

## RESEARCH

# Estimation of flow trajectories in a multi-lines transportation network

Guillaume Guex<sup>1\*</sup>, Romain Loup<sup>2</sup> and François Bavaud<sup>1,2</sup>

\*Correspondence: gguex@unil.ch

<sup>1</sup>Department of Language and Information Sciences, University of Lausanne, Lausanne, Switzerland  
Full list of author information is available at the end of the article

## Abstract

Characterizing a public transportation network, such as an urban multi-lines bus network, requires the origin-destination trip counts during a given period. Yet, if automatic counting makes the embarkment (boarding) and disembarkment (alighting) counts at each bus stop known, it often happens that pedestrian transfers between stops are unknown, and this contribution proposes a three-steps procedure for estimating the missing information, involving maximum entropy and iterative fitting. \*\* à poursuivre \*\*\*

**Keywords:** multiline bus network; origin-destination flows; boarding and alighting counts; transit flows; maximum entropy estimation

## 1 Introduction

Transportation networks determine our mobility, require a considerable amount of planning and resources, and elicit much public hopes and critics. They also constitute an endless source of inspiration in formal modeling and optimization, as attested in operations research (classical optimal transportation, maximum flow problem), quantitative geography and spatial econometrics (spatial navigation, multimodality, gravity models for flows), and machine learning (recent developments in regularized optimal transportation, such as color transfer or images interpolation; see e.g. [1]).

This contribution addresses a straightforward, yet central question in public transportation networks: given a network made of many bus lines, how can one estimate the real trips made by the travelers, on the sole basis of the embarkment (boarding) counts and disembarkment (alighting) counts for each bus stop? Although estimating origin-destination flows is a much addressed issue in transportation modeling (see e.g [2] [3] [4] [5] and references therein), the specific problem addressed in this contribution seems, to the best of our knowledge, original.

Pedestrian transfers of travelers between bus lines here constitute the missing information, whose principled evaluation require some methodological reflexion and experimentation. Section 2 introduces the notations and the formalism, as well as the statement of the problem and the iterative solution method, which consist of three consecutive steps: a maximum-entropy computation of the trip distributions, obeying marginal constraints and with a given prior (section 2.4.1); an update of the marginal flows to avoid transfer overflow (section 2.4.3); and an update of the prior distribution (section 2.4.4) by shrinking the components responsible for overflow.

The first step only is required for solving the single line case (section 2.5), naturally much simpler but yet not trivial, and exhibiting a disembarking probability independent of the embarking stop (Markov property).

Cases studies are presented in section 3 \*\*\* à poursuivre \*\*\*

## 2 Notations and formalism

### 2.1 Lines, stops and junctions

Consider a *transportation network* made of *lines* numbered  $\ell = 1, \dots, q$ , of respective lengths (number of stops)  $l_\ell$ . Opposite lines, that is parallel lines running in the back and forth directions are considered as distinct.

The  $l = \sum_{\ell=1}^q l_\ell$  stops constitute the nodes of the transportation network. Each stop  $i = 1, \dots, l$  belongs to a single line, and defines a unique next or forward stop  $F(i)$  (unless  $i$  is the line terminus) and a unique backward stop  $B(i)$  (unless  $i$  is the line start), both on the same line.

Let  $S_i$  denotes the *set of stops which can be reached from stop  $i$  outside lines connection* (with, e.g., an acceptable walking distance), excluding  $i$  itself. A stop  $i$  is referred to as an *isolated stop* if  $S_i = \emptyset$ , and to as a *junction* otherwise.

### 2.2 Edges, trips, and the incidence matrix

Two sorts of oriented edges are involved in the transportation network:

- *intra-line edges*  $(i, j) = (i, F(i))$  belonging to a single line  $\ell(i) = \ell(j)$
- *inter-line or transfer edges*  $(i, j)$  connecting different lines  $\ell(i) \neq \ell(j)$ , involving walks from junction  $i$  to  $j \in S_i$ . The *set of transfer edges* is denoted by  $T$ .

A *st-trip*, noted  $[s, t]$ , consists of entering into the network at stop  $s$ , and leaving the network at  $t$ , by following the shortest-path (i.e. achieving the minimum distance, minimum time, or minimum cost), supposed unique, leading to  $s$  from  $t$ .

The succession of edges  $(i, j)$  belonging to the *st-trip*, noted  $(i, j) \in [s, t]$ , is unique. Define the *edge-trip incidence matrix* as

$$\chi_{ij}^s = \begin{cases} 1 & \text{if } (i, j) \in [s, t], \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Note that we can also forbid some aberrant trips across the network, for example, trips  $[s, t]$  where  $(s, t)$  is a transfer edge (making this trip do not actually use the line network). The *set of permitted trips* across the network is denoted by  $P$ , and can be defined regarding some conditions.

### 2.3 Transportation flows

Let  $x_{ij}$  count the *number of travelers using edge  $(i, j)$*  in a given period, such as a given hour, day, week or year. The edge flow  $x_{ij}$  is denoted by  $y_{ij}$  for an intra-line edge  $(i, j)$ , and by  $z_{ij}$  for a transfer edge  $(i, j)$ . By construction,  $x_{ij} = y_{ij} + z_{ij}$ , where  $y_{ij}z_{ij} = 0$ .

Let  $a_i$ , respectively  $b_i$ , the *number of passengers embarking*, respectively *disembarking* at stop  $i$ . By construction,

$$\begin{cases} y_{i, F(i)} = a_i \text{ and } b_i = 0 & \text{if } i \text{ is a line start,} \\ y_{B(i), i} = b_i \text{ and } a_i = 0 & \text{if } i \text{ is a line terminus,} \\ y_{i, F(i)} = y_{B(i), i} + a_i - b_i & \text{otherwise.} \end{cases} \quad (2)$$

Also,  $\mathbf{a}$  and  $\mathbf{b}$  must be consistent, in the sense that  $A_i \geq B_i$ , where  $A_i$  (respectively  $B_i$ ) is the *cumulated number of embarked (resp. disembarked) passengers* on the line under consideration, recursively defined as  $A_{F(i)} = A_i + a_i$  (resp.  $B_{F(i)} = B_i + b_i$ ). Moreover,  $A_i = B_i$  at a terminal line stop  $i$ . This common value yields the total number of passengers transported by the line.

