

# Spark Solution for Rank Product

Mahmoud Parsian  
Ph.D in Computer Science

Senior Architect @ illumina<sup>1</sup>

August 25, 2015

---

<sup>1</sup>[www.illumina.com](http://www.illumina.com)

# Table of Contents

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions
- 4 Input Data
- 5 Rank Product Algorithm
- 6 Moral of Story

# Outline

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions
- 4 Input Data
- 5 Rank Product Algorithm
- 6 Moral of Story

# Who am I?

- Name: Mahmoud Parsian
- Education: Ph.D in Computer Science
- Work: Senior Architect @Illumina, Inc
  - Lead Big Data Team @Illumina
  - Develop scalable regression algorithms
  - Develop DNA-Seq and RNA-Seq workflows
  - Use Java/MapReduce/Hadoop/Spark/HBase
- Author: of 3 books
  - Data Algorithms (O'Reilly:  
<http://shop.oreilly.com/product/0636920033950.do/>)
  - JDBC Recipes (Apress: <http://apress.com/>)
  - JDBC MetaData Recipes (Apress: <http://apress.com/>)

# Outline

- 1 Biography
- 2 What is Rank Product?**
- 3 Basic Definitions
- 4 Input Data
- 5 Rank Product Algorithm
- 6 Moral of Story

# Problem Statement: Data Algorithms Book

- Bonus Chapter: <http://shop.oreilly.com/product/0636920033950.do>
- Source code: <https://github.com/mahmoudparsian/data-algorithms-book>



# Problem Statement

- Given a large set of studies, where each study is a set of (Key, Value) pairs
- The goal is to find the Rank Product of all given Keys for all studies.
- Magnitude of this data is challenging to store and analyze:
  - Several hundreds of studies
  - Each study has billions of (Key, Value) pairs
  - Find "rank product" for all given Keys

# Problem Statement

- Given a large set of studies, where each study is a set of (Key, Value) pairs
- The goal is to find the Rank Product of all given Keys for all studies.
- Magnitude of this data is challenging to store and analyze:
  - Several hundreds of studies
  - Each study has billions of (Key, Value) pairs
  - Find "rank product" for all given Keys



# Magnitude of Data per Analysis

- 100's of studies
- Each study may have billions of  $(K, V)$  pairs
- Example: 300 studies
- Each study: 2000,000,000  $(K, V)$  pairs
- Analyze: 600,000,000,000  $(K, V)$  pairs

# Magnitude of Data per Analysis

- 100's of studies
- Each study may have billions of  $(K, V)$  pairs
- Example: 300 studies
- Each study: 2000,000,000  $(K, V)$  pairs
- Analyze: 600,000,000,000  $(K, V)$  pairs

# Outline

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions**
- 4 Input Data
- 5 Rank Product Algorithm
- 6 Moral of Story

# Some Basic Definitions

- Study is a set of  $(K, V)$  pairs
- Example-1:  $K=\text{geneID}$ ,  $V=\text{geneValue}$
- Example-2:  $K=\text{userID}$ ,  $V=\text{user's followers in social networks}$
- Example-3:  $K=\text{bookID}$ ,  $V=\text{book rating by users}$

# Some Basic Definitions

- Study is a set of  $(K, V)$  pairs
- Example-1:  $K=\text{geneID}$ ,  $V=\text{geneValue}$
- Example-2:  $K=\text{userID}$ ,  $V=\text{user's followers in social networks}$
- Example-3:  $K=\text{bookID}$ ,  $V=\text{book rating by users}$

# What is a Ranking?

- Let  $S = \{(K_1, 40), (K_2, 70), (K_3, 90), (K_4, 80)\}$
- Then  $\text{Rank}(S) = \{(K_1, 4), (K_2, 3), (K_3, 1), (K_4, 2)\}$
- Since  $90 > 80 > 70 > 40$
- Ranks are assigned as: 1, 2, 3, 4, ..., N

# What is Rank Product?

- Let  $\{A_1, \dots, A_k\}$  be a set of (key-value) pairs where keys are unique per dataset.
- Example of (key-value) pairs:
  - $(K,V) = (\text{item}, \text{number of items sold})$
  - $(K,V) = (\text{user}, \text{number of followers for the user})$
  - $(K,V) = (\text{gene}, \text{test expression})$
- Then the ranked product of  $\{A_1, \dots, A_k\}$  is computed based on the ranks  $r_i$  for key  $i$  across all  $k$  datasets. Typically ranks are assigned based on the sorted values of datasets.

