

# Confidence Intervals for OD Demand Estimation

Yingying Chen, Fernando Ordóñez\*, and Kurt Palmer

Industrial and Systems Engineering, University of Southern California,  
3715 McClintock Ave. GER-240, Los Angeles, CA 90089-0193,  
juphia@gmail.com, fordon@usc.edu, kpalmer@usc.edu

January, 2006

## Abstract

Representative origin-destination (OD) demand tables are a crucial part of making many transportation models relevant to practice. However estimating these OD tables is a challenging problem, even more so determining the confidence intervals on these OD estimates. In this work we propose a method to construct estimates and confidence intervals of OD demand tables from link flow data. Our method separates the uncertainty in estimating OD tables into the statistical uncertainty of link flow data and the possibility of multiple feasible OD demand solutions for the same link flow data. A confidence interval is constructed from concise representations of the uncertainty in each part. We illustrate our estimation method through examples, including one with real data from an intersection in the Los Angeles freeway system.

**Keywords:** Origin-destination Demand Estimation; Confidence Intervals; Link Flow Data.

---

\*corresponding author

# 1 Introduction

A fundamental part in modeling a transportation network is an accurate estimate of the traffic demand. This demand is typically represented by the number of trips between specific origin and destination pairs (OD pairs), which forms the origin-destination demand table or matrix (or OD tables, for short). This demand information can further indicate the types of flow (or commodities) between OD pairs and also capture the dynamic nature of demand throughout the day.

For example, Intelligent Transportation Systems (ITS) are identified as the means to achieve sustainable and environment friendly transportation for the 21st Century. An ITS system could collect and process real-time traffic conditions data, along with OD demand estimation, to control and manage traffic. Having accurate estimates of demand, with confidence intervals, can help develop policies for reduction of traffic congestion, enhanced safety, and mitigation of environmental impacts of transportation systems.

In most applications the estimates of demand will be subject to significant uncertainty, either due to the fact that the demand is being estimated in advance for planning purposes, or that the information available to infer the OD flows is not complete. In this work, we assume that the data used to estimate these OD trip tables originates from incomplete surveys or economic studies and is adjusted with link flow information, available from loop detectors in the freeway system. In particular we are interested in situations where the data used to estimate the OD flows is insufficient to determine OD flows unambiguously. In other words, we focus on the under-determined case when the number of observations is less than the number of parameters to estimate. It turns out that this is a common occurrence when estimating OD flows, see Bierlaire and Crittin (2003) for a discussion. Due to this uncertain nature of future OD demands, confidence intervals of the demand estimates are in fact more relevant than a specific estimate of demand.

There are a number of demand estimation models in the literature, each making its own assumptions to select a single estimate from the underdetermined estimation problem. However, to our knowledge, the only analytical method to obtain confidence intervals for an OD demand estimate from link flow data is to repeatedly use a demand estimation model. This would make the mean of the OD demand estimates obtained a multivariate normal random variable for which an ellipsoidal confidence level set is readily available, see for example (Morrison 1976). Such a confidence interval has two important drawbacks: it is influenced by the estimation model assumptions, rather than only by the data; and, the volume of the confidence interval can be reduced arbitrarily by increasing the number of times  $q$  the estimation model is used, since the variance of the sample mean is proportional to  $\frac{1}{\sqrt{q}}$ .

We introduce a method to estimate static OD demand and construct confidence intervals that only depends on the uncertainty present in the data. We focus on static demand for a single flow type, however the extension to multicommodity flows is straightforward leading simply to a problem that is further underdetermined as loop detectors are not able to discriminate between flow types from the link flow data. Our demand estimates and confidence intervals are obtained by explicitly and concisely representing the possible solutions to the underdetermined system and the statistical uncertainty of our estimate which lies in the complement space. Each part is represented with ellipsoids which yield easy answers for coordinate wise projected confidence intervals. The demand estimate corresponds to the analytic center of the set of possible solutions.

The next section introduces the estimation model and discusses prior demand estimation methods. In Section 3 we present an analytic center based estimate of OD-flows and describe how to construct a confidence interval for this demand estimate. We describe our numerical experiments and present their results in Section 4. Finally, we present some closing remarks and conclusions in Section 5.

## 2 OD pairs estimation and prior work

Estimating OD pair demands for a place like Los Angeles is truly a challenging problem not only because of the sheer size of the transportation network, but also because of the underdetermined nature of the problem. However in spite of this difficulty, or possibly because of it, many different estimation models have been proposed. Distinguishing features of these models are whether they address a dynamic or a static traffic model, whether the models obtain demands that satisfy equilibrium conditions or that optimize some objective, and also whether the assignment matrix  $X$  is stochastic to account for user path choice or not.

We assume that the observed link flows at  $m$  different loop detectors, denoted  $L \in \mathfrak{R}^m$ , are related to the true OD flows,  $f \in \mathfrak{R}^n$ , in the following linear model

$$L_l = \sum_{r \in R} X_{lr} f_r + \varepsilon_l, \quad l = 1, \dots, m. \quad (1)$$

The matrix  $X \in \mathfrak{R}^{m \times n}$ , referred to as the assignment matrix, indicates whether OD-flow  $r$  passes through the loop detector  $l$ . In addition we assume that dynamic traffic conditions, data measurement errors, and unmodeled random events are jointly represented by the error term  $\varepsilon \sim N(0, \sigma^2 I)$ . In this work we consider a deterministic assignment matrix, which takes only the values  $X_{lr} = 1$  or  $= 0$  when OD-flow  $r$  passes through  $l$  or not, respectively. This model requires one OD flow variable for each possible route between each origin and destination. Thus  $n$  is the total number of paths between OD pairs. An alternative model is to consider a stochastic assignment matrix, where  $X_{lr} \in [0, 1]$  now represents the probability that OD pair  $r$  routes its flow through loop detector  $l$ , reducing the number of flow variables,  $n$ , to the number of OD pairs. For illustration consider Figure 1 to compare deterministic and stochastic assignment matrix models. This network has three OD flows: from node 1 to node 3, either straight or through node 2, and from node 2 to node 3. We have link flow data from every arc.

This network leads to the following two linear models depending on which assignment

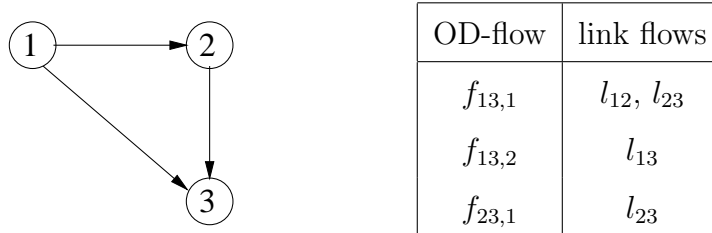


Figure 1: A network example with three OD flows and the link flows each impacts.

matrix is used:

$$\begin{pmatrix} l_{12} \\ l_{13} \\ l_{23} \end{pmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{pmatrix} f_{13,1} \\ f_{13,2} \\ f_{23,1} \end{pmatrix} \qquad \begin{pmatrix} l_{12} \\ l_{13} \\ l_{23} \end{pmatrix} = \begin{bmatrix} \alpha_{13} & 0 \\ 1 - \alpha_{13} & 0 \\ \alpha_{13} & 1 \end{bmatrix} \begin{pmatrix} f_{13} \\ f_{23} \end{pmatrix}$$

deterministic matrix stochastic matrix .

