

JAVA 课程大作业

学院	软件与微电子学院	
组员	潘江 2301210367	
组员	于畅 2301210509	
组员	常逸歌 2301210174	

2024年1月4日

目录

1	作业	地址	3		
2	任务一: 基本算法实现				
	2.1	问题描述	3		
		算法设计思路			
	2.3	测试方案	9		
	2.4	部分结果如下所示	9		
3	任务	·二:基于多线程的预测	10		
		改进内容			
	3.2	测试方案	12		
	3.3	部分结果如下所示	12		

1 作业地址

- Gitee: https://gitee.com/jiang000/java_homework
- Overleaf: https://www.overleaf.com/read/qzpzjpnpkytg#316605
- 结果: output.txt

2 任务一: 基本算法实现

2.1 问题描述

对于给定一个数据图和一个查询需求(多个关键词),通过图搜索算法在给定数据图中给出一组合适的 Web APIs, 这组 APIs 覆盖所有的关键词,而且能形成连通路径(API 具有兼容性)。

2.2 算法设计思路

- 数据预处理:根据题目要求,该任务只需要 api.csv 和 mashup.csv 中的部分数据。
 - api.csv 文件中仅保留 Name、PrimaryCategory、SecondaryCategories 三列,得到新的 api.csv 文件。
 - mashup.csv 文件中仅保留的 Name、RelatedAPIs 和 Categories 三列,得
 到新的 mashup.csv 文件。
- **构建图**: 读取 *api.csv* 和 *mashup.csv* 中的数据,从而构建出图。使用了 *OpenCsv* 和 *JgraphT* 工具包来完成 *csv* 文件的读取,以及整个图的构建。其中 *OpenCsv* 对于读写 *csv* 文件提供了便利的操作, *JGraphT* (*Java* 图工具包)是一个用于处理图和图算法的 *Java* 库。
 - 读入图的顶点和顶点的类别属性,该部分核心代码如下:

```
1 //定义一个无向图,读取 api.csv 文件,构建图数据结构
2 Graph<String,DefaultWeightedEdge> undirectedGraph = new

→ SimpleWeightedGraph<>(DefaultWeightedEdge.class);
3 String apiCsvFile = "api.csv";
4 try {
5 CSVParser csvParser = new CSVParser(new

→ FileReader(apiCsvFile), CSVFormat.DEFAULT.withHeader());
6 for (CSVRecord record : csvParser) {
```

String name = record.get("Name");

```
String category = record.get("Primary Category");
          undirectedGraph.addVertex(name);
          PrimaryCategory.put(name, category);
      }
   } catch (IOException e) {
      e.printStackTrace();
- 生成边的权值。读入 mashup.csv 文件, 读入其中的 Related APIs 列, 每一个
   单元里面出现的 api, 说明他们曾经被一起调用过,则两两之间存在一条边。
   边的权值对应着一起出现的次数。该部分核心代码如下:
  // 读取 mashup.csv 文件, 计算边的权值
  String mashupCsvFile = "mashup.csv";
  try {
      CSVParser csvParser = new CSVParser(new
          FileReader(mashupCsvFile),
          CSVFormat.DEFAULT.withHeader());
      for (CSVRecord record : csvParser) {
          String[] relatedApis = record.get("Related
           → APIs").split(",");
          for(String a :relatedApis){
              for(String b :relatedApis){
                  if (undirectedGraph.containsVertex(a) &&
                     undirectedGraph.containsVertex(b)){
                      if(a != b){
                          if(!undirectedGraph.containsEdge(a,b)){
11
12
                                 addWeightedEdge(undirectedGraph,a,b,1);
                         }
13
                         else{
14
                             DefaultWeightedEdge targetEdge =
15
                                 findEdge(undirectedGraph,a,b);
16
                                 setEdgeWeight(undirectedGraph,targetEdge,
                                 undirectedGraph.getEdgeWeight(targetEdge)
                                 + 1);
                         }
17
                      }
```

