

结合Spark源码分析, combineByKey, aggregateByKey, foldByKey, reduceByKey

置顶 2016年06月16日 09:42:03

阅读数 : 5243

转载请标明出处：小帆的帆的专栏

combineByKey

```
1 def combineByKey[C](  
2   createCombiner: V => C,  
3   mergeValue: (C, V) => C,  
4   mergeCombiners: (C, C) => C): RDD[(K, C)] = self.withScope {  
5   combineByKeyWithClassTag(createCombiner, mergeValue, mergeCombiners)(null)  
6 }
```

- createCombiner : 当combineByKey第一次遇到值为k的Key时，调用createCombiner函数，将v转换为c
- mergeValue : combineByKey不是第一次遇到值为k的Key时，调用mergeValue函数，将v累加到c中
- mergeCombiners : 将两个c，合并成一个

```
1 // 实例  
2 SparkConf conf = new SparkConf().setAppName("test").setMaster("local");  
3 JavaSparkContext sc = new JavaSparkContext(conf);  
4  
5 List<Tuple2<Integer, String>> list = new ArrayList<>();  
6  
7 list.add(new Tuple2<>(1, "www"));  
8 list.add(new Tuple2<>(1, "iteblog"));  
9 list.add(new Tuple2<>(1, "com"));  
10 list.add(new Tuple2<>(2, "bbs"));  
11 list.add(new Tuple2<>(2, "iteblog"));  
12 list.add(new Tuple2<>(2, "com"));  
13 list.add(new Tuple2<>(3, "good"));  
14  
15 JavaPairRDD<Integer, String> data = sc.parallelizePairs(list);  
16  
17 JavaPairRDD<Integer, List<String>> result = data.combineByKey(v -> {  
18   ArrayList<String> strings = new ArrayList<>();  
19   strings.add(v);  
20   return strings;
```

```
21  }, (c, v) -> {  
22      c.add(v);  
23      return c;  
24  }, (c1, c2) -> {  
25      c1.addAll(c2);  
26      return c1;  
27  });  
28  
29  result.collect().foreach(System.out::println);
```

aggregateByKey

```
1  def aggregateByKey[U: ClassTag](zeroValue: U, partitioner: Partitioner)(seqOp: (U,  
2    combOp: (U, U) => U): RDD[(K, U)] = self.withScope {  
3  
4    // 中间代码省略, 主要看最后一个, 调用combineByKey  
5  
6    val cleanedSeqOp = self.context.clean(seqOp)  
7    // seqOp, 同时是, createCombiner, mergeValue。而combOp是mergeCombiners  
8    combineByKeyWithClassTag[U]((v: V) => cleanedSeqOp(createZero(), v),  
9      cleanedSeqOp, combOp, partitioner)  
10 }
```

- createCombiner : cleanedSeqOp(createZero(), v)是createCombiner, 也就是传入的seqOp函数, 只不过其中一个值是传入的zeroValue
- mergeValue : seqOp函数同样是mergeValue, createCombiner和mergeValue函数相同是aggregateByKey函数的关键
- mergeCombiners : combOp函数

因此, 当createCombiner和mergeValue函数的操作相同, aggregateByKey更为合适

```
1  // 例子与combineByKey相同, 只是改用aggregateByKey实现  
2  SparkConf conf = new SparkConf().setAppName("test").setMaster("local");  
3  JavaSparkContext sc = new JavaSparkContext(conf);  
4  
5  List<Tuple2<Integer, String>> list = new ArrayList<>();  
6  
7  list.add(new Tuple2<>(1, "www"));  
8  list.add(new Tuple2<>(1, "iteblog"));  
9  list.add(new Tuple2<>(1, "com"));  
10 list.add(new Tuple2<>(2, "bbs"));  
11 list.add(new Tuple2<>(2, "iteblog"));  
12 list.add(new Tuple2<>(2, "com"));
```

```

13 list.add(new Tuple2<>(3, "good"));
14
15 JavaPairRDD<Integer, String> data = sc.parallelizePairs(list);
16
17 JavaPairRDD<Integer, List<String>> result = data.aggregateByKey(new ArrayList<Stri
18     c.add(v);
19     return c;
20 }, (Function2<List<String>, List<String>, List<String>>) (c1, c2) -> {
21     c1.addAll(c2);
22     return c1;
23 });
24
25 result.collect().foreach(System.out::println);

```

foldByKey

```

1 def foldByKey(
2     zeroValue: V,
3     partitioner: Partitioner)(func: (V, V) => V): RDD[(K, V)] = self.withScope {
4
5     // 中间代码省略, 主要看最后一个, 调用combineByKey
6
7     val cleanedFunc = self.context.clean(func)
8     // 传入的func函数, 同时是, createCombiner, mergeValue, mergeCombiners
9     // createCombiner函数传入了零值, 首次遇到一个key时, 根据零值进行初始化
10    combineByKeyWithClassTag[V]((v: V) => cleanedFunc(createZero(), v),
11        cleanedFunc, cleanedFunc, partitioner)
12 }

```

- createCombiner : cleanedFunc(createZero(), v)是createCombiner, 也就是func函数, 只不过其中一个值是传入的zeroValue
- mergeValue, mergeCombiners : func函数也是mergeValue和 mergeCombiners

当createCombiner , mergeValue和mergeCombiners函数操作都相同, 唯独需要一个zeroValue时, 适用

```

1 // 根据Key把Value相加, 但是不从0开始, 设置初始值为100
2 val conf = new SparkConf().setAppName("test").setMaster("local")
3 val sc = new SparkContext(conf)
4 val sqlContext = new SQLContext(sc)
5
6 var rdd = sc.makeRDD(Array(("A",0),("A",2),("B",1),("B",2),("C",1)))
7

```

8

```
rdd.foldByKey(100)(_+_).collect.foreach(println)
```

reduceByKey

```
1 def reduceByKey(partitioner: Partitioner, func: (V, V) => V): RDD[(K, V)] = self.w
2   combineByKeyWithClassTag[V]((v: V) => v, func, func, partitioner)
3 }
```

- createCombiner : 与foldByKey相比, reduceByKey没有初始值, createCombiner也没有调用func函数, 而是直接将参数作为返回值返回了,
- mergeValue, mergeCombiners : func函数同时是mergeValue和 mergeCombiners

当不需要createCombiner , 且mergeValue和mergeCombiners函数操作都相同时, 适用

```
1 val conf = new SparkConf().setAppName("test").setMaster("local")
2 val sc = new SparkContext(conf)
3 val sqlContext = new SQLContext(sc)
4
5 var rdd = sc.makeRDD(Array(("A", 0), ("A", 2), ("B", 1), ("B", 2), ("C", 1)))
6
7 rdd.reduceByKey(_ + _).collect.foreach(println)
```

总结

这几个算子, 核心就要弄明白combineByKey, 其他三个都是调用它. 上文主要也是从combingByKey传入的三个函数的角度在分析.

而在实际运用中, 最先要考虑的应该是类型. combiningByKey和aggregateByKey输入跟输出的类型可以不一致, 而foldByKey和reduceByKey不行. 类型确定后再根据自己的业务选择最简洁的算子.