

Quality Assessment of Volunteered Geographic Information

Roya Esmaili^{1*}, Farzin Naseri², Ali Esmaili³

¹GIS Engineering, Graduate University of Advanced Technology, Kerman, Iran

²Dep. of GIS Engineering, Graduate University of Advanced Technology, Kerman, Iran

³Dep. of RS Engineering, Graduate University of Advanced Technology, Kerman, Iran

Abstract Recent advances in spatial data collection technologies and online services dramatically increase the contribution of ordinary people to produce, share and use geographic information. Collecting spatial data and disseminating them on the internet by citizens has led to a huge source of spatial data termed as Volunteered Geographic Information (VGI) by Mike Goodchild. Despite the advantages of VGI, there are a lot of doubts about its quality. This article examines the early literature on this phenomenon and illustrating the current methods for quality assessment and assurance of VGI. Almost all the existing researches are based on comparing the VGI data with the accurate official data. However in many cases there is no access to correct data, so looking for an alternative way to determine the quality of VGI data is essential. In this article a method for positional accuracy assessment of VGI is suggested based on comparing the existing data of the same place with each other according to the metadata that their creators have obtained. The proposed method was implemented for the different maps that were produced by various methods from our case study.

Keywords Crowd-sourcing, Volunteered Geographic Information (VGI), Quality Assessment

1. Introduction

Historically, professional surveyors, cartographers, geographers and governmental agencies endeavoured to provide geographic data from the Earth in an authoritative manner[1] through several methods including surveying, photogrammetry and remote sensing[2], which could be later ordered by the users in paper or digital format. However, during the last decade, Internet was increasingly deployed for providing the users with geographic information, resulted in “Geographic World Wide Web”[3]. In this period consumers could just download maps free or by paying to use them and they couldn’t make changes, edit or add something to maps on the Internet.

Improvement of web-based mapping, invention of cell phones and devices that are equipped with Global Positioning System (GPS), PDAs and digital cameras have made it possible for ordinary people to collect spatial data, which are then shared and disseminated on the internet using web map services and specially Web 2.0[3]. Producing maps and geographic Information by non-expert people without academic studies and with local knowledge about their environment and generally the world, is

preparing a phenomena that is named by different terms in researches such as Neogeographic[4], Public Participation GIS[5], Ubiquitous cartography[6] and Goodchild named this phenomena Volunteered Geographic information (VGI)[7]. Web sites such as OpenStreetMap (OSM) and Wikimapia, aiming to produce a free and editable map of the world, are examples of VGI[8].

VGI possess many advantages such as free access, quick acquisition and provides types of data of places as well. However, accuracy and quality assessment of these data still are challenge for researchers.

In this paper we reviewed the existing methods for quality assessment of VGI data and suggest a method for determining the quality of data to the users that the majority of them are the ordinary people and have no academic knowledge. Section 2 provides the background of crowdsourced data and specially describes the VGI. Section 3 explains the quality of crowdsourced data and presents the existing methods for quality assessment of the user generated content, particularly in spatial domain. In Section 4 our method for spatial quality assessment of VGI data is explained. Section 5 presents the results of applying the proposed method for a case study. Finally, section 6 contains some concluding remarks and ideas for future work.

* Corresponding author:

r.esmaily@student.kgut.ac.ir (Roya Esmaili)

Published online at <http://journal.sapub.org/ajgis>

Copyright © 2013 Scientific & Academic Publishing. All Rights Reserved

2. User Generated Content

There are different definitions for user-generated content, however Goodchild expressed that user-generated content refers to “the ability of users to create content that is then integrated and made available through Web sites”[9]. The Organization for Economic Cooperation and Development (OECD) defined UGC as a creative content that is created by public without any professional practices and is disseminated on the internet[10].

