

1 **A TWO-STAGE STOCHASTIC PROGRAMMING APPROACH FOR DYNAMIC OD**
2 **ESTIMATION**

3

4

5

6 **Qing-Long Lu**

7 Chair of Transportation Systems Engineering, Technical University of Munich, Germany

8 Email: qinglong.lu@tum.de

9

10 **Moeid Qurashi, Corresponding Author**

11 Chair of Transport Modeling and Simulation, Technical University of Dresden, Germany

12 Email: moeid.qurashi@tu-dresden.de

13

14 **Constantinos Antoniou, Ph.D.**

15 Chair of Transportation Systems Engineering, Technical University of Munich, Germany

16 Email: c.antoniou@tum.de

17

18

19 Word Count: 5770 words + 0 table(s) \times 250 = 5770 words

20

21

22

23

24

25

26 Submission Date: October 9, 2022

1 ABSTRACT

2 Estimating origin-destination (OD) demand is indispensable for urban transport management and
3 traffic control systems. While the existing estimation methods rely on data sources like house-
4 hold travel surveys and traffic network detection, they incur very high costs and are still either less
5 frequent or low in coverage density triggering lower observability and indeterminacy issues for
6 OD estimation. With ubiquity of smartphones, Location based social networks (LSBN) data has
7 emerged as a new rich data source with broad urban spatial and temporal coverage highly suitable
8 for OD estimation. However, thus far, most LSBN-based estimation models only focus on static
9 (day-level) OD estimation. This paper establishes a two-stage stochastic programming (TSSP)
10 framework integrating the activity chains to model activity-level mobility flows using LBSN data.
11 The first stage model aims to minimize the errors introduced by the inter-zone OD flows alongside
12 the expected errors of the check-in patterns. The second stage model attempts to minimize the
13 errors produced by the considered check-in pattern scenarios. A generalized Benders decomposi-
14 tion algorithm is presented to solve the two-stage stochastic programming model. We conduct the
15 experiments employing generalized least squares (GLS) estimator on the case study of Tokyo city.
16 The results depict that the algorithm convergence can be guaranteed within several steps. The al-
17 gorithm shows satisfactory performance in check-in pattern estimation, OD flows estimation, and
18 activity share estimation. Further, the implementation of the model in practical applications is also
19 specifically discussed.

20

21 *Keywords:* OD estimation, Demand estimation, Large networks, Stochastic programming, LBSN
22 data, Generalized Benders Decomposition

1 INTRODUCTION

2 An accurate origin-destination (OD) matrix, as a typical representative of mobility demand pattern,
3 is indispensable in urban transportation management and traffic control systems (1, 2). Integrating
4 with a traffic assignment model, it can be used to reproduce the traffic flow and network state in a
5 detailed manner, helping to design practical measures to improve the transport system efficiency.

6 Existing OD estimation methods primarily includes three travel data sources, i.e., tradi-
7 tional household surveys, traffic measurements, and positioning technology based data (3). The
8 first data source type of traditional household surveys are time-consuming, labor-intensive and
9 expensive, therefore are normally restrained within a limited area at low frequencies (e.g., once
10 or twice a decade). Similarly, OD estimation using the second data type of traffic measurements
11 also relies on dense detection infrastructure distributed over the network, requiring high cost for
12 both installation and maintenance. The traffic measurement based methods also structurally suffer
13 from the issue of indeterminateness in estimating realistic OD flows patterns (i.e., multiple sets of
14 varying OD matrix patterns can satisfy the constraints imposed by the traffic measurements and
15 optimize the system objective at the same time) due to the high dimensionality of the OD matrix
16 (4–6). Therefore, methods using the third data source type have attracted much attention in re-
17 cent years. The ubiquity of smartphones equipped with positioning technologies, such as GPS and
18 Bluetooth, has resulted in regular real-time generation of large sets of well-distributed data that
19 also provides unprecedented opportunities for the implementation and application of OD estima-
20 tion methods.

21 Generally, people travel for specific purposes, which enables the integration of travel pur-
22 poses and activity chains into OD estimation. Location based social networks (LBSN) data has
23 been used to develop such models, attributed to its broad urban spatial and temporal coverage and
24 confirmed trip purposes (7). LBSN services generate a large amount of anonymous check-in data
25 of venues and users, making it a natural “host” of urban mobility patterns (3, 8). More specifically,
26 check-in time series of venues record the travel destination distribution in both spatial and temporal
27 dimensions, while check-ins of users reflect the activity chains of individuals. It is thus probable to
28 develop an activity-based OD estimator to simulate the urban mobility using the patterns extracted
29 from LBSN data.

30 To this end, Yang et al. (3) proposed a singly constrained gravity model to estimate the
31 non-commuting OD flows using LBSN check-in data. The model was further improved by Jin
32 et al. (8) who replaced the singly constrained gravity model with a doubly constrained one to re-
33 duce the sampling bias of check-in data. The performance of other conventional trip distribution
34 models, such as radiation model, rank-based model, and population-weighted opportunities model,
35 calibrated with LBSN data have also been compared and evaluated in Kheiri et al. (9). However,
36 all these models can only provide a static (day-level) solution to the OD estimation problem. Ac-
37 cordingly, inspired by the promising performance of the application of the Hawkes process to
38 self-reinforcing behavior modeling in Cho et al. (10), Hu and Jin (7) presented a time-of-day zonal
39 arrival estimation model by integrating the Hawkes process and a LBSN check-in observation
40 model into a state-space modeling framework. Such an approach can reduce the sampling bias in
41 OD estimation caused by the difference between the social behaviors and the real travel patterns.
42 As per Jin et al. (8), the accuracy of trip arrival estimation is significant for the performance of OD
43 estimator using LBSN data.

44 There is still a significant scope left to explore the usage of LBSN data and construct a
45 dynamic OD estimator that can thoroughly utilize the trip purposes and activity chains informa-

1 tion. Therefore, as a step in that direction, this paper establishes a two-stage stochastic program-
2 ming (TSSP) framework integrating the activity chains to model activity-level mobility flows using
3 LBSN data. It is worth mentioning that, stochastic programming has been applied to optimize the
4 allocation of traffic sensors considering the uncertainty in the path flow distribution (11), and the
5 OD reconstruction problem based on traffic counts (12). However, to the best of our knowledge,
6 this is the first effort to apply it to model the dynamic OD estimation problem based on LBSN data.

7 In this study, we assume that similar check-in patterns are generated by the same OD flow
8 pattern, and these check-in patterns are treated as scenarios in the stochastic programming frame-
9 work. The first stage minimizes the errors introduced by the inter-zone flows (refer to OD flows
10 hereafter) alongside the expected errors of the check-in patterns. The second stage is to minimize
11 the errors produced by each check-in pattern scenario separately. Note that, a scenario is defined
12 as a realization of the second stage problem state, i.e., the check-in patterns. Finally, the proposed
13 two-stage stochastic programming model is addressed by the generalized Benders decomposition
14 (GBD) algorithm. The idea is to construct a master problem and a series of subproblems (one per
15 scenario) with respect to the first stage and second stage decision variables, respectively. These
16 problems are then optimized alternately and iteratively until the global optimum is found (11).