# What is Rank Product?

- Let  $A_1 = \{(K_1, 30), (K_2, 60), (K_3, 10), (K_4, 80)\}$ ,  
then  $Rank(A_1) = \{(K_1, 3), (K_2, 2), (K_3, 4), (K_4, 1)\}$   
since  $80 > 60 > 30 > 10$   
Note that 1 is the highest rank (assigned to the largest value).
- Let  $A_2 = \{(K_1, 90), (K_2, 70), (K_3, 40), (K_4, 50)\}$ ,  
 $Rank(A_2) = \{(K_1, 1), (K_2, 2), (K_3, 4), (K_4, 3)\}$   
since  $90 > 70 > 50 > 40$
- Let  $A_3 = \{(K_1, 4), (K_2, 8)\}$   
 $Rank(A_3) = \{(K_1, 2), (K_2, 1)\}$   
since  $8 > 4$

The **rank product** of  $\{A_1, A_2, A_3\}$  is expressed as:

$$\{(K_1, \sqrt[3]{3 \times 1 \times 2}), (K_2, \sqrt[3]{2 \times 2 \times 1}), (K_3, \sqrt[2]{4 \times 4}), (K_4, \sqrt[2]{1 \times 3})\}$$



# Calculation of the Rank Product

- Given  $n$  genes and  $k$  replicates,
- Let  $e_{g,i}$  be the fold change and  $r_{g,i}$  the rank of gene  $g$  in the  $i$ 'th replicate.
- Compute the rank product (RP) via the geometric mean:

$$RP(g) = \left( \prod_{i=1}^k r_{g,i} \right)^{1/k}$$

$$RP(g) = \sqrt[k]{\left( \prod_{i=1}^k r_{g,i} \right)}$$

# Outline

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions
- 4 Input Data**
- 5 Rank Product Algorithm
- 6 Moral of Story

# Input Data Format

- Set of  $k$  studies  $\{S_1, S_2, \dots, S_k\}$
- Each study has billions of (Key, Value) pairs
- Sample Record:  
`<key-as-string><,><value-as-double-data-type>`

# Input Data Persistence

- Data persists in HDFS
- Directory structure:

```
/input/study-001/file-001-1.txt  
/input/study-001/file-001-2.txt  
/input/study-001/file-001-3.txt  
...  
/input/study-235/file-235-1.txt  
/input/study-235/file-235-2.txt  
/input/study-235/file-235-3.txt  
...
```

# Formalizing Rank Product

- Let  $S = \{S_1, S_2, \dots, S_k\}$  be a set of  $k$  studies, where  $k > 0$  and each study represent a micro-array experiment
- Let  $S_i (i = 1, 2, \dots, k)$  be a study, which has an arbitrary number of assays identified by  $\{A_{i1}, A_{i2}, \dots\}$
- Let each assay (can be represented as a text file) be a set of arbitrary number of records in the following format:  
`<gene_id><,><gene_value_as_double_data-type>`
- Let `gene_id` be in  $\{g_1, g_2, \dots, g_n\}$  (we have  $n$  genes).

# Rank Product: in 2 Steps

Let  $S = \{S_1, S_2, \dots, S_k\}$  be a set of  $k$  studies:

- STEP-1: find the mean of values per study per gene
  - you may replace the "mean" function by your desired function
  - finding mean involves `groupByKey()` or `combineByKey()`
- STEP-2: perform the "Rank Product" per gene across all studies
  - finding "Rank Product" involves `groupByKey()` or `combineByKey()`

# Formalizing Rank Product

The last step will be to find the rank product for each gene per study:

$$S_1 = \{(g_1, r_{11}), (g_2, r_{12}), \dots\}$$

$$S_2 = \{(g_1, r_{21}), (g_2, r_{22}), \dots\}$$

...

$$S_k = \{(g_1, r_{k1}), (g_2, r_{k2}), \dots\}$$

then Ranked Product of  $g_j =$

$$RP(g_j) = \left( \prod_{i=1}^k r_{i,j} \right)^{1/k}$$

or

$$RP(g_j) = \sqrt[k]{\left( \prod_{i=1}^k r_{i,j} \right)}$$

# Spark Solution for Rank Product

- 1 Read  $k$  input paths (each path represents a study, which may have any number of assay text files).
- 2 Find the mean per gene per study
- 3 Sort the genes by value per study and then assign rank values; To sort the dataset by value, we will swap the key with value and then perform the sort.
- 4 Assign ranks from 1, 2, ...,  $N$  (1 is assigned to the highest value) we use `JavaPairRDD.zipWithIndex()`, which zips the RDD with its element indices (these indices will be the ranks). Spark indices will start from 0, we will add 1 when computing the ranked product.
- 5 Finally compute the Rank Product per gene for all studies. This can be accomplished by grouping all ranks by the key (we may use `JavaPairRDD.groupByKey()` or `JavaPairRDD.combineByKey()` – note that, in general, `JavaPairRDD.combineByKey()` is more efficient than `JavaPairRDD.groupByKey()`).