UGC is not limited to the internet and people can generate content everywhere for example when they tell someone something they’re generating content[10]. However internet made the creation and distribution of the content easier and quicker. Web 2.0 converted internet from one-directional, that users could just see the content of web sites, to a bi-directional one that users in addition to using and seeing the contents of the internet could create and generate data by themselves and share it on the Internet[11]. Social networks such as Facebook, Tweeter and creation of weblogs and wikis and the photo and video sharing websites such as Flickr and YouTube are all using the Web 2.0 technology and made it possible for ordinary people to create data and disseminate it on the Internet. The most famous example of crowd sourcing websites is Wikipedia that people can create articles and edit them.

The ability of individuals to create their own digital geographic information and acquiring spatial data such as place names, topographic features and transport networks by ordinary people through up to five functioning senses and by means of tools that are equipped with GPS brought out this idea that people can be imagined as mobile sensors[11]. They can collect data about everywhere and mostly free and share it with others. Goodchild named the user generated content in spatial domain “Volunteered Geographic Information” or in brief VGI[7]. VGI has changed the way of producing and sharing and the characteristics of digital spatial data[12].

By traditional spatial data, collected by official organization of map provider, we cannot answer to the questions in detail such as the name of a local restaurant in a specific street, maybe because these data are not important to be collected or change rapidly and cannot be updated quickly[13]. While one of the advantages of VGI data is that we can be aware of local spatial events by citizens that would be time consuming to be collected officially[14].

Ostlaender named some of VGI applications such as: finding new bike road, monitoring forest fires, or analyzing tweeter messages to create landing map of planes[15]. The almost famous examples of VGI are websites such as OpenStreetMap and Wikimapia that provided tools for users to add geographic data on the map by uploading their GPS tracks or digitizing the background satellite image or simply naming the streets or places all around the world.

3. Quality Assessment of Crowd-Sourced Data

Using improper data for an application can cause many

socio-economic problems. Crowd sourced data are not the exception and their quality should be determined before usage.

In spite of doubts about the user generated contents which are collected by non-experts and without any quality standards, the quality gap is balanced by volume, as user generated content is much larger than editorial content[16].

As the volume of such collaborative information increases, the problem of assessing its quality, preventing vandalism and spam grows simultaneously[17].

Several studies have evaluated the accuracy of crowdsourced data such as quality of Flickr tags[18] or quality of Wikipedia articles[19]. The majority of researches were about quality assessment by reputation systems.

Reputation systems are in two categories; user-driven and content-driven. User-driven which means giving reputation to the article or product by users according to their own factors such as websites like Baywhere the users themselves will assess the process of selling and buying of each other; on the other hand content-driven reputation is giving reputation without the interference of users and according to the content of an article or giving reputation to contributors according to the lifespan of their edits on articles[19] the more the article has words, the more reputation is gained.

3.1. Quality Assessment of VGI

Awareness of the quality of spatial data as Tveite said is important because of two reasons; economic and legal importance of spatial data in decision making processes and the possibility of combining multiple spatial data set for different purposes[20]. For example OSM maps are used in different commercial projects as background maps so awareness of their quality is essential[21].

Data quality can be defined as fitness for purpose, or how suitable some data is in satisfying particular needs or fulfilling certain requirements to solve a problem[22]. When considering the quality of spatial data, most naive users consider only the positional accuracy of data. However, the spatial quality is more than that and has different aspects: positional accuracy, attribute accuracy, currency, completeness, logical consistency, lineage, accuracy and resolution[10].

VGI data are sometimes called “asserted” because there is no standard for checking their quality and there is no reference or citation for them, in divers the official data are called “authoritative” because their quality is checked with standards[14]. Although the quality of VGI data might not be clear but in emergency situation such as forest fires where we have no official data, using volunteered data with quality vagueness is better than waiting for better data to arrive[23]. The most significant advantage of VGI is that they can be up to date in less time than traditional data so for the projects that we have time limitations such as updating the streets of a city in few weeks the local citizens are the best source to collect data and update the

information about the streets[9].

Cooper mentioned some of the challenges that VGI has for quality assessment, one of them is that the user cannot assess the quality of data in isolation because the quality depends on some parameters such as the data user, purpose and the content in which it is used the users are not forced to create metadata for the data and also all aspects of quality cannot be assessed quantitatively and they depend on the language (e.g. about the completeness degree of a data set)[10].

