

1 (Article)

A Method for Analysis and Visualization of Similar 2

Hotspot Flow Patterns between Different Regional 3

Groups 4

- 5 Haiping Zhang 1,2,3#, Xingxing Zhou 1,2,3#, Xin Gu 5, Genlin Ji 1,2,4 and Guoan Tang 1,2,3,*
- 6 ¹ Key laboratory of Virtual Geographic Environment, Ministry of Education, Nanjing Normal University; 7 Nanjing 210023, China; gissuifeng@163.com
- 8 State Key Laboratory Cultivation Base of Geographical Environment Evolution (Jiangsu Province), Nanjing 9 210023, China;
- 10 Jiangsu Center for Collaborative Innovation in Geographical Information Resource Development and 11 Application, Nanjing 210023, China;
- 12 College of Computer Science and Technology, Nanjing Normal University; Nanjing 210023, China
- 13 ⁵ Department of Geography Sciences, University of Maryland College Park, MD 20742, USA; 14 xgu12347@terpmail.umd.edu
- 15 * Correspondence: tangguoan@njnu.edu.cn; Tel.: +8613 776623891
- 16 #The contribution to the article is the same
- 17

18 **Abstract:** The interaction between different regions normally is reflected by the form of the stream. 19 For example, the interaction of the flow of people and flow of information between different regions 20 can reflect the structure of cities' network, and also can reflect how the cities function and connect 21 to each other. Since big data has become increasingly popular, it is much easier to acquire flow data 22 for various types of individuals. Currently, it is a hot research topic to apply the regional interaction 23 model, which is based on the summary level of individual flow data mining. So far, previous 24 research on spatial interaction methods focused on point-to-point and area-to-area interaction 25 patterns. However, there are a few scholars who study the hotspot interaction pattern between two 26 regional groups with some predefined neighborhood relationship by starting with two regions. In 27 this paper, a method for identifying a similar hotspot interaction pattern between two regional 28 groups has been proposed, and the Geo-Information-Tupu methods are applied to visualize the 29 interaction patterns. For an example of an empirical analysis, we discuss China's air traffic flow 30 data, so this method can be used to find and analyze any hotspot interaction patterns between 31 regional groups with adjoining relationships across China. Our research results indicate that this 32 method is efficient in identifying hotspot interaction flow patterns between regional groups. 33 Moreover, it can be applied to any analysis of flow space that is used to excavate regional group 34 hotspot interaction patterns.

35 Keywords: regional group interaction; similar hotspot flow patterns; spatial interaction; visual 36 analytics; Geo-Information-Tupo; GIS.

37

38 1. Introduction

39 Our society is built based on mobility, such as the flow of people, the flow of goods and the flow 40

of information technology. And these elements of flow form a flow space[1]. Compared to traditional 41

- local space, the flow space pays more attention on the interaction of elements and their interaction 42 relationship[2,3]. In the past, geographers were focus on physical space[4-6]. Nowadays, with the
- 43
- increasing development of economic globalization and Internet technology, geography researchers 44

transfer their sights to flow space[7-10]. On one hand, the outstanding change of economic

45 globalization is: people have strengthened their exchanges in tourism, trade and technology from all 46 over the world, thus directly leading the advanced enhancement on the flow of people, logistics and 47 technology; On the other hand, information flow has further weakened the distance between places 48 with the development of internet technology. One apparently fact is that: the distance is no longer 49 suitable to apply for the metric space when the time required to transmit information for 1 kilometer 50 is almost the same as the time required to transmit information for 100,000 kilometers, that is to say, 51 the connection of Internet has realized the change on transmission of spatial information. In fact, for 52 geographers, the main points not only should be the flow space itself, but it is also about how these 53 flow elements reconstruct the spatial organization structure, how to make the organization works, 54 and what kind of flow patterns emerge[11]. Based on all these facts above, it is important for us to 55 use the quantitative analysis method to excavate the interaction patterns and define the interaction 56 patterns, because they are the basis methods to solve spatial relationships between two regional 57 groups.

58 Over the past few decades, many methods have already been proposed to find out interaction 59 patterns of flow space. In terms of spatial interaction model, some scholars have built many spatio-60 temporal interaction pattern mining algorithms from a summary sight[12-16]. However, the spatial 61 dependence of interactive nodes among all these methods are lacked. Some of them apply complex 62 network methods to discover the spatial interaction patterns[17-20]. The conception of interaction 63 regions model based on the idea of complex networks has been proposed by some other scholars^[21]. 64 To some extent, the dependencies and similarities between flowing nodes are considered, especially 65 the method of interaction relation which is proposed by Kira are able to identify areas with strong 66 interactivity. But the limitation is that it only recognizes all the individual regions that are similar in 67 interaction, rather than the interaction modes between different regions. Kwanho Kim et al. has 68 proposed a regional mobile pattern recognition algorithm(MZP) based on the aggregation of metro 69 nodes recently. Based on his idea[22], Chen et al. expanded the proximity relationship and realized 70 the mobile pattern recognition based on taxi OD data, and MPFZ algorithm is proposed[23]. All these 71 methods are mainly focus on point data and its adjacency relationship, the main disadvantage is that 72 their algorithm is inefficient and have not given a visually well-resolved solution to excavate 73 interaction model. So in terms of visualizing the interaction patterns results, none of the above 74 methods can solve the interaction pattern between two regional groups.

75 A basic characteristic of the existing models is: it lacks a more interactive flow pattern 76 recognition method that may exist between one regional group and another, and it's hard to define 77 the adjacency relationship between two regional groups. We can refer to the adjacency matrix related 78 literature for the definition of regional adjacency relations. For example, there is a strong interaction 79 between region A and region B (B does not have a predefined proximity relationship). Not only that, 80 but also between several regions around A, and regions around B. We assume there is a strong 81 interaction relationship, and a pre-defined adjacency relationship is satisfied between region A and 82 its surroundings, and also a pre-defined adjacency relationship between region B and its 83 surroundings, then regional group A and regional group B are located. So we can conclude that there 84 is a strong interaction between A or B and surroundings, more importantly, a regional group 85 interaction flow pattern is formed between regional group A and regional group B.

