



## Apache Spark Developer Cheat Sheet

### Transformations (return new RDDs – Lazy)

Where	Function	DStream API	Description
<a href="#">RDD</a>	<a href="#">map(function)</a>	<a href="#">Yes</a>	Return a new distributed dataset formed by passing each element of the source through a function.
<a href="#">RDD</a>	<a href="#">filter(function)</a>	<a href="#">Yes</a>	Return a new dataset formed by selecting those elements of the source on which function returns true.
<a href="#">OrderedRDD Functions</a>	<a href="#">filterByRange(lower, upper)</a>	No	Returns an RDD containing only the elements in the the inclusive range lower to upper.
<a href="#">RDD</a>	<a href="#">flatMap(function)</a>	<a href="#">Yes</a>	Similar to map, but each input item can be mapped to 0 or more output items (so function should return a Seq rather than a single item).
<a href="#">RDD</a>	<a href="#">mapPartitions(function)</a>	<a href="#">Yes</a>	Similar to map, but runs separately on each partition of the RDD.
<a href="#">RDD</a>	<a href="#">mapPartitionsWithIndex(function)</a>	No	Similar to mapPartitions, but also provides function with an integer value representing the index of the partition.
<a href="#">RDD</a>	<a href="#">sample(withReplacement, fraction, seed)</a>	No	Sample a fraction of the data, with or without replacement, using a given random number generator seed.
<a href="#">RDD</a>	<a href="#">union(otherDataset)</a>	<a href="#">Yes</a>	Return a new dataset that contains the union of the elements in the datasets.
<a href="#">RDD</a>	<a href="#">intersection(otherDataset)</a>	No	Return a new RDD that contains the intersection of elements in the datasets.
<a href="#">RDD</a>	<a href="#">distinct([numTasks])</a>	No	Return a new dataset that contains the distinct elements of the source dataset.
<a href="#">PairRDD Functions</a>	<a href="#">groupByKey([numTasks])</a>	<a href="#">Yes</a>	Returns a dataset of (K, Iterable<V>) pairs. Use reduceByKey or aggregateByKey to perform an aggregation (such as a sum or average).
<a href="#">PairRDD Functions</a>	<a href="#">reduceByKey(function, [numTasks])</a>	<a href="#">Yes</a>	Returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function.
<a href="#">PairRDD Functions</a>	<a href="#">aggregateByKey(zeroValue)(seqOp, combOp, [numTasks])</a>	No	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.
<a href="#">OrderedRDD Functions</a>	<a href="#">sortByKey([ascending], [numTasks])</a>	No	Returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean ascending argument.
<a href="#">PairRDD Functions</a>	<a href="#">join(otherDataset, [numTasks])</a>	<a href="#">Yes</a>	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.
<a href="#">PairRDD Functions</a>	<a href="#">cogroup(otherDataset, [numTasks])</a>	<a href="#">Yes</a>	When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (Iterable<V>, Iterable<W>)) tuples.
<a href="#">RDD</a>	<a href="#">cartesian(otherDataset)</a>	No	When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).
<a href="#">RDD</a>	<a href="#">pipe(command, [envVars])</a>	No	Pipe each partition of the RDD through a shell command, e.g. a Perl or bash script.
<a href="#">RDD</a>	<a href="#">coalesce(numPartitions)</a>	No	Decrease the number of partitions in the RDD to numPartitions. Useful for running operations more efficiently after filtering down a large dataset.
<a href="#">RDD</a>	<a href="#">repartition(numPartitions)</a>	<a href="#">Yes</a>	Reshuffle the data in the RDD randomly to create either more or fewer partitions and balance it across them. This always shuffles all data over the network.
<a href="#">OrderedRDD Functions</a>	<a href="#">repartitionAndSortWithinPartitions(partitioner)</a>	No	Repartition the RDD according to the given partitioner and, within each resulting partition, sort records by their keys. More efficient than calling repartition and then sorting.