As Goodchild mentioned[24] quality assurance in traditional data is done in two parts; one of them is assuring and checking the quality of data during the contribution and creating them and the second one is checking the quality of data after making them by comparing them with reference data and save the results in form of metadata. He discussed three approaches to quality assurance of VGI data follows the first procedures of quality assurance; Crowd-sourcing, social and geographic approaches. Here they will be explained briefly.

● **Crowd-sourcing approach:** It expresses the number of people even without qualification can solve the problem that is caused by individual and also the answer of a group is converge to the truth in compare to the answer of an individual.

● **Social approach:** In this approach there will be a gate-keeper that will check the entrance of data to the data set to avoid vandalism, like the Data Working Group (DWG) of OpenStreetMap who will decide the solution in critical situations such as arguing between users about a certain place or any vandalism and violation that may be happened.

● **Geographic approach:** This approach suggests checking the geographic data according to some rules such as language structures. For example the data that are near each other should be consistence and obey some rules.

Goodchild also mentioned that the first approach may not be so useful for geographic domain because there will be many spatial error despite of many volunteers in an area. But the last approach can be more useful because it can be automated and it's specifically for geographic data.

As mentioned above there are two processes for quality assurance that Goodchild defined three approaches for the first one. Here the existing methods for quality assessment that followed the second part are explained.

3.1.1. Metadata

One method to determine the quality of data set is using metadata. Despite the importance of metadata no metadata exists for projects such as OpenStreetMap or Google Earth[25], for example even there is no information about the date that the image has been obtained in Google Earth[14]but the user that contributes in volunteered projects may be volunteer to give more data about the data[9].

Goodchild argued for a binary user-centric metadata that in addition to the single quality of the data can describe the

ability of two dataset to work together because in geospatial web relative quality of data sets that are being integrated is so important. He called it metadata 2.0[9].

The quality of spatial data collected by a contributor can be determined by factors such as the experience of the contributor in GIS projects and specially crowdsourced projects or how they apply metadata to their contributed data[26].

The vagueness of crowdsourced data can be determined by two types of metadata one of them is "user-encoded vagueness metadata" for example the number that the user himself give to the data that has been created and the other type is "system-created vagueness metadata" in which the system itself save to define the quality of data such as saving the scale in which the data has been added automatically[27].

3.1.2. Comparison

As mentioned in the previous section one of the spatial quality parameters is positional accuracy. To measure this parameter, the data should be compared with an accurate exiting data. The main problem for measuring positional accuracy is choosing proper dataset that prepared data be compared with proper dataset[28]. However on of the problem with this method specially about volunteered geographic data is that in many cases we do not have access to a suitable ground-truth datasets for comparison[26]. Generally evaluating positional accuracy is based on a buffer technique that was developed by Goodchild and Hunter in 1997; illustrated by Figure 1. Ather in 2009 compared OSM data with OS data in England[22] and Kounadi compared OSM data with HMGS data in Athens, Greece using the mentioned buffer method[20].

In the absence of availability of a ground truth dataset a simple alternative solution is to analyze the density of data points within the grid squares (Figure 2)[26].

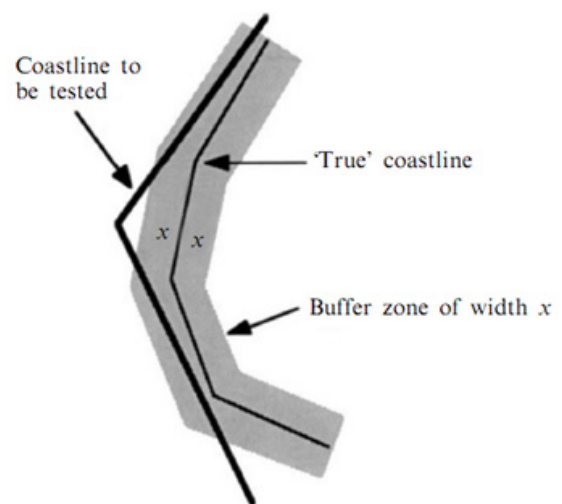


Figure 1. Width x is created around the reference object and the proportion of the tested source that lines within the buffer is then calculated. Depending on the size of the buffer chosen, the level of accuracy of the tested source can then be determined[29]

