

DS256:Jan17 (3:1)

Spark Programming with RDD

Yogesh Simmhan

21 Mar, 2017

<http://spark.apache.org/docs/latest/programming-guide.html>

<http://spark.apache.org/docs/latest/api/java/index.html?org/apache/spark/api/java/JavaRDD.html>

<http://spark.apache.org/docs/latest/api/java/index.html?org/apache/spark/api/java/JavaPairRDD.html>

©Yogesh Simmhan & Partha Talukdar, 2016

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/)

Copyright for external content used with attribution is retained by their original authors





Spark Language

- Much more flexible than MapReduce
- Compose complex dataflows
- Data (RDD) centric... ~object oriented



Creating RDD

- Load external data from distributed storage
- Create logical RDD on which you can operate
- Support for different input formats
 - HDFS files, Cassandra, Java serialized, directory, gzipped
- Can control the number of partitions in loaded RDD
 - Default depends on external DFS, e.g. 128MB on HDFS

```
JavaRDD<String> distFile = sc.textFile("data.txt");
```



RDD Operations

■ Transformations

- From one RDD to one or more RDDs
- Lazy evaluation...*use with care*
- Executed in a distributed manner

■ Actions

- Perform aggregations on RDD items
- Return single (or distributed) results to “driver” code

- `RDD.collect()` brings RDD partitions to single driver machine



Anonymous Classes

- Data-centric model allows functions to be passed
 - Functions applied to items in the RDD
 - Typically, on individual partitions in data-parallel
- Anonymous class implements interface

```
JavaRDD<String> lines = sc.textFile("data.txt");
JavaRDD<Integer> lineLengths = lines.map(new Function<String, Integer>() {
    public Integer call(String s) { return s.length(); }
});
int totalLength = lineLengths.reduce(new Function2<Integer, Integer, Integer>() {
    public Integer call(Integer a, Integer b) { return a + b; }
});
```

```
class GetLength implements Function<String, Integer> {
    public Integer call(String s) { return s.length(); }
}
class Sum implements Function2<Integer, Integer, Integer> {
    public Integer call(Integer a, Integer b) { return a + b; }
}
```

```
JavaRDD<String> lines = sc.textFile("data.txt");
JavaRDD<Integer> lineLengths = lines.map(new GetLength());
int totalLength = lineLengths.reduce(new Sum());
```



Anonymous Classes & Lambda Expressions

- Or Java 8 functions are short-forms for simple code fragments to iterate over collections

```
JavaRDD<String> lines = sc.textFile("data.txt");
JavaRDD<Integer> lineLengths = lines.map(s -> s.length());
int totalLength = lineLengths.reduce((a, b) -> a + b);
```

- Caution: Cannot pass “local” driver variables to lambda expressions/anonymous classes....only **final**
 - Will fail when distributed

```
int counter = 0;
JavaRDD<Integer> rdd = sc.parallelize(data);

// Wrong: Don't do this!!
rdd.foreach(x -> counter += x);

println("Counter value: " + counter);
```



RDD and PairRDD

- RDD is logically a collection of items with a generic type
- PairRDD is like a “Map”, where each item in collection is a <key,value> pair, each a generic type
- Transformation functions use RDD or PairRDD as input/output
- E.g. Map-Reduce

```
JavaRDD<String> lines = sc.textFile("data.txt");  
JavaPairRDD<String, Integer> pairs = lines.mapToPair(s -> new Tuple2(s, 1));  
JavaPairRDD<String, Integer> counts = pairs.reduceByKey((a, b) -> a + b);
```



Transformations

Transformation	Meaning
<code>map(func)</code>	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<code>filter(func)</code>	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<code>flatMap(func)</code>	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).

- **JavaRDD<R> map(Function<T,R> f) : 1:1** mapping from input to output. Can be different types.
- **JavaRDD<T> filter(Function<T,Boolean> f) : 1:0/1** from input to output, same type.
- **JavaRDD<U> flatMap(FlatMapFunction<T,U> f) : 1:N** mapping from input to output, different types.



Transformations

mapPartitions(*func*)

Similar to map, but runs separately on each partition (block) of the RDD, so *func* must be of type `Iterator<T> => Iterator<U>` when running on an RDD of type T.

- Earlier Map and Filter operate on one item at a time. No state across calls!
- **JavaRDD<U>**
mapPartitions(FlatMapFunc<Iterator<T>,U> f)
- mapPartitions has access to iterator of values in entire partition, not just a single item at a time.



Transformations

sample(*withReplacement, fraction, seed*)

Sample a fraction *fraction* of the data, with or without replacement, using a given random number generator *seed*.

union(*otherDataset*)

Return a new dataset that contains the union of the elements in the source dataset and the argument.