# Spark Solution for Rank Product

- 1 Read  $k$  input paths (each path represents a study, which may have any number of assay text files).
- 2 Find the mean per gene per study
- 3 Sort the genes by value per study and then assign rank values; To sort the dataset by value, we will swap the key with value and then perform the sort.
- 4 Assign ranks from 1, 2, ...,  $N$  (1 is assigned to the highest value) we use `JavaPairRDD.zipWithIndex()`, which zips the RDD with its element indices (these indices will be the ranks). Spark indices will start from 0, we will add 1 when computing the ranked product.
- 5 Finally compute the Rank Product per gene for all studies. This can be accomplished by grouping all ranks by the key (we may use `JavaPairRDD.groupByKey()` or `JavaPairRDD.combineByKey()` – note that, in general, `JavaPairRDD.combineByKey()` is more efficient than `JavaPairRDD.groupByKey()`).

# Spark Solution for Rank Product

- 1 Read  $k$  input paths (each path represents a study, which may have any number of assay text files).
- 2 Find the mean per gene per study
- 3 Sort the genes by value per study and then assign rank values; To sort the dataset by value, we will swap the key with value and then perform the sort.
- 4 Assign ranks from 1, 2, ...,  $N$  (1 is assigned to the highest value) we use `JavaPairRDD.zipWithIndex()`, which zips the RDD with its element indices (these indices will be the ranks). Spark indices will start from 0, we will add 1 when computing the ranked product.
- 5 Finally compute the Rank Product per gene for all studies. This can be accomplished by grouping all ranks by the key (we may use `JavaPairRDD.groupByKey()` or `JavaPairRDD.combineByKey()` – note that, in general, `JavaPairRDD.combineByKey()` is more efficient than `JavaPairRDD.groupByKey()`).

# Spark Solution for Rank Product

- 1 Read  $k$  input paths (each path represents a study, which may have any number of assay text files).
- 2 Find the mean per gene per study
- 3 Sort the genes by value per study and then assign rank values; To sort the dataset by value, we will swap the key with value and then perform the sort.
- 4 Assign ranks from 1, 2, ...,  $N$  (1 is assigned to the highest value) we use `JavaPairRDD.zipWithIndex()`, which zips the RDD with its element indices (these indices will be the ranks). Spark indices will start from 0, we will add 1 when computing the ranked product.
- 5 Finally compute the Rank Product per gene for all studies. This can be accomplished by grouping all ranks by the key (we may use `JavaPairRDD.groupByKey()` or `JavaPairRDD.combineByKey()` – note that, in general, `JavaPairRDD.combineByKey()` is more efficient than `JavaPairRDD.groupByKey()`).

# Spark Solution for Rank Product

- 1 Read  $k$  input paths (each path represents a study, which may have any number of assay text files).
- 2 Find the mean per gene per study
- 3 Sort the genes by value per study and then assign rank values; To sort the dataset by value, we will swap the key with value and then perform the sort.
- 4 Assign ranks from 1, 2, ...,  $N$  (1 is assigned to the highest value) we use `JavaPairRDD.zipWithIndex()`, which zips the RDD with its element indices (these indices will be the ranks). Spark indices will start from 0, we will add 1 when computing the ranked product.
- 5 Finally compute the Rank Product per gene for all studies. This can be accomplished by grouping all ranks by the key (we may use `JavaPairRDD.groupByKey()` or `JavaPairRDD.combineByKey()` – note that, in general, `JavaPairRDD.combineByKey()` is more efficient than `JavaPairRDD.groupByKey()`).

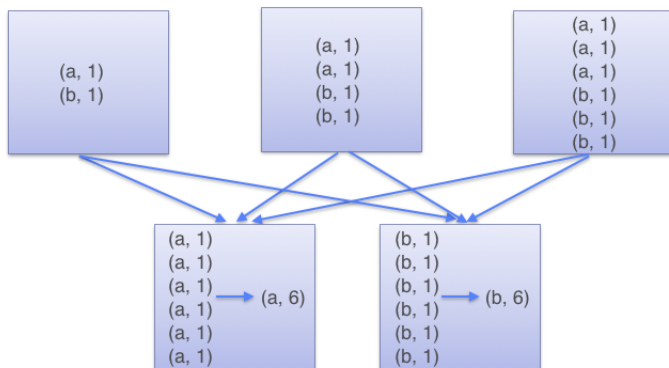
## Two Spark Solutions: `groupByKey()` and `combineByKey()`

Two solutions are provided using Spark-1.4.0:

- `SparkRankProductUsingGroupByKey`
  - uses `groupByKey()`
- `SparkRankProductUsingCombineByKey`
  - uses `combineByKey()`