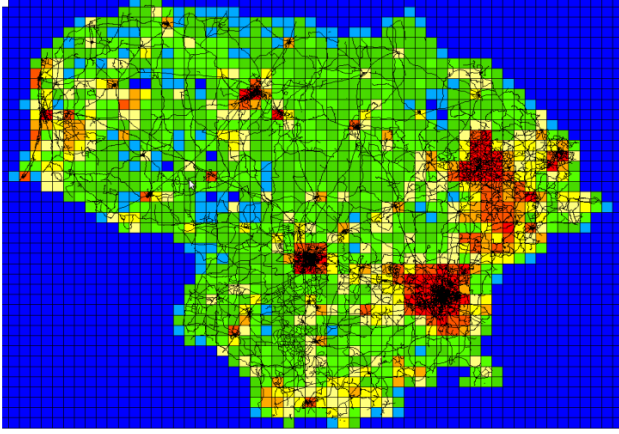


Figure 2. Density of points in OSM Lithuania using 5KM grid squares[26]

3.1.3. Reputation System

Some VGI websites for controlling their contributors and avoiding vandalism have ranked their contributors according to different parameters. Wikimapia has ranked its users according to the level of experience, number of edits, and the number of objects or pictures and etc. that has added. For example a new comer will ranked in lower level(level zero) that means he can only add objects and cannot delete any objects on the map but after a period of time and gaining experience of mapping he will go to level 1 that can delete objects and will have more authority. And the highest level belongs to the administrators that have the most authority for instance they can block the user that make vandalism.

4. Our Proposed Method

The most existing methods for quality assessment of volunteered data as mentioned are based on comparing with correct data. However, in many cases there is no access to the correct data. Therefore, we would be forced to look for an alternative method. Among different parameters of spatial quality we suggest a method for positional accuracy.

To determine the positional accuracy of data sets that belong to the same area we can follow the below procedure:

1. Ask the user to get information about the way that he or she has collected the data and the precision of the tools that were used. Determine whether he used GPS or digitizing satellite image or etc. and if he used GPS determine the precision of the device. And also information about his age, being citizen of the place that he collects the data from or not and etc.

2. Set an initial positional accuracy for each data set according to the information that its creator has given to us.

3. Determine in common points between maps and label them one by one (the same points in different maps will have the same label). Then, points of the same label will be grouped.

4. The weighted mean and standard deviation will be calculated for each group.

5. In each group the distance of each point is calculated with the mean value of the group. Thereafter the calculated distances compared with the triple of standard deviation and if it was lower the point will be kept and otherwise it will be omitted.

6. The process 4 and 5 will be repeated for the group that a point has been left out from. It will be repeated until no point is omitted.

7. The mean value of each group is determined, and then the distance of each point (except the ones that have been omitted) with the mean value of its group is calculated. This distance is the error for each point.

8. The final error of each data set is determined by calculating the mean value of the errors of the points in each data set that was calculated in the previous level.

9. The positional accuracy of each dataset is the inverse of the error of it.

5. Implementation

5.1. Study area

Kerman Graduate University of Technology was chosen to implement the proposed approach (figure 3).



