

```

1  1.回顾
2      Spark通用计算引擎，速度快、通用性、运行方式多，易用性。
3
4      core
5      sql
6      streaming
7      mllib
8      graphx
9
10     集群部署模式
11     -----
12         local
13         standalone
14         yarn
15         mesos
16
17     RDD
18     -----
19     resilient distributed dataset,弹性分布式数据集。
20     java集合(List)。
21
22     [scala]
23     val conf = new SparkConf()
24     conf.setAppName("wc")
25     conf.setMaster("spark://s101:7077")
26     val sc = new SparkContext(conf)
27
28     val rdd1 = sc.textFile("hdfs://mycluster/user/centos/1.txt")
29     val rdd2 = rdd1.flatMap(_.split(" "))
30     val rdd3 = rdd2.map((_,1))
31     val rdd4 = rdd3.reduceByKey(_+_ )
32     val res = rdd4.collect();
33
34     [java]
35     JavaSparkContext
36     JavaRDD
37     JavaPairRDD
38
39
40     List<Tuple2<String, Tuple2<String, Integer>>> list = new ArrayList<Tuple2<String,
41     Tuple2<String, Integer>>>() ;
42     String busid = t._1._1 ;
43     String tag = t._1._2 ;
44     Integer count = t._2() ;
45
46     list.add(new Tuple2<String, Tuple2<String, Integer>>(busid ,new Tuple2<String,
47     Integer>(tag, count)));
48     return null ;
49
50     Spark实现标签生成
51     -----
52     1.scala版
53         import com.oldboy.spark.Util.TagUtil
54         import org.apache.spark.{SparkConf, SparkContext}
55
56         import scala.collection.mutable.ArrayBuffer
57
58         /**
59          * Created by Administrator on 2018/5/10.
60          */
61         object TaggenScala {
62             def main(args: Array[String]): Unit = {
63                 //1.创建SparkContext
64                 val conf = new SparkConf
65                 conf.setAppName("taggen")
66                 conf.setMaster("local")
67                 val sc = new SparkContext(conf)
68
69                 //2.加载文件
70                 val rdd1 = sc.textFile("d:/mr/temptags.txt")

```

```

71
72 //3.压扁成RDD[(busid,tag)]
73 val rdd2 = rdd1.flatMap(line=>{
74     val buf = ArrayBuffer[Tuple2[Tuple2[String, String], Int]]()
75     val arr = line.split("\t")
76     val busid = arr(0)
77     val json = arr(1)
78     //
79     import scala.collection.JavaConversions._
80     val tags = TagUtil.extractTags(json)
81     for(tag <- tags){
82         buf.+=((busid , tag),1))
83     }
84     buf
85 })
86
87 //4.聚合
88 val rdd3 = rdd2.reduceByKey(_+_ )
89
90 //5.变换组合成(busid, List((tag,count)))
91 val rdd4 = rdd3.map(t=>{
92     (t._1._1, (t._1._2 , t._2)::Nil)
93 })
94
95 //6.聚合
96 val rdd5 = rdd4.reduceByKey((a,b)=> a ::: b)
97
98 //7.商家内评论倒排序并取前三
99 val rdd6 = rdd5.mapValues(list=>{
100     list.sortBy(_._2).take(3)
101 })
102
103 val rdd7 = rdd6.sortBy(t=> -t._2(0)._2)
104
105 rdd7.foreach(println)
106 }
107 }
108

```

## 2.java版

```

110 package com.oldboy.spark.java;
111
112 import com.oldboy.spark.Util.TagUtil;
113 import org.apache.spark.SparkConf;
114 import org.apache.spark.api.java.JavaPairRDD;
115 import org.apache.spark.api.java.JavaRDD;
116 import org.apache.spark.api.java.JavaSparkContext;
117 import org.apache.spark.api.java.function.*;
118 import scala.Tuple2;
119 import scala.Tuple4;
120
121 import java.util.ArrayList;
122 import java.util.Comparator;
123 import java.util.Iterator;
124 import java.util.List;
125
126 /**
127  * 统计气温数据
128  */
129 public class TaggenJava {
130     public static void main(String[] args) {
131         SparkConf conf = new SparkConf();
132         conf.setAppName("TaggenJava");
133         conf.setMaster("local") ;
134
135         JavaSparkContext sc = new JavaSparkContext(conf);
136
137         //1.加载文件
138         JavaRDD<String> rdd1 = sc.textFile("d:/mr/temptags.txt");
139
140         //2.压扁
141         JavaPairRDD<Tuple2<String, String>, Integer> rdd2 =
142             rdd1.flatMapToPair(new PairFlatMapFunction<String,
143                 Tuple2<String,String>, Integer>() {

