

MOBIN

离开舒适区，坚持不懈，持续学习！！！

博客园 首页 新随笔 联系 管理 订阅 XML

随笔- 45 文章- 0 评论- 47

Github

昵称：MOBIN  
园龄：2年11个月  
粉丝：86  
关注：0  
[+加关注](#)

Spark函数详解系列之RDD基本转换

摘要：

**RDD**：弹性分布式数据集，是一种特殊集合，支持多种来源，有容错机制，可以被缓存，支持并行操作，一个**RDD**代表一个分区里的数据集

**RDD**有两种操作算子：

**Transformation**（转换）：**Transformation**属于延迟计算，当一个**RDD**转换成另一个**RDD**时并没有立即进行转换，仅仅是记住了数据集的逻辑操作

**Ation**（执行）：触发**Spark**作业的运行，真正触发转换算子的计算

本系列主要讲解Spark中常用的函数操作：

- 1.RDD基本转换
- 2.键-值RDD转换
- 3.Action操作篇

本节所讲函数

- 1.map(func)
- 2.flatMap(func)
- 3.mapPartitions(func)
- 4.mapPartitionsWithIndex(func)
- 5.simple(withReplacement,fraction,seed)
- 6.union(orthorDataset)
- 7.intersection(otherDataset)
- 8.distinct([numTasks])
- 9.cartesian(otherDataset)
- 10.coalesce(numPartitions, shuffle)
- 11.repartition(numPartition)
- 12.glom()
- 13.randomSplit(weight:Array[Double],seed)

基础转换操作：

**1.map(func)**：数据集中的每个元素经过用户自定义的函数转换形成一个新的RDD，新的RDD叫MappedRDD

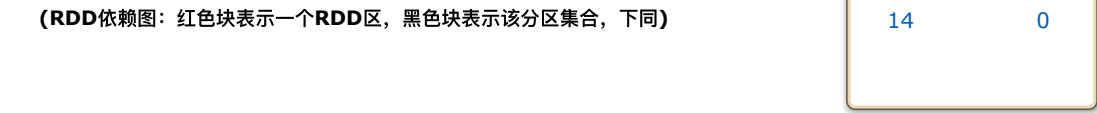
(例1)

```
1 object Map {
2   def main(args: Array[String]) {
3     val conf = new SparkConf().setMaster("local").setAppName("map")
4     val sc = new SparkContext(conf)
5     val rdd = sc.parallelize(1 to 10) //创建RDD
6     val map = rdd.map(_*2)           //对RDD中的每个元素都乘于2
7     map.foreach(x => print(x+" "))
8     sc.stop()
9   }
10 }
```

输出：

2 4 6 8 10 12 14 16 18 20

(RDD依赖图：红色块表示一个RDD区，黑色块表示该分区集合，下同)



<2018年5月>

日	一	二	三	四	五	六
29	30	1	2	3	4	5
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31	1	2
3	4	5	6	7	8	9

搜索

常用链接

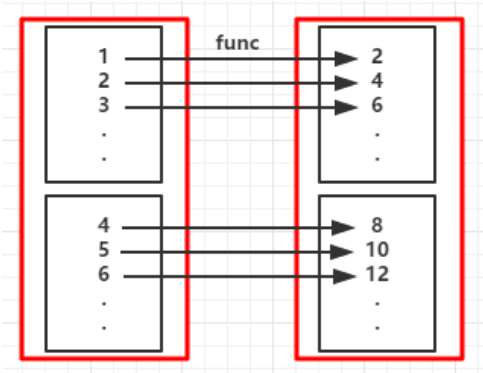
- 我的随笔
- 我的评论
- 我的参与
- 最新评论
- 我的标签

我的标签

- JAVA(9)
- hbase(7)
- Spark(7)
- 数据结构(6)
- 算法(6)
- Scala(5)
- Phoenix(4)
- Spring(4)
- Hive(3)
- Java并发(3)
- 更多