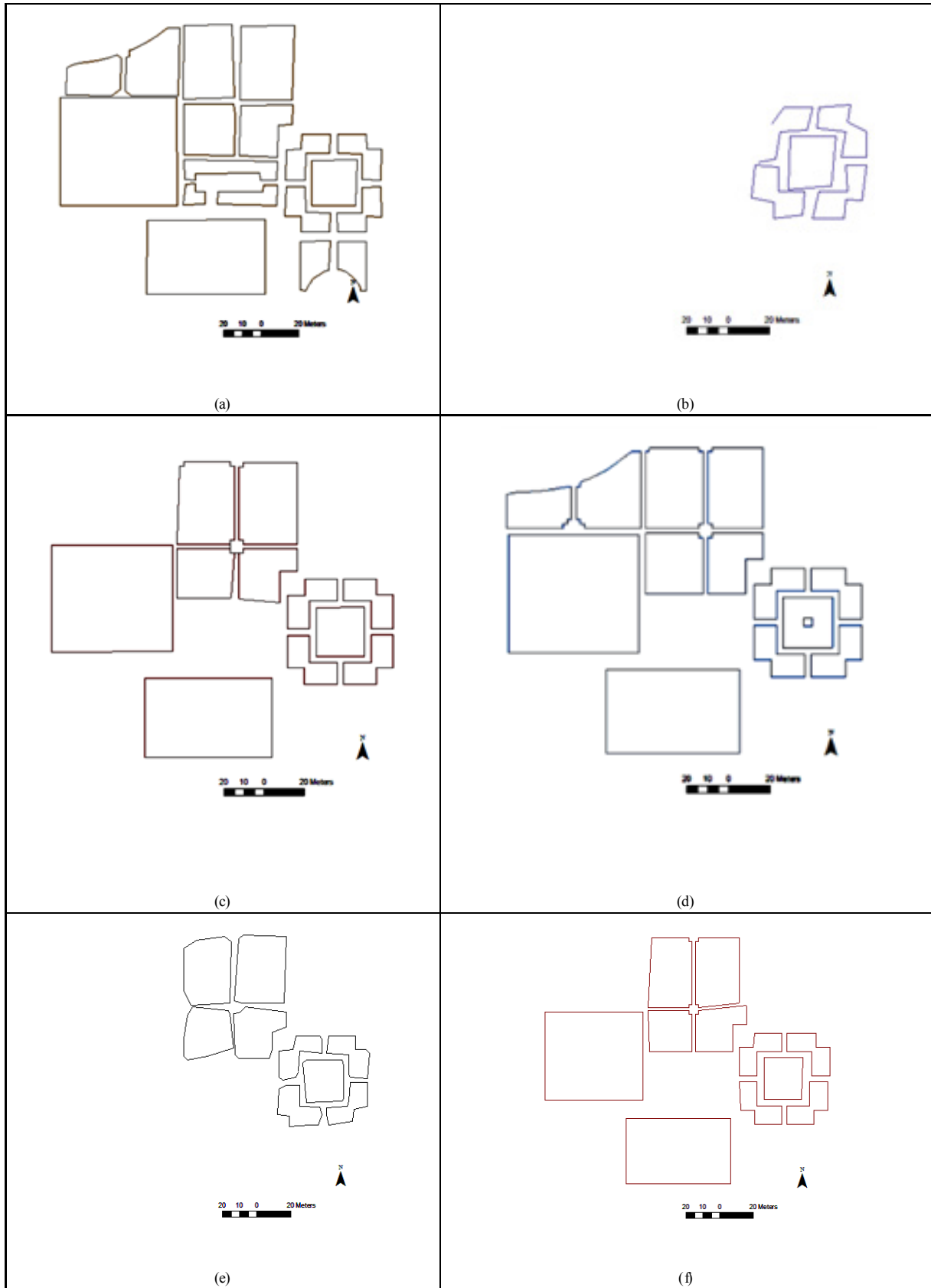
Figure 3. The Google Earth image of study area (Kerman Graduate University of Technology)

16 versions of 2D map were produced using different data collection methods in the same local coordinate system. The methods and the number of each dataset are shown in table 1. Among these maps one from each method as sample is shown in table 2. Also we can see that these maps have differences with each other in figure 4.