Let the *transportation flow*  $n_{st}$  denotes the number of passengers following an  $st$ -trip, that is entering the network at  $s$  and leaving the network at  $t$  by using the shortest-path. One gets from (1)

$$x_{ij} = \sum_{st} \chi_{ij}^{st} n_{st} \quad (3)$$

Among the passengers embarking in  $i$ , some transfer from another line, and some others enter into the network:

$$a_i = z_{\bullet i} + n_{i\bullet} \quad (4)$$

where “ $\bullet$ ” denotes the summation over the replaced index, as in  $n_{i\bullet} = \sum_{j=1}^l n_{ij}$ . Similarly, among the passengers disembarking in  $i$ , some transfer to another line, and some others leave the network:

$$b_i = z_{i\bullet} + n_{\bullet i} \quad (5)$$

By construction

$$a_{\bullet} = b_{\bullet} = z_{\bullet\bullet} + n_{\bullet\bullet}$$

where  $n_{\bullet\bullet}$  counts the number of passengers, and  $z_{\bullet\bullet}$  counts the number of transfers.  $z_{\bullet\bullet}/n_{\bullet\bullet}$  is the average number of transfers per passenger.

As explained in section 2.1, transfers can only occur at junctions, that is  $z_{ij} > 0$  implies  $(i, j) \in T$ . In particular,  $z_{ii} = 0$  : no traveller is supposed to disembark and re-embark later at the same stop.

## 2.4 Statement of the problem and solution method

*Automatic passenger counters* measure the number of passengers entering and leaving lines at each stop [Boyle, 1998], that is  $\mathbf{a}$  and  $\mathbf{b}$ , which provide the basic raw data of the present study, kindly provided by the Lausanne Transportation Agency (tl) for the case study in section 3.3. We will suppose here that this data obeys the necessary consistency condition  $a_{\bullet}^{\ell} = b_{\bullet}^{\ell}$  (where the latter quantities denote the total embarkments and disembarkments on line  $\ell$ ), even if, in real case studies, a rescaling must usually be performed to balance in and out-flows on each lines.

Intra-line edge flows  $\mathbf{Y} = (y_{ij})$  can be determined by (2), but transfer edge flows  $\mathbf{Z} = (z_{ij})$  are, here and typically, unknown. The objective is to estimate the  $l \times l$  transportation flow matrix  $\mathbf{N} = (n_{st})$ . Many consistent solutions coexist in general, even for a single line with no transferts (section 2.5). This issue of incompletely observed data can be tackled by the maximum entropy formalism [6], which has often been the case in transportation modelling researches [7] [8].

Let  $f_{st} = n_{st}/n_{\bullet\bullet}$  be the *distribution of st-trips* (empirical distribution) and let  $g_{st}$  be some prior guess on its shape (theoretical distribution). Assuming some reasonable initial prior  $g_{st}$ ,

- (1) we shall first suppose that the empirical margins  $\alpha_s = f_{s\bullet}$  and  $\beta_t = f_{\bullet t}$  are known. Then  $f_{st}$  can be determined as the maximum entropy solution (section 2.4.1), i.e. as the distribution closest to  $g_{st}$  in the Kullback-Leibler divergence sense under the margin constraints, to be calibrated by an iterative fitting inner loop
- (2) then (section 2.4.3), the margins will be updated to  $\tilde{\alpha}_s$  and  $\tilde{\beta}_t$  by requiring a *minimum proportion*  $\theta \in [0, 1)$  of passengers entering/leaving the network at each stop, as well as avoiding transfer overflow exceeding the embarking and disembarking counts at each stop
- (3) finally (section 2.4.4), the prior will be updated to  $\tilde{g}_{st}$  by shrinking, if necessary, the priors  $g_{st}$  associated to overflows.

With the new prior distribution  $\tilde{g}_{st}$  and the new margin distributions  $\tilde{\alpha}_s, \tilde{\beta}_t$ , we can iterate the the above steps, until convergence. The only free parameter is  $\theta$ , whose effect is studied on toy examples in section 3.2.

The above iterative solution method is somewhat reminiscent of the EM algorithm. As a matter of fact, the first step (maximum entropy) exactly correspond to the “expectation step” of the EM algorithm (see e.g. [9] [10]), but steps two and three, aiming at calibrating parameters  $\alpha_s, \beta_t$  and  $g_{st}$ , do not follow the maximum likelihood rationale of the “maximisation step”. Pseudocode of the algorithm is shown with Algorithm 1.

