

作业2

姓名：吴双

学号：10164102141

作业2

1 软件体系结构

2 详细设计

1 软件体系结构

Map

要求

班级学生成绩的随机生成

输入：本班同学的学号 输出：<学号，成绩>

数据

`stuID.csv`文件，每一列为一个学号，在 hdfs 上，shell input:

```
1 | hdfs dfs put stuID.csv input
```

reduceByKey

要求

求平均成绩：将全班同学每隔5号分为一组，求每组的平均成绩

输入：<学号，成绩> 输出：<组号，平均分>

数据

一个 `score.csv`文件，每一列为学号和学生成绩，放入 `HDFS` 中，shell input:

```
1 | hdfs dfs put score.csv input
```

NaturalJoin

要求：

匹配出person中每个人所在的位置信息；

每条记录各个字段之间以空格为分隔符。

数据:

person.txt

address.txt

Kmeans

数据

k-means.dat

第一行标明K的值和数据个数N, 均为整形, 由","隔开 (如 3,10 表示K=3, N=10)。

之后N行中每行代表一个二维向量, 向量元素均为整形, 由","隔开 (如 1,2 表示向量(1, 2))。

输出:

K行, 每行是一个聚类图心的二维向量, 向量元素均为浮点型 (如 1.1,2.3)。

2 详细设计

Map:

code:

```
1  import org.apache.spark.SparkConf;
2  import org.apache.spark.api.java.*;
3  import org.apache.spark.api.java.function.Function;
4
5  import java.util.ArrayList;
6  import java.util.List;
7  import java.util.Random;
8
9  public class StuScore {
10
11      private static Random rand = new Random();
12
13      public static void main(String[] args) {
14          SparkConf conf = new SparkConf().setAppName("StuScore");
15          JavaSparkContext sc = new JavaSparkContext(conf);
16          String logFile = "hdfs:///user/hadoop/input/stuID.csv";
17          JavaRDD<String> stuID = sc.textFile(logFile);
18
19          JavaRDD<String> stuScore = stuID.map(
20              new Function<String, String>() {
21                  @Override
22                  public String call(String s) throws Exception {
23                      String tmp = s + " " + String.valueOf( rand.nextInt(100)
24                          +1);
25                      return tmp;
26                  }
27              }
28          );
```

```
28  
29     stuScore.saveAsTextFile("hdfs:///user/hadoop/output/spark/StuRandomScore");  
30     }  
31 }
```

运行

本地运行

shell input:

```
1 | spark-submit --class StuScore StuScore-1.0.jar
```

查看结果,shell input:

```
1 | hdfs dfs -cat output/spark/StuRandomScore/*
```

集群上运行

同样，需要将HDFS地址更改为集群的地址，例如：

```
hdfs:///user/hadoop/input/stuID.csv => hdfs://10.11.6.91:9000/user/hadoop23/input/stuID.csv
```

shell input:

```
1 | spark-submit --class StuScore StuScore-1.0.jar
```

查看运行结果,shell input:

```
1 | hdfs dfs -cat output/spark/StuRandomScore/*
```

reducebykey

code:

```
1  import org.apache.spark.SparkConf;  
2  import org.apache.spark.api.java.*;  
3  import org.apache.spark.api.java.function.Function;  
4  import org.apache.spark.api.java.function.Function2;  
5  import org.apache.spark.api.java.function.PairFunction;  
6  import scala.Tuple2;  
7  
8  
9  import java.util.Iterator;  
10 import java.util.Random;  
11  
12 public class AVGScore {
```

```

13     private static Integer groupSize = 5;
14
15     public static void main(String[] args) {
16         SparkConf conf = new SparkConf().setAppName("AVGScore");
17         JavaSparkContext sc = new JavaSparkContext(conf);
18         sc.setLogLevel("WARN"); //http://stackoverflow.com/questions/27781187/how-to-
stop-messages-displaying-on-spark-console
19
20         String logFile = "hdfs:///user/hadoop/input/score.csv";
21
22         JavaRDD<String> fileRDD = sc.textFile(logFile);
23
24
25         /**
26          * map string to (id, score) and convert to (group_id, (score,1))
27          * reduceByKey => (group_id,(sumScore, count)
28          * and then mapValues to avg score
29          */
30
31         JavaPairRDD<Integer, Double> stuScore = fileRDD.mapToPair(
32             line -> new Tuple2<>(
33                 (Integer.parseInt(line.split(",")[0]) + 1 )/5, new Tuple2<>
(Double.parseDouble(line.split(",")[1]),1))
34             ).reduceByKey(
35                 (x,y) -> new Tuple2<>(x._1 + y._1,x._2+y._2))
36             ).mapValues(x -> x._1/x._2);
37
38         stuScore.saveAsTextFile("hdfs:///user/hadoop/output/spark/AVGscore");
39     }
40 }

