

# 作业4

姓名： 吴双

学号： 10164102141

## 《高级编程》课程项目实践报告

1 软件体系结构

2 详细设计

2.4 编码实现及问题分析

Map部分：

ReduceByKey部分：

NaturalJoin部分：

KMeans部分：

结果：

## 1 软件体系结构

person.txt

address.txt

README.md

map.java

ReduceByKey.java

NaturalJoin.java

KMeans

k-means.dat

README\_k-means.md

point.java

kmeansRun.java

## 2 详细设计

Map

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>

## 数据准备

1. 首先需要有一个 stuID.csv 文件，每一列为一个学号：

```
1 2182863732
2 2439766117
3 3266939511
4 2366826895
5 3045445297
6 2936280284
7 2650797670
8 2308375640
9 2272287061
```

2. 然后将文件放入 HDFS 中：

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

## ReduceByKey

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>

## 数据准备

1. 首先需要有一个 score.csv 文件，每一列为学号和学生成绩：

```
1 1,55
2 2,12
3 3,48
4 4,98
5 5,31
6 6,81
7 7,19
8 8,54
9 9,65
```

## Natural join

person.txt

```
Aaron 210000
Abbott 214000
Abel 221000
Abner 215000
Abraham 226000
Adair 225300
Adam 223800
Addison 224000
Adolph 223001
```

address.txt

```
210000 Nanjing
214000 Wuxi
221000 Xuzhou
213000 Changzhou
```

要求以code为连接属性，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。

## KMeans

输入：

- 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)。

## 2.4 编码实现及问题分析

Map部分：

```
1 import org.apache.flink.api.common.functions.MapFunction;
2 import org.apache.flink.api.java.DataSet;
3 import org.apache.flink.api.java.ExecutionEnvironment;
4 import org.apache.flink.api.java.tuple.Tuple2;
5 import org.apache.flink.api.java.utils.ParameterTool;
6 import java.util.Random;
7
```

```

8 public class StuScore {
9     private static Random rand = new Random();
10
11     public StuScore(){}
12
13     public static void main(String[] args) throws Exception {
14         ParameterTool params = ParameterTool.fromArgs(args);
15         ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
16         env.getConfig().setGlobalJobParameters(params);
17         DataSet<String> text;
18
19         if(params.has("input")){
20             text = env.readTextFile(params.get("input"));
21         }else{
22             System.out.println("Please confirm input keywords!");
23             return;
24         }
25
26
27         DataSet<Tuple2<String,Integer>> stuscore = text.map(new MapFunction<String,
Tuple2<String, Integer>>() {
28             @Override
29             public Tuple2<String, Integer> map(String s) throws Exception {
30                 return new Tuple2<>(s,rand.nextInt(100) +1);
31             }
32         });
33
34         //如果没有指定输出，则默认打印到控制台
35         if(params.has("output")){
36             stuscore.writeAsCsv(params.get("output"),"\n", ",");
37             env.execute();
38         }else{
39             System.out.println("Printing result to stdout. Use --output to specify
output path.");
40             stuscore.print();
41         }
42
43     }
44 }

```

## ReduceByKey部分:

```

1 import org.apache.commons.compress.archivers.dump.DumpArchiveEntry;
2 import org.apache.flink.api.common.typeinfo.Types;
3 import org.apache.flink.api.java.DataSet;
4 import org.apache.flink.api.java.ExecutionEnvironment;
5 import org.apache.flink.api.java.tuple.Tuple2;
6 import org.apache.flink.api.java.utils.ParameterTool;
7
8 public class AVGscore {
9     private static Integer groupSize = 5;
10
11     public static void main(String[] args) throws Exception {

```

```

12 ParameterTool params = ParameterTool.fromArgs(args);
13 ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
14 env.getConfig().setGlobalJobParameters(params);
15 DataSet<Tuple2<Integer, Double>> fileDataSet;
16
17 if (params.has("input")) {
18     fileDataSet = env.readCsvFile(params.get("input"))
19         .types(Integer.class, Double.class);
20 } else {
21     System.out.println("Please confirm input keywords!");
22     return;
23 }
24
25 /**
26  * map string to (id, score) and convert to (group_id, (score,1))
27  * GroupBy and reduce == reduceByKey
28  * and then map to avg score
29  */
30 DataSet<Tuple2<Integer, Double>> stuAVGscore = fileDataSet
31     .map(line -> Tuple2.of(
32         (line.f0-1)/ 5, Tuple2.of(line.f1, 1)))
33     .returns(Types.TUPLE(Types.INT, Types.TUPLE(Types.DOUBLE,
Types.INT)))
34     .groupBy(0)
35     .reduce(
36         (kv1, kv2) -> Tuple2.of(kv1.f0, Tuple2.of(kv1.f1.f0 +
kv2.f1.f0, kv1.f1.f1 + kv2.f1.f1)))
37     .returns(Types.TUPLE(Types.INT, Types.TUPLE(Types.DOUBLE,
Types.INT)))
38     .map(
39         line -> Tuple2.of(line.f0, line.f1.f0 / line.f1.f1)
40     ).returns(Types.TUPLE(Types.INT, Types.DOUBLE));
41
42 //如果没有指定输出，则默认打印到控制台
43 if (params.has("output")) {
44     stuAVGscore.writeAsCsv(params.get("output"), "\n", ",");
45     env.execute();
46 } else {
47     System.out.println("Printing result to stdout. Use --output to specify
output path.");
48     stuAVGscore.print();
49 }
50
51 }
52 }

