

# Apache Spark Developer Cheat Sheet 9

## Transformations (return new RDDs – Lazy)

| Where                        | Function   | DStream API | Description  |
|------------------------------|--|-------------|--|
| RDD                          | <code>map(function)</code>                           | Yes         | Return a new distributed dataset formed by passing each element of the source through a function.  |
| RDD                          | <code>filter(function)</code>                        | Yes         | Return a new dataset formed by selecting those elements of the source on which function returns true.                                    |
| Order-<br>edRDD<br>Functions | <code>filterByRange(lower, upper)</code>             | No          | Returns an RDD containing only the elements in the inclusive range lower to upper.   |
| RDD                          | <code>flatMap(function)</code>                       | Yes         | 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). |
| RDD                          | <code>mapPartitions(function)</code>                 | Yes         | Similar to map, but runs separately on each partition of the RDD.  |
| RDD                          | <code>mapPartitionsWithIndex(function)</code>        | No          | Similar to mapPartitions, but also provides function with an integer value representing the index of the partition.                      |
| RDD                          | <code>sample(withReplacement, fraction, seed)</code> | No          | Sample a fraction of the data, with or without replacement, using a given random number generator seed.                                  |

| Where                       | Function  | DStream API | Description   |
|-----------------------------|---|-------------|---|
| <b>RDD</b>                  | <b>union(otherDataset)</b>                                  | <b>Yes</b>  | Return a new dataset that contains the union of the elements in the datasets.   |
| <b>RDD</b>                  | <b>intersection(otherDataset)</b>                           | No          | Return a new RDD that contains the intersection of elements in the datasets.  |
| <b>RDD</b>                  | <b>distinct([numTasks])</b>                                 | No          | Return a new dataset that contains the distinct elements of the source dataset.   |
| <b>PairRDD Functions</b>    | <b>groupByKey([numTasks])</b>                               | <b>Yes</b>  | Returns a dataset of (K, Iterable<V>) pairs. Use reduceByKey or aggregateByKey to perform an aggregation (such as a sum or average).  |
| <b>PairRDD Functions</b>    | <b>reduceByKey(function, [numTasks])</b>                    | <b>Yes</b>  | Returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function.   |
| <b>PairRDD Functions</b>    | <b>aggregateByKey(zeroValue)(seqOp, combOp, [numTasks])</b> | 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. |
| <b>OrderedRDD Functions</b> | <b>sortByKey([ascending], [numTasks])</b>                   | No          | Returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean ascending argument.  |
| <b>PairRDD Functions</b>    | <b>join(otherDataset, [numTasks])</b>                       | <b>Yes</b>  | 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.    |
| <b>PairRDD Functions</b>    | <b>cogroup(otherDataset, [numTasks])</b>                    | <b>Yes</b>  | When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (Iterable<V>, Iterable<W>)) tuples.   |
| <b>RDD</b>                  | <b>cartesian(otherDataset)</b>                              | No          | When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).  |

| Where                | Function   | DStream API | Description  |
|----------------------|--|-------------|--|
| RDD                  | <code>pipe(command, [envVars])</code>                        | No          | Pipe each partition of the RDD through a shell command, e.g. a Perl or bash script.  |
| RDD                  | <code>coalesce(numPartitions)</code>                         | No          | Decrease the number of partitions in the RDD to numPartitions. Useful for running operations more efficiently after filtering down a large dataset.                                |
| RDD                  | <code>repartition(numPartitions)</code>                      | Yes         | 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.                       |
| OrderedRDD Functions | <code>repartitionAndSortWithinPartitions(partitioner)</code> | 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  |
|---------------------------|--|-------------|--|
| RDD                       | <code>reduce(function)</code>                          | Yes         | Aggregate the elements of the dataset using a function (which takes two arguments and returns one).  |
| RDD                       | <code>collect()</code>                                 | No          | Return all the elements of the dataset as an array at the driver program. Best used on sufficiently small subsets of data.                             |
| RDD                       | <code>count()</code>                                   | Yes         | Return the number of elements in the dataset.  |
| RDD                       | <code>countByValue()</code>                            | Yes         | Return the count of each unique value in this RDD as a local map of (value, count) pairs.  |
| RDD                       | <code>first()</code>                                   | No          | Return the first element of the dataset (similar to <code>take(1)</code> ).  |
| RDD                       | <code>take(n)</code>                                   | No          | Return an array with the first <code>n</code> elements of the dataset.   |
| RDD                       | <code>takeSample(withReplacement, num, [seed])</code>  | No          | Return an array with a random sample of <code>num</code> elements of the dataset.  |
| RDD                       | <code>takeOrdered(n, [ordering])</code>                | No          | Return the first <code>n</code> elements of the RDD using either their natural order or a custom comparator.   |
| RDD                       | <code>saveAsTextFile(path)</code>                      | Yes         | Write the elements of the dataset as a text. Spark will call <code>toString</code> on each element to convert it to a line of text in the file.        |
| SequenceFileRDD Functions | <code>saveAsSequenceFile(path) (Java and Scala)</code> | 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.          |
| RDD                       | <code>saveAsObjectFile(path) (Java and Scala)</code>   | Yes         | Write the elements of the dataset in a simple format using Java serialization, which can then be loaded using <code>SparkContext.objectFile()</code> . |