```

```

142         public Iterator<Tuple2<Tuple2<String, String>, Integer>>
call(String s) throws Exception {
143             List<Tuple2<Tuple2<String, String>, Integer>> list = new
ArrayList<Tuple2<Tuple2<String, String>, Integer>>();
144             String[] arr = s.split("\t");
145             String busid = arr[0];
146             String json = arr[1];
147             List<String> tags = TagUtil.extractTags(json);
148             for(String tag : tags){
149                 list.add(new Tuple2<Tuple2<String, String>, Integer>(new
Tuple2<String, String>(busid,tag) , 1));
150             }
151             return list.iterator();
152         }
153     });
154
155     //3. 聚合((busid,tag),count)
156     JavaPairRDD<Tuple2<String, String>, Integer> rdd3
=rdd2.reduceByKey(new Function2<Integer, Integer, Integer>() {
157         public Integer call(Integer v1, Integer v2) throws Exception {
158             return v1 + v2;
159         }
160     });
161
162     //4. 交换元组形成新元组 (busid , List((tag,count)))
163     JavaPairRDD<String, List<Tuple2<String, Integer>>> rdd4 =
rdd3.mapToPair(new
PairFunction<Tuple2<Tuple2<String,String>,Integer>, String,
List<Tuple2<String,Integer>>>() {
164         public Tuple2<String, List<Tuple2<String, Integer>>>
call(Tuple2<Tuple2<String, String>, Integer> t) throws Exception {
165             List<Tuple2<String,Integer>> list = new
ArrayList<Tuple2<String, Integer>>();
166             list.add(new Tuple2<String, Integer>(t._1._2() , t._2()));
167             return new Tuple2<String, List<Tuple2<String,
Integer>>>(t._1._1 , list);
168         }
169     });
170
171     //5. 聚合, (busid,List((tag,count),(tag1,count1)))
172     JavaPairRDD<String, List<Tuple2<String, Integer>>> rdd5
=rdd4.reduceByKey(
173         new Function2<List<Tuple2<String, Integer>>,
List<Tuple2<String, Integer>>, List<Tuple2<String,
Integer>>>() {
174             public List<Tuple2<String, Integer>>
call(List<Tuple2<String, Integer>> v1,
List<Tuple2<String, Integer>> v2) throws Exception {
175                 v1.addAll(v2);
176                 return v1;
177             }
178         });
179
180     //6. 排序
181     JavaPairRDD<String, List<Tuple2<String, Integer>>> rdd6 =
rdd5.mapValues(new Function<List<Tuple2<String,Integer>>,
List<Tuple2<String, Integer>>>() {
182         public List<Tuple2<String, Integer>> call(List<Tuple2<String,
Integer>> v1) throws Exception {
183             v1.sort(new Comparator<Tuple2<String, Integer>>() {
184                 public int compare(Tuple2<String, Integer> o1,
Tuple2<String, Integer> o2) {
185                     return -(o1._2() - o2._2());
186                 }
187             });
188             return new ArrayList<Tuple2<String, Integer>>(v1.size() > 4
? v1.subList(0,5) : v1);
189         }
190     });
191     //7. 全排序
192     JavaRDD<Tuple2<String, List<Tuple2<String, Integer>>>>
rdd7=rdd6.map(new
Function<Tuple2<String,List<Tuple2<String,Integer>>>, Tuple2<String,

```

```

193         List

```

## Spark创建RDD方式

- 1.textFile  
    文本加载
- 2.通过RDD进行特定变换  
    rdd.map()  
    rdd.flatMap()  
    rdd.filter()  
    ...
- 3.parallelize  
    创建静态rdd.
- 4.