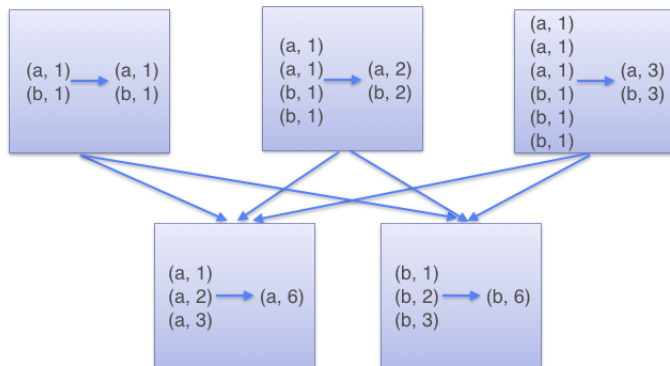
# How does groupByKey() work

## GroupByKey



# How does reduceByKey() work

## ReduceByKey



# Outline

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions
- 4 Input Data
- 5 Rank Product Algorithm**
- 6 Moral of Story



# Rank Product Algorithm in Spark

Algorithm: High-Level Steps	
Step	Description
STEP-1	import required interfaces and classes
STEP-2	handle input parameters
STEP-3	create a Spark context object
STEP-4	create list of studies (1, 2, ..., K)
STEP-5	compute mean per gene per study
STEP-6	sort by values
STEP-7	assign rank
STEP-8	compute rank product
STEP-9	save the result in HDFS

# Main Driver

## Listing 1: performRrankProduct()

```

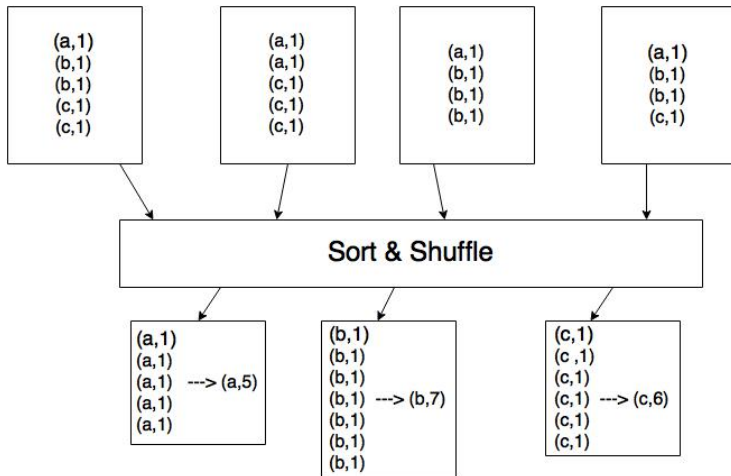
1  public static void main(String[] args) throws Exception {
2      // args[0] = output path
3      // args[1] = number of studies (K)
4      // args[2] = input path for study 1
5      // args[3] = input path for study 2
6      // ...
7      // args[K+1] = input path for study K
8      final String outputPath = args[0];
9      final String numOfStudiesAsString = args[1];
10     final int K = Integer.parseInt(numOfStudiesAsString);
11     List<String> inputPathMultipleStudies = new ArrayList<String>();
12     for (int i=1; i <= K; i++) {
13         String singleStudyInputPath = args[1+i];
14         inputPathMultipleStudies.add(singleStudyInputPath);
15     }
16     performRrankProduct(inputPathMultipleStudies, outputPath);
17     System.exit(0);
18 }

```

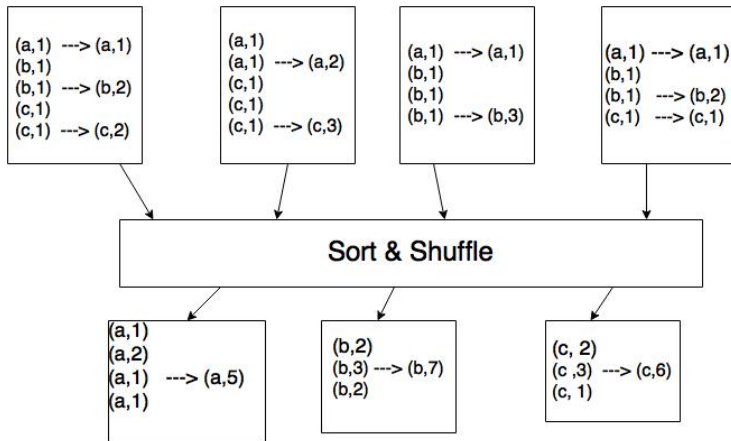
## groupByKey() vs. combineByKey()

- Which one should we use? `combineByKey()` or `groupByKey()`? According to their semantics, both will give you the same answer. But `combineByKey()` is more efficient.
- In some situations, `groupByKey()` can even cause out of disk problems. In general, `reduceByKey()`, and `combineByKey()` are preferred over `groupByKey()`.
- Spark shuffling is more efficient for `reduceByKey()` than `groupByKey()` and the reason is this: in the shuffle step for `reduceByKey()`, data is combined so each partition outputs at most one value for each key to send over the network, while in shuffle step for `groupByKey()`, all the data is wastefully sent over the network and collected on the reduce workers.
- To understand the difference, the following figures show how the shuffle is done for `reduceByKey()` and `groupByKey()`

