# Visual Analysis of Public Utility Service Problems in a Metropolis



Jiawan Zhang, *Member, IEEE*, Yanli E, Jing Ma, Yahui Zhao, Binghan Xu, Liting Sun, Jinyan Chen, and Xiaoru Yuan, *Member, IEEE* 

Fig. 1. Visual Analytics Approach to Public Utility Service Problems. a)Dotmap for visualizing the spatial distribution of issues; b)Heatmap for visualizing the spatial distribution of issues; c)Visualization of one event; d)Visualization of events; e)Changes in spatial distributions over months; f)Multiple temporal scale analysis; g)Animation for visualizing time sequence of issues in one event.

**Abstract**—Issues about city utility services reported by citizens can provide unprecedented insights into the various aspects of such services. Analysis of these issues can improve living quality through evidence-based decision making. However, these issues are complex, because of the involvement of spatial and temporal components, in addition to having multi-dimensional and multivariate natures. Consequently, exploring utility service problems and creating visual representations are difficult. To analyze these issues, we propose a visual analytics process based on the main tasks of utility service management. We also propose an aggregate method that transforms numerous issues into legible events and provide visualizations for events. In addition, we provide a set of tools and interaction techniques to explore such issues. Our approach enables administrators to make more informed decisions.

Index Terms—utility services, evidence-based decision making, visual analytics, aggregate

## **1** INTRODUCTION

City utility service management has become an emergent topic in the era of rapid city expansion. Normal city operations, particularly those in a large-scale city, are guaranteed through the provision of proper utility services by multiple city agencies or non-profit organizations,

- Jiawan Zhang is with School of Computer Science and Technology, and School of Computer Software(SCS), Tianjin University. E-mail: jwzhang@tju.edu.cn.
- Yanli E is with School of Computer Science and Technology, Tianjin University. E-mail: yanlitju@163.com.
- Jing Ma is with SCS, Tianjin University. E-mail: majingtju@163.com.
- Yahui Zhao is with SCS, Tianjin University. E-mail: zhaoyahui@tju.edu.cn.
- Binghan Xu is with SCS, Tianjin University. E-mail: iceuia@hotmail.com.
- Liting Sun is with SCS, Tianjin University. E-mail: slttju2010@sina.cn.
- Jinyan Chen is with SCS, Tianjin University. E-mail: chenjinyan@tju.edu.cn.
- Xiaoru Yuan is with School of EECS, Peking University. E-mail:xiaoru.yuan@pku.edu.cn.

Manuscript received 31 Mar. 2014; accepted 1 Aug. 2014. Date of publication 11 Aug. 2014; date of current version 9 Nov. 2014. For information on obtaining reprints of this article, please send e-mail to: tvcg@computer.org.

Digital Object Identifier 10.1109/TVCG.2014.2346898

such as the natural gas supply, water supply, drainage, roads or bridges and heat supply divisions. These divisions guarantee the quality of the daily lives of citizens by providing crucial public infrastructures. However, ensuring the orderliness and safety of a city at all times can be challenging, because the various daily issues or emergencies related to utility services arise due to the increasing rate of city construction, vandalism and infrastructure failures. Moreover, city governments have their internal demands to grasp whether utility services meet citizen requirements, as well as identify and solve existing problems instantly. An equally important factor is the surveillance of utility services because normal city operations are closely related to democratic self-government [24, 28, 30]. Inhabitants are the main service subjects and the issues they reported can reflect the city utility service problems in a more straightforward manner. Administrators must thus place importance on these issues.

Therefore, a mechanism that enables issues to be properly and effectively received, recorded, documented, reflected to related divisions, solved and followed up must be created. This mechanism should also be capable of archiving the handling process for each issue. To realize the mechanism, a hierarchical system has to be established. This system should include city-level centers and sub-centers, which are set up by the subordinate local districts. Such systems have been set up in several countries by providing a single non-emergency number, such as 311 in America [26], 101 in England, 115 in Germany,

1077-2626 © 2014 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information. and 12319 in China. They have integrated the previous help lines of multiple divisions into a single call center. These centers enable residents to contribute in addressing the issues that they witnessed in their own neighborhoods. During emergencies, these centers can also interact with emergency call centers, such as the fire department and police department, to address issues promptly. City administrators can improve service delivery and service request traceability, as well as citizen satisfaction with utility services through these centers. Thus, non-emergency call centers have been attached more importance by governments all over the world, although they vary from country to country because of different business management approaches.

In order to present our research more clearly, it is necessary to define some terms used in our paper. We cooperated with and used data from a non-emergency municipal service call center in China. This call center contains information from nine agencies or non-profit organizations, which includes the water supply, heat supply, natural gas supply, drainage, roads or bridges, bus, metro, taxi and light rail divisions. When citizens get into troubles due to utility service failures, they will call the center and report the issues they encountered. However, many issues reflect the same event. For example, a water or gas pipe burst event will cause many citizens to report the issues. Thus, one event in our study includes all the correlative issues reported by citizens. Moreover, some issues have a relatively low timeliness requirement, such as trivial daily life and common utility service management issues. Nevertheless, several issues may directly result in casualties and must thus be solved immediately. These issues include pipeline bursts, gas leaks, and stolen or damaged manhole covers. We utilize issue manifestation to distinguish the difference among these issues. On the one hand, issues vary from division to division. For instance, the natural gas, water and heat supply divisions encounter pipeline leak issues, but other divisions do not. On the other hand, every division has various issue manifestations.

Our research context is visual analysis of the issues reported by citizens to identify the utility service problems for administrators. To the best of our knowledge, such analysis has not yet been conducted. In addition, a method that improves the quality of utility service through evidence-based decisions is yet to be proposed, although a number of decision-making studies exist in the field of visual analytics [20, 38]. However, given the increasing number of issues have been captured and are now available, visual analysis and understanding of the issues have become challenging. The data are high-dimensional because of different divisions and are multivariate because of various issue manifestations. In addition, the data are also temporal-spatial because these issues differ in either spatial or temporal nature. The analysis of these issues can be difficult because of the finite feasibilities available in analyzing the location, time and category simultaneously. Moreover, city administrators need to analyze utility service problems from different perspectives because their mission requirements vary in the role of management. For instance, decision-makers are more concerned with the events leading to anomalies rather than every discrete issue. City analysts focus on various trends and patterns of issues. Personnel in the operation and maintenance departments focus on the locations of infrastructures that need to be maintained. Therefore, we provide tools and methodological support for administrators.