#### 2.4.1 Maximum entropy estimate of st-trips

As announced, the proportion of st-trips  $f_{st} = n_{st}/n_{\bullet\bullet}$  (empirical distribution) will be estimated from some prior guess  $g_{st}$  (theoretical distribution) and margin constraints  $\alpha_s$  and  $\beta_t$  for  $f_{st}$  by maximum entropy, i.e. by solving the problem

$$\begin{aligned} \min_{\mathbf{f} \in \mathcal{F}} \quad & \sum_{st} f_{st} \log \frac{f_{st}}{g_{st}}, \\ \text{s.t.} \quad & \sum_t f_{st} = \alpha_s, \\ & \sum_s f_{st} = \beta_t. \end{aligned} \tag{6}$$

The Lagragian is

$$L = \sum_{st} f_{st} \log \frac{f_{st}}{g_{st}} - \sum_s \lambda_s (\alpha_s - \sum_t f_{st}) - \sum_t \mu_t (\beta_t - \sum_s f_{st}),$$

which gives, after deriving and setting to zero,

$$f_{st} = \phi_s \psi_t g_{st} \quad \text{with } \phi_s := \exp(-1 - \lambda_s), \psi_t := \exp(-\mu_t). \tag{7}$$

Using constraints in (6), we find

$$\phi_s = \frac{\alpha_s}{\sum_t \psi_t g_{st}}, \quad \psi_t = \frac{\beta_t}{\sum_s \phi_s g_{st}}, \tag{8}$$

which yields the following *iterative fitting algorithm*: starting with some  $\psi_t^{(0)} > 0$ , one performs the iteration

$$\phi_s^{(t)} = \frac{\alpha_s}{\sum_t \psi_t^{(t)} g_{st}}, \quad \psi_t^{(t+1)} = \frac{\beta_t}{\sum_s \phi_s^{(t)} g_{st}}, \quad (9)$$

until convergence to  $\phi_s$  and  $\psi_t$  obeying (8).

In view of (4) and (5), the postulated margins must satisfy, for each isolated stop  $i$

$$\alpha_i = \frac{a_i}{n_{\bullet\bullet}}, \quad \beta_i = \frac{b_i}{n_{\bullet\bullet}} \quad (10)$$

permitting to determine the total flow as  $n_{\bullet\bullet} = \frac{a_i}{\alpha_i}$ , or  $n_{\bullet\bullet} = \frac{b_i}{\beta_i}$  for any isolated stop  $i$ , and thus the  $st$ -flow itself as

$$n_{st} = n_{\bullet\bullet} f_{st} = n_{\bullet\bullet} \phi_s \psi_t g_{st} \quad (11)$$

whose plugging into (3) yields the intra-line edge flows  $\mathbf{Y} = (y_{ij})$  and the transfer edge flows  $\mathbf{Z} = (z_{ij})$ .

#### 2.4.2 Initialization of the prior and the margins

The geometry of the network permits to define the set of permitted  $st$ -trips across the network denoted by  $P$ . The initial prior was chosen as the uniform distribution on admissible paths, that is as

$$g_{st} = \begin{cases} \frac{1}{|P|} & [s, t] \in P, \\ 0 & \text{otherwise.} \end{cases}$$

The initial margins were chosen to initially match  $g_{st}$ , namely  $\alpha_s = g_{s\bullet}$  and  $\beta_t = g_{\bullet t}$  for all stops.

#### 2.4.3 Updating the margin distributions

Let us define the hyperparameter  $\theta \in [0, 1)$  as the *minimum proportion of passengers (among  $a_i$  and  $b_i$ ) entering/leaving the network at each stop*, that is  $n_{s\bullet} \geq \theta a_s$  and  $n_{\bullet t} \geq \theta b_t$ . Note that we could set a different hyperparameter for each node, and differing for embarkments and disembarkments, but without addition information, we will restrain to this simpler case. Identities (4) and (5) then imply the inequalities

$$z_{\bullet s} \leq (1 - \theta) a_s \quad z_{t\bullet} \leq (1 - \theta) b_t$$

the violation of which constitutes transfer overflow. Hence requiring a maximal transfer yet avoiding overflow can be granted with the following updating of margins

$$\tilde{\alpha}_s = \frac{\min(\theta a_s, a_s - z_{\bullet s})}{\sum_{s'} \min(\theta a_{s'}, a_{s'} - z_{\bullet s'})} \quad \tilde{\beta}_t = \frac{\min(\theta b_t, b_t - z_{t\bullet})}{\sum_{t'} \min(\theta b_{t'}, b_{t'} - z_{t'\bullet})} . \quad (12)$$

#### 2.4.4 Updating the prior distribution

Overflow occurs in transfer edge  $(i, j)$  if  $z_{i\bullet} > (1 - \theta)b_i$  or  $z_{\bullet j} > (1 - \theta)a_j$ . To avoid it, components  $g_{st}$  of the prior distribution will be shrunk by a suitable ratio whenever edge flows  $(i, j) \in [s, t]$  exhibit overflow. For any edge  $(i, j)$ , let us compute the *flow ratio*  $r_{ij}$  as

$$r_{ij} = \max \left( 1, \frac{z_{i\bullet}}{(1 - \theta)b_i}, \frac{z_{\bullet j}}{(1 - \theta)a_j} \right) \geq 1, \quad (13)$$

where  $r_{ij} > 1$  denotes an overflow through edge  $(i, j)$ . For a given origin-destination  $[s, t]$ , define the *origin-destination flow ratio*  $\bar{r}_{st}$  as the largest  $r_{ij}$  among edge flows  $(i, j) \in [s, t]$ , that is as

$$\bar{r}_{st} = \max_{ij} \chi_{ij}^{st} r_{ij} \geq 1. \quad (14)$$

By construction,  $\bar{r}_{st} > 1$  denotes an overflow on some transfer edge between  $s$  and  $t$ . To adjust the flow, we shall divide the previous flow by this ratio

$$\tilde{n}_{st} = \frac{n_{st}}{\bar{r}_{st}} \quad (15)$$

and define the new prior distribution as

$$\tilde{g}_{st} = \frac{\left( \frac{\tilde{n}_{st}}{\phi_s \psi_t} \right)}{\sum_{s', t'} \left( \frac{\tilde{n}_{s', t'}}{\phi_{s'} \psi_{t'}} \right)}. \quad (16)$$

where  $\phi_s$  and  $\psi_t$  are the values (8) obtained in the previous maximum entropy step.