# Understanding groupByKey()



# Understanding reduceByKey() or combineByKey()



## Listing 2: performRrankProduct()

```

1 public static void performRrankProduct(
2     final List<String> inputPathMultipleStudies,
3     final String outputPath) throws Exception {
4     // create a context object, which is used
5     // as a factory for creating new RDDs
6     JavaSparkContext context = Util.createJavaSparkContext(useYARN);
7
8     // Spark 1.4.0 requires an array for creating union of many RDDs
9     int index = 0;
10    JavaPairRDD<String, Double>[] means =
11        new JavaPairRDD[inputPathMultipleStudies.size()];
12    for (String inputPathSingleStudy : inputPathMultipleStudies) {
13        means[index] = computeMeanByGroupByKey(context, inputPathSingleStudy);
14        index++;
15    }
16
17    // next compute rank
18    ...
19 }

```

## Listing 3: performRrankProduct()

```

1 public static void performRrankProduct(
2     ...
3     // next compute rank
4     // 1. sort values based on absolute value of mean value
5     // 2. assign rank from 1 to N
6     // 3. calculate rank product for each gene
7     JavaPairRDD<String,Long>[] ranks = new JavaPairRDD[means.length];
8     for (int i=0; i < means.length; i++) {
9         ranks[i] = assignRank(means[i]);
10    }
11    // calculate ranked products
12    // <gene, T2<rankedProduct, N>>
13    JavaPairRDD<String, Tuple2<Double, Integer>> rankedProducts =
14        computeRankedProducts(context, ranks);
15
16    // save the result, shuffle=true
17    rankedProducts.coalesce(1,true).saveAsTextFile(outputPath);
18
19    // close the context and we are done
20    context.close();
21 }

```

## STEP-3: create a Spark context object

```

1 public static JavaSparkContext createJavaSparkContext(boolean useYARN) {
2     JavaSparkContext context;
3     if (useYARN) {
4         context = new JavaSparkContext("yarn-cluster", "MyAnalysis"); // YARN
5     }
6     else {
7         context = new JavaSparkContext(); // Spark cluster
8     }
9     // inject efficiency
10    SparkConf sparkConf = context.getConf();
11    sparkConf.set("spark.kryoserializer.buffer.mb", "32");
12    sparkConf.set("spark.shuffle.file.buffer.kb", "64");
13    // set a fast serializer
14    sparkConf.set("spark.serializer",
15                  "org.apache.spark.serializer.KryoSerializer");
16    sparkConf.set("spark.kryo.registrator",
17                  "org.apache.spark.serializer.KryoRegistrator");
18    return context;
19 }

```

---



## Listing 4: STEP-4: compute mean

```

1 static JavaPairRDD<String, Double> computeMeanByGroupByKey(
2     JavaSparkContext context,
3     final String inputPath) throws Exception {
4     JavaPairRDD<String, Double> genes =
5         getGenesUsingTextFile(context, inputPath, 30);
6
7     // group values by gene
8     JavaPairRDD<String, Iterable<Double>> groupedByGene = genes.groupByKey();
9
10    // calculate mean per gene
11    JavaPairRDD<String, Double> meanRDD = groupedByGene.mapValues(
12        new Function<
13            Iterable<Double>,      // input
14            Double                 // output: mean
15        >() {
16            @Override
17            public Double call(Iterable<Double> values) {
18                double sum = 0.0;
19                int count = 0;
20                for (Double v : values) {
21                    sum += v;
22                    count++;
23                }
24                // calculate mean of samples
25                double mean = sum / ((double) count);
26                return mean;
27            }
28        });
29    return meanRDD;
30 }

```

## Listing 5: getGenesUsingTextFile()

```

1 static JavaPairRDD<String, Double> getGenesUsingTextFile(
2     JavaSparkContext context,
3     final String inputPath,
4     final int numberOfPartitions) throws Exception {
5
6     // read input and create the first RDD
7     // JavaRDD<String>: where String = "gene,test_expression"
8     JavaRDD<String> records = context.textFile(inputPath, numberOfPartitions);
9
10    // for each record, we emit (K=gene, V=test_expression)
11    JavaPairRDD<String, Double> genes
12        = records.mapToPair(new PairFunction<String, String, Double>() {
13        @Override
14        public Tuple2<String, Double> call(String rec) {
15            // rec = "gene,test_expression"
16            String[] tokens = StringUtils.split(rec, ",");
17            // tokens[0] = gene
18            // tokens[1] = test_expression
19            return new Tuple2<String, Double>(
20                tokens[0], Double.parseDouble(tokens[1]));
21        }
22    });
23    return genes;
24 }