随笔分类

- Akka(1)
- Docker(1)
- Hadoop(2)
- HBase(7)
- Hive(3)
- Java(5)
- Java并发编程(3)
- Maven(2)
- Phoenix(4)
- Scala(5)
- Scalatra(1)
- Spark(6)
- Spring Boot(1)



Struts2(2)  
数据结构(6)  
微服务(1)

随笔档案

- 2018年3月 (1)
- 2017年11月 (1)
- 2017年7月 (1)
- 2017年3月 (1)
- 2016年12月 (2)
- 2016年9月 (1)
- 2016年7月 (2)
- 2016年6月 (4)
- 2016年5月 (1)
- 2016年4月 (8)
- 2016年3月 (8)
- 2016年1月 (3)
- 2015年12月 (1)
- 2015年11月 (1)
- 2015年10月 (1)
- 2015年8月 (1)
- 2015年7月 (8)

文章分类

Ruby

积分与排名

积分 - 87178  
排名 - 3955

最新评论

- 1. Re:Actor模型原理  
二楼的例子更简单  
--nuc093
- 2. Re:Spark常用函数讲解之键值RDD转换  
讲的很好，受教了。  
--Super<John
- 3. Re:Hive2.0函数大全(中文版)  
收藏了多谢  
--哈士奇说喵
- 4. Re:Hive2.0函数大全(中文版)  
非常方便，感谢  
--柏原森森
- 5. Re:图解堆排序  
for(int i = len/2 - 1; i >=0; i --){ //堆构造  
这里的应该用arr.length / 2，而不应该再拿arr.length减1了吧  
--cumtli

阅读排行榜

- 1. Hive2.0函数大全(中文版)(43023)
- 2. Spark函数详解系列之RDD基本转换(38519)
- 3. 图解快速排序(34090)
- 4. 图解插入排序--直接插入排序(23331)
- 5. java并发编程--Executor框架(22829)

评论排行榜

- 1. java并发编程--Executor框架(7)
- 2. 图解堆排序(5)
- 3. 图解快速排序(4)
- 4. Java对象在JVM中的生命周期(3)
- 5. Actor模型原理(3)

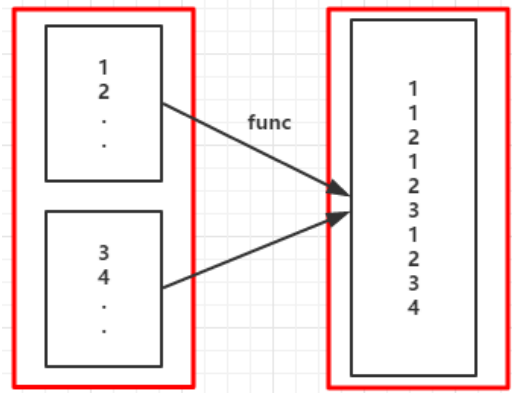
2.flatMap(func):与map类似，但每个元素输入项都可以被映射到0个或多个的输出项，最终将结果“扁平化”后输出  
(例2)

```
1 //...省略sc
2 val rdd = sc.parallelize(1 to 5)
3 val fm = rdd.flatMap(x => (1 to x)).collect()
4 fm.foreach( x => print(x + " "))
```

输出:  
1 1 2 1 2 3 1 2 3 4 1 2 3 4 5

如果是map函数其输出如下:  
Range(1) Range(1, 2) Range(1, 2, 3) Range(1, 2, 3, 4) Range(1, 2, 3, 4, 5)

(RDD依赖图)



