

作业2

09_茜

August 2019

1 RDD算子实现PAT Ranking

```

1 //spark
2 info = [(("00002", 2, 12), ("00007", 4, 17), ("00005", 1, 19), ("00007", 2, 25), ("00005", 1, 20), ("00002", 2, 2), ("00005", 1, 15), ("00001", 1, 18), ("00005", 3, 22), ("00006", 4, -1), ("00001", 2, 18), ("00002", 1, 20), ("00004", 1, 15), ("00002", 4, 18), ("00001", 3, 4), ("00001", 4, 2), ("00005", 2, -1), ("00004", 2, 0
3 ))
4 high_marks = sc.parallelize(info).map(lambda line : ((line[0], line[1]), line[2])).reduceByKey(max)
5 marks = sc.parallelize(high_marks.collect()).map(lambda line : (line[0][0], (line[0][1], line[1])))
6 zeroValue = []
7 mergeVal = (lambda aggregated, el: aggregated + [el] if el[1][0] != -1 else aggregated + [(el[0], '-')]
8 )
9 y = marks.aggregateByKey(zeroValue, mergeVal, mergeComb).collect()
10 for element in y:
11     full_set = set([1, 2, 3, 4])
12     s = set()
13     if (len(element[1]) < 4):
14         for i in range(0, len(element[1])):
15             s.add(element[1][i][0])
16             missing = full_set - s
17             for e in missing:
18                 element[1].append((e, '-'))
19
20 def takeFirst(elem):
21     return elem[0]
22
23 for element in y:
24     element[1].sort(key=takeFirst)
25
26 def add(ele):
27     sum = 0
28     for i in range(0, 4):
29         if (ele[1][i][1] != '-'):
30             sum += ele[1][i][1]
31     return sum
32
33 for ele in y:
34     sum = add(ele)
35     ele[1].insert(0, (0, sum))
36
37 result = sc.parallelize(y).filter(lambda line : line[1][0][1] != 0).collect()
38 for ele in result:
39     values = [str(x[1]) for x in ele[1]]
40     print(ele[0]+' '+ ' '.join(values))
41
42 00002 63 20 25 _ 18
43 00007 42 _ 25 _ 17
44 00005 42 20 _ 22 _
45 00001 42 18 18 4 2
46 00004 40 15 0 25 _

```

Figure 1: RDD实现PAT Ranking