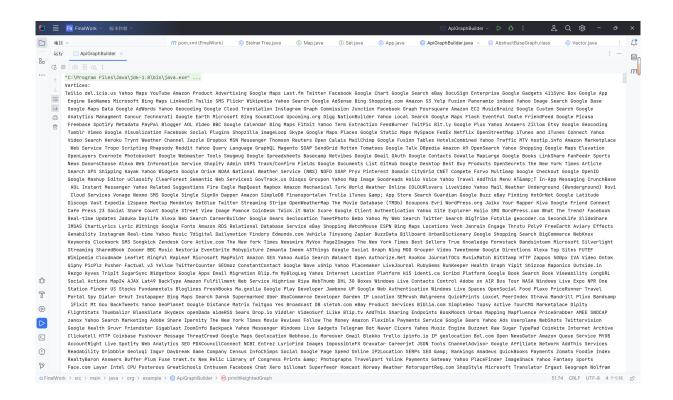


图 1: 构建好的图

- **以上构建了初始的图**,输出图的结点和权值,部分截图如上图所示:
- **生成 api 推荐算法**: 使用的 api 推荐算法主要是论文中的方法,下面是一些核心代码。
- **定义树的结构**: 首先定义一个树的数据结构 MidTree, 用于生成树, 树的一些属性如下:
- 1 private Graph<String, DefaultWeightedEdge> graph;//存储从属的图
- 2 private String root; //存储该树的跟借点
- private Vector<String> vertexSetofMidTree = new Vector<>();//存储该

 → 树存在的结点
 - private Vector<DefaultWeightedEdge> edgeSetofMidTree = new
 - → Vector<>();//存储该树的根

```
private Set<String> KeywordsofTree = new HashSet<>();//存储该树结点
   → 包含的类别
6 private double totleweight = 0;//存储该树的权值
• 两个队列和对队列的操作
  Vector<MidTree> Q = new Vector<>();
  Vector<MidTree> RQ = new Vector<>();
  public static void enqueue(Vector<MidTree> WQ,MidTree t){
      if(WQ.contains(t)) return;
      WQ.add(t);
  public static MidTree dequeue(Vector<MidTree> WQ){
      if(WQ.size() == 0){
          return null;
      }
10
      MidTree now = WQ.get(0);
      WQ.remove(0);
12
      return now;
13
  public void updatequeue(Vector<MidTree> WQ){
      sortVector(WQ,
16
       }
17
• Tree Merging 算法: 树的合并
  public static Vector<MidTree> mergeTreesWithSameRoot(Vector<MidTree>
   → Q) {
      Vector<MidTree> mergedTrees = new Vector<>();
      HashSet<Integer> processedIndices = new HashSet<>();
      for (int i = 0; i < Q.size(); i++) {</pre>
          if (processedIndices.contains(i)) continue;
          MidTree tree1 = Q.get(i);
          for (int j = i + 1; j < Q.size(); j++) {</pre>
              if (processedIndices.contains(j)) continue;
              MidTree tree2 = Q.get(j);
              if (tree1.getRoot().equals(tree2.getRoot())) {
10
                 // 合并 tree1 和 tree2
```

```
MidTree mergedTree = mergeTwoTrees(tree1, tree2);
12
                   mergedTrees.add(mergedTree);
                   processedIndices.add(i);
14
                   processedIndices.add(j);
15
                   break;
               }
17
           }
18
           if (!processedIndices.contains(i)) {
               mergedTrees.add(tree1);
20
           }
21
       }
       return mergedTrees;
23
24
25
   private static MidTree mergeTwoTrees(MidTree tree1, MidTree tree2) {
26
   // 创建一个新的 MidTree 实例来合并 tree1 和 tree2
27
       MidTree mergedTree = new MidTree(tree1.getGraph());
28
       mergedTree.setRoot(tree1.getRoot());
29
       // 合并顶点集
30
       for (String vertex : tree1.getVertexSetofMidTree()) {
31
           mergedTree.addVertexToMidTree(vertex);
       }
33
       for (String vertex : tree2.getVertexSetofMidTree()) {
34
           mergedTree.addVertexToMidTree(vertex);
       }
       // 合并边集
37
       for (DefaultWeightedEdge edge : tree1.getEdgeSetofMidTree()) {
38
           mergedTree.addEdgeToMidTree(edge);
       }
40
       for (DefaultWeightedEdge edge : tree2.getEdgeSetofMidTree()) {
41
           mergedTree.addEdgeToMidTree(edge);
       }
43
       // 合并关键字集
44
       for (String keyword : tree1.getKeywordsofTree()) {
           mergedTree.setKeywordsofTree(keyword);
46
       }
47
       for (String keyword : tree2.getKeywordsofTree()) {
48
```

```
mergedTree.setKeywordsofTree(keyword);
49
       }
       return mergedTree;
  }
  Tree growth 算法: 树的生长
   //tree growth 算法
   Set<String> N = new HashSet<>();//找到所有邻居
  N.addAll(Graphs.neighborListOf(graph,midTree.getRoot()));
   Set<String> keyofmidtree = midTree.getKeywordsofTree();
   for(String n : N){//遍历邻接点
       DefaultWeightedEdge e = graph.getEdge(n,midTree.getRoot());//对
       → 应边
       if(e == null) continue;
       String NprimaryCategory = PrimaryCategory.get(n);
       Set<MidTree> Willremove = new HashSet<>();
       Iterator<MidTree> iterator = Q.iterator();
10
       Vector<MidTree> increaseQ = new Vector<>();
11
       while(iterator.hasNext()){
           MidTree q = iterator.next();
13
           if(q!=midTree){
14
               if(q.getRoot().equals(n)){
15
                   //判断两者包含的 key 是是否完全一致
16
                   keyofmidtree.add(NprimaryCategory);
17
                   if(q.coverKeywords(keyofmidtree)){
18
                       if(q.getTotleweight()>
                           midTree.getTotleweight()+graph.getEdgeWeight(e)){
                           iterator.remove();
20
                           midTree.setKeywordsofTree(NprimaryCategory);
21
                           midTree.setRoot(n);
22
                           midTree.addEdgeToMidTree(e);
23
                           midTree.addVertexToMidTree(n);
                           increaseQ.add(midTree);
25
                           continue;
26
                       }
27
28
                   }
29
```

```
midTree.setKeywordsofTree(NprimaryCategory);
30
                     midTree.setRoot(n);
                     midTree.addEdgeToMidTree(e);
32
                     midTree.addVertexToMidTree(n);
33
                     increaseQ.add(midTree);
34
                }
35
            }
36
        }
        for(MidTree newq : increaseQ){
38
            Q.add(newq);
39
        }
40
        increaseQ.clear();
41
   }
42
```

2.3 测试方案

- 读取 mashup.csv 中的 relatedApis 作为验证结果,然后将查询这些 api 所属的种类作为查询的关键词,使用这些关键词调用算法,将 relatedApis 与算法返回的结点做比较,一致则说明命中。
- 一共命中了 1434 个任务, 命中率为 22.17573221757
- 若返回的 api 数少于所给的 api 数也算命中,则命中率为 80.04029133736246
- 每个任务的平均耗时为 1157.6011157601115760(ms)
- 总用时 2 小时 11 分钟

2.4 部分结果如下所示

- 第 329 个任务返回 api 共耗时 4296 毫秒
- 第 329 个任务返回的群 Steiner Tree 的权重和为: 0.5
- 第 329 个任务返回的群 Steiner Tree 的顶点个数为:2
- 第 329 个任务未命中
- 第 367 个任务返回 api 共耗时 1459 毫秒
- 第 367 个任务返回的群 Steiner Tree 的权重和为: 0.5
- 第 367 个任务返回的群 Steiner Tree 的顶点个数为:2

- 第367个任务未命中
- 第 392 个任务返回 api 共耗时 26 毫秒
- 第 392 个任务返回的群 Steiner Tree 的权重和为: 0.0
- 第 392 个任务返回的群 Steiner Tree 的顶点个数为:1
- 第 392 个任务命中这是第 88 个命中
- 第 394 个任务返回 api 共耗时 26 毫秒
- 第 394 个任务返回的群 Steiner Tree 的权重和为: 0.0
- 第 394 个任务返回的群 Steiner Tree 的顶点个数为:1
- 第 394 个任务命中这是第 90 个命中

3 任务二:基于多线程的预测

3.1 改进内容

使用 ExecutorService 创建多线程算法,通过 ExecutorService 能够更方便地管理线程,提高程序的并发性能,并最终减少了算法的运行时间,加快了算法的收敛。核心代码如下:

```
Set<CSVRecord> csvSet = new HashSet<>(csvParser.getRecords());
      ExecutorService executor =
2
      System.out.println(" 一共" + csvSet.size()+" 任务,任务开始");
      for (CSVRecord record : csvSet) {
         executor.submit(() -> {
             int now = totle.incrementAndGet();
             System.out.println("一共" + csvSet.size()+"任务,第"+now+"
             → 个任务开始");
             String[] relatedApis = record.get("Related APIs").split(",");
             //String[] Categories = record.get("Categories").split(",");
10
             Set<String> comparedres = new HashSet<>();
11
             Set<String> Query = new HashSet<>();
12
13
             for (String a : relatedApis) {
14
```

```
if(undirectedGraph.containsVertex(a)){
15
                       comparedres.add(a);
16
                       Query.add(PrimaryCategory.get(a));
17
                   }
18
19
               }
20
              /* for (String b : Categories) {
21
                   Query.add(b);
               7*/
23
               SteinerTree steinerTree1 = new
24
                   SteinerTree(undirectedGraph,PrimaryCategory);
               long startTime = System.currentTimeMillis();
25
               MidTree ress = steinerTree1.bulidSteinerTree(Query);
26
               long endTime = System.currentTimeMillis();
27
               if (ress != null) {
                   System.out.println(" 第"+now+" 个任务返回 api 共耗时"+
29
                      (endTime - startTime)+" 毫秒");
                   System.out.println(" 第"+now+" 个任务返回的群 Steiner Tree
30
                       的权重和为: "+ ress.getTotleweight());
                   Vector<String> resV = ress.getVertexSetofMidTree();
31
                   System.out.println(" 第"+now+" 个任务返回的群 Steiner Tree
                     的顶点个数为:"+ resV.size());
                   if(resV.size() <= comparedres.size())</pre>
33
                       lessthan.incrementAndGet();
                   if (comparedres.containsAll(resV)) {
                       int nowright = right.incrementAndGet();
35
                       System.out.println(" 第"+now+" 个任务命中"+" 这是第" +
36
                       → nowright+" 个命中");
                   }
37
                   else{
38
                       System.out.println(" 第"+now+" 个任务未命中");
                   }
40
               }
41
               else{
                   System.out.println(" 第"+now+" 个任务未命中");
43
               }
44
           });
45
```

```
46
       }
47
48
       executor.shutdown();
49
       try {
50
           executor.awaitTermination(Long.MAX_VALUE, TimeUnit.NANOSECONDS);
51
       } catch (InterruptedException e) {
52
           e.printStackTrace();
       }
54
       double hitrate = right.doubleValue()/totle.doubleValue();
55
       System.out.println(" 一共命中了"+right.get()+" 个任务, "+" 命中率为" +
56
       → hitrate*100+"%");
       System.out.println(" 若返回的 api 数少于所给的 api 数也算命中,则命中
57
       → 率为"+(lessthan.doubleValue()/totle.doubleValue())*100 + "%");
       System.out.println("任务结束");
   } catch (IOException e) {
59
       e.printStackTrace();
60
  }
61
```

3.2 测试方案

- 读取 mashup.csv 中的 relatedApis 作为验证结果,然后将查询这些 api 所属的种类作为查询的关键词,使用这些关键词调用算法,将 relatedApis 与算法返回的结点做比较,一致则说明命中。
- 一共命中了 1457 个任务, 命中率为 22.58
- 若返回的 api 数少于所给的 api 数也算命中,则命中率为 80.04029133736246
- 每个任务的平均耗时为 4054.669929669923(ms)
- 总用时 20 分钟,显著的减少了运行时间。

3.3 部分结果如下所示

- 第 4818 个任务返回 api 共耗时 336 毫秒
- 第 4818 个任务返回的群 Steiner Tree 的权重和为: 0.5
- 第 4818 个任务返回的群 Steiner Tree 的顶点个数为:2
- 第 4818 个任务命中这是第 1056 个命中

- 第 4578 个任务返回 api 共耗时 35202 毫秒
- 第 4578 个任务返回的群 Steiner Tree 的权重和为: 6.5621212121212125
- 第 4578 个任务返回的群 Steiner Tree 的顶点个数为:18
- 第 4578 个任务未命中
- 第 4840 个任务返回 api 共耗时 85 毫秒
- 第 4840 个任务返回的群 Steiner Tree 的权重和为: 0.5
- 第 4840 个任务返回的群 Steiner Tree 的顶点个数为:2
- 第 4840 个任务命中这是第 1063 个命中
- 第 4849 个任务返回 api 共耗时 8 毫秒
- 第 4849 个任务返回的群 Steiner Tree 的权重和为: 0.0
- 第 4849 个任务返回的群 Steiner Tree 的顶点个数为:1
- 第 4849 个任务命中这是第 1066 个命中
- 第 4708 个任务返回 api 共耗时 21958 毫秒
- 第 4708 个任务返回的群 Steiner Tree 的权重和为: 4.75
- 第 4708 个任务返回的群 Steiner Tree 的顶点个数为:12
- 第 4708 个任务未命中

```
第6050个任务返回api共耗时58127毫秒
```

第6050个任务返回的群Steiner Tree的权重和为: 86.85119047619048

第6050个任务返回的群Steiner Tree的顶点个数为:115

第6050个任务未命中

第6294个任务返回api共耗时36939毫秒

第6294个任务返回的群Steiner Tree的权重和为: 14.1249999999998

第6294个任务返回的群Steiner Tree的顶点个数为:28

第6294个任务未命中

第5848个任务返回api共耗时96702毫秒

第5848个任务返回的群Steiner Tree的权重和为: 26.32784250063662

第5848个任务返回的群Steiner Tree的顶点个数为:47

第5848个任务未命中

第6451个任务未命中

第6438个任务返回api共耗时66549毫秒

第6438个任务返回的群Steiner Tree的权重和为: 19.98053391053391

第6438个任务返回的群Steiner Tree的顶点个数为:49

第6438个任务未命中

超时未收敛

第6244个任务未命中

图 2: 预测结果输出

```
第6450个任务返回api共耗时30毫秒
第6450个任务返回的群Steiner Tree的权重和为: 0.0
第6450个任务返回的群Steiner Tree的顶点个数为:1
第6450个任务未命中
一共6453任务,第6451个任务开始
第6404个任务返回api共耗时3100毫秒
第6404个任务返回的群Steiner Tree的权重和为: 0.5
第6404个任务返回的群Steiner Tree的顶点个数为:2
第6404个任务未命中
一共6453任务,第6452个任务开始
第6452个任务返回api共耗时23毫秒
第6452个任务返回的群Steiner Tree的权重和为: 0.0
第6452个任务返回的群Steiner Tree的顶点个数为:1
第6452个任务命中这是第1456个命中
一共6453任务,第6453个任务开始
第6453个任务返回api共耗时13毫秒
第6453个任务返回的群Steiner Tree的权重和为: 0.0
第6453个任务返回的群Steiner Tree的顶点个数为:1
第6453个任务未命中
第6416个任务返回api共耗时2851毫秒
第6416个任务返回的群Steiner Tree的权重和为: 0.5
第6416个任务返回的群Steiner Tree的顶点个数为:2
第6416个任务未命中
```

图 3: 预测结果输出

```
一共6453任务,第5739个任务开始
```

第5739个任务返回api共耗时2毫秒

第5739个任务返回的群Steiner Tree的权重和为: 0.0

第5739个任务返回的群Steiner Tree的顶点个数为:1

第5739个任务命中这是第1276个命中

一共6453任务,第5740个任务开始

第5676个任务返回api共耗时6956毫秒

第5676个任务返回的群Steiner Tree的顶点个数为:3

第5676个任务未命中

一共6453任务,第5741个任务开始

第5741个任务返回api共耗时12毫秒

第5741个任务返回的群Steiner Tree的权重和为: 0.0

第5741个任务返回的群Steiner Tree的顶点个数为:1

第5741个任务未命中

一共6453任务,第5742个任务开始

第5686个任务返回api共耗时5788毫秒

第5686个任务返回的群Steiner Tree的权重和为: 0.125

第5686个任务返回的群Steiner Tree的顶点个数为:2

第5686个任务未命中

一共6453任务,第5743个任务开始

第5743个任务返回api共耗时24毫秒

第5743个任务返回的群Steiner Tree的权重和为: 0.0

第5743个任务返回的群Steiner Tree的顶点个数为:1

第5743个任务命中这是第1277个命中

一共6453任务,第5744个任务开始

第5744个任务返回api共耗时18毫秒

第5744个任务返回的群Steiner Tree的权重和为: 0.0

第5744个任务返回的群Steiner Tree的顶点个数为:1

第5744个任务未命中

图 4: 预测结果输出