Comparative Spatial Analysis of Positional Accuracy of OpenStreetMap and Proprietary Geodata

Marco HELBICH, Christof AMELUNXEN, Pascal NEIS and Alexander ZIPF

Abstract

The emergence and ubiquitary availability of geotechnologies yield a rapid increase of user generated geographical data, utilized for mapping, modeling etc. On the example of a well mapped German city this paper analyzes the positional accuracy of OpenStreetMap and TomTom data by means of a statistical comparative approach using official survey data as the reference dataset. The results show that OpenStreetMap and TomTom data feature similar spatial deviations while both do not coincide with the survey data. Furthermore, OpenStreetMap data show spatial heterogeneity in the positional error distribution, leading to significant clusters of high and low positional accuracy.

1 Introduction

Profound changes occurred in GIScience lately (e.g. GOODCHILD 2007, ELWOOD 2008, SUI 2008). Until now, the generation, maintenance, and distribution of geographic data had been almost solely the domain of either official land surveying offices or commercial companies. This was due to the immense costs related to the actual surveying and data maintenance as well as the efforts involved to share and distribute spatial data. What can be seen nowadays is a massive increase of geographic data collected and shared by volunteers, working in a collaborative fashion via the Web. The dramatically reduced costs of modern handheld devices equipped with satellite navigation have enabled people to privately collect geographic data with ease of use and in precision levels which had formerly been simply beyond reach for the masses. Furthermore, the progress of the internet to the participatory "Web 2.0" approach has made collaborative efforts to generate and share content of various kinds very common. This phenomenon is widely known as Volunteered Geographic Information (VGI, GOODCHILD 2007, ELWOOD 2008). In combination with nowadays ubiquitous available software (e.g., Google Earth) and miscellaneous Web services, Sui (2008:1) calls this revolutionary development the "wikification of GIS".

Among a broad list of initiatives working with VGI, OpenStreetMap (OSM) is one of the most promising crowd sourced projects. When the project started, its primary goal was simply to generate a free map of the world through volunteered participation. Nowadays, the objective is to build up a geodatabase, where the collected spatial data are made publicly available and may thus be used for other individual purposes, such as regional planning (HAGENAUER & HELBICH 2012), disaster management (NEIS et al. 2010), among others. Additionally, OSM serves as a platform for location based services, including routing (NEIS & ZIPF 2008), geocoding (AMELUNXEN 2010), accessibility analysis and spatial searches.

However, using these data means accepting their limitations, particularly concerning the data quality. For a comprehensive overview of several data quality aspects we refer to VAN OORT (2006). This paper instead focuses exclusively on positional accuracy. It denotes the coordinate deviation of a spatial object compared to its real location (HAKLAY 2010). The positional accuracy of the collected data is affected by different influences, e.g. the technological bias like the accuracy of the GPS-receiver used, different data acquisition techniques (e.g., digitizing) or subjective knowledge about the data gathering process. In order to assess the usability of VGI in varying cases of application the positional accuracy of the data to be used has thus to be thoroughly evaluated, because missing and imprecise data effect spatial modeling. For example, Burra et al. (2002) concludes that the well-known Moran's *I* autocorrelation statistic is affected by locational errors (e.g., through inaccurate geocoding), which may lead to false conclusions in the worst case.

Hence the main purpose of this research is the statistical analysis of the positional accuracy of three different data sources for the year 2011, namely OSM, TomTom (TT), and survey data (SD), for a well mapped medium size German city, comprising an enthusiastic and active OSM community. For this purpose, we have to make the crucial assumption that our official SD features the highest accuracy as it is based on precise surveying techniques (e.g., triangulation). Therefore, SD serve as the reference data, to which the other data sets are relatively evaluated. In this context it must be mentioned that, because routing being its main application, the primarily intention of the TT dataset is topological correctness and positional accuracy is just a secondary but of course still essential issue. Thus, comparisons can be, but must not be, biased.

2 Related Work

Research concerning different kinds of accuracies (e.g., positional or topological) of VGI has unfortunately not gained much interest yet. A first descriptive attempt was conducted by HAKLAY (2010) who analyzed the positional accuracy of OSM compared to commercial data (OS Meridian 2) for the United Kingdom. He analyzed the percentage of overlaps between both data vendors within a buffer distance, as proposed by GOODCHILD & HUNTER (1997). The methodology has been adapted by ZIELSTRA & ZIPF (2010) for Germany, comparing the completeness of OSM to TT. Both studies concluded that OSM is a viable alternative data source, but emphasize that there are some limitations in usage concerning its completeness in rural areas. LUDWIG et al. (2010) alluded to a related issue, criticizing the lack of specific attributes, like maximum speed limits and street names. Similarly, SCHMITZ et al. (2008) analyzed the routing capabilities of OpenRouteService, based on OSM. Because of topological errors (e.g., unconnected street segments) within the street graph, 3-5 percent of all routing requests were not executable. CHEN (2010) deals with topology correctness and completeness of digital maps through the integration of different usergenerated (OSM) and commercial data sources (NavTeq, TT), comparing the correspondence of road crossings. His findings, among others, clarify that NavTeq and TT have a higher topological similarity than TT and OSM. Furthermore, urban areas show a higher similarity than rural areas. HAGENAUER & HELBICH (2012) dedicated their empirical work to urban areas as well. They evaluated the use of OSM to map urban regions in Europe by means of machine learning approaches and concluded the fitness for use of OSM strongly depends on location. OVER et al. (2010) extended the range of application of OSM data to

the third dimension. In combination with free elevation data (SRTM), the usefulness for 3D visualizations (e.g., buildings) is shown. Further research has addressed VGI for the purpose of geocoding (AMELUNXEN 2010). He concluded that the positional accuracy of geocoding results based on OSM data is equal to or even better than the accuracy provided by the commercial geocoding service offered by Google Maps. Nevertheless, these accuracy levels could only be achieved when OSM data were available on house number level which, at the time of the research, had been the case for only about 5 percent of the sample requests within the study area, but is increasing fast. Finally, NEIS et al. (2012) analyzed the development of the German OSM total street network from 2007 to 2011 and predicted that the discrepancy between TomTom data and OSM will disappear in late 2012.

This brief literature review highlights, mostly in a descriptive fashion, some limitations as well as the potential of OSM data. Further, it clarifies the need of statistical analysis of the positional accuracy of OSM compared to proprietary geographical data, like SD and TT. The present research tackles this issue.

3 Methodology

3.1 Data Processing

Datasets containing line strings of road segments (center lines) from all three sources are semantically aligned and loaded into a PostgreSQL/PostGIS spatially enabled relational database. The spatial reference system of the survey data has been transformed from the German "GK3" (EPSG 31467) to WGS 84 in this process in order to provide a common reference system for all datasets (OSM and TT data had been present in WGS 84 already).

As a first preprocessing step for each dataset, separate road segments sharing the same street name are merged in order to provide a single line string for either street. The junctions within the datasets are then extracted by determining all point coordinates where exactly two distinct line strings cross each other. This approach admittedly rules out junctions where three or more streets cross but has been preferred for the sake of clarity.

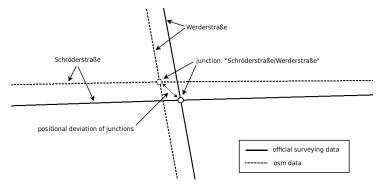


Fig. 1: Extraction and comparison of road junctions

The concatenated names of the streets crossing each other serve as an identifier for a given junction. These identifiers are then used to select and spatially compare corresponding

junctions among the datasets. As the identifier has to be unique, this approach additionally requires to rule out those cases where two streets cross each other more than once. Figure 1 illustrates this approach.

The deviation of the junction point coordinates from the corresponding points in the defined reference data set is then used as a measure of positional accuracy. Based on this, a scatter diagram of positional errors can investigated to inspect their spatial distribution and in order to detect potential systematic errors.

3.2 Geometrical Evaluation of the Distortion

To get insights of this geometrical distortion of the point patterns a global bidimensional regression is calculated (TOBLER 1994, FRIEDMAN & KOHLER 2003). This method allows assessing the transformation parameters between two (plain) maps and point patterns, respectively. Contrary to TOBLER'S (1965) remark, that bidimensional regression could be particularly useful for geographical analysis, it is rarely applied till these days. Primarily, spatial positional accuracies in cognitive maps are analyzed (LLOYD 1989). Other scope of applications are concerned with the lineage of historical maps (SYMINGTON et al. 2002), rubber sheeting as well as corrections of remote sensing images (TOBLER 1994).

The present research uses bidimensional regression models as descriptive statistics to determine the correspondence between OSM, TT, and SD. In Euclidean bidimensional regression the vectors of the regression equation are extended to be two-dimensional Cartesian coordinates pairs $(x_i, y_i; u_i, v_i)$, where x_i, y_i are the estimated coordinates from the of OSM and TT data, respectively, and u_i, v_i are the associated dependent reference coordinates of the surveying data. Scaling, translation, and rotation parameters reflect how the estimated point pattern must be transformed in order to fit back into the reference point pattern. This allows to quantify the geometrical relationship between two point patterns. The resulting bidimensional regression equation has following notation:

$$\begin{pmatrix} u_j \\ v_j \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} + \begin{pmatrix} b_{11}b_{12} \\ b_{21}b_{22} \end{pmatrix} \begin{pmatrix} x_j \\ y_j \end{pmatrix} + \begin{pmatrix} e_j \\ f_j \end{pmatrix}$$
 (1)

The parameters a_1 and a_2 correspond to the ordinary least squares (OLS) intercept term and carry out the translation. The b_{ij} values conduct the scaling as well as rotation and can be understood as the slope coefficient in OLS regression. Parameters e_j and f_j are the errors terms.

First, the magnitude of the horizontal (a_1) and vertical (a_2) translation between the reference pattern and the independent pattern is estimated, determining a least squares solution. A positive value of a_1 indicates a west-to-east shift and a negative value indicates an east-to-west shift. Likewise, positive values of a_2 are in accordance with a south-to-north shift and vice versa. Second, b_1 and b_2 are used to derive a scale parameter ϕ and angle parameter θ . Former causes the magnitude of contraction or expansion, whereas a ϕ value < 1 indicates a contraction and a ϕ value > 1 means an expansion relative to the reference pattern. The direction of the rotation necessary to get the best fit is determined by the angle parameter θ . A positive θ value indicates a counterclockwise rotation and a negative θ a clockwise one (LLOYD 1989, FRIEDMAN & KOHLER 2003). An overall "quality criterion" for the bidimensional regression is the Distortion Index (DI) introduced by WATERMAN & GORDON (1984) and extended in FRIEDMAN & KOHLER (2003). This index "can be thought

of as a standardized measure of relative error" (LLOYD 1989:110) and has a range between 0 and 100, where a lower value means less distortion.

3.3 Local Spatial Association of Positional Errors

Because bidimensional regression is a global statistic, it seems necessary to explore spatial heterogeneity in the positional errors as well. Therefore, an appealing method, among others, to detect local patterns of spatial association is the G^* -statistic (Getis & ORD 1992). Compared to the Moran's I, Burra et al. (2002) concludes that Getis-Ord G^* -statistic is more robust against spatial inaccuracies and thus preferable in our case. The G^* -statistic yields the proportion of the weighted sum of the variable within a distance d from location i as a proportion of the variable aggregated over the entire study region:

$$G_i^*(d) = \frac{\sum_{j=1}^n w_{ij}(d)x_j}{\sum_{j=1}^n x_j}$$
 (2)

where x_j correspond to the value of the observation at j, $w_{ij}(d)$ is the ij element of the spatial weight matrix and n is the number of observations. As a result spatial clusters of high and low values can be evaluated. In our case, a cluster of high values (z-scores) means a clustering of high positional errors and low values (z-scores) are related to an accumulation of low errors, always compared to SD. Significance is tested via a randomization approach.

4 Results

The preprocessing algorithm was able to extract 121 identical road junctions within our three datasets. The resulting point pattern is visualized in Figure 2. It can be seen that the junctions are spatially bounded to urban areas. Taking into account, that SD serve as a spatially precise reference dataset, the spatial deviation between SD and OSM and TT, respectively, was evaluated.

The mean deviation error has been found to be approximately one meter smaller in the OSM dataset, compared to TT (Table 1). A two sample Welch's t-test confirms significant differences between both mean values (t = -3.037, p = 0.003). The Fligner-Killeen-statistic is used to test homogeneity of both variances. The result clearly rejects the null hypotheses (FK = 57.644, p < 0.001) and it can be concluded that there are significant differences between the OSM and TT error variances. Moreover, OSM scatters more around the mean than TT, but comprising the directional scattering around the true position of the road junctions, as shown in Figure 3, it is noticeable that TT error clearly scatters more westward around the "true" position, than OSM does. The two varying mean centers of each point pattern support this finding and refer to a possible systematically variation.

Table 1: Descriptive statistics of the error deviation (in meters) between reference data and OSM and TT

	OSM	TT
Min.	0.220	2.759
Max.	18.694	13.607
Mean	5.229	6.145
Std. dev.	3.037	1.300

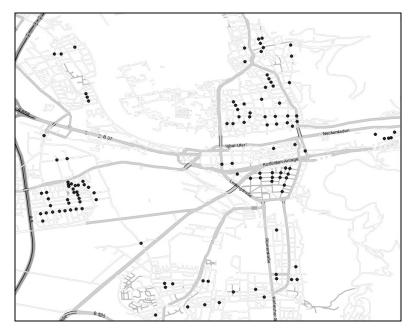


Fig. 2: Study site and identical road junctions (point signatures). Line signatures represent the road network.

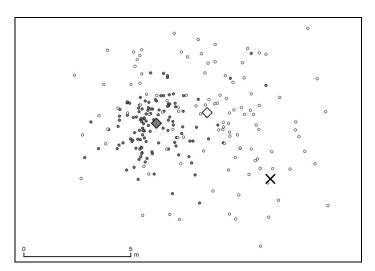


Fig. 3: Directional scattering around the "true" position of the road junctions. Darker points represents TT junctions, brighter ones are OSM junctions, the cross marks the "true" position, and the rectangles show the spatial means of the error distributions for TT (darker square) and OSM (white square).

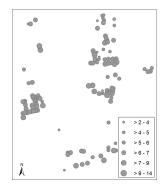
Next the geometrically distortion of the OSM and TT point pattern are compared to the one of SD on a global level, using the bidimensional regression framework. Thus, the OSM and TT pattern are regressed on the SD reference pattern, leading to the estimated parameters shown in Table 2. These parameters indicate how the OSM and TT pattern, respectively, must be transformed to get the SD pattern. Overall the OSM and TT pattern have the same geometrical distortion, hence having the same parameter signs. Both patterns are shifted east-to-west as well as south-to-north. Furthermore, the contraction or expansion parameter is negligible, because differences occur only after the fifth decimal place. θ refers to a clockwise rotation of the OSM and TT pattern, whereas OSM is marginally more rotated. The DI suggests that the relative error is slightly lower and thus the TT pattern corresponds more to the reference pattern. Nevertheless, the conclusions concerning the positional accuracy are not overwhelming and henceforth local statistics are used.

Table 2: Estimated Parameters (rounded) of the bidimensional regression (SD dependent variable, OSM or TT independent variable)

	a_{I}	a_2	b_I	b_2	φ	θ	DI
OSM	-35.844	0.690	1.000	-0.000	1.000	-0.011	0.178
TT	-2.561	19.351	1.000	-0.000	1.000	-0.009	0.095

Mapping the positional errors (Fig. 4) gives a first indication of spatial heterogeneity, but needs some statistical validation. Therefore, to explore areas with high and low accuracy, the G^* -statistic is calculated. We applied the zone of indifference option for conceptualization of the spatial relationships between the points, which is a combination of the inverse distance and fixed distance band model, leading to a neighborhood search threshold of 967 meters. Points with high z-scores and p-values below 0.05 indicate a spatial clustering of high positional errors (approx. beyond +/-2 standard deviations) and vice versa. Values between +/-2 standard deviations suggest no significant clustering.





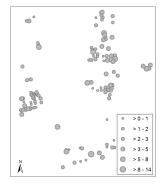
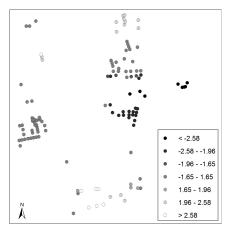


Fig. 4: Absolute deviation in meters between OSM and SD (left) and TT and SD (middle). Absolute value of deviation (in meters) between OSM and TT (right).



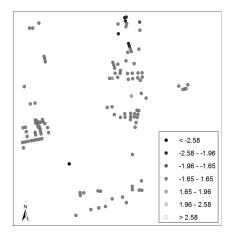


Fig. 5: Results of the G^* -statistics of OSM (left) and TT (right) in standard deviations. Values beyond +/-2 standard deviations have significant p-values (p<0.05).

Both maps in Figure 5 show some significant clusters of low values, corresponding to clusters with high positional accuracy. In the case of OSM (Fig. 5 left), this cluster is primarily located in the center of the map. 27 out of 121 observations have a significant negative z-score (p < 0.05). The opposite is valid for TT (Fig. 5 right), whereas these areas are located in the northern part of the map (7 significant observations). Positive values beyond 2 standard deviations are interpreted as badly mapped areas. In this regard, OSM shows some limitation, because such areas are present in the northern as well as southern part of the study site, whereas TT is not affected by limited position accuracy, compared to SD. Comparing the amount of such observations confirms this, OSM has 10 times more significantly imprecise mapped observations (OSM: 21, p < 0.05; TT: 2, p < 0.1). In general, the TT map gives a more homogeneous impression of the position accuracy errors.

5 Conclusion

The present paper is devoted to the comparison of positional accuracy of volunteered geographic information and proprietary geospatial data, using the case study of a German city. On the one hand bidimensional regression analysis is applied to evaluate the global geometries of the patterns and on the other hand clusters of high and low precision are detected by means of local autocorrelation statistics. The results show that both data sets, OSM and TT, have a highly positional accuracy and may be used for small and medium scale mapping applications. However, the bidimensional regression estimates shows highest correlation between OSM/TT and their true position, but the TT dataset has less distortion than OSM. The G^* -statistic results in some clusters with high and low positional accuracy, interpretable as spatial heterogeneity. Furthermore, the OSM areas of high accuracy are primarily located in the highly populated urban centers, leading to the conclusion that these areas are subject to a higher validation rate and consequently, errors are corrected more quickly than in rural areas. These findings are similar to those reported by Chen (2010), where urban areas have a higher (topological) accuracy. Hence, future comparisons between urban and

rural areas seems fruitful, because rural areas are mapped with significantly less completeness (ZIELSTRA & ZIPF 2010) but the tremendous growth of OSM data may shrink this disparity. There are, however, some limitations to this research. OSM as well as TT show similar spatial distortion, which raises the question whether the SD are affected by inaccuracy. Hence, future research is needed to get confidence, particularly other reference datasets and more case studies must be analyzed on different scales.

References

- AMELUNXEN, C. (2010), An approach to geocoding based on volunteered spatial data. In ZIPF, A. et al. (Eds.), Geoinformatik 2010. Die Welt im Netz, 7-12.
- Burra, T., Jerrett, M., Burnett, R. & Anderson, M. (2002), Conceptual and practical issues in the detection of local disease clusters: A study of mortality in Hamilton, Ontario. Canadian Geographer, 46, 160-171.
- CHEN H. (2010), Entwicklung von Verfahren zur Beurteilung und Verbesserung der Qualität von digitalen Karten. PhD Thesis, University of Stuttgart, Germany.
- FRIEDMAN, A. & KOHLER, B. (2003), Bidimensional regression: A method for assessing the configural similarity of cognitive maps and other two-dimensional data. Psychological Methods, 8, 468-491.
- ELWOOD, S. (2008), Volunteered geographic information: Future research directions motivated by critical, participatory, and feminist GIS. GeoJournal, 72, 173-183.
- GOODCHILD, M. F. (2007), Citizens as sensors: The world of volunteered geography. Geo-Journal, 69, 211-221.
- GOODCHILD, M. F. & HUNTER, G. J. (1997), A simple positional accuracy measure for linear features. International Journal of Geographical Information Science, 11, 299-306.
- GETIS, A. & ORD, J. K. (1992), The analysis of spatial association by use of distance statistics. Geographical Analysis, 24, 189-206.
- HAGENAUER, J. & HELBICH, M. (2012) Mining urban land-use patterns from volunteered geographic information by means of genetic algorithms and artificial neural networks. International Journal of Geographical Information Science, DOI:10.1080/13658816. 2011.619501 (online first).
- HAKLAY, M. (2010), How good is volunteered geographical information? A comparative study of OpenStreetMap and Ordnance Survey datasets. Environment and Planning B, 37, 682-703.
- LLOYD, R. (1989), Cognitive maps: Encoding and decoding information. Annals of the Association of American Geographers, 79, 101-124.
- NEIS, P. & ZIPF, A. (2008), OpenRouteService.org is three times "open": Combining OpenSource, OpenLS and OpenStreetMaps. GIS Research UK, Manchester.
- NEIS, P., SINGLER, P. & ZIPF, A. (2010) Collaborative mapping and emergency routing for disaster logistics – Case studies from the Haiti earthquake and the UN portal for Afrika. In CAR, A. et al. (Eds.), Geospatial crossroads @ GI_Forum 2010. Proceedings of the Geoinformatics Forum Salzburg, 2010, 239-248.
- NEIS, P., ZIELSTRA, D. & ZIPF, A. (2012), The street network evolution of crowdsourced maps: OpenStreetMap in Germany 2007-2011. Future Internet 2012, 4, 1-21.
- OpenStreetMap (2010), The free wiki world map. http://www.openstreetmap.org/ (last date accessed Feb. 1st, 2012).

- OVER, M., SCHILLING, A., NEUBAUER, S. & ZIPF, A. (2010), Generating web-based 3D city models from OpenStreetMap: The current situation in Germany. Computers, Environment and Urban Systems, 36, 496-507.
- SCHMITZ, S., NEIS, P. & ZIPF A. (2008), New applications based on collaborative geodata The case of routing. XXVIII INCA International Congress on Collaborative Mapping and SpaceTechnology, Gandhinagar, Gujarat, India.
- STRUNCK, A. (2010), Raumzeitliche Qualitätsuntersuchung von OpenStreetMap. Master Thesis, University of Bonn, Germany.
- SUI, D. (2008), The wikification of GIS and its consequences: Or Angelina Jolie's new tattoo and the future of GIS. Computers, Environment and Urban Systems, 32, 1-5.
- SYMINGTON, A., CHARLTON, M. & BRUNSDON, C. (2002), Using bidimensional regression to explore map lineage. Computers, Environment and Urban Systems, 26, 201-218.
- TOBLER, W. (1965), Computation of the correspondence of geographical patterns. Papers in Regional Science Association, 15, 131-139.
- TOBLER, W. (1994), Bidimensional regression. Geographical Analysis, 26, 187-212.
- VAN OORT, P. (2006), Spatial data quality: From description to application. PhD Thesis, Wageningen University, The Netherlands.
- WATERMAN, S. & GORDON, D. (1984), A quantitative-comparative approach to the analysis of distortion in mental maps. The Professional Geographer, 36, 326-337.
- ZIELSTRA, D. & ZIPF, A. (2010) A comparative study of proprietary geodata and volunteered geographic information for Germany. 13th AGILE International Conference on Geographic Information Science. Guimaraes, Portugal.