## RDD的算子

RDD算子包括变换和动作两种类型，只要返回值是rdd就是变换，反之大部分action. 最终通过是否调用SparkContext的runJob方法来界定。 变换不会触发job的执行，只有action方法调用时才会触发job的执行。

[transformation]	
map	//变换
flatMap	//压扁
filter	//过滤
reduceByKey()	//聚合
groupByKey	//分组
mapPartitions	//*****对分区进行变换.*****
mapPartitionsWithIndex	//*****对分区进行变换.*****
union	//同型操作,等价于sql union,不需要
intersect	//交集,需要shuffle ,
distinct	//交集,需要shuffle
groupByKey()	//变换
aggregateByKey()	//按照key进行聚合，可以改变类型。
join	//shuffle，按照key连接，等价于sql的连接操作
leftOuterJoin	
rightOuterJoin	
fullOuterJoin	
cogroup	//协分组，对两个rdd的value进行聚合成二元组
cartesian	//笛卡尔积，交叉连接。不需要shuffle
repartition	
//再分区,不论是增加还是减少分区，都需要跨网络进行数据shuffle	
	//处理，负载均衡。一旦涉及shuffle，性能会有影响。
coalesce	//*****

//改变分区，需要携带shuffle=true参数来控制是否需要shuffle。  
。

//默认false，不进行shuffle，只能用于减少分区，repartition调用  
//调用的就是该方法，但是shuffle设置为true。

## 解决数据倾斜

-----

### 1.重新设计key

```
import org.apache.spark.{SparkConf, SparkContext}

import scala.util.Random

/**
 * Created by Administrator on 2018/5/8.
 */
object DataLeanScala {
  def main(args: Array[String]): Unit = {
    //1.创建spark配置对象
    val conf = new SparkConf()
    conf.setAppName("wcApp")
    conf.setMaster("local")

    val sc = new SparkContext(conf)

    val rdd1 = sc.textFile("d:/mr/1.txt")
    val rdd2 = rdd1.flatMap(line=>{
      var arr = line.split(" ")
      arr
    })
    val rdd3= rdd2.map((_,1))

    val rdd4 = rdd3.map(t=>{
      val n = Random.nextInt(3)
      (t._1 + "_" + n , t._2)
    })

    val rdd5 = rdd4.reduceByKey(_+_ )
    val rdd6 = rdd5.map(t=>{
      var arr = t._1.split("_")
      (arr(0),t._2)
    })
    val rdd7 = rdd6.reduceByKey(_+_ )
    rdd7.collect().foreach(println)

    while(true){
      Thread.sleep(1000)
    }
  }
}
```

### 2.重新设计分区类。

```
import org.apache.spark.{HashPartitioner, Partitioner, SparkConf,
SparkContext}

import scala.util.Random

/**
 * Created by Administrator on 2018/5/8.
 */
object DataLean2Scala {
  def main(args: Array[String]): Unit = {
    //1.创建spark配置对象
    val conf = new SparkConf()
    conf.setAppName("wcApp")
    conf.setMaster("local")

    val sc = new SparkContext(conf)
```

```

326
327         val rdd1 = sc.textFile("d:/mr/1.txt")
328         val rdd2 = rdd1.flatMap(line=>{
329             var arr = line.split(" ")
330             arr
331         })
332         val rdd3= rdd2.map((_,1))
333
334         class MyPartitioner(n:Int) extends Partitioner{
335             override def numPartitions: Int = n
336
337             override def getPartition(key: Any): Int = {
338                 Random.nextInt(n)
339             }
340         }
341
342         val rdd4 = rdd3.reduceByKey(new MyPartitioner(4) , (a,b)=>{a + b})
343         //需要重新指定分区类对象，否则使用之前的分区对象
344         val rdd5 = rdd4.reduceByKey(new HashPartitioner(3) , (a,b)=>a + b)
345
346         rdd5.foreach(println)
347
348         while(true){
349             Thread.sleep(1000)
350         }
351     }
352 }
353
354
355
356