3.mapPartitions(func):类似与map，map作用于每个分区的每个元素，但mapPartitions作用于每个分区工  
func的类型: Iterator[T] => Iterator[U]  
假设有N个元素，有M个分区，那么map的函数的将被调用N次,而mapPartitions被调用M次,当在映射的过程中不断的创建对象时就可以使用mapPartitions比map的效率要高很多，比如当向数据库写入数据时，如果使用map就需要为每个元素创建connection对象，但使用mapPartitions的话就需要为每个分区创建connetcion对象  
(例3): 输出有女性的名字:

```
1 object MapPartitions {
2 //定义函数
3 def partitionsFun(*index : Int,*iter : Iterator[(String,String)]) : Iterator[String] = {
4     var woman = List[String]()
5     while (iter.hasNext){
6         val next = iter.next()
7         next match {
8             case (_, "female") => woman = /*"+index+"*/+*/next._1 :: woman
9             case _ =>
10         }
11     }
12     return woman.iterator
13 }
14
15 def main(args: Array[String]) {
```

## 推荐排行榜

1. 图解快速排序(36)
2. 深度剖析JDK动态代理机制(30)
3. java并发编程--Executor框架(25)
4. Hive2.0函数大全(中文版)(20)
5. Spark函数详解系列之RDD基本转换(14)

```

16  val conf = new SparkConf().setMaster("local").setAppName("mappartitions")
17  val sc = new SparkContext(conf)
18  val l = List(("kpop", "female"), ("zorro", "male"), ("mobin", "male"), ("lucy", "female"))
19  val rdd = sc.parallelize(l, 2)
20  val mp = rdd.mapPartitions(partitionsFun)
21  /*val mp = rdd.mapPartitionsWithIndex(partitionsFun)*/
22  mp.collect.foreach(x => (print(x + " "))) //将分区中的元素转换成Array再输出
23  }
24  }

```

输出:

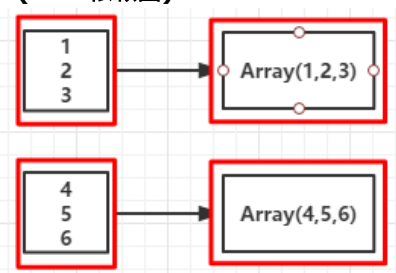
```
kpop lucy
```

其实这个效果可以用一条语句完成

```
1 val mp = rdd.mapPartitions(x => x.filter(_._2 == "female")).map(x => x._1)
```

之所以不那么做是为了演示函数的定义

### (RDD依赖图)



**4.mapPartitionsWithIndex(func):**与mapPartitions类似，不同的时函数多了个分区索引的参数

func类型: (Int, Iterator[T]) => Iterator[U]

(例4): 将例3橙色的注释部分去掉即是

输出: (带了分区索引)

```
[0]kpop [1]lucy
```

**5.sample(withReplacement,fraction,seed):**以指定的随机种子随机抽样出数量为fraction的

数据，withReplacement表示是抽出的数据是否放回，true为有放回的抽样，false为无放回的抽样

(例5): 从RDD中随机且有放回的抽出50%的数据，随机种子值为3 (即可能以1 2 3的其中一个起始值)

```

1  //省略
2  val rdd = sc.parallelize(1 to 10)
3  val sample1 = rdd.sample(true, 0.5, 3)
4  sample1.collect.foreach(x => print(x + " "))
5  sc.stop

```

**6.union(otherDataset):**将两个RDD中的数据集进行合并，最终返回两个RDD的并集，若RDD中存在相同的元素也不会去重

```

1  //省略sc
2  val rdd1 = sc.parallelize(1 to 3)
3  val rdd2 = sc.parallelize(3 to 5)
4  val unionRDD = rdd1.union(rdd2)
5  unionRDD.collect.foreach(x => print(x + " "))
6  sc.stop

```

输出:

```
1 2 3 3 4 5
```

**7.intersection(otherDataset):**返回两个RDD的交集

```

1  //省略sc
2  val rdd1 = sc.parallelize(1 to 3)
3  val rdd2 = sc.parallelize(3 to 5)
4  val unionRDD = rdd1.intersection(rdd2)

```

14

0

```

5 unionRDD.collect.foreach(x => print(x + " "))
6 sc.stop

```

输出:

3 4

## 8.distinct([numTasks]):对RDD中的元素进行去重

```

1 //省略sc
2 val list = List(1,1,2,5,2,9,6,1)
3 val distinctRDD = sc.parallelize(list)
4 val unionRDD = distinctRDD.distinct()
5 unionRDD.collect.foreach(x => print(x + " "))

```

输出:

1 6 9 5 2

## 9.cartesian(otherDataset):对两个RDD中的所有元素进行笛卡尔积操作

```

1 //省略
2 val rdd1 = sc.parallelize(1 to 3)
3 val rdd2 = sc.parallelize(2 to 5)
4 val cartesianRDD = rdd1.cartesian(rdd2)
5 cartesianRDD.foreach(x => println(x + " "))

```

输出:

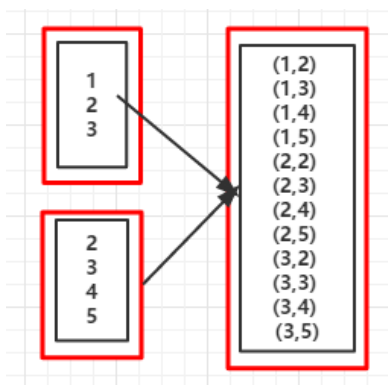
```

(1,2)
(1,3)
(1,4)
(1,5)
(2,2)
(2,3)
(2,4)
(2,5)
(3,2)
(3,3)
(3,4)
(3,5)

```



(RDD依赖图)



## 10.coalesce(numPartitions, shuffle):对RDD的分区进行重新分区, shuffle默认值为false,当

shuffle=false时, 不能增加分区数  
目,但不会报错, 只是分区个数还是原来的

(例9:) **shuffle=false**

```

1 //省略
2 val rdd = sc.parallelize(1 to 16,4)
3 val coalesceRDD = rdd.coalesce(3) //当shuffle的值为false时, 不能增加分区数(即

```

14

0

```
4 | println("重新分区后的分区个数:"+coalesceRDD.partitions.size)
```

输出:

```
重新分区后的分区个数:3
//分区后的数据集
List(1, 2, 3, 4)
List(5, 6, 7, 8)
List(9, 10, 11, 12, 13, 14, 15, 16)
```

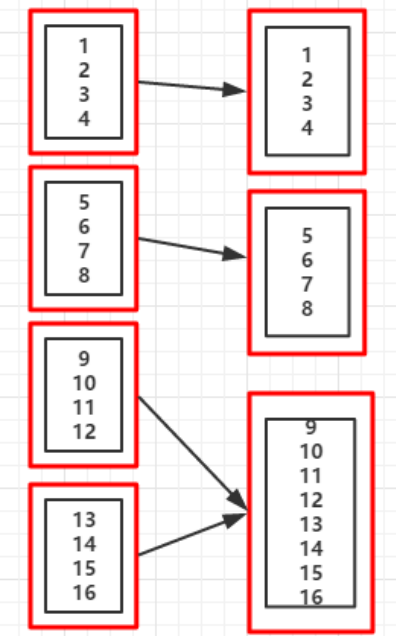
(例9.1:) shuffle=true

```
1 | //...省略
2 | val rdd = sc.parallelize(1 to 16,4)
3 | val coalesceRDD = rdd.coalesce(7,true)
4 | println("重新分区后的分区个数:"+coalesceRDD.partitions.size)
5 | println("RDD依赖关系:"+coalesceRDD.toDebugString)
```