17 In the remainder of this paper, we first briefly introduce the LBSN check-in data. Then, the
18 mathematical model of the proposed OD estimator based on LBSN data is constructed, followed
19 with the solution algorithm – generalized Benders decomposition. Later on, case studies are
20 elaborated and model performance is evaluated. Finally, we draw some conclusions and suggest
21 future directions for research.

22 **LBSN CHECK-IN DATA DESCRIPTION**

23 This section provides a brief introduction on the characteristics of LBSN check-in data and relevant
24 concepts. An LBSN check-in event is automatically recorded when a user posts with geo-location
25 information or visits a venue (a point-of-interest). Each check-in is described by a user ID, a
26 venue ID, and the time of the check-in. In this regard, we can treat venues as detectors of such
27 events, while users are the objects or counts being detected. Overall, venues and users participate
28 in the services actively as venues can interact with customers in a creative and convenient manner
29 and customers can get awarded (e.g., discounts or "badges") from the social networking sites.
30 Compared to the conventional household surveys, such check-in data can be collected at a very
31 low cost with much higher frequency, and compared to the traffic measurements, detectors (i.e.,
32 venues) of check-in events are "deployed" much denser over the urban area.

33 Combining with the pre-registered location and category information of venues, the check-
34 in data has become a carrier of activity-oriented urban mobility patterns and can thus be used to
35 model the urban travel demand after appropriate aggregation. Normally, venue-side data and user-
36 side data are distinguished in the site server (3). Venue-side data contains the check-in statistics
37 with respect to the venue, while user-side data preserves the check-in history of the user. Conse-
38 quently, one can aggregate the venue-side check-in data based on the categorical hierarchy adopted
39 by the site to model the activity-based mobility flows. Likewise, the activity chains of individuals
40 can be extracted from user-side data. Inspired by this basic idea, in next section we develop a
41 mathematical model for OD estimation using LBSN check-in data which integrates the aggregated
42 check-in patterns of venue-side data and the activity chains extracted from user-side data.

1 METHODOLOGIES

2 Given that the OD patterns do not change dramatically within a short period without any disruptive
 3 events, we do a plausible assumption that similar check-in patterns at a specific time interval in
 4 different days during the reference period are generated by the same OD pattern. Similarly, it is also
 5 reasonable to say that OD flows generate when people travel for various activities across different
 6 regions within the network. In other words, OD flows are aggregated results of activity flows.
 7 Based on the said basic assumption and conceptual analysis, we developed an activity-oriented
 8 OD estimator in this section, leveraging the two-stage stochastic programming framework.

9 In particular, the OD estimator is built upon the graph model shown in Figure 1. For conve-
 10 nience, we define an *activity node* as an aggregating representative for a specific type of activities,
 11 e.g., “Food”. An *activity flow* is then the movements of people between two types of activities. For
 12 each traffic analysis zone (TAZ), we define a virtual *source* and a virtual *sink* to: (i) “memorize”
 13 the sum of in- and out-flows; (ii) counteract the noise in the check-in data collection; (iii) bridge the
 14 first-stage and the second-stage model decisions. Noteworthy, for a specific TAZ both the source
 15 node and the sink node are connecting to all activity nodes in the TAZ. The proposed approach is
 16 to optimize the OD pattern in the first-stage, fulfilling the specific constraints on OD flows and the
 17 constraints imposed by the expected cost from the second-stage problem. In the second stage, the
 18 check-in pattern will be optimized conditional on a specific problem state and OD pattern. Clearly,
 19 the optimal OD pattern and check-in pattern are interdependent.

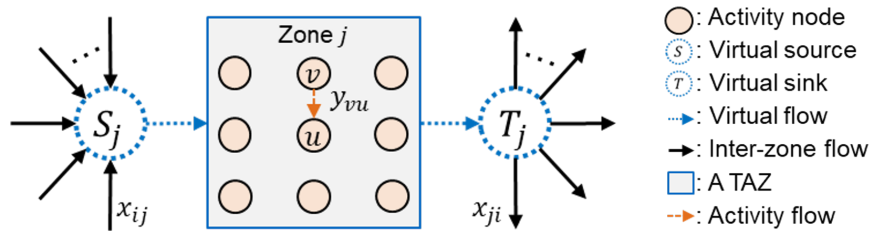


FIGURE 1: Graphical illustration of the model.

20 Two-Stage Stochastic Programming Model

21 The check-in patterns are the scenarios considered in the second stage, incurring by a common
 22 OD pattern. We formulate the generic two-stage stochastic programming model for OD estimation
 23 using LBSN check-in data as follows:

$$\min_{\mathbf{x}} f_1(\mathbf{x}, \mathbf{x}^{(p)}) + \kappa f_c(\Phi(\mathbf{x}), \Phi_c(\mathbf{c})) + \omega \mathbb{E}_{\xi} [Q(\mathbf{x}, \xi)] \quad (1)$$

$$\text{s.t. } \underline{\epsilon}_b x_{ij}^{(p)} \leq x_{ij} \leq \bar{\epsilon}_b x_{ij}^{(p)} \quad \forall i, j \in \mathbb{Z} \quad (2)$$

24

25 where \mathbf{x} is the decision variable of the first-stage problem, i.e., the vector of OD flows, $\mathbf{x}^{(p)}$ is
 26 the given prior OD flows. $f_1(\cdot)$ is the function measuring the difference between the estimated
 27 posterior OD flows and the prior OD flows. Similar to traffic measurement based OD estimators,
 28 the idea of including $f_1(\cdot)$ in the objective function is to help restrict the search space of the
 29 posterior OD flows. $\Phi(\mathbf{x})$ is the vector of out-flows of zones which is obtained by aggregating
 30 the estimated OD flows correspondingly, $\Phi_c(\mathbf{c})$ is the given out-flows estimated with the observed

1 check-in statistics \mathbf{c} . Note that $f_c(\cdot)$ measures the distance between the modeled and the measured
 2 out-flows. The addition of $f_c(\cdot)$ is inspired by the linear relationship observed from ten-month
 3 empirical LBSN check-in dataset. Figure 2 compares the observed out-flows and the out-flows
 4 estimated by a simple linear regression model based on the number of check-ins, i.e., $\Phi_c(\mathbf{c}) = \hat{\theta}^T \mathbf{c}$,
 5 where $\hat{\theta} = (C^T C)^{-1} C^T \Phi_0$, C is the matrix of the number of check-ins aggregated by TAZs, and
 6 Φ_0 is the vector of observed out-flows. This out-flows estimator is an input to the proposed OD
 7 estimator and will be adopted in the following experiments.

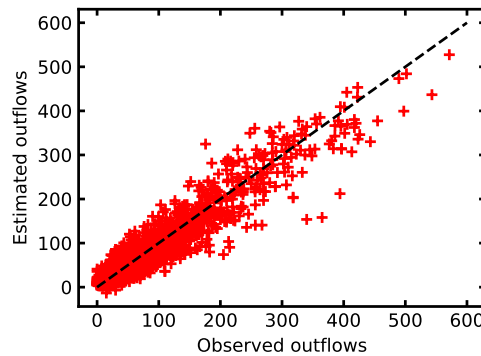


FIGURE 2: Linear relationship between out-flows and the number of check-ins.