```

#### Spark aggregateByKey()

```

357 -----
358
359     aggregateByKey(zeroValue: U)(seqOp: (U, V) => U, combOp: (U, U) => U): RDD[(K, U)]
360     1.(zeroValue: U)
361         初始值
362     2.(seqOp: (U, V) => U, combOp: (U, U) => U)
363
364     3.实现每个年度下，各个分区内的气温最大值列表
365         import java.sql.DriverManager
366
367         import org.apache.spark.{SparkConf, SparkContext}
368
369         /**
370          * Created by Administrator on 2018/5/10.
371          */
372         object RDDTestScala {
373
374             //返回当前线程
375             def tname() = {
376                 Thread.currentThread().getName
377             }
378
379             def main(args: Array[String]): Unit = {
380                 //1.创建SparkContext
381                 val conf = new SparkConf
382                 conf.setAppName("taggen")
383                 conf.setMaster("local[4]")
384                 val sc = new SparkContext(conf)
385
386                 //         val rdd1 = sc.parallelize(Array(1,2,2,3) , 3)
387                 //         val rdd2 = sc.parallelize(Array(2, 2, 3, 4), 3)
388                 //         val rdd3 = rdd1.intersection(rdd2 , 2)
389                 //         //         val rdd3 = rdd1.union(rdd2)
390                 //         println("3 ): " + rdd3.partitions.length)
391                 //         //         val rdd4 = rdd3.distinct(3)
392                 //         println("4 ): " + rdd4.partitions.length)
393                 //         println(rdd4.partitions.length)
394                 //         rdd3.collect().foreach(println)
395
396
397                 //         val rdd1 =
398                 sc.parallelize(Array("1.1"->"tom","1.1"->"tomas","1.2"->"bob","1.2"->"alice"))

```

```

398         //      val rdd2 = rdd1.aggregateByKey()
399         //      rdd2.collect()
400
401         val rdd1 = sc.textFile("d:/mr/temp15.txt", 3)
402         val rdd2 = rdd1.map(line=>{
403             val arr = line.split(" ")
404             (arr(0).toInt , arr(1).toInt)
405         })
406
407         //初始值,在每个分区内按key聚合,
408         val zero:List[Int] = Nil
409         //分区内按key聚合的聚合函数 (U,V) => U
410         def seqOp(a:List[Int] , b:Int) = {
411             if(a.isEmpty){
412                 b :: a
413             }
414             else{
415                 import scala.math._
416                 max(a(0),b) :: Nil
417             }
418         }
419         //shuffle之后,对key下的shuffle过来的U进行聚合, (U,U)=>U
420         def combOp(a:List[Int],b:List[Int]) ={
421             a:::b
422         }
423         val rdd3 = rdd2.mapPartitionsWithIndex((index, it)=>{
424             for(t <- it){
425                 println(index + " : " + t )
426             }
427             it
428         })
429         rdd3.collect().foreach(println)
430         println("=====")
431         val rdd4 = rdd2.aggregateByKey(zero)(seqOp , combOp )
432         //练习: 提取每个年度下,每个分区内的最高气温和最低温形成的元组列表。
433         //加强题:0-(20,-19),1-(,)
434         rdd4.collect().foreach(println)
435
436         while(true){
437             Thread.sleep(1000)
438         }
439     }
440 }

```

## 作业

-----

- 1.aggregateByKey改成java版实现气温统计。
- 2.统计实现每年内,在每个分区下最高最低气温形成列表
- 3.统计实现每年内,在每个分区下最高最低气温形成列表并携带分区索引信息。
- 4.java版实现上述2,3两题。

```

[action]
collect           //Array
first             //take(1)
take              //
foreach()         //循环
reduce()          //T
save...           //保存.

