# Validation of Origin-Destination Matrix by Open Data

Evgeny Mayorov
*Department of Transportation Organization and Management*
*Samara National Research University*
Samara, Russia
benjamin1437@mail.ru

Oleg Saprykin
*Department of Transportation Organization and Management*
*Samara National Research University*
Samara, Russia
saprykinon@ssau.ru

*Abstract*—**Many global factors can influence the socio-economic situation in the country. One such factor is the transport systems of cities. It depends on the state of the transport infrastructure of the city: business development, investment, quality of life, employment growth. However, in larger cities observed the deterioration of the situation with the transport systematic in connection with the increasing number of cars and other factors. To solve these problems resort to modelling. The dynamic modelling allows to determine the characteristics of the future project (capacity, number of road lanes, traffic light cycle, etc.) without high capital investments. Modeling of transport processes has been used in practice for a long time, but there is a problem in checking the compliance of the constructed model with the real transport situation. Existing methods, such as the manual method of collecting information about traffic flows, are expensive and labor-intensive. Authors propose the method of validation of origin-destination matrix by open data provided by Uber company.**

*Keywords—Transport process, origin-destination matrix, SUMO, gravity model, Uber Movement, transport simulation, OpenStreetMap, transport analysis zones.*

## I. INTRODUCTION

In connection with the active development and growth of the urban population, the problem of matching the transport infrastructure with the needs of citizens' mobility is acute. The number of residents of the city is growing as well as the size of the cities. Transport systems remain one of the main problems of modern megacities. According to statistics, residents of large cities every 365 days spend 3 days waiting in traffic jams, which is an average of about 16% of the total driving time due to the poor condition of urban road networks. These statistics show the existing problem of urban road networks and the need to address it. In this regard, there are several modern approaches. The first way is to improve the road network, its reconstruction [1]. Another approach to solving the problem is to create comfortable conditions in public transport to attract citizens. This makes sense, since the desire to transfer people from personal to public transport will certainly lead to the unloading of the street and road network of the city.

Transport systems and their research become more complex over time and require precise knowledge in this area. Currently, there are many studies devoted to this area. Speaking about the transport system, it is worth noting that it consists of many interrelated parts. That is why in the world there are a large number of scientific groups engaged in research in a certain area.

Most attention should be given to the construction of origin-destination matrices (OD), which are designed to contain information on population movements throughout the city and to be able to apply these data in modelling. The data-driven approach and the model-driven approach are the most popular approaches to constructing OD matrices.

Application of GPS data for scientific research in transport systems is very popular nowadays. Therefore, the data approach is often used because it is based on the fact that the model is created based on dynamic data of residents' movement throughout the city using GPS, social networks, etc. This approach is very modern due to the widespread use of GPS navigation in our lives. In addition, the installation of GPS-Navigator in cars is now very common, and the matrix OD can be formed from the data obtained from the GPS-device. It is worth noting that the use of this method significantly reduces the possibility of obtaining an inadequate model. This area of research is relevant in our time [2-6].

Information from social networks is another source of data on population movements. Today, Foursquare and Facebook are among the best for this purpose. This method is based on the processing of geolocation data included in messages and photos. Having processed a sufficiently large amount of data, it is possible to make statistics showing the intensity of population movement in certain areas of the city at certain intervals. In this area, a lot of research is carried out on the creation and use of new social networks [7-10].

The most accurate method of obtaining data on mobility of the inhabitants were Clastres [11]. The high accuracy of the method is provided by using the data obtained from the cellular communication of residents of the city. This method provides the most up-to-date data, but processing large sets of data from different operators is very complex. The difficulties of using this method include the fact that it may be limited by the level of legislation of some countries, as well as the disagreement of subscribers with the processing of their personal data.

The model-driven approach is also popular with researchers and is used when mobile data or GPS is not available. Filippovskiy [12], which uses a universal model for mobility and migration, use the gravity method to obtain data on population movements across the country. This method makes it possible to monitor trade flows both within and between countries, migration flows and long-distance telephone calls. The gravitational method was developed by W. Reilly [13]. The idea of constructing a model of

gravitational type is based on the universal law of gravity. Applying this law to the transport system, it turns out that the objects are points that generate and absorb traffic flows [14], and the mass of the object is the total volume of incoming and outgoing flow. It is worth noting that it is possible to replace the physical distance by any costs [15] associated with the movement.

The entropy model [16] is specific for the second law of thermodynamics and was proposed by physicists. The law states that any closed physical system tends to achieve a stable equilibrium, which is characterized by the maximum entropy of this system. The first application of the method was carried out to the transport processes of A. Wilson [17]. The principle of the entropy method is that the system of movement of residents on the road network of the city has a sufficiently large number of uncontrolled elements. Due to the fact that our study is devoted to urban transport, it is necessary to consider the transport system as closed.

Michael Ballmer compares two types of modeling in his work [18]: the micro-simulation model and the traditional destination model. The traditional distribution model is based on the use of od matrices to reflect the movement of large population flows [19]. The main advantage of the micro-simulation model is the ability to simulate the movement of each agent separately [20].

Modern transport modeling systems allow creating a visual representation of traffic flows. Currently there is a wide range of software for transport modeling, among which are more popular programs such as T7F/TSIS, SUMO [21], TRANSYT, VISUM, MATsim, Aimsun. They allow to carry out verification of transport infrastructure changes without special economic and time costs. In addition, a study of the construction of od matrices based on the results of microscopic modeling was carried out [22-24].

The main problem of modeling is the lack of data with which to compare the resulting model. The method of full-scale measurements is quite labor-intensive and takes a large amount of time. Recently, however, sources with open statistics have begun to emerge. For example, Uber has created its own statistical data system based on GIS – Uber Movement. Uber engineers took data from GPS trackers in a taxi of their company, depersonalized it, and then created a statistical database on the time spent on the movement of the city districts. These data can be used for comparative analysis of the experimental correspondence matrix.

The article is devoted to the development of methods and software for automated creation of transport microscopic models for any urban area as well as comparative analysis of the received origin-destination matrix based on open data of Uber Movement. The developed system allows getting up-to-date information from open Internet resources, process it and use it to create or update the transport model. As an example of an open resource, consider OpenStreetMap, which provides information about urban infrastructure collected by volunteers. Using OpenStreetMap information, we have identified supply and demand for transportation to build a microscopic transportation model. The proposed approach allows to obtain models with acceptable quality using only data from open sources.

## II. A METHOD FOR ACTIVITY CHAINS CONSTRUCTION

One of the main objectives of this research is to develop an automated method of activity chains construction (fig. 1).

The first step in the method of activity chains construction is to obtain a map of the city of Amsterdam (one of the cities that available in Uber Movements by this day) from an open source OpenStreetMap.

The further stage is the division of the city into transport analysis zones [25]. In this study, it is necessary to take the division of the city into transport analysis zones is the same as in Uber Movement, in order to further be able to conduct a comparative analysis (fig. 2).

In this case, the division of the city is as follows: the closer to the city center, in which the number of trips as much as possible, the smaller the size of transport areas (in some cases up to the size of the quarter), and the farther from the city center, closer to the suburban areas the size of transport areas more. According to this principle, the city is divided into 181 transport areas.
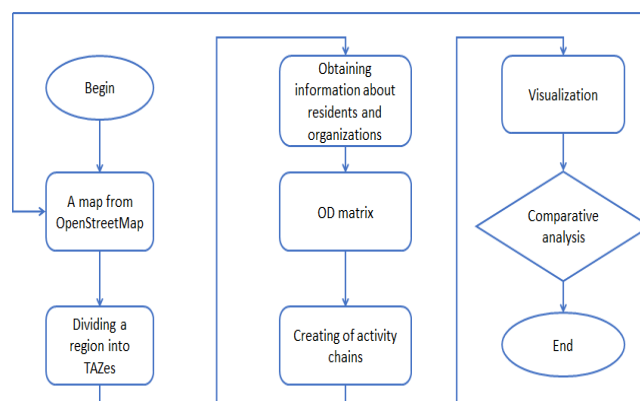


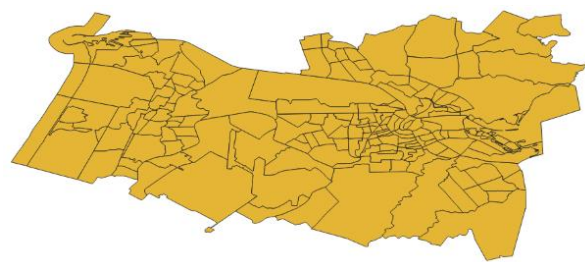Fig. 1. Schema of method for activity chains construction



Fig. 2. Dividing the city of Amsterdam by transport analysis zones

The next stage of the developed method of activity chains construction is to collect information about the number of residents of the city and assign them to a specific transport area. In this particular case, it was possible to obtain only General information on the number of inhabitants in the city of Amsterdam.

To divide the total population of the city by transport analysis zones, it was decided to use the coefficient method. This is because the size of the transport areas in this study is different and some areas are very different from the General trends, especially those far from the center. Therefore, it is necessary to assign a certain coefficient to such areas, which will equalize the size of the area and the number of inhabitants in it.

The coefficients were assigned to the administrative division of the city as some administrative areas have a large number of small areas, namely the center of Amsterdam – has 100 and Haarlem - 22 transport areas. It was therefore necessary to assign coefficients ranging from 0 to 2 to all transport areas so that the total number of inhabitants would remain the same. The smaller the transport area, the closer its coefficient is to the value 2, and vice versa, the larger the transport area, the value is closer to 0, but 0 will never be.

The next step is to build an origin-destination matrix. A gravity model is used to distribute correspondence by calculating cost matrices representing the cost of travel between each pair of zones. The transport gravity model relates to the intensity of the flow. Between the total number of departures from the *i*-th zone and arrivals to the *j*-th zone and the cost of travel between zones *i* and *j* (1). Applying the formula of the gravitational method was obtained origin-destination matrix which describes the movement of the population in transport analysis zones:

$$T_{ij} = \frac{Q_i * D_j}{c_{ij}^2} \, i = 1, \dots, N, j = 1, \dots, M, \qquad (1)$$

where *N* is the total number of departure zones, and *M* is the total number of arrival zones. In this model, the distance between areas is considered the distance between the centres of these areas.

However, for the possibility of further comparative analysis of semi-cluded matrix with Uber open data, the data about the intensity of vehicle should be converted to traffic time (here the time spent moving from one area to another):

$$t_{ij} = L_{ij} * \frac{p}{T_{ij}}, \iota, j \in R, \qquad (2)$$

where *t* is the time taken to move from one area to another

*L* is distance between transport areas;

*p* is average traffic density in the city.

The figure 3 shows the transformed origin-destination matrix in time characteristics.
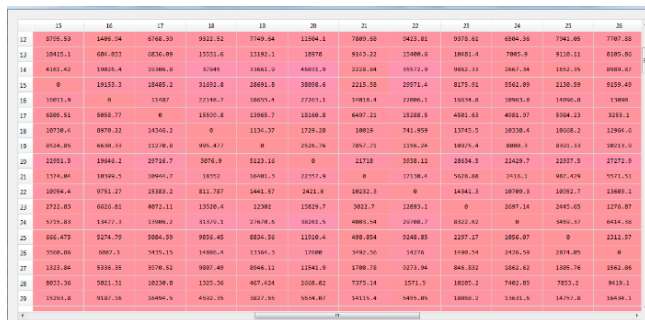


Fig. 3. Origin-destination matrix in time characteristics

The whole developed method of constructing a correspondence matrix is universal. In order to apply this method to any other city, you only need to download a new map of the city, determine the division of the city into transport areas, as well as add new information about the number of residents in the city.

The next step is to build an activity chain. The activity chain is one of the components of visualization in the SUMO simulation program. The activity chain is obtained by converting the origin-destination matrix into the format required for modeling. To transform a matrix, you must use applications SUMO – OD2TRIPS and DUAROUTER. These programs include all the necessary configurations for accurate modeling of transport processes. OD2TRIPS determines the start and end points of the path of each simulated correspondent based on the origin-destination matrix loaded into it, but it does not include the function of building a route along the roads. And DUAROUTER already directly takes the data obtained from OD2TRIPS and converts them in accordance with the roads laid down in the downloadable map of the city. Thus, we get the connected elements of the correspondence matrix and the map of the city of Amsterdam.

The final stage of creating a model of traffic flows of the city of Amsterdam is the visualization in the program of microscopic simulation of SUMO. To do this, you need to download the correspondence matrix processed in the OD2TRIPS and DUAROUTER programs, as well as the city map in xml format.

## III. COMPARATIVE ANALYSIS OF OD MATRIX WITH OPEN DATA

After obtaining an adequate model of transport processes in the city in the SUMO program, it is necessary to determine how the resulting model corresponds to the real situation on the roads of the city of Amsterdam.

To obtain information about the origin-destination matrix with open data on the time spent moving between transport areas taken from The Uber Movement website, it is necessary to compare them. Comparative analysis in this paper is performed using regression analysis.

TABLE I.     STATISTICAL DATA FROM UBER MOVEMENT RESOURCE

| source_id | dstid | hod | mean_travel_time | standard_deviation_travel_time | geometric_mean_travel_time | geometric_standard_deviation_travel_time |
|---|---|---|---|---|---|---|
| 143 | 141 | 7 | 263,89 | 233,27 | 200,94 | 2,2 |
| 144 | 131 | 7 | 527,19 | 513,77 | 394,14 | 2,04 |
| 61 | 108 | 7 | 472,58 | 250,77 | 425,29 | 1,56 |
| 156 | 16 | 12 | 1723,8 | 335,02 | 1691,98 | 1,21 |
| 153 | 46 | 12 | 1468,7 | 354,5 | 1425,79 | 1,28 |
| 74 | 127 | 6 | 980,88 | 488,84 | 889,75 | 1,52 |
| 38 | 171 | 1 | 1252,2 | 252,61 | 1230,38 | 1,2 |
| 42 | 27 | 18 | 814,19 | 314,18 | 763,15 | 1,42 |
| 40 | 47 | 18 | 1087,9 | 420,01 | 1018,93 | 1,43 |
| 86 | 109 | 7 | 537,3 | 302,23 | 448,22 | 1,89 |
| 53 | 87 | 0 | 519,95 | 194,15 | 489,4 | 1,4 |
| 59 | 27 | 0 | 594,47 | 243,35 | 554,03 | 1,46 |
| 57 | 47 | 0 | 724,27 | 211,11 | 698,47 | 1,3 |

From an open source Uber Movement were taken data on the average time of movement between transport areas for the third quarter of 2018. Table 1 shows the format of the downloaded data where: "sourceid" is the number of the departure area, "dstid" is the number of the arrival area, "mean_travel_time" is the average travel time between the departure area and the arrival area, the other elements of this table are not relevant to this study. Average time obtained by statistically anonymized data with GPS the taxi industry for Uber in Amsterdam.

Then it is necessary to load data on the areas of departure and arrival, and the average time spent on correspondence. These data, as well as the origin-destination matrix should be obtained in the form of a matrix shown in figure 4. A small disadvantage of the downloaded open data is the lack of some elements of the matrix, this is due to the lack of statistical data on the movement of population from certain areas in Uber specialists.



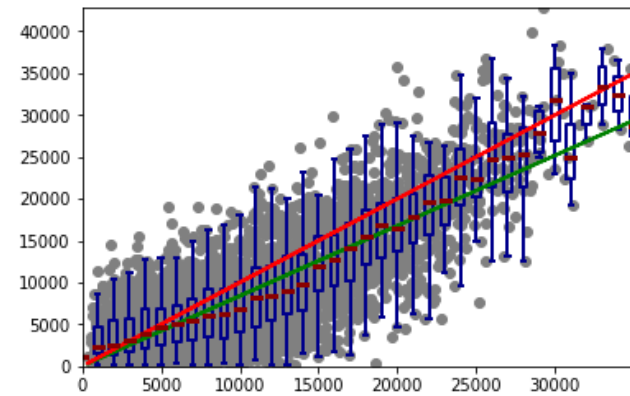Fig. 4.   Preprocessed Uber Movement data



Fig. 5.   Comparison modeled data with real data using regression model

The linear regression is used, the result of which will be a two-dimensional graph. The linear regression method was also calculated using the Python programming language.

Since this study uses only one comparison parameter, the linear regression formula will look like this:

$$Y = bX \qquad (3)$$

where $X$ is an independent variable, in this case it is a matrix obtained from an open source;

$Y$ is dependent variable, in this case it is a matrix of correspondence calculated by the gravitational method and transformed into a time matrix.

To obtain a more accurate model it is necessary to conduct machine learning by loading the test data X and Y. The test data were the last 200 values of each temporary matrix. After the test, a linear regression of all remaining data can be run. The obtained two-dimensional graph of the linear regression presented in figure 5. Coefficient of the model is 0.84. R-squared metric is equal to 0.6.

## IV. Conclusion

After a comparative analysis of the obtained model of the correspondence matrix with open data taken from Uber Movement, we can say that the resulting model is quite suitable for creating a simulation model of any city. The deviation from the ideal model is that there are no data on points of attraction, namely data on jobs, the number of students in schools and universities. When these data are obtained, it is possible to create an ideal model of any city.

Further stages of work on this project is the search for additional open data, as well as the creation of its own principle of dividing the city into transport areas, not based on data from Uber.

## References

[1]  "Bundesministeriumfür Verkehr. Bauund Wohnungswesen" (BMVBW).: Neubauvon Bundesautobahnen, 2004.

[2]  A.V. Gasnikov, S.L. Klenov, E.A. Nurminskiy, Y.A. Kholodov and N.B. Shamrai, "Vvedeniye v matematicheskoye modelirovanie transportnykh potokov," 2010, MFTI. (In Russian).

[3]  L. Moreira-Matias, "Time-evolving O-D matrix estimation using high-speed GPS data streams," Expert Systems With Applications, 2016, 44 pp.275–288.

[4]  J.S. Russell, M. Ye, B.D.O. Anderson, H. Hmam and P. Sarunic, "Cooperative Localisation of a GPS-Denied UAV in 3-Dimensional Space Using Direction of Arrival Measurements," IFAC, vol. 50, 2017, issue 1 pp.8019-8024.

[5]  S. Singh and P.B. Sujit, "Landmarks based path planning for UAVs in GPS-denied areas," IFAC, vol 49, 2016, issue 1 pp.396-400.

[6]  B. Kim, D. Kim, S. Park, Y. Jung and K. Yi, "Automated Complex Urban Driving based on Enhanced Environment Representation with GPS/map," Radar, Lidar and Vision, IFAC, vol. 49, 2016, issue 11 pp.190-195.

[7]  A.A. Kheir, "Intra-urban movement flow estimation using location based social networking data," The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, Volume XL-1/W5, 2015 International Conference on Sensors & Models in Remote Sensing & Photogrammetry, Kish Island, Iran.

[8]  Z. Zhang and Z. Wang, "The data-driven null models for information dissemination tree in social networks," Physica A: Statistical Mechanics and its Applications, 2017, pp.394-411.

[9]  C. Yu, B. Xiao, D. Yao, X. Ding and H. Jin, "Using check-in features to partition locations for individual users in location based social network," Information Fusion, 2017, pp.86-97.

[10]  Z. Jiao, L. Ran, J. Chen, H. Meng and C. Li, "Data-Driven Approach to Operation and Location Considering Range Anxiety of One-Way Electric Vehicles Sharing System," Energy Procedia, 2017, pp.2287-2294.

[11]  K. Friso, "Enriching the transport model of the Rotterdam region by cell phone data," Models and Technologies for Intelligent Transportation Systems (MT-ITS), 2015.

[12]  F. Simini, "A universal model for mobility and migration patterns," Macmillan Publishers Limited, 2012.

[13]  W.J. Reilly, "The law of retail gravitation," New York, 1931.

[14] R. Wiedemann, "Simulation des Straßenverkehrs flusses," PhD-thesis, University of Karlsruhe, Germany, 1974.

[15] M. Beckmann, C.B. McGuire and C.B. Winsten, "Studies in the economics of transportation," RM 1488. Santa Monica: RAND Corporation, 1955.

[16] S.C. Fang, J.R. Rajasekera and H-S.J. Tsao, "Entropy optimization and mathematical programming," Kluwer Academic Publisher, 1997.

[17] A.G. Wilson, "Entropy in urban and regional modelling," London: Pion, 1970.

[18] M. Balmer, "Generating Day Plans Based on Origin-Destination," 2005.

[19] T. Roughgarden and E. Tardos, "How bad is selfish routing," Journal of the ACM, 2002.

[20] D. Krajzewicz, J. Erdmann, M. Behrisch and L. Bieker, "Recent Development and Applications of SUMO – Simulation of Urban Mobility," International Journal On Advances in Systems and Measurements,2012.

[21] C. Hofer, G. Jäger, M. and Füllsack, "Generating Realistic Road Usage Information and Origin-Destination Data for Traffic Simulations: Augmenting Agent-Based Models with Network Techniques," In: C. Cherifi, H. Cherifi, M. Karsai and M. Musolesi (eds) Complex Networks & Their Applications VI. COMPLEX NETWORKS 2017. Studies in Computational Intelligence, 689, pp.1223-1233. Springer, Cham, 2018.

[22] O. Saprykin and O. Saprykina, "Multilevel Modelling of Urban Transport Infrastructure," In Proceedings of the 1st International Conference on Vehicle Technology and Intelligent Transport Systems (VEHITS-2015). Portugal, Lisbon: SCITEPRESS, 2015 pp.78-82.

[23] A.E. Gorev, K. Buttger, A.V. Prokhorov and R.R. Gizatullin, "Transport modeling fundamentals," St. Petersburg, Kosta, 2015, 168 p.

[24] O. Saprykin and O. Saprykina, "Validation of Transport Infrastructure Changes via Microscopic Simulation: A Case Study for the City of Samara," Russia, In Proceedings of the 5th IEEE International Conference on Models and Technologies for Intelligent Transportation Systems (MT-ITS 2017) in Naples, Italy, 26-28 June, 2017 pp.788-793.

[25] L.M. Martinez, J.M. Viegas and E.A. Silva, "A traffic analysis zone definition: a new methodology and algorithm," Springer Science+Business Media, LLC, 2009.