| Where             | Function          | DStream API | Description   |
|-------------------|-------------------|-------------|---|
| PairRDD Functions | countByKey()      | No          | Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.                |
| RDD               | foreach(function) | Yes         | 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   |
|-------|---------------------------------------|-------------|---|
| RDD   | <code>cache()</code>                  | Yes         | Don't be afraid to call <code>cache</code> on RDDs to avoid unnecessary recomputation.<br>NOTE: This is the same as <code>persist(MEMORY_ONLY)</code> . |
| RDD   | <code>persist([Storage Level])</code> | Yes         | Persist this RDD with the default storage level.  |
| RDD   | <code>unpersist()</code>              | No          | Mark the RDD as non-persistent, and remove its blocks from memory and disk.   |
| RDD   | <code>checkpoint()</code>             | Yes         | 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  |
|---------------------------|---|--|
| <b>SparkCon-<br/>text</b> | <b>doubleRDDToDou-<br/>bleRDDFunctions</b>  | Extra functions available on RDDs of Doubles               |
| <b>SparkCon-<br/>text</b> | <b>numericRDDToDou-<br/>bleRDDFunctions</b> | Extra functions available on RDDs of Doubles               |
| <b>SparkCon-<br/>text</b> | <b>rddToPairRDDFunc-<br/>tions</b>          | Extra functions available on RDDs of (key, value) pairs    |
| <b>SparkCon-<br/>text</b> | <b>hadoopFile()</b>                         | Get an RDD for a Hadoop file with an arbitrary InputFormat |
| <b>SparkCon-<br/>text</b> | <b>hadoopRDD()</b>                          | Get an RDD for a Hadoop file with an arbitrary InputFormat |
| <b>SparkCon-<br/>text</b> | <b>makeRDD()</b>                            | Distribute a local Scala collection to form an RDD         |
| <b>SparkCon-<br/>text</b> | <b>parallelize()</b>                        | Distribute a local Scala collection to form an RDD         |
| <b>SparkCon-<br/>text</b> | <b>textFile()</b>                           | Read a text file from a file system URI                    |
| <b>SparkCon-<br/>text</b> | <b>wholeTextFiles()</b>                     | Read a directory of text files from a file system URI      |

## Extended RDDs w/ Custom Transformations and Actions

| RDD Name            | Description  |
|---------------------|--|
| <b>CoGroupedRDD</b> | A RDD that cogroups its parents. For each key <i>k</i> in parent RDDs, the resulting RDD contains a tuple with the list of values for that key.    |
| <b>EdgeRDD</b>      | Storing the edges in columnar format on each partition for performance. It may additionally store the vertex attributes associated with each edge. |
| <b>JdbcRDD</b>      | An RDD that executes an SQL query on a JDBC connection and reads results. For usage example, see test case JdbcRDDSuite.                           |
| <b>ShuffledRDD</b>  | The resulting RDD from a shuffle.  |
| <b>VertexRDD</b>    | 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   |
|------------------------|---|---|
| DStream                | <code>window(windowLength, slideInterval)</code>  | Return a new DStream which is computed based on windowed batches of the source DStream.   |
| DStream                | <code>countByWindow(windowLength, slideInterval)</code>                                       | Return a sliding window count of elements in the stream.  |
| DStream                | <code>reduceByWindow(function, windowLength, slideInterval)</code>                            | Return a new single-element stream, created by aggregating elements in the stream over a sliding interval using function.   |
| PairD-Stream Functions | <code>reduceByKeyAndWindow(function, windowLength, slideInterval, [numTasks])</code>          | 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.  |
| PairD-Stream Functions | <code>reduceByKeyAndWindow(function, invFunc, windowLength, slideInterval, [numTasks])</code> | A more efficient version of the above <code>reduceByKeyAndWindow()</code> . Only applicable to those reduce functions which have a corresponding “inverse reduce” function. Checkpointing must be enabled for using this operation. |
| DStream                | <code>countByValueAndWindow(windowLength, slideInterval, [numTasks])</code>                   | Returns a new DStream of (K, Long) pairs where the value of each key is its frequency within a sliding window.  |
| DStream                | <code>transform(function)</code>  | The transform operation (along with its variations like <code>transformWith</code> ) allows arbitrary RDD-to-RDD functions to be applied on a Dstream.  |
| PairD-Stream Functions | <code>updateStateByKey(function)</code>   | The <code>updateStateByKey</code> operation allows you to maintain arbitrary state while continuously updating it with new information.   |

## RDD Persistence

| Storage Level                                   | Meaning  |
|---|--|
| <b>MEMORY_ONLY (default level)</b>              | 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. |
| <b>MEMORY_AND_DISK</b>                          | 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.   |
| <b>MEMORY_ONLY_SER</b>                          | Store RDD as serialized Java objects. Generally more space-efficient than deserialized objects, but more CPU-intensive to read.                              |
| <b>MEMORY_AND_DISK_SER</b>                      | 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.           |
| <b>DISK_ONLY</b>                                | Store the RDD partitions only on disk.   |
| <b>MEMORY_ONLY_2, MEMORY_AND_DISK_2, etc...</b> | 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});  broadcastVar.value();</pre> |
| Python   | <pre>broadcastVar = sc.broadcast([1, 2, 3])  broadcastVar.value</pre>                                      |

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

## Other References

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