A dynamic traffic model has to explicitly represent the evolution of traffic through the network by keeping track that the current flow at a loop detector originated at different locations at different previous time periods. This would modify our static model by indexing it on different times, yielding something like

$$L_{l,t} = \sum_{p=t-p'}^t X_{lrt}^p x_{rp} + \varepsilon_{l,t} ,$$

where  $L_{l,t}$  contains the observed traffic flow for time interval  $t$ , the assignment matrix  $X_{lrt}^p$  represents the likelihood that OD flows departing during interval  $p$  are observed during interval  $t$  at loop detector  $l$ .

For our work, the most relevant distinguishing feature of previous estimation models is the method used to resolve the under-determinedness of the problem. Below we classify prior estimation models following the method used to distinguish an estimate out of the potentially many solutions that satisfy the link flow data.

Some classic models incorporate concepts from physics or information theory to select a single OD demand table with maximum likelihood out of all the feasible ones.

In *gravity models* (Robillard 1975) the likelihood of a trip between an origin and destination is considered inversely proportional to the square of the distance between them. In *entropy models* (Van Zuylen and Willumsen 1980) the maximum likelihood estimate is obtained by selecting the OD demand table that provides the least amount of additional information, minimizing an information measure. There are additional maximum likelihood models, which differ in the functions and algorithms used and the method of incorporating prior information, see for example (Spiess 1987; Van Aerde et al. 2003).

Statistical features of dynamic link flows are also used to identify single OD demand. For example, Cremer and Keller (1987) consider using the correlations between entering and exiting flows to estimate demand, and Hazelton (2003) uses covariance information of the link flows data to resolve the problem of indeterminacy.

Additionally some models provide an OD demand estimate which satisfies flow equilibrium constraints while others solve some optimization problem, thus the estimated demand is optimal with respect to a given objective. Examples of models which include equilibrium constraints, and thus lead to bilevel programming, include (Fisk and Boyce 1983; Yang et al. 1992; Sherali et al. 1994; Bell et al. 1997; Cascetta and Postorino 2001; Nie et al. 2005), where these models differ on the objective and type of network considered.

Typically, optimization based models minimize a least squares expression of estimation errors. Examples on dynamic traffic models with a fixed deterministic assignment matrix include (Cascetta 1984; Cascetta et al. 1993; Sherali and Park 1999; Ashok and Ben-Akiva 2000; Bierlaire and Crittin 2003). These models differ in whether they use survey or historical data and on the specific algorithms proposed. Lo et al. (1996) and Ashok and Ben-Akiva (2002) minimize least squares and consider a stochastic assignment matrix. Additional optimization models include (Brenninger-Göthe et al. 1989; Bell 1991; Doblus and Benitez 2005).

These models invariably bias the OD estimate to a single point estimate using some additional criterion or assumption. In our work we propose a concise method to characterize all solutions implied by the link flow data, explicitly identifying the uncertainty in the estimate due to the under-determinedness of the problem. In addition these models do not discuss how to construct confidence intervals for their OD demand estimate. We note that Van Aerde et al. (1993) present confidence intervals for OD demand estimates obtained from probe vehicle data. Their approach however does not translate to a method for demand estimates from link flow data.

### 3 Estimation model

Recall from (1) our linear model relating the observed link flows  $L \in \mathfrak{R}^m$  and the true OD-flows  $f \in \mathfrak{R}^n$ , which in matrix notation is

$$L = Xf + \varepsilon . \tag{2}$$

where the matrix  $X$  corresponds to the deterministic assignment matrix described in the previous section. As was discussed in the introduction, it turns out that the number of possible OD-pairs far exceeds the number of loop detectors that are gathering link flow data for most transportation networks. Hence the system (2) is in general underdetermined, which implies that it is likely there are multiple solutions. For example, consider an intersection as in Figure 2. The letters in uppercase stand for inflows while those in lowercase stand for outflows. Such a freeway intersection has 12 different OD pairs, since every inflow to the freeway intersection can exit in three different directions, and there are 4 such inflows; and if we gather link flows in every possible segment of road there are a total of 8 link flows. Thus in this system there are 12 unknowns with only 8 data points. Our approach explicitly models these multiple optimal solutions with the additional assumption that OD flows are positive and bounded. This is the only additional assumption used, which is reasonable considering that capacities on roads

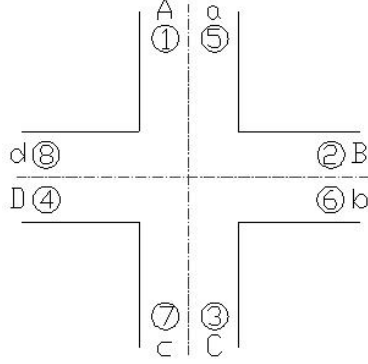


Figure 2: Two-way traffic street intersection.

and total population limit the total OD flow. We concisely express the set of optimal solutions to (2), which turns out to be bounded, through ellipses that are contained in the likely feasible region and build the confidence intervals that depend only on the error uncertainty.

### 3.1 Characterizing all minimal error solutions

The least squares estimates of Model (2) are the solutions which minimize the sums of squares of the residual  $Xf - L$ . From optimality conditions we know that the estimates satisfy the normal equations  $X^T X f = X^T L$  and equivalently belong to

$$S = X^+ L + \text{Ker}(X^T X) = (X^T X)^- X^T L + \{v \mid X^T X v = 0\} , \quad (3)$$

where  $(A)^+$  denotes the Moore-Penrose inverse of a possibly singular matrix  $A$  and  $(A)^-$  denotes the 1-matrix inverse, see (Campbell and Meyer 1991). For example, when  $X^T X$  is invertible (and  $\text{Ker}(X^T X) = \{0\}$ ), then  $(X^T X)^- = (X^T X)^{-1}$  and  $X^+ = (X^T X)^{-1} X^T$  is the regular least squares projection matrix. We now provide a short proof of this result in the general case, with a possibly singular matrix, for completeness and to set the notation that will be used later in this section.

**Proposition 1.** For  $X \in \mathbb{R}^{m \times n}$  with  $\text{rank}(X) = q \leq \min\{m, n\}$ , the solution to the problem  $\min_f (Xf - L)^T (Xf - L)$  is given by the set  $S$  defined in (3).

