

1 **TRANSIT ROUTE ORIGIN-DESTINATION MATRIX ESTIMATION USING**  
2 **COMPRESSED SENSING**

7 **Pramesh Kumar**

8 Department of Civil, Environmental, and Geo-Engineering  
9 University of Minnesota - Twin Cities  
10 500 Pillsbury Drive S.E.  
11 Minneapolis, MN 55455, USA  
12 Tel: (716) 903-2366  
13 Fax: (612) 626-7750  
14 Email: [kumar372@umn.edu](mailto:kumar372@umn.edu)

16 **Alireza Khani, Corresponding Author**

17 Department of Civil, Environmental, and Geo-Engineering,  
18 University of Minnesota - Twin Cities,  
19 500 Pillsbury Drive S.E.,  
20 Minneapolis, MN 55455, USA  
21 Tel: (612) 624-4411  
22 Fax: (612) 626-7750;  
23 Email: [akhani@umn.edu](mailto:akhani@umn.edu)

25 **Gary A. Davis**

26 Department of Civil, Environmental, and Geo-Engineering,  
27 University of Minnesota - Twin Cities,  
28 500 Pillsbury Drive S.E.,  
29 Minneapolis, MN 55455, USA  
30 Tel: (612) 625-2598  
31 Fax: (612) 626-7750;  
32 Email: [drtrips@umn.edu](mailto:drtrips@umn.edu)

38 Word count: 4,985 words text + 1 table x 250 words (each) = 5,235 words

40 Submitted for publication in Transportation Research Record

41 Submission Date: November 15, 2018

1 **ABSTRACT**

2 Development of an origin-destination (OD) demand matrix is crucial for transit planning. With the  
3 help of automated data, it is possible to estimate a stop-level OD matrix. We propose a novel  
4 method for estimating transit route origin-destination (OD) matrix using Automatic Passenger  
5 Count (APC) data. The method uses  $l_0$  norm regularizer, which leverages the sparsity in the actual  
6 OD matrix. The technique is popularly known as compressed sensing (CS). We also discuss the  
7 mathematical properties of the proposed optimization program and the complexity of solving it.  
8 We use simulation to assess the accuracy and efficiency of the method and found that the proposed  
9 method is able to recover the actual matrix within small errors. With increased sparsity in the  
10 actual OD matrix, the solution gets closer to the actual value of the matrix. The method was found  
11 to perform more efficiently even for different demand patterns. We also present a real numerical  
12 example of OD estimation of A Line BRT route in Twin Cities, MN.

13  
14

15 *Keywords:* origin-destination (OD) matrix, transit, compressed sensing, Lasso, sparsity,  $l_0$  norm,  $l_1$   
16 norm, Automatic Passenger Count (APC), automated data

17  
18  
19  
20  
21  
22  
23

24 **The authors confirm contribution to the paper as follows: study conception and design:**  
25 **Pramesh Kumar, Alireza Khani, Gary A. Davis; data collection: Pramesh Kumar, Alireza**  
26 **Khani, Gary A. Davis; analysis and interpretation of results: Pramesh Kumar, Alireza**  
27 **Khani, Gary A. Davis; draft manuscript preparation: Pramesh Kumar, Alireza Khani, Gary**  
28 **A. Davis, All authors reviewed the results and approved the final version of the manuscript.**

29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

## 1 INTRODUCTION AND LITERATURE REVIEW

2 To understand the travel pattern of passengers, transit agencies require an origin-destination (OD)  
3 flow matrix of the passengers. This is a fundamental element of interest that helps in designing  
4 new routes and schedules, understanding and forecasting demand on transit network, adjusting  
5 marketing strategy, etc. The transit OD flow matrix is the quantification of the flow of passengers  
6 from one transit stop to another. To evaluate such a matrix, the agencies conduct on-board surveys,  
7 which collect data about passenger boarding and alighting stops, the purpose of travel, etc. These  
8 surveys are expensive to conduct and cover only a small sample of passengers (1). However, with  
9 recent advancement in automated data collection systems (ADCS), it is possible to mine the full  
10 origin-destination matrix. The automated data such as Automatic Fare Collection (AFC) and  
11 Automatic Passenger Count (APC) data are a rich source about information of passenger travel  
12 over a continuous period, making it possible to estimate OD matrix more frequently.

