spark Developer Hands-on Certification available with total 75 solved problem scenarios. Click for More Detail)**

1. reduceByKey()
2. groupByKey()
3. combineByKey()
4. foldByKey()
5. aggregateByKey()
6. Comparision Between Function

reduce Bykey()

K1, (1,2) ——— reduceBykey ((K1,(1,2)),(K1,(1,7)))

K1, (1,7)                                    ↑           ↑
                                             X           Y
k2, (1,8)

k2, (2,9)        ① K1 = [(1+1),(2+7)]

k7, (1,8)           K1 = [(2,9)]
K1, (2,6)

Part-1         ② K1  reduceBykey((K1,(2,9),(K1,(2,6)))

                                        ↑              ↑
                                        X              Y

k3,(1,4)           K1 = [(2+2),(9+6)]
k4, (1,4)
                      = [(4,15)
k5, (1,6)

k3, (1,9)

k3,(1,6)        P1 ⟹ [(K1,(4,15)),(k2,(3,17)),(k7,(2,6))

k4, (2,7)
                P2⟹ [(k3,(3,19)),(k4,(3,11)),(k5,(1,6))
Part-2
                P3 ⟹ [(k6,(2,12)),(k9(2,14))]

k6,(1,4)              Locally each Partition

k9,(1,7)

k6,(1,8)

k9,(1,7)

dataRDD

dataRDD. reduceBykey ((x,y) ⟹

                      (x.-1+y.-1, x.-2+y.-2))

# groupByKey()

Val words = Array ("one", "two", "two", hadoop, hadoop, hadoop)

Val generatePairRDD = se.parallelize (words).map( word => (word, 1))

⇓

( one, 1), (two, 1), (two, 1), (hadoop, 1), (hadoop, 1), (hadoop, 1)

Val wordCountWithReduce

= generatePairRDD. reduceByKey(_ + _). collect ()

(one, 1), => (one, 1)

(two, 1), (two, 1) => two =>(1+1) => (two, 2)

(hadoop, 1), (hadoop, 1), (hadoop, 1)

↳ (hadoop, (1+1)), (hadoop, 1)

↳ (hadoop, 3)

Val wordWithGroup = generatePairRDD. groupByKey ()

. map (t => (t._1, t._2.sum)

. collect ()

$$
\begin{array}{|l}
(one, 1) \\
(two, (1,1)) \\
(hadoop, (1,1,1))
\end{array}
$$

=> .map

(one, 1)

(two => (1,1).sum)

(hadoop => (1,1,1).sum)

⇓

(one, 1)

(two, 2)

(hadoop, 3)

=) Both will produce same result.

=) reduceByKey works better ( local Aggregations)

=) graepByKey => first shuffle and then aggregate
( Hence, huge new traffic involved, which will cause
performance issue)

=> ~~$th~~ Explaination to below image: —

> ┌─────────────────────────────────┐
> │ ReduceBykey  v/s   GroupByKey. │
> └─────────────────────────────────┘

Here are more functions to prefer over groupByKey()

CombineByKey(): — Can be used when you are combining
elements but your return type is different
from your input.

foldBykey():— merges the values for each key using an
associative function and a neutral "zero value"

# foldBykey() & fold()

## fold(): Example.

$$sc.\ parallelize(1\ to\ 10)$$
$$.fold(0)\ \{\ (acc,\ element) =>$$

$$acc + element)$$

first pass $\longrightarrow$ [(0,1) $\Rightarrow$ (0+1)] $\Rightarrow$ 1

2nd pass $\longrightarrow$ [(1,2) $\Rightarrow$ (1+2)] $\Rightarrow$ 3

3rd pass $\longrightarrow$ [(3,3) $\Rightarrow$ (3+3)] $\Rightarrow$ 6

4th pass $\longrightarrow$ [(6,4) $\Rightarrow$ (6+4)] $\Rightarrow$ 10

$\vdots$

1th pass $\longrightarrow$ [(45,10) $\Rightarrow$ (45+10)] $\Rightarrow$ 55

Another Example: Count element in list

$$sc.\ parallelize\ (1\ to\ 10).\ fold(0)\{(acc,\ element) =>$$

$$acc + 1\}$$

output would be = 10

## foldByKey() :-

```
Val depEmployee = List (
                        ( CS, ( Amit, 1000)),
                        ( CS, (Rahul, 1200)),
                        ( ECE, (Rakesh, 1500)),
                        ( ECE, (Ankit, 1200))
                    )

Val empRDD =  sc. makeRDD ( depEmployee)
// find max score by dept

Val maxdept = empRDD. foldByKey (("dummy", 0.0) )
                ((acc, element) => if (acc.-2 > element.-2)
                                    acc
                                    else
                                    element )
```

for CS

firstPass

CS ⟶ [( dummy, 0.0), (Amit, 1000)]
        => ( Amit, 1000)

CS ⟶ [ (Amit, 1000), (Rahul, 1200)]
        => (Rahul, 1200)

ECE ⟶ [(dummy, 0.0), (Rakesh, 1500)]
        => (Rakesh, 1500)

ECE ⟶ [(Rakesh, 1500), (Ankit, 1200)]
        => (Rakesh, 1500)

maxdept => [ (CS, (Rahul, 1200)), (ECE, (Rakesh, 1500))]

CombineByKey():-

CombineByKey ( CreateCombiner,
                MergeValue,
                mergeCombiners,
                partitioner )

Example Code : Calculate Running Avarage

val result = inputRDD.combineByKey (

①————————.(V) => (V,1)

② ←———————— (acc : (Int, Int), V) => (acc._1 + V, acc._2 + 1)

③ ←——— (acc1 : (Int, Int), acc2 : (Int, Int)
       => ( acc._1 + acc2._1 , acc1._2 + acc2._2 ))

      • map { case (key, value) => (key, value._1/

                                    value._2.tofloat)



InputRDD          P1
                 ┌──────────┐      => for k1
                 │ (k1, 12) │      Part1 first time  k1 → (12,1)
                 │ (k2, 13) │      first time        k2 → (13,1)
Part-1 ────→     │ (~~~~~~) │      2nd time          k1 → ((12,1),16)
                 │ (k1, 16) │                          → (28,2)
                 │ (k1, 12) │      3rd time           k1 → ((28,2),12)
                 │ (k2, 13) │                          → (~~36,3~~) (40,3)
                 │ (k2, 14) │
Part-2 ────      │ (k3,11)  │      Part-2 for k1
                 │ (k3, 12) │        first time  k1 → (12,1)
                 │ (k3,11)  │        2nd time    k1 → ((12,1),14)
                 │ k2,14    │                      → (26,2)
                 │ k1, 12   │
                 │ k1,14    │      finally operation3 (Using all Accumulator)
                 └──────────┘      across partition ((36,3)(26,2)) => (62,5)

Similarly it applies to all keys (k1, k2, k3 ---etc.)
across all partitions.

$$k1 \rightarrow 62/5$$

$$k2 \rightarrow ((26,2), (28,2)) \Rightarrow (54,4)$$
$$\Rightarrow 54/4$$

$$k3 \rightarrow ((34,3)) \Rightarrow 34/3$$

~~result = [ k1, 12.4~~

result = [ (k1, 12.4), (k2, 13.5), (k3, 11.33) ]

Running Average for all the keys.

Aggregate ByKey() :- This function also requires three
parameters.

① An initial zero value, which will not affect the total
values to be collected.

   e.g. Summation = 0 + 5 = 5
                    0 + 6 = 6   etc
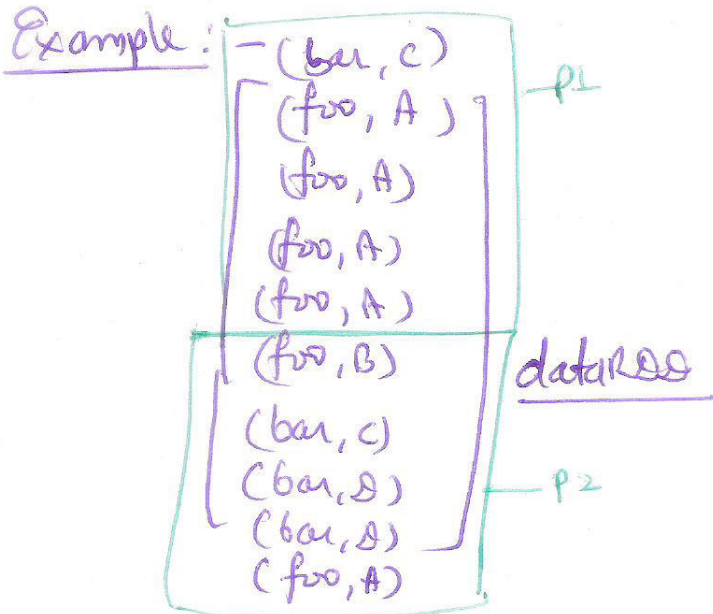
   Collecting unique elements from set. - Empty set
            empty Set.add ("hadoop") ⟹ empty Set [hadoop]

② A combiner function, which accept two parameters
   ⟹ Second parameter will be merged into the first
      parameter
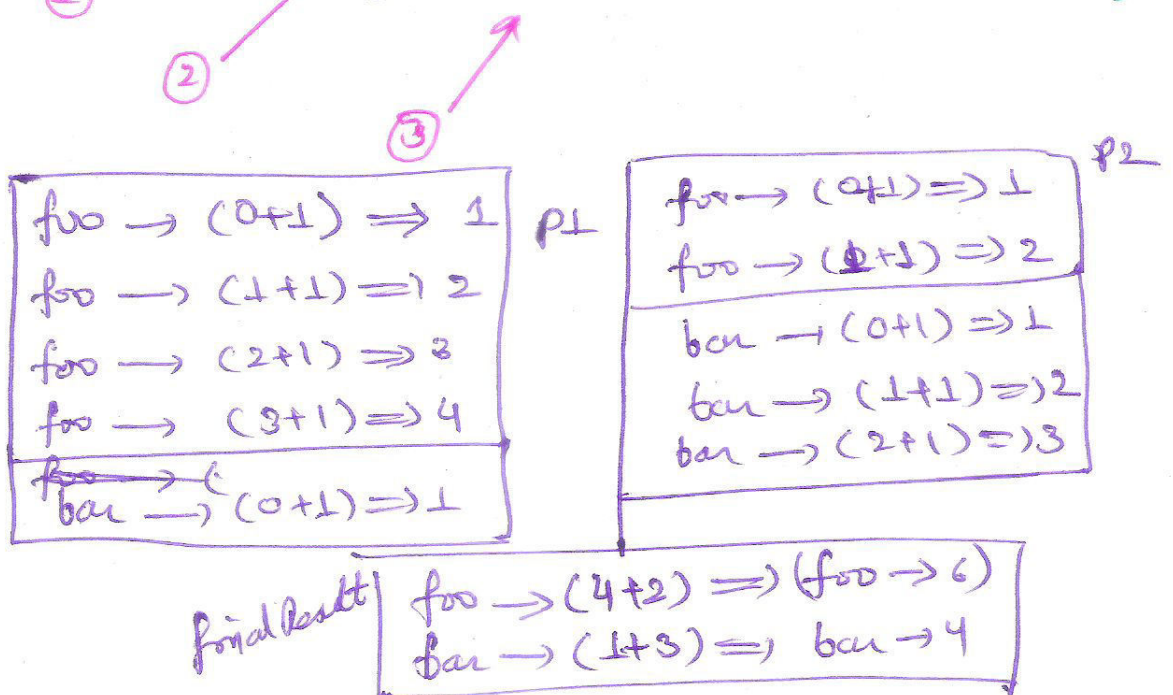
⇒ Combining function will work on local partition only.