**Proof:** Let  $\text{rank}(X^T X) = q \leq \min\{m, n\}$ , then let  $W$  be the  $m \times q$  matrix of orthonormal eigenvectors associated to positive eigenvalues of  $X^T X$  that form the  $q \times q$  diagonal matrix  $D^2$ . Note that this makes  $D$  the matrix of singular values of  $X$ . We also let  $V$  be the (possibly empty)  $m \times (n - q)$  matrix of orthonormal eigenvectors with 0 eigenvalue. Then the singular value decomposition is  $X = U[D \ 0; 0^T \ 0][WV]^T$  for some unitary matrix  $U$ , and the eigenvalue decomposition of  $X^T X = [W \ V][D^2 \ 0; 0^T \ 0][W \ V]^T$ , and the 1-matrix is given by  $(X^T X)^- = [W \ V][D^{-2} \ 0; 0^T \ 0][W \ V]^T = WD^{-2}W^T$ . Substituting these singular value and eigenvector decompositions our problem is simply  $\min_f L^T L - 2L^T U[D; 0^T]W^T f + f^T W D^2 W^T f$ . Which through the change of variables  $z = W^T f$  leads to a strictly convex problem in  $q$  variables with a single optimal solution:  $z^* = [D^{-1} \ 0]U^T L$ . Since for any  $z \in \mathbb{R}^q$ , the vector  $f = Wz + V\alpha$ ,  $\alpha \in \mathbb{R}^{n-q}$  satisfies  $W^T f = z$ , we have that  $f^* \in W[D^{-1} \ 0]U^T L + \text{Ker}(X^T X)$ . It is easy to check that  $W[D^{-1} \ 0]U^T L = (X^T X)^- X^T L$ . ■

A reasonable assumption on OD flows is that the flow between each OD pair is non-negative and bounded, in other words  $0 \leq f_r \leq U_r$  for any OD flow  $r$ . Hence, out of the potentially multiple solutions to the normal equations above, we are interested only in those that in addition satisfy these upper and lower bounds. In fact we assume that the system is such that  $S \cap [0, U] \neq \emptyset$ , in other words, for some  $v \in \text{Ker}(X^T X)$ ,  $0 \leq (X^T X)^- X^T L + v \leq U$ . If this were not the case, and  $S \cap [0, U] = \emptyset$ , then it means that the least square estimate of the link flow data obtained leads to flows that are either negative or larger than  $U_r$ , which suggests we are missing something. A least squares estimate is still given by solving the following problem

$$\begin{aligned} \min \quad & (L - Xf)^T (L - Xf) \\ \text{s.t.} \quad & 0 \leq f_r \leq U_r \quad r \in \{1, \dots, OD\} \end{aligned}$$

which provides a biased estimator of the true flow under Model (2).

We now describe our concise representation of the set of solutions  $S$ . Let  $\bar{f} = (X^T X)^- X^T L = X^+ L$  be the least square solution to  $Xf = L$ . In general the solution set  $S$  is formed by the subspace  $\text{Ker}(X^T X)$  passing through the point  $\bar{f}$ . If in addition we consider the upper and lower bounds, we want to efficiently represent the set of solutions  $S' = \{\bar{f} + Vz \mid 0 \leq \bar{f} + Vz \leq U\}$ , where recall we denote by  $V$  the  $m \times (n - q)$  matrix of eigenvectors of  $X^T X$  associated with eigenvalue 0. For instance in the crossroads example, the vectors  $f \in \Re^{12}$  and the null space of  $X^T X$  is of dimension 5. If we denote by  $f^*$  the true flow we are trying to estimate, then we show a schematic drawing in Figure 3 the true flow and least squares estimate  $\bar{f}$  along with the bounded set  $S'$  of possible solutions. We approximate the bounded set  $S'$  by inscribed and circumscribed ellipsoids centered at the analytic center of  $S'$ . The analytic center  $\hat{f} = \bar{f} + V\hat{z}$  of  $S'$  is given as the solution to

$$\min_z - \sum_{r=1}^{OD} (\log(U_r - \bar{f}_r - (Vz)_r) + \log(\bar{f}_r + (Vz)_r)) .$$

Let  $H$  be the Hessian of the objective function for the analytic center problem, and denote by  $B(\hat{z}, r) = \{z \mid (z - \hat{z})^T H^{-1} (z - \hat{z}) \leq r\}$  the ball centered at  $\hat{z}$  of radius  $r$  with the norm defined by  $H$ . We know from (Renegar 2001) that

$$B(\hat{z}, 1) \subseteq \{z \mid 0 \leq \bar{f} + Vz \leq U\} \subseteq B(\hat{z}, 4\vartheta + 1)$$

and hence

$$\bar{f} + VB(\hat{z}, 1) \subseteq S' \subseteq \bar{f} + VB(\hat{z}, 4\vartheta + 1) .$$

Note that  $\bar{f} + VB(\hat{z}, r) = \hat{f} + VB(0, r)$ . We present a schematic figure of how these different solutions are related in Figure 3, including both ellipsoidal sets.

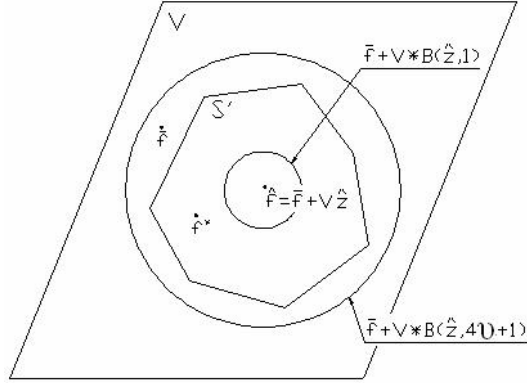


Figure 3: Schematic representation of the null space  $V = \text{Ker}(X^T X)$ , the solution set  $S'$ , the inscribed and circumscribed ellipsoids, and true flow  $f^*$ , least square estimate  $\bar{f}$ , and analytic center estimate  $\hat{f}$ .

### 3.2 Confidence intervals

In this section we describe how to construct confidence intervals on our OD flow estimate, given by  $\hat{f}$ , (note that it is possible that  $\bar{f} \notin [0, U]$ ).

$$\begin{aligned}
 \hat{f} &= \bar{f} + V\hat{z} \\
 &= (X^T X)^{-1} X^T L + V\hat{z} \\
 &= (X^T X)^{-1} X^T (X f^* + \varepsilon) + V\hat{z} \\
 &= f^* - VV^T f^* + V\hat{z} + (X^T X)^{-1} X^T \varepsilon .
 \end{aligned}$$

Here we consider that  $f^*$  represents our true flow, and the term  $V\hat{z}$  amounts to a change in  $\text{Ker}(X^T X)$  to move our unconstrained least squares estimate  $\bar{f}$  to the analytic center of the feasible estimates of  $f^*$ .

