



HAL
open science

Reconstructing activity locations from zone-based trip data for discrete choice modeling

Milos Balac, Sebastian Hörl

► **To cite this version:**

Milos Balac, Sebastian Hörl. Reconstructing activity locations from zone-based trip data for discrete choice modeling. 101st Annual Meeting of the Transportation Research Board (TRB), Jan 2022, Washington D.C., United States. hal-03405572

HAL Id: hal-03405572

<https://hal.science/hal-03405572>

Submitted on 27 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Reconstructing activity locations from zone-based trip data for discrete choice modeling**

2

3 **Milos Balac**

4 Institute for Transport Planning and Systems, ETH Zürich

5 Stefano-Francini-Platz 5, 8093 Zürich, Switzerland

6 ORCID number: 0000-0002-6099-7442

7 milos.balac@ivt.baug.ethz.ch

8

9 **Sebastian Hörl**

10 Institut de Recherche Technologique SystemX

11 9 avenue de la Vauve, 91220 Palaiseau, France

12 ORCID number: 0000-0002-9018-432X

13 sebastian.horl@irt-systemx.fr

14

15

16 Word Count: 4555 words + 3 table(s) \times 250 = 5305 words

17

18 Submission Date: July 29, 2021

1 ABSTRACT

2 This paper presents a methodology to disaggregate activity locations from zone-based activity
3 chain data usually reported in the anonymized travel surveys. We propose an algorithm that aims
4 to find a feasible sequence of activity locations, for each individual, that minimizes the maximum
5 error of each trip's Euclidean distance within the activity chain. The reconstructed activity loca-
6 tions are then used to create unchosen alternatives within the choice set for each individual. This
7 is followed by the mode-choice model estimation. We test our approach on three large-scale travel
8 surveys conducted in Switzerland, Île-de-France and São Paulo. We find that with our approach we
9 can reconstruct activity locations that accurately match trip Euclidean distances, but with location
10 errors that still provide location protection. The models estimated on the reconstructed locations
11 perform similarly, in terms of goodness of fit and prediction, to the ones obtained on the original
12 activity locations.

13

14 *Keywords:* anonymization, data privacy, travel survey, choice model, discrete choice

1 INTRODUCTION

2 One of the most critical parts of transport planning is transport modeling. It should be able to
3 support transport planners in anticipating the impacts of policies and infrastructure projects. The
4 collection of various transport-related data supports transport modeling. While today information
5 can be collected through smartphone applications, transit tap-in/tap-out data, or mobile phone
6 data, the traditional approach is to utilize (household) travel surveys. These surveys, also referred
7 to as revealed preference (RP) surveys, usually collect detailed sociodemographic information
8 on individuals living in the area of interest together with their activity and trip behavior on one
9 or multiple days of the week. The activities can frequently be identified by a GPS coordinate
10 or detailed address. Typically, the gathered information on mobility behavior is enriched with
11 unchosen alternatives for each trip based on the choice set for each individual. This serves as a
12 preparatory step for further mode-choice modeling.

13 Due to privacy concerns and governing laws in many countries, the information in travel
14 surveys has to be anonymized at a level that protects the identity of individuals and their link to
15 the survey data. For this reason, identifying information like first and last name, home address, or
16 coordinates of activities are removed. The location of activities in publicly available versions of
17 surveys is usually published on a zonal level (i.e., traffic analysis zone, census zone). While this
18 protects the interviewed individuals, it is unknown how this aggregation affects the generation of
19 the unchosen alternatives, and subsequently, the modeling of the data and the forecasting power of
20 the created models. Therefore, in this paper, we aim to answer these questions.

21 The paper is organized as follows. Section 2 goes over the current literature in data
22 anonymization and its application to the field of transportation. Section 3 proposes a heuristic
23 to reconstruct activity locations based on zone-based trip data and explains the subsequently used
24 mode-choice modeling approach. Section 4 explains the used data sets, and Section 5 presents the
25 results. After, Sections 6 and 7 provide discussion and closing remarks.

26 BACKGROUND

27 With the increasing popularity of the open-data concept, the need to protect the privacy of indi-
28 viduals that provided their data has increased. One of the most usual pieces of information that
29 needs to be anonymized is location. Techniques used to provide location protection aim to obscure
30 the location of activities of individuals. Some of these techniques involve aggregation, spatial
31 cloaking, or random perturbation (for a detailed overview of different mechanisms, please refer to
32 (1)). A typical example is perturbation of residential locations of surveyed individuals, where the
33 anonymization procedure aims to maintain the usefulness of the data (2). The authors of (2) focus
34 on analyzing the performance of different perturbation mechanisms for protecting the privacy of
35 survey respondents. They also point out that current methods mainly deal with the anonymization
36 of single points and that further research is needed in developing methods for multi-point data.

37 Travel surveys that collect the mobility behavior of respondents over a day or week have
38 to deal with such multi-location data. Since each respondent reports multiple activities, a suitable
39 technique needs to be utilized that protects the privacy of individuals while still maintaining the
40 usefulness of the data. Most surveys utilize zone aggregation mechanisms (i.e., activity locations
41 are provided on a zone level). In the United States, each activity is usually aggregated to the census
42 tract (i.e., California Household Travel Survey (3), or My Daily Travel Survey conducted in the
43 Chicago Metropolitan Region (4)). In the case of France, multiple surveys exist. The publicly
44 accessible national survey has a high degree of aggregation on the level of departments, which

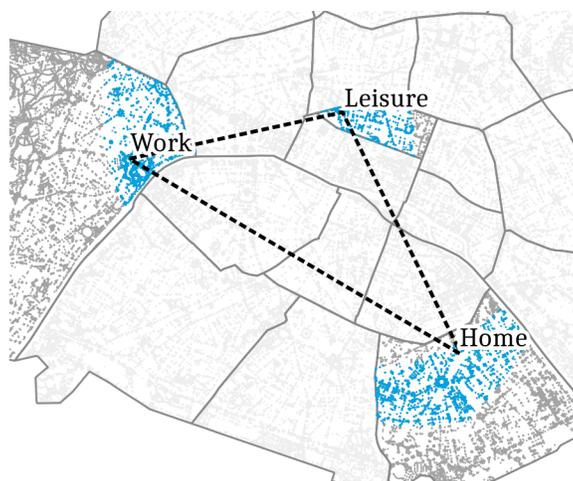


FIGURE 1: Example of a feasible set of candidate points

1 cover thousands or millions of residents. More local surveys, such as the one for the Île-de-France
 2 region around Paris, are only accessible on request and provide locations aggregated to a grid of
 3 100x100 meters. A commonly used aggregation level in French data sets are municipalities with
 4 thousands to tens of thousands inhabitants. In São Paulo, the publicly available travel survey does
 5 not provide location protection. In contrast, publicly available Brazilian census data is aggregated
 6 to a census zone containing between 20 and 55 thousand people.

7 Even when privacy protection techniques are used, confidential data can be at risk if addi-
 8 tional information obtained from other sources can uniquely identify individuals. For example, (5)
 9 show that mobile-phone traces provided in hourly intervals and with the spatial resolution provided
 10 by antennas can be uniquely identified in 95% of the cases with only four spatio-temporal points.
 11 (6) show that by revealing home and work census tract information, the anonymity set (i.e., the
 12 number of potential matching individuals) has a median size of 21 for the case of the U.S. working
 13 population. This raises a potential privacy concern for anonymized travel or commuting surveys.
 14 Nevertheless, identifying the level of privacy that the location protection techniques bring to the re-
 15 spondents in these surveys is not a direct aim of this paper, even though we provide some insights.
 16 We, however, aim to show how much the level of aggregation provided by the travel surveys could
 17 affect the prediction power of downstream models.

18 Therefore, to the best of our knowledge, we provide a first documented effort of the fol-
 19 lowing aspects:

- 20 • We propose a heuristic that, based on anonymized and aggregated zone-based trip data,
 21 creates disaggregated activity locations for all trips conducted by interviewed individuals.
- 22 • We perform analyses on the prediction accuracy of discrete choice models estimated on the
 23 basis of non-anonymized location information versus reconstructed locations.
- 24 • We show the universality of our findings based on survey data from three different countries.

25 **METHODOLOGY**

26 **Problem statement**

27 Figure 1 shows a motivating example for our approach. It shows an activity chain with four activ-
 28 ities, where a person starts the daily travels at home in the 13th arrondissement in Paris, then goes

1 to work close to the Eiffel tour which is located in the 16th arrondissement, continues to the Opera
 2 (2nd arrondissement) in the evening and then goes back home. In an anonymized travel survey, we
 3 may only know the Euclidean (and/or routed) distances between the activities, but also the zones in
 4 which the activities occur, represented by the arrondissements in this example. In dark gray, a set
 5 of possible activity locations in the zones has been obtained (here based on OpenStreetMap data).
 6 Furthermore, the Euclidean distances between all activities are known (exemplified by the dotted
 7 lines). If one now starts to move the locations of the four activities under the two conditions that (1)
 8 both “home” activities need to be at the same place, (2) Euclidean distances between the locations
 9 need to deviate no more than 50 meters from the reference distances, we arrive at a feasible set
 10 of locations which is colored in blue. The smaller the allowed deviation gets (e.g., 10 meters, 5
 11 meters), the smaller the feasible set of locations will become. Ideally, if our set of possible activity
 12 locations represents well the locations used in the survey, one would find the exact locations by
 13 reducing the deviation to zero.

14 **Location search problem**

15 The algorithm to find locations for the activities in a chain of a specific person is described in the
 16 following. As input, we know the number of activities in the chain N , as well as whether each of
 17 the activities $i \in \{1, \dots, N\}$ is a “home” activity. The indices of those activities are noted down in
 18 the index set \mathcal{H} . Furthermore, reference Euclidean distances are given as $r_i \in \mathbb{R}$.

19 The potential locations for the i th activity correspond to the potential locations in the re-
 20 spective zone. We denote the set of those locations as \mathcal{L}_i and the set of all potential locations in the
 21 activity chain is $\mathcal{L} = \mathcal{L}_1 \cap \dots \cap \mathcal{L}_N$. Let $k \in \{1, |\mathcal{L}|\}$ reference the elements of \mathcal{L} , then $y_{k,i}$ indicates
 22 whether location k is a potential location for the zone of activity i . The Euclidean distance between
 23 location k and k' is denoted as $d(k, k')$.

24 The aim of the algorithm is then to find a sequence $l = (l_1, \dots, l_N)$ with $l_i \in \mathcal{L}_i$ such that (1)
 25 the location for each activity is located in the respective zone, and (2) “home” activities always
 26 take place at the same location. To select among the feasible locations, the maximum deviation of
 27 the generated distances along the chain, compared to the reference distances, is minimized. The
 28 optimization problem is defined by the following objective function

$$\text{minimize}_{(l_1, \dots, l_N)} \max_{i \in \{1, \dots, N-1\}} \{ |d(l_i, l_{i+1}) - r_i| \} \quad (1)$$

29 with the following constraints:

$$\begin{aligned} y_{l_i, i} &= 1 & \forall i \in \{1, \dots, N\} \\ l_i &= l_{\min \mathcal{H}} & \forall i \in \mathcal{H} \end{aligned} \quad (2)$$

30 The first constraint makes sure that activities along the sequence only take place in locations
 31 that belong to the respective zone. The second constraint requires that all home activities take place
 32 at the same location.

33 **Solution strategy**

34 The solution strategy aims to find a feasible and optimal sequence (l_1, \dots, l_N) for each person.
 35 The most straightforward approach would use a depth-first branch-and-bound algorithm, where
 36 we would start a chain at any location in the first zone, then extend these chains with locations
 37 from the second zone and after with succeeding zones until one complete chain is found. The

1 maximum deviation along this chain can then be used to bound further exploration steps of the
 2 graph. Additionally, locations for home activities are set to the first occurrence of a home location
 3 along the constructed chain.

4 Our experiments have shown that such an approach causes very long run times if multiple
 5 times hundreds of potential locations need to be examined, especially for long activity chains.
 6 Hence, we perform a directed search where candidates in the following zones are chosen such
 7 that the local error is minimized. While the solutions of such an algorithm are not optimal, they
 8 perform well for the following modeling steps, as will be shown further below. Formally, the
 9 following depth-first branch-and-bound algorithm is proposed:

ALGORITHM 1: Chain-based location assignment

Input:

Location sets $\mathcal{L}_1, \dots, \mathcal{L}_N$ and \mathcal{L}

Home activity index set \mathcal{H}

Initialize:

$C = []$ $l^* = \emptyset$ $q^* = \infty$

For each $l_1 \in \mathcal{L}_1$

$C \leftarrow ((l_1), 0)$

Continue

While $|C| > 0$

$(l_1, \dots, l_n), q_n \leftarrow \text{pop } C$

If $q_n < q^*$ **Then**

If $n = N$ **Then**

$q^*, l^* = q_n, l$

Else

If $n \in \mathcal{H}$ and $n > \min \mathcal{H}$

$l_{n+1} = l_{\min \mathcal{H}}$

Else

$l_{n+1} = \arg \min_{l_u} \{|d(l_n, l_u) - r_i| \mid l_u \in \mathcal{L}_{n+1}\}$

End

$q_{n+1} = \max\{q_n, |d(l_n, l_{n+1}) - r_i|\}$

$C \leftarrow ((l_1, \dots, l_n, l_{n+1}), q_{n+1})$

End

End

Continue

Return l^*

10 Note that location sequences are only extended in a best-response fashion using the closest
 11 successor in terms of minimizing the Euclidean distance error, rather than enumerating all possible
 12 options. However, the algorithm can be easily modified to perform a complete enumeration if
 13 necessary.

1 Choice model

2 To test the impacts of location error on mode-choice model estimates, we make use of a straightfor-
 3 ward logistic regression model. We model the mode-choice for trips where car or public transport
 4 were a chosen mode. Therefore, the choice set includes only public transport and private car. To
 5 obtain relevant characteristics of the two alternatives, we perform a minimum cost path routing for
 6 all car trips, based on road networks obtained from OpenStreetMap data and free flow speeds. For
 7 public transport, we use an implementation of the RAPTOR algorithm (7) to find routes through
 8 the public transport network provided in GTFS format which minimize the total travel time of the
 9 trips. The data sets are documented in the scope of the development of synthetic populations for
 10 agent-based transport simulation for the three cases of São Paulo (8), Switzerland (9) and Île-de-
 11 France (10). As for some trips a public transport route cannot be found (i.e., the trip is too short, or
 12 public transport is not accessible), those trips are filtered out, which creates some differences in the
 13 size of the data set for reconstructed and original coordinates (see also Table 1). The mathematical
 14 formulation of the model is as follows:

$$\log\left(\frac{p_{car}}{1-p_{car}}\right) = \alpha +$$

$$\beta_{hascar} \cdot \delta_{car} + \beta_{haslicense} \cdot \delta_{license}$$

$$\beta_{invehicle,car} \cdot tt_{invehicle,car} + \beta_{invehicle,pt} \cdot tt_{invehicle,pt} +$$

$$\beta_{access,pt} \cdot tt_{access,pt} + \beta_{egress,pt} \cdot tt_{egress,pt} +$$

$$\beta_{transfer,pt} \cdot tt_{transfer,pt}$$
(3)

15 where p_{car} is the probability of choosing a car. All independent variables are continu-
 16 ous except δ_{car} and $\delta_{license}$, which are dummy variables representing whether a person has a car
 17 or driver's license, respectively. $tt_{invehicle,car}$ represents the travel time by car, and $tt_{invehicle,pt}$,
 18 $tt_{access,pt}$, $tt_{egress,pt}$, and $tt_{transfer,pt}$ represent the in-vehicle travel time, access time, egress time,
 19 and transfer time of public transport alternative. In São Paulo, driver's license information was not
 20 collected and, therefore, is not used in the models for São Paulo.

21 For each of the three case studies denoted by i , we estimate two models, one based on
 22 the original coordinates M_i^o and one based on the reconstructed coordinates M_i^r . To compare the
 23 predictive power of these two models, we split both data sets into a training set containing 70%
 24 (T_i^o and T_i^r) and a test set containing 30% (V_i^o and V_i^r) of the data by ensuring that the same trips
 25 are contained in both (i.e., T_i^o and T_i^r contain the same trips, but with different routing data). We
 26 train both M_i^o and M_i^r on the respective training set T_i^o and T_i^r . Finally, we analyze the predictive
 27 accuracy of the trained models on V_i^o data.

28 All models are estimated using the `scikit-learn` package in Python (11).

29 CASE STUDY

30 We make use of the already existing travel surveys from Switzerland (12), Île-de-France (13), and
 31 Greater São Paulo Metropolitan Region (14) to create the inputs for the reconstruction algorithm
 32 and the downstream mode-choice model estimation.

33 Switzerland

34 The *Mikrozensus Mobilität und Verkehr* (12) is a national travel survey conducted every five years
 35 in Switzerland. For the last edition conducted in 2015, about 56 000 persons ($\simeq 0.6\%$ of the total

1 Swiss population) are asked questions about their mobility behavior and their socio-demographic
2 attributes. Disaggregated, coordinate-level information about activities is available to the research
3 community upon request. The aggregated zonal information used in this study comes from the
4 National transport Model (15).

5 **Île-de-France**

6 The *Enquête globale de transport* (EGT, 13) is a household travel survey conducted in the Île-de-
7 France region, mainly during the year 2010. The EGT contains the trip chains of around 35 000
8 respondents in 15 000 households in the Île-de-France region. These numbers translate to a sample
9 of around 0.3% of people living in the region. Within Île-de-France, around 122 000 trips are
10 reported of all the members in each household. Unfortunately, EGT is only available on request
11 from the regional authorities and therefore not publicly available. Activity locations are reported
12 on a grid of 100x100 meters. As zoning data, French municipalities are used.

13 **São Paulo**

14 The last household travel survey in the Greater São Paulo Metropolitan Region was conducted in
15 2017 and is publicly available (16). It contains 84 889 weighted samples. For each sample, both
16 person and household-level information is provided. Unfortunately, no driver's license information
17 is available. Locations of activities performed by the respondents are reported with coordinate
18 accuracy. The dataset also provides a traffic zone for each of the activities, which are then used to
19 test the performance of the disaggregation algorithm.

20 **Candidates**

21 For the three cases, multiple sets of candidate points are created, among which the locations of the
22 activities can be chosen. Two different ways of generating such points are looked at.

23 First, we sample points at random for each zone in the three use cases. To do so, we obtain
24 the bounding box of each zone, sample N points within the bounding box, and then keep those
25 points that fall inside the zone boundaries. The number of points is defined as $N = A \cdot \eta$ with A
26 being the bounding box area and η a configurable density. In the experiments below, densities of
27 1, 5, 10, and 20 km^{-2} are used.

28 Second, we obtain OpenStreetMap data for each case. We filter for all road geometries that
29 are included or intersect with the case study area and use the *nodes* of the remaining road shapes
30 as location candidates.

31 **RESULTS**

32 **Reconstruction process**

33 First, the results of the reconstruction algorithm are presented. We examine the *distance errors*
34 and the *location errors* produced by the reconstruction algorithm. The *distance error* is defined
35 as the absolute difference between the Euclidean distance of a trip from the original data set and
36 the Euclidean distance between the selected location candidates. It is, hence, a measure of how
37 well the algorithm can recover the reference distances. The *location error* represents the distance
38 between an activity's location in the reference data set and its location. Therefore, it is a measure
39 of how well the algorithm reconstructs the original locations. Note that it is a validation measure,
40 as in the general case (with an anonymized data set), the original locations would not be available.

41 Figure 2 shows the cumulative distribution function of both error types for the three use

1 cases. In all cases, we observe that the distance error decreases strongly with an increased density
 2 of the location candidates, as more options allow a more fine-grained assignment. Furthermore, the
 3 OSM-based assignment performs the best in terms of reducing the distance error. For the location
 4 error, the same effects can be observed.

5 Interestingly, using the OSM candidates, the distance error is reduced to zero for almost all
 6 trips, i.e., point sequences that match the actual distances can be found in almost every case. The
 7 Euclidean distances are, hence, replicated almost perfectly.

8 The results on the location error are essential in terms of identifying specific activity lo-
 9 cations. Even with the high-density OSM data, locations can not be reconstructed perfectly. For
 10 Switzerland, however, 90% of activities are located within 1km of the original location. For Île-
 11 de-France and São Paulo, this threshold is reached at about 2km. On the contrary, more than 50%
 12 of locations in Switzerland can be reconstructed with an accuracy of 300m.

13 While Figure 2 gives a general impression on the matching performance of the algorithm,
 14 it is interesting to analyze how errors are distributed spatially. Figure 3 shows the location error,
 15 capped at 2km, for the three use cases. A high matching performance can be observed for Switzer-
 16 land for the finely zoned and highly populated areas around Zurich in the North and along the
 17 Geneva lake in the South-West. On the contrary, the sparsely populated and coarsely zoned areas
 18 in the Alps can be identified clearly as a strip of high location errors. For Île-de-France, errors are
 19 distributed somewhat randomly across space, especially no increase in accuracy can be observed
 20 for Paris and its metropolitan region, which would otherwise stick out in the center of the map.
 21 For São Paulo, accuracy is very low in the outer regions, where enormous zones contain large,
 22 unpopulated areas. Accuracy, however, increases towards the city center of São Paulo.

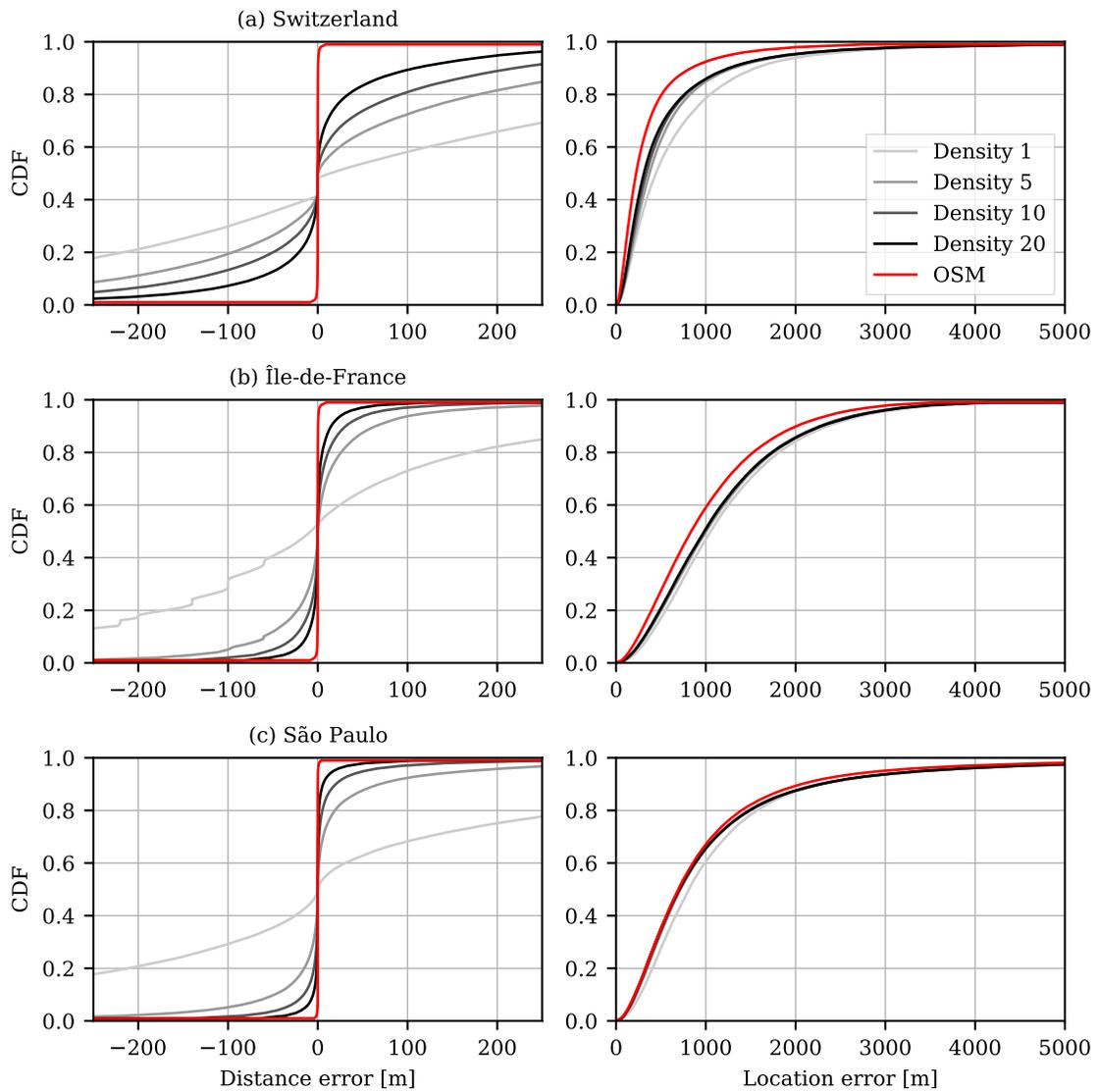
23 **Model estimation**

24 Table 1 presents the models estimated for different study areas and activity locations source. All
 25 parameters have the expected sign and are significant at 0.1% level. The parameters are in most
 26 cases very similar between models estimated on reconstructed and original activity locations. How-
 27 ever, some differences are observable, with the most prominent being for $\beta_{car, invehicle}$ and $\beta_{pt, waiting}$
 28 in Île-de-France.

29 Table 2 shows the prediction accuracy of the models, an evaluation mechanism that is fre-
 30 quently used in machine learning. For this measure, the systematic utilities of the two alternatives
 31 are calculated, and the better one is chosen. After, it is evaluated how many choices have been
 32 predicted correctly this way. Interestingly, both models perform similarly.

33 Prediction accuracy assumes that we have perfect knowledge of the individuals and their
 34 decision behavior. However, Train (17) argues that, given the taste variations within the population,
 35 it might be more suitable to compare mode-share predictions by sampling from the obtained choice
 36 probabilities. The results of this approach can be seen in Table 3. Both models predict mode-shares
 37 quite well. However, the models based on the original coordinates perform slightly better.

38 Figure 4 shows the car mode share in 1km distance bins for two models and the observed
 39 data. Once more, both models show similar patterns and forecasting quality. Towards longer
 40 distances, both models start to deviate from the observed mode-share. This could be accredited to
 41 the small number of observations for large distances leading higher likelihood of error.

**FIGURE 2:** Distance and location errors after the matching process

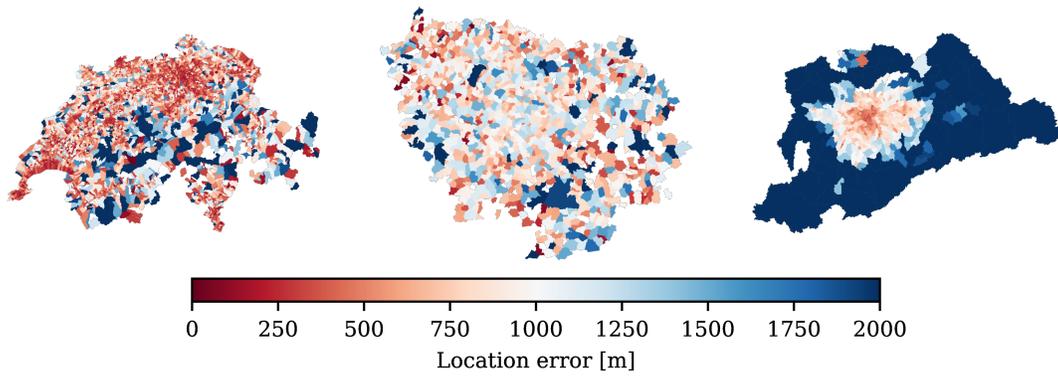


FIGURE 3: Spatial distribution of the location error (from left to right: Switzerland, Île-de-France, São Paulo)

TABLE 1: Models estimated for three study areas for original and reconstructed coordinates.

Parameter	Switzerland		Île-de-France		São Paulo	
	Rec.	Orig.	Rec.	Orig.	Rec.	Orig.
α	-2.954	-2.748	-7.989	-7.771	-2.368	-2.379
$\beta_{access,pt}$	[min^{-1}] 0.040	0.036	0.039	0.033	0.009	0.010
$\beta_{regress,pt}$	[min^{-1}] 0.042	0.040	0.039	0.037	0.008	0.008
$\beta_{waiting,pt}$	[min^{-1}] 0.035	0.045	0.040	0.058	0.012	0.013
$\beta_{invehicle,pt}$	[min^{-1}] 0.006	0.004	0.017	0.025	0.012	0.016
$\beta_{invehicle,car}$	[min^{-1}] -0.052	-0.052	-0.102	-0.123	-0.077	-0.089
β_{hascar}	1.781	1.757	4.691	4.720	2.933	2.933
$\beta_{license}$	2.642	2.592	4.588	4.491	-	-
Observations:	57589	59329	53869	55506	76867	77764
<i>Pseudo R-squared:</i>	0.339	0.338	0.411	0.422	0.205	0.208

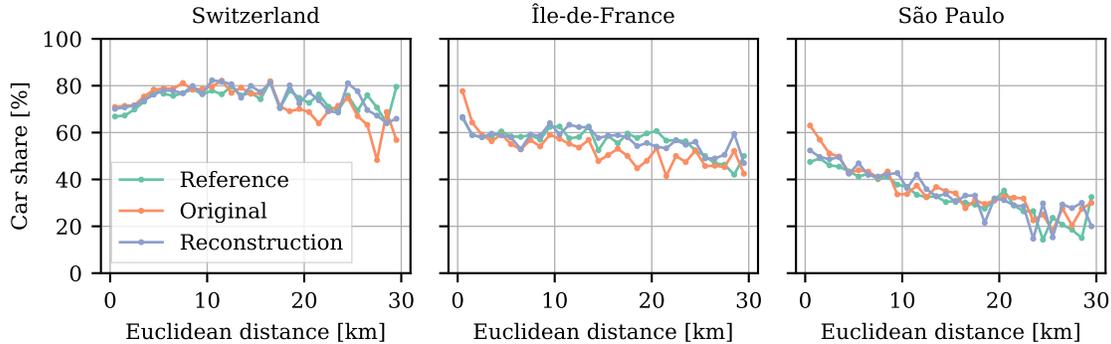
Note: All parameters are significant at the 0.1% level.

TABLE 2: Prediction accuracy of models estimated based on reconstructed vs. original locations

	Reconstructed	Original
Switzerland	0.854	0.853
Île-de-France	0.821	0.822
São Paulo	0.699	0.700

TABLE 3: Forecasted car mode share for models based on reconstructed vs. original location data in comparison to survey reference shares

	Reconstructed	Original	Reference
Switzerland	0.726	0.721	0.716
Île-de-France	0.576	0.570	0.572
São Paulo	0.414	0.409	0.409

**FIGURE 4:** Car mode-share in 1km distance bins for reference data from the surveys, and the models based on original and reconstructed locations

1 DISCUSSION

2 Based on the three data sets, the results show that the proposed location reconstruction algorithm
 3 generates activity locations that match Euclidean trip distances well. Furthermore, models based
 4 on reconstructed location provide a prediction quality very similar to the original data. While we
 5 only model the binary choice between a car and public transport, the results are promising. Future
 6 work should show if the models can still be estimated with a good fit when additional transport
 7 modes are added, or more complex models are estimated.

8 Some of the additional ways that the reconstruction of locations can be improved are:

- 9 • For trips made with public transport, origin or destination activity locations with reasonable
 10 access to public transport could be sampled within the zones. Consequently, unrealistic
 11 locations can be avoided, and higher location precision may be obtained.
- 12 • Currently, we only consider Euclidean distances between consecutive activities. Taking
 13 into account network distances could potentially improve the accuracy of the algorithm.
 14 Even (congested) network travel times could be used to reconstruct activity-to-activity travel
 15 times, if available.
- 16 • In the current approach, we extract all road nodes from the OSM network. In areas where
 17 OSM data has good quality, like in Switzerland or France, one could sample from potential
 18 locations based on the origin and destination activity. This way, possible locations for shop-
 19 ping activities would come from the location of shopping facilities present in OSM. More
 20 importantly, this could speed up the reconstruction algorithm. On the other hand, it could po-
 21 tentially increase the chances of precisely identifying activity locations of individuals, which
 22 would violate the anonymity requirement. If this is the case, suitable measures would need

- 1 to be taken to further anonymize the data.
- 2 • During location reconstruction, we only restrict home activities to happen at the same lo-
3 cation. Similarly, we could impose restrictions on education and work activities. However,
4 some individuals perform work activities in different places during the day. If this is the
5 case, we could identify this change in the activity chain by the change of the zone where the
6 work activity is performed.
 - 7 • Finally, from the location protection perspective, it would be interesting to investigate how
8 knowing the exact location of one of the activities would affect the knowledge about the
9 other activity locations in the chain, which would give insights on the potential vulnerability
10 of the data to outside attacks.

11 CONCLUSION

12 This paper demonstrates that discrete choice models estimated from disaggregated zone-based trip
13 data obtained with the proposed reconstruction methodology exhibit similar goodness of fit as those
14 based on non-anonymized data. These results are encouraging as they imply that by using spatial
15 cloaking on the level employed in the three datasets described for Switzerland, Île-de-France, and
16 São Paulo, the usefulness of the data sets for mode-choice modeling can be maintained. The
17 reconstruction algorithm presented in this paper can easily be applied to other data sets (such as
18 California Household Travel Survey (3)), which are spatially anonymized by default.

19 We observe that anonymity of individuals is not endangered by the methodology we em-
20 ploy. We have highlighted some essential future investigations that can help answer whether ad-
21 ditional data could potentially endanger the privacy of the surveyed individuals. As different en-
22 tities are increasingly collecting data from their users, the possibility to identify individuals from
23 anonymized surveys is increasing, which could have consequences on how future datasets should
24 be anonymized. Therefore, future work should focus on finding the potential weak points of cur-
25 rent anonymization techniques, especially when combined with other data sources, to inform on
26 potential risks and vulnerabilities.

REFERENCES

1. John Krumm. A survey of computational location privacy. *Personal and Ubiquitous Computing*, 13(6):391–399, 2009.
2. Godwin Badu-Marfo, Bilal Farooq, and Zachary Patterson. Perturbation methods for protection of sensitive location data: Smartphone travel survey case study. *Transportation Research Record*, 2673(12):244–255, 2019.
3. Transportation Secure Data Center. National Renewable Energy Laboratory. URL www.nrel.gov/tsdc. (Accessed: 20.07.2021).
4. Chicago Metropolitan Agency for Planning. My daily travel survey. URL <https://www.cmap.illinois.gov/data/transportation/travel-survey>. (Accessed: 20.07.2021).
5. Yves-Alexandre De Montjoye, César A Hidalgo, Michel Verleysen, and Vincent D Blondel. Unique in the crowd: The privacy bounds of human mobility. *Scientific reports*, 3(1): 1–5, 2013.
6. Philippe Golle and Kurt Partridge. On the anonymity of home/work location pairs. In Hideyuki Tokuda, Michael Beigl, Adrian Friday, A. J. Bernheim Brush, and Yoshito Tobe, editors, *Pervasive Computing*, pages 390–397, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
7. Daniel Delling, Thomas Pajor, and Renato F Werneck. Round-based public transit routing. *Transportation Science*, 49(3):591–604, 2015.
8. Aurore Sallard, Milos Balac, and Sebastian Hörl. Synthetic travel demand for the Greater São Paulo Metropolitan Region, based on open data. *Regional Studies, Regional Science*, In Press.
9. Christopher Tchervenkov, Aurore Sallard, Grace Kagho, Sebastian Hörl, Milos Balac, and Kay W. Axhausen. Synthetic travel demand for Switzerland. *Working Paper*, 2021.
10. Sebastian Hörl and Milos Balac. Synthetic population and travel demand for Paris and Île-de-France based on open and publicly available data. *Transportation Research Part C: Emerging Technologies*, 130:103291, 2021.
11. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
12. Swiss Federal Office of Statistics (BFS) and Federal Office for Spatial Development (ARE). Mikrozensus Mobilität und Verkehr, 2018. Neuchâtel.
13. Île-de-France Mobilités, OMNIL, and DRIEA. Enquête Globale Transport 2010, 2010.
14. Secretaria Estadual dos Transportes Metropolitanos and Companhia do Metropolitanano de São Paulo – METRÔ. Pesquisa Origem Destino 2017, 2019.
15. Bundesamt für Raumentwicklung (ARE). Modelletablierung Nationales Personenverkehrsmodell (NPVM) 2017, 2020.
16. Transportes Metropolitanos. Resultados finais da pesquisa origem e destino 2017 (final results of the 2017 origin-destination survey), 2017. URL <http://www.metro.sp.gov.br/pesquisa-od/>.
17. Kenneth E Train. *Discrete choice methods with simulation*. Cambridge university press, 2009.