

# <sup>1</sup> Outlier Detection and Comparison of <sup>2</sup> Origin-Destination Flows Using Data Depth

<sup>3</sup> **Myeong-Hun Jeong<sup>1</sup>**

<sup>4</sup> Department of Civil Engineering, Chosun University, Gwangju, Republic of Korea

<sup>5</sup> mhjeong@chosun.ac.kr

<sup>6</sup>  [0000-0003-4850-8121]

<sup>7</sup> **Junjun Yin<sup>2</sup>**

<sup>8</sup> Social Science Research Institute; Institute for CyberScience, Penn State University, PA, USA

<sup>9</sup> jyin@psu.edu

<sup>10</sup>  [0000-0002-4196-2439]

<sup>11</sup> **Shaowen Wang<sup>3</sup>**

<sup>12</sup> CyberGIS Center for Advanced Digital and Spatial Studies; Department of Geography and

<sup>13</sup> Geographic Information Science, University of Illinois at Urbana-Champaign, IL, USA

<sup>14</sup> shaowen@illinois.edu

---

## <sup>15</sup> — Abstract —

<sup>16</sup> Advances in location-aware technology have resulted in massive trajectory data. Origin-destination  
<sup>17</sup> (OD) trajectories provide rich information on urban flow and transport demand. This study de-  
<sup>18</sup> scribes a new method for detecting OD flows outliers and conducting hypothesis testing between  
<sup>19</sup> two OD flow datasets in terms of the variations of spatial extent, that is, spread. The proposed  
<sup>20</sup> method is based on data depth, which measures the centrality and outlyingness of a point with  
<sup>21</sup> respect to a given dataset in  $\mathbb{R}^d$ . Based on the center-outward ordering property, the proposed  
<sup>22</sup> method analyzes the underlying characteristics of OD flows, such as location, outlyingness, and  
<sup>23</sup> spread. The ability of the method to detect OD anomalies is compared with that of the Ma-  
<sup>24</sup> halanobis distance approach, and an F-test is used to verify the difference in scale. Empirical  
<sup>25</sup> evaluation has demonstrated that our method effectively identifies OD flows outliers in an in-  
<sup>26</sup> teractive way. Furthermore, the method can provide new perspectives such as spatial extent by  
<sup>27</sup> considering the overall structure of data when comparing two different OD flows in terms of scale.

<sup>28</sup> **2012 ACM Subject Classification** Computing methodologies → Anomaly detection

<sup>29</sup> **Keywords and phrases** Movement Analysis, Trajectory Data Mining, Data Depth, Outlier De-  
<sup>30</sup> tection

<sup>31</sup> **Digital Object Identifier** 10.4230/LIPIcs.GIScience.2018.6

---

## <sup>32</sup> **1 Introduction**

<sup>33</sup> With ubiquitous geolocation-aware sensors, knowledge discovery is greatly enhanced by  
<sup>34</sup> extracting and mining interesting patterns from spatiotemporal big data in various domains.  
<sup>35</sup> Massive movement data are collected to track people, animals, vehicles, and even natural  
<sup>36</sup> phenomena. Such data help us better model moving objects and reveal hidden patterns that

---

<sup>1</sup> [This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. 2018R1C1B5043892).]

<sup>2</sup> [This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by the National Science Foundation grant number ACI-1548562]

<sup>3</sup> [This work was supported by the U.S. National Science Foundation (grant numbers: 1047916 and 1443080)]

## 6:2 Outlier Detection and Comparison of Origin-Destination Flows using Data Depth

37 are important to urban planning [17], understanding human mobility [30, 11], achieving the  
38 sustainability of urban systems [1, 3] and the environment [4], and improving public security  
39 and safety [2].

40 This paper a new method that identifies origination-destination (OD) flow anomalies  
41 and conducts hypothesis testing between two sets of different OD flows. In this study, the  
42 OD flow data represents a particular type of trajectory data, which records the origin and  
43 destination of each movement while ignoring the exact trajectory route [9]. The method was  
44 applied to OD flows derived from New York City taxi trip records, in which each record  
45 contains the origin and destination of each trip, without intermediate locations of the actual  
46 routes.

47 In recent years, researchers have investigated a variety of approaches to trajectory data  
48 mining. Most contemporary trajectory mining methods can be classified into four categories:  
49 clustering, classification, frequent/group pattern mining, and outlier detection [18, 33]. These  
50 methods can be used independently or together for trajectory mining applications. This study  
51 focuses on outlier detection of OD flows. Outlier detection aims to identify trajectories that  
52 do not follow the typical flows of trajectory that characterize the connectivity between regions  
53 [18]. Euclidean distance is employed by [7, 13] to find outlier patterns from trajectories.  
54 Studies by [20, 14] question the Euclidean distance approach because of the loss of local  
55 features and unavailability when external factors, such as topography, land cover or weather  
56 condition, affect the trajectories. In their research, [20, 14] addresses this by using robust  
57 distance measurements, e.g., Mahalanobis distance [20] and relative distance [14]. Instead of  
58 using distance or density, anomalous trajectories are detected by exploiting comparisons of  
59 the structural features of each trajectory segment [31] and an isolation tree of trajectories  
60 [32]. Most of these methods are related to trajectory data analysis, and thus, it is reasonable  
61 to extend the application of these approaches to the identification of OD flow anomalies.  
62 To overcome the sensitivity of Euclidean distance-based approaches to non-normal data  
63 distribution and the difficulty of selecting parameters for anomaly detection techniques based  
64 on distance or density, this study employs robust statistics, such as data depth, to detect  
65 OD flow outliers.

66 Flow mapping, a type of visual analytics, is a common approach to analyzing OD flow data.  
67 Visual representations of massive movement data facilitate comprehensive exploration of  
68 data, in turn enabling interpretation and understanding of complex flow trends. Aggregation  
69 and generalization of movement data are frequently utilized to resolve visual clutter [9, 29].  
70 While visual analytics can help to extract inherent patterns from massive data, it is difficult  
71 to quantitatively compare two sets of different OD flows based on hypothesis testing. In  
72 other words, it is complicated to comprehend how two OD flows differ and, more importantly,  
73 the magnitude of the difference, using a test of statistical significance. Recently published  
74 articles employ multidimensional spatial scan statistics [8] and local Ripley's K-function [23]  
75 to identify clusters of flow data based on statistical significance testing. In a similar vein,  
76 this paper applies bivariate hypothesis testing methods based on data depth to understand  
77 the difference between two OD flow datasets in the context of different spatial extents.