We note that all the uncertainty due to errors  $\varepsilon$  is concentrated in the last term. If we assume that  $\varepsilon \sim N(0, \sigma^2 I)$  we get that  $(X^T X)^{-1} X^T \varepsilon \sim N(0, \Sigma)$  with  $\Sigma = \sigma^2 (X^T X)^{-1}$ . If we assume that  $\sigma$  is known, we have that

$$\frac{1}{\sigma^2} \varepsilon^T X (X^T X)^{-1} X^T \varepsilon = \frac{1}{\sigma^2} (\hat{f} - f^*)^T (X^T X) (\hat{f} - f^*) \sim \chi_q^2 , \quad (4)$$

i.e. follows a Chi-square distribution with  $q$  degrees of freedom. The first equality in (4) uses the fact that  $X^T X V = 0$ . A question for further study is how to construct the confidence intervals when we assume  $\sigma$  is not known. The typical procedure is to approximate  $\sigma$  with the sample standard deviation  $\hat{\sigma}$ . The distribution of this estimator and what it means for the distribution of expression (4) must be investigated.

We now show that  $(X^T X)^{-1} X^T \varepsilon$  is not contained in the subspace  $V$  by showing that it makes positive inner products with vectors in  $V^\perp$ .

**Proposition 2.** *If  $W = [w_1 \dots w_q]$ , then we have that  $w_i^T (X^T X)^{-1} X^T \neq 0$  for all  $i = 1, \dots, q$ .*

**Proof:**  $w_i^T (X^T X)^{-1} X^T = w_i^T W D^{-2} W^T X^T = \frac{1}{d_i^2} e_i^T W^T X^T = \frac{1}{d_i^2} w_i^T X^T$ , which must be  $\neq 0$ , since otherwise it would imply that  $X^T X w_i = 0$ , a contradiction. ■

Confidence intervals on our estimate  $\hat{f}$  are then given by

$$S(w) = \{f \mid (\hat{f} - f)^T (X^T X) (\hat{f} - f) \leq c_q(w) \sigma^2\} .$$

Taking into account that we have in fact multiple optimal estimates, represented by an ellipse, our confidence intervals are given by the following intersection of ellipses

$$S(w) \cap W ,$$

where  $W = VB(0, 1) \times (\text{Ker}(X^T X))^\perp$ .

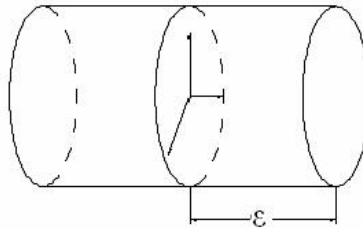


Figure 4: Schematic representation of the construction of a confidence interval.

Here we use the ball inside the ellipsoid to represent the  $\text{Ker}(X^T X)$ . We try to construct a ball with a small radius, let the center of the ball represent the true flow. The space with the circles are the balls in  $z$ . The distance of the center to the surface of ellipsoid is  $\varepsilon$ . Thus the parallel movement of these ellipses is given by the error terms. And we try to verify that most of  $f$  would be in the ellipse. Our aim is to calculate the confidence interval of  $\bar{f}$ .

### 3.3 Coordinate-wise confidence intervals

To build confidence intervals on each coordinate of  $\hat{f}$ , we construct the box that contains the confidence level ellipse  $S(w) \cap W$  found above.

This ellipsoidal set can also be represented as:

$$\left\{ f \mid f = \hat{f} + Vz + \delta, z^T H z \leq 1, \delta = (X^T X)^{-1} X^T \varepsilon, \varepsilon^T X (X^T X)^{-1} X^T \varepsilon \leq c_q(w) \sigma^2 \right\} .$$

So to find the upper (and lower) bound on coordinate  $\hat{f}_i$  we need to solve

$$\begin{aligned} & \max(\min) \quad e_i^T f \\ & \text{s.t.} \quad f = \hat{f} + Vz + (X^T X)^{-1} X^T \varepsilon \\ & \quad \quad z^T H z \leq 1 \\ & \quad \quad \varepsilon^T X (X^T X)^{-1} X^T \varepsilon \leq c_q(w) \sigma^2 \end{aligned} \tag{5}$$

or equivalently

$$\begin{aligned} & \max_{z, \varepsilon}(\min) \quad e_i^T Vz + e_i^T (X^T X)^{-1} X^T \varepsilon \\ & \text{s.t.} \quad z^T H z \leq 1 \\ & \quad \quad \varepsilon^T X (X^T X)^{-1} X^T \varepsilon \leq c_q(w) \sigma^2 \end{aligned} \tag{6}$$

which can be separated in  $z$  and  $\varepsilon$ :

$$\begin{aligned} & \max_z(\min) \quad e_i^T Vz \quad + \quad \max_\varepsilon(\min) \quad e_i^T (X^T X)^{-1} X^T \varepsilon \\ & \text{s.t.} \quad z^T H z \leq 1 \quad \quad \quad \text{s.t.} \quad \varepsilon^T X (X^T X)^{-1} X^T \varepsilon \leq c_q(w) \sigma^2 . \end{aligned} \tag{7}$$

The above quadratic problems have closed form solutions which yield an objective function value of  $\sqrt{e_i^T V H V^T e_i} + \sqrt{\sigma^2 c_q(w) e_i^T (X^T X)^{-1} e_i}$  where each term corresponds to the solution of each problem. In the case of minimization and lower bound the optimal objective function value is:  $-\sqrt{e_i^T V H V^T e_i} - \sqrt{\sigma^2 c_q(w) e_i^T (X^T X)^{-1} e_i}$ .

## 4 Numerical Experiments

We conduct three types of experiments: first, through controlled experiments, we examine how accurate our estimates are compared to the true flows; second, we construct an artificial example with data generated that matches real data; finally, we present examples of estimation of OD flows from real network flow data.

### 4.1 Controlled Experiment

The goal of our first experiment is to quantify how accurate are the ellipsoidal estimates constructed. We do this by randomly generating true OD flows and using these true flows to determine the observed link flows and the OD flow estimates obtained from these link flows. We are interested both in the distance between the analytic center estimate and the true OD flows and on how often is the true OD flow contained in the inscribed ellipse. This experiment therefore does not explore whether the true OD flow is contained in the confidence interval but whether it is contained in the set estimate obtained in the subspace of the null-space  $\text{Ker}(X^T X)$ .