---

**Algorithm 1** Compute the transportation flow matrix  $\mathbf{N} = (n_{st})$  knowing the edge-trip incidence matrix  $\chi = (\chi_{ij}^{st})$ , the set of transfer edges  $T$ , the set of permitted trips  $P$ , the embarking flow  $\mathbf{a}$ , the disembarking flow  $\mathbf{b}$ , the index of an isolated source node  $\tilde{s}$ , and the minimum proportion of passengers entering/leaving the network  $\theta$ .

---

```

1:  $g_{st} \leftarrow I([s, t] \in P) / |P|, \forall s, t$  ▷ Initialize the prior distribution
2:  $\alpha_s \leftarrow a_s / a_{\bullet}, \forall s$  ▷ Initialize the network ingoing distribution
3:  $\beta_t \leftarrow b_t / b_{\bullet}, \forall t$  ▷ Initialize the network outgoing distribution
4:  $\varepsilon \leftarrow 10^{-40}$  ▷ Fix a small quantity
5: while  $\mathbf{N} = (n_{st})$  has not converge do ▷ Main loop
6:    $\psi_t \leftarrow 1, \forall t$ 
7:   while  $\psi = (\psi_t)$  has not converge do ▷ Iterative fitting loop
8:      $\phi_s \leftarrow \alpha_s / (\sum_t \psi_t g_{st} + \varepsilon), \forall s$ 
9:      $\psi_t \leftarrow \beta_t / (\sum_s \phi_s g_{st} + \varepsilon), \forall t$ 
10:  end while
11:   $n_{st} \leftarrow \frac{a_{\tilde{s}}}{\alpha_{\tilde{s}}} \phi_s \psi_t g_{st}, \forall s, t$  ▷ Compute the transportation flow
12:   $z_{ij} \leftarrow I((i, j) \in T) \sum_{st} \chi_{ij}^{st} n_{st}, \forall i, j$ 
13:   $\alpha_s \leftarrow \frac{\min(\theta a_s, a_s - z_{s\bullet})}{\sum_{s'} \min(\theta a_{s'}, a_{s'} - z_{s'\bullet})}, \forall s$  ▷ Update the network ingoing distribution
14:   $\beta_t \leftarrow \frac{\min(\theta b_t, b_t - z_{\bullet t})}{\sum_{t'} \min(\theta b_{t'}, b_{t'} - z_{\bullet t'})}, \forall t$  ▷ Update the network outgoing distribution
15:   $r_{ij} \leftarrow \max \left( 1, \frac{z_{i\bullet}}{(1 - \theta)b_i}, \frac{z_{\bullet j}}{(1 - \theta)a_j} \right), \forall i, j$ 
16:   $\bar{r}_{st} \leftarrow \max_{ij} \chi_{ij}^{st} r_{ij}, \forall s, t$ 
17:   $\tilde{g}_{st} \leftarrow \frac{\left( \frac{n_{st}}{\phi_s \psi_t \bar{r}_{st} + \varepsilon} \right)}{\sum_{s', t'} \left( \frac{n_{s', t'}}{\phi_{s'} \psi_{t'} \bar{r}_{s', t'} + \varepsilon} \right) + \varepsilon}, \forall s, t$  ▷ Update the prior distribution
18: end while
19: return  $\mathbf{N} = (n_{st})$ 

```

---

## 2.5 Markov property for a single line

A “network” made of a single line contains no transfers, and flow estimates can be obtained at once by the maximum entropy step only.

Let  $i = 1, \dots, l$  enumerate the bus stops in increasing order, i.e.  $F(i) = i + 1$ . The initial prior is simply  $g_{st} = c I(s < t)$  and captures solely the unidirectional nature of trips, where  $I(\cdot)$  denotes the 0/1 indicator function and  $c = \frac{1}{(l-1)(l-2)}$ . The margins of the empirical distribution  $f_{st}$ , as well as the total flow, are here known :

$$\alpha_s = \frac{a_s}{a_\bullet} \quad \beta_t = \frac{b_t}{b_\bullet} \quad n_{\bullet\bullet} = a_\bullet = b_\bullet .$$

Following (7) maximum entropy flows are of the form

$$n_{st} = n_{\bullet\bullet} c I(s < t) \phi_s \psi_t \quad (17)$$

where (setting  $\Psi_s := \sum_{t>s} \psi_t$  and  $\Phi_t := \sum_{s<t} c \phi_s$ ) the constraints (8) equivalently read

$$\phi_s = \frac{\alpha_s}{c \sum_{t>s} \psi_t} = \frac{a_s}{n_{\bullet\bullet} c \Psi_s} \quad \psi_t = \frac{\beta_t}{c \sum_{s<t} \phi_s} = \frac{b_t}{n_{\bullet\bullet} c \Phi_t} \quad (18)$$

to be solved by iterative fitting.

Interestingly enough, the form (17) for the flows is reminiscent of the *gravity flows* of quantitative Geography [7] [8] [11] [12], where  $\phi_s$  is the *push factor*,  $\psi_t$  is the *pull factor*, and  $I(s < t)$  the *distance deterrence function*. Yet, instead of being symmetric in  $s, t$  and decreasing with the distance  $|s - t|$ , the distance deterrence function is here asymmetric due to the line orientation, but otherwise constant.

This constancy entails the following Markovian behaviour for flows: let  $m_{st}$  be the number of travelers embarking at stop  $s$  and still inside the bus at stop  $t > s$ , and let  $\rho_{st}$  the probability that travelers embarking at  $s$  will disembark at  $t$ . By (17),

$$m_{st} = \sum_{u \geq t} n_{su} = n_{\bullet\bullet} c \phi_s \sum_{u \geq t} I(s < u) \psi_u = n_{\bullet\bullet} c \phi_s (\psi_t + \Psi_t)$$

The empirical estimate of  $\rho_{st}$  is given by the proportion, among the travelers embarking at  $s$  and still present at  $t > s$ , of travelers disembarking at  $t$ , that is

$$\rho_{st} = \frac{n_{st}}{m_{st}} = \frac{n_{\bullet\bullet} c \phi_s \psi_t}{n_{\bullet\bullet} c \phi_s (\psi_t + \Psi_t)} = \frac{\psi_t}{\psi_t + \Psi_t} \leq 1$$

which depends on  $t$  only: it appears that the disembarkment probability  $\rho_t = \frac{\psi_t}{\psi_t + \Psi_t}$  at  $t$  is *independent* of the embarkment stop  $s$ . Said otherwise, a traveler embarking at any stop  $s$  (and thus necessarily in the bus at  $F(s) = s + 1$ ) experiences the *same disembarkment probability* at each further stop  $t > s$ .

This Markov property, enjoyed by maximum-entropic flows, contrasts other possible solutions, such as the “first in, first out” (FIFO) flows (homogenizing the traveled distances among users) or the “last in, first out” (LIFO) flows (tending to generate maximally contrasted traveled distances).

### 3 Case Studies

Case studies are divided in two sections. In the first section, we test the algorithm on toy examples, which are artificial networks where the transportation flow  $n_{st}$  is randomly drawn. These examples enable some kind of validation of the algorithm, as the "real" transportation flow is known and can be compared to solutions given by our method. This setup differs from the second section, which is dedicated to applying the algorithm to the real case of the public transportation network of the city of Lausanne (tl), where embarkment and disembarkment flows are measured but the real transportation flow is unknown. This second case study shows that the algorithm is applicable on large, real datasets and can give insights about passengers probable routes in the network.

#### 3.1 Error measurements

In all case studies, we obtain an estimation of the transportation flow with the algorithm, noted  $\mathbf{N} = (n_{st})$ , starting from the real embarkment flow  $\mathbf{a}_{\text{ref}}$  and disembarkment flow  $\mathbf{b}_{\text{ref}}$ . In toy examples, we also have access to the real transportation flow  $\mathbf{N}_{\text{ref}}$ . There are two types of dissimilarity measures between the data and the solution proposed by the algorithm: (1) if we have access to  $\mathbf{N}_{\text{ref}}$ , how much  $\mathbf{N}$  differs from it, and (2) how well constraints defined by  $\mathbf{a}_{\text{ref}}$  and  $\mathbf{b}_{\text{ref}}$  are respected. The first dissimilarity is measured through the *mean transportation error*, denoted by  $\text{MTE}(\mathbf{N})$ , and computed with

$$\text{MTE}(\mathbf{N}) = \sum_{st} \frac{n_{st}^{\text{ref}}}{n_{\bullet\bullet}^{\text{ref}}} \frac{|n_{st} - n_{st}^{\text{ref}}|}{n_{st}^{\text{ref}}} = \frac{\sum_{st} |n_{st} - n_{st}^{\text{ref}}|}{n_{\bullet\bullet}^{\text{ref}}} \quad (19)$$

and the second one with the *mean margin error*, noted  $\text{MME}(\mathbf{N})$ , obtained with

$$\begin{aligned} \text{MME}(\mathbf{N}) &= \frac{1}{2} \sum_i \frac{a_i^{\text{ref}}}{a_{\bullet\bullet}^{\text{ref}}} \frac{|z_{\bullet i} + n_{i\bullet} - a_i^{\text{ref}}|}{a_i^{\text{ref}}} + \frac{1}{2} \sum_i \frac{b_i^{\text{ref}}}{b_{\bullet\bullet}^{\text{ref}}} \frac{|z_{i\bullet} + n_{\bullet i} - b_i^{\text{ref}}|}{a_i^{\text{ref}}} \\ &= \frac{\sum_i (|z_{\bullet i} + n_{i\bullet} - a_i^{\text{ref}}| + |z_{i\bullet} + n_{\bullet i} - b_i^{\text{ref}}|)}{2n_{\bullet\bullet}^{\text{ref}}} \end{aligned} \quad (20)$$

where  $z_{ij} = I((i, j) \in T) \sum_{st} \chi_{ij}^{st} n_{st}$  is the flow on transfer edges  $T$ . Both errors can be interpreted as a weighted mean deviation percentage.

Note that, by construction, the MME should be null when the algorithm converges. However, it can be informative to track down this error along iterations and, in some practical cases where the network is large, the algorithm convergence criterion is reached without having margin constraints perfectly respected.

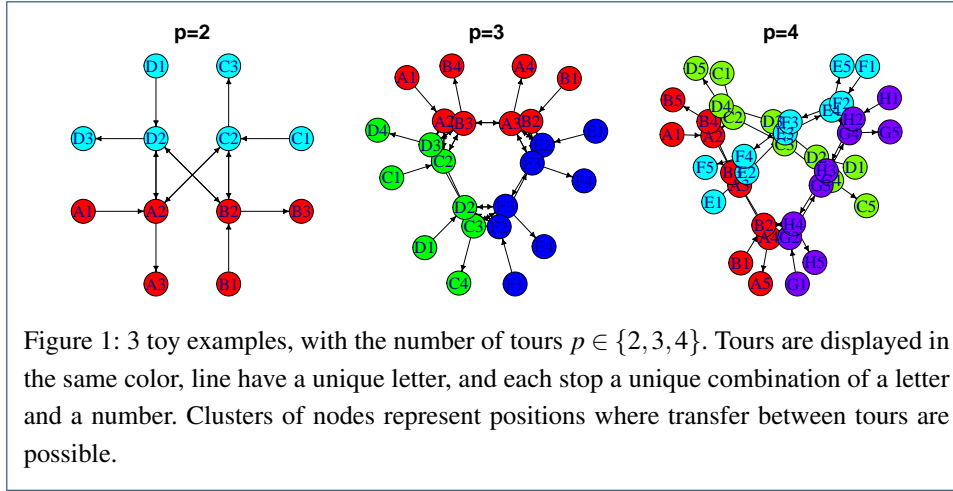
#### 3.2 Toy Examples

##### 3.2.1 Construction

All constructed toy examples are built following the same approach, which aims at being simple but somewhat realistic. We fix a number of *line tours*  $p \geq 2$ , each of which is constituted of a forward line and a backward line, for a total of  $q = 2p$  lines. Every line has a starting and ending node, which are isolated nodes, and possesses  $p - 1$  intermediary nodes which allows transfers to the other tour lines, giving a total of  $n = 2p(p + 1)$  nodes in the network. Examples of these toy networks can be found in Figure 1.

In order to be realistic, permitted *st*-trips set  $P$  is constructed considering all shortest-path between pair of nodes, excluding :





- $s$  and  $t$  that are on the same line but with  $t$  preceding  $s$  in the line order.
- $s$  and  $t$  that are on the same tour but opposite line.
- $s$  and  $t$  whose shortest-path starts with a transfer edge, ends with a transfer edge, or possesses two (or more) consecutive transfer edges.

A transportation flow  $\mathbf{N}_{\text{ref}} = (n_{st}^{\text{ref}})$  is drawn by setting a fixed number of passengers  $n_{\bullet\bullet}^{\text{ref}}$ , and each passenger is assigned randomly to a  $(s, t)$  pair drawn uniformly among  $P$ . From this reference transportation flow  $\mathbf{N}_{\text{ref}}$ , using the edge-trip incidence matrix  $\chi$  and equation (3), we can compute flow on edges  $\mathbf{X}_{\text{ref}}$  and, in turn, the number of passengers embarking  $\mathbf{a}_{\text{ref}}$  and the number of passengers disembarking  $\mathbf{b}_{\text{ref}}$  at each stop.

### 3.2.2 Algorithm iterations

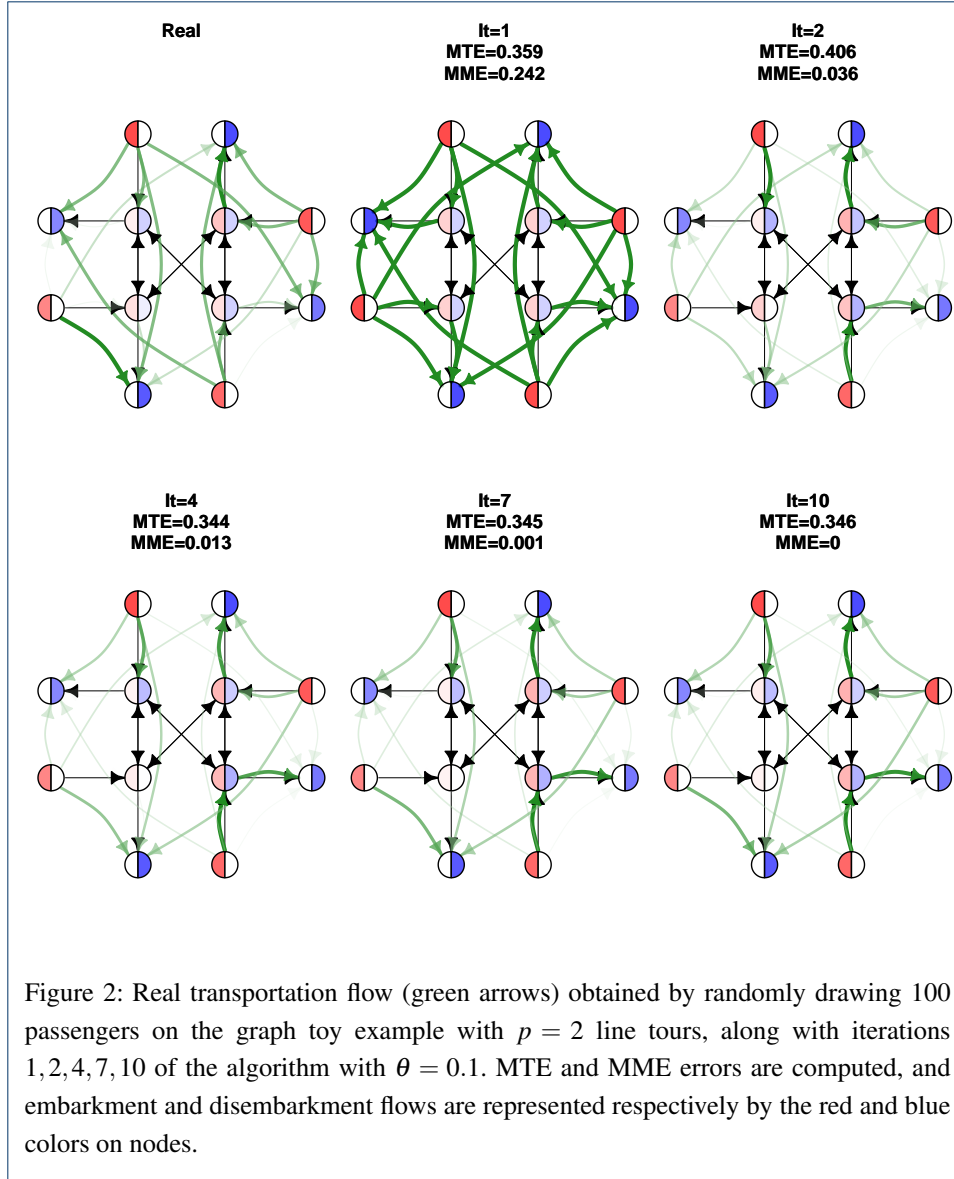
First, we show some algorithm iterations on a toy example with  $p = 2$ , where 100 passengers were drawn uniformly across the  $|P| = 20$  possible  $st$ -trips. Some iterations of the algorithm with  $\theta = 0.1$ , along with MTE and MME errors, are shown in Figure 2. On this small example, we can see that the algorithm quickly finds an estimation giving small MME error, but still gives an MTE of 0.252. This result is due to the fact that only 50 passengers are drawn, giving a large deviation compared to the optimally found solution which maximizes the entropy. Interestingly, we see that a better result is found on iteration, but margin constraints at this point are not perfectly respected yet.

### 3.2.3 MTE study

The main goal of toy examples, since we have access to the real transportation flow, is to study how the resulting MTE behave regarding passenger  $st$ -trips distribution and hyperparameter  $\theta$ , on different network sizes.

The randomness of the  $st$ -trips distribution is controlled by the number of drawn passengers and we can see that the algorithm performs better if this number increases, as seen in Figure 3. This behavior can be expected as the method is constructed to find the solution maximizing the entropy while respecting embarkment and disembarkment constraint. In fact, all MTE should converge to 0 eventually, however this rate of convergence seems to be low.

As for the hyperparameter study, found in Figure 4, we see that the optimal parameter seems to decrease as network size increases. However,

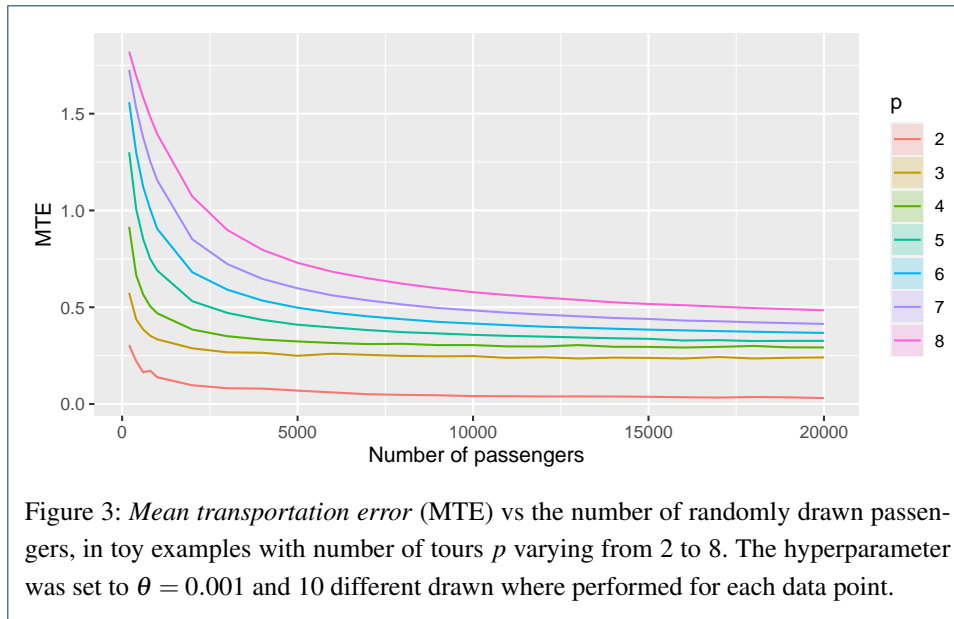


### 3.3 Real Data

after some preliminary, undocumented corrections (i.e. the components of  $\mathbf{a}$  and  $\mathbf{b}$  can be non-integer). It may also happen that, on some lines  $\ell$ , raw data do not obey the necessary consistency condition  $a_{\bullet}^{\ell} = b_{\bullet}^{\ell}$  (where the latter quantities denote the total embarkments and disembarkments on line  $\ell$ ), in which case we did rescale the embarking and disembarking line counts as

$$\hat{a}_i = \left(1 - \frac{a_{\bullet}^{\ell} - b_{\bullet}^{\ell}}{a_{\bullet}^{\ell} + b_{\bullet}^{\ell}}\right) a_i \quad \hat{b}_i = \left(1 + \frac{a_{\bullet}^{\ell} - b_{\bullet}^{\ell}}{a_{\bullet}^{\ell} + b_{\bullet}^{\ell}}\right) b_i$$

ensuring  $\hat{a}_{\bullet}^{\ell} = \hat{b}_{\bullet}^{\ell} = 2a_{\bullet}^{\ell}b_{\bullet}^{\ell}/(a_{\bullet}^{\ell} + b_{\bullet}^{\ell})$ . However, strongly unbalanced lines such that  $|a_{\bullet}^{\ell} - b_{\bullet}^{\ell}|/a_{\bullet}^{\ell} > 0.3$  or  $|a_{\bullet}^{\ell} - b_{\bullet}^{\ell}|/b_{\bullet}^{\ell} > 0.3$  (which always turned out to be temporary lines with small counts) were simply disregarded and line  $\ell$  removed from the network.

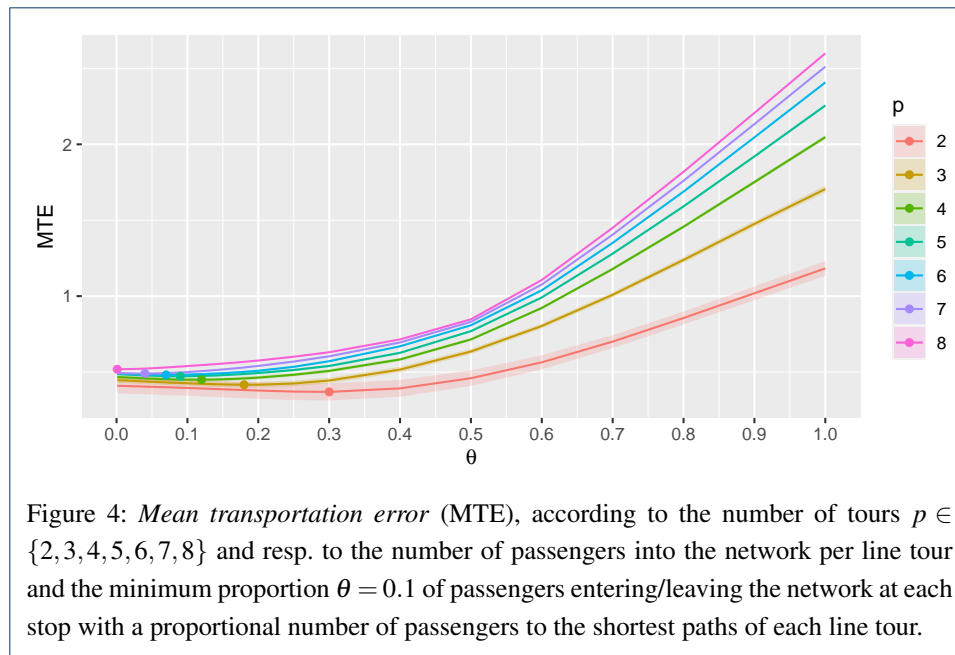


Also, the geometry of the network permits to derive the edge-trip incidence matrix  $\chi$  defined in (1).

### 3.3.1 Dataset construction

This research presents a dataset derived from the public transportation agency in Lausanne (tl), Switzerland. Automatic passenger counters were used to collect data on the number of passengers entering and leaving buses and metros at each stop. The dataset comprises 44 lines of metro and bus transportation, 1361 stops, and over 115 million passengers per year. Each round trip forms two separate lines, which often have stops that correspond to both the inbound and outbound directions, but this is not always the case. Some stops may be unique to either the inbound or outbound direction. The initial dataset contained 13 million data rows, which were aggregated to the year 2019, and the passenger data was filtered to include only passengers who boarded and disembarked at each stop. Lines with the most errors were removed, resulting in a final dataset of 35 transportation lines and 1216 stops. The threshold for invalid lines was set at a difference of 15% between boarded and disembarked passenger counters.

## 4 Conclusion



## Appendix

Text for this section. . .

### Acknowledgements

Text for this section. . .

### Funding

Text for this section. . .

### Abbreviations

Text for this section. . .

### Availability of data and materials

Text for this section. . .

### Ethics approval and consent to participate

Text for this section. . .

### Competing interests

The authors declare that they have no competing interests.

### Consent for publication

Text for this section. . .

### Authors' contributions

Text for this section. . .

### Authors' information

Text for this section. . .

### Author details

<sup>1</sup>Department of Language and Information Sciences, University of Lausanne, Lausanne, Switzerland. <sup>2</sup>Institute of Geography and Sustainability, University of Lausanne, Lausanne, Switzerland.

### References

1. Peyré, G., Cuturi, M., et al.: Computational optimal transport: With applications to data science. Foundations and Trends® in Machine Learning **11**(5-6), 355–607 (2019)
2. Bell, M.G., Lida, Y.: Transportation Network Analysis. Wiley, Chichester (1997)
3. Hazelton, M.L.: Estimation of origin–destination matrices from link flows on uncongested networks. Transportation Research Part B: Methodological **34**(7), 549–566 (2000)
4. Ashok, K., Ben-Akiva, M.E.: Estimation and prediction of time-dependent origin-destination flows with a stochastic mapping to path flows and link flows. Transportation science **36**(2), 184–198 (2002)

5. Cui, A.: Bus passenger origin-destination matrix estimation using automated data collection systems. PhD thesis, Massachusetts Institute of Technology (2006)

6. Jaynes, E.T.: Information theory and statistical mechanics. *Physical review* **106**(4), 620 (1957)

7. Wilson, A.: A statistical theory of spatial distribution models. *Transportation Research* **1**(3), 253–269 (1967)

8. Erlander, S., Stewart, N.F.: *The Gravity Model in Transportation Analysis: Theory and Extensions* vol. 3. VSP, Leiden (1990)

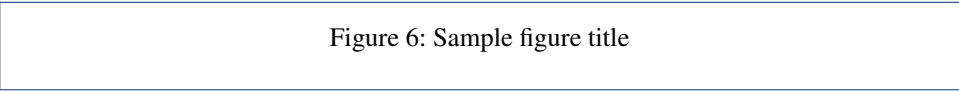
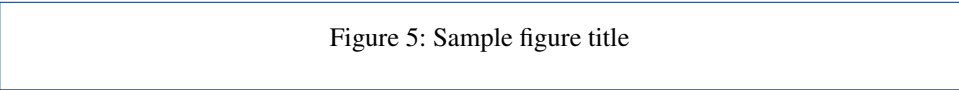
9. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society: series B (methodological)* **39**(1), 1–22 (1977)

10. Bavaud, F.: Information theory, relative entropy and statistics. In: Sommaruga, G. (ed.) *Formal Theories of Information: From Shannon to Semantic Information Theory and General Concepts of Information*. LNCS, vol. 5363, pp. 54–78. Springer, Berlin (2009)

11. Bavaud, F.: The quasi-symmetric side of gravity modelling. *Environment and Planning A* **34**(1), 61–79 (2002)

12. Thomas-Agnan, C., LeSage, J.P.: In: Fischer, M.M., Nijkamp, P. (eds.) *Spatial Econometric OD-Flow Models*, pp. 2179–2199. Springer, Berlin, Heidelberg (2021)

Figures



Tables

Table 1: Sample table title. This is where the description of the table should go

	B1	B2	B3
A1	0.1	0.2	0.3
A2	...	..	.
A3	..	.	.

Additional Files

Additional file 1 — Sample additional file title  
Additional file descriptions text (including details of how to view the file, if it is in a non-standard format or the file extension). This might refer to a multi-page table or a figure.

Additional file 2 — Sample additional file title  
Additional file descriptions text.