8 While $f_1(\cdot)$ force the posterior OD flows to follow a similar OD pattern as the prior, $f_c(\cdot)$
 9 adjust the OD demand level based on check-in observations. κ is a weight factor that quantifies the
 10 relative reliability of the prior OD estimate and the prior out-flow estimate. \mathbb{Z} is the set of TAZs
 11 within the study area, and x_{ij} is the OD flow from TAZ i to j . Equation (2) represents the bound
 12 constraints on the OD flows, and bounds are defined as a multiple of the prior OD estimate. $\underline{\epsilon}_b$
 13 (< 1) and $\bar{\epsilon}_b$ (> 1) are threshold parameters.

14 Two-stage stochastic programming framework provides an opportunity for further restrict-
 15 ing the search space of the OD flows. This is achieved by introducing a batch of check-in pattern
 16 scenarios of the second-stage problem state, as expressed by the third term of Equation (1). \mathbb{E}_ξ
 17 calculates the expectation with respect to a random vector ξ , defined on the probability space
 18 $(\Omega, \mathcal{F}, \mathbb{P})$, with Ω being the sample space, \mathcal{F} being the event space, and \mathbb{P} being a probability
 19 distribution defined on \mathcal{F} . ξ is a random variable to describe the problem state at the second stage.
 20 $\mathbb{E}_\xi [Q(x, \xi)]$ is also called the recourse function. ω is a weight factor that quantifies the trade-off
 21 between the optimization of OD flow pattern and check-in patterns. $Q(\mathbf{x}, \xi)$ is the optimal value
 22 of the second stage problem given by

$$Q(\mathbf{x}, \xi) = \min_{\mathbf{y}} \sum_{z \in \mathbb{Z}} f_2(\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi)) \quad (3)$$

$$\text{s.t.} \quad \sum_{u \in \{\mathbb{V}_z - v\}} y_{vu,z} - \left(\sum_{u \in \{\mathbb{V}_z - v\}} y_{uv,z} + q_{v,z}^{\tau-1}(\xi) \right) \leq 0 \quad \forall v \in \mathbb{V}_z, \forall z \in \mathbb{Z} \quad (4)$$

$$(1 - \underline{\varepsilon}_a) \rho_{vu} q_{v,z}^{\tau-1}(\xi) \leq y_{vu,z} \leq (1 + \bar{\varepsilon}_a) \rho_{vu} q_{v,z}^{\tau-1}(\xi) \quad \forall v, u \in \mathbb{V}_z, \forall z \in \mathbb{Z} \quad (5)$$

$$y_{vu,z} \geq 0 \quad \forall v, u \in \mathbb{V}_z, \forall z \in \mathbb{Z} \quad (6)$$

$$(1 - \underline{\varepsilon}_s) \sum_{m \in \{\mathbb{Z} - z\}} x_{mz} \leq \sum_{v \in \mathbb{V}_z} y_{sv,z} \leq (1 + \bar{\varepsilon}_s) \sum_{m \in \{\mathbb{Z} - z\}} x_{mz} \quad \forall z \in \mathbb{Z} \quad (7)$$

$$(1 - \underline{\varepsilon}_t) \sum_{m \in \{\mathbb{Z} - z\}} x_{zm} \leq \sum_{v \in \mathbb{V}_z} y_{vt,z} \leq (1 + \bar{\varepsilon}_t) \sum_{m \in \{\mathbb{Z} - z\}} x_{zm} \quad \forall z \in \mathbb{Z} \quad (8)$$

1

2 where $\hat{\Delta}_z(\xi) = \mathbf{q}_z^{\tau-1}(\xi) - \mathbf{q}_z^{\tau}(\xi)$ is the ‘‘check-in pattern’’ of TAZ z in scenario ξ with \mathbf{q}_z^{τ} being
 3 the vector of the number of check-ins at different activity nodes in time interval τ , $\Delta_z(\mathbf{y}_z)$ is the
 4 estimated check-in pattern derived from the optimized activity flows \mathbf{y}_z . $f_2(\cdot)$ is a goodness-of-fit
 5 function measuring the distance between the observed and estimated check-in patterns. \mathbb{V}_z is the
 6 set of activity nodes in TAZ z . In practice, only the main venue categories in the TAZ will be
 7 selected for the sake of: (i) reducing the noise in the statistics caused by insufficient venues of a
 8 specific category; (ii) distinguishing different TAZs with respect to the land-use functionality and
 9 characteristics.

10

Equation (4) (denominated as inventory constraints) expresses that for a specific activity
 11 node v the sum of leaving flows cannot be greater than the sum of the coming flows and the number
 12 of check-ins recorded at the previous interval. Equation (5) (denominated as activity share con-
 13 straints) incorporates the activity chain information extracted from user-side data into the model,
 14 which is used to restrict the search space of activity flows and prevent the optimization from over-
 15 fitting issue to some extent. ρ_{vu} is the activity share of u in the flow out from v , which can be
 16 estimated from the historical check-in data. Note, the activity share can be aggregated at either the
 17 network level or TAZ level. All activity flows should be non-negative as expressed by Equation (6).

18

Moreover, recall that we define a source node and a sink node for each zone in the graph
 19 model. Thus, it is natural to have the in- and out-flow balance constraints on the source and sink
 20 nodes, which connect the decision variables of the first- and second-stage problems. However, due
 21 to the randomness and incompleteness of the activity information, we allow a subtle deviation in
 22 both constraints. The in-flow balance constraints are then given by Equation (7), representing that
 23 for a specific TAZ the sum of activity flows from the source node should not deviate too much
 24 from the sum of inter-zone flows to that TAZ. Similarly, the out-flow balance constraints are given
 25 by Equation (8), indicating that for a specific TAZ the sum of activity flows to the sink node should
 26 not deviate significantly from the sum of inter-zone flows from that TAZ.

27

$\underline{\varepsilon}_a, \bar{\varepsilon}_a, \underline{\varepsilon}_s, \bar{\varepsilon}_s, \underline{\varepsilon}_t$ and $\bar{\varepsilon}_t$ are predefined threshold parameters in the range (0,1). We note that
 28 the optimal activity flows \mathbf{y}^* depends on the first-stage OD pattern \mathbf{x} and the second-stage problem
 29 state ξ .

30

By comparison, the proposed LBSN data based OD estimator resembles the generic traf-
 31 fic measurement based OD estimator to a certain extent: (i) the proposed model also relies on
 32 a prior OD flow estimate; (ii) $f_c(\Phi(x), \Phi_c)$ in the objective function of the first-stage problem
 33 (Equation (1)) and the activity share constraints in Equation (5) play a similar role as the traffic

1 assignment method; (iii) both models need to handle the stochasticity of the observed data. The
 2 difference between the two method lies in that LBSN data are collected when the travel is finished
 3 but traffic measurements are collected during the travel. In other words, LBSN contains the end-
 4 to-end information, while traffic measurements reflects the situation between ends. As a result, the
 5 application of LBSN data based OD estimator usually demands no network structure, but needs
 6 the information of activity preference of travelers. More importantly, different from most traffic
 7 measurement based OD estimators in the existing literature, the proposed model can be used for
 8 dynamic OD estimation by only using the LBSN data for estimated check-in patterns without the
 9 need of otherwise running computationally expensive dynamic traffic simulation to generate sim-
 10 ulated traffic measurements. This puts the proposed methodology on a significant computational
 11 advantage against most dynamic OD estimation approaches.