86 This paper presents an advanced method for discovering, analyzing and visualizing the 87 interaction hotspot flow patterns between two different regional groups. During the next section, a 88 review of related work has been written, and the expected results of the method will also be 89 expressed. After the second section, a new method which is used to mine regional group flow 90 patterns has also been proposed. During that part, we mainly include the definition of regional 91 adjacent relationship, the structuration of flow pattern mining algorithm, the introduction of flow 92 pattern visualization and methodological issues. In the end, an air flow volume data will be shown 93 during the case study part by using the section three methods.

95 2. Related work

96 In most cases, individual flow data will be modeled as a flow pattern from node to node [24, 25]. 97 It is also for the above reasons that many of the macro mode summary or interactive mode discovery 98 methods for individual flow data are based on node flow data [26-28]. In contrast, there is few flow 99 data modeling and analysis between regions, moreover, the interactions between regions can also be 100 abstracted as point-to-point interactions. It is easy to use basic spatial analysis methods to achieve the 101 goal even if the flow data from point to point is aggregated to the region-to-region flow data. 102 However, interaction modeling and analysis between regional groups will involve many issues such 103 as how to identify and determine the regional adjacency relationship, and better visualize the 104 expression. Most of the existing researches are based on the first two cases. Today we will have a 105 brief introduction to the existing related research below. In order to better understand the limitations 106 of the research objectives and the existing methods, we will also discuss point-to-point flow pattern, 107 area-to-area flow pattern and flow pattern of two different regional groups, but we only mention the 108 one having a strong relationship.

109 As we mentioned above, most of the flow data exists in the form of point-to-point with arrow. 110 Related interaction analysis methods mainly include point-to-point interaction pattern mining [29-111 33], interactive pattern mining in between multiple points, and a model analysis of adjacent points in 112 a same community [34-37]. We can see clearly from Figure 1(a) that the interaction between the three 113 nodes in the northwest corner and the two nodes in the southeast corner are significantly stronger 114 than other flow data. A similar situation exists between the neighboring points in the southwest and 115 northeast corner. As is shown on Figure 1(b), the MZP algorithm proposed by Publication et al. can 116 discover a strong interaction pattern between a set of adjacent nodes in a network structure data to 117 another set of adjacent nodes. Then, the two modes shown on Figure 1(c) could be identified. The 118 MZP algorithm mainly represents a prefect method for solving such problems, and provides valuable 119 reference value for related researches. However, the time complexity of this algorithm process is too 120 high, and the visualization of the analysis results has not yet proposed a good solution. Based on that, 121 Chen et al. proposed the MPFZ method, but Chen's method only extended the data which was 122 applied by the network MZP algorithm from the network node flow to other analysis of the arbitrary 123 node flow data. No other major changes have been improved in other areas.



Figure 1. An example for point-to-point flow data and its analysis methods.(a)point-to-point flow data;(b)points-to-points flow data;(c)points-to-points flow patterns.

We may pay more attention to the interaction pattern between different areas for flow data in some cases. For example, for the point-to-point with arrow data shown in Figure 1(a), we can easily obtain the area based area-to-area flow data through the basic spatial superposition and statistical analysis methods. As is shown in Figure 2(a), the arrow must also contain an attribute to indicate the size of the interaction value for each area-to-area flow data. Based on the results shown on Figure 2(a), we can easily identify the regional interaction shown on Figure 2(b), and thus obtaining the

133 region interaction pattern shown on Figure. 2(d).



134

Figure 2. An example for region flow data and its analysis methods.(a)area-to-area flow data with high
 interaction values;(b);(c)area-to-area flow patterns.

137 Concerning area-to-area model, the obvious disadvantage is that each area interaction mode 138 does not consider the correlation characteristics of the starting and ending area with other existing 139 adjacent areas, which means the spatial autocorrelation of any area interaction mode and the 140 surrounding area in interaction directions and sizes. As shown in Figure 3(a), it is much more 141 significant that the interaction between several adjacent areas in the northwest and southeast among 142 the area-to-area flow data. Also similar patterns are applied on the southwest and northeast sides. As 143 shown in Figure 3(b), the goal of this paper is to identify the existence of regional group interaction 144 flow patterns by defining specific area adjacency relationships, Figure 3(c) shows the results and 145 visualization of the flow pattern that is expected. And then do further research on the interaction 146 strength, value size and significant level of each regional group based on the results of the analysis.





150 3. Methodology

151 The entire research framework includes the input of node-based flow data, data processing, and 152 mining of regional group flow patterns, flow pattern output, and visualization. Since most of the flow 153 data is counted and then stored by nodes, this study supports node-based flow data input during the 154 design process. Firstly, the input node-to-node flow data is converted according to a certain regional 155 unit and then converted into regional-to-regional flow data. This process can be realized by using the 156 common GIS overlay and statistics functions. Then determine the adjacency relationship of the 157 regional units (Section 3.1.1), and based on this adjacency relationship, merge the adjacent areas 158 where the interaction value reaches a certain threshold before being constructed into regional groups. 159 After that, we are able to identify all similar hotspot flow patterns among different regional groups 160 (Section 3.1.2). In the end, the Geo-Information-Tupu visualization method is used to present the 161 regional groups with similar hotspot flow patterns and visual variables are used to visualize the 162 evaluation results of their own characteristics in each flow pattern. The rest of the writing, we refer

163 to the similar hotspot flow pattern between regional groups as RG-Flow-Pattern.



164

Figure 4 Overview of the framework for analysis and visualization the similar hotspot flow patternsbetween regional groups.

168 3.1 Building algorithm for similar hotspot pattern between regional groups

169 In this study, the regional hotspot interaction model algorithm mainly includes three aspects. 170 They are 1) Defining the regional neighborhood relationship. 2) Constructing a regional hotspot 171 interaction pattern recognition algorithm based on the defined neighborhood relationships. 3) 172 Multiple test parameters are used to evaluate the results of the identified area hotspot interaction 173 models.

174 3.1.1 Regional adjacency relationship modeling

175 In order to identify the hotspot interaction pattern, we must clearly define the regional adjacency 176 relationship and its merger principle. In this method, four ways are defined for determining the 177 adjacency relationship of the area. As shown in figure 3, if each grid is used as a region, the adjacency 178 relationship between regions can be expressed as the following four ways: figure3 (b), 3(c), 3(d) and 179 3(e). As shown in figure 3(a), if we assume the target area is the red one, the specific meanings of the 180 four adjacency relationships are briefly described as below.

181 1. Adjacent edges

182 In figure 3(b), there are four areas have common edges with target area, and these four areas are 183 defined as adjacent areas of the target area. The adjoining relationship in this manner is called an 184 edge-adjacent relationship. In an actual partition, under this rule, a target area may have more 185 adjacent areas or less than four adjacent areas.

186 2. Adjacent edges and corners

Figure 3(c) shows a similar adjacency relationship to figure 3(b). However, except that the area having a common edge with the target area belongs to the adjacent edge of the target area, it also includes an area having a common node with the target area. This kind of adjoining relationship is called edge-corner adjoining.

191 3. Customized adjacent range

In figure 3(d), a circular buffer area is defined with the center of mass of the target area as the origin. When other areas are within or intersect the buffer area, they are defined as the adjacent areas of the target area. In this method, the adjacency relationship is called the adjoining relationship of customized adjacent range.

196 4. Logical adjacent relationship

In addition to the above three methods to define the adjacency relationship, we can also determine whether the target area and the other areas are adjacent by customizing the logical relationship that is independent of the spatial position. In figure 3(e), there are some logical relations between the three blue areas and the target area. Therefore, even though these areas do not coincide with the target area or coincide with the vertices, these three areas are defined as the adjacent areas of the target area.

Basically, the above four are the typical modeling methods for the spatial relationship of surface
 features. Other adjacencies include k-nearest, custom based on spatial adjacency matrix, etc.



207

Figure 4. similar interaction hotspot flow pattern and it's visualizaiton among regional groups.(a);(b);(c);(d);(e).

208 3.1.2 Region merge and similar hotspot flow pattern recognition

209 1. Definitions of similar hotspot flow pattern between regional groups

210 In this research, a set of data sets containing n planar area units were given, Rset =211 $\{R_1, R_2, ..., R_n\}$ (i = 1,2, ..., n), R_i represents the nth region. In a regional group interactive hotspot flow 212 pattern, the origin area group is defined as $RGOset = \{R_1, R_2, ..., R_u\}$, and the destination area group 213 is defined as RGDset = { $R_1, R_2, ..., R_v$ }. In addition, a pair of origin areas and destination areas that 214 have interactions in a regional group interaction hotspot pattern are called regional flow. A data set 215 RFset was given to store all the regional flow in Regional group interaction hotspot flow pattern, (216 RIH-FP), $RFset = \{RF_1, RF_2, \dots, RF_m\}(j = 1, 2, \dots, m)$, The *jth* regional flow can be represented as 217 $RF_i = RF_{i_o} \rightarrow RF_{i_d}, RF_{i_o} \in RGOset$, it indicates the origin area of the regional flow. $RF_{i_d} \in RGDset$, it 218 represents the destination area of this regional flow. In some situations, for ease of exposition, we use 219 the term flow pattern instead of RIH-FP in the remainder of the paper. There are some definitions 220 about the flow pattern.

221 **Definition1:** A regional group interactive hotspot flow model consists of three parts. They are 222 the starting regional group RGOset, the destination regional group RGDset, and the interaction 223 direction indicating the interaction relationship. A regional group hotspot interaction pattern has the 224 same direction as any $RF_j = RF_{j_0} \rightarrow RF_{j_d}$ in the RFset.

225 Definition2: Given an area-adjacent relationship defined in 3.1.1, a single region in the RGOset 226 of the origin region group must satisfy such a adjacent relationship, and a single region in the 227 destination region group RGDset must also satisfy this adjacent relationship.

228 **Definition3:** The number of regions in the origin and destination regional group of a regional 229 group interactive hot spot flow mode cannot be 1 at the same time, that is, at least more than one 230 region is included in the start or termination regional group.

231 Definition4: The interaction value of the regional flow refers to the interaction value from one 232 region to another, which is represented by InterVal. This value has different meanings in different 233 applications, but the following conditions must be required:

Given a threshold θ , the interaction strength value $P(RF_j)$ of the j-th regional stream RF_j must satisfy the following conditions:

$$P(RF_j) = \frac{InterVal(RF_{j_o} \to RF_{j_d})}{InterVal(RF_{j_o} \to RF_{s_d})*InterVal(RF_{s_o} \to RF_{j_d})} (P(RF_j) \ge \theta)$$
(1)

237 InterVal $(RF_{j_o} \rightarrow RF_{j_d})$ represents the interaction value which is from origin area RF_{j_o} to 238 destination area RF_{j_d} . InterVal $(RF_{j_o} \rightarrow RF_{*_d})$ represents the sum of the interaction values of the 239 origin region RF_{j_o} to all other destination regions. InterVal $(RF_{*_o} \rightarrow RF_{j_d})$ represents the sum of 240 the interaction values of all the origin regions to the destination region RF_{j_d} .

Definition5: The RFset, which contains all regional flow in the same flow pattern, is no pre defined adjacent relationship from the starting region(s) to the ending region(s) in any regional flow
 RF.

244 2. Region merge

245 Firstly, we randomly select a group of regional flow data that satisfy: $P(RF_i) \ge \theta$, $RF_i = RF_i \circ A$ 246 RF_{id} and make $RF_i = RF_{id} \rightarrow RF_{id}$ as the first region flow of a new regional interactive hotspot 247 flow pattern, and the interaction value size is expressed as Inter Val($RF_{i,o} \rightarrow RF_{i,d}$). Then using $RF_{i,o}$ 248 as the starting regional group elements of the new region interactive hotspot flow mode, and satisfy 249 $RF_{i,o} \in RGOset$. Using $RF_{i,d}$ as the new regional interactive hotspot flow mode, which is the 250 termination elements of regional group, it should satisfy $RF_{id} \in RGDset$. Search for all regions 251 adjacent to RF_{j_o} , whose set is defined as ARGOset = { $RF_{j_o_1}, R_{j_o_2}, \dots, R_{j_o_u}$ }(m = 1,2, ..., u), the mth 252 adjacent region of RF_{j_0} is $RF_{j_0,m}$; all regions adjacent to RF_{j_d} are searched in the same way, and 253 the set is defined as ARGDset = $\{RF_{j_d_1}, R_{j_d_2}, \dots, R_{j_d_v}\}$ (n = 1,2, ..., v),. The nth adjacent of RF_{j_d} is 254 $RF_{j_{d_n}}$. For $RF_{j_{d_n}}$ in any ARGOset, if $RF_{j_{d_n}}$ interacts with the area $RF_{j_{d_n}}$ in the ARGDset, it 255 constitutes the regional flow $RF_{j_m_n} = RF_{j_{o_m}} \rightarrow RF_{j_{d_n}}$, then:

$$P(RF) =$$

$$\frac{InterVal(RF_{j_o} \rightarrow RF_{j_d}) + InterVal(RF_{j_o} \rightarrow RF_{j_d}, n)}{(InterVal(RF_{j_o} \rightarrow RF_{*_d}) + InterVal(RF_{j_o} \rightarrow RF_{*_d}, n)) * (InterVal(RF_{*_o} \rightarrow RF_{j_d}) + InterVal(RF_{*_o} \rightarrow RF_{j_d}, n))}$$
(2)