```

运行

本地运行

运行程序, shell input:

```
1 | spark-submit --class AVGScore AVGScore-1.0.jar
```

查看结果, 同时查看行号, shell input:

```
1 | hdfs dfs -cat output/spark/AVGscore/* | wc -l
```

集群上运行

需要将HDFS地址更改为集群的地址, 例如:

```
hdfs:///user/hadoop/input/score.csv => hdfs://10.11.6.91:9000/user/hadoop23/input/score.csv
```

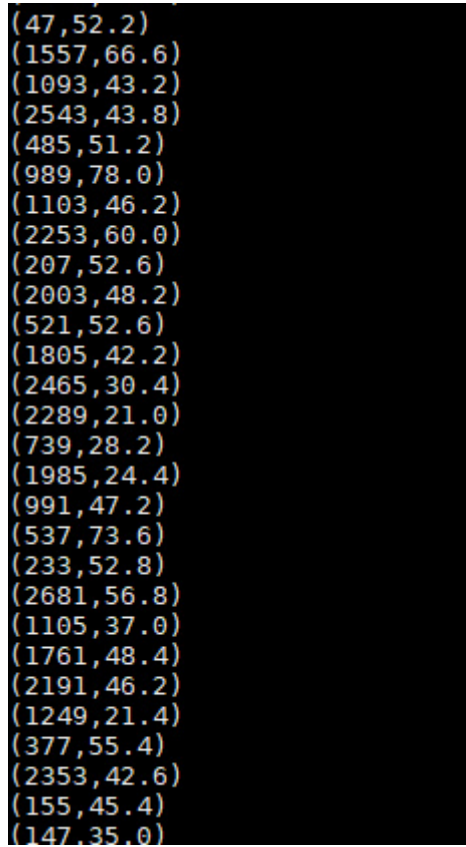
提交作业并运行, shell input:

```
1 | spark-submit --master spark://10.11.6.91:7077 --class AVGScore AVGScore-1.0.jar
```

查看结果, shell input:

```
1 | hdfs dfs -cat output/spark/AVGscore/*
```

output:



```
(47,52.2)
(1557,66.6)
(1093,43.2)
(2543,43.8)
(485,51.2)
(989,78.0)
(1103,46.2)
(2253,60.0)
(207,52.6)
(2003,48.2)
(521,52.6)
(1805,42.2)
(2465,30.4)
(2289,21.0)
(739,28.2)
(1985,24.4)
(991,47.2)
(537,73.6)
(233,52.8)
(2681,56.8)
(1105,37.0)
(1761,48.4)
(2191,46.2)
(1249,21.4)
(377,55.4)
(2353,42.6)
(155,45.4)
(147,35.0)
```

Natural join

code:

```
1 | import org.apache.commons.lang.StringUtils;
2 | import org.apache.hadoop.hdfs.protocol.DirectoryListing;
3 | import org.apache.spark.SparkConf;
4 | import org.apache.spark.api.java.JavaPairRDD;
5 | import org.apache.spark.api.java.JavaRDD;
6 | import org.apache.spark.api.java.JavaSparkContext;
7 | import org.apache.spark.api.java.function.Function;
8 | import org.apache.spark.api.java.function.PairFunction;
9 | import scala.Tuple2;
10 |
11 | public class NaturalJoin {
12 |
13 |
```

```

14     public static void main(String args[]) {
15         SparkConf conf = new SparkConf().setAppName("NaturalJoin");
16         JavaSparkContext sc = new JavaSparkContext(conf);
17         sc.setLogLevel("WARN"); //http://stackoverflow.com/questions/27781187/how-to-
stop-messages-displaying-on-spark-console
18
19
20         String addFile = "hdfs:///user/hadoop/input/address.txt";
21         String personFile = "hdfs:///user/hadoop/input/person.txt";
22
23
24         /**
25          * return code,city
26          */
27         JavaPairRDD<Integer, String> addRDD = sc.textFile(addFile).mapToPair(
28             line -> new Tuple2<>(
29                 Integer.parseInt(line.split(" ")[0]), line.split(" ")[1]));
30
31         /**
32          * return return code,{ id + name }
33          */
34         JavaPairRDD<Integer, String> personRDD = sc.textFile(personFile).mapToPair(
35             new PairFunction<String, Integer, String>() {
36                 @Override
37                 public Tuple2<Integer, String> call(String s) throws Exception {
38                     String[] splitLines = StringUtils.split(s, " ");
39                     if (splitLines.length < 3)
40                         return null;
41                     return new Tuple2<>(Integer.parseInt(splitLines[2]),
splitLines[0] + splitLines[1]);
42                 }
43             }
44         );
45
46         /**
47          * return code, [{id + name}, city]
48          */
49         JavaPairRDD<Integer, Tuple2<String,String>> resultRDD =
personRDD.join(addRDD);
50
51         resultRDD.saveAsTextFile("hdfs:///user/hadoop/output/spark/NaturalJoin");
52     }
53 }

```

运行

本地运行

运行程序并查看结果,shell input:

```
1 | spark-submit --class NaturalJoin NaturalJoin-1.0.jar
2 | hdfs dfs -cat output/spark/NaturalJoin/*
```

集群上运行

需要将HDFS地址更改为集群的地址，例如：

```
hdfs:///user/hadoop/input/address.txt => hdfs://10.11.6.91:9000/user/hadoop23/input/address.txt
```

提交作业并运行,shell input:

```
1 | spark-submit --master spark://10.11.6.91:7077 --class NaturalJoin NaturalJoin-1.0.jar
```

查看结果,ishell input:

```
1 | hdfs dfs -cat output/spark/NaturalJoin/*
```

```
226000,(74Conrad,Nantong))
226000,(5Abraham,Nantong))
226000,(35August,Nantong))
225300,(43Baldwin,Taizhou))
225300,(54Berger,Taizhou))
225300,(6Adair,Taizhou))
225300,(13Alan,Taizhou))
225300,(29Andy,Taizhou))
224000,(75Cornelius,Yancheng))
224000,(8Addison,Yancheng))
224000,(30Angelo,Yancheng))
224000,(36Aubrey,Yancheng))
224000,(39Asa,Yancheng))
223001,(49Beck,Huaian))
223001,(53Berg,Huaian))
223001,(58Bowen,Huaian))
223001,(9Adolph,Huaian))
hadoop7@ubuntu16g-1:~/Spark/Natural join/targets$
```

Kmeans

code:

Kmeans.java

```
1 | import org.apache.spark.api.java.JavaPairRDD;
2 | import org.apache.spark.api.java.JavaRDD;
3 | import org.apache.spark.api.java.JavaSparkContext;
4 | import org.apache.spark.api.java.function.Function2;
5 | import scala.Serializable;
6 | import scala.Tuple2;
7 |
8 | import java.util.ArrayList;
9 | import java.util.List;
10 |
11 | public class Kmeans implements KmeansInterface, Serializable {
12 |
```

```

13 // every point has a cluster number and point(x,y)
14 private List<Tuple2<Integer, Point>> oldCenterList = new ArrayList<>();
15 private List<Tuple2<Integer, Point>> newCenterList = new ArrayList<>();
16 private double threshold = 0.000001;
17
18
19 /**
20  * @param point
21  * @return cluster belonged
22  * @Method get the closest cluster for the point
23  */
24 public int findClosest(Point point) {
25     int argmin = -1;
26     double minimalDistance = Double.MAX_VALUE;
27     for (Tuple2<Integer, Point> i : oldCenterList) {
28         double distance = point.EuclideanDis(i._2);
29         if (distance < minimalDistance) {
30             minimalDistance = distance;
31             argmin = i._1;
32         }
33     }
34     return argmin;
35 }
36
37
38 /**
39  * @param outFile string
40  * @Method save center to txt
41  */
42 public void saveToFile(String outFile, JavaSparkContext sc) {
43     List<String> outCenterList = new ArrayList<>();
44     // format center points
45     for (Tuple2<Integer, Point> tmp : newCenterList) {
46         outCenterList.add(String.valueOf(tmp._2.getX()) + " " +
String.valueOf(tmp._2.getY()));
47     }
48
49     JavaRDD<String> center = sc.parallelize(outCenterList);
50     center.saveAsTextFile(outFile);
51 }
52
53
54 /**
55  * @return False for not stable
56  * @Method compare two cluster center with threshold
57  */
58 public boolean clusterCompare() {
59
60     for (Tuple2<Integer, Point> oldCenter : oldCenterList) {
61         int clusterNum = oldCenter._1;
62         for (Tuple2<Integer, Point> newCenter : newCenterList) {
63             if (newCenter._1 == clusterNum) {
64                 double dis = oldCenter._2.EuclideanDis(newCenter._2);

```



```

65         if (dis > threshold)
66             return false;
67         break;
68     }
69 }
70 }
71 return true;
72 }
73
74
75 /**
76  * @param kmeansRDD
77  * @return init pointsRDD
78  * @Method prepare Points RDD and select clusters randomly
79  */
80 public JavaPairRDD<Integer, Point> Prepare(JavaRDD<String> kmeansRDD) {
81     // get the number of cluster
82     String fisrtLine = kmeansRDD.first();
83     int clusterCount = Integer.parseInt(fisrtLine.split(",")[0]);
84
85     // filter first line and convert to <Point,clusternum>, init set all cluster
86     // number 1
87     JavaPairRDD<Integer, Point> pointsRDD = kmeansRDD.filter(line ->
88 !line.equals(fisrtLine)).mapToPair(
89     line -> {
90         String[] splitLine = line.split(",");
91         double X = Double.parseDouble(splitLine[0]);
92         double Y = Double.parseDouble(splitLine[1]);
93         return new Tuple2<>(0, new Point(X, Y));
94     }
95 );
96
97 // init center list
98 oldCenterList.addAll(pointsRDD.take(clusterCount));
99
100 for (int i = 0; i < clusterCount; i++) {
101     Tuple2<Integer, Point> tmp = oldCenterList.get(i);
102     oldCenterList.set(i, new Tuple2<>(i, tmp._2));
103 }
104
105 newCenterList.addAll(oldCenterList);
106
107 return pointsRDD;
108 }
109
110 /**
111  * @param pointsRDD to cluster
112  * @return new classify PointsRDD
113  * @method cluster and update new cluster center
114  */

```

```

115     public JavaPairRDD<Integer, Point> cluster(JavaPairRDD<Integer, Point>
pointsRDD) {
116
117         JavaPairRDD<Integer, Point> newPointsRDD = pointsRDD.mapToPair(
118             kv -> new Tuple2<>(findClosest(kv._2), kv._2)
119         );
120
121         JavaPairRDD<Integer, Point> newClusterRDD = newPointsRDD
122             .mapValues(
123                 value -> new Tuple2<>(value, 1))
124             .reduceByKey(
125                 new Function2<Tuple2<Point, Integer>, Tuple2<Point,
Integer>, Tuple2<Point, Integer>>() {
126                     @Override
127                     public Tuple2<Point, Integer> call(Tuple2<Point,
Integer> value1, Tuple2<Point, Integer> value2) throws Exception {
128                         Point tmp = new Point(value1._1.getX() +
value2._1.getX(), value1._1.getY() + value2._1.getY());
129                         int count = value1._2 + value2._2;
130                         return new Tuple2<>(tmp, count);
131                     }
132                 }
133             ).mapValues(
134                 v -> new Point(v._1.getX() / v._2, v._1.getY() / v._2)
135             );
136
137         oldCenterList.clear();
138         oldCenterList.addAll(newCenterList);
139         // convert to list to store
140         newCenterList.clear();
141         newCenterList.addAll(newClusterRDD.collect());
142
143         return newPointsRDD;
144     }
145
146 }