We explore whether the mean and variance of the true OD flows or the geometry of the network have an influence on the quality of the estimate. We considered three different networks in this experiment: an intersection, an intersection with additional secondary flow along one of the directions, and a network in which users have path choice. The intersection network, depicted in Figure 2, considers that any inbound flow

can exit in any of the three other directions out of the intersection, since there are four different directions into the intersection this gives a total of 12 OD pairs. We also assume that we can observe the traffic flow at any of the 8 segments of road into and out of the intersection, Note that this example assumes we do not have access to turn counts and that there is no flow leaving or originating at the intersection itself. In the augmented intersection network, Figure 5, we consider two additional inflows and two outflows along one direction out of the intersection. This generates a problem with a total of 20 OD flows (by removing unreasonable flows such as (A,a) or (B,d) for example, see Figure 5), and 12 traffic flow observations. We consider that the flow to and from the secondary origins and destinations is significantly lower than the flow on the main roads to the intersection, this explores the effect of having very disparate flows in the network. Our

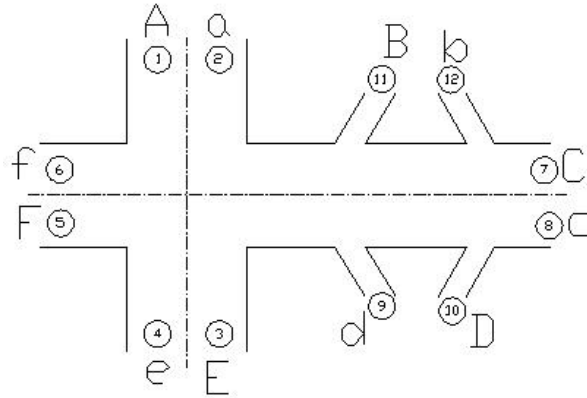


Figure 5: A multi-exit highway segment with a major intersection.

third example considers a network that allows more than one path for some OD pairs. For example, as it is shown in Figure 6, to travel from A to c, you can either pass the link flows 1, 11, 17 and 4, or link flows 1, 15, 13 and 4. This example considers a total of 20 OD flows and has 18 locations where traffic is quantified.

For each of the network examples above, we generate 100 random instances of true OD flow. We consider that the OD flow between each OD pair is independent and

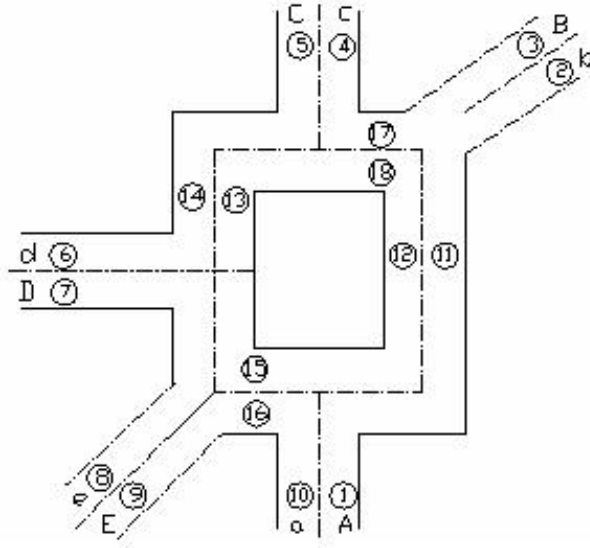


Figure 6: A network with multiple paths per OD pair.

identically distributed following a Normal distribution, with a given mean  $\mu$  and standard deviation  $\sigma$ . The only exception is in the augmented intersection, where the mean and standard deviation to the secondary, or branch, flows were assumed lower with mean  $b\mu$  and standard deviation  $b\sigma$ . We considered the above experiment for different mean values (from 500 to 2500) and standard deviations (from 100 to 500). In Table 1 we report the average results for all three networks for the different combinations of mean and standard deviation used to generate the true flow. We provide the relative distance between the analytic center estimate  $\hat{f}$  and the true flow  $f^*$ , given by  $\hat{d} = \frac{\|f^* - \hat{f}\|}{\|f^*\|}$ , averaged over the 100 repetitions. We also present for comparison, the relative distance of the regular unconstrained least squares estimate  $\bar{f}$ , given by  $\frac{\|f^* - \bar{f}\|}{\|f^*\|}$ , averaged over the 100 repetitions. Finally we also present in IN the number of times out of the 100 repetitions that the true flow was contained in the inscribed ellipsoid.

The results in Table 1 show that for all three networks as the standard deviation increases, the distance between the true flow and the estimated OD flows increases, also the likelihood that the true flow is contained in the inscribed ellipse decreases.

Table 1: Simulation results for three network examples. True flow generated with a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ .

OD flow		Intersection			Augmented Intersection					Multiple Paths		
parametrs.		$\hat{d}$	$\bar{d}$	IN	branch flow		$\hat{d}$	$\bar{d}$	IN	$\hat{d}$	$\bar{d}$	IN
$\mu$	$\sigma$				$b\mu$	$b\sigma$						
500	100	0.122	0.121	100	100	20	0.123	0.131	81	0.105	0.102	100
500	300	0.327	0.317	90	100	60	0.310	0.315	37	0.242	0.235	68
500	500	0.437	0.421	71	100	100	0.349	0.364	27	0.291	0.269	54
1500	100	0.043	0.042	100	300	20	0.043	0.047	100	0.036	0.037	100
1500	300	0.117	0.126	100	300	60	0.133	0.132	87	0.108	0.094	100
1500	500	0.203	0.192	100	300	100	0.216	0.213	49	0.184	0.170	95
2500	100	0.026	0.026	100	500	20	0.025	0.028	100	0.022	0.022	100
2500	300	0.075	0.071	100	500	60	0.075	0.082	100	0.066	0.061	100
2500	500	0.125	0.124	100	500	100	0.128	0.131	88	0.107	0.108	100

We also observed that when the true flow was not contained in the inscribed ellipse it was so because the estimates of a few OD flows were inaccurate while most OD flows estimates were close to the true values. As expected, even in this case the true OD flow is always contained in the circumscribed ellipse. We also note from Table 1 that the relative distance of the analytic center estimate  $\hat{d}$  is similar to the relative distance of the unconstrained least squares estimate  $\bar{d}$  in all three networks. So there is no sacrifice in accuracy in forcing the solution to satisfy the upper and lower-bound constraints. In summary, regardless of network structure or mean OD flow value, the relative distance between the estimate and true OD flows increases with the standard deviation of the true OD flow values. This observation also held on experiments with different distributions (uniform, lognormal) on the same three networks. In Figure 7 we show the increase in the relative distance from the true OD flow,  $\hat{d}$ , as a function of the standard deviation for the Intersection network and for different mean values. The graphs for the other networks considered are similar. We also observe in Figure 7 that as the mean value

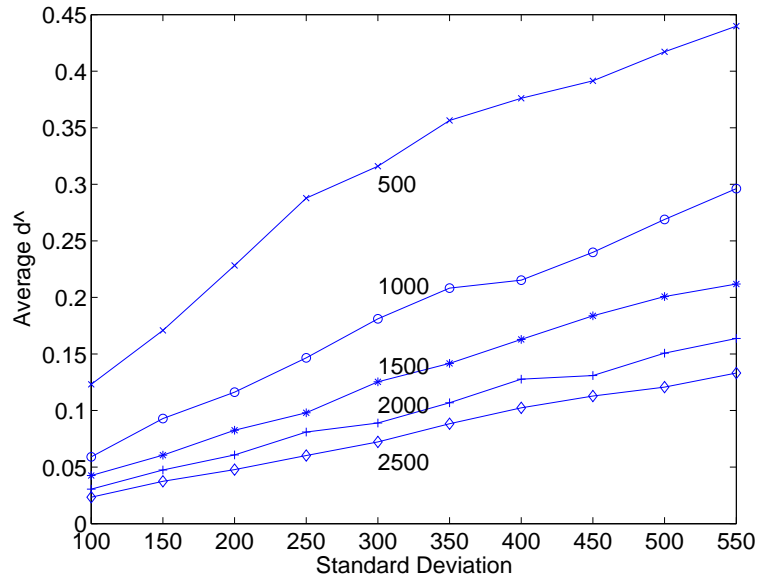


Figure 7: Average relative distance to true flow  $\hat{d}$  as a function of the standard deviation used to generate  $f^*$ , for different means.

increases, the relative error of the estimate decreases for the same standard deviation. This suggests that the relative error is indeed related to the coefficient of variation of the true OD flow, that is  $\frac{\sigma}{\mu}$  the standard deviation divided by the mean. We explore this relationship below.

We observe from Table 1 that the estimates for the Multiple Paths network typically lead to lower relative distances  $\hat{d}$ . Although slightly further away from the true flow, the estimates for the Intersection network were most effective at containing the true OD flow in the inscribed ellipse. Finally, although the relative distance  $\hat{d}$  obtained for Augmented Intersection network is similar to the Intersection Network, these estimates are considerably less accurate in containing the true OD flow in the circumscribed ellipsoid. We note that in the Augmented Intersection network, the OD flows to and from a secondary (or branch) segment had a significantly lower mean and standard deviation. This suggests that overall the OD flows in this network had significantly more variation than the other two examples. In Figure 8 we summarize the effect of the mean and standard deviation on the relative error by plotting the mean  $\hat{d}$  versus the coefficient of variation used to generate the OD flows. We observe that there is indeed a difference depending on which network example we are considering. We note that as the coefficient of variance increases so does the relative distance of the estimate  $\hat{d}$  for all three networks. This increase behaves linearly for small values of the coefficient of variation and tapers off for larger values.

Finally we note that the definition of  $\hat{f}$  involves an arbitrary upper bound on the true OD flow,  $U$ , that can be safely set to the maximum observed link flow value. This value however is a conservative estimate as all link flows are composed of multiple, and positive, OD flows. Which of the many  $f(z) = Vz + \bar{f}$  flows is the best estimate is difficult to predict. We explored the effect of reducing the upper-bound  $U$  on the quality of the analytic center estimate and found no significant trend if the value  $U$  is not set too close to the real upper bound of the OD flow.

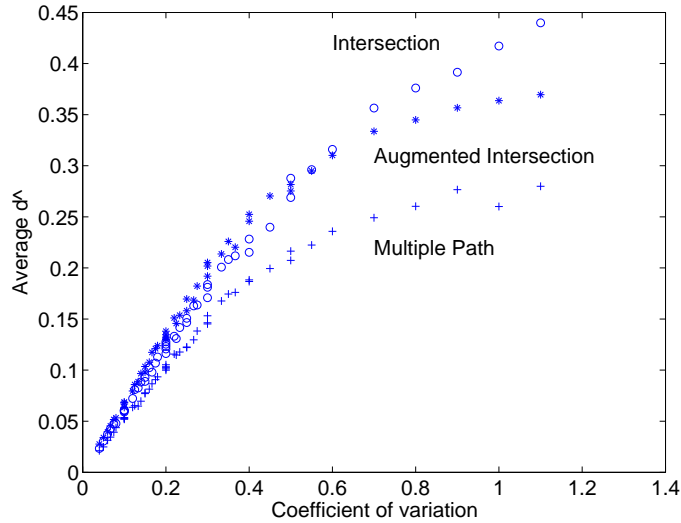


Figure 8: Average  $\hat{d}$  as a function of the coefficient of variation of the normal distribution used to generate  $f^*$ , for the three networks.

## 4.2 Experiments validating model with real data

Our second set of experiments investigates how accurate is our estimation model for data that behaves similarly to real traffic flow data, and what are the coordinate-wise confidence intervals we construct for this data. To achieve this we consider link flow data from a real intersection and construct randomly generated link flow data that matches the distribution of this real traffic data. We then use this synthetic data to study the accuracy of the estimation model and the form of the confidence intervals.

We obtained true link flow data from loop detectors under the Los Angeles freeway system from the PeMS website (<http://pems.eecs.berkeley.edu/Public>). This site provides data on traffic volume at a number of loop detectors in the Los Angeles metropolitan region throughout the day, and maintains repositories of this data that can be queried. Although there is a lot of data to work with, there are several issues with the data. Nevertheless, we found that the intersection between highways 405 and 10 in Los Angeles County has enough reliable loop detector data to compile traffic flow

data on every link flow in and out of the intersection. We downloaded link flow data for the morning commute, specifically data between 8am and 9am from Monday to Friday excluding holidays for dates between Fall 2003 and Summer 2004. We preprocessed this data removing the outliers (possibly due to bad weather or traffic accidents) in each link flow. After addressing these issues we still have between 78 to 150 reliable data observations for each link. In the first row of plots in Figure 9 we present the histogram of the true link flows for three representative link flows. This data is found to approximately follow a Weibull distribution, also depicted in the graphs.

To assess the efficiency of our estimation technique, we simulated OD flows to generate link flows that closely approximate the true link flows. We ensured the simulated link flows matched the mean and standard deviation of the true link flows through linear constraints on the mean and variance of the OD flows. These constraints represent the fact that each link flow is the sum of three OD flows in an intersection. The second row of plots in Figure 9 presents the simulated link flows for each of the true link flows presented in the first row of plots.

To stress the accuracy of our simulation, we present the summary statistics for all observed and simulated link flows in Table 2.

Making use of the simulated flows, we conduct the same type of experiment as before. We generate synthetic OD flows, which lead to link flows. These link flows are used to determine the ellipsoidal estimate of the OD flows. In Table 3 we present the relative distance between the analytic center estimate  $\hat{f}$  and the true flow simulated  $f^*$ ,  $\hat{d}$ ; the relative distance of the regular unconstrained least squares estimate  $\bar{f}$  and the true flow,  $\bar{d}$ ; we also present the percent of experiments where the true OD flow belongs to the inscribed ellipsoid.

We now explore the significance of the coordinate-wise confidence intervals that can be constructed from this link flow data. For this we construct a true flow  $f^*$  and

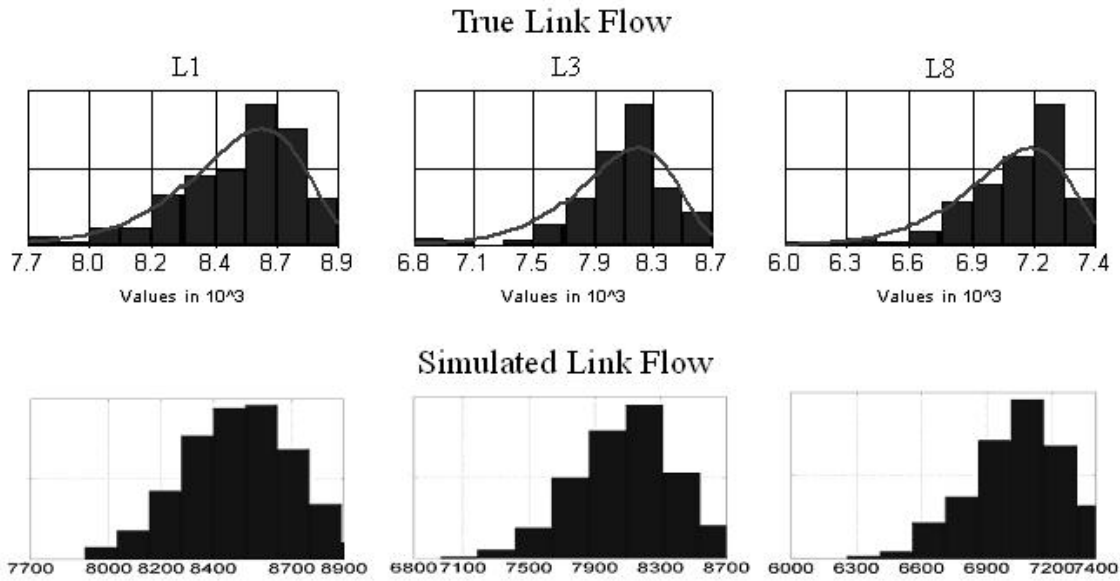


Figure 9: Sample comparison of true and simulated link flows for some loop detectors in the 405/10 intersection.

Table 2: Summary statistics of true and simulated link flows for the 405/10 intersection.

Link flow	Mean of True $L$	$\sigma$ of True $L$	Mean of Simulated $L$	$\sigma$ of Simulated $L$
1	8512	228	8502	215
2	7258	321	7273	331
3	8104	326	8110	334
4	7586	414	7597	410
5	10066	254	10065	261
6	7645	367	7639	378
7	7133	278	7131	279
8	7058	245	7058	236
Average	7920	304	7922	306

Table 3: Simulation results for OD flow estimation for data approximating the flow on the 405/10 intersection.

No. of experiment	Average $\hat{d}$	Average $\bar{d}$	% in $B(\hat{z}, 1)$
1000	0.0395	0.0445	100

coordinate-wise confidence intervals from link flow data. For link flow data we use the mean link flow data generated from 1000 simulation experiments to avoid extreme behavior. The true flow  $f^*$  is constructed by assuming that one third of the inflow to the intersection goes straight for every inbound direction. The remaining OD flows are then found solving a 8 by 8 linear system of equations. The coordinate-wise confidence intervals are obtained following Equation (7) which shows that the coordinate-wise limits are indeed made up by the sum of two different limits. The first due to the inscribed ellipse in  $\text{Ker}(X^T X)$  which is  $\omega_i^1 = \sqrt{e_i^T V H V^T e_i}$  for coordinate  $i$ , and the second due to the uncertainty from the actual data which is  $\omega_i^2 = \sqrt{\sigma^2 c_7(\omega) e_i^T (X^T X)^{-1} e_i}$  for coordinate  $i$ .

In Table 4 we present both the constructed true flow  $f^*$  and the analytic center estimate  $\hat{f}$  for each of the different coordinates, or origin destination pairs, in the intersection problem. For a 95% and a 80% confidence level, we also present the sum of the coordinate-wise limits  $\omega^1 + \omega^2$ , and the relative size of this limit with respect to the estimate  $\hat{f}$ . Note that the confidence interval for coordinate  $i$  is  $[\hat{f}_i - \omega_i^1 - \omega_i^2, \hat{f}_i + \omega_i^1 + \omega_i^2]$ . Finally we indicate whether the true flow is contained in each coordinate confidence interval in the column *Valid in*.

We observe that the difference between the true OD flows and the estimated analytic center solution is about 2.7%. It is noteworthy that all OD flows are inside the range of 80% confidence interval. Therefore to be certain with 80% probability that our confidence intervals contain the true flow we use an estimate with intervals that are less than 43% of the estimated value in every coordinate. The coefficient of variation of our

Table 4: Comparison of estimation and coordinate-wise confidence intervals to an artificial true flow  $f^*$  with mean link flows  $L$ .

OD pairs	$f^*$	$\hat{f}$	$\omega = 95\%$			$\omega = 80\%$		
			$\omega_1 + \omega_2$	$\frac{\omega_1 + \omega_2}{f_i} (\%)$	Valid in	$\omega_1 + \omega_2$	$\frac{\omega_1 + \omega_2}{f_i} (\%)$	Valid in
1-6	2937	2957	519	17.5	yes	433	14.6	yes
1-7	2904	2910	519	17.8	yes	433	14.9	yes
1-8	2753	2727	519	19.0	yes	433	15.9	yes
2-5	3214	3160	519	16.4	yes	433	13.7	yes
2-7	2132	2132	519	24.3	yes	433	20.3	yes
2-8	1974	2027	519	25.6	yes	433	21.4	yes
3-5	3491	3564	519	14.6	yes	433	12.1	yes
3-6	2436	2392	519	21.7	yes	433	18.1	yes
3-8	2252	2225	519	23.3	yes	433	19.5	yes
4-5	3262	3243	519	16.0	yes	433	13.4	yes
4-6	2217	2242	519	23.1	yes	433	19.3	yes
4-7	2184	2179	519	23.8	yes	433	19.9	yes

assumed true OD flow is 0.1735, which according to Figure 8, would lead to a  $\hat{d}$  smaller than 0.1, similar to what was observed in this case.

### 4.3 Estimation of OD flows from real link flow data

We now present results on the estimation and confidence intervals of OD demand for real link flow data. Here we use the observed link flow data for the intersection of highways 405 and 10 in Los Angeles County to construct the confidence intervals on the OD flows. We note that because of the quality of the data available from the PeMS website we had to consider this simple section of the highway system, as two factors have conspired to prevent us from obtaining good data for larger pieces of the highway system. First, the use of large sections of highway provides more opportunities for faulty loop detectors with missing data to be included in the example dataset. Second, the use of large highway sections exacerbates the effects of traffic dynamics. Since our model is static, these dynamic effects have the appearance of significant inconsistencies in the data.

We consider the same real link flow data described in the previous subsection; in fact we use the mean link flow (second column in Table 2) as a surrogate of representative link flow and as the input to the estimation model. In Table 5 we present the analytic center estimate obtained for each coordinate or OD pair and the width of the coordinate-wise confidence interval with 95% and 80% confidence. We also provide the percent of the estimate that each interval width represents. We note that these results are similar to the ones obtained with the simulated link flow data.

Table 5: Estimated flow given link flows equal to the mean of true link flows

OD pairs	$\hat{f}$	$\omega = 95\%$		$\omega = 80\%$	
		$\omega_1 + \omega_2$	$\frac{\omega_1 + \omega_2}{\hat{f}_i} (\%)$	$\omega_1 + \omega_2$	$\frac{\omega_1 + \omega_2}{\hat{f}_i} (\%)$
1-6	2987	526	17.6	439	14.7
1-7	2837	526	18.5	439	15.5
1-8	2743	526	19.2	439	16.0
2-5	3182	526	16.5	439	13.8
2-7	2093	526	25.1	439	21.0
2-8	2038	526	25.8	439	21.5
3-5	3559	526	14.8	439	12.3
3-6	2378	526	22.1	439	18.5
3-8	2222	526	23.7	439	19.8
4-5	3269	526	16.1	439	13.4
4-6	2224	526	23.6	439	19.7
4-7	2148	526	24.5	439	20.4

## 5 Conclusions

In this work, we propose an analytic center estimate and ellipsoidal confidence interval of OD pair flow obtained from link flow data for general transportation networks. We show that for this generally underdetermined problem there are naturally multiple estimation solutions, which we represent through the analytic center of the set of estimation solutions and an inscribed ellipse. We show that the confidence interval due to the data uncertainty is not contained in the subspace of the multiple estimation solutions, yielding an ellipsoidal confidence interval for the estimation. Our computational experiments show that the ellipsoidal estimate of the multiple estimation solutions is accurate, in particular for OD flows which have small coefficient of variation. The proposed estimation method provided tight coordinate-wise OD flow confidence intervals for the 405/10 highway intersection for real link flow data. We illustrate the dependency of this estimation method on the quality of link flow data through another real data example.

The proposed estimation method does not make use of additional assumptions on the behavior of the traffic flows and attempts to represent all possible realistic OD flows. Its use as part of planning or operational models would correctly represent the uncertainty on OD flows in a system. For example, OD flow estimation models are an important part of models to decide capacity expansion of a transportation network and traffic simulation. An estimate of OD demand identifies the actual origin and destination for each type of customer, as opposed to the manifestation of these trips constrained to the existing network present on the link flows. The presence of confidence intervals on these estimates allow modelers to enhance planning and simulation models to include the uncertainty present in these estimates. One such approach is to develop robust-optimization based planning or simulation models.

The main obstacle in developing confidence intervals and estimates of OD flows for

a large area is the availability of enough representative data for the area in question. In our study of the Los Angeles region we noticed that although the PeMS website provides substantial loop detector data, it is difficult to find in this data large swaths of the Los Angeles region with all loop detectors continuously providing accurate and reliable data.

There are two modeling enhancements that are planned as future work that could lead to more accurate estimates of the OD flows: the first is to incorporate dynamic aspects of the traffic flow in the estimation process. The current estimates are constructed off-line assuming all the flow traverses the network instantaneously. This approximation of what the link flow data represents likely increases the inaccuracies of the estimates. Second, the estimation procedure should be integrated with the decision of where to gather link flow data to improve the estimation accuracy.

## 6 Acknowledgements

The authors are grateful for the financial support for this research provided by the METTRANS Transportation Center through research grant #04-15.

## References

- Ashok, K. and M. Ben-Akiva (2000). Alternate approaches for real-time estimation and prediction of time-dependent origin-destination flows. *Transportation Science* 34(1), 21–36.
- Ashok, K. and M. Ben-Akiva (2002). 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.
- Bell, M. (1991). The estimation of origin-destination matrices by constrained generalized least squares. *Transportation Research Part B* 25, 13–22.

- Bell, M. G. H., C. M. Shield, F. Busch, and K. Kruse (1997). A stochastic user equilibrium path flow estimator. *Transportation Research Part C* 5, 197–210.
- Bierlaire, M. and F. Crittin (2003). An efficient algorithm for real-time estimation and prediction of dynamic OD tables. *Operations Research* 52(1), 116–127.
- Brenninger-Göthe, M., K. O. Jörnsten, and J. T. Lundgren (1989). Estimation of origin-destination matrices from traffic counts using multi-objective programming formulations. *Transportation Research Part B* 23, 257–269.
- Campbell, S. L. and C. D. Meyer (1991). *Generalized Inverses of Linear Transformations*. Dover Publications.
- Cascetta, E. (1984). Estimation of trip matrices from traffic counts and survey data: a generalized least squares estimator. *Transportation Research Part B* 18, 289–299.
- Cascetta, E., D. Inaudi, and G. Marquis (1993). Dynamic estimators of origin-destination matrices using traffic counts. *Transportation Science* 27(4), 363–373.
- Cascetta, E. and M. N. Postorino (2001). Fixed point approaches to the estimation of O/D matrices using traffic counts on congested networks. *Transportation Science* 35(2), 134–147.
- Doblas, J. and F. G. Benitez (2005). An approach to estimating and updating origin-destination matrices based upon traffic counts preserving the prior structure of a survey matrix. *Transportation Research Part B* 39, 565–591.
- Fisk, C. S. and D. E. Boyce (1983). A note on trip matrix estimation from link traffic count data. *Transportation Research Part B* 17, 245–250.
- Lo, H. P., N. Zhang, and W. H. K. Lam (1996). Estimation of an origin-destination matrix with random link choice proportions: a statistical approach. *Transportation Research Part B* 30, 309–324.
- Morrison, D. F. (1976). *Multivariate Statistical Methods*. McGRAW-HILL Book Company.

- Nie, Y., H. M. Zhang, and W. W. Recker (2005). Inferring origin-destination trip matrices with a decoupled gls path flow estimator. *Transportation Research Part B* 39, 497–518.
- Renegar, J. (2001). *A Mathematical View of Interior-Point Methods in Convex Optimization*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia.
- Robillard, P. (1975). Estimating the O-D matrix from observed link volumes. *Transportation Research* 9, 123–128.
- Sherali, H. D. and T. Park (1999). Estimation of dynamic origin-destination trip tables for a general network. *Transportation Research Part B* 35(3), 217–235.
- Sherali, H. D., R. Sivanandan, and A. G. Hobeika (1994). A linear programming approach for synthesizing origin-destination trip tables from link traffic volumes. *Transportation Research Part B* 28, 213–233.
- Spiess, H. (1987). A maximum likelihood model for estimating origin-destination matrices. *Transportation Research Part B* 21, 395–412.
- Van Aerde, M., B. Hellenga, L. Yu, and H. Rakha (1993). Vehicle probes as real-time ATMS sources of dynamic O-D and travel time data. In *Conference on Advanced Traffic Management Systems (ATMS)*, St. Petersburg, Florida, pp. 207–230.
- Van Aerde, M., H. Rakha, and H. Paramahamsan (2003). Estimation of O-D matrices: The relationship between practical and theoretical considerations. *Transportation Research Record No. 1831*, 122–130.
- Van Zuylen, J. J. and L. G. Willumsen (1980). The most likely trip matrix estimated from traffic counts. *Transportation Research Part B* 14, 281–293.
- Yang, H., T. Sasaki, Y. Iida, and Y. Asakura (1992). Estimation of origin-destination matrices from link traffic counts on congested networks. *Transportation Research Part B* 26, 417–434.