We utilize task-driven visual analytics approach [27] based on the city utility service management tasks [29], which include visually analyzing issues, finding and understanding anomalies and temporal-spatial patterns, identifying events, and performing evidence-based decision making [8, 12, 23].

The main contributions of this research are summarized as follows: 1) Proposed a task-driven visual analytics process for city adminis-

trators to identify and solve city utility service problems.

2) Proposed a method that can identify events by aggregating numerous issues into legible events and provided visualizations to represent these events.

3) Developed a set of visual analytics tools and interaction techniques that will enable city administrators to explore issue trends and patterns from different perspectives.

The rest of this paper is structured as follows. The previous studies

are reviewed in the second section. The proposed method is presented in the third section. The data and experimental results are discussed in the fourth section and finally we draw our conclusions.

#### 2 RELATED WORK

Previous studies that are relevant to our research are divided into three categories, namely, study on city utility services, city utility service study in the information technology context, and study for city life improvement in a visual analytics context. In this section, we briefly summarize the studies in each category.

#### 2.1 Study on City Utility Services

Extensive studies have been conducted on utility services. However, the majority of these studies focused on infrastructures. Some of these studies were conducted for infrastructure management. [1] provided an optimization model for sewer network management. [32] presented an infrastructure interdependency analysis. Moreover, some studies were conducted for infrastructure deployment optimization. Shukla [36] proposed an optimization framework that utilized mathematical programming to determine the best location for establishing alternative transportation fuel stations. Agustin-Blas [4] addressed the issue of deploying metropolitan wireless networks by using a hybrid grouping genetic algorithm. A few studies were conducted for infrastructure failure detection. Oliveira [11] presented a novel approach that extended the existing spatial scan statistics approach to detect and locate clusters of break points in a water distribution network. [41] reported a novel data mining method based on the requirements of end users in water utilities to predict pipe failures.

A considerable number of studies were conducted to investigate utility services with different approaches. Tanimoto [37] proposed a new methodology to calculate likely utility load profiles in a dwelling. Campisi-Pinto [7] addressed the problems of forecasting urban water demand by using back-propagation artificial neural network(ANN) coupled with wavelet-denoising. Adamowski [2] demonstrated that the ANN approach provided better prediction of peak daily summer water demand than multiple linear regression and time series analysis. Many researchers [15, 22, 18] focused on an electric load forecast model and system based on ANN and several other methods.

## 2.2 City Utility Service Study in the Information Technology Context

A number of studies have adopted information technology in the study of city utility services. Some studies were designed to help administrators in the emergency management of unexpected infrastructure incidents. [14] presented a framework to assess risks and uncertainties in the domain of utility services for up to 100 years in the future. Armbruster, etc [6] illustrated that the threats associated with aging infrastructures had been continuously increasing in intensity and were at par with catastrophic events encountered in the context of critical infrastructure protection.

Some studies focused on infrastructures through the Geographic Information System(GIS). Vanier [39] provided an overview of the stateof-practice of GIS to manage and optimize municipal infrastructures. [21] presented a web-based GIS for the assessment and visualization of critical infrastructure and its hazards. [34] proposed an integrated infrastructure management information system based on GIS. [13] predicted the location of future emergency events by combining GIS and neural networks.

However, all the previous studies for city utility services were based on infrastructure network operating data, which had no necessary connection with utility service problems. On the one hand, most of these data were based on the entire network. Data granularity is usually insufficient to reflect every network terminal, such as every family, company or school. That is, normal network operation does not imply that no issues exist. But each reported issue indicates that the network has some problems. On the other hand, no existing channel can cause a single-issue as feedback to adjust the infrastructure network operating parameters, and the current channel can only work when massive issues demand adjustment.

## 2.3 Study for City Life Improvement in a Visual Analytics Context

An increasing number of studies that aim to improve city life have focused on the visual analytics and visualization domain. Bristle maps [19] can be used to visualize spatio-temporal data, such as city burglary rates. [40] established a framework to analyze provenance in human terrain and meet the requirements of the intelligence community. [9] visualized spatio-temporal social media data to detect abnormal events and their severity. In addition, a considerable number of studies [12, 23, 42] visually analyzed spatio-temporal urban data to process city taxi trips and conducted traffic jam analysis. [33] designed AlVis for the surveillance of road tunnels to recognize complex disaster scenarios in advance. [43] employed GIS and KDE to analyze the spatial distribution of public housing households. [16] used GIS to locate and rank suitable storm-water harvesting sites in urban areas. [17] developed a WebGIS system function that provided real-time information on road conditions. [35] presented several visualizations of integrated data on local municipalities for the use of municipality managers. Collins, etc [10] introduced Bubble Sets, which refer to a spatial layout that minimizes cluster overlap, and is used to visualize some city problems.

A large number of overlaps among points and intersections among lines in the geo-spatial view could occur and make the visualization illegible. Thus, an appropriate aggregation method must be applied to solve this problem. A number of studies have proposed aggregation methods and conducted visual analysis of aggregation results by highlighting significant information. [25] proposed a predictive analytics approach by geographically visualizing an aggregated distribution to enable analysts to establish prevention measures corresponding to the perceived threats. [3] suggested a method for spatial generalization and aggregation of massive movement data. This method aggregated numerous trajectories into legible flow maps. [5] defined spatiotemporal aggregation methods suitable for movement data. However, these aggregation methods only considered spatial and temporal factors and were thus unsuitable for the aggregation of citizen-reported issues. Thus, in this study, we propose a new aggregate method that considers three aspects, namely, spatial, temporal and semantic.

#### **3 OUR METHOD**

## 3.1 Visual Analytics Process Overview

All city utility service administrators are considered as the main users in this study. The main tasks of administrators are analyzing the issues and determining all the events that are caused by utility services promptly. However, the routine work-flow is insufficient to satisfy administrator requirements gradually because utility service management is becoming increasingly complex, which is why we conduct this study. Based on this motivation, both the requirement and design methods of our visual analytics system are clarified through user discussion. According to user opinions, we not only retained the previous work-flow that they have been used to, but also addressed its drawbacks. To guarantee the usability and adoptability of our developed system and to enable users to fulfill their tasks more efficiently, we interacted with the users from the earliest design stages through the intermediate and final evaluations.