256

Among them, InterVal($RF_{j_o} \rightarrow RF_{j_d}$) is the interaction value of the regional flow RF_j , InterVal($RF_{j_o} \rightarrow RF_{*_d}$) indicates the sum of the interaction values of the starting area RF_{j_o} to all other termination areas RF_{*_d} , InterVal($RF_{*_o} \rightarrow RF_{j_d}$) represents the sum of the interaction values of all other starting regions RF_{*_o} to the ending region RF_{j_d}). Similarly, InterVal $InterVal(RF_{j_o_m} \rightarrow RF_{j_{d_n}})$ is the interaction value of the regional flow $RF_{j_{m_n}}$, $InterVal(RF_{j_{o_m}} \rightarrow RF_{*_{d_n}})$ indicates the sum of the interaction values of the starting area $RF_{j_{o_m}}$ to all other areas. $InterVal(RF_{*_{o_m}} \rightarrow RF_{*_{d_n}})$ indicates $RF_{j_d_n}$ indicates all other areas to the interaction value of $RF_{j_d_n}$.

264 After calculating the P(RF) value, if $P(RF) \ge \theta$, then RF_{i_0} is also included at the origin 265 regional group of the regional group interaction mode, $RF_{i o m} \in RGOset$ is satisfied, The $RF_{i d n}$ is 266 included in the termination zone group of the regional group interaction mode, $RF_{i,d,n} \in RGD$ set is 267 satisfied. After all the above is completed, statistical analysis is performed on other adjacent areas by 268 the same method, and it is known that an area does not meet the merge threshold and the merge 269 operation is ended. The newly included start and end regions are then searched for their adjacent 270 regions, and the above operations are iterated until no region satisfies the merge threshold. Finally, 271 a complete regional interaction hotspot flow mode start zone group and termination zone group are 272 obtained.

273 3. Regional interaction hotspot flow pattern recognition

Through the merging of the upper part of the region, a starting zone group and an ending zone group of several regional interactive hotspot modes are formed. For an area interaction hotspot flow pattern RIH-FP if the set of start area groups is defined as: $RGOset = \{R_1, R_2, ..., R_u\}$, the ending regional group is defined as $RGDset = \{R_1, R_2, ..., R_v\}$, the set of regional flow is defined as $RFset = \{RF_1, RF_2, ..., RF_m\}(j = 1, 2, ..., m)$. RF_p represents the p-th region flow, and RF_q represents the q-th region flow. The initial regional group RGOset, the termination area group RGDset, and the interaction stream set RFset between the two regional groups constitute a complete regional hotspot interaction flow mode. The direction of interaction between the regional groups is indicated by the directional arrows. Thus, the start region group, the termination region group, and the direction arrow constitute a basic visualization element of an area hotspot interaction flow pattern and form a feature structure of the flow pattern. Based on a complete regional interaction hotspot flow pattern, in addition to the visual elements and feature structure, some evaluation values are needed to distinguish the strength of each flow pattern. If the variable P is used to indicate the strength of a certain RIH-FP, then:

$$P = \sum_{j=1}^{m} P(RF_j) \tag{3}$$

288

 $P(RF_j)$ represents the interaction strength value of the jth regional flow in the regional flow set RFset. The interaction strength of the entire RIH-FP is the sum of the values of all the regional flow interaction strengths in the RFset.

If the V denote the size of the interaction value of a certain RIH-FP, then V should satisfy the following formula:

$$V = \sum_{j=1}^{m} Interval(RF_j)$$
(4)

294

295 $Interval(RF_j)$ represents the interaction value of the jth region flow in the regional flow set 296 RFset. The interaction value of the entire RIH-FP is the sum of all the regional flow interaction value 297 in the RFset.

In addition, it is also necessary to separately calibrate the contribution of each of the start regional group and the termination regional group to the current flow mode interaction value in a complete mode. For the i-th region R_i in the starting regional group RGOset:

$$DO(R_j) = \frac{\sum_{j=1}^{m} Interval(RF_j)}{InterVal(RF_{j_0} \to RF_{*_d})} (R_j \in \text{RGOset}, RF_j \in \text{RFset}, RF_{j_0} = R_i)$$
(5)

301

302 For the i-th region R_i in the termination regional group RGDset:

$$DD(R_j) = \frac{\sum_{j=1}^{m} Interval(RF_j)}{InterVal(RF_{*_o} \to RF_{j_d})} (R_j \in \text{RGOset}, RF_j \in \text{RFset}, RF_{j_d} = R_i)$$
(6)

303

304 3.2 SHFP-RG visualization method based on Geo-information Tupo theory

305 3.2.1 Visualization of single RG-Flow-Pattern

306 In the RG-Flow-Pattern method of this paper, the analysis results are evaluated and investigated 307 by using different flow pattern variables. These variables have both an assessment of the starting and 308 ending regional groups as well as an overall assessment of the interaction model. Viewing these 309 evaluation variables that match a particular pattern in a table is not desirable for spatial pattern 310 analysis. It also loses the advantage of visualizing spatial data analysis results based on maps and 311 further visual analysis. Therefore, it is very important to design a scientific and reasonable RG-Flow-312 Pattern visualization method. So based on the above facts, the RG-Flow-Pattern visualization method 313 is designed as shown in Figure6 (a) and (b). Figure6 (a) and (b) are two basic examples of RG-Flow-314 Pattern visualization. The basic meanings and expression purposes of the two model examples are 315 described in detail below.

