

Python For Data Science Cheat Sheet

PySpark - RDD Basics

Learn Python for data science [Interactively](https://www.datacamp.com) at [www.DataCamp.com](https://www.datacamp.com)



Spark

PySpark is the Spark Python API that exposes the Spark programming model to Python.



Initializing Spark

SparkContext

```
>>> from pyspark import SparkContext
>>> sc = SparkContext(master = 'local[2]')
```

Inspect SparkContext

<pre>>>> sc.version >>> sc.pythonVer >>> sc.master >>> str(sc.sparkHome) >>> str(sc.sparkUser()) >>> sc.appName >>> sc.applicationId >>> sc.defaultParallelism >>> sc.defaultMinPartitions</pre>	<p>Retrieve SparkContext version</p> <p>Retrieve Python version</p> <p>Master URL to connect to</p> <p>Path where Spark is installed on worker nodes</p> <p>Retrieve name of the Spark User running SparkContext</p> <p>Return application name</p> <p>Retrieve application ID</p> <p>Return default level of parallelism</p> <p>Default minimum number of partitions for RDDs</p>
--	--

Configuration

```
>>> from pyspark import SparkConf, SparkContext
>>> conf = (SparkConf()
           .setMaster("local")
           .setAppName("My app")
           .set("spark.executor.memory", "1g"))
>>> sc = SparkContext(conf = conf)
```

Using The Shell

In the PySpark shell, a special interpreter-aware SparkContext is already created in the variable called `sc`.

```
$ ./bin/spark-shell --master local[2]
$ ./bin/pyspark --master local[4] --py-files code.py
```

Set which master the context connects to with the `--master` argument, and add Python `.zip`, `.egg` or `.py` files to the runtime path by passing a comma-separated list to `--py-files`.

Loading Data

Parallelized Collections

```
>>> rdd = sc.parallelize([('a',7), ('a',2), ('b',2)])
>>> rdd2 = sc.parallelize([('a',2), ('d',1), ('b',1)])
>>> rdd3 = sc.parallelize(range(100))
>>> rdd4 = sc.parallelize(["a", "x", "y", "z"],
                        ["b", "p", "r"])
```

External Data

Read either one text file from HDFS, a local file system or any Hadoop-supported file system URI with `textFile()`, or read in a directory of text files with `wholeTextFiles()`.

```
>>> textFile = sc.textFile("/my/directory/*.txt")
>>> textFile2 = sc.wholeTextFiles("/my/directory/")
```

Retrieving RDD Information

Basic Information

<pre>>>> rdd.getNumPartitions() >>> rdd.count() 3 >>> rdd.countByKey() defaultdict(<type 'int'>, {'a':2, 'b':1}) >>> rdd.countByValue() defaultdict(<type 'int'>, {'b':2}:1, ('a',2):1, ('a',7):1) >>> rdd.collectAsMap() {'a': 2, 'b': 2} >>> rdd3.sum() 4950 >>> sc.parallelize([]).isEmpty() True</pre>	<p>List the number of partitions</p> <p>Count RDD instances</p> <p>Count RDD instances by key</p> <p>Count RDD instances by value</p> <p>Return (key,value) pairs as a dictionary</p> <p>Sum of RDD elements</p> <p>Check whether RDD is empty</p>
---	--

Summary

<pre>>>> rdd3.max() 99 >>> rdd3.min() 0 >>> rdd3.mean() 49.5 >>> rdd3.stdev() 28.866070047722118 >>> rdd3.variance() 833.25 >>> rdd3.histogram(3) ([0, 33, 66, 99], [33, 33, 34]) >>> rdd3.stats()</pre>	<p>Maximum value of RDD elements</p> <p>Minimum value of RDD elements</p> <p>Mean value of RDD elements</p> <p>Standard deviation of RDD elements</p> <p>Compute variance of RDD elements</p> <p>Compute histogram by bins</p> <p>Summary statistics (count, mean, stdev, max & min)</p>
---	--

Applying Functions

<pre>>>> rdd.map(lambda x: x+(x[1],x[0])) .collect() [('a',7,7,'a'), ('a',2,2,'a'), ('b',2,2,'b')] >>> rdd5 = rdd.flatMap(lambda x: x+(x[1],x[0])) >>> rdd5.collect() ['a',7,7,'a','a',2,2,'a','b',2,2,'b'] >>> rdd4.flatMapValues(lambda x: x) .collect() [('a','x'),('a','y'),('a','z'),('b','p'),('b','r')]</pre>	<p>Apply a function to each RDD element</p> <p>Apply a function to each RDD element and flatten the result</p> <p>Apply a flatMap function to each (key,value) pair of rdd4 without changing the keys</p>
---	---

