

# SPARK RDD Cheat Sheet

## Transformations:

Transformations are kind of operations which will transform your RDD data from one form to another. And when you apply this operation on any RDD, you will get a new RDD with transformed data (RDDs in Spark are immutable). Operations like map, filter, flatMap are transformations.

Now there is a point to be noted here and that is when you apply the transformation on any RDD it will not perform the operation immediately. It will create a DAG(Directed Acyclic Graph) using the applied operation, source RDD and function used for transformation. And it will keep on building this graph using the references till you apply any action operation on the last lined up RDD. That is why the transformation in Spark are lazy.

Spark has certain operations which can be performed on RDD. An operation is a method, which can be applied on a RDD to accomplish certain task. RDD supports two types of operations, which are Action and Transformation. An operation can be something as simple as sorting, filtering and summarizing data.

## Actions

Transformations create RDDs from each other, but when we want to work with the actual dataset, at that point action is performed. When the action is triggered after the result, new RDD is not formed like transformation. Thus, actions are RDD operations that give non-RDD values. The values of action are stored to drivers or to the external storage system. It brings laziness of RDD into motion.

Spark drivers and external storage system store the value of action. It brings laziness of RDD into motion.

An action is one of the ways of sending data from Executer to the driver. Executors are agents that are responsible for executing a task. While the driver is a JVM process that coordinates workers and execution of the task.

With RDD's we have access to a large collection of useful transformations and actions. Here are some of the more important ones:

Transformation	Meaning
<b>map</b> ( <i>func</i> )	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<b>filter</b> ( <i>func</i> )	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<b>flatMap</b> ( <i>func</i> )	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).
<b>sample</b> ( <i>withReplacement</i> , <i>fraction</i> , <i>seed</i> )	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator seed.
<b>union</b> ( <i>otherDataset</i> )	Return a new dataset that contains the union of the elements in the source dataset and the argument.
<b>intersection</b> ( <i>otherDataset</i> )	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
<b>distinct</b> ([ <i>numPartitions</i> ])	Return a new dataset that contains the distinct elements of the source dataset.
<b>groupByKey</b> ([ <i>numPartitions</i> ])	<p>When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable&lt;V&gt;) pairs.</p> <p><b>Note:</b> If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using <code>reduceByKey</code> or <code>aggregateByKey</code> will yield much better performance.</p> <p><b>Note:</b> By default, the level of parallelism in the output depends on the number of partitions of the parent RDD.</p>

	<p>You can pass an optional <code>numPartitions</code> argument to set a different number of tasks.</p>
<p><b>reduceByKey</b>(<i>func</i>, <i>[numPartitions]</i>)</p>	<p>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 <i>func</i>, which must be of type (V,V) =&gt; V. Like in <code>groupByKey</code>, the number of reduce tasks is configurable through an optional second argument.</p>
<p><b>aggregateByKey</b>(<i>zeroValue</i>)(<i>s</i> <i>eqOp</i>, <i>combOp</i>, <i>[numPartitions]</i>)</p>	<p>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 <code>groupByKey</code>, the number of reduce tasks is configurable through an optional second argument.</p>
<p><b>sortByKey</b>(<i>[ascending]</i>, <i>[numPartitions]</i>)</p>	<p>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 <code>ascending</code> argument.</p>
<p><b>join</b>(<i>otherDataset</i>, <i>[numPartitions]</i>)</p>	<p>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 <code>leftOuterJoin</code>, <code>rightOuterJoin</code>, and <code>fullOuterJoin</code>.</p>
<p><b>cogroup</b>(<i>otherDataset</i>, <i>[numPartitions]</i>)</p>	<p>When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (Iterable&lt;V&gt;, Iterable&lt;W&gt;)) tuples. This operation is also called <code>groupWith</code>.</p>
<p><b>cartesian</b>(<i>otherDataset</i>)</p>	<p>When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).</p>

Action	Meaning
<b>reduce</b> ( <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.
<b>collect</b> ()	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.
<b>count</b> ()	Return the number of elements in the dataset.
<b>first</b> ()	Return the first element of the dataset (similar to take(1)).
<b>take</b> ( <i>n</i> )	Return an array with the first <i>n</i> elements of the dataset.
<b>takeSample</b> ( <i>withReplacement</i> , <i>num</i> , [ <i>seed</i> ])	Return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.
<b>takeOrdered</b> ( <i>n</i> , [ <i>ordering</i> ])	Return the first <i>n</i> elements of the RDD using either their natural order or a custom comparator.
<b>saveAsTextFile</b> ( <i>path</i> )	Write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call toString on each element to convert it to a line of text in the file.
<b>saveAsSequenceFile</b> ( <i>path</i> )	Write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. This is available on RDDs of key-value pairs that implement Hadoop's Writable interface. In

---

(Java and Scala)	Scala, it is also available on types that are implicitly convertible to Writable (Spark includes conversions for basic types like Int, Double, String, etc).
------------------	--

---

<b>saveAsObjectFile</b> ( <i>path</i> )	Write the elements of the dataset in a simple format using Java serialization, which can then be loaded using <code>SparkContext.objectFile()</code> .
---	--

(Java and Scala)

---

<b>countByKey()</b>	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.
---------------------	--