78 It is worth noting that flow mapping approaches frequently suffer from the modifiable  
79 areal unit problem (MAUP). Essentially, MAUP reflects the influence of different aggregations  
80 determined by location on the identification and representation of coherent patterns. Kernel-  
81 based flow estimation and smoothing are used to overcome different spatial resolutions [9].  
82 Instead of attempting to find the best areal unit by which to partition urban space and  
83 aggregate the OD flows, this study adopted the established traffic analysis zones of New  
84 York City as a base unit. That said, the proposed method can be adapted to other areal

units. In this study, New York City taxi trip data includes origins and destinations within traffic analysis zones, while ignoring the intermediate locations of the actual routes. Note that it is not necessary to reconstruct individual movements for flow estimation (see [5]).

In summary, this paper presents a new algorithm which conducts outlier detection as well as hypothesis testing on OD flow data. Our approach investigates the central regions of OD flows, based on data depth, to detect OD flow anomalies and conduct hypothesis testing between two different OD flow datasets. We believe that our method for analyzing taxi trip data has the potential to aid administrative authorities to better understand crowd patterns for improving urban planning activities such as determining transportation investments.

The remainder of this paper is organized as follows: Section 2 overviews how to detect OD flow outliers and conduct hypothesis testing between two different OD flow datasets using the concept of data depth. Experimental design and the evaluation of the proposed method are presented in Section 3. These results are discussed in Section 4. Section 5 concludes this paper with a summary and future work perspectives.

## 2 Methods

### 2.1 Data Depth

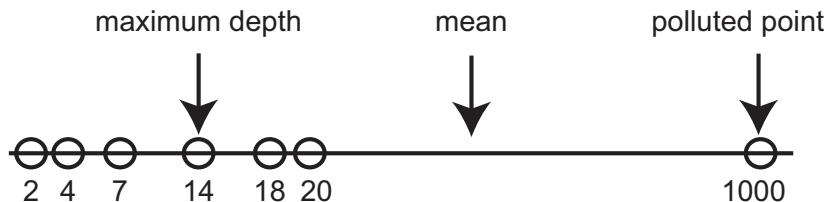
Data depth measures the centrality of a point with regard to a given dataset in  $\mathbb{R}^d$ . Originally developed by [24], the notion of data depth (i.e., halfspace depth) generalizes the univariate concept of ranking to multivariate data. Halfspace depth represents how deeply a point is located within a given dataset by ordering all points according to their degree of centrality.

Generally, the halfspace depth (HD) of point  $x$  in  $\mathbb{R}^d$  is defined as the minimum probability,  $P$  on  $\mathbb{R}^d$ , associated with any closed halfspace containing  $x$  [34].

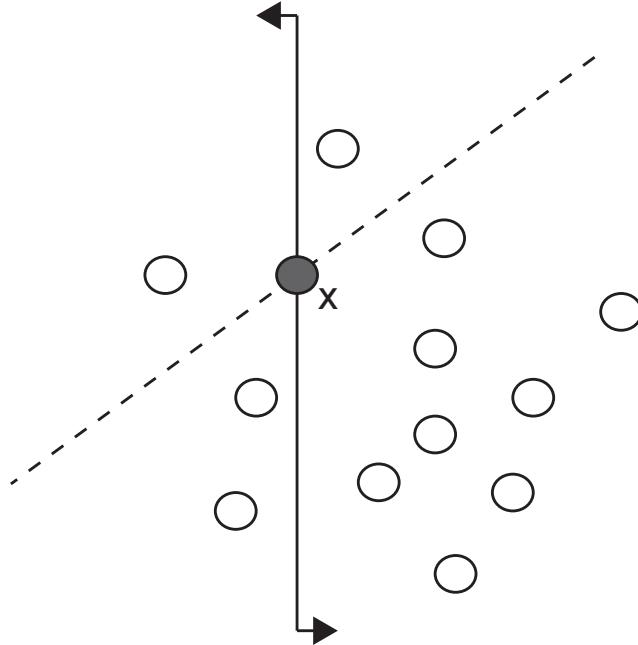
$$HD(x; P) = \inf\{P(H) : H \text{ is a closed halfspace}, x \in H\}, x \in \mathbb{R}^d.$$

For the univariate case, all values less than or equal (greater than or equal) to  $x$  form a closed halfspace. All values less (greater) than  $x$  are an open halfspace. The smallest probability associated with two closed halfspaces developed by  $x$  is the halfspace depth of point  $x$ . In Figure 1, the probability of values less than or equal to 4 is 2/7 and the probability of values greater than or equal to 4 is 6/7. Thus, the halfspace depth of 4 is 2/7, which is the minimum probability carried by any closed halfspace containing 4. Furthermore, as the sample median, 14 has the largest halfspace depth. Note that the polluted point inflates the standard error of the sample mean, thereby distorting the view of the data.

Similarly, the halfspace depth of  $x$  for the bivariate case is defined by the minimal number of data points in any closed halfspace, which is determined by a hyperplane through  $x$  [21]. In Figure 2, the solid line through  $x$  is rotated by 180°. The halfspace depth of  $x$  is determined by the smallest portion of data separated by such a hyperplane. For example, the halfspace



**Figure 1** Robustness of halfspace depth for the univariate case



■ **Figure 2** Halfspace depth for the bivariate case

depth of  $x$  is  $3/13$ , as determined by the dotted line. However, the halfspace depth of  $x$  determined by the solid line is  $4/13$ . Therefore, the halfspace depth of  $x$  is  $3/13$ , which is the minimal number of data points in any closed halfspace through  $x$ .

The property of halfspace depth is a center-outward ordering of points in  $\mathbb{R}^d$  and is affine invariant [19]. These features make halfspace depth a useful tool in nonparametric inference, which leads to various applications such as data classification and cluster analysis [12, 10]. There are multiple approaches to calculating data depth, including halfspace depth [21], projection depth [25], and simplicial depth [15]. While the computational complexity of the projection approach is  $\mathcal{O}(n^2)$  (where  $n$  is the number of points), the computational complexity of simplicial depth is  $\mathcal{O}(n^3)$ . This can significantly increase computing time when  $n$  is large. Thus, this paper uses the more efficient method proposed by [21], in which the computational complexity for both approaches is  $\mathcal{O}(n \log n)$ .