Our proposed process is an iterative work-flow based on the main tasks, as described in Fig.2. City administrators need to know what happened to citizens, as well as when, where and why these issues occurred. Therefore, they must understand the temporal-spatial patterns and keep track of the handling information for these issues. Data mining methods are necessary in determining significant knowledge, such as the relation between issues and notifications released ahead of schedule by the concerned divisions. For the full use of knowledge, we build a data model of them. The visualization based on data model can assist administrators in discovering more patterns. Moreover, given that the quantity of events is far smaller than the quantity of issues and that events better demonstrate the reasons behind anomalies than discrete issues, the visual analytics of events would be more significant to administrators. The administrators should understand the events by determining their severity, influence area and evolution trend. Thus, aggregating correlative issues into one event is essential. If the available analysis results of events are sufficient to support administrators in making decisions, the concerned division can take action to promote utility service quality, such as improving infrastructures and solving emergencies. Otherwise, after the administrators reviewed the tasks and current analysis situation, they had to analyze from issues to events again.

We developed a method and a set of tools for the administrators or related divisions to understand city utility service problems better from different perspectives and hierarchies. First, the common visualization tool is essential to administrators. We visualized the distribution of issues by utilizing heatmap and dotmap on the basis of their geographical coordinates. Administrators can employ zoom interaction techniques to analyze spatial information at different scales(Section 3.2.1). We provided a time-filtered line chart to assist users in identifying temporal patterns or anomalies. For ease of analysis, the line charts provide multiple time scale support and enable the comparison of spatial distributions with linked views (Section 3.2.2). City administrators often focus on the efficiency in solving issues. We thus provided a bar chart that presents changes in resolved and unresolved issue quantities over time. We combined dotmap and heatmap to show solution situation in two-dimensional geographical space(Section 3.2.3). Second, we also proposed a method that could identify events by aggregating correlative issues based on a clustering algorithm(Section 3.3.1). In addition, we presented visualization design and visual analytics tools to assist administrators in analyzing events(Section 3.3.2). Incorporating other interactive functionalities to drill down further from issues to events according to user requirements is also essential.



Fig. 2. Proposed Analytics Process of City Utility Service Problems.

#### 3.2 Visual Analytics of Issues

## 3.2.1 Geo-spatial Visual Analytics

Issues in different locations are of different quantities and manifestations, so the spatial distribution of issues should be understood by city administrators. We first utilized a dotmap based on Scalable Vector Graphics map from OpenStreetMap to display all the issues, each of which is represented by a dot. Different issue manifestations in one division are distinguished by different colors. Moreover, we utilized the blue noise sampling algorithm [31] to avoid overlaps between dots. Given the possibility that links exist among different issue manifestations, one manifestation might give rise to another manifestation. Thus, when one manifestation is spatially surrounded consistently by another manifestation, a certain link exists between the two manifestations(Fig.1(a)). We then employed the heatmap technique to show the quantity of issues from a macroscopic perspective. The higher the issue quantity is, the closer to red the color is, and vice versa. The spatial distribution of issues is clearly visualized in Fig.1(b).

## 3.2.2 Temporal Visual Analytics

The number of issues produced by infrastructure failures shows patterns over time. These patterns include periodical decline or increase



Fig. 3. Temporal and Spatial Analysis of Multiple Divisions.



Fig. 4. a)Temporal comparison of the same periods in 2011 and 2012 for the heat supply division; b)Spatial comparison of November in 2011 and 2012 for the heat supply division.

in issue quantity and the anomalies that numerous issues suddenly occurred. Furthermore, infrastructures often exhibit an aging trend because of internal and external factors. Thus, city administrators should analyze the temporal patterns of issues.

Traditional statistical approaches cannot adequately support the analysis of characteristic values. Therefore, we combined line charts and dotmap or heatmap to provide administrators with a multi-timescale analysis support, including month, day and hour. First, multiple divisions or multiple issue manifestations in one division can be analyzed on the basis of lines with different colors and heatmaps from temporal and spatial perspectives. For example, five divisions are shown in Fig.3. The number of issues from heat and water supply is relatively large, whereas the spatial distributions of the other three divisions are relatively sparse. Second, it can display the maximum and minimum quantities of previous years as two grey lines, and the quantity in same period of last year(Fig.4(a)) or all the other years with different color lines. This procedure guides the administrators in identifying anomalies. In addition, administrators can compare spatial distributions in the same period over multiple years(Fig.4(b)) or in the twelve months of one year with linked view(Fig.1(e)). Furthermore, characteristic or abnormal values can be analyzed on any time scale(Fig.1(f)).

Given that weather is closely related to day-to-day life, administrators should consider the study of weather. We provided a more detailed time-scale line chart to determine the relationship between issue quantity and rainfall. We divided 24h into two parts, namely, 8:00 to 20:00 and 20:00 to 8:00(the next day). In the line chart, one color represents the issue quantity and the other color represents rainfall. The horizontal axis represents the time series in which one day is divided into two parts. The two vertical axes represent the issue quantity and rainfall respectively. This visualization can assist administrators in identifying not only temporal patterns of issue quantity and rainfall but also their relationship(Fig.10(a) and Fig.10(c)).

## 3.2.3 Visual Analytics of Efficiency

To guarantee citizen life quality, administrators aim to solve more issues by inspecting the efficiency of every division. Thus, resolved and unresolved issues should be monitored. From a spatial perspective, all the discrete issues on the map are shown by using dotmap. The problem is serious in the areas where the dots are denser. Meanwhile, we also utilized heatmap technique on the basis of a geographical map to visualize the spatial distribution of unresolved issues. The closer to red the color is, the worse the solution situation of the corresponding location is, and vice versa. Thus, administrators can obtain the spatial distribution of solution situation(Fig.5(a)). From a temporal perspective, we used a bar chart in which the light color represents the quantity of resolved issues and the dark color represents the quantity of unresolved issues. The horizontal axis represents time, whereas the vertical axis represents quantity(Fig.5(b)). A higher bar indicates a larger quantity, and vice versa. In addition, we used different colors to distinguish different divisions to help administrators identify which has the worst efficiency.



Fig. 5. a)Spatial distribution of solution situation; b)Temporal trend of resolved and unresolved issues.



Fig. 6. Process for Identifying Events. a)Discrete issues; b)Clustering result without human interaction; c)Process of human interaction; d)Clustering result after human interaction.