```

## Listing 6: Assign Rank

```

1 // result is JavaPairRDD<String, Long> = (gene, rank)
2 static JavaPairRDD<String,Long> assignRank(JavaPairRDD<String,Double> rdd){
3     // swap key and value (will be used for sorting by key); convert value to abs(value)
4     JavaPairRDD<Double,String> swappedRDD = rdd.mapToPair(
5         new PairFunction<Tuple2<String, Double>,           // T: input
6             Double,                                         // K
7             String>>(){                                     // V
8                 public Tuple2<Double,String> call(Tuple2<String, Double> s) {
9                     return new Tuple2<Double,String>(Math.abs(s._2), s._1);
10                }
11            });
12     // we need 1 partition so that we can zip numbers into this RDD by zipWithIndex()
13     JavaPairRDD<Double,String> sorted = swappedRDD.sortByKey(false, 1); // sort means descending
14     // JavaPairRDD<T,Long> zipWithIndex()
15     // Long values will be 0, 1, 2, ...; for ranking, we need 1, 2, 3, ..., therefore, we will add 1
16     JavaPairRDD<Tuple2<Double,String>,Long> indexed = sorted.zipWithIndex();
17     // next convert JavaPairRDD<Tuple2<Double,String>,Long> into JavaPairRDD<String,Long>
18     // JavaPairRDD<Tuple2<value, gene>,rank> into JavaPairRDD<gene,rank>
19     JavaPairRDD<String, Long> ranked = indexed.mapToPair(
20         new PairFunction<Tuple2<Tuple2<Double,String>,Long>, // T: input
21             String,                                           // K: mapped_id
22             Long>(){                                           // V: rank
23                 public Tuple2<String, Long> call(Tuple2<Tuple2<Double,String>,Long> s) {
24                     return new Tuple2<String,Long>(s._1._2, s._2 + 1); // ranks are 1, 2, ..., n
25                 }
26            });
27     return ranked;
28 }

```

## Listing 7: Compute Rank Product using groupByKey()

```

1 static JavaPairRDD<String, Tuple2<Double, Integer>> computeRankedProducts(
2     JavaSparkContext context,
3     JavaPairRDD<String, Long>[] ranks) {
4     JavaPairRDD<String, Long> unionRDD = context.union(ranks);
5
6     // next find unique keys, with their associated values
7     JavaPairRDD<String, Iterable<Long>> groupedByGeneRDD = unionRDD.groupByKey();
8
9     // next calculate ranked products and the number of elements
10    JavaPairRDD<String, Tuple2<Double, Integer>> rankedProducts = groupedByGeneRDD.mapValues(
11        new Function<
12            Iterable<Long>,           // input: means for all studies
13            Tuple2<Double, Integer> // output: (rankedProduct, N)
14        >() {
15        @Override
16        public Tuple2<Double, Integer> call(Iterable<Long> values) {
17            int N = 0;
18            long products = 1;
19            for (Long v : values) {
20                products *= v;
21                N++;
22            }
23            double rankedProduct = Math.pow((double) products, 1.0/((double) N));
24            return new Tuple2<Double, Integer>(rankedProduct, N);
25        }
26    });
27    return rankedProducts;
28 }

```

## Next FOCUS on combineByKey()

We do need to develop 2 functions:

- `computeMeanByCombineByKey()`
- `computeRankedProductsUsingCombineByKey()`

## combineByKey(): how does it work?

- `combineByKey()` is the most general of the per-key aggregation functions. Most of the other per-key combiners are implemented using it.
- Like `aggregate()`, `combineByKey()` allows the user to return values that are not the same type as our input data.
- To understand `combineByKey()`, it is useful to think of how it handles each element it processes.
- As `combineByKey()` goes through the elements in a partition, each element either has a key it has not seen before or has the same key as a previous element.

# combineByKey() definition

```

1 public <C> JavaPairRDD<K,C> combineByKey(
2   Function<V,C> createCombiner,      // V -> C
3   Function2<C,V,C> mergeValue,      // C+V -> C
4   Function2<C,C,C> mergeCombiners   // C+C -> C
5 )
6 Description: Generic function to combine the elements for each key
7 using a custom set of aggregation functions. Turns a JavaPairRDD[(K, V)]
8 into a result of type JavaPairRDD[(K, C)], for a "combined type" C. Note
9 that V and C can be different -- for example, one might group an RDD of
10 type (Int, Int) into an RDD of type (Int, List[Int]).
11 Users provide three functions:
12 - createCombiner, which turns a V into a C (e.g., creates a one-element list)
13 - mergeValue, to merge a V into a C (e.g., adds it to the end of a list)
14 - mergeCombiners, to combine two Cs into a single one.