13  
14 The OD estimation problem has attracted the attention of many researchers over several decades.  
15 More recently, the use of automated data such as AFC, APC, or cell phone data has become  
16 popular for OD estimation. For example, the AFC data can be used to estimate a stop level OD  
17 matrix. It usually lacks the passenger alighting stop, which can be inferred using a trip-chaining  
18 algorithm based on several assumptions (2–5). The inference rate depends on the quality of data,  
19 the percentage of passengers using the smart card, and assumptions involved in the trip-chaining  
20 algorithm. On the other hand, APC systems collect information about the number of passenger  
21 boarding and alighting at each transit stop. OD estimation using the boarding and the alighting  
22 counts is a classic problem, which is hard to solve. The problem requires solving an  
23 underdetermined system of equations, in which case the number of unknowns to solve is far more  
24 than the number of equations available. Usually, multiple solutions are possible for this problem,  
25 which satisfy the given equations. Other information is supplemented to produce the accurate OD  
26 flows. To deal with this underdetermined problem, various methods have been proposed in the  
27 literature, which is summarized below:

- 28  
29 1. *Iterative Proportional Fitting (IPF) method*: This is a popular and easy to apply method to  
30 evaluate the OD matrix using count data (6, 7). The method starts with a base matrix,  
31 which is improved iteratively by multiplying the columns and rows of the matrix by a  
32 constant factor. The base matrix can be taken as a null matrix or any other seed matrix.  
33 Mishalani et al. found that using onboard survey data as a base matrix gives more accurate  
34 results than using null base matrix (8). The method has several issues such as the problem  
35 of non-structural zeros (6), due to which a zero entry remains zero in every iteration. The  
36 method also fails to converge if the number of zero entries become large in the matrix.
- 37 2. *Bayesian inference methods*: These methods use Bayesian approach to evaluate an OD  
38 matrix by formulating the problem as a partially observed Markov chain and utilizing prior  
39 information along with current observations of count data (9–12).
- 40 3. *Optimization methods*: As there are multiple solutions possible for this system of  
41 equations, these methods try to find the one, which optimizes an objective. The objective  
42 can be maximizing entropy (13) or the likelihood (14–16) function. With isotropic  
43 Gaussian noise, the maximum likelihood estimation turns into a classic least squares  
44 problem.

1 Another class of optimization methods consider above objectives along with a regularizer. The  
2 regularizer helps to mitigate the ill-posedness of the system of equations. The regularization can be  
3 included as a least square term between the unknown and a prior OD matrix obtained from a  
4 survey or from domain knowledge. This technique is quite popular in the literature. For example,  
5 Cascetta and Nguyen minimized generalized least square objective with a prior matrix (7), Van  
6 Zuylen and Willumsen maximized the relative entropy or minimized the Kullback-Leibler (KL)  
7 divergence of unobserved and observed flow distributions (13). This approach tries to force the  
8 solution, as close to the prior matrix as possible which may result in poor estimates if the prior or  
9 seed matrix used is not reliable.

10  
11 In this research, we evaluate the transit route OD matrix using APC data. The problem is the  
12 estimation of the flow of passengers between stops for a single trip. The route matrix problem has  
13 a special structure that provides an extra piece of information to reduce the ill-posedness of the  
14 system of equations. The estimation requires the selection of the correct estimate out of the  
15 multiple solutions. We use an estimation method that encourages the sparse OD matrix using  $l_0$   
16 norm regularizer. This helps in mitigating the ill-posedness of the system and offers interpretability  
17 (17) as there is only a subset of the origin-destination pairs which carries flow in an actual OD  
18 matrix. The method is popularly known as compressed sensing (18) and can also be viewed as  
19 least absolute shrinkage and selection operator (LASSO) regression proposed by (19).