## 3.3 Visual Analytics of Events

## 3.3.1 Event Identification

Decision-makers need to focus on events to make more informed decisions. Thus, we proposed a method that could identify events by using a clustering algorithm to aggregate correlative issues to an event.

Every issue has occurrence time, spatial location including text description(semantic space) and geographical coordinate(absolute space), and detailed issue content. These attributes describe issues from different perspectives and must be considered in the proposed algorithm. We communicated with the city utility service administrators to obtain additional information. They mentioned that three aspects must be included in our algorithm. First, the time span of issues caused by the same event should be five days at most. Second, the locations of correlative issues should be spatially close. However, "close" involves two aspects. On the one hand, their locations are close when their absolute space distance is less than 1 km. On the other hand, the locations of issues that occur in the same neighborhood or street are close even if the absolute space distance is greater than 1 km(the longest distance of two locations in the same neighborhood). Therefore, both the text description and geographical coordinate of the location are essential. Third, the issue content describes what happens to citizens in detail. We summarized the keywords that often appeared in the content. Given that many different keywords express the same meaning, we divided them into several categories as  $O_i$ , i = 1, 2, ..., n, that is, one category may contain many keywords. Thus,  $O_i = \{P_{ij}\},\$  $j = 1, 2, ..., m, P_{ij}$  represents a keyword in category  $O_i$ . Because some citizens report reasons and some report results of one event,  $O_i$  is either a reason or result. The categories and keywords vary from division to division.

In our algorithm, the correlation is used to determine whether two issues belong to the same event. We considered four factors comprehensively in our clustering algorithm to calculate the correlation: time span, absolute space distance, semantic space relevance and issue content similarity. When the time span of two issues exceeded five days, we defined such issues as uncorrelated. Absolute space is two-dimensional. We thus calculated the actual geographical distance between two issues by using longitude and latitude. When the distance was greater than 1 km, we defined the issues as uncorrelated. With regard to semantic space, we compared the text descriptions of the spatial locations of the two issues. If two issues occurred in the same neighbourhood or street, we defined them as relevant and did not consider their absolute space distance; otherwise, we did not consider their semantic space relevance. Issue content similarity is complicated. We thus quantified every issue content to a high dimension vector. The dimension of the vector was equal to the number of categories summarized above. The particular steps are as follows. If the content contains any keyword in O<sub>i</sub>, the *i*th dimension is 1; otherwise, the *i*th dimension is 0.1. If all dimensions representing reason are 0.1 and at least one dimension representing result is 1, we provided an empirical value for every reason and the value was decided by the possibility that it would lead to this result. The larger the empirical value is, the higher the priority of the reason is. However, the sum of all empirical values in one issue should be 1. Thus, issue content becomes a vector. We quantified content similarity of two issues by calculating the Euclidean distance of two vectors.

The calculations of the above four factors have an equal contribution to the correlation of every two issues. Thus, we should normalize the calculations linearly with Formula (1), except for semantic space relevance, where  $D_{ii}$  represents the calculation between the *i*th and *j*th issues;  $D_{max}$  and  $D_{min}$  represent the maximum and minimum of all calculations of every two issues, respectively; and  $d_{ij}$  is the normalized result. However, several exceptions should be specifically noted.  $t_{ij}$ ,  $a_{ij}$  and  $c_{ij}$  respectively represent normalized results of time span, absolute space distance, and issue content similarity between the *i*th and *j*th issues, which are calculated by using Formula (1). T is calculated with Formula (2), which represents critical value of time span. A is calculated with Formula (3), which represents critical value of absolute space distance. If  $t_{ij}$  is less than T,  $t_{ij} = 0.1$  (the two issues are uncorrelated). Meanwhile, if  $a_{ij}$  is less than A,  $a_{ij} = 0.1$ . For semantic space relevance  $s_{ij}$  and critical value S, when two issues occur in the same community or street,  $s_{ii} = A$  and  $a_{ii} = 1$ , S = A, A = 1.

Otherwise,  $s_{ij} = 1$  and S = 1. C represents the critical value of issue content similarity and varies from division to division. In our method, C is assigned with the second-largest value of all the  $c_{ii}$ , which is the optimum value for determining whether two issues have similar contents. Then, F = T \* A \* S \* C; F is the critical value of correlation. The correlation  $f_{ij}$  between the *i*th and *j*th issues is calculated with Formula (4). As for the *i*th issue, if  $f_{ij}$  (j = 1, 2, ..., n) is greater than F, the *j*th and *i*th issues are correlative. All issues correlative to the *i*th issue are aggregated into one event. Thus, discrete issues(Fig.6(a)) are aggregated into different events(clusters) marked by rectangles with different colors representing different event types(Fig.6(b)).

$$d_{ij} = \frac{D_{max} - D_{ij}}{D_{max} - D_{min}} \tag{1}$$

$$T = \frac{T_{max} - 5}{T_{max} - T_{min}} \tag{2}$$

$$A = \frac{A_{max} - 1}{A_{max} - A_{min}} \tag{3}$$

$$f_{ij} = t_{ij} a_{ij} s_{ij} c_{ij} \tag{4}$$

In our method, one issue may belong to multiple events, because the keywords we extracted from issue content may contain more than one reason. We considered every reason that appeared in the issue content according the priority of reasons. The quantity of events that one issue can belong to is decided by administrators.

However, no algorithm can deal with issue content more accurately than humans. Every method that combines temporal-spatial and semantic variables is often accompanied by errors, such as the issues marked by darkgoldenrod circles in the shaded rectangle in Fig.6(b). Our algorithm introduced human interaction. The clustering result of the first iteration is verified by eliminating or affirming several issues from the clusters. The interaction in our algorithm then continues to execute according to the following steps(Fig.6(c)). First, we analyzed the content of issues that need to be removed and found all of them to be described irregularly. When users removed issues, they must choose the right event type that can generalize issue content. Second, the above method determines issues that have similar content to the removed issues. These issues are distinguished with red font and should be manually inspected. Next, the removed issues are added to other clusters or become a lone cluster according to the event type chosen by users. Our algorithm keeps iterating until the result has met the user requirements. The clustering result after iteration improves the accuracy of our algorithm(the shaded rectangle in Fig.6(d)).

## 3.3.2 Event Visualization

On the basis of the analysis results of a large number of experiments, we concluded that issues in one event often presented linearly along infrastructure networks of their divisions. Thus, to visualize events better, the visualization of one event must include all issues in the event and must not change their original locations on the one hand. On the other hand, overlaps among events must be minimized because of numerous events in the same time quantum and space. Thus, we utilized a visualization based on Minimal Spanning Tree(MST) algorithm(kruskal) to link all points representing issues in one event. In this algorithm, we took the spatial distance between two issues as the weight of the edge. Moreover, after obtaining the MST, we utilized lines to draw this MST, with the width of the lines being slightly greater than the diameter of the point(Fig.7(a)). However, this measure made MST visually discontinuous due to intersections among lines. In order to avoid intersections, we proposed an algorithm shown in Algorithm 1, which draws all the edges of MST with a smooth stroke and guarantees that the length of the stroke is the shortest(Fig.7(b)). However, we can not identify different events by using opaque strokes, as shown in the top figure of Fig.7(c). Thus, we utilized transparent strokes to identify overlaps among events. For example, three events were identified in the bottom figure of Fig.7(c).

## Algorithm 1 Drawing the MST with a stroke

Input: The set of edges of MST,  $S = [s_0, s_1, s_2...]$ , for each  $s_k \in S$ ,  $s_k =$ [pointi, pointj], edgeflag[k] = false, k = 0, 1, 2...;Initialize a set of points for stroke,  $L = [s_0[0], s_0[1]]$  and edgeflag[0] = true;**Output:** The stroke for drawing the MST, L; repeat for each  $s_k \in S$  and edgeflag[k] = false do if only *i* in *L* then insert *j* after the first location *i* in *L*; insert *i* after *j*; edgeflag[k] = true;end if if only *j* in *L* then insert *i* after the first location *j* in *L*; insert *j* after *i*; edgeflag[k] = true;end if if *i* and *j* in *L* then edgeflag[k] = true;end if end for until each edgeflag[k] is true return L;

As mentioned above, different issues contain different timeliness requirements. Therefore, the events are also handled with different priorities. An event contains several features, such as issue quantity, influence area, and event type. The urgency degree of an event is affected by the event type and issue quantity. Administrators often focus on emergencies and wish to discover and solve them as soon as possible. Thus, in our design, the priority that events are displayed is affected by their degree of urgency. If the event is more urgent, its display is much prior. On the one hand, different event types are distinguished with different colors. On the other hand, the higher the quantity of issues, the larger the opacity of color. Through this design, administrators can identify emergencies easily(Fig.8(a)). We also provided a filter tool for administrators to analyze any event type further.

After that, administrators wish to analyze further the time sequence of issues in one event and discover the reason or occurrence pattern of issues. When users focus on one event, other events on the map are hid to better present this event(Fig.1(c)). Meanwhile, we provided an animation based on geo-spatial coordinate for every event, which visualized the occurrence time and spatial location of every issue. A line chart that visualizes the changes in issue quantity over time is drawn(Fig.1(g)). We also displayed a list including content of issues to assist administrators in monitoring issues.

However, when we visualized events over a long period, the quantity of events became large and a mass of overlaps among events caused visualization to become illegible. Users could not easily distinguish between different overlapped events that have similar opacities(Fig.8(b)). Thus, to visualize large quantities of events better, we utilized a design that used the root of the event to stand for the event. The root is an issue that occurred earlier than other issues of the event. We gave preference to the earlier issue that described the essential reason of the event. If all the issues do not contain reasons, we considered the earlier issue that only described the results of the event. Thus, the root is the first issue that reflects the event. The root of event is visualized by a circle, and the larger the influence area of the event, the larger the circle. The influence area is defined as the length of the stroke that draws MST. A larger issue quantity results in a lower transparency of the circle. We also highlighted circles that represented emergencies with yellow strokes to help administrators to identify emergencies quickly(Fig.8(c)).



Fig. 7. a)Visualization of one event before applying algorithm; b)Visualization of one event after applying algorithm; c)Identifying multiple events when intersections between transparent strokes appear.

## 4 EXPERIMENT RESULT

## 4.1 Our Data

The data in our study are a real set of issues collected by the call center in Tianjin, China. A system that supports the normal operations of the call center was developed by our laboratory in 2002. We operate and maintain this system continuously and are thus very familiar with its business processes. The calls are answered 24h a day for 365 days a year. Upon receiving a call, the center will record the issue and hand it off to the concerned division, dispatch concerned division to investigate, solve the issue, and call back after several hours or days to follow up. The citizen will then provide the degree of satisfaction based on the solution situation of the issue. In our system, we established three datasets to track each issue including the characteristic information, hand-off information, and follow-up information of issue. Our datasets contain nine million issues from all divisions from 2007 to 2012. However, most issues pertain to the bus, metro, taxi and light rail divisions, which are queries rather than complaints or issue reports. Hand-off and follow-up on these issues are unnecessary. Thus, these issues were removed from our datasets. We then obtained approximately 500,000 issues from five divisions including heat, water, and natural gas supply, drainage system, and road or bridge divisions.

Characteristic information of every issue contains issue ID (a unique identifier), contact information of the citizen, location of the issue including the geographical coordinate and text description, the time the issue occurred, and detailed issue content. The hand-off and follow-up information include the concerned division, the time issue was solved, the call-back time and the degree of satisfaction.

We also collected notification data from five divisions in this city from 2007 to 2012, including data on cutting off water, gas, or heat supply. Each notification belonging to a particular division contains the effective time quantum, the influence area as well as the detailed notification content. Given that a close relationship exists between weather and the daily lives of people, rainfall data in this city from 2007 to 2012 were also included in our study. The datasets in our study will be released to the public at a later time.

## 4.2 Temporal-Spatial Clustering of Issues

Different divisions present different temporal-spatial clustering patterns. In this section, we present how our system can assist users in identifying the temporal and spatial patterns of different divisions. Fig.9(a) shows that the issues from heat supply mostly appear in November, December, January, February and March of every year. This pattern is obvious because the heat supply system only operates during the heating season. We further analyzed the peaks on a more detailed time scale. We discovered that the peak of every year mostly appeared in November, and the quantity of issues on "No heat supply(NHS)" was almost equal to the total quantity of the issues(Fig.9(a)).