## 132 2.2 OD Flow Outlier Detection Based on Data Depth

133 The center-outward ordering in data depth is closely related to the detection of outliers. The  
 134 upper level sets of data depth in  $\mathbb{R}^2$  form the central regions. The most central region can  
 135 be regarded as a median. Conversely, the lower level sets of data depth, which coincide with  
 136 larger distances from the center, can be regarded as outlyingness. This concept was utilized  
 137 by [22, 28] to generate bag plots, which are analogous to one-dimensional box plots based  
 138 on data depth. This paper uses the bag plot to identify the outliers of OD flows. Before  
 139 explaining the method of outlier detection, we first introduce a basic definition of OD flow.

140 ▶ **Definition 1.** Origin-destination (OD) flow. The OD flow  $OD_i = (o_i, d_i, c_i, ts_i, te_i)$  is the  
 141 number of trips ( $c_i$ ) from the origin ID ( $o_i$ ) to destination ID ( $d_i$ ) of traffic analysis zones  
 142 between the start time ( $ts_i$ ) and the end time ( $te_i$ ), where  $ts_i < te_i$ .

143 Based on this basic definition, Figure 3 depicts bag plots representing the OD flows  
 144 of New York City taxi data collected on May 21, 2014 and July 1, 2014 respectively. We  
 145 exploited taxi data on May 21, 2014 because the National September 11 Memorial Museum  
 146 and Pavilion was opened to the public on this date. We also randomly selected another  
 147 data set on July 1, 2014. In Figure 3a, the deepest depth of OD flows, depth median, is  
 148 represented by a star symbol. This point is surrounded by a dark blue bag, which contains  
 149 the half of OD flows. This region is regarded as a central region of OD flows. The OD flows in  
 150 the bag are the dominant patterns. Magnifying the bag by a factor of three, relative to depth  
 151 median, constructs a fence, as indicated by the light-blue area. The fence is comparable to  
 152 the whiskers of a one-dimensional boxplot. The OD flows outside the fence, represented by  
 153 red circles, are outliers. Every OD pair is represented by a point in Figure 3. The x-axis  
 154 indicates the counts of forward OD flows (e.g., the number of OD flows from origin ID 2 to  
 155 destination ID 10), and the y-axis indicates the counts of reverse OD flows (e.g., the number  
 156 of OD flows from origin ID 10 to destination ID 2) in Figure 3a.

157 The bag plot presents the data using the following attributes: location is represented by  
 158 the depth median; spread or the spatial extent of bag; correlation or the orientation of the  
 159 bag; and skewness, as represented by the shape of the bag and the fence [22]. In Figure 3a,  
 160 we observe that some forward OD flows have higher counts than their paired reverse OD  
 161 flows. We also note the relatively linear correlation between forward OD flows and reverse  
 162 OD flows and the skewness of forward (reverse) OD flows.

163 It is also possible to detect the outliers of OD flows of two different time stamps. In  
 164 Figure 3c, we visualize the OD flows recorded on two different days. Comparing the two  
 165 sets of OD flows not only indicates the central region of OD flows, it also distinguishes the  
 166 significantly different OD flows.

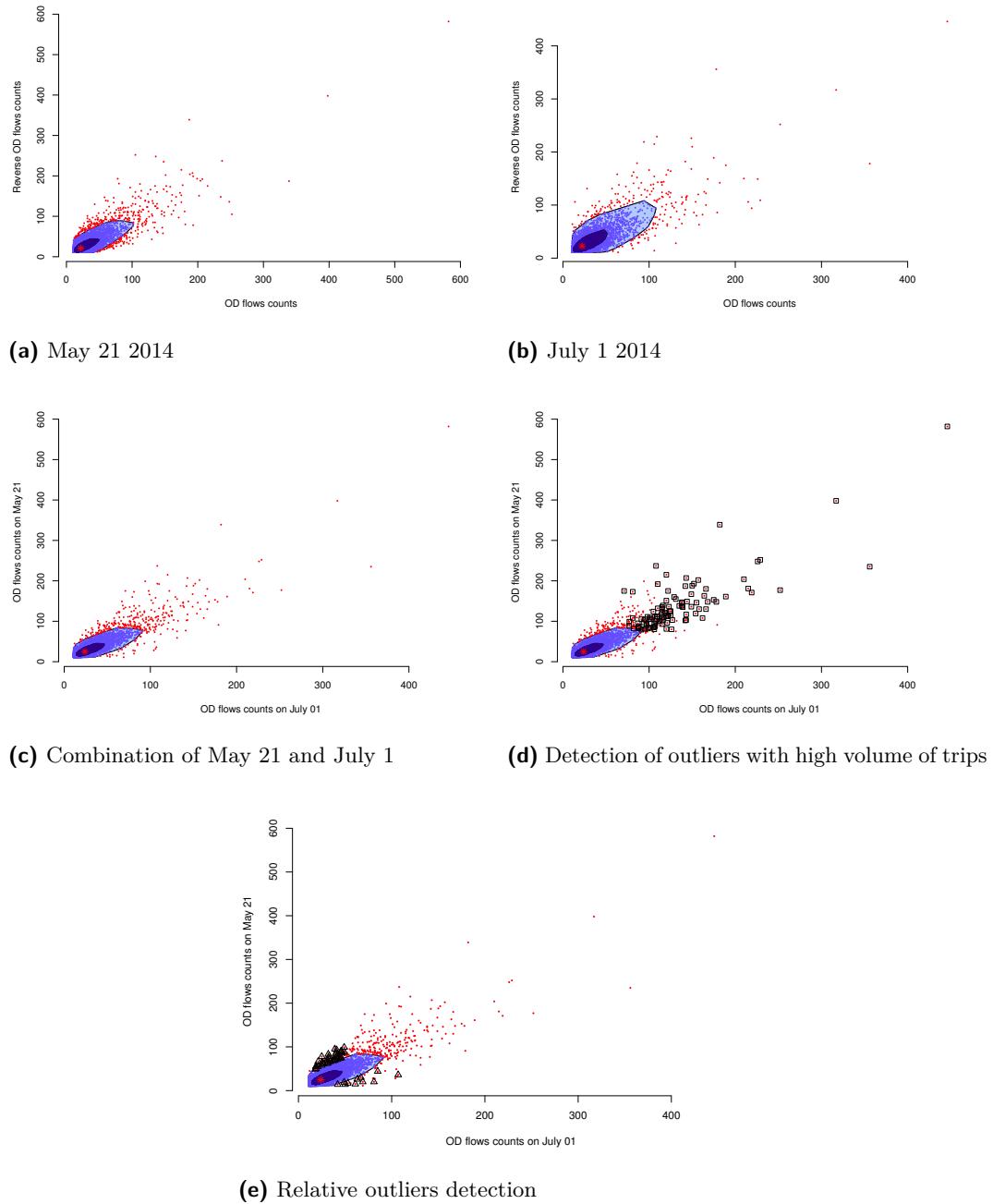
167 The OD flows in high activity areas of a city are more likely to have large trip volumes.  
 168 We use set operations to detect such outliers. We regard OD flows on July 1 as the control  
 169 dataset (*control*); OD flows on May 21 as test dataset (*test*); and the combination of two  
 170 OD flows as combination dataset (*combination*) in Figure 3. Then we can calculate the  
 171 intersection of three outliers sets (*control*  $\cap$  *test*  $\cap$  *combination*), which are represented as  
 172 rectangle symbols in Figure 3d.