316 As we mentioned earlier, a complete RG-Flow-Pattern contains three basic constructs, namely 317 the start regional group, the termination regional group, and the directional arrows. In order to

325

326

327

318 visualize the results of each RG-Flow-Pattern, the interaction value size, and the contribution rate of

each RG-Flow-Pattern in each of the start and termination regional group, some Visual variables such
as color and size are expressed. As shown in Figure. 6(a) and Figure.6(b), if one proceeds from the
basic definition, it is obvious that the basic requirements of the RG-Flow-Pattern structure are
satisfied.



Figure 6. Two simple example for single RG-Flow-Pattern visualization and instrument its meaning.(a) A regional interaction hot spot flow pattern with low interaction value. (b)A regional interaction cold spot flow pattern with high interaction value.

328 Comparing the two findings, there are significant differences in the overall color design of the 329 regional group. Figure 6(a) shows a warm tone, while Figure 6(b) shows a cool tone. The purpose of 330 this design is to express the strength of each RG-Flow-Pattern by means of cool and warm colors. The 331 warm tone indicates that the RG-Flow-Pattern behaves in a strong interactive mode, and the cool 332 color represents the performance of the RG-Flow-Pattern behaves in a weak interactive mode. The 333 degree of strength is measured by the P value in equation (2). The critical value of strength is divided 334 according to the overall distribution of P values of all models by using natural discontinuity method, 335 quantile method, etc., and the user of the model can definite it by themselves. Obviously in the two 336 examples given in this paper, Figure 6(a) belongs to the strong regional interaction flow mode, where 337 the hot spot flow mode is further defined. Figure 6(b) belongs to the weaker interactive flow mode, 338 which is further defined as the cold spot flow mode. In addition to the differences in the cool and 339 warm tones of the regional groups as a whole, there are also differences among the inner regions of 340 each RG-Flow-Pattern. This represents the contribution rate of a single region to the current RG-Flow-341 Pattern interaction value. The darker the color, the greater the contribution rate of the region to the 342 RG-Flow-Pattern interaction value, and vice versa. The contribution rate is measured by Equation (4) 343 and Equation (5). The former measures each mode. The contribution of a single zone in the starting 344 regional group, which is used to measure the contribution rate of a single region in the termination 345 regional group for each flow mode. This rule can be applied to both hot flow pattern and cold spot 346 flow pattern. The first two parts of the legend shown in Figure 6(c) illustrate the specific meanings 347 and corresponding relationships between the expression flow pattern strength and the contribution 348 rate of interaction values in each region in the visualization results.

In addition, RG-Flow-Pattern also needs to evaluate the value of the overall model interaction value through the value of V, so as to make up for the inadequacy of the interaction value that can be used to evaluate the strength of the interaction model. In the visualization, the size of the V value is expressed by the thickness of the arrow, which indicates the current RG-Flow-Pattern interaction value. Comparing Figure. 6(a) and Figure. 6(b), although RG-Flow-Pattern in Figure. 6(a) shows a
strong flow pattern, the interaction value is smaller than that in Figure. 6(b). The flow pattern
direction portion of Figure. 6(c) is a legend of the interaction value size relationship.

We can conclude that in addition to the directional arrows including the starting regional group and the ending regional group, the group of cooling and heating tone variables representing the strong and weak P value of the interaction mode, a saturation vision variable of a single region contribution rate V value to the current mode interaction value, and an arrow size vision variable representing the size of the flow pattern interaction value are also included in a complete visualization result of RG-Flow-Pattern.

362 3.2.2 Visualization and classification of multiple RG-Flow-Patterns based on Geo-information Tupu

363 In the traditional spatial data distribution and visualization mode, the distribution pattern of the 364 same topic and region can be presented on a map. For example, the classic analysis method Local 365 moran's I and General G index for analyzing the local spatial autocorrelation, the analysis of the 366 models are easy to present on the same map. However, it is difficult to present on the same map for 367 the regional group interactive hotspot flow mode. As shown in Figure. 6, pattern-01 and pattern-02 368 belong to two different flow patterns in the same region, but both patterns have a single repeating 369 unit in both the real regional group and the termination regional group, which means, it is difficult 370 for such situations to express two modes on the same map.

371 In the 90 decade of the 20th century, the theory and method of Geo-Information-Tupu put 372 forward by Chen can be used to solve this problem[38]. In Chen's Geo-Information-Tupu theory, it 373 emphasizes the structuring, abstraction, type, and relevance features of geographic laws, and uses 374 these principles in a map sequence. Since in many cases, it is difficult to present multiple RG-Flow-375 Patterns in the same map, and different RG-Flow-Patterns of the same topic can also be type-divided, 376 the map sequence can be adopted by the Geo-information-Tupu method. The RG-Flow-Pattern map 377 sequence can be arranged according to types, and can also be arranged according to interaction 378 strength, interaction value size. Since the interaction strength and interaction values can be directly 379 organized by P value and Z value, thus only the type division of the RG-Flow-Pattern map is 380 introduced in this paper.

- In fact, for RG-Flow-Patterns, the type division is also a relatively simple task. In this paper, RG Flow-Patterns is divided into two basic types and complex types. The basic types mainly include the
- 383 five types shown in figure 8.



Figure 7. An example of the same region belong to different patterns



Figure 8. Basic categories of RG-Flow-Pattern based on Geo-information-Tupu.(a) many-to-many
 regions and single direction RG-Flow-Pattern.(2) many-to-many regions and double directions RG-Flow Pattern.(c) one-to-many single direction RG-Flow-Pattern. (d) many-to-one single direction RG-FP.(e) one
 and many double direction RG-Flow-Pattern.

391 4. Case Study: the national scale migration flow data of China

392 *4.1 Study area and data descriptions*

393 Due to work, leisure travel and other purposes, a large number of people travel from one place 394 to another every day. Human mobility can reflect lots of issues, such as urban attractiveness, tourism 395 resources and so on. China has a population of 1.3 billion and there are significant differences in 396 economic, political, cultural and resource characteristics in different regions. The huge imbalance in 397 population size and regional disparities further promotes population movements. In terms of 398 transportation, China's national-wide cross-regional transportation includes three types of 399 transportation: automobiles, trains, and aircraft. A car is more suitable for short trips, the train is 400 more suitable for people with short-to-medium-distance or low- and middle-income groups, while 401 the airplane is mainly for long-distance or high-income travel. Because the method proposed in this 402 paper is more effective in the analysis of flow data across regions, this paper uses the migratory flow 403 data of the Chinese mainland as the main data source, and the prefecture-level city as the smallest 404 research unit. We adopt the RG-Flow-Pattern method to develop the empirical analysis. Figure 8 405 shows the distribution of the population migration routes (by airplane) for the main study area on 406 April 1, 2017. It should be noted that only the top ten data inflows and relocations from each 407 prefecture-level city are used here.

408 The demographic data provided by the Tencent location big data platform was used in this 409 research. Tencent is a major Internet company in China that provides nationwide location-based real-410 time migration big data services. On this platform, daily migration data from mainland China are 411 provided. The migration types include aircraft, trains, and automobiles. Also, the top ten regions by 412 rank of flow data was included, and the degree of hotspot flow value of moving in and out was 413 calculated. Among the three modes of transportation migration data, the flight data has the longest 414 distance, and the RG-Flow-Pattern method is better for this analysis. Therefore, the population 415 migration data of flights was analyzed in this paper. The data used in this study involves 315 cities,

- 416 and the total number of data points for all of the cities is approximately 6300, including flow data
- 417 with original city, destination city, and hot value as the main attributes.



Figure 9. Study area and visualization of flow data

420 4.2 Result

418 419

421 The RG-Flow-Pattern method proposed in this paper was adopted, and set the prefecture-level

- 422 city was a regional unit and the modal method of spatial relationship shown in Figure 5(c) was used, 423
- then set the θ value of $P(RF_i) \ge \theta$ to 0.00001. The partial patterns in the analysis result are shown in
- 424 figure 10 below:





Figure 10. Four examples of RG-Flow-Patterns Geo-Infomaion-Tupu by threshold of 0.00001.(a) a
coldspot RG-Flow-Pattern.(b) a hotspot RG-Flow-Pattern. (c) a cold spot RG-Flow-Pattern. (d) a hotspot
RG-Flow-Pattern. (e) legend for RG-Flow-Patterns.

430 As shown in Figure 10(a), through the RG-Flow-Pattern algorithm analysis, it found that some 431 regions in the southwestern part of China (the red part) and the eastern part of the coastal area (the 432 blue part) form the regional group interaction flow model. Figure 10(a) shows the geographical 433 distribution of the flow pattern on the left, and figure 10(a) shows the pattern representation of this 434 pattern on the right. As can be seen from the latter, the pattern belongs to the cold spot flow pattern, 435 and the direction of flow pattern is from the southwest area to the eastern coastal area. The color of a 436 single area represents the contribution of the flow of that area to the entire pattern. The southwest 437 area used as the starting regional group of the flow pattern, in which the color depth of each area 438 represents the contribution of the sum of the values of the area flowing out to the termination area 439 group to the outflow value (also called the outdegree) of the entire model; The coastal area is the most 440 frequent end-of-flow model, in which the color depth of each area indicates the contribution rate of 441 the inflow value of the area to the inflow value of the entire model. The darker the color, the greater 442 the contribution rate. Refer to the legend shown in Figure 10(e) for details.

Figure 10(b) and figure 10(d) are interactive cold spot flow patterns recognized by the RG-Flow-Pattern algorithm. Figure. 10(c) is another set of identified regional group interaction hotspot flow patterns.

446 **5.** Discussion and conclusions

This section may be divided by subheadings. It should provide a concise and precise description
of the experimental results, their interpretation as well as the experimental conclusions that can be
drawn.

451 5.1 Discussion

452 5.1.1 Selection pricinple of region ajacency relationship and region merge threshold

453 In this case, the adjacent edges and corners approach is used for the adjacency of the area, which 454 means that this approach is considered as the adjoining area of the target area as long as there is an 455 edge or corner adjacent to the target area. When we model the area's adjacency, other methods 456 mentioned in 3.2.1 section can be chosen. However, based on RG-Flow-Pattern analysis, using different 457 regional adjacency relationships, there may also be differences in the models. This is the impact of 458 regional adjoining relationships on the model. Among the specific issues, it is recommended to refer 459 to the selection principles of regional spatial relationships in spatial statistical methods such as 460 Moran's I, the Geary index, and Geographically Weighted Regression (GWR). Another problem is 461 that when the value of θ in $P(RF_i) \ge \theta$ is different, the resulting flow pattern may also be different. 462 The larger the value of θ , the smaller the number of flow patterns to be formed. The number of areas 463 in the flow pattern that make up the start area group and the termination area group also decreases. 464 To solve this problem, the recommended practice is to first obtain the $P(RF_i) \ge \theta$ values for all 465 regional flows, and then use the bar histogram to evaluate the distribution of all regional flow 466 $P(RF_i) \ge \theta$ values and select them according to the analysis target. A reasonable threshold is taken 467 as the value of θ . This method can control the number and strength of flow patterns to a certain 468 extent. So in this case study, figure 11 shows the plot distribution of P values for all regional flows.

469 5.1.2 result evaluation

470 In a complete flow pattern, both the basic elements of the flow pattern (starting regional group, 471 termination regional group, and interaction arrows) are included, as well as the interaction strength, 472 interaction value size, each individual flow pattern and the rate of contribution of the area's traffic to 473 the interaction value of the entire flow pattern. Although this design makes it possible for each flow 474 pattern to contain enough information to evaluation itself, the disadvantages are also obvious. First 475 of all, these assessments are for a single flow model and lack the assessment of the overall 476 characteristics of all models. For a single flow mode, starting from the strength of the mode and the 477 size of the interaction value, there are four situations: firstly, a strong interaction mode with a large 478 interaction value; secondly, a weak interaction mode with a small interaction value; thirdly, a strong 479 interaction mode with a small interaction value; and fourthly, a weak interaction mode with a large 480 interaction value. For all the overall characteristics of the model, it is obviously very useful for 481 subsequent analysis to understand these four scenarios. If the strength and interaction values of each 482 flow pattern can be described by XY coordinate system, the four cases can be expressed clearly and 483 transparently through the four-quadrant diagram.

484 5.1.3 shortcomings and future improvements

485 The RG-Flow-Pattern method realizes that all flow patterns satisfying a certain intensity are 486 recognized from the mass flow data, and a plurality of visual variables are used to better express the 487 pattern and the related evaluation amount. However, there still exist some deficiencies. First of all, 488 although the goal of this method is to analyze any type of flow data such as people flow, logistics, 489 and traffic flow, for some flow data with a shorter interaction distance, it is difficult to find two cross-490 regional regional groups by this algorithm. This means that this method is more suitable for the 491 mining of regional group interaction patterns between regions with long interaction distances. 492 Although one can solve this problem by setting smaller partitions, more often than not, the interactive 493 areas used for analysis are predefined and show some geographic significance, and cannot be 494 customized for their size. In subsequent studies, we will try to build a flow data model mining model 495 that is suitable for short interaction distances based on this method. Secondly, in a complete regional 496 group interaction flow model, a strong self-interactive mode may exist between a single region of a 497 starting regional group and a single region of an ending regional group, and the RG-Flow-Pattern 498 method cannot recognize their self-interactive mode in this case Although it is not considered in the

499 RG-Flow-Pattern method, this self-interactive pattern mining method is relatively simple. Its main 500 challenge is how this kind of self-interactive mode plays a role in the flow model of this article and 501 how it can be improved. The expression is performed in a visual manner to facilitate subsequent 502 visual analysis. These are the tasks that need to be further improved.

503 5.2 Conclusion

504 With globalization and the development of the Internet, geographers have turned their attention 505 from physical space to flow space. Spatial analysis methods also have been extended from spatial 506 pattern analysis to spatial interaction pattern discovery. Although spatial interaction has always been 507 the focus of the GIS field, with the advent of big data technologies, spatial interactions and even 508 space-time interactions have successfully attract the attention of scholars nowadays. Many 509 researchers mainly focus on point-to-point, area-to-area, or interaction-based research on regional 510 convergence or diffusion. Few people consider the interaction patterns that may exist between 511 regional groups that have some sort of adjoining relationship. In fact, the interaction of most flow 512 data does not only exist between two separate areas, but the interaction always happens between a 513 group of areas and another regional group.

514 We assume that two different regions, the relationship of one area to another is formed since an 515 imbalance in certain resources. Furthermore, since this kind of imbalance, the surrounding area of 516 one certain region has similarity demand of this resource, so it leads the target area and its 517 surroundings with limited sources (we call it regional groups) interacts with other regional groups 518 having abundant resources. The area and the surrounding areas that also have such resources interact 519 with the surrounding areas that require such resources but lack them, forming an interaction between 520 the two regional groups. In this paper, the RG-Flow-Pattern analysis and visualization method 521 proposed can effectively mine the possible interaction patterns between two regional groups under 522 such scenarios. In this analysis method, not only can all the regional groups having such interaction 523 relationships which satisfy a specific traffic threshold be identified, but also the level of the strength 524 of each group of interaction flow modes, the size of the interaction of the modes, and each of the 525 interaction variables and the extent to which the area contributes to the overall model interaction 526 volume can be measured by some outcome variables.

527 The first law of geography is the basic principle of the GIS spatial analysis model, that is, the 528 spatial unit has spatial correlation characteristics. In the past, spatially-distributed characteristics tent 529 to be considered in analytical models in spatial distribution models and spatial relationship 530 modeling. Concomitant with the "interactive" turn of the GIS analysis model, and under the 531 perspective of flow space, the spatial flow model or spatial interaction model should also be 532 considered as spatial correlation. However, describing the spatial flow model is more complex than 533 the spatial distribution model and the spatial relationship modeling, and it is difficult to visualize all 534 the patterns through a single map. In this paper, based on the consideration of the relevance of 535 neighboring regional units, we proposed a spatial group interaction model analysis method, and at 536 the same time, geo-information maps was used to express the analysis results model and to deal with 537 the difficulty of single diagram visualization. This analysis method can be extended to mine regional 538 data interaction relationship in any other flow data forms.

539 6. Patents

Acknowledgments: This work is supported by National Natural Science Foundation of China under Grants No.
 41471371, and supported by National Natural Science Foundation of China, No.41671389 We would like to
 express our sincere appreciation to the anonymous reviewers for their insightful comments that have greatly

543 aided us in improving the quality of this paper.

544 References

Marty, P.F. An introduction to digital convergence: Libraries, archives, and museums in the information
 age. *Libr Quart* 2010, *80*, 1-5.

547	2.	Andris, C.; Liu, X.; Ferreira, J. Challenges for social flows. Computers, Environment and Urban Systems 2018,
548		70, 197-207.
549	3.	Andris, C. Integrating social network data into gisystems. Int J Geogr Inf Sci 2016, 30, 2009-2031.
550	4.	Midler, J.C. Non-euclidean geographic spaces: Mapping functional distances. <i>Geogr Anal</i> 2010, 14, 189-203.
551 552	5.	Alamri, S.; Taniar, D.; Safar, M.; Al-Khalidi, H. A connectivity index for moving objects in an indoor cellular
553	6	Wang LE: Li X H: Christakos C: Liao XI: Zhang T: Cu X: Zhang X X Coographical detectors based
554	0.	health rick assessment and its application in the neural tube defects study of the bechun region, china. <i>Lut</i>
555		L Construct Study of the residur region, china. Int
556	7	J Geogr III) Set 2010 , 24, 107-127.
557	7.	Developments in the dutch urban system on the basis of nows.
558	0	Reg Stuu 2009, 45, 179-196.
550	о.	McKenzie, G.; Janowicz, K.; Gao, S.; Gong, L. How where is when? On the regional variability and
560	0	resolution of geosocial temporal signatures for points of interest. <i>Comput Environ Urban</i> 2015, 54, 336-346.
561	9.	Tao, K.; Thill, J.C. Spatial cluster detection in spatial flow data. <i>Geogr Anal</i> 2016 , 48, 355-372.
562	10.	Seto, K.C.; Reenberg, A.; Boone, C.G.; Fragkias, M.; Haase, D.; Langanke, T.; Marcotullio, P.; Munroe, D.K.;
502 5(2		Olah, B.; Simon, D. Urban land teleconnections and sustainability. <i>P Natl Acad Sci USA</i> 2012 , <i>109</i> , 7687-7692.
563	11.	Zhu, X.; Guo, D.S. Mapping large spatial flow data with hierarchical clustering. <i>Transactions In Gis</i> 2014 , <i>18</i> ,
564		421-435.
363	12.	Adams, P.C. A taxonomy for communication geography. <i>Prog Hum Geog</i> 2011 , <i>35</i> , 37-57.
566	13.	Mesbah, M.; Currie, G.; Lennon, C.; Northcott, T. Spatial and temporal visualization of transit operations
567		performance data at a network level. <i>Journal Of Transport Geography</i> 2012 , 25, 15-26.
568	14.	Fonte, C.C.; Fontes, D.; Cardoso, A. A web gis-based platform to harvest georeferenced data from social
569		networks: Examples of data collection regarding disaster events. <i>Int J Online Eng</i> 2018 , <i>14</i> , 165-172.
570	15.	Hale, M.L.; Ellis, D.; Gamble, R.; Walter, C.; Lin, J. Secuwear: An open source, multi-component
571		hardware/software platform for exploring wearable security. <i>Ieee Int Conf Mo</i> 2015, 97-104.
572	16.	Li, M.; Sun, Y.R.; Fan, H.C. Contextualized relevance evaluation of geographic information for mobile users
573		in location-based social networks. Isprs Int J Geo-Inf 2015, 4, 799-814.
574	17.	Li, J.W.; Ye, Q.Q.; Deng, X.K.; Liu, Y.L.; Liu, Y.F. Spatial-temporal analysis on spring festival travel rush in
575		china based on multisource big data. Sustainability-Basel 2016, 8.
576	18.	Rosvall, M.; Bergstrom, C.T. Maps of random walks on complex networks reveal community structure. P
577		Natl Acad Sci USA 2008 , 105, 1118-1123.
578	19.	Esquivel, A.V.; Rosvall, M. Compression of flow can reveal overlapping-module organization in networks.
579		<i>Phys Rev X</i> 2011 , <i>1</i> , 1668-1678.
580	20.	Zhou, M.; Yue, Y.; Li, Q.Q.; Wang, D.G. Portraying temporal dynamics of urban spatial divisions with
581		mobile phone positioning data: A complex network approach. Isprs Int J Geo-Inf 2016, 5.
582	21.	Kempinska, K.; Longley, P.; Shawe-Taylor, J. Interactional regions in cities: Making sense of flows across
583		networked systems. Int J Geogr Inf Sci 2018, 32, 1348-1367.
584	22.	Kim, K.; Oh, K.; Lee, Y.K.; Kim, S.; Jung, J.Y. An analysis on movement patterns between zones using smart
585		card data in subway networks. Int J Geogr Inf Sci 2014 , 28, 1781-1801.
586	23.	Chen, Z.L.; Gong, X.; Xie, Z. An analysis of movement patterns between zones using taxi gps data.
587		<i>Transactions In Gis</i> 2017 , <i>21</i> , 1341-1363.
588	24.	Liu, L.A.; Hou, A.Y.; Biderman, A.; Ratti, C.; Chen, J. Understanding individual and collective mobility
589		patterns from smart card records: A case study in shenzhen. 2009 12th International leee Conference on
590		Intelligent Transportation Systems (Itsc 2009) 2009 , 1-6.
591	25.	Munizaga, M.A.; Palma, C. Estimation of a disaggregate multimodal public transport origin-destination
592		matrix from passive smartcard data from santiago, chile. Transport Res C-Emer 2012, 24, 9-18.
593	26.	Ghasemzadeh, M.; Fung, B.C.M.; Chen, R.; Awasthi, A. Anonymizing trajectory data for passenger flow
594		analysis. Transport Res C-Emer 2014, 39, 63-79.
595	27.	Zhang, Y.P.; Martens, K.; Long, Y. Revealing group travel behavior patterns with public transit smart card
596		data. Travel Behav Soc 2018 , 10, 42-52.
597	28.	Chu, K.K.A.; Chapleau, R. Enriching archived smart card transaction data for transit demand modeling.
598		<i>Transp Res Record</i> 2008 , 63-72.

599	29.	Higuchi, T.; Shimamoto, H.; Uno, N.; Shiomi, Y. A trip-chain based combined mode and route choice
600		network equilibrium model considering common lines problem in transit assignment model. State Of the
601		Art In the European Quantitative Oriented Transportation And Logistics Research 2011, 20.

- 602 30. Concas, S.; DeSalvo, J.S. The effect of density and trip-chaining on the interaction between urban form and
 603 transit demand. *Journal Of Public Transportation* 2014, 17, 16-38.
- Shou, L.; Ji, Y.X.; Wang, Y.Z. Analysis of public transit trip chain of commuters based on mobile phone data and gps data. 2017 4th International Conference on Transportation Information And Safety (Ictis) 2017, 635-606
 639.
- 607 32. Blythe, P.T. Improving public transport ticketing through smart cards. *P I Civil Eng-Munic* **2004**, *157*, 47-54.
- 608 33. Pelletier, M.P.; Trepanier, M.; Morency, C. Smart card data use in public transit: A literature review.
 609 *Transport Res C-Emer* 2011, *19*, 557-568.
- 610 34. Wang, Y.; Lim, E.P.; Hwang, S.Y. Efficient algorithms for mining maximal valid groups. *Vldb Journal* 2008, 611 17, 515-535.
- 612 35. Aung, H.H.; Tan, K.L. Discovery of evolving convoys. *Scientific And Statistical Database Management* 2010, 6187, 196-213.
- 614 36. Li, Y.X.; Bailey, J.; Kulik, L. Efficient mining of platoon patterns in trajectory databases. *Data & Knowledge*615 *Engineering* 2015, 100, 167-187.
- 616 37. Williams, H.J.; Holton, M.D.; Shepard, E.L.C.; Largey, N.; Norman, B.; Ryan, P.G.; Duriez, O.; Scantlebury,
 617 M.; Quintana, F.; Magowan, E.A., *et al.* Identification of animal movement patterns using tri-axial
 618 magnetometry. *Mov Ecol* 2017, 5.
- 8. Ye, Q.; Tian, G.; Liu, G.; Ye, J.; Yao, X.; Liu, Q.; Lou, W.; Wu, S. Tupu methods of spatial-temporal pattern on land use change. Journal of Geographical Sciences 2004, 14, 131-142.