- **JavaRDD<T> sample(boolean withReplacement, double fraction)**: fraction between [0,1] without replacement, >0 with replacement
- **JavaRDD<T> union(JavaRDD<T> other)**: Items in other RDD added to this RDD. Same type. Can have duplicate items (i.e. not a 'set' union).



Transformations

intersection(*otherDataset*)

Return a new RDD that contains the intersection of elements in the source dataset and the argument.

distinct(*[numTasks]*)

Return a new dataset that contains the distinct elements of the source dataset.

- **JavaRDD<T> intersection(JavaRDD<T> other):** Does a set intersection of the RDDs. Output *will not* have duplicates, even if inputs did.
- **JavaRDD<T> distinct():** Returns a new RDD with unique elements, eliminating duplicates.



Transformations: PairRDD

`groupByKey([numTasks])`

When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs.

Note: If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using `reduceByKey` or `aggregateByKey` will yield much better performance.

Note: By default, the level of parallelism in the output depends on the number of partitions of the parent RDD. You can pass an optional `numTasks` argument to set a different number of tasks.

`reduceByKey(func, [numTasks])`

When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function `func`, which must be of type `(V,V) => V`. Like in `groupByKey`, the number of reduce tasks is configurable through an optional second argument.

- **`JavaPairRDD<K,Iterable<V>> groupByKey()`**: Groups values for each key into a single iterable.
- **`JavaPairRDD<K,V> reduceByKey(Function2<V,V,V> func)`** : Merge the values for each key into a single value using an associative and commutative reduce function. Output value is of same type as input.
- **For aggregate that returns a different type?**
- **`numPartitions`** can be used to generate output RDD with different number of partitions than input RDD.



Transformations

aggregateByKey(zeroValue)(seqOp, combOp, [numTasks])

When called on a dataset of (K, V) pairs, returns a dataset of (K, U) pairs where the values for each key are aggregated using the given combine functions and a neutral "zero" value. Allows an aggregated value type that is different than the input value type, while avoiding unnecessary allocations. Like in `groupByKey`, the number of reduce tasks is configurable through an optional second argument.

sortByKey([ascending], [numTasks])

When called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean `ascending` argument.

- **JavaPairRDD<K,U> aggregateByKey(U zeroValue, Function2<U,V,U> seqFunc, Function2<U,U,U> combFunc) :** Aggregate the values of each key, using given combine functions and a neutral "zero value".
 - **SeqOp** for merging a V into a U within a partition
 - **CombOp** for merging two U's, within/across partitions
- **JavaPairRDD<K,V> sortByKey(Comparator<K> comp) :** Global sort of the RDD by key
 - Each partition contains a sorted range, i.e., output RDD is range-partitioned.
 - Calling `collect` will return an ordered list of records



Transformations

join(*otherDataset*, [*numTasks*])

When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key. Outer joins are supported through `leftouterJoin`, `rightouterJoin`, and `fullouterJoin`.

cartesian(*otherDataset*)

When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).

- **JavaPairRDD<K, Tuple2<V,W>>**
join(JavaPairRDD<K,W> other, int numParts):
Matches keys in *this* and *other*. Each output pair is (k, (v1, v2)). Performs a hash join across the cluster.
- **JavaPairRDD<T,U> cartesian(JavaRDDLike<U,?> other):** Cross product of values in each RDD as a pair



Actions

reduce(<i>func</i>)	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
collect()	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
count()	Return the number of elements in the dataset.
first()	Return the first element of the dataset (similar to <code>take(1)</code>).
take(<i>n</i>)	Return an array with the first <i>n</i> elements of the dataset.



RDD Persistence & Caching

- RDDs can be reused in a dataflow
 - Branch, iteration
- But it will be re-evaluated each time it is reused!
- Explicitly persist RDD to reuse output of a dataflow path multiple times
- Multiple storage levels for persistence
 - Disk or memory
 - Serialized or object form in memory
 - Partial spill-to-disk possible
 - *Cache* indicates “persist” to memory



RePartitioning

repartition

```
public JavaRDD<T> repartition(int numPartitions)
```

Return a new RDD that has exactly `numPartitions` partitions.

Can increase or decrease the level of parallelism in this RDD. Internally, this uses a shuffle to redistribute data.

If you are decreasing the number of partitions in this RDD, consider using `coalesce`, which can avoid performing a shuffle.

coalesce

```
public JavaRDD<T> coalesce(int numPartitions,  
                           boolean shuffle)
```

Return a new RDD that is reduced into `numPartitions` partitions.



From DAG to RDD lineage