```

## NaturalJoin部分:

```

1 import org.apache.flink.api.common.functions.FlatJoinFunction;
2 import org.apache.flink.api.common.typeinfo.Types;
3 import org.apache.flink.api.java.DataSet;
4 import org.apache.flink.api.java.ExecutionEnvironment;
5 import org.apache.flink.api.java.tuple.Tuple2;

```

```

6  import org.apache.flink.api.java.tuple.Tuple3;
7  import org.apache.flink.api.java.tuple.Tuple4;
8  import org.apache.flink.api.java.utils.ParameterTool;
9  import org.apache.flink.util.Collector;
10 import scala.Int;
11
12 import java.lang.reflect.Type;
13
14 public class NaturalJoin {
15
16     public static void main(String args[]) throws Exception{
17         ParameterTool params = ParameterTool.fromArgs(args);
18         ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
19         env.getConfig().setGlobalJobParameters(params);
20
21         // code, city
22         DataSet<Tuple2<Integer, String>> addDataSet;
23         // id ,name, code
24         DataSet<Tuple3<Integer, String, Integer>> personDataSet;
25
26         if (params.has("addinput")) {
27             addDataSet = env.readCsvFile(params.get("addinput"))
28                 .fieldDelimiter(" ")
29                 .ignoreInvalidLines()
30                 .types(Integer.class, String.class);
31         } else {
32             System.out.println("Please confirm input keywords!");
33             return;
34         }
35
36         if (params.has("personinput")) {
37             personDataSet = env.readCsvFile(params.get("personinput"))
38                 .fieldDelimiter(" ")
39                 .ignoreInvalidLines()
40                 .types(Integer.class, String.class, Integer.class);
41         } else {
42             System.out.println("Please confirm input keywords!");
43             return;
44         }
45
46         DataSet<Tuple4<Integer, String, Integer, String>> result =
47         personDataSet.join(addDataSet)
48             .where(2)
49             .equalTo(0)
50             .with(
51                 (x, y) -> Tuple4.of(x.f0, x.f1, x.f2, y.f1)
52             ).returns(Types.TUPLE(Types.INT,Types.STRING,Types.INT,Types.STRING));
53
54         personDataSet.print();
55
56         //如果没有指定输出, 则默认打印到控制台
57         if (params.has("output")) {

```

```

57         result.writeAsCsv(params.get("output"), "\n", ",");
58         env.execute();
59     } else {
60         System.out.println("Printing result to stdout. Use --output to specify
output path.");
61         result.print();
62     }
63
64 }
65 }

```

## KMeans部分:

point.java:用于自定义point类 (POJO对象)

```

1  import java.io.Serializable;
2
3  public class Point implements Serializable {
4      public double x,y;
5
6      public Point() {
7
8      }
9
10     public Point(double x, double y) {
11         this.x = x;
12         this.y = y;
13     }
14
15     public Point add(Point other) {
16         x += other.x;
17         y += other.y;
18         return this;
19     }
20
21     public Point div(long val) {
22         x /= val;
23         y /= val;
24         return this;
25     }
26
27     public double euclideanDistance(Point other) {
28         return Math.sqrt((x - other.x) * (x - other.x) + (y - other.y) * (y -
other.y));
29     }
30
31     public String toString() {
32         return x + " " + y;
33     }
34
35 }

```

KMeansRun.java:

```

1  import org.apache.flink.api.common.functions.*;
2  import org.apache.flink.api.java.DataSet;
3  import org.apache.flink.api.java.ExecutionEnvironment;
4  import org.apache.flink.api.java.functions.FunctionAnnotation;
5  import org.apache.flink.api.java.operators.IterativeDataSet;
6  import org.apache.flink.api.java.tuple.Tuple2;
7  import org.apache.flink.api.java.tuple.Tuple3;
8  import org.apache.flink.api.java.utils.ParameterTool;
9  import org.apache.flink.configuration.Configuration;
10
11 import java.io.BufferedReader;
12 import java.io.FileReader;
13 import java.util.ArrayList;
14 import java.util.Collection;
15
16 public class Kmeans {
17
18     public static void main(String[] args) throws Exception {
19         // Checking input parameters
20         final ParameterTool params = ParameterTool.fromArgs(args);
21         // set up execution environment
22         ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
23         env.getConfig().setGlobalJobParameters(params); // make parameters available
24         in the web interface
25
26         // get input data:
27         DataSet<Point> points = getPoint(params, env);
28         if (points == null)
29             return;
30         DataSet<Centroid> centroids = getCentroid(params, env);
31
32         // set number of bulk iterations for KMeans algorithm
33         IterativeDataSet<Centroid> loop =
34             centroids.iterate(params.getInt("iterations", 100));
35
36         DataSet<Centroid> newCentroids = points
37             // compute closest centroid for each point
38             .map(new SelectNearestCenter()).withBroadcastSet(loop, "centroids")
39             // count and sum point coordinates for each centroid
40             .map(new CountAppender())
41             .groupBy(0).reduce(new CentroidAccumulator())
42             // compute new centroids from point counts and coordinate sums
43             .map(new CentroidAverager());
44
45         // feed new centroids back into next iteration
46         DataSet<Centroid> finalCentroids = loop.closeWith(newCentroids,
47             newCentroids.filter(new thresholdFilter()).withBroadcastSet(loop, "centroids"));
48
49         // DataSet<Tuple2<Integer, Point>> clusteredPoints = points
50         // // assign points to final clusters

```



```

50 // .map(new SelectNearestCenter()).withBroadcastSet(finalCentroids,
    "centroids");
51
52 // emit result
53 if (params.has("output")) {
54     finalCentroids.writeAsCsv(params.get("output"), "\n", " ");
55     env.execute();
56 } else {
57     System.out.println("Printing result to stdout. Use --output to specify
output path.");
58     finalCentroids.print();
59 }
60 }
61
62 private static DataSet<Point> getPoint(ParameterTool params,
ExecutionEnvironment env) {
63     DataSet<Point> points;
64     if (params.has("input")) {
65         // read points from CSV file
66         points = env.readCsvFile(params.get("input"))
            .ignoreFirstLine()
            .pojoType(Point.class, "x", "y");
67     } else {
68         System.out.println("Use --input to specify file input.");
69         return null;
70     }
71     return points;
72 }
73
74 private static DataSet<Centroid> getCentroid(ParameterTool params,
ExecutionEnvironment env) throws Exception {
75
76     ArrayList<Centroid> centroidArrayList = new ArrayList<>();
77     BufferedReader br = new BufferedReader(new FileReader(params.get("input")));
78     String text = br.readLine();
79     int k = Integer.parseInt(text.split(",")[0]);
80
81     while (k != 0) {
82         text = br.readLine();
83         double x = Double.parseDouble(text.split(",")[0]);
84         double y = Double.parseDouble(text.split(",")[1]);
85         centroidArrayList.add(new Centroid(k, x, y));
86         k--;
87     }
88
89     DataSet<Centroid> centroids = env.fromCollection(centroidArrayList);
90     return centroids;
91 }
92
93 /**
94  * Determines the closest cluster center for a data point.
95  */
96 @FunctionAnnotation.ForwardedFields("*->1")
97

```

```

98     public static final class SelectNearestCenter extends RichMapFunction<Point,
Tuple2<Integer, Point>> {
99         private Collection<Centroid> centroids;
100
101         /**
102          * Reads the centroid values from a broadcast variable into a collection.
103          */
104         @Override
105         public void open(Configuration parameters) throws Exception {
106             this.centroids = getRuntimeContext().getBroadcastVariable("centroids");
107         }
108
109         @Override
110         public Tuple2<Integer, Point> map(Point p) throws Exception {
111
112             double minDistance = Double.MAX_VALUE;
113             int closestCentroidId = -1;
114
115             // check all cluster centers
116             for (Centroid centroid : centroids) {
117                 // compute distance
118                 double distance = p.euclideanDistance(centroid);
119
120                 // update nearest cluster if necessary
121                 if (distance < minDistance) {
122                     minDistance = distance;
123                     closestCentroidId = centroid.id;
124                 }
125             }
126
127             // emit a new record with the center id and the data point.
128             return new Tuple2<>(closestCentroidId, p);
129         }
130     }
131
132     /**
133      * Appends a count variable to the tuple.
134      */
135     @FunctionAnnotation.ForwardedFields("f0;f1")
136     public static final class CountAppender implements MapFunction<Tuple2<Integer,
Point>, Tuple3<Integer, Point, Long>> {
137         @Override
138         public Tuple3<Integer, Point, Long> map(Tuple2<Integer, Point> t) {
139             return new Tuple3<>(t.f0, t.f1, 1L);
140         }
141     }
142
143     /**
144      * Sums and counts point coordinates.
145      */
146     @FunctionAnnotation.ForwardedFields("0")
147     public static final class CentroidAccumulator implements
ReduceFunction<Tuple3<Integer, Point, Long>> {

```

```

148
149     @Override
150     public Tuple3<Integer, Point, Long> reduce(Tuple3<Integer, Point, Long>
val1, Tuple3<Integer, Point, Long> val2) {
151         return new Tuple3<>(val1.f0, val1.f1.add(val2.f1), val1.f2 + val2.f2);
152     }
153 }
154
155 /**
156  * Computes new centroid from coordinate sum and count of points.
157  */
158 @FunctionAnnotation.ForwardedFields("0->id")
159 public static final class CentroidAverager implements
MapFunction<Tuple3<Integer, Point, Long>, Centroid> {
160     @Override
161     public Centroid map(Tuple3<Integer, Point, Long> value) {
162         // id, x/num y/num
163         return new Centroid(value.f0, value.f1.div(value.f2));
164     }
165 }
166
167 /**
168  * Filter that filters vertices where the centroid difference is below a
threshold.
169  */
170 public static final class thresholdFilter extends RichFilterFunction<Centroid> {
171     private Collection<Centroid> centroids;
172     private double threshold = 1e-5;
173
174     /**
175      * Reads the centroid values from a broadcast variable into a collection.
176      */
177     @Override
178     public void open(Configuration parameters) throws Exception {
179         this.centroids = getRuntimeContext().getBroadcastVariable("centroids");
180     }
181
182     @Override
183     public boolean filter(Centroid centroid) {
184         for (Centroid oldcentroid : centroids) {
185             if (centroid.id == oldcentroid.id) {
186                 // compute distance
187                 double distance = centroid.euclideanDistance(oldcentroid);
188                 return (distance > this.threshold);
189             }
190         }
191         return true;
192     }
193 }
194 }

```

**结果:**

Map部分:

```
1 | flink run -c StuScore StuScore.jar --input  
  | /home/hadoop/Documents/distribution/Flink/StuScore/stuID.csv
```

output:

```
(3076722373,20)  
(407826062,3)  
(2115951107,75)  
(3116395451,48)  
(651866235,8)  
(245285626,29)  
Program execution finished  
Job with JobID 6fcaba071f7e23671ad6c7ba4e836ec5 has finished.  
Job Runtime: 475 ms  
Accumulator Results:  
- 8ab08a881ddb93127a7f75ca2aba7270 (java.util.ArrayList) [14094
```

ReduceByKey部分:

```
1 | flink run -c AVGscore AVGscore.jar --input  
  | /home/hadoop/Documents/distribution/Flink/AVGscore/score.csv
```

output:

```
(2813,38.2)  
(2814,51.6)  
(2815,71.0)  
(2816,45.2)  
(2817,39.4)  
(2818,55.0)  
Program execution finished
```

Natural join部分:

```
1 | flink run -c NaturalJoin NaturalJoin.jar --addinput  
  | /home/hadoop/Documents/distribution/Flink/NaturalJoin/address.txt --personinput  
  | /home/hadoop/Documents/distribution/Flink/NaturalJoin/person.txt
```

output:

```
(70,Cliff,222000,Lianyungang)
(71,Cleveland,214000,Wuxi)
(72,Clyde,226000,Nantong)
(73,Colin,222000,Lianyungang)
(74,Conrad,226000,Nantong)
(75,Cornelius,224000,Yancheng)
Program execution finished
Job with JobID 6ab7c20915728a60c808fca6
Job Runtime: 2176 ms
```

Kmeans部分:

```
1 | flink run -c Kmeans Kmeans.jar --input
   | /home/hadoop/Documents/distribution/Flink/kmeans/k-means.dat
```

output:

```
Printing result to stdout. Use --output to specify output path.
1 91.05208333333333 206.19791666666666
2 493.2277227722772 798.3267326732673
3 496.11650485436894 207.3398058252427
4 98.31 803.89
Program execution finished
Job with JobID 1dbe3cd013615653e999d27047f8971e has finished.
Job Runtime: 1405 ms
Accumulator Results:
```