20  
21 The rest of the paper is structured as follows. Section 2 presents the methodology for sparse matrix  
22 recovery followed by results of the experiments in section 3, then limitations of this research and  
23 directions for future research are discussed in section 4. Finally, conclusions are presented in  
24 section 5.

25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

## 1 METHODOLOGY

2 In this section, we present the method to estimate the route level OD matrix using boarding and  
3 alighting counts available from APC data. We use the following notations throughout the paper  
4 (Table 1):

5  
6  
7  
8  
9

10 Let  $N$  be the set of stops along a transit route at which passenger board or alight. We consider the  
11 boarding and alighting in a single direction. Let  $b_i$  and  $a_i$  be the observed number of passengers  
12 who board and alight at stop  $i = (1, 2, \dots, |N|)$  respectively. The values of  $b_i$  and  $a_i$  are obtained  
13 from APC data. Let  $X = \{x_{ij}\} \in R^{|N| \times |N|}$  be the origin-destination flow matrix, where  $x_{ij}$  denotes  
14 the number of passengers boarding at stop  $i$  and alighting at stop  $j$ . The overall setup is shown in  
15 Figure 1. Let  $x \in R^{|N|^2}$  be the vectorized form of matrix  $X$  i.e.,  $x = Vect(X)$ . The estimation  
16 procedure is subject to the following constraints, which are taken from (11, 20).

17  
18  
19  
20  
21

### 22 Constraints

23 1. If we sum the values of  $x_{ij}$  along all the columns, then we get the total number of passengers  
24 boarding at stop  $i$  i.e.  $b_i$ ,

$$25 \quad \sum_{j=1}^{|N|} x_{ij} = b_i \quad \forall i \in N \quad (1)$$

26 2. Similarly, if we sum the values of  $x_{ij}$  along all the rows, then we get the total number of  
27 passengers alighting at stop  $j$ , i.e.,  $a_j$

$$28 \quad \sum_{i=1}^{|N|} x_{ij} = a_j \quad \forall j \in N \quad (2)$$

29

30 3. The total number of boarding at all the stops should be equal to the total number of alighting.

$$31 \quad \sum_{j=1}^{|N|} b_j = \sum_{i=1}^{|N|} a_i \quad (3)$$

32

33 4. The number of boarding and alighting at the same stop is zero, which means the diagonal  
34 elements of the matrix  $X$  should be equal to zero.

35

$$36 \quad x_{ii} = 0 \quad \forall i \in N \quad (4)$$

37

38 5. As the transit vehicle runs in a single direction, a passenger boarding at one stop cannot alight at  
39 the previous stops that vehicle has already visited. This means,

40

$$41 \quad x_{ij} = 0 \quad \forall i > j, \quad \forall i, j \in N \quad (5)$$

42

43 7. The total load on a link between two stops is equal to the passengers boarding between those  
44 stops.

$$45 \quad \sum_{i=1}^k (b_i - a_i) = \sum_{i=1}^k \sum_{j=k+1}^n x_{ij} \quad (6)$$

1  
2 By imposing these constraints, the structure of the matrix will look as in Figure 2.  
3  
4  
5  
6  
7  
8

9 We can express the linear constraints (1)-(6), in form of a matrix as

$$10 \quad \mathcal{A}(x) = b \quad (7)$$

11  
12  
13 Where,  $\mathcal{A} \in R^p \times |N|^2$  is the linear map (which is a matrix in this case) for  $p$  number of  
14 constraints and  $b \in R^p$  represents the constant vector for these constraints. In the following  
15 subsections, we describe the proposed solution to the given problem.  
16