```

## spark分区

-----

hadoop中的切片 == spaark.rdd的分区

## 启动spark-shell

-----

```

spark-shell --master local[4] //启动4个线程模拟worker
spark-shell --master local[*] //*匹配CPU内核数

```

## Spark常用变换

-----

map

```

471     flatMap
472     mapValues()
473     filter() //
474
475
476 spark运行
477 -----
478     spark-submit --class ... --master ... xxx.jar
479
480
481 List<Integer> = 12345
482 list2
483 for(Integer n : list){
484     list.add(n * 2);
485 }
486
487 list.map(_ * 2)
488
489 spark相关概念
490 -----
491     1.RDD
492         resilient distributed dataset,弹性(容错)分布式数据集,spark的基本抽象。
493         代表不可变、分区化Partition的元素集合,可以进行并行计算。
494         包含可用于所有RDD的基本操作,例如map、filter、persist等。
495         PairRDDFunctions包含针对kv对类型的RDD的操作。
496         轻量级。
497         [内部表现为5个主要属性]
498         a.分区列表
499
500         b.作用于每个切片的计算函数
501
502         c.到其他RDD的依赖列表
503             lineage , 血统
504
505         e.(可选)针对kv的RDD的分区类
506
507         f.(可选)计算每个切片的首选位置列表(hdfs)。
508
509         g.HadoopRDD
510             该RDD为用于读取hdfs文件提供了核心功能。
511
512         h.hadoop.textFile()
513             textFile() -> HadoopRDD(inputFormat,k,v,minParts)
514     2.SparkContext
515         spark程序的主入口点,代表到cluster的连接,可以创建RDD、累加器、广播变量、
516         每个JVM只有一个active的SparkContext。在创建新的上下文之前必须调用stop停止
517         当前的上下文。
518         Spark调度机制是三级调度,DagScheduler -> TaskScheduler -> BackendSheduler
519
520         创建SparkContext -> createTaskScheduler -> {new TaskSchedulerImpl +
521         match(local) => SchedulerBackend}
522
523         LocalSchedulerBackend -> local[n]
524         //local模式下
525         sc.textFile(path , minPartitions = defaultMinPartitons=>
526         min((get("spark.default.parallelism" , totleCores)) ,2))
527
528
529
530     3.DAGScheduler
531         Direct acycline graph,有向无环图调度器。
532         面向stage的高级调出层。计算每个job的stage DAG,跟踪RDD和stage,寻找最短路径。
533
534         提交taskset给下层的调度器(task调度器),taskset包含的是完全独立的任务(可以在clus
535         ter直接运行)。
536         stage划分以shuffle边界,具有窄依赖的RDD操作(map |
537         filter)进入同一taskset的管线中。
538         Shuffle依赖的操作需要多个Stage,上一个的输出做下一个的输入。
539         每个Stage只有一次依赖于其他stage的shuffle操作。
540         DAG调度器决定task的首选运行位置,

```



DAG调度器处理因为shuffle输出文件丢失导致的故障。不是由shuffle输出文件丢失引发的故障由task调度器处理。  
取消整个stage之前重试少量次数。

job :ActiveJob表示, 最顶层的工作单元。action发生时, 提交的activeJob  
stage  
:task集合, 计算中间结果。同一RDD的每个分区都执行相同的计算过程。stage以shuffle  
边界作为划分标准。  
task :工作单元, 每个task发送一个machine。

#### 4.ActiveJob

DAG调度器中运行的job, job类型有两种。  
一种是result job, 执行action, 对ResultStage进行计算。  
另一个是MapStage job, 计算shuffleMapStage的map输出。  
使用finalStage字段来区分是哪种类型。

#### 5.Stage

并行执行的task集合, 所有task运行同一函数。所有task都有相同的shuffle依赖, DAG调度器以拓扑顺序执行stage。

stage类型有两种,  
一种是ShuffleMapStage, 的输出是下一个stage的输入。  
另一种是ResultStage, 通过执行一个rdd的函数运行action。  
每个stage都有一个firstJobid, 标识第一个提交的stage, 每个阶段都对应一个rdd。  
Stage类型有两种:

5.1) ShuffleMapStage  
在RDD的某些分区上执行函数计算结果

5.2) ResultStage  
在RDD的某些分区上执行函数计算结果

#### 6.Task

Spark执行单元, 两种类型:ShuffleMapTask + ResultTask。

Job由多个Stage构成, 每个job的最后stage由多个ResultTask构成, 之前的stage由多个ShuffleMapTask构成。

ResultTask执行task并将结果回传到Driver (执行入口点)。  
ShuffleMapTask执行task并按照分区类将task的输出划分到多个bucket (分区) 中。

#### 7.TaskScheduler

底层任务调度器接口, 当前只有TaskSchedulerImpl实现。  
该接口可插拔, 可以使用不同的实现类。每个任务调度器为一个SparkContext调度任务。  
该调度器从上层调度器 (DagScheduler) 得到提交过来的task set, 并发送任务给集群。  
运行在driver端,

通过后台调度器为多种集群类型调度任务。是过渡阶段。是以taskset为单位进行调度。

#### 8.Executor

spark执行程序, 通过线程池运行task。

#### 9.依赖

Dependency。  
RDD的分区和上级RDD分区之间的对应关系。

[窄依赖]  
子RDD的每个分区依赖于父RDD的少量分区。  
OneToOneDependency  
RangeDependency  
PruneDependency

[宽依赖]  
需要shuffle,  
ShuffleDependency

#### 10.driver

驱动。  
client : 入口程序。

spark job提交流程[local]

并发度 :

```

604 -----
605     最大并发度，等价于切片数(分区数)
606     sc.textFile("",10) ;
607
608
609 并发能力
610 -----
611     并行线程数决定。
612     conf.setMaster("local[4]")
613
614
615 spark的textFile() 默认并发度计算
616 -----
617     1.sc.textFile(,,)
618         def defaultMinPartitions: Int = math.min(defaultParallelism, 2)
619
620     2.sc.parallelize(1 to 10,3)
621         defaultParallelism = scheduler.conf.getInt("spark.default.parallelism",
622             totalCores)
623
624 后台调度器的创建时间
625 -----
626     sc.createTaskScheduler()
627         -->new LocalSchedulerBackend()
628
629 再分区
630 -----
631     shuffle动作可以重新指定分区。
632     reduceByKey(f, n) ;
633     groupByKey(n) ;
634
635 map变换的操作
636 -----
637     1.mapPartitions
638         //映射分区
639         val rdd1 = sc.parallelize(1 to 10 , 2)
640         val rdd2 = rdd1.mapPartitions(it => {
641             println("start")
642             var list:List[Int] = Nil
643             for(e <- it){
644                 list = (e * 2)::list
645             }
646             println("end")
647             list.iterator
648         } )
649
650         //数据入库问题,可以使用foreachPartition
651         val rdd2 = rdd1.mapPartitions(it=>{
652             Class.forName("com.mysql.jdbc.Driver")
653             val url = "jdbc:mysql://localhost:3306/big9"
654             val user = "root"
655             val pass = "root"
656             val conn = DriverManager.getConnection(url,user,pass)
657             conn.setAutoCommit(false);
658             val sql ="insert into tt(id) values(?)" ;
659             val ppst = conn.prepareStatement(sql)
660             for(e <- it){
661                 ppst.setInt(1, e)
662                 ppst.executeUpdate()
663             }
664             conn.commit()
665             it
666         })
667
668     2.union
669         联合的分区数是父RDD的分区数总和.
670         rdd1.union(rdd2) ;
671
672     3.intersection
673         交集,需要shuffle。采用协分组。
674         val rdd3 = rdd1.intersection(rdd2,4)
675         println(rdd3.getNumPartitions)