We conducted explorations to determine the underlying reasons. Firstly, we checked whether notifications were released in November. Indeed, the heat supply division usually provide pre-heating to determine whether heating infrastructures are normal and thus, a number of notifications on the start of pre-heating supply are released in early November. Secondly, we compared the spatial distributions of twelve months in every year and found more red areas appearing in the heatmap of November than in other months(Fig.1(e)). We then observed the spatial distribution of different issue manifestations in November(Fig.9(b)). The spatial distribution of issues on NHS was the most severe and the quantities of issues on other manifestations were relatively smaller, such as, the second heatmap of issues on "Pipeline leak(PL)" in Fig.9(b). Our observations proved that the quantity of issues on NHS was always larger than other manifestations. Finally, we checked the content of issues on NHS that occur in the red area of the first heatmap in Fig.9(b) and arrived at the following explanations. On the one hand, more problems occur in the pre-heating period than in other periods. Such problems include low temperature of the boiler and infrastructure failure because they have not been used for a long time. On the other hand, citizens focus on the quality of heat supply to guarantee normal heat supply later. Thus, the peaks appear in this period and the quantity of issues on NHS is the largest.

In Fig.1(e), a small number of issues also occur during the nonheating period. We utilized the line chart and heatmap to visualize the temporal trends and spatial distributions of different manifestations of



Fig. 8. Visualization of Numerous Events. a)Visualization of several events with MST; b)Visualization of numerous events with MST; c)Visualization of numerous events with roots of events.

heat supply in this period. We found that these issues can mostly be attributed to construction by the heat supply division.

We then communicated with the administrators from the heat supply division. They were satisfied with our discoveries and deemed that our explanations were reasonable. They also thought that our discoveries could provide guidance for taking precautions against numerous issues in the pre-heating period. For example, administrators can process some troubleshooting before pre-heating period to discover problems in advance of citizens.



Fig. 9. a)Temporal clustering with season and peaks appearing in November of every year; b)Spatial distributions of issues on "No heat supply" and "Pipeline leak" in November of 2007.

## 4.3 Hotspot Prediction of Drainage System

Rainstorms often impose considerable load on the drainage system of the metropolis. Occasionally, a rainstorm renders the drainage system invalid. From another perspective, city administrators can predict the locations where more issues occur with previous multiple rainfall. Firstly, we used a line chart to observe the temporal trend of rainfall and the quantity of issues from the drainage division(Fig.10(a) and Fig.10(c)). We found that the trend of issue quantity was nearly consistent with that of rainfall. However, the peaks of issue quantity lagged slightly behind the peaks of rainfall. This pattern indicates that we can determine the days with a large amount of rainfall on the basis of the quantity of issues from the drainage division.

Then we used heatmap to determine the spatial distribution of issues from the drainage system when a peak of rainfall occurred(Fig.10(b) and Fig.10(d)). Through the second heatmap in Fig.10(b), we found that several locations had more issues. The system predicted that these locations would still have more issues if rainfall was large as in the second heatmap in Fig.10(b). This prediction can assist the operation and maintenance department in deciding on the locations where repairs on infrastructure need to be conducted. However, related divisions may promote infrastructures given the long-term numerous issues in the same location. Therefore, despite the rainstorms on August 16, 2011 and August 31, 2011, location A in the third and fourth heatmaps of Fig.10(b) contains almost no issue. Furthermore, it proved that infrastructure promotions were effective and city utility service quality increased. The same reason can explain the change in spatial distributions of issues after several rainstorms in 2012(Fig.10(d)).

Administrators from the drainage system provided their views on this discovery. Before using our system, they only became aware of various issues that citizen complained during rains. It was difficult to determine accurately where have more issues by using merely simple statistics and work experiences. Through our prediction method, the operation and maintenance department can repair drainage infrastructures in a more timely and effective manner.

#### 4.4 Causal Relationship of Different Issue Manifestations

Issue manifestations differ in every division. Thus, to analyze associations between these manifestations, we selected every two manifestations in one division and presented their spatial distributions with dots in different colors. Here, issues on PL and NHS are shown on the map



Fig. 10. a)Lag relationship between the quantity of issues from the drainage division and rainfall in 2011; b)Spatial distributions of issues after rainstorms in 2011; c)Lag relationship between the quantity of issues from the drainage division and rainfall in 2012; d)Spatial distributions of issues after rainstorms in 2012.

in November for three years(Fig.11(a)) and different zoom scales for November 2012 are shown in Fig.11(b). The blue circles represent the issues on PL, whereas the turquoise circles represent the issues on NHS. The results of several of our experiments consistently showed that blue circles were surrounded by turquoise circles. Based on the data, during the heating season, more than 90% of the regions that had issues on PL also experienced issues on NHS.

We then determined the underlying causes to confirm our discovery. Indeed, PL was caused by the poor status of heating facilities including valves, radiators, and heating pipelines. On the one hand, the normal operation of the heating system is guaranteed through momentby-moment cycle of cold and hot water. However, when pipelines leak, the cycle is interrupted and cold water dominates the pipelines. Therefore, an increasing number of citizens report issues on NHS. On the other hand, PL often results in rust, which may block the pipelines and cause nonuniform hot water flow. Thus, PL also often leads to NHS. This discovery can provide a warning to administrators that the issues on PL indicate that facilities should be repaired to avoid more serious issues, such as NHS that might pose ill effects on the living conditions of the citizens. However, various other factors may cause NHS. Thus, NHS may possibly be irrelevant to PL. In addition, the issues on PL may not be associated with NHS. Because pipelines may suffer from different types of problems, and not all of them could lead to the issues on NHS.

We presented this causal relationship to administrators from the heat supply division. They were very interested in the results. Then, they confirmed and used our system to discover other relationships among manifestations. After that, they applied the valuable discoveries they obtained by using our system into their work-flow. For example, when citizens report issues on NHS, administrators will firstly check whether there are issues on PL nearby.

## 4.5 Emergency Analysis

Several emergencies occur because of old infrastructures or vandalism and all of them cause a sudden increase in issue quantity. Fig.12(a) shows that the quantity of issues from water supply division changes over time. An anomaly appeared in January 2010. Similarly, we firstly confirmed that no notification existed in this time quantum.

To analyze what happened to the citizens, we used our aggregation method to identify events and visualized them on the map(Fig.12(b)). Among the numerous clusters, a blue one was more severe than others. The cluster was the event that the water pipeline cracked. Finally, we checked the detailed content of issues in the right view of Fig.12(b). Most of the issues pertained to "No water supply(NWS)" and "Water leaking(WL)" when the anomaly appeared because of



Fig. 11. a)Spatial distributions of two manifestations in November 2010, November 2011 and November 2012; b)Different zoom scales of November 2012.

frozen pipelines. To further verify whether the increase in the quantity of issues resulted from the cracked water pipeline rather than other reasons, we analyzed the weather data and found that heavy snowfall occurred on January 3, 2010 and January 4, 2010. The temperatures on these days were so low that outdoor or underground pipelines lost efficiency. Thus, the anomaly was due to the pipeline failures, which was effectively identified by our system and analysis process.

In order to illustrate the spatial distribution of issues caused by this event, the MST of the event is shown in Fig.12(c). All the issues in this event were caused by cracked pipelines, and almost occurred along relatively long water pipelines. Some of issues in this event occurred far from other issues. Therefore, in our aggregation method, the space distance between two issues must be based not only on absolute space distance but also on semantic space relevance.

After that, we used the method mentioned in section 3.2.3 to assist users in tracking the solution situation(Fig.12(d)). We found that most issues were solved during the day, but some serious issues which occurred in the dark blue areas in the heatmap of Fig.12(d) needed to be handled in a span of several days.

The discovery of emergencies was beyond the expectation of administrators, because we had little awareness of utility service management. Administrators confirmed our finding because that emergency was indeed caused by heavy snowfall. What's more, our method provided accurate locations and solution situation for every event.

## 5 CONCLUSION AND FUTURE WORK

In this study, we proposed a visual analytics process for city utility service administrators to discover and solve problems. To identify events that caused numerous issues, we proposed an aggregation method based on our clustering algorithm. We also integrated various visual analytics tools and interaction techniques to analyze the utility service problems. In addition, our system was used in Tianjin, China. Some patterns and events that help administrators in making informed decisions were identified. In the future, we will continue our endeavors by investigating additional factors that may be associated with issue quantity, such as underground network data of infrastructures.

#### REFERENCES

- D. M. Abraham, R. Wirahadikusumah, T. Short, and S. Shahbahrami. Optimization modeling for sewer network management. *Journal of con*struction engineering and management, 124(5):402–410, 1998.
- [2] J. F. Adamowski. Peak daily water demand forecast modeling using artificial neural networks. *Journal of Water Resources Planning and Man*agement, 134(2):119–128, 2008.
- [3] N. Adrienko and G. Adrienko. Spatial generalization and aggregation of massive movement data. *Visualization and Computer Graphics, IEEE Transactions on*, 17(2):205–219, 2011.
- [4] L. E. Agustín-Blas, S. Salcedo-Sanz, P. Vidales, G. Urueta, and J. A. Portilla-Figueras. Near optimal citywide wifi network deployment using a hybrid grouping genetic algorithm. *Expert Systems with Applications*, 38(8):9543–9556, 2011.
- [5] G. Andrienko and N. Andrienko. Spatio-temporal aggregation for visual analysis of movements. In *Visual Analytics Science and Technology*, 2008. VAST'08. IEEE Symposium on, pages 51–58. IEEE, 2008.
- [6] G. Armbruster, B. Endicott-Popovsky, and J. Whittington. Threats to municipal information systems posed by aging infrastructure. *International Journal of Critical Infrastructure Protection*, 6(3):123–131, 2013.
- [7] S. Campisi-Pinto, J. Adamowski, and G. Oron. Forecasting urban water demand via wavelet-denoising and neural network models. case study: city of syracuse, italy. *Water resources management*, 26(12):3539–3558, 2012.
- [8] N. Cao, Y.-R. Lin, X. Sun, D. Lazer, S. Liu, and H. Qu. Whisper: Tracing the spatiotemporal process of information diffusion in real time. *Visualization and Computer Graphics, IEEE Transactions on*, 18(12):2649– 2658, 2012.
- [9] J. Chae, D. Thom, H. Bosch, Y. Jang, R. Maciejewski, D. S. Ebert, and T. Ertl. Spatiotemporal social media analytics for abnormal event detection and examination using seasonal- trend decomposition. In *Visual Analytics Science and Technology (VAST), 2012 IEEE Conference on*, pages 143–152. IEEE, 2012.
- [10] C. Collins, G. Penn, and S. Carpendale. Bubble sets: Revealing set relations with isocontours over existing visualizations. *Visualization and Computer Graphics, IEEE Transactions on*, 15(6):1009–1016, 2009.
- [11] D. P. de Oliveira, D. B. Neill, J. H. Garrett Jr, and L. Soibelman. Detection of patterns in water distribution pipe breakage using spatial scan statistics for point events in a physical network. *Journal of Computing in Civil Engineering*, 25(1):21–30, 2010.



Fig. 12. Emergency Discovery and Analysis. a)Temporal clustering; b)Discovering event that pipelines cracked with our method; c)Spatial distribution of issues caused by this event with MST; d)Solution situation of issues that happened between December 19, 2009 and February 21, 2010.