### 17 **Transit route OD estimation using compressed sensing technique**

18 As discussed, (7) is usually an ill posed problem for which one can expect multiple solutions. A  
19 generic regularizer can help in mitigating the ill posedness of the problem. One such regularizer is  
20 the generalized least square of the difference between the unknown and a prior matrix obtained from  
21 the survey data. The quality of the solution depends upon the availability of a good prior matrix as  
22 the optimal solution is forced to be as close as possible to the prior matrix. We can use other  
23 regularizers based on the domain knowledge on the space of the plausible OD flows in the network  
24 (17). To use one such regularizer, we make the following assumption:  
25

26 **Assumption** *The planted OD matrix in the set of linear equations is sparse which means that the*  
27 *flow between many of the OD pairs should be equal to zero. The observed flow is only due to a*  
28 *small subset of  $\frac{N(N-1)}{2}$  pairs.*

29  
30 The intuition behind the above assumption is that there is a large number of OD pairs for a transit  
31 route, but the travel happens only along few pairs. For example, during the morning peak hours,  
32 there are only a few popular origin stops such as residential locations and few destination stops  
33 such as central business areas, park and rides, etc. Moreover, it is unlikely that passengers boarding  
34 at initial stops of the route will alight at all the following stops. This makes the flow between most  
35 of the OD pairs equal to zero. This is opposite to the solution evaluated using entropy  
36 maximization, which tries to achieve the solution, as uniform as possible to minimize the errors.  
37 The sparsity as a regularizer has been used before for highway network OD estimation and has  
38 found promising results (17, 21–24). For example, (17) leverages sparsity in highway OD matrix  
39 to estimate a set of suitable traffic analysis zones (TAZs) and use those zones to evaluate an OD  
40 matrix. The method proposed in (17) has a bi-level structure with sparse OD estimation on upper  
41 level and traffic assignment using user equilibrium at lower level. The use of non-negativity  
42 constraints for improving the solution is also emphasized. This paper uses similar optimization for  
43 the transit route OD estimation problem, which has a special structure as we get an extra set of  
44 constraints because of transit movement in one direction. We also describe the conditions under  
45 which sparse recovery is possible.  
46  
47

1 *Using sparsity as the regularizer for OD estimation*

2 To achieve the sparsity in the solution, we minimize the number of non-zero entries in the solution,  
 3 which can be done by minimizing  $l_0$  norm of the vector  $x$ . We can state the problem as the  
 4 minimization of  $l_0$  norm of  $x$  subject to linear constraints. The optimization formulation is given  
 5 below:

$$\begin{aligned} & \text{minimize } \|x\|_0 \\ & \text{s. t. } \mathcal{A}(x) = b \\ & \quad x \geq 0 \end{aligned} \tag{8}$$

11 Where,  $\|x\|_0$  is the  $l_0$  norm of vector  $x$ , which is defined as  $\lim_{p \rightarrow 0} \sum_j |x_j|^p$ . The non-negativity  
 12 should not be dropped from (8) as it helps to mitigate the ill-posedness of the problem (17). Using  
 13 Lagrangian relaxation, the linear constraints can be included in the objective function as a least  
 14 square term and formulated as following:

$$\text{minimize}_{x \geq 0} \|\mathcal{A}(x) - b\|_2 + \mu \|x\|_0 \tag{9}$$

18 The problem (9) tries to find the sparse vector  $x$  planted in the given ill-posed system of linear  
 19 equations. The regularization parameter  $\mu$  controls the sparsity of the vector and requires tuning to  
 20 get the best results. A higher value of the  $\mu$  will impose more sparsity in the solution. When  $\mu = 0$ ,  
 21 (9) reduces to an ordinary least squares problem. The optimization program (9) is useful for the  
 22 APC data when the total number of boarding and alighting do not match as the least square term  
 23 will try to find a solution which best explains the observed flows. This happens quite often in the  
 24 APC systems due to the errors in recording data. The given problem (9) is an NP-hard as the  
 25 minimization of  $l_0$  norm cannot be done in polynomial time. Recent work in compressed sensing  
 26 has proposed a tightest convex relaxation of the  $l_0$  norm which is  $l_1$  norm (25). The problem (9)  
 27 can be restated as follows.