12 **Sampling and Sample Average Approximation**

13 Apparently, the proposed model is very complicated and non-convex as the expectation \mathbb{E}_ξ is
 14 usually an integral of a complex function. Accordingly, in practice, we often need to assume ξ
 15 has a finite number of possible realizations with a known probability distribution, such that we can
 16 estimate \mathbb{E}_ξ by:

$$\mathbb{E}_\xi [Q(\mathbf{x}, \xi)] = \sum_n^N p_n f_2 (\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi_n)) \quad (9)$$

17 where N is the total number of realizations. Applying some sampling techniques and sample
 18 average approximation (SAA) method, the expectation can then be approximated by:

$$\mathbb{E}_\xi [Q(\mathbf{x}, \xi)] \approx \frac{1}{N_s} \sum_n^{N_s} f_2 (\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi_n)) \quad (10)$$

19 where N_s is the total number of scenario samples selected. In this paper, we apply the k -nearest
 20 neighbors algorithm (k -NN) to search for similar check-in patterns in the historical data to compose
 21 the set of check-in pattern scenarios of the OD pattern, in which each activity node pair is regarded
 22 as one dimension of the observation.

23 **Generalized Benders Decomposition Algorithm**

24 The generalized benders decomposition algorithm (GBD) was first proposed in Geoffrion (13) for
 25 addressing the mathematical programming problems with complicating variables (i.e., variables
 26 that if fixed to given values render a simple or decomposable problem). Obviously, in two-stage
 27 stochastic programming, the first-stage decision variables are the complicating variables of the
 28 problem. The idea behind GBD is of decomposing the original problem into a master problem
 29 and a series of subproblems (one per scenario). In the master problem, the first-stage decisions
 30 (i.e., OD flows, \mathbf{x}) are optimized. In the subproblems, the second-stage decisions (i.e., activity
 31 flows, \mathbf{y}) are optimized. They are solved iteratively until convergence. At a specific iteration k ,
 32 the subproblems are solved first separately resulting in the optimum $\mathbf{y}^k(\xi)$ given \mathbf{x}^{k-1} and scenario
 33 ξ . An optimality cut (or feasibility cut) is generated based on the dual solutions of subproblems
 34 (or feasibility problem), which is added to the master problem as a new constraint. Given all cut
 35 constraints created through a pass-back mechanism from subproblems in previous iterations, the
 36 master problem is solved with respect to \mathbf{x} resulting in \mathbf{x}^k . Note, these cuts gradually shrink the
 37
 38

1 feasible space of the complicating variable.

2 To describe the algorithm, we sequentially provide the formulations of the subproblem
3 (SP), the feasibility problem (FP), and the master problem (MP). At the k -th iteration, for a given
4 scenario ξ_n and \mathbf{x}^{k-1} , the SP is formulated as follows:

$$\min_{\mathbf{y}} \sum_{z \in \mathbb{Z}} f_2(\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi_n)) \quad (11)$$

$$\text{s.t. Constraints (4)-(8)} \quad (12)$$

$$\mathbf{x} = \mathbf{x}^{k-1} : \lambda_n^k \quad (13)$$

5

6 The solution of the SP provides values for the activity flows \mathbf{y}^k in different scenarios, as well as
7 the corresponding optimal Lagrange multipliers vector associated with Constraints (13), i.e., the
8 optimal dual variables vector, λ^k . Note, we can have N_s SPs solved in parallel. If the SP is feasible,
9 the Lagrangian function can be written as:

$$\mathcal{L}_o(\mathbf{x}, \mathbf{y}^k(\xi_n), \lambda_n^k) = \sum_{z \in \mathbb{Z}} f_2(\Delta_z(\mathbf{y}_z^k(\xi_n)), \hat{\Delta}_z(\xi_n)) + (\lambda_n^k)^T (\mathbf{x} - \mathbf{x}^{k-1}) \quad (14)$$

10

11 However, if the SP is infeasible, the following FP will be solved.

$$\min_{\mathbf{y}, \eta} \eta \quad (15)$$

$$\text{s.t. Constraints (4)-(8)} \quad (16)$$

$$\eta \geq 0 \quad (17)$$

$$\mathbf{x} - \mathbf{x}^{k-1} \leq \eta : \mu_n^k \quad (18)$$

12

13 Similarly, we can get the Lagrangian multiplier vector μ_n^k for Constraints (18). The Lagrangian
14 function of the FP is given by:

$$\mathcal{L}_f(\mathbf{x}, \mu_n^k) = (\mu_n^k)^T (\mathbf{x} - \mathbf{x}^{k-1} - \eta) \quad (19)$$

15

16 Then, the MP is formulated as follows:

$$\min_{\mathbf{x}, \alpha} f_1(\mathbf{x}, \mathbf{x}^{(p)}) + \kappa f_c(\Phi(\mathbf{x}), \Phi_c(\mathbf{c})) + \alpha \quad (20)$$

$$\text{s.t. } \underline{\epsilon}_b x_{ij}^{(p)} \leq x_{ij} \leq \bar{\epsilon}_b x_{ij}^{(p)} \quad \forall i, j \in \mathbb{Z} \quad (21)$$

$$\omega / N_s \sum_{n=1}^{N_s} \mathcal{L}_o(\mathbf{x}, \mathbf{y}^t(\xi_n), \lambda_n^t) \leq \alpha \quad \forall t \in \mathbb{I}_o \quad (22)$$

$$\mathcal{L}_f(\mathbf{x}, \mu_j^t) \leq 0 \quad \forall l \in \mathbb{S}_f^t, \forall t \in \mathbb{I}_f \quad (23)$$

17

18 where \mathbb{I}_o is the set of indices of the iterations at which all SPs are feasible, \mathbb{I}_f is the set of indices
19 of the iterations at which at least one of the SPs is infeasible, and \mathbb{S}_f^t is the set of scenarios whose
20 associated SPs are infeasible at iteration t . Constraints (22) are denominated as optimality cuts,
21 while Constraints (23) are feasibility cuts. From MP, we can get the values of the first-stage

1 decision variables \mathbf{x}^k .

2 If the original objective function is convex on the complicating variable, GBD can guar-
3 antee the strong optimality condition, i.e., the optimal solution from the decomposed problems is
4 equivalent to the original problem. For convenience, further details on the procedure of the GBD
algorithm are presented in Algorithm 1.