```

```

676         val rdd4 = rdd3.mapPartitionsWithIndex((n,it)=>{
677             for(e <- it){
678                 println(n + " : " + e)
679             }
680             it
681         })
682
683 4.distinct
684     去重，内部通过reduceByKey实现。
685     map(x => (x, null)).reduceByKey((x, y) => x, numPartitions).map(_._1)
686
687 5.groupByKey
688     需要shuffle。
689     val rdd3 = rdd1.map(e=>{
690         (if(e%2==0) "even" else "odd" , e)
691     })
692     val rdd4 = rdd3.groupByKey(4);
693
694 6.reduceByKey()
695     可以指定分区数。
696
697
698 7.aggregateByKey
699     按key聚合，
700     zeroU是初始值。
701     该函数现在分区内进行预聚合，使用函数f1。
702     分区间聚合使用f2函数。
703
704     val zeroU = "X"
705     def f1(a:String,b:Int) :String = a + "(f1)" + b
706     def f2(a:String,b:String) = a + "(f2)" + b
707     val rdd4 = rdd3.aggregateByKey(zeroU) (f1,f2)
708
709 8.join
710     leftOuterJoin
711     rightOuterJoin
712     fullOuterJoin
713
714 9.cogroup
715     协分组。
716     对当前的rdd和传递的rdd参数进行联合分组，将相同的key的v组合成元组。
717     rdd1.cogroup(rdd2)
718     (K, (List[V1],List[V2]))
719
720 10.groupByKey
721     对当前的rdd进行按key分组。
722
723 11.cartesian
724     不需要shuffle。
725     for (x <- rdd1.iterator(currSplit.s1, context);
726         y <- rdd2.iterator(currSplit.s2, context)) yield (x, y)
727
728 12.coalesce
729     再分区。
730     推荐不进行shuffle来降低分区数，产生新的rdd。
731     如果增加分区数，必须shuffle=true，否则无法实现增加目的。
732
733 13.repartition
734     恒使用shuffle操作，对于减少分区的动作，推荐使用coalesce。
735     coalesce(numPartitions, shuffle = true)
736
737 14.repartitionAndSortWithinPartitions
738     分区带排序
739     rdd1.repartitionAndSortWithinPartitions(new HashPartitioner(3))
740
741 groupByKey和reduceByKey
742 -----
743     groupByKey不会进行map端combine操作。
744     reduceByKey() 默认进行map端combine操作，减少网络负载。
745
746
747
748

```

1.使用aggregateByKey函数实现groupByKey功能，将所有v放入到集合List中.

```
val zeroU:List[Int] = Nil
def f1(a:List[Int] , b:Int) = b :: a
def f2(a:List[Int],b:List[Int]) = a ++ b
```

2.准备cust.txt和orders.txt,  
使用spark实现如下计算

2.1) 查询每个客户的订单编号集合

集合形式:(1,tomas,List(No001,No002))

2.2) 查询每个客户订单价格的总和

集合形式:(1,tomas,300.4)

2.3) 查询没有订单的客户

```
import org.apache.spark.{SparkConf, SparkContext}
```

```
/**
```

```
 * Created by Administrator on 2018/3/2.
```

```
 */
```

```
object SimpleApp3 {
```

```
  case class Customer(id:Int,name:String,age:Int)
```

```
  case class Order(id:Int,orderno:String,price:Float,cid:Int)
```

```
  def main(args: Array[String]): Unit = {
```

```
    val conf = new SparkConf()
```

```
    conf.setAppName("app")
```

```
    conf.setMaster("local[3]")
```

```
    val sc = new SparkContext(conf)
```

```
    val crdd1 = sc.textFile("file:///d:/java/custs.txt")
```

```
    val ordd1 = sc.textFile("file:///d:/java/orders.txt")
```

```
    val crdd2 = crdd1.map(e=>{
```

```
      val arr = e.split(",")
```

```
      (arr(0).toInt , Customer(arr(0).toInt ,arr(1),arr(2).toInt))
```

```
    })
```

```
    val ordd2 = ordd1.map(e=>{
```

```
      val arr = e.split(",")
```

```
      (arr(3).toInt, Order(arr(0).toInt, arr(1), arr(2).toFloat ,  
        arr(3).toInt))
```

```
    })
```

```
    //连接
```

```
    val ardd1 = crdd2.join(ordd2)
```

```
      .map(t=>(t._2._1,t._2._2.orderno))
```

```
      .groupByKey().map(t=>(t._1.id,t._1.name,t._2.toList))
```

```
    //ordd2.foreach(println)
```

```
    ardd1.foreach(println)
```

```
    while(true){
```

```
      Thread.sleep(1000)
```

```
    }
```

```
  }
```

```
}
```