```

kmeansRun.java

```

1  import org.apache.spark.SparkConf;
2  import org.apache.spark.api.java.JavaPairRDD;
3  import org.apache.spark.api.java.JavaRDD;
4  import org.apache.spark.api.java.JavaSparkContext;
5  import scala.Tuple2;
6
7  import java.util.ArrayList;
8  import java.util.List;
9
10 public class kmeansRun {
11
12     public static void main(String args[]) {
13         Kmeans kmeans = new Kmeans();
14

```

```

15 SparkConf conf = new SparkConf().setAppName("Kmeans");
16 JavaSparkContext sc = new JavaSparkContext(conf);
17
18 String kmeansFile = "hdfs://10.11.6.91:9000/user/hadoop23/input/k-means.dat";
19 String outFile = "hdfs://10.11.6.91:9000/user/hadoop23/output/spark/kmeans";
20
21 int counter = 1, maxIteration = 500;
22 boolean result = false;
23
24 JavaRDD<String> kmeansRDD = sc.textFile(kmeansFile).cache();
25
26 //      init centerList and Points RDD
27 JavaPairRDD<Integer, Point> PointsRDD = kmeans.Prepare(kmeansRDD);
28
29 while (!result && counter < maxIteration) {
30     PointsRDD = kmeans.cluster(PointsRDD);
31
32     result = kmeans.clusterCompare();
33     System.out.println("*****KMEANS finished iteration:>> " + counter + "
|| means stable: " + result);
34     counter++;
35 }
36
37 kmeans.saveToFile(outFile,sc);
38
39 }
40 }

```

运行

本地运行

运行程序, shell input:

```
1 | spark-submit --class kmeansRun Kmeans-1.0.jar 2>&1 | grep "KMEANS finished iteration:"
```

output:

```

*****KMEANS finished iteration:>> 1 || means stable: false
*****KMEANS finished iteration:>> 2 || means stable: false
*****KMEANS finished iteration:>> 3 || means stable: false
*****KMEANS finished iteration:>> 4 || means stable: false
*****KMEANS finished iteration:>> 5 || means stable: false
*****KMEANS finished iteration:>> 6 || means stable: false
*****KMEANS finished iteration:>> 7 || means stable: false
*****KMEANS finished iteration:>> 8 || means stable: true

```

查看结果, shell input:

```
1 | hdfs dfs -cat output/spark/kmeans/*
```

output:

```
98.31,803.89
496.11650485436894,207.3398058252427
493.2277227722772,798.3267326732673
91.05208333333333,206.19791666666666
```

集群上运行

需要将HDFS地址更改为集群的地址，例如：

```
hdfs:///user/hadoop/input/k-means.dat => hdfs://10.11.6.91:9000/user/hadoop23/input/k-means.dat
```

提交作业并运行, shell input:

```
1 | spark-submit --master spark://10.11.6.91:7077 --class kmeansRun Kmeans-1.0.jar 2>&1 |
  grep "KMEANS finished iteration:"
```

output:

```
*****KMEANS finished iteration:>> 1 || means stable: false
*****KMEANS finished iteration:>> 2 || means stable: false
*****KMEANS finished iteration:>> 3 || means stable: false
*****KMEANS finished iteration:>> 4 || means stable: false
*****KMEANS finished iteration:>> 5 || means stable: false
*****KMEANS finished iteration:>> 6 || means stable: false
*****KMEANS finished iteration:>> 7 || means stable: false
*****KMEANS finished iteration:>> 8 || means stable: true
```

查看结果：

```
18/11/01 22:22:50 WARN util.NativeCodeLoader: Unable to load native
. using builtin-java classes where applicable
98.31 803.89
493.2277227722772 798.3267326732673
496.11650485436894 207.3398058252427
91.05208333333333 206.19791666666666
```