```

---

# computeMeanByCombineByKey():

## Define C data structure

```
1  /**
2   * AverageCount is used by combineByKey() to hold
3   * the total values and their count.
4   */
5  static class AverageCount implements Serializable {
6      double total;
7      int count;
8
9      public AverageCount(double total, int count) {
10         this.total = total;
11         this.count = count;
12     }
13
14     public double average() {
15         return total / (double) count;
16     }
17 }
```



# computeMeanByCombineByKey()

```

1 static JavaPairRDD<String, Double> computeMeanByCombineByKey(
2     JavaSparkContext context,
3     final String inputPath) throws Exception {
4
5     JavaPairRDD<String, Double> genes = getGenesUsingTextFile(context, inputPath, 30);
6
7     // we need 3 function to be able to use combineByKey()
8     Function<Double, AverageCount> createCombiner = ...
9     Function2<AverageCount, Double, AverageCount> addAndCount = ...
10    Function2<AverageCount, AverageCount, AverageCount> mergeCombiners = ...
11
12    JavaPairRDD<String, AverageCount> averageCounts =
13        genes.combineByKey(createCombiner, addAndCount, mergeCombiners);
14
15    // now compute the mean/average per gene
16    JavaPairRDD<String, Double> meanRDD = averageCounts.mapToPair(
17        new PairFunction<
18            Tuple2<String, AverageCount>, // T: input
19            String,                       // K
20            Double,                       // V
21            >() {
22
23        @Override
24        public Tuple2<String, Double> call(Tuple2<String, AverageCount> s) {
25            return new Tuple2<String, Double>(s._1, s._2.average());
26        }
27    });
28    return meanRDD;
29 }

```

# 3 functions for combineByKey()

```

1  Function<Double, AverageCount> createCombiner = new Function<Double, AverageCount>() {
2      @Override
3      public AverageCount call(Double x) {
4          return new AverageCount(x, 1);
5      }
6  };
7
8  Function2<AverageCount, Double, AverageCount> addAndCount =
9      new Function2<AverageCount, Double, AverageCount>() {
10     @Override
11     public AverageCount call(AverageCount a, Double x) {
12         a.total += x;
13         a.count += 1;
14         return a;
15     }
16 };
17
18 Function2<AverageCount, AverageCount, AverageCount> mergeCombiners =
19     new Function2<AverageCount, AverageCount, AverageCount>() {
20     @Override
21     public AverageCount call(AverageCount a, AverageCount b) {
22         a.total += b.total;
23         a.count += b.count;
24         return a;
25     }
26 };

```

# Per-key average using combineByKey() in Python

```
1 sumCount = nums.combineByKey(  
2     (lambda x: (x,1)),  
3     (lambda x, y: (x[0] + y, x[1] + 1)),  
4     (lambda x, y: (x[0] + y[0], x[1] + y[1])))  
5  
6  
7 mean = sumCount.map(lambda key, xy: (key, xy[0]/xy[1])).collectAsMap()  
8 print(*mean)
```

---

# Sample Run using groupByKey(): script

```

1 # define input/output for Hadoop/HDFS
2 ## args[0] = output path
3 ## args[1] = number of studies (K)
4 ## args[2] = input path for study 1
5 ## args[3] = input path for study 2
6 ## ...
7 ## args[K+1] = input path for study K
8 OUTPUT=/rankproduct/output
9 NUM_OF_INPUT=3
10 INPUT1=/rankproduct/input1
11 INPUT2=/rankproduct/input2
12 INPUT3=/rankproduct/input3
13
14 # remove all files under input
15 $HADOOP_HOME/bin/hadoop fs -rmr $OUTPUT
16 #
17 # remove all files under output
18 driver=org.dataalgorithms.bonus.rankproduct.spark.SparkRankProductUsingGroupByKey
19 $SPARK_HOME/bin/spark-submit --class $driver \
20     --master yarn-cluster \
21     --jars $OTHER_JARS \
22     --conf "spark.yarn.jar=$SPARK_JAR" \
23     $APP_JAR $OUTPUT $NUM_OF_INPUT $INPUT1 $INPUT2 $INPUT3

```

## Sample Run using groupByKey(): input

```
1 # hadoop fs -cat /rankproduct/input1/rp1.txt
2 K_1,30.0
3 K_2,60.0
4 K_3,10.0
5 K_4,80.0
6
7 # hadoop fs -cat /rankproduct/input2/rp2.txt
8 K_1,90.0
9 K_2,70.0
10 K_3,40.0
11 K_4,50.0
12
13 # hadoop fs -cat /rankproduct/input3/rp3.txt
14 K_1,4.0
15 K_2,8.0
```

---

## Sample Run using groupByKey(): output

```
1 # hadoop fs -cat /rankproduct/output/part*  
2 (K_2,(1.5874010519681994,3))  
3 (K_3,(4.0,2))  
4 (K_1,(1.8171205928321397,3))  
5 (K_4,(1.7320508075688772,2))
```

---

# computeRankedProductsUsingCombineByKey(): Define C data structure

```
1  /**
2   * RankProduct is used by combineByKey() to hold
3   * the total product values and their count.
4   */
5  static class RankProduct implements Serializable {
6      long product;
7      int count;
8
9      public RankProduct(long product, int count) {
10         this.product = product;
11         this.count = count;
12     }
13
14     public double rank() {
15         return Math.pow((double) product, 1.0/ (double) count);
16     }
17 }
```