③ Merging function = Works across the partition.

It also accepts two parameters.

Example :
$$
\begin{array}{l}
- (bar, c) \\
(foo, A) \\
(foo, A) \\
(foo, A) \\
(foo, A) \\
(foo, B) \\
(bar, c) \\
(bar, 8) \\
(bar, 8) \\
(foo, A)
\end{array}
$$

P1

dataRDD

P2

dataRDD . aggregateByKey ( 0 ) ( ( A : Int, V : String ) ⇒ A+1 ),

(P1

dataRDD . aggregateByKey ( 0 ) ( ( acc : Int, V : String ) ⇒ acc+1 ),

( ( acc : Int, acc : Int ) ⇒ acc+acc ) )

①
②
③

P1
$$
\begin{array}{l}
foo \rightarrow (0+1) \Rightarrow 1 \\
foo \rightarrow (1+1) \Rightarrow 2 \\
foo \rightarrow (2+1) \Rightarrow 3 \\
foo \rightarrow (3+1) \Rightarrow 4 \\
\cancel{foo \rightarrow (} \\
bar \rightarrow (0+1) \Rightarrow 1
\end{array}
$$

P2
$$
\begin{array}{l}
foo \rightarrow (0+1) \Rightarrow 1 \\
foo \rightarrow (1+1) \Rightarrow 2 \\
bar \rightarrow (0+1) \Rightarrow 1 \\
bar \rightarrow (1+1) \Rightarrow 2 \\
bar \rightarrow (2+1) \Rightarrow 3
\end{array}
$$

final Result
$$
\begin{array}{l}
foo \rightarrow (4+2) \Rightarrow (foo \rightarrow 6) \\
bar \rightarrow (1+3) \Rightarrow bar \rightarrow 4
\end{array}
$$

# function Comparision =>

① You can replace groupByKey() with reduceByKey() to improve performance.

② reduceByKey() performs map side combine which can reduce n/w IO and shuffle size.

③ groupByKey() will not perform map side combine.

④ CombineByKey() is more general than aggregateByKey()

⑤ Implementation of aggregateByKey, reduceByKey, and groupByKey is achieved by combineByKey().

⑥ AggregateByKey() is similar to reduceByKey(), but you can provide initial values when performing aggregation.

=> As name suggests, aggregateByKey is suitable for compute aggregations for keys such as sum, avg. etc.

=> CombineByKey(): is more general and you have the flexibility to specify whether you would like to perform map side combine

=> However, eg CombineByKey() is more complex, at the minimum you need to implement three functions,
  ① create Combiner
  ② merge Value
  ③ merge Combiners.

# cloudera CCA175   Is Now Available , with Hands On Sessions

## CCA Spark and Hadoop Developer Certification

HadoopExam Learning Resource provides the following material for the Advanced Technologies.
Please visit www.HadoopExam.com for more detail this is just a few products from portfolio.

**Price start for training with Just $79/3500INR**

| | | | | |
|---|---|---|---|---|
| **Spark** | **hadoop** | **APACHE HBASE** | **ANDROID** | **ORACLE Java** |
| Apache Spark Professional Training with HandsOn Session **+ Certification Material** | Hadoop Professional Training with HandsOn Sessioon **+ Certification Material** | HBase Professional Training with HandsOn Session **+ Certification Material** | **Certification Material** | **Certification Material** |
| **amazon** webservices | **S.sas** | **EMC²** | **Microsoft Azure** | **EMC² Data Scientist** |
| **Certification Material** | **Certification Material** | **Certification Material** | **Certification Material** | **Certification Material** |