- [12] N. Ferreira, J. Poco, H. T. Vo, J. Freire, and C. T. Silva. Visual exploration of big spatio-temporal urban data: A study of new york city taxi trips. *Visualization and Computer Graphics, IEEE Transactions on*, 19(12):2149–2158, 2013.
- [13] G. Grekousis and Y. N. Photis. Analyzing high-risk emergency areas with gis and neural networks: The case of athens, greece. *The Professional Geographer*, 66(1):124–137, 2014.
- [14] T. Grubic, L. Varga, and S. Varga. Future utility services(un) knowns framework: Knowledge existence and knowledge reach. *Futures*, 54:68– 86, 2013.
- [15] L. Hernandez, C. Baladrón, J. M. Aguiar, B. Carro, A. J. Sanchez-Esguevillas, and J. Lloret. Short-term load forecasting for microgrids based on artificial neural networks. *Energies*, 6(3):1385–1408, 2013.
- [16] P. Inamdar, S. Cook, A. Sharma, N. Corby, J. O'Connor, and B. Perera. A gis based screening tool for locating and ranking of suitable stormwater harvesting sites in urban areas. *Journal of environmental management*, 128:363–370, 2013.
- [17] S. Jiao, Y. Qu, Z. Liu, Q. Feng, J. Ren, and X. Chen. Webgis application based on real-time traffic flow network analysis. In *Geoinformatics 2007*, pages 67542K–67542K. International Society for Optics and Photonics, 2007.
- [18] A. Khotanzad, R. Afkhami-Rohani, T.-L. Lu, A. Abaye, M. Davis, and D. J. Maratukulam. Annstlf-a neural-network-based electric load forecasting system. *Neural Networks, IEEE Transactions on*, 8(4):835–846, 1997.
- [19] S. Kim, R. Maciejewski, A. Malik, Y. Jang, D. Ebert, and T. Isenberg. Bristle maps: A multivariate abstraction technique for geovisualization. 2013.
- [20] S. Ko, R. Maciejewski, Y. Jang, and D. S. Ebert. Marketanalyzer: an interactive visual analytics system for analyzing competitive advantage using point of sale data. In *Computer Graphics Forum*, volume 31, pages 1245–1254. Wiley Online Library, 2012.
- [21] M. Kulawiak, Z. Lubniewski, K. Bikonis, and A. Stepnowski. Geographical information system for analysis of critical infrastructures and their hazards due to terrorism, man-originated catastrophes and natural disasters for the city of gdansk. In *Information Fusion and Geographic Information Systems*, pages 251–262. Springer, 2009.
- [22] N. Kunwar, R. Kumar, et al. Area-load based pricing in dsm through ann and heuristic scheduling. 2013.
- [23] H. Liu, Y. Gao, L. Lu, S. Liu, H. Qu, and L. M. Ni. Visual analysis of route diversity. In Visual Analytics Science and Technology (VAST), 2011 IEEE Conference on, pages 171–180. IEEE, 2011.
- [24] L. Liu, M. Zhang, L. Chen, and C. Du. The research of city community safety management system. In *Computer Science and Information Technology (ICCSIT), 2010 3rd IEEE International Conference on*, volume 2, pages 243–247. IEEE, 2010.
- [25] R. Maciejewski, R. Hafen, S. Rudolph, S. G. Larew, M. A. Mitchell, W. S. Cleveland, and D. S. Ebert. Forecasting hotspotsła predictive analytics approach. *Visualization and Computer Graphics, IEEE Transactions on*, 17(4):440–453, 2011.
- [26] I. MACRO GROUP and U. S. of America. Building a 311 system: A case study of the city of minneapolis. 2008.
- [27] A. Malik, R. Maciejewski, B. Maule, and D. S. Ebert. A visual analytics process for maritime resource allocation and risk assessment. In *Visual Analytics Science and Technology (VAST)*, 2011 IEEE Conference on, pages 221–230. IEEE, 2011.
- [28] A. Mostashari, F. Arnold, M. Maurer, and J. Wade. Citizens as sensors: The cognitive city paradigm. In *Emerging Technologies for a Smarter World (CEWIT), 2011 8th International Conference & Expo on*, pages 1–5. IEEE, 2011.
- [29] T. Munzner. A nested model for visualization design and validation. *Visualization and Computer Graphics, IEEE Transactions on*, 15(6):921–928, 2009.
- [30] N. Okada, L. Fang, and D. M. Kilgour. Community-based decision making in japan. *Group Decision and Negotiation*, 22(1):45–52, 2013.
- [31] V. Ostromoukhov, C. Donohue, and P.-M. Jodoin. Fast hierarchical importance sampling with blue noise properties. In ACM Transactions on Graphics (TOG), volume 23, pages 488–495. ACM, 2004.
- [32] P. Pederson, D. Dudenhoeffer, S. Hartley, and M. Permann. Critical infrastructure interdependency modeling: a survey of us and international research. *Idaho National Laboratory*, pages 1–20, 2006.
- [33] H. Piringer, M. Buchetics, and R. Benedik. Alvis: Situation awareness in the surveillance of road tunnels. In *Visual Analytics Science and Tech*-

nology (VAST), 2012 IEEE Conference on, pages 153-162. IEEE, 2012.

- [34] A. R. Pradhan, D. F. Laefer, and W. J. Rasdorf. Infrastructure management information system framework requirements for disasters. *Journal* of computing in civil engineering, 21(2):90–101, 2007.
- [35] S. Savoska, S. Loskovska, and I. Dimitrovski. Information visualization from the public utilities databases of local municipality for municipalities managers. In *Information Technology Interfaces, 2008. ITI 2008. 30th International Conference on*, pages 237–242. IEEE, 2008.
- [36] A. Shukla, J. Pekny, and V. Venkatasubramanian. An optimization framework for cost effective design of refueling station infrastructure for alternative fuel vehicles. *Computers & Chemical Engineering*, 35(8):1431– 1438, August 2011.
- [37] J. Tanimoto and A. Hagishima. Total utility demand prediction system for dwellings based on stochastic processes of actual inhabitants. *Journal* of Building Performance Simulation, 3(2):155–167, 2010.
- [38] D. Teng, S. Song, H. Yang, C. Ma, H. Wang, and G. Dai. An approach to visual analysis for task flow management. *Science China Information Sciences*, 56(5):1–12, 2013.
- [39] D. J. Vanier. Geographic information systems (gis) as an integrated decision support tool for municipal infrastructure asset management. In *Proceedings of the CIB 2004 Triennial Congress, Toronto, Ont*, pages 2–9, 2004.
- [40] R. Walker, A. Slingsby, J. Dykes, K. Xu, J. Wood, P. H. Nguyen, D. Stephens, B. Wong, and Y. Zheng. An extensible framework for provenance in human terrain visual analytics. *Visualization and Computer Graphics, IEEE Transactions on*, 19(12):2139–2148, 2013.
- [41] R. Wang, W. Dong, Y. Wang, K. Tang, and X. Yao. Pipe failure prediction: A data mining method. In *Data Engineering (ICDE), 2013 IEEE* 29th International Conference on, pages 1208–1218. IEEE, 2013.
- [42] Z. Wang, M. Lu, X. Yuan, J. Zhang, and H. v. d. Wetering. Visual traffic jam analysis based on trajectory data. *Visualization and Computer Graphics, IEEE Transactions on*, 19(12):2159–2168, 2013.
- [43] Z. Zhang, Y. Liu, B. Chen, and K. Chen. Using gis and kde analysis spatial distribution on public housing households: A case study. In *Computer Science & Education (ICCSE), 2013 8th International Conference* on, pages 925–930. IEEE, 2013.