Algorithm 1 Generalized Benders decomposition algorithm for OD estimation

- 1: Initialize the OD flows \mathbf{x}_0 .
 - 2: Initialize the iteration index $k = 1$, the complicating variables $\mathbf{x}^k = \mathbf{x}_0$, error tolerance ε , the maximum number of iterations M .
 - 3: Set the lower bound of the objective function $\underline{z}^k = 0$, and the upper bound $\bar{z}^k = \infty$.
 - 4: **while** $|\bar{z}^k - \underline{z}^k|/|\underline{z}^k| \geq \varepsilon$ and $k < M$ **do**
 - 5: Set $k := k + 1$.
 - 6: Solve the subproblems by fixing \mathbf{x} as \mathbf{x}^{k-1} .
 - 7: **if** all subproblems are feasible **then**
 - 8: Obtain solution \mathbf{y}^k and the dual variables of those constraints that fix the complicating variables to given values λ^k .
 - 9: Calculate $z = f_1(\mathbf{x}^{k-1}, \mathbf{x}^{(p)}) + \kappa f_c(\Phi(\mathbf{x}^{k-1}), \Phi_c(\mathbf{c})) + \omega/N_s \sum_n^{N_s} \sum_z f_2(\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi_n))$.
 - 10: Update the upper bound $\bar{z}^k = \min\{\bar{z}^{k-1}, z\}$.
 - 11: Set $\mathbb{I}_o := \mathbb{I}_o \cup \{k\}$.
 - 12: **else**
 - 13: Solve the feasibility problems associated with the infeasible subproblems.
 - 14: Obtain solution \mathbf{y}^k , dual variable vector μ^k and the set of infeasible subproblems \mathbb{S}_f^k .
 - 15: Set $\mathbb{I}_f := \mathbb{I}_f \cup \{k\}$.
 - 16: **end if**
 - 17: Add the new optimality cut (or feasibility cuts) to the master problem.
 - 18: Solve the master problem to get \mathbf{x}^k and α^k .
 - 19: Update the lower bound $\underline{z}^k = f_1(\mathbf{x}^k, \mathbf{x}^{(p)}) + \kappa f_c(\Phi(\mathbf{x}^k), \Phi_c(\mathbf{c})) + \alpha^k$.
 - 20: **end while**
-

5

6 EXPERIMENTAL DESIGN

7 In this section, we verify the proposed OD estimator using the Foursquare check-in data. Foursquare
8 was launched in 2009 and has provided the leading LBSN service for more than a decade. As of
9 2016, it has more than 60 million registered users, with over 50 million monthly active. At the
10 same time, more than 95 million venues from over 190 countries or regions are registered on the
11 site. Their real-world images and consumer reviews are being updated constantly. It indicates that
12 the Foursquare data has a broad spatial coverage and therefore can somewhat capture the human
13 behavior in urban areas.

14 Case study setup

15 The Foursquare check-in data (*I4*) of Tokyo city, Japan, from April 2012 to February 2013 are
16 used in the following experiments. Figure 3a shows the map of the study area and the delineation
17 of TAZs. It is clear that the study area (1,302 km²) is divided into 17 TAZs. Figure 3b exhibits
18 a heatmap of 10,000 check-in records randomly sampled from the entire dataset, which contains

1 57,3703 records. The heatmap has a clear center and the color intensity gradually fades from the
 2 center outward. We note that TAZs are devised based on the density of check-ins for the sake of
 statistical significance, i.e., denser area has more TAZs.

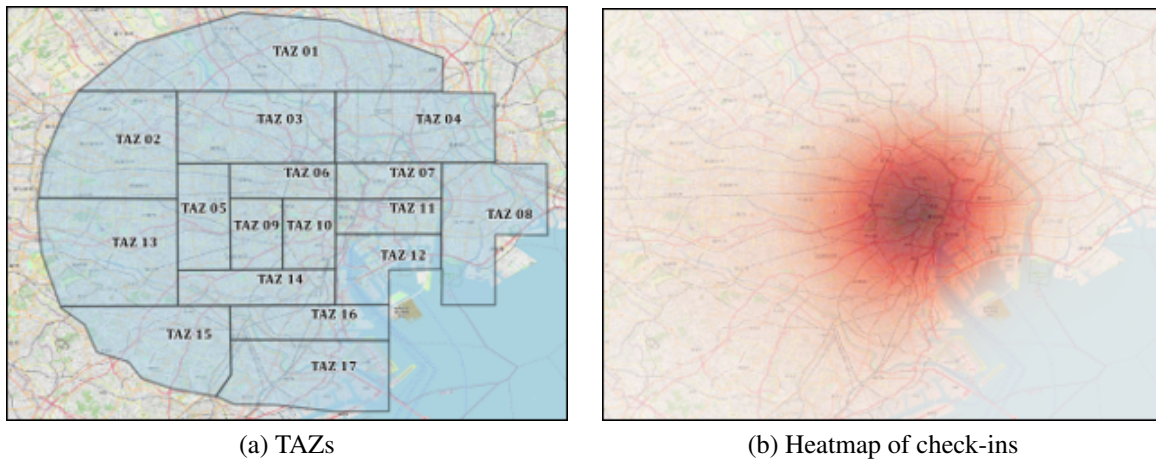


FIGURE 3: Study area: City of Tokyo.

3

4

Due to the lack of venue-side data, we aggregate user-side data hourly for each parent venue
 5 category¹, each TAZ, and each day to reconstruct the venue-side dataset. Categories with fewer
 6 than five check-ins are not further defined as an activity node of the TAZ. This can help identify the
 7 functionality of TAZs in land use and the generator of their attractiveness over time. For instance,
 8 the TAZ that has a huge number of “College & University” check-ins but a negligible quantity
 9 of “Professional & Other Places” check-ins is more likely an area including higher-educational
 10 institutions. Furthermore, we apply the moving average (seven days) technique to cancel the ran-
 11 domness of check-in behavior. In terms of the user-side dataset, for a specific time interval, we first
 12 extract the activity chain of each user. An activity share matrix can then be derived by counting
 13 the number of transfers between every two activities followed with normalization.

14 **Algorithm setup**

15

16

17

18

19

20

21

22

as follows:

$$f_1(\mathbf{x}, \mathbf{x}^{(p)}) = (\mathbf{x} - \mathbf{x}^{(p)})^T \Lambda_1 (\mathbf{x} - \mathbf{x}^{(p)}) \quad (24)$$

¹Foursquare leverages its own proprietary taxonomy of more than 1000 categories. According to the hierarchical taxonomy of categories (version 2012), ten parent categories are defined, including: Arts & Entertainment, College & University, Event, Food, Nightlife Spot, Outdoors & Recreation, Professional & Other Places, Residence, Shops & Service, Travel & Transport.

$$f_c(\Phi(\mathbf{x}), \Phi_c(\mathbf{c})) = (\Phi(\mathbf{x}) - \Phi_c(\mathbf{c}))^T \Lambda_c (\Phi(\mathbf{x}) - \Phi_c(\mathbf{c})) \quad (25)$$

$$f_2(\Delta_z(\mathbf{y}_z), \hat{\Delta}_z(\xi_n)) = (\Delta_z(\mathbf{y}_z) - \hat{\Delta}_z(\xi_n))^T \Lambda_2 (\Delta_z(\mathbf{y}_z) - \hat{\Delta}_z(\xi_n)) \quad (26)$$

1

2 where Λ_1, Λ_c and Λ_2 are the dispersion matrix of the prior OD estimate, of the out-flows distribu-
3 tion, and of the check-in pattern, respectively. For simplicity, we set $\Lambda_1 = \Lambda_c = \Lambda_2 = \text{diag}(\mathbf{1})$.

4 Demand scenario setup

5 To conduct the experiments, we chose the morning peak (7 am - 10 am) of February 1st, 2013,
6 divided in three estimation time intervals each for one hour. The entire check-in dataset is used for
7 estimating the relationship between the number of check-ins and the out-flows, i.e., $\hat{\theta}$, as afore-
8 mentioned. Further to generate the demand estimation scenarios, Antoniou et al. (5) points out that
9 the quality of the prior OD estimate, in terms of both demand level and patterns, is a key element
10 affecting the performance of the OD estimator. Following the suggestions therein, we perturb the
11 true OD flows to derive the historical OD flow estimates to be provided as inputs for the OD es-
12 timator. Due to space limitations, here we only test the performance of the proposed approach in
13 low-demand scenarios, i.e., the prior OD estimate is out-of-date and is lower than the true demand
14 level. More specifically, we create the prior OD estimate using the following equation:

$$\mathbf{x}^{(p)} = (0.7 + 0.3\delta)\mathbf{x} \quad \delta \sim \mathcal{N}(0, 1/3) \quad (27)$$

15 RESULTS

16 In this section, we first analyze the convergence performance of the GBD algorithm. Then, the
17 fit of the estimated OD flows, to the true OD demand, to the check-in pattern, and to the activity
18 share is presented. Finally, we show that the LBSN OD matrix can be easily scaled up to approxi-
19 mate the network OD matrix, illustrating the potential of the proposed OD estimator for practical
20 applications.

21 Algorithm Convergence Analysis

22 Figure 4 depicts the convergence results for estimating the OD matrices of the three experiment
23 intervals. Since the initial upper bound is infinity, it is not visible in the figure. Recalled that the
24 upper bound is updated by solving the subproblems, while the lower bound is updated by solving
25 the relaxed master problem. As expected, the algorithm converges within only a few iterations in
26 all three experiments, due to the convexity of the problem (GLS estimator). At each iteration, the
27 algorithm needs to solve the relaxed master problem once, and all the N_s subproblems in parallel.
28 Note, the relaxed master problem can be solved efficiently (in seconds) as it has limited number
29 of constraints and all constraints share the same format. We also found that subproblems are
30 always feasible if \mathbf{x}_0 is feasible, which means all Benders type cuts are optimality cuts, and thus
31 no feasibility problems needed to solve and no feasibility cuts are added to the master problem.
32 In consequence, the three experiments are solved in rather cheap computational efforts, indicating
33 that the proposed modeling framework has the potential for estimating the dynamic OD matrices
34 for large scale networks.

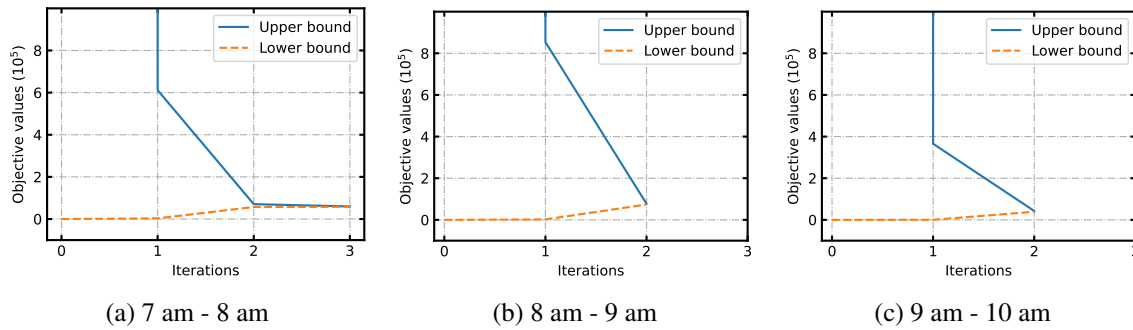


FIGURE 4: Convergence performance of the GBD algorithm.

1 Estimation Quality Evaluation

2 Figure 5 illustrates the quality of estimation with respect to the check-in patterns by comparing the
 3 empirical and estimated check-in pattern using 45° plots. Recall that check-in pattern is defined
 4 as the map of the difference of check-in counts between successive time stamps. We can see that
 5 all points are aligned closer to the “ $y = x$ ” line at interval 7 am - 8 am, confirming the capability
 6 of the model in recreating the check-in pattern. However, in intervals 8 am - 9 am and 9 am - 10
 7 am, though most points are also located near the 45° line, some deviate relatively further from the
 8 line. It is caused by the usage of k -NN in scenario selection as described in the previous section.
 9 k -NN would lead to biased results under the situation of limited candidate set. As a result, OD
 10 estimates that are promising to some check-in scenarios accidentally incur a biased estimate to
 11 some activity nodes in the others. The problem could be eliminated by incorporating the model
 12 with better sampling methods. Moreover, we can also see that some venues observe large negative
 13 check-in difference in the interval 7 am - 8 am. These venues may represent the residence places
 14 given that people are more likely to leave their home to work at this time. In all three intervals, we
 15 can find that most points are located in the range $[-20, 20]$, which represent the regular movements
 16 between different activities.

17 Similarly, Figure 6 visualizes the quality of the estimated OD flows by comparing it with
 18 the target OD flows using 45° plots. Overall, the model reaches to an acceptable estimate with
 19 a slight underestimation in the high demand OD pairs. We note that the prior OD estimate is
 20 also underestimated as 70% of the target values in average. Our model can somehow improve the
 21 situation attributed to the inclusion of $f_c(\cdot)$ and a batch of check-in pattern scenarios.

22 Further, Figure 7 compares the theoretic activity shares and the estimated activity shares.
 23 The difference has been restricted by the activity share constraints expressed in Equation (5). Due
 24 to the large number of points, we add the heatmap effect in the figure to represent the density of
 25 points. Brighter colors mean greater density, vice versa. Overall, there are more points in the range
 26 of smaller values. Since the constraints are defined based on percentage values, it is plausible
 27 to see the points with larger values are more scattered. In addition, we also note that more active
 28 movements (larger activity shares present) can be observed in 8 am - 9 am and 9 am - 10 am. Hence,
 29 together with the label of these activity flows, Figure 7 can provide useful auxiliary information for
 30 dynamic traffic management, and help venues design working schedules and “production” plans.

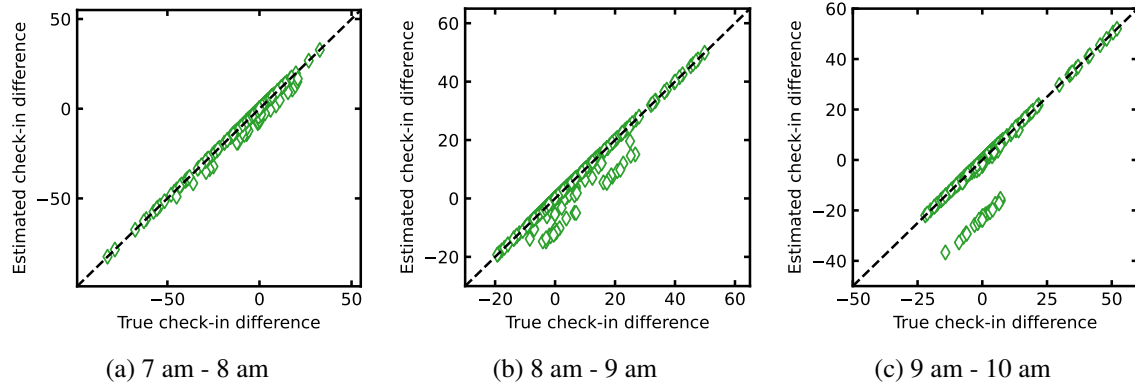


FIGURE 5: Comparison of true and estimated check-in patterns.