$$\text{minimize}_{x \geq 0} \|\mathcal{A}(x) - b\|_2 + \mu \|x\|_1 \tag{10}$$

31 where,  $\mu \|x\|_1 = \sum_i |x_i|$ . (10) is a convex optimization program as the absolute value of  $x_i$  can be  
 32 written as a set of linear inequality constraints. The use of  $l_1$  norm is better than the  $l_2$  norm (also  
 33 called ridge regression) to achieve sparsity. This is because the  $l_1$  norm ball has corner points that  
 34 can intersect the given plane at the sparsest solutions, unlike  $l_2$  norm ball. The problem can also be  
 35 viewed as least absolute shrinkage and selection operator (or Lasso regression) proposed by (19)  
 36 as given a set of observations, we try to estimate the coefficients which satisfies the given  
 37 equations. However, there is a key difference between compressed sensing and LASSO. The  
 38 former provides conditions under which the linear map  $\mathcal{A}$  is nicely behaved and the uniqueness of  
 39 the solution can be proved (these conditions are discussed in the next subsection). In other words,  
 40 we can design  $\mathcal{A}$  in such a way that it can guarantee to recover the actual solution. On the other  
 41 hand, LASSO is a regression method in which we have no control over the data and we try to find  
 42 the best coefficients which are sparse and satisfy the equations obtained from data. We can also  
 43 interpret these estimates as a Bayesian posterior mode estimate when the regression parameters  
 44 have independent Laplace (i.e., double exponential) priors (26). Now the natural question which  
 45 arises is that when does solving (10) gives a good solution to (9). In other words, what natural  
 46 conditions can be applied on linear map  $\mathcal{A}$  so that we can say that the solution is unique. Candès  
 47 and Tao, 2005 proposed the idea of restricted isometry property (RIP) of the matrices, which states

1 that if  $\mathcal{A}$  satisfies the isometry property, then there exists a unique solution to the problem (10)  
 2 which is equal to the solution of (9).

### 3 *Restricted Isometry Property (RIP)*

4 The linear map  $\mathcal{A}$  has RIP with constants  $k$  and  $\delta_k$ , if  $\forall \|x\|_0 \leq k$ ,  $\mathcal{A}$  behaves almost as an  
 5 isometry in following sense i.e.  $l_2$  norm of  $\mathcal{A}(x)$  is close to the  $l_2$  norm of vector  $x$ :  
 6

$$7 \quad (1 - \delta_k)\|x\|_2^2 \leq \|\mathcal{A}(x)\|_2^2 \leq (1 + \delta_k)\|x\|_2^2 \quad (11)$$

8  
 9  
 10 RIP matrices are extremely common in practice and most of the random matrices satisfy this  
 11 property. Based on the above definition, a theorem is proposed by Candès and Tao, 2005 (25).  
 12

### 13 **Theorem (Candès and Tao, 2005 (25))**

14 If  $\mathcal{A}(x) = b$  and  $b$  is constructed using a sparse solution with  $\|x\|_0 \leq k$ , and the RIP condition is  
 15 satisfied with constants  $\delta_{2k}$  and  $\delta_{3k}$ , satisfying  $\delta_{2k} + \delta_{3k} < 1$ , then (10) can obtain a unique  
 16 solution to the problem (9) with as few as  $\mathcal{O}\left(k \log\left(\frac{|N|^2}{k}\right)\right)$  number of equations  
 17

18 As the passenger flow cannot be negative, we can replace the  $l_1$  norm with sum of the components  
 19 of vector  $x$ , which allow us to use the gradient-based approaches to solve the optimization  
 20 program (10) efficiently. If we have some idea about the number of non-zero entries (say less than  
 21  $k$ ), we can constraