

# Python数据科学速查表

## PySpark - RDD 基础

### Spark

PySpark 是 Spark 的 Python API，允许 Python 调用 Spark 编程模型。



### 初始化 Spark

### SparkContext

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

### 核查 SparkContext

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; sc.version &gt;&gt;&gt; sc.pythonVer &gt;&gt;&gt; sc.master &gt;&gt;&gt; str(sc.sparkHome) &gt;&gt;&gt; str(sc.sparkUser())  &gt;&gt;&gt; sc.appName &gt;&gt;&gt; sc.applicationId &gt;&gt;&gt; sc.defaultParallelism &gt;&gt;&gt; sc.defaultMinPartitions</pre> | <p>获取 SparkContext 版本<br/>获取 Python 版本<br/>要连接的 Master URL<br/>Spark 在工作节点的安装路径<br/>获取 SparkContext 的 Spark 用户名</p> <p>返回应用名称<br/>获取应用程序ID<br/>返回默认并行级别<br/>RDD默认最小分区数</p> |
|--|--|

### 配置

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

### 使用 Shell

PySpark Shell 已经为 SparkContext 创建了名为 sc 的变量。

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

用 --master 参数设定 Context 连接到哪个 Master 服务器，通过传递逗号分隔列表至 --py-files 添加 Python.zip、.egg 或 .py 文件到 Runtime 路径。

### 加载数据

### 并行集合

```
>>> 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"])]])
```

### 外部数据

使用 `textFile()` 函数从HDFS、本地文件或其它支持 Hadoop 的文件系统里读取文本文件，或使用 `wholeTextFiles()` 函数读取目录里文本文件。

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

### 提取 RDD 信息

#### 基础信息

|   |   |
|---|---|
| <pre>&gt;&gt;&gt; rdd.getNumPartitions() &gt;&gt;&gt; rdd.count() 3 &gt;&gt;&gt; rdd.countByKey() defaultdict(&lt;type 'int'&gt;, {'a': 2, 'b': 1}) &gt;&gt;&gt; rdd.countByValue() defaultdict(&lt;type 'int'&gt;, {'b': 2}: 1, ('a', 2): 1, ('a', 7): 1) &gt;&gt;&gt; rdd.collectAsMap() {'a': 2, 'b': 2} &gt;&gt;&gt; rdd3.sum() 4950 &gt;&gt;&gt; sc.parallelize([]).isEmpty() True</pre> | <p>列出分区数<br/>计算 RDD 实例数量</p> <p>按键计算 RDD 实例数量</p> <p>按值计算 RDD 实例数量</p> <p>以字典形式返回键值</p> <p>汇总 RDD 元素</p> <p>检查 RDD 是否为空</p> |
|---|---|

#### 汇总

|   |   |
|---|---|
| <pre>&gt;&gt;&gt; rdd3.max() 99 &gt;&gt;&gt; rdd3.min() 0 &gt;&gt;&gt; rdd3.mean() 49.5 &gt;&gt;&gt; rdd3.stdev() 28.866070047722118 &gt;&gt;&gt; rdd3.variance() 833.25 &gt;&gt;&gt; rdd3.histogram(3) ([0, 33, 66, 99], [33, 33, 34]) &gt;&gt;&gt; rdd3.stats()</pre> | <p>RDD 元素的最大值</p> <p>RDD 元素的最小值</p> <p>RDD 元素的平均值</p> <p>RDD 元素的标准差</p> <p>计算 RDD 元素的方差</p> <p>分箱 (Bin) 生成直方图</p> <p>综合统计<br/>包括：计数、平均值、标准差、最大值和最小值</p> |
|---|---|

### 应用函数

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; rdd.map(lambda x: x+(x[1],x[0])) ...      .collect() [('a', 7, 7, 'a'), ('a', 2, 2, 'a'), ('b', 2, 2, 'b')] &gt;&gt;&gt; rdd5 = rdd.flatMap(lambda x: x+(x[1],x[0]))  &gt;&gt;&gt; rdd5.collect() ['a', 7, 7, 'a', 'a', 2, 2, 'a', 'b', 2, 2, 'b'] &gt;&gt;&gt; rdd4.flatMapValues(lambda x: x) ...      .collect() [('a', 'x'), ('a', 'y'), ('a', 'z'), ('b', 'p'), ('b', 'r')]</pre> | <p>对每个 RDD 元素执行函数</p> <p>对每个 RDD 元素执行函数，并拉平结果</p> <p>不改变键，对 rdd4 的每个键值对执行 flatMap 函数</p> |
|--|--|

### 选择数据

|  |   |
|--|---|
| <p><b>获取</b></p> <pre>&gt;&gt;&gt; rdd.collect() [('a', 7), ('a', 2), ('b', 2)] &gt;&gt;&gt; rdd.take(2) [('a', 7), ('a', 2)] &gt;&gt;&gt; rdd.first() ('a', 7) &gt;&gt;&gt; rdd.top(2) [('b', 2), ('a', 7)]</pre> <p><b>抽样</b></p> <pre>&gt;&gt;&gt; rdd3.sample(False, 0.15, 81).collect() [3, 4, 27, 31, 40, 41, 42, 43, 60, 76, 79, 80, 86, 97]</pre> <p><b>筛选</b></p> <pre>&gt;&gt;&gt; rdd.filter(lambda x: "a" in x) ...      .collect() [('a', 7), ('a', 2)] &gt;&gt;&gt; rdd5.distinct().collect() ['a', 2, 'b', 7] &gt;&gt;&gt; rdd.keys().collect() ['a', 'a', 'b']</pre> | <p>返回包含所有 RDD 元素的列表</p> <p>提取前两个 RDD 元素</p> <p>提取第一个 RDD 元素</p> <p>提取前两个 RDD 元素</p> <p>返回 rdd3 的采样子集</p> <p>筛选 RDD</p> <p>返回 RDD 里的唯一值</p> <p>返回 RDD 键值对里的键</p> |
|--|---|

### 迭代

|   |                   |
|---|-------------------|
| <pre>&gt;&gt;&gt; def g(x): print(x) &gt;&gt;&gt; rdd.foreach(g) ('a', 7) ('b', 2) ('a', 2)</pre> | <p>为所有RDD应用函数</p> |
|---|-------------------|

### 改变数据形状

|   |  |
|---|--|
| <p><b>规约</b></p> <pre>&gt;&gt;&gt; rdd.reduceByKey(lambda x,y : x+y) ...      .collect() [('a', 9), ('b', 2)] &gt;&gt;&gt; rdd.reduce(lambda a, b: a + b) ('a', 7, 'a', 2, 'b', 2)<p><b>分组</b></p><pre>&gt;&gt;&gt; rdd3.groupBy(lambda x: x % 2) ...      .mapValues(list) ...      .collect() &gt;&gt;&gt; rdd.groupByKey() ...      .mapValues(list) ...      .collect() [('a', [7, 2]), ('b', [2])]<p><b>聚合</b></p><pre>&gt;&gt;&gt; seqOp = (lambda x,y: (x[0]+y,x[1]+1)) &gt;&gt;&gt; combOp = (lambda x,y:(x[0]+y[0],x[1]+y[1])) &gt;&gt;&gt; rdd3.aggregate((0,0),seqOp,combOp) (4950,100) &gt;&gt;&gt; rdd.aggregateByKey((0,0),seqOp,combOp) ...      .collect() [('a', (9, 2)), ('b', (2, 1))] &gt;&gt;&gt; rdd3.fold(0,add) 4950 &gt;&gt;&gt; rdd.foldByKey(0, add) ...      .collect() [('a', 9), ('b', 2)] &gt;&gt;&gt; rdd3.keyBy(lambda x: x+x) ...      .collect()</pre></pre></pre> | <p>合并每个键的 RDD 值</p> <p>合并 RDD 的值</p> <p>返回 RDD 的分组值</p> <p>按键分组 RDD</p> <p>汇总每个分区里的 RDD 元素，并输出结果<br/>汇总每个 RDD 的键的值</p> <p>汇总每个分区里的 RDD 元素，并输出结果<br/>合并每个键的值</p> <p>通过执行函数，创建 RDD 元素的元组</p> |
|---|--|

### 数学运算

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; rdd.subtract(rdd2) ...      .collect() [('b', 2), ('a', 7)] &gt;&gt;&gt; rdd2.subtractByKey(rdd) ...      .collect() [('d', 1)] &gt;&gt;&gt; rdd.cartesian(rdd2).collect()</pre> | <p>返回在 rdd2 里没有匹配键的 rdd 键值对</p> <p>返回 rdd2 里的每个 (键, 值) 对，rdd中没有匹配的键</p> <p>返回 rdd 和 rdd2 的笛卡尔积</p> |
|--|--|

### 排序

|   |   |
|---|---|
| <pre>&gt;&gt;&gt; rdd2.sortBy(lambda x: x[1]) ...      .collect() [('d', 1), ('b', 1), ('a', 2)] &gt;&gt;&gt; rdd2.sortByKey() ...      .collect() [('a', 2), ('b', 1), ('d', 1)]</pre> | <p>按给定函数排序 RDD</p> <p>按键排序 RDD 的键值对</p> |
|---|---|

### 重分区

|   |  |
|---|--|
| <pre>&gt;&gt;&gt; rdd.repartition(4) &gt;&gt;&gt; rdd.coalesce(1)</pre> | <p>新建一个含4个分区的 RDD<br/>将 RDD 中的分区数缩减为1个</p> |
|---|--|

### 保存

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

### 终止 SparkContext

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

### 执行程序

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

原文作者

**DataCamp**  
Learn Python for Data Science Interactively