Table 1. Number of maps created in each method

Method	Iteration
Total Station	2
Walking	3
Meter	2
Digitizing	3
GPS Track	4
GPS Mark	2

Table 2. Sample maps of different method (a) Digitize-(b) Mark GPS-(c) Meter-(d) Total Station-(e) Track GPS-(f) Walking



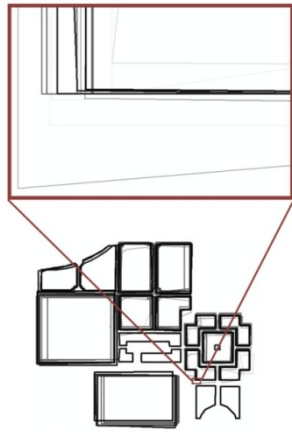


Figure 4. Overlay of 16 Maps. We can see the difference between maps

5.2. Positional Accuracy

To determine the positional accuracy of each map the following processes should be passed:

1. Gathering information about map from its creator such as the method that the data was produced by.
2. The error was set for each map according to its method, and then the initial accuracy based on the error was calculated. The accuracy is the inverse of the error (table 3).

Table 3. Error and accuracy of each method is shown. The calculated accuracy was multiplied to 70 to have the first accuracy start from 1

Method	Error	Accuracy
Mark GPS	70 cm	1
Tack GPS	50 cm	1.4
Digitize	50 cm	1.4
Walk	25 cm	2.8
Meter	5 cm	14
Total Station	5 cm	14

3. In common points between maps and their coordinates were defined. There were 23 in common points between the 16 maps (figure 5). The points were labelled from 1 to 23 the same points in maps have the same number. Therefore 23 groups of points that each of them has 16 members were created.

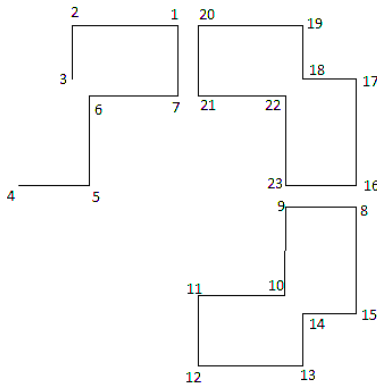


Figure 5. In common points between 16 maps

4. The weighted mean value and standard deviation were calculated for each of these 23 groups.

$$X_i = \frac{\sum_{j=1}^{16} (X_{ij} * W_j)}{\sum_{j=1}^{16} W_j} \quad i = 1 \text{ to } 23$$

$$Y_i = \frac{\sum_{j=1}^{16} (Y_{ij} * W_j)}{\sum_{j=1}^{16} W_j}$$

$$std X_i = \sqrt{\frac{\sum_{j=1}^{16} (X_{ij} - X_i)^2}{16}}$$

$$std Y_i = \sqrt{\frac{\sum_{j=1}^{16} (Y_{ij} - Y_i)^2}{16}}$$

$$std com_i = \sqrt{\frac{2}{var X_i + var Y_i}}$$

y_{ij} = The y of point i in map j
 X_{ij} = The x of point i in map j
 \bar{X}_i = The mean value of x of points i
 \bar{Y}_i = The mean value of y of points i
 $std com_i$ = complete standard deviation
 $std X_i$ = standard deviation of x of points i
 $std Y_i$ = standard deviation of y of points i

5. In each group the distance of each point is calculated with the mean value of the group (error of each point). Thereafter the calculated distance is compared with the $3 * standard deviation$ and if it was lower the point will be kept and otherwise it will be omitted.

Point number 12 from GPS-T2, points 8, 10 and 15 from GPS_W2 and point 13 from GPS_T1 were omitted and again the levels 2 and 3 were repeated for and the points 8, 10, 12, 13, 15 (The distance for point 8 in each map and its mean is shown in table 4 as example).

Table 4. The distance between point 8 in each 16 maps and the mean value (d8)

3 times of Standard deviation of point 8 = 4.985553	
Map ID	$ p_8 - \bar{p}_8 $
Digitize1	0.6640022
Digitize2	1.0370212
Digitize3	0.6828414
Walk1	0.9578549
Walk2	0.7410992
Walk3	0.5673364
GPS_T1	0.8496667
GPS_T2	0.7947039
GPS_T3	1.6886566
GPS_T4	1.4142113
GPS_W1	2.8686217
GPS_W2	5.0346549
Metr1	0.361555
Metr2	0.3424417
Total1	0.5098871
Total2	0.3564448

6. The process of 4 and 5 were repeated for the group of points 8, 10, 12, 13 and 15 until no points was leaved out any group. After one repeat point 10 from the map GPS_W1 was omitted and after that there was no need to repeat.