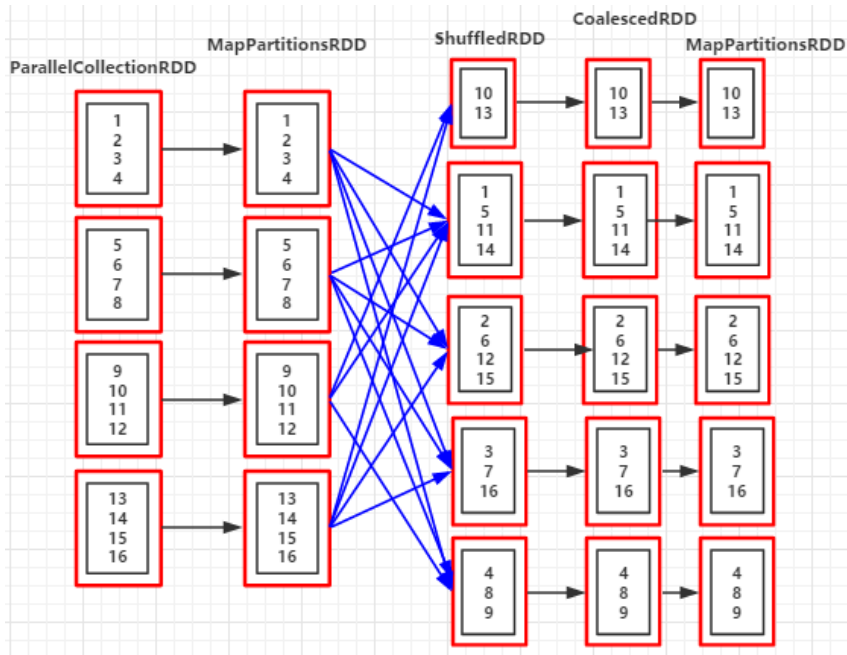
输出:

```
重新分区后的分区个数:5
RDD依赖关系: (5) MapPartitionsRDD[4] at coalesce at Coalesce.scala:14 []
| CoalescedRDD[3] at coalesce at Coalesce.scala:14 []
| ShuffledRDD[2] at coalesce at Coalesce.scala:14 []
+-(4) MapPartitionsRDD[1] at coalesce at Coalesce.scala:14 []
| ParallelCollectionRDD[0] at parallelize at Coalesce.scala:13 []
//分区后的数据集
List(10, 13)
List(1, 5, 11, 14)
List(2, 6, 12, 15)
List(3, 7, 16)
List(4, 8, 9)
```

(RDD依赖图:coalesce(3,flase))



(RDD依赖图:coalesce(3,true))



**11.repartition(numPartition):**是函数coalesce(numPartition,true)的实现，效果和例9.1的coalesce(numPartition,true)的一样

**12.glom():**将RDD的每个分区中的类型为T的元素转换成数组Array[T]

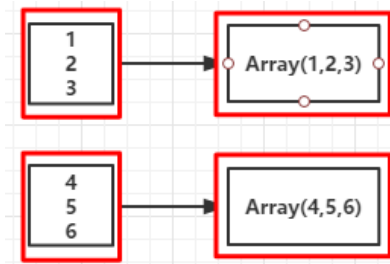
```

1 //省略
2 val rdd = sc.parallelize(1 to 16,4)
3 val glomRDD = rdd.glom() //RDD[Array[T]]
4 glomRDD.foreach(rdd => println(rdd.getClass.getSimpleName))
5 sc.stop

```

输出:

```
int[] //说明RDD中的元素被转换成数组Array[Int]
```



**13.randomSplit(weight:Array[Double],seed):**根据weight权重值将一个RDD划分成多个RDD,权重越高划分得到的元素较多的几率就越大

```

1 //省略sc
2 val rdd = sc.parallelize(1 to 10)
3 val randomSplitRDD = rdd.randomSplit(Array(1.0,2.0,7.0))
4 randomSplitRDD(0).foreach(x => print(x + " "))
5 randomSplitRDD(1).foreach(x => print(x + " "))
6 randomSplitRDD(2).foreach(x => print(x + " "))
7 sc.stop

```

输出:

```
2 4
3 8 9
1 5 6 7 10
```

14

0

以上例子源码地址: <https://github.com/Mobin-F/SparkExample/tree/master/src/main/scala/com/mobin/SparkRDDFun/Transformation/KVRDD>