# computeRankedProductsUsingCombineByKey()

```

1 // JavaPairRDD<String, Tuple2<Double, Integer>> = <gene, T2(rankedProduct, N)>
2 // where N is the number of elements for computing the rankedProduct
3 static JavaPairRDD<String, Tuple2<Double, Integer>> computeRankedProductsUsingCombineByKey(
4     JavaSparkContext context,
5     JavaPairRDD<String, Long>[] ranks) {
6     JavaPairRDD<String, Long> unionRDD = context.union(ranks);
7
8     // we need 3 function to be able to use combineByKey()
9     Function<Long, RankProduct> createCombiner = ...
10    Function2<RankProduct, Long, RankProduct> addAndCount = ...
11    Function2<RankProduct, RankProduct, RankProduct> mergeCombiners = ...
12
13    // next find unique keys, with their associated copa scores
14    JavaPairRDD<String, RankProduct> combinedByGeneRDD =
15        unionRDD.combineByKey(createCombiner, addAndCount, mergeCombiners);
16
17    // next calculate ranked products and the number of elements
18    JavaPairRDD<String, Tuple2<Double, Integer>> rankedProducts = combinedByGeneRDD.mapValues(
19        new Function<
20            RankProduct,                               // input: RankProduct
21            Tuple2<Double, Integer>                       // output: (RankedProduct.count)
22        >() {
23
24            @Override
25            public Tuple2<Double, Integer> call(RankProduct value) {
26                double theRankedProduct = value.rank();
27                return new Tuple2<Double, Integer>(theRankedProduct, value.count);
28            }
29        });
30    return rankedProducts;
31 }

```



# 3 functions for computeRankedProductsUsingCombineByKey()

```

1  Function<Long, RankProduct> createCombiner = new Function<Long, RankProduct>() {
2      @Override
3      public RankProduct call(Long x) {
4          return new RankProduct(x, 1);
5      }
6  };
7
8  Function2<RankProduct, Long, RankProduct> addAndCount =
9      new Function2<RankProduct, Long, RankProduct>() {
10     @Override
11     public RankProduct call(RankProduct a, Long x) {
12         a.product *= x;
13         a.count += 1;
14         return a;
15     }
16 };
17
18 Function2<RankProduct, RankProduct, RankProduct> mergeCombiners =
19     new Function2<RankProduct, RankProduct, RankProduct>() {
20     @Override
21     public RankProduct call(RankProduct a, RankProduct b) {
22         a.product *= b.product;
23         a.count += b.count;
24         return a;
25     }
26 };

```

# REVISED computeRankedProductsUsingCombineByKey():

## Define C data structure

```
1  static class RankProduct implements Serializable {
2      long product;
3      int count;
4
5      public RankProduct(long product, int count) {
6          this.product = product;
7          this.count = count;
8      }
9
10     public product(long value) {
11         this.product *= value;
12         this.count++;
13     }
14
15     public product(RankProduct pr) {
16         this.product *= pr.value;
17         this.count += pr.count;
18     }
19
20     public double rank() {
21         return Math.pow((double) product, 1.0/ (double) count);
22     }
23 }
```

# REVISED 3 functions for computeRankedProductsUsingCombineByKey()

```

1  Function<Long, RankProduct> createCombiner = new Function<Long, RankProduct>() {
2      public RankProduct call(Long x) {
3          return new RankProduct(x, 1);
4      }
5  };
6
7  Function2<RankProduct, Long, RankProduct> addAndCount =
8      new Function2<RankProduct, Long, RankProduct>() {
9      public RankProduct call(RankProduct a, Long x) {
10         a.product(x);
11         return a;
12     }
13 };
14
15 Function2<RankProduct, RankProduct, RankProduct> mergeCombiners =
16     new Function2<RankProduct, RankProduct, RankProduct>() {
17     public RankProduct call(RankProduct a, RankProduct b) {
18         a.product(b);
19         return a;
20     }
21 };

```

---

# Outline

- 1 Biography
- 2 What is Rank Product?
- 3 Basic Definitions
- 4 Input Data
- 5 Rank Product Algorithm
- 6 Moral of Story**

# Moral of Story...

- Understand data requirements
- Select proper platform for implementation: Spark
- Partition your RDDs properly
  - number of cluster nodes
  - number of cores per node
  - amount of RAM
- Avoid unnecessary computations  $(g1, g2)$ ,  $(g2, g1)$
- Use `filter()` often to remove non-needed RDD elements
- Avoid unnecessary `RDD.saveAsTextFile(path)`
- Run both on YARN and Spark cluster and compare performance
- use `groupByKey()` cautiously
- use `combineByKey()` over `groupByKey()`
- Verify your test results against R language

# Thank you!

Questions?