Selecting Data

<pre>>>> rdd.collect() [('a', 7), ('a', 2), ('b', 2)] >>> rdd.take(2) [('a', 7), ('a', 2)] >>> rdd.first() ('a', 7) >>> rdd.top(2) [('b', 2), ('a', 7)] >>> rdd3.sample(False, 0.15, 81).collect() [3, 4, 27, 31, 40, 41, 42, 43, 60, 76, 79, 80, 86, 97] >>> rdd.filter(lambda x: "a" in x) .collect() [('a',7), ('a',2)] >>> rdd5.distinct().collect() ['a',2,'b',7] >>> rdd.keys().collect() ['a', 'a', 'b']</pre>	<p>Return a list with all RDD elements</p> <p>Take first 2 RDD elements</p> <p>Take first RDD element</p> <p>Take top 2 RDD elements</p> <p>Return sampled subset of rdd3</p> <p>Filter the RDD</p> <p>Return distinct RDD values</p> <p>Return (key,value) RDD's keys</p>
---	--

Iterating

<pre>>>> def g(x): print(x) >>> rdd.foreach(g) ('a', 7) ('b', 2) ('a', 2)</pre>	<p>Apply a function to all RDD elements</p>
---	---

Reshaping Data

<pre>>>> rdd.reduceByKey(lambda x,y : x+y) .collect() [('a',9), ('b',2)] >>> rdd.reduce(lambda a, b: a + b) ('a',7,'a',2,'b',2) >>> rdd3.groupBy(lambda x: x % 2) .mapValues(list) .collect() >>> rdd.groupByKey() .mapValues(list) .collect() [('a', [7,2]), ('b', [2])]</pre>	<p>Merge the rdd values for each key</p> <p>Merge the rdd values</p> <p>Return RDD of grouped values</p> <p>Group rdd by key</p>
<pre>>>> seqOp = (lambda x,y: (x[0]+y,x[1]+1)) >>> combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1])) >>> rdd3.aggregate((0,0),seqOp,combOp) (4950,100) >>> rdd.aggregateByKey((0,0),seqOp,combOp) .collect() [('a', (9,2)), ('b', (2,1))] >>> rdd3.fold(0,add) 4950 >>> rdd.foldByKey(0, add) .collect() [('a',9), ('b',2)] >>> rdd3.keyBy(lambda x: x+x) .collect()</pre>	<p>Aggregate RDD elements of each partition and then the results</p> <p>Aggregate values of each RDD key</p> <p>Aggregate the elements of each partition, and then the results</p> <p>Merge the values for each key</p> <p>Create tuples of RDD elements by applying a function</p>

Mathematical Operations

<pre>>>> rdd.subtract(rdd2) .collect() [('b',2), ('a',7)] >>> rdd2.subtractByKey(rdd) .collect() [('d', 1)] >>> rdd.cartesian(rdd2).collect()</pre>	<p>Return each rdd value not contained in rdd2</p> <p>Return each (key,value) pair of rdd2 with no matching key in rdd</p> <p>Return the Cartesian product of rdd and rdd2</p>
--	--

Sort

<pre>>>> rdd2.sortBy(lambda x: x[1]) .collect() [('d',1), ('b',1), ('a',2)] >>> rdd2.sortByKey() .collect() [('a',2), ('b',1), ('d',1)]</pre>	<p>Sort RDD by given function</p> <p>Sort (key, value) RDD by key</p>
---	---

Repartitioning

<pre>>>> rdd.repartition(4) >>> rdd.coalesce(1)</pre>	<p>New RDD with 4 partitions</p> <p>Decrease the number of partitions in the RDD to 1</p>
---	---

Saving

```
>>> rdd.saveAsTextFile("rdd.txt")
>>> rdd.saveAsHadoopFile("hdfs://namenodehost/parent/child",
                        'org.apache.hadoop.mapred.TextOutputFormat')
```

Stopping SparkContext

```
>>> sc.stop()
```

Execution

```
$ ./bin/spark-submit examples/src/main/python/pi.py
```