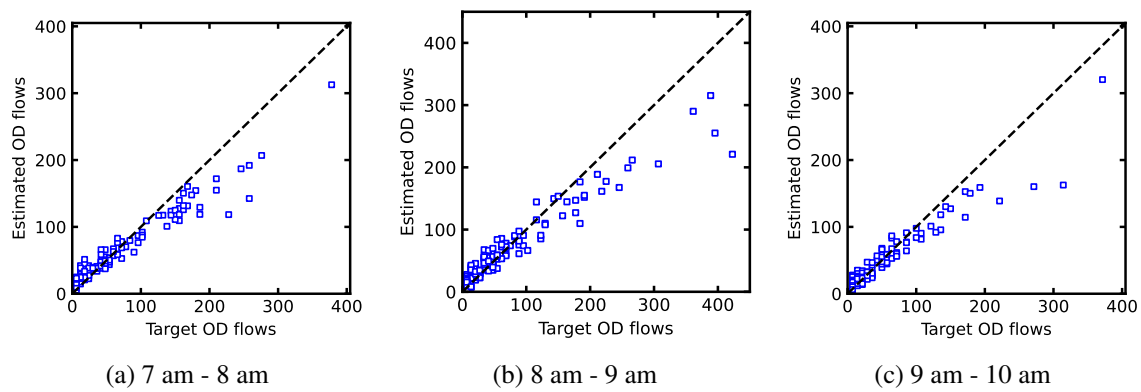


FIGURE 6: Comparison of target and estimated OD flows.

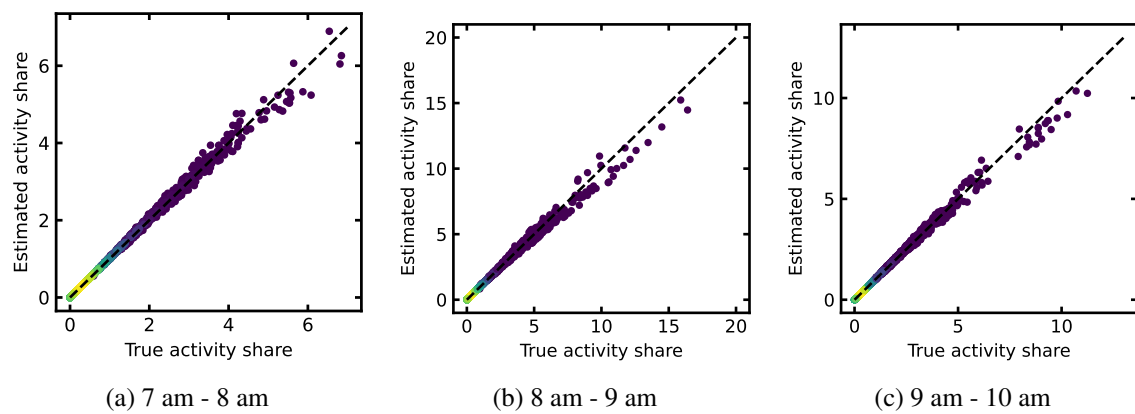


FIGURE 7: Comparison of true and estimated activity share.

1 **Scaling Towards Network OD Flows**

2 The OD flows compared in the previous sections only represent a partial set of the overall network
 3 OD flows that are observed in the LBSN check-in dataset (refer to LBSN OD matrix/flows here-
 4 after). Whilst, we usually need the actual network OD matrix in real practical applications instead.
 5 It is a critical input to the traffic simulation models for evaluating traffic management and policy
 6 measures (5). However, we argue that due to the random nature of check-in behaviors the LBSN
 7 OD matrix follows the same structural pattern as the network OD matrix, except for the deviation
 8 in the magnitude of demand level. Therefore, Figure 8 compares the network OD flows and the
 9 magnified LBSN OD flows using a similar method as in Figure 2. More specifically, the “true”
 10 network OD flows here are the average of TomTom OD flows of all Wednesdays in January 2021.
 11 The estimated network OD flows are obtained by scaling up the observed LBSN OD flows using a
 12 common scaling matrix for morning time and another matrix for afternoon time. Despite a rough
 13 estimation, the estimated ODs has a close pattern to the true ODs. It is reasonable to say that the
 14 Scatteredness of points is the joint result of the combination of the long distance in time (that the
 15 two datasets are collected) and the bias in data collection methods.

16 TomTom is a Dutch corporation launched in 1991 that specializes in the production of car
 17 navigation systems of all types. Many enterprises use TomTom’s positioning technologies such as
 18 Microsoft and Uber, due to their high precision. The TomTom data are collected from the probe
 19 vehicles which are equipped with TomTom’s positioning devices. Given the high penetration rate
 20 of TomTom positioning devices, we believe that TomTom data can capture the real traffic state to
 21 a large extent. Therefore, Figure 8 confirms that the LBSN OD flows estimated by the proposed
 22 model can approximate the real network OD flows after appropriate scaling. Theoretically, the
 23 scaling matrix can be either used after the LBSN OD estimation or directly integrated into the
 24 estimation framework. Also note that, exploration of more suitable scaling method which can
 25 incorporate demand information such as trip length frequency distributions to improve the network
 26 OD estimation are part of our future work.

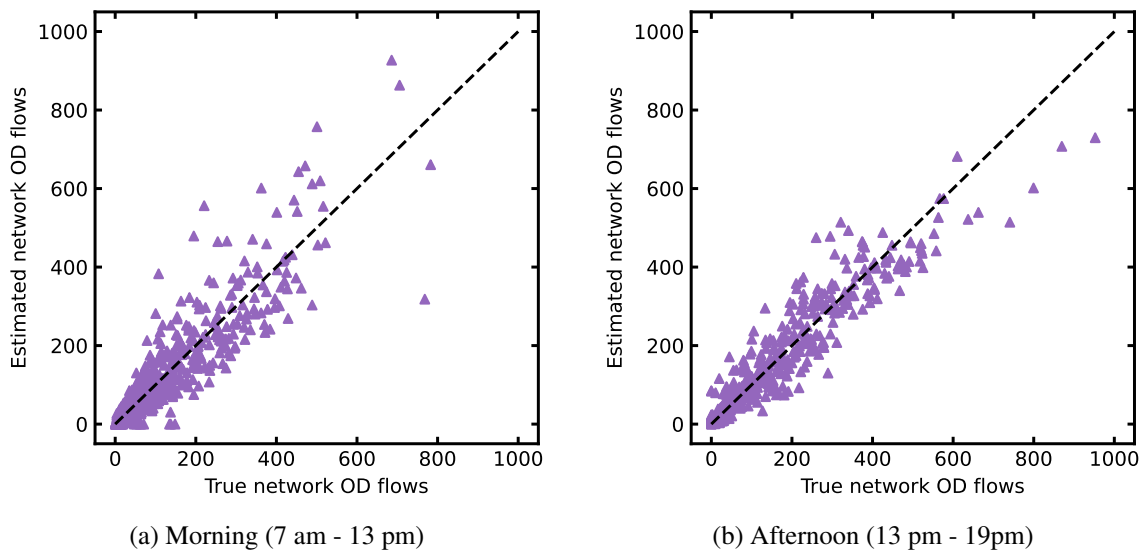


FIGURE 8: Comparison of true and estimated network OD flows.