### Actions (return values – NOT Lazy)



Where	Function	DStream API	Description
<a href="#">RDD</a>	<a href="#">reduce(function)</a>	<a href="#">Yes</a>	Aggregate the elements of the dataset using a function (which takes two arguments and returns one).
<a href="#">RDD</a>	<a href="#">collect()</a>	No	Return all the elements of the dataset as an array at the driver program. Best used on sufficiently small subsets of data.
<a href="#">RDD</a>	<a href="#">count()</a>	<a href="#">Yes</a>	Return the number of elements in the dataset.
<a href="#">RDD</a>	<a href="#">countByValue()</a>	<a href="#">Yes</a>	Return the count of each unique value in this RDD as a local map of (value, count) pairs.
<a href="#">RDD</a>	<a href="#">first()</a>	No	Return the first element of the dataset (similar to take(1)).
<a href="#">RDD</a>	<a href="#">take(n)</a>	No	Return an array with the first n elements of the dataset.
<a href="#">RDD</a>	<a href="#">takeSample(withReplacement, num, [seed])</a>	No	Return an array with a random sample of num elements of the dataset.
<a href="#">RDD</a>	<a href="#">takeOrdered(n, [ordering])</a>	No	Return the first n elements of the RDD using either their natural order or a custom comparator.
<a href="#">RDD</a>	<a href="#">saveAsTextFile(path)</a>	<a href="#">Yes</a>	Write the elements of the dataset as a text. Spark will call toString on each element to convert it to a line of text in the file.
<a href="#">SequenceFileRDD Functions</a>	<a href="#">saveAsSequenceFile(path) (Java and Scala)</a>	No	Write the elements of the dataset as a Hadoop SequenceFile in a given path. For RDDs of key-value pairs that use Hadoop's Writable interface.
<a href="#">RDD</a>	<a href="#">saveAsObjectFile(path) (Java and Scala)</a>	<a href="#">Yes</a>	Write the elements of the dataset in a simple format using Java serialization, which can then be loaded using SparkContext.objectFile().
<a href="#">PairRDD Functions</a>	<a href="#">countByKey()</a>	No	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.
<a href="#">RDD</a>	<a href="#">foreach(function)</a>	<a href="#">Yes</a>	Run a function on each element of the dataset. This is usually done for side effects such as updating an Accumulator.

## Persistence Methods

Where	Function	DStream API	Description
<a href="#">RDD</a>	<a href="#">cache()</a>	<a href="#">Yes</a>	Don't be afraid to call cache on RDDs to avoid unnecessary recomputation. NOTE: This is the same as persist(MEMORY_ONLY).
<a href="#">RDD</a>	<a href="#">persist([Storage Level])</a>	<a href="#">Yes</a>	Persist this RDD with the default storage level.
<a href="#">RDD</a>	<a href="#">unpersist()</a>	No	Mark the RDD as non-persistent, and remove its blocks from memory and disk.
<a href="#">RDD</a>	<a href="#">checkpoint()</a>	<a href="#">Yes</a>	Save to a file inside the checkpoint directory and all references to its parent RDDs will be removed.

## Additional Transformation and Actions

Where	Function	Description
<a href="#">SparkContext</a>	<a href="#">doubleRDDToDoubleRDDFunctions</a>	Extra functions available on RDDs of Doubles
<a href="#">SparkContext</a>	<a href="#">numericRDDToDoubleRDDFunctions</a>	Extra functions available on RDDs of Doubles
<a href="#">SparkContext</a>	<a href="#">rddToPairRDDFunctions</a>	Extra functions available on RDDs of (key, value) pairs
<a href="#">SparkContext</a>	<a href="#">hadoopFile()</a>	Get an RDD for a Hadoop file with an arbitrary InputFormat
<a href="#">SparkContext</a>	<a href="#">hadoopRDD()</a>	Get an RDD for a Hadoop file with an arbitrary InputFormat
<a href="#">SparkContext</a>	<a href="#">makeRDD()</a>	Distribute a local Scala collection to form an RDD
<a href="#">SparkContext</a>	<a href="#">parallelize()</a>	Distribute a local Scala collection to form an RDD
<a href="#">SparkContext</a>	<a href="#">textFile()</a>	Read a text file from a file system URI
<a href="#">SparkContext</a>	<a href="#">wholeTextFiles()</a>	Read a directory of text files from a file system URI

## Extended RDDs w/ Custom Transformations and Actions

RDD Name	Description
<a href="#">CoGroupedRDD</a>	A RDD that cogroups its parents. For each key k in parent RDDs, the resulting RDD contains a tuple with the list of values for that key.
<a href="#">EdgeRDD</a>	Storing the edges in columnar format on each partition for performance. It may additionally store the vertex attributes associated with each edge.
<a href="#">JdbcRDD</a>	An RDD that executes an SQL query on a JDBC connection and reads results. For usage example, see test case JdbcRDDSuite.
<a href="#">ShuffledRDD</a>	The resulting RDD from a shuffle.
<a href="#">VertexRDD</a>	Ensures that there is only one entry for each vertex and by pre-indexing the entries for fast, efficient joins.

## Streaming Transformations

Where	Function	Description
<a href="#">DStream</a>	<a href="#">window(windowLength, slideInterval)</a>	Return a new DStream which is computed based on windowed batches of the source DStream.
<a href="#">DStream</a>	<a href="#">countByWindow(windowLength, slideInterval)</a>	Return a sliding window count of elements in the stream.
<a href="#">DStream</a>	<a href="#">reduceByWindow(function, windowLength, slideInterval)</a>	Return a new single-element stream, created by aggregating elements in the stream over a sliding interval using function.
<a href="#">PairDStream Functions</a>	<a href="#">reduceByKeyAndWindow(function, windowLength, slideInterval, [numTasks])</a>	Returns a new DStream of (K, V) pairs where the values for each key are aggregated using the given reduce function over batches in a sliding window.
<a href="#">PairDStream Functions</a>	<a href="#">reduceByKeyAndWindow(function, invFunc, windowLength, slideInterval, [numTasks])</a>	A more efficient version of the above reduceByKeyAndWindow(). Only applicable to those reduce functions which have a corresponding "inverse reduce" function. Checkpointing must be enabled for using this operation.
<a href="#">DStream</a>	<a href="#">countByValueAndWindow(windowLength, slideInterval, [numTasks])</a>	Returns a new DStream of (K, Long) pairs where the value of each key is its frequency within a sliding window.
<a href="#">DStream</a>	<a href="#">transform(function)</a>	The transform operation (along with its variations like transformWith) allows arbitrary RDD-to-RDD functions to be applied on a Dstream.
<a href="#">PairDStream Functions</a>	<a href="#">updateStateByKey(function)</a>	The updateStateByKey operation allows you to maintain arbitrary state while continuously updating it with new information.

## RDD Persistence

Storage Level	Meaning
<a href="#">MEMORY_ONLY (default level)</a>	Store RDD as deserialized Java objects. If the RDD does not fit in memory, some partitions will not be cached and will be recomputed on the fly when needed.
<a href="#">MEMORY AND DISK</a>	Store RDD as deserialized Java objects. If the RDD does not fit in memory, store the partitions that don't fit on disk, and load them when they're needed.
<a href="#">MEMORY_ONLY_SER</a>	Store RDD as serialized Java objects. Generally more space-efficient than deserialized objects, but more CPU-intensive to read.
<a href="#">MEMORY AND DISK SER</a>	Similar to MEMORY_ONLY_SER, but spill partitions that don't fit in memory to disk instead of recomputing them on the fly each time they're needed.
<a href="#">DISK_ONLY</a>	Store the RDD partitions only on disk.
<a href="#">MEMORY_ONLY_2</a> <a href="#">MEMORY AND DISK 2, etc...</a>	Same as the levels above, but replicate each partition on two cluster nodes.

## Shared Data

[Broadcast Variables](#) Broadcast variables allow the programmer to keep a read-only variable cached on each machine rather than shipping a copy of it with tasks.

Language	Create, Evaluate
Scala	<pre>val broadcastVar = sc.broadcast(Array(1, 2, 3))  broadcastVar.value</pre>
Java	<pre>Broadcast&lt;int[]&gt; broadcastVar = sc.broadcast(new int[] {1, 2, 3});</pre>

Language	Create, Add, Evaluate
Python	<pre>broadcastVar.value(); broadcastVar = sc.broadcast([1, 2, 3])</pre>
	broadcastVar.value

[Accumulators](#) Accumulators are variables that are only "added" to through an associative operation and can therefore be efficiently supported in parallel.

Language	Create, Add, Evaluate
Scala	<pre>val accum = sc.accumulator(0, My Accumulator)  sc.parallelize(Array(1, 2, 3, 4)).foreach(x =&gt; accum += x)  accum.value</pre>
Java	<pre>Accumulator&lt;Integer&gt; accum = sc.accumulator(0);  sc.parallelize(Arrays.asList(1, 2, 3, 4)).foreach(x -&gt; accum.add(x))  accum.value();</pre>
Python	<pre>accum = sc.accumulator(0)</pre>

## MLlib Reference

Topic	Description
<a href="#">Data types</a>	Vectors, points, matrices.
<a href="#">Basic Statistics</a>	Summary, correlations, sampling, testing and random data.
<a href="#">Classification and regression</a>	Includes SVMs, decision trees, naïve Bayes, etc...
<a href="#">Collaborative filtering</a>	Commonly used for recommender systems.
<a href="#">Clustering</a>	Clustering is an unsupervised learning approach.
<a href="#">Dimensionality reduction</a>	Dimensionality reduction is the process of reducing the number of variables under consideration.
<a href="#">Feature extraction and transformation</a>	Used in selecting a subset of relevant features (variables, predictors) for use in model construction.
<a href="#">Frequent pattern mining</a>	Mining is usually among the first steps to analyze a large-scale dataset.
<a href="#">Optimization</a>	Different optimization methods can have different convergence guarantees.
<a href="#">PMML model export</a>	MLlib supports model export to Predictive Model Markup Language.

## Other References

- [Launching Jobs](#)
- [SQL and DataFrames Programming Guide](#)
- [GraphX Programming Guide](#)
- [SparkR Programming Guide](#)