7. After omitting 12 from GPS-T2, points 8, 10 and 15

from GPS_W2 and point 13 from GPS_T1 and 10 from GPS_W1, the mean value of these groups is determined again: 8, 10, 12, 13 and 15. Thereafter the distance of each point (except the ones that have been omitted) with the mean value of its group is calculated (error of point).

8. The final error of each data set is determined by calculating the mean value of the errors of the points in each data set that was calculated in previous levels.

9. The positional accuracy of each dataset is the inverse of the error of it. After defining accuracy they have been scaled to the range of 0 to 100 to be easier for interpretation (table 5).

Table 5. The final positional accuracy of each map

Map ID	Positional Accuracy	Map ID	Positional Accuracy
Digitize1	47.1	GPS_T3	9.9
Digitize2	41.1	GPS_T4	16.7
Digitize3	40.3	GPS_W1	8.6
Walk1	20.9	GPS_W2	7.59
Walk2	32.5	Metr1	59.8
Walk3	44	Metr2	71.1
GPS_T1	20.9	Total1	64.1
GPS_T2	17.2	Total2	100

6. Conclusions

Crowdsourced data in general and volunteered geographic information in particular, are becoming the huge source of data. VGI has enormous advantages such as: it's free, has the ability to produce large amount of data in short period of time and collect local data that are in some cases impossible to obtain them by traditional methods of mapping. Despite its benefits, volunteered geographic information cannot be used in many applications because its quality is not determined and there is vagueness about it. Therefore many researches have done to determine the quality of VGI.

In this article the existing methods for quality assessment of crowdsourced data and VGI is explained. The majority of the methods are based on comparing the VGI data with an accurate official data but in most cases there is no access to accurate data. We looking for an alternative way and we suggested a method for assessing positional accuracy based on comparing the existing data of the same place with each other according to the metadata that their creators have given. The proposed method was implemented for the different maps that were produced by various methods from our case study.

In this paper just a method for positional accuracy assessment was suggested. Determining the other spatial quality parameters and also display the quality to the users of volunteered geographic information with the methods that can be easily interpreted by ordinary people can be the aim of future researches.

REFERENCES

- [1] Andrew Flanagin and Miriam Metzger, "The credibility of volunteered geographic information", *GeoJournal*, Vol.72, no.3, pp.137-148, 2008.
- [2] W. T. Castelein, L. Grus, J. W. H. C. Crompvoets, and A. K. Bregt, "A Characterization of Volunteered Geographic Information", in *13th AGILE International Conference on Geographic Information Science*, pp.10, 2010.
- [3] Muki Haklay, Alex Singleton, and Chris Parker, "Web Mapping 2.0: The Neogeography of the GeoWeb", *Geography Compass*, Vol.2, no.6, pp.2011-2039, 2008.
- [4] Sara Elwood, "Volunteered geographic information: key questions, concepts and methods to guide emerging research and practice", *GeoJournal*, Vol.72, pp.133-135, 2008.
- [5] Geisa Bugs, Oscar Fonts, Joaquín Huerta, and Marco Painho, "An assessment of Public Participation GIS and Web2.0 technologies in urban planning practice in Canela, Brazil", *Cities* Vol.27, no.3, pp.172-181, 2010.
- [6] Georg Gartner, David A. Bennett, and Takashi Morita, "Towards Ubiquitous Cartography", *Cartography and Geographic Information Science*, Vol.34, no.4, pp.247-257, 2007.
- [7] Michael F. Goodchild, "Citizens as Sensors: Web 2.0 and The Volunteering of Geographic Information", *GeoFocus*, Vol.7, pp.8-10, 2007.
- [8] M. Haklay, P. Weber, "OpenStreetMap: User-Generated Street Maps", *PERVASIVE computing*, Vol.7, no.4, pp.12-18, 2008.
- [9] Michael F. Goodchild, "Spatial Accuracy 2.0", in *8th international symposium on spatial accuracy assessment in natural resources and environmental sciences*, 2008.
- [10] Antony Cooper, Serena Coetzee, and Derrick Kourie, "Volunteered geographical information – the challenges", *PoPositionIT*, 2012.
- [11] Michael F. Goodchild, "Citizens as Voluntary Sensors: Spatial Data Infrastructure in the World of Web 2.0", *International Journal of Spatial Data Infrastructures Research*, Vol.2, pp.24-32, 2007.
- [12] Sara Elwood, "Volunteered geographic information: future research directions motivated by critical, participatory, and feminist GIS", *GeoJournal*, Vol.72, pp.173-183, 2008.
- [13] X. and L. Di Qian, "Data cleaning approaches in Web2.0 VGI application", in *Geoinformatics, 2009 17th International Conference* pp.1-4, 2009.
- [14] Michael F. Goodchild, "Citizens as sensors: the world of volunteered geography", *GeoJournal*, Vol.69, no.4, pp.211-221, 2007.
- [15] Nicole Ostlaender and Robin S. Smith, "What Volunteered Geographic Information is (good for) - designing a methodology for comparative analysis of existing applications to classify VGI and its uses", *Geoscience and Remote Sensing Symposium (IGARSS)*, IEEE, pp.1422-1425, 2010.
- [16] Ricardo Baeza-Yates, "User generated content: how good is it?", in *3rd workshop on Information credibility on the web*,

- pp.1-2, 2009.
- [17] Krishnendu Chatterjee, Luca de Alfaro, and Ian Pye, "*Robust Content-Driven Reputation*", in *1st ACM workshop on Workshop on AISec*, pp.33-42, 2008.
- [18] D. Liu, M. Wang, et al, "*Tag quality improvement for social images*", Multimedia and Expo, 2009. ICME 2009. IEEE International, 2009.
- [19] B. Thomas Adler and Luca de Alfaro, "*A Content-Driven Reputation System for the Wikipedia*", in *16th international conference on World Wide Web*, pp.261-270, 2007.
- [20] Ourania Kounadi, "Assessing the quality of OpenStreetMap data", University College of London, 2009.
- [21] M. van Exel, E. Dias, and S. Fruijtjer, "*The impact of crowdsourcing on spatial data quality indicators*", in *GIScience*, 2010.
- [22] Aamer Ather, "A Quality Analysis of OpenStreetMap Data", University College London, 2009.
- [23] Sara Elwood, Michael F. Goodchild, and Daniel Z. Sui, "*Researching Volunteered Geographic Information: Spatial Data, Geographic Research, and New Social Practice*", *Annals of the Association of American Geographers*, Vol.102, no.3, pp.571-590, 2012.
- [24] Michael F. Goodchild and L. Li, "*Assuring the quality of volunteered geographic information*", *Spatial Statistics*, Vol.1, pp.110-120, 2012.
- [25] Michael F. Goodchild, "*NeoGeography and the nature of geographic expertise*", *Journal of Location Based Services*, Vol.3, no.2, pp.82-96, 2009.
- [26] Błażej Ciepluch, Peter Mooney, and Adam C. Winstanley, "*Building Generic Quality Indicators for OpenStreetMap*", in *19th annual GIS Research UK (GISRUK)*, 2011.
- [27] Bertrand De Longueville, Nicole Ostländer, and Carina Keskitalo, "*Addressing vagueness in Volunteered Geographic Information (VGI) – A case study*", *International Journal of Spatial Data Infrastructures Research*, Vol.5, 2010.
- [28] Mordechai Haklay, "*How good is volunteered geographical information? A comparative study of OpenStreetMap and Ordnance Survey datasets*", *Environment and Planning B: Planning and Design*, Vol.37, no.4, pp.682-703, 2010.
- [29] Michael F. Goodchild and Gary J. Hunter, "*A simple positional accuracy measure for linear features*", *International Journal of Geographical Information Science*, Vol.11, no.3, pp.299-306, 1997.