1 CONCLUSION

2 Origin-destination (OD) estimation methods have long relied on the two traditional data sources of
3 household travel surveys and traffic network detection. On the one hand, travel surveys are time-
4 consuming and labor-intensive, restraining their coverage and frequencies, while on the other, the
5 traffic detection infrastructure is expensive to install and maintain, restraining its density and trig-
6 gering indeterminateness issues for OD estimation methods. Therefore, OD estimation methods
7 that utilize inexpensive and widespread data sources

8 Considering the stochastic nature of human behaviors and transportation systems, we pro-
9 pose a dynamic OD estimator by utilizing the scenario-based two-stage stochastic programming
10 framework, which integrates the activity chains extracted from LBSN check-in data to model
11 activity-level mobility flows. Given that OD flows are aggregated results of activity flows, the
12 OD matrix can be derived from these activity flows. Within the framework, the first stage model
13 aims to minimize the errors introduced by the inter-zone OD flows alongside the expected errors
14 of the check-in patterns. At the same time, the second stage model attempts to minimize the errors
15 produced by the considered check-in pattern scenarios. To solve the two-stage stochastic program-
16 ming model, a generalized Benders decomposition is presented, which seeks the optimal solution
17 by solving a relaxed master problem and a series of subproblems iteratively.

18 To evaluate the approach, we employ the case study of Tokyo city and use the generalized
19 least squares (GLS) estimator to measure the goodness of fit of both check-in and OD patterns. The
20 experiment results show that the convergence of the algorithm can be guaranteed within several
21 steps. Note, since our model is simulation free, the computational efficiency can be significantly
22 improved. More importantly, the model leads to a good fit for check-in patterns, OD flows, and
23 activity share distributions. Furthermore, we also present the method to scale up the LBSN OD
24 matrix to approximate the real network OD matrix, which can inspire the implementation of the
25 proposed model in practical applications.

26 Future directions for research would be to integrate appropriate sampling methods (e.g.,
27 importance sampling) into the proposed estimation framework for generating significant scenarios
28 instead of simply applying k -NN for scenario selection. On the other hand, embedding a suitable
29 scaling method into the model will also be useful.

30 ACKNOWLEDGEMENTS

31 This work was supported by the European Interest Group CONCERT-Japan DARUMA project
32 (Grant Number: 01DR21010) funded by the German Federal Ministry of Education and Re-
33 search (BMBF). This research was also partially funded by the German Research Foundation DFG
34 (TRAMPA Project, Grant 415208373).

1 REFERENCES

- 2 1. Ren, J. and Q. Xie, Efficient OD trip matrix prediction based on tensor decomposition.
3 *Proceedings - 18th IEEE International Conference on Mobile Data Management, MDM*
4 *2017*, 2017, pp. 180–185.
- 5 2. Xiong, X., K. Ozbay, L. Jin, and C. Feng, Dynamic Origin–Destination Matrix Prediction
6 with Line Graph Neural Networks and Kalman Filter. *Transportation Research Record*,
7 Vol. 2674, No. 8, 2020, pp. 491–503.
- 8 3. Yang, F., P. J. Jin, Y. Cheng, J. Zhang, and B. Ran, Origin-destination estimation for
9 non-commuting trips using location-based social networking data. *International Journal*
10 *of Sustainable Transportation*, Vol. 9, No. 8, 2015, pp. 551–564.
- 11 4. Cascetta, E., A. Papola, V. Marzano, F. Simonelli, and I. Vitiello, Quasi-dynamic estima-
12 tion of o–d flows from traffic counts: Formulation, statistical validation and performance
13 analysis on real data. *Transportation Research Part B: Methodological*, Vol. 55, 2013, pp.
14 171–187.
- 15 5. Antoniou, C., J. Barceló, M. Breen, M. Bullejos, J. Casas, E. Cipriani, B. Ciuffo, T. Djukic,
16 S. Hoogendoorn, V. Marzano, L. Montero, M. Nigro, J. Perarnau, V. Punzo, T. Toledo, and
17 H. van Lint, Towards a generic benchmarking platform for origin–destination flows estima-
18 tion/updating algorithms: Design, demonstration and validation. *Transportation Research*
19 *Part C: Emerging Technologies*, Vol. 66, 2016, pp. 79–98.
- 20 6. Qurashi, M., Q.-L. Lu, G. Cantelmo, and C. Antoniou, Dynamic demand estimation on
21 large scale networks using Principal Component Analysis: The case of non-existent or
22 irrelevant historical estimates. *Transportation Research Part C: Emerging Technologies*,
23 Vol. 136, 2022, p. 103504.
- 24 7. Hu, W. and P. J. Jin, An adaptive hawkes process formulation for estimating time-of-day
25 zonal trip arrivals with location-based social networking check-in data. *Transportation*
26 *Research Part C: Emerging Technologies*, Vol. 79, 2017, pp. 136–155.
- 27 8. Jin, P. J., M. Cebelak, F. Yang, J. Zhang, C. M. Walton, and B. Ran, Location-based social
28 networking data: exploration into use of doubly constrained gravity model for origin–
29 destination estimation. *Transportation Research Record*, Vol. 2430, No. 1, 2014, pp. 72–
30 82.
- 31 9. Kheiri, A., F. Karimipour, and M. Forghani, Intra-urban movement flow estimation us-
32 ing location based social networking data. *The International Archives of Photogrammetry,*
33 *Remote Sensing and Spatial Information Sciences*, Vol. 40, No. 1, 2015, p. 781.
- 34 10. Cho, Y.-S., G. Ver Steeg, and A. Galstyan, Where and Why Users" Check In". In *Proceed-*
35 *ings of the AAAI Conference on Artificial Intelligence*, 2014, Vol. 28.
- 36 11. Fu, C., N. Zhu, and S. Ma, A stochastic program approach for path reconstruction oriented
37 sensor location model. *Transportation Research Part B: Methodological*, Vol. 102, 2017,
38 pp. 210–237.
- 39 12. Jeong, I.-J. and D. Park, Stochastic programming approach for static origin–destination
40 matrix reconstruction problem. *Computers & Industrial Engineering*, Vol. 157, 2021, p.
41 107373.
- 42 13. Geoffrion, A. M., Generalized benders decomposition. *Journal of optimization theory and*
43 *applications*, Vol. 10, No. 4, 1972, pp. 237–260.

- 1 14. Yang, D., D. Zhang, V. W. Zheng, and Z. Yu, Modeling user activity preference by lever-
2 aging user spatial temporal characteristics in LBSNs. *IEEE Transactions on Systems, Man,*
3 *and Cybernetics: Systems*, Vol. 45, No. 1, 2014, pp. 129–142.