173 In addition, it is interesting to detect the outliers of OD flows which are typical patterns  
 174 at time  $t_1$  but atypical behaviors at time  $t_2$ . We define the union of points in the bag, the  
 175 central region, at time  $t_1$  and  $t_2$ . Then we calculate the intersection of two sets, the outliers  
 176 of the combination set and the previous union set. These outliers are represented as triangle  
 177 symbols in Figure 3e. These outliers are typical OD flows at time  $t_1$ , located in the central  
 178 regions in the bag plot. When we consider two OD flows together, they become unusual OD  
 179 flows, some have more trips and some have fewer trips, relative to the control dataset. Thus,  
 180 we can detect and treat outliers interactively based on data depth.

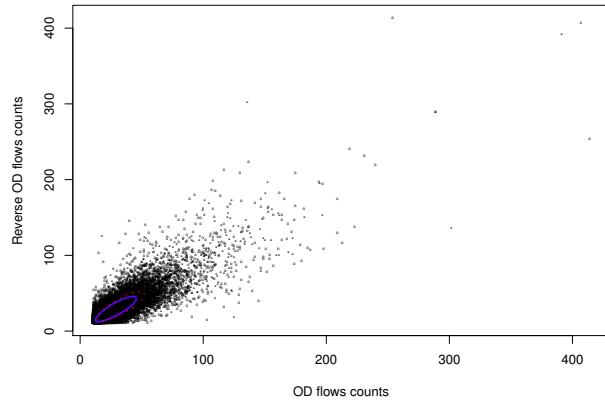
### 181 2.3 OD Flow Comparisons Based on Data Depth

182 Data depth can compare bivariate data from two independent groups. A *t*-test can be used  
 183 to compare means from two independent groups. For example, the *t*-test reveals whether the  
 184 means of two OD flows are different between two different temporal ranges. However, it is  
 185 also worth examining how groups differ in terms of scale, which is also referred to as spread.  
 186 Comparisons of central regions in data depth evaluate the marginal distribution, thereby  
 187 considering the overall structure of the data [26].

188 Let  $X$  and  $Y$  be the random variables having distributions  $F$  and  $G$  for two independent  
 189 groups. The quality index proposed by [16] is the probability that the depth of  $Y$  is greater



**Figure 3** Outliers detection of OD flows using a bag plot



**Figure 4** Central regions of two OD flows:  $\circ$  indicates the OD flows for Saturday, March 29 2014 and  $*$  indicates the OD flows for a list of Saturdays; blue line presents the central region of the OD flows for the list of Saturdays and red dotted line presents the central region of the OD flows on March 29.

190 than or equal to depth of  $X$ .

$$191 \quad Q(F, G) = P[D(X; F) \leq D(Y; F)],$$

192 where  $P$  is the probability and  $D(X; F)$  is the depth of randomly sampled observations  
193 according to distribution  $F$ . The range of  $Q$ , as presented by [16], is  $[0, 1]$  and  $Q(F, G) = 0.5$   
194 if and only if  $F = G$ . If  $Q < 0.5$  or if  $Q > 0.5$ , the scale increases or decreases from  $F$  to  $G$ .  
195 Therefore, it is possible to detect differences in scale using a bootstrap method.

196 Let  $X_1, \dots, X_a$  be a random sample from  $F$ , and  $Y_1, \dots, Y_b$  be a random sample from  $G$ .  
197 The estimate of  $Q(F, G)$  is calculated as shown below.

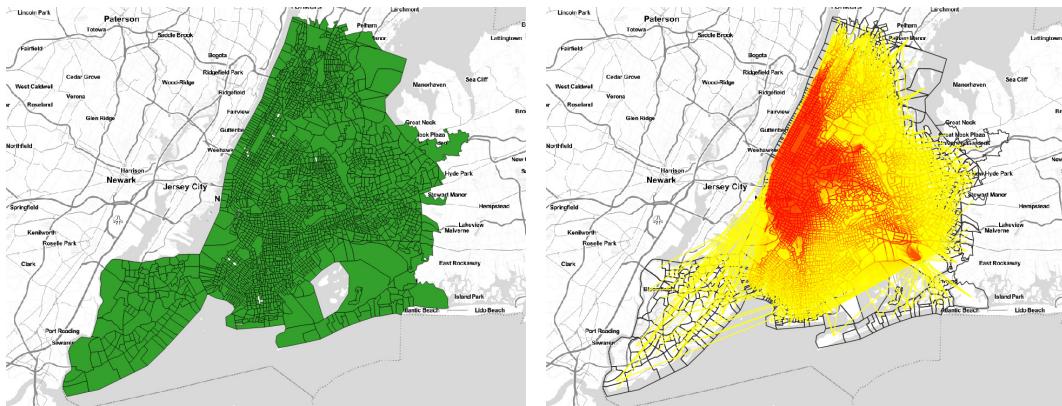
$$198 \quad \hat{Q}(F, G) = \frac{1}{b} \sum_{i=1}^b R(Y_i; F_a),$$

199 where  $R(Y_i; F_a)$  indicates the proportion of  $X_j$  which has  $D(X_j; F_a) \leq D(Y_i; F_a)$ . Simil-  
200 arly, the estimate of  $Q(G, F)$  can be defined as follows:

$$201 \quad \hat{Q}(G, F) = \frac{1}{a} \sum_{i=1}^a R(X_i; G_b).$$

202 Bootstrap samples are obtained by resampling from the two groups ( $F$  and  $G$ ). Under the  
203 null hypothesis ( $H_0 : Q(F, G) = Q(G, F)$ ), the difference of the resulting bootstrap estimates  
204 is  $Q^*(F, G) - Q^*(G, F)$ . Thus, if the confidence interval of  $Q(F, G) - Q(G, F)$  does not  
205 contain zero, we can reject the null hypothesis,  $H_0$  [16, 26].

206 For ease of understanding, Figure 4 presents the central regions of two OD flows. One  
207 dataset is OD flows for Saturday, March 29, 2014, and the other dataset includes multiple  
208 Saturdays, those of March 1, 8, 15, 22, and April 5. At 552,064 taxi trips, the day of March  
209 29 had the highest number of taxi trips for the year of 2014. The dataset for the other five



(a) 2,250 traffic analysis zones in New York City (b) OD flows on July 1 2014

**Figure 5** Experimental data: New York City taxi data

210 Saturdays comprised 2,621,703 taxi trips. The bootstrap method reveals that the confidence  
 211 interval is 0.0247 and 0.0596. This confidence interval does not include zero, thus rejecting  
 212 the  $H_0$  null hypothesis. This indicates that scale range is significantly changed between two  
 213 OD flow datasets. Furthermore, the OD flows from the group of Saturdays are nested within  
 214 the OD flows corresponding to March 29. This additional perspective was based on data  
 215 depth comparisons.

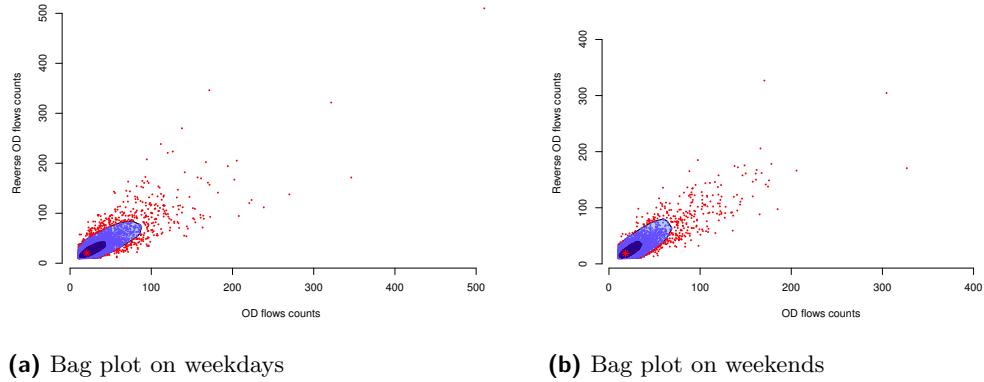
216 The bootstrap method is a time consuming process. For this study, we generate 2,000  
 217 bootstrap samples. To improve the efficiency of the bootstrap computation, we distributed  
 218 the work across multiple computing nodes and cores by implementing an embarrassingly  
 219 parallel R code.

### 220 3 Experiments

#### 221 3.1 Data

222 This study uses New York City taxi data collected in 2014 to evaluate the effectiveness  
 223 of the proposed approach. Figure 5a presents traffic analysis zones in New York City  
 224 which indicate the origin and the destination IDs of the OD flows. A traffic analysis zone  
 225 (TAZ) is the most commonly adopted basic geographic unit in transportation planning  
 226 models. The geographic areas of TAZ are delineated by transportation officials for tabulating  
 227 traffic-related data. The size of TAZ varies because it accounts the underlying popula-  
 228 tion in each zone, which consists of one or more census blocks, block groups, or census  
 229 tracts. The shapes of the TAZs in this study are derived from the cartographic bound-  
 230 ary shapefiles developed by the U.S. Census Bureau in conjunction with the 2010 census  
 231 (<https://www2.census.gov/geo/tiger/TIGER2010/TAZ/2010/>). Considering the TAZs are  
 232 particularly useful for journey-to-work and place-of-work statistics, we employed them as the  
 233 basic units for accounting the taxi trips. Figure 5b shows OD flows on July 1. Red lines  
 234 indicate the dominant OD flows.

235 As a case study, this paper examined OD flows recorded on weekdays and weekends in  
 236 June 2014. The weekday dataset includes taxi trajectories collected on June 3, 10, 17, and 24,  
 237 and represents 1,721,655 taxi trips. The weekend dataset includes taxi trajectories collected  
 238 on June 8, 15, 22, and 29, and describes 1,593,480 trips.



**Figure 6** Outliers detection of OD flows: X-axis indicates forward OD flows counts and Y-axis indicates reverse OD flows counts.

### 239 3.2 Workflow

240 The performance of the proposed method was compared with alternative methods. Trajectory  
 241 anomaly detection based on Mahalanobis distance [20] was used to evaluate the performance  
 242 of outliers detection by the proposed method. The Mahalanobis distance is distinguished  
 243 from Euclidean distance by its consideration of the correlations of the data, in this case, the  
 244 two OD flow datasets. According to [20], the anomaly detection threshold can be defined as  
 245 follows:

$$246 \quad d_M(OD_{t_1}, \mu_{[t_0, t_1]}) \geq 3 \cdot \sqrt{\frac{1}{N} \sum_{t \in [t_0, t_1]} (OD_t - \mu_{[t_0, t_1]})^2}$$

247 where  $OD_{t_1}$  is the current OD flow, and  $\mu_{[t_0, t_1]}$  is the median of all OD flows during  $[t_0, t_1]$ .  
 248 In addition, we visualized the results in order to compare them and make the difference  
 249 easier to understand. The difference of scale was evaluated using standard statistics, such as  
 250 F-test, to compare the variance of two datasets.

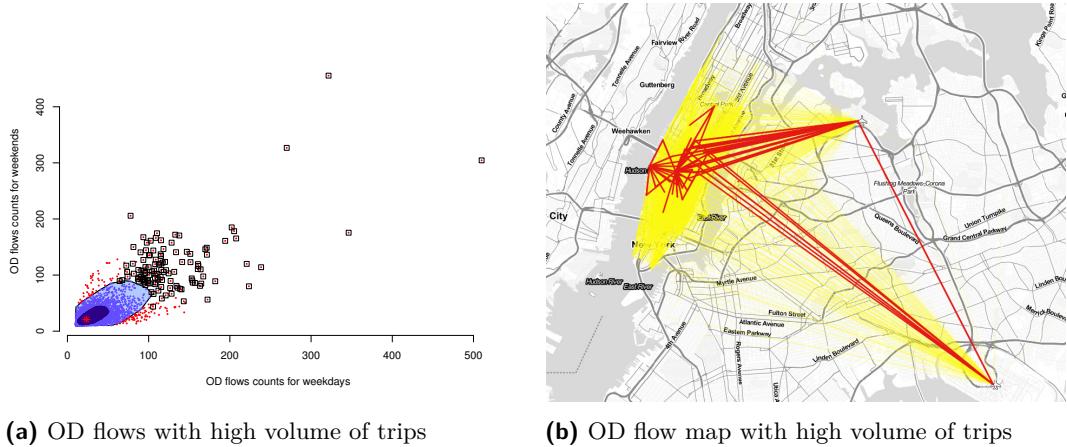
251 For data cleaning process, this study used Hadoop with Pig. We developed a Hadoop  
 252 program to resolve large data volume, which was composed of 173 million taxi trip records,  
 253 remove trips with invalid OD coordinates, and assign each OD locations into the corresponding  
 254 traffic analysis zone. To implement the OD flow outliers detection, this study used R. The  
 255 computing environment used Amazon Web Service and the Bridges supercomputer at the  
 256 Pittsburgh Supercomputing Center. This study only evaluated OD flows more than 10 trips,  
 257 as the low trip number OD flows could have distorted the view of the data. All the code will  
 258 be released as open source (the link to the code is available upon request).

### 259 3.3 Case study: weekdays vs weekends

#### 260 3.3.1 Outlier Detection

261 The bag plots presented OD flow outliers on weekdays and weekends in Figures 6a and 6b,  
 262 respectively. The outliers are detected by considering forward OD flows and reverse OD  
 263 flows together.

## 6:10 Outlier Detection and Comparison of Origin-Destination Flows using Data Depth



**Figure 7** Outliers with high volume of trips on weekdays and weekends: Rectangles in Figure 7a coincide with red lines in Figure 7b.

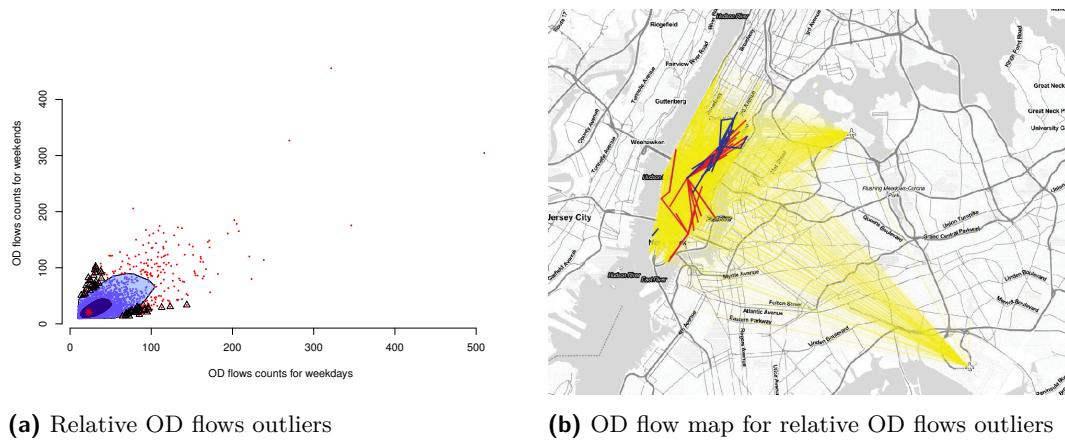
To find the difference between two datasets, we considered two forward OD flows together with the bag plot. Then, we identified the outliers OD flows in Figure 7a. The outliers with rectangle symbols indicate OD flows with large volumes of taxi trips during weekdays and weekends. Figure 7b depicts these outliers superimposed on a map with red lines. The yellow lines represent the other OD flows, excluding the large volume OD flows on weekdays and weekends. This case clearly demonstrates that most OD flows occurred in three broad areas: within Manhattan, between the center of Manhattan and the two major airports (J.F.K International Airport and LaGuardia Airport), and between the two airports.

In addition, we investigated abnormal weekend OD flows that are typical weekday OD flows. These abnormal weekend OD flows exhibited substantial variance in number of taxi trips relative to their weekday counterparts. Figure 8a presents these OD flows outliers with triangle symbols. In Figure 8b, red lines indicate the substantial increases in weekend trip volumes. Conversely, blue lines indicate the decreases in trip volumes. Figure 8b reveals that OD flows between the center of Manhattan and the two airports or between the two airports were not significantly different during weekdays and weekends. However, we did observe some meaningful decrease in OD flows during the weekends in business district, as depicted by the blue lines in Figure 8b.

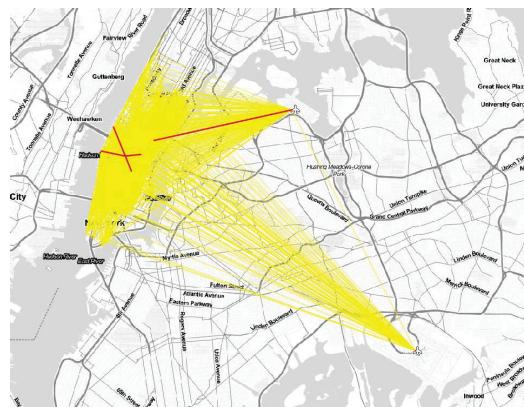
We also detected outlier OD flows using Mahalanobis distance. The results are presented in Figure 9. Far fewer outlier OD flows were detected using Mahalanobis distance than by our method. The Mahalanobis method only considers the forward OD flows of the two datasets. It identified OD flow outliers with high volume of trips because Mahalanobis distance considers the correlations between two OD flows. Thus, Mahalanobis distance is more likely to identify outliers when two OD flows have large trip volumes. In fact, the OD flows outliers from Mahalanobis distance are a subset of the outliers identified by our method, as depicted in Figure 7b. Furthermore, the Mahalanobis distance approach could not detect the outliers detected by our method in Figure 8 because the Mahalanobis distance approach cannot compare two flows to evaluate significant increases or decreases.

### 291 3.3.2 Scale Comparisons

292 We further investigated how two OD flows differ. Our approach is sensitive to the difference  
293 in scale. Hypothesis testing of the differences between two central regions in Figure 10



**Figure 8** Relative OD flows outliers on weekdays and weekends: Triangles in Figure 8a coincide with red and blue lines in Figure 8b.



**Figure 9** Outlier OD flows on weekdays and weekends based on Mahalanobis distance.

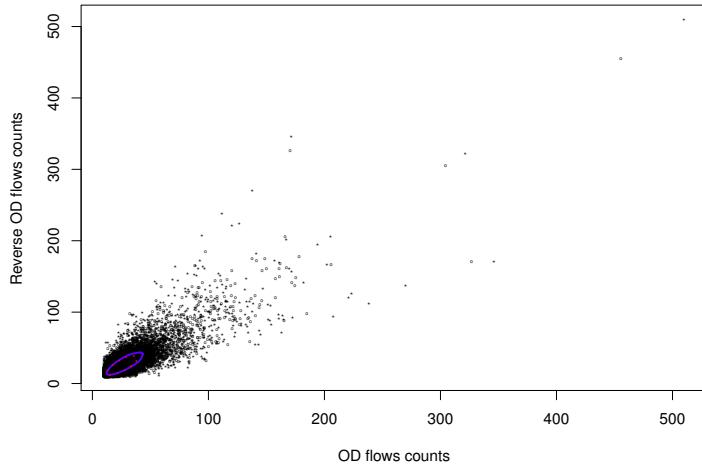
inadvertently revealed that the confidence interval was -0.0277 and 0.0157, which includes zero. Thus, it failed to reject the null hypothesis. The two central regions were similar in terms of the spread.

Interestingly, the standard statistic F-test was significant,  $F(9530, 7637) = 1.1786, p \leq 0.05$ . The variances of two groups were significantly different. The result of F-test directly opposed that of our method.

#### 4 Discussion

The results demonstrate that the method effectively identifies outlier OD flows based on data depth. It is also feasible to detect outlier OD flows by querying with conditional clauses, such as which outlier OD flows always have high trip volumes during time  $t_1$  and time  $t_2$ .

As an alternative, the state-of-the-art Mahalanobis distance approach detected similar outlier OD flows. However, the number of outliers detected was different. This occurred because the proposed method's OD flows data had heavy tail distributions, which means many of the OD flows with a long distance from the depth median depicted in Figure 8a. Mahalanobis distance is known to be inadequate when the underlying data have heavy tail distributions [27]. Thus, the presence of outliers may mask the detection of other outliers



**Figure 10** OD flows comparisons based on data depth:  $\circ$  indicates the OD flows on weekdays and  $*$  indicates the OD flows on weekends; blue line presents the central region of the OD flows for the weekdays and red dotted line presents the central region of the OD flows on weekends.

310 in Mahalanobis distance approach. Furthermore, it can only detect OD flow outliers with  
 311 high numbers of trips during time  $t_1$  and time  $t_2$ . It is difficult to detect OD flows outliers  
 312 that have different properties, such as substantial differences in the number of trips when  
 313 comparing between time  $t_1$  and time  $t_2$ .

314 In terms of the difference in spread, our method used a bootstrap technique to compare  
 315 the central regions of data depth. This technique investigated the difference in scale as well  
 316 as the structure of data. It can provide information about how deeply points from group 1,  
 317 OD flows at  $t_1$ , tend to be located within group 2, OD flows at  $t_2$ . General statistics such as  
 318 F-test only provide their difference in variation and do not further specify how groups differ.

319 Interestingly, the F-test results revealed a statistically significant difference in terms of  
 320 variation of OD flows on weekdays and weekends. Our approach showed no statistically  
 321 significant differences. This contrast may be caused by the sensitivity of F-test to non-  
 322 normality [6], which increases the Type-I error rate. Conversely, data depth makes no  
 323 assumptions about the distributions of the underlying dataset.

## 324 5 Conclusions and Future Work

325 This paper describes a new method for identifying outlier OD flows and the difference in scale  
 326 between two different OD flows at  $t_1$  and  $t_2$ . The new method is based on the concept of data  
 327 depth. Data depth is robust statistics, which is suitable to non-Gaussian distribution of the  
 328 underlying datasets. Compared with standard statistics, our method enhances understanding  
 329 of the differences and the magnitude of the differences between two OD flow datasets.

330 This study made no attempt to incorporate geographic contexts such as locational  
 331 circumstances or surrounding environment in understanding OD flows. Ultimately, further  
 332 research should focus on integrating the analysis of OD flows with appropriate geographic  
 333 contexts. Such research will lead to desirable knowledge discovery and better understanding  
 334 of movement dynamics.

335 

---

 References

- 336 1 Marina Alberti, John M Marzluff, Eric Shulenberger, Gordon Bradley, Clare Ryan, and  
337 Craig Zumbrunnen. Integrating humans into ecology: Opportunities and challenges for  
338 studying urban ecosystems. *AIBS Bulletin*, 53(12):1169–1179, 2003.
- 339 2 Maike Buchin, Somayeh Dodge, and Bettina Speckmann. Similarity of trajectories taking  
340 into account geographic context. *Journal of Spatial Information Science*, 2014(9):101–124,  
341 2014.
- 342 3 Chao Chen, Daqing Zhang, Zhi-Hua Zhou, Nan Li, Tülin Atmaca, and Shijian Li. B-  
343 planner: Night bus route planning using large-scale taxi GPS traces. In *2013 IEEE Interna-*  
344 *tional Conference on Pervasive Computing and Communications (PerCom)*, pages  
345 225–233. IEEE, 2013.
- 346 4 Srinivas Devarakonda, Parveen Sevusu, Hongzhang Liu, Ruilin Liu, Liviu Iftode, and Badri  
347 Nath. Real-time air quality monitoring through mobile sensing in metropolitan areas. In  
348 *Proc. 2nd ACM SIGKDD International Workshop on Urban Computing*, page 15. ACM,  
349 2013.
- 350 5 Matt Duckham, Marc van Kreveld, Ross Purves, Bettina Speckmann, Yaguang Tao, Kevin  
351 Verbeek, and Jo Wood. Modeling checkpoint-based movement with the earth mover’s  
352 distance. In *International Conference on Geographic Information Science*, pages 225–239.  
353 Springer, 2016.
- 354 6 Andy Field, Jeremy Miles, and Zoë Field. *Discovering statistics using R*. Sage, London,  
355 UK, 2012.
- 356 7 Vitor Cunha Fontes, Lucas Andre de Alencar, Chiara Renso, and Vania Bogorny. Discov-  
357 ering trajectory outliers between regions of interest. In *Proc. XIV GeoInfo*, pages 49–60,  
358 2013.
- 359 8 Yizhao Gao, Ting Li, Shaowen Wang, Myeong-Hun Jeong, and Kiumars Soltani. A multi-  
360 dimensional spatial scan statistics approach to movement pattern comparison. *International*  
361 *Journal of Geographical Information Science*, 0(0):1–22, 2018.
- 362 9 Diansheng Guo and Xi Zhu. Origin-destination flow data smoothing and mapping. *IEEE*  
363 *Transactions on Visualization and Computer Graphics*, 20(12):2043–2052, 2014.
- 364 10 Myeong-Hun Jeong, Yaping Cai, Clair J Sullivan, and Shaowen Wang. Data depth based  
365 clustering analysis. In *Proc. 24th ACM SIGSPATIAL International Conference on Ad-*  
366 *vances in Geographic Information Systems*, page 29. ACM, 2016.
- 367 11 Mei-Po Kwan. Space-time and integral measures of individual accessibility: A comparative  
368 analysis using a point-based framework. *Geographical Analysis*, 30(3):191–216, 1998.
- 369 12 Tatjana Lange, Karl Mosler, and Pavlo Mozharovskyi. Fast nonparametric classification  
370 based on data depth. *Statistical Papers*, 55(1):49–69, 2014.
- 371 13 Jae-Gil Lee, Jiawei Han, and Xiaolei Li. Trajectory outlier detection: A partition-and-  
372 detect framework. In *IEEE 24th International Conference on Data Engineering*, pages  
373 140–149. IEEE, 2008.
- 374 14 Liangxu Liu, Shaojie Qiao, Yongping Zhang, and JinSong Hu. An efficient outlying trajec-  
375 tories mining approach based on relative distance. *International Journal of Geographical*  
376 *Information Science*, 26(10):1789–1810, 2012.
- 377 15 Regina Y Liu. On a notion of data depth based on random simplices. *The Annals of*  
378 *Statistics*, pages 405–414, 1990.
- 379 16 Regina Y Liu and Kesar Singh. A quality index based on data depth and multivariate rank  
380 tests. *Journal of the American Statistical Association*, 88(421):252–260, 1993.
- 381 17 Jean Damascène Mazimpaka and Sabine Timpf. Exploring the potential of combining  
382 taxi GPS and flickr data for discovering functional regions. In *AGILE 2015*, pages 3–18.  
383 Springer, 2015.

- 384   **18** Jean Damascène Mazimpaka and Sabine Timpf. Trajectory data mining: A review of  
385   methods and applications. *Journal of Spatial Information Science*, 2016(13):61–99, 2016.
- 386   **19** Karl Mosler. *Robustness and Complex Data Structures*, chapter Depth Statistics, pages  
387   17–34. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.
- 388   **20** Bei Pan, Yu Zheng, David Wilkie, and Cyrus Shahabi. Crowd sensing of traffic anomalies  
389   based on human mobility and social media. In *Proc. 21st ACM SIGSPATIAL International  
390   Conference on Advances in Geographic Information Systems*, pages 344–353. ACM, 2013.
- 391   **21** Peter J Rousseeuw and Ida Ruts. Algorithm AS 307: Bivariate location depth. *Journal of  
392   the Royal Statistical Society. Series C (Applied Statistics)*, 45(4):516–526, 1996.
- 393   **22** Peter J Rousseeuw, Ida Ruts, and John W Tukey. The bagplot: A bivariate boxplot. *The  
394   American Statistician*, 53(4):382–387, 1999.
- 395   **23** Ran Tao and Jean-Claude Thill. Spatial cluster detection in spatial flow data. *Geographical  
396   Analysis*, 48(4):355–372, 2016.
- 397   **24** John W Tukey. Mathematics and the picturing of data. In *Proc. International Congress  
398   of Mathematicians*, volume 2, pages 523–531, 1975.
- 399   **25** Rand R Wilcox. Approximating Tukey’s depth. *Communications in Statistics-Simulation  
400   and Computation*, 32(4):977–985, 2003.
- 401   **26** Rand R Wilcox. Two-sample, bivariate hypothesis testing methods based on Tukey’s depth.  
402   *Multivariate Behavioral Research*, 38(2):225–246, 2003.
- 403   **27** Rand R Wilcox. *Introduction to robust estimation and hypothesis testing*. Academic Press,  
404   2012.
- 405   **28** Hans Peter Wolf and Uni Bielefeld. aplpack: Another Plot PACKage: stem.leaf, bagplot,  
406   faces, spin3r, plotssummary, plot hulls, and some slider functions, 2014. R package version  
407   1.3.0. URL: <https://CRAN.R-project.org/package=ap1pack>.
- 408   **29** Junjun Yin, Yizhao Gao, Zhenhong Du, and Shaowen Wang. Exploring multi-scale spati-  
409   otemporal twitter user mobility patterns with a visual-analytics approach. *ISPRS Interna-  
410   tional Journal of Geo-Information*, 5(10):187, 2016.
- 411   **30** Junjun Yin, Aiman Soliman, Dandong Yin, and Shaowen Wang. Depicting urban bound-  
412   aries from a mobility network of spatial interactions: A case study of great britain  
413   with geo-located twitter data. *International Journal of Geographical Information Science*,  
414   31(7):1293–1313, 2017.
- 415   **31** Guan Yuan, Shixiong Xia, Lei Zhang, Yong Zhou, and Cheng Ji. Trajectory outlier detec-  
416   tion algorithm based on structural features. *Journal of Computational Information Systems*,  
417   7(11):4137–4144, 2011.
- 418   **32** Daqing Zhang, Nan Li, Zhi-Hua Zhou, Chao Chen, Lin Sun, and Shijian Li. iBAT: Detecting  
419   anomalous taxi trajectories from GPS traces. In *Proc. 13th International Conference on  
420   Ubiquitous Computing*, pages 99–108. ACM, 2011.
- 421   **33** Yu Zheng. Trajectory data mining: An overview. *ACM Transactions on Intelligent Systems  
422   and Technology*, 6(3):29, 2015.
- 423   **34** Yijun Zuo and Robert Serfling. General notions of statistical depth functions. *The Annals  
424   of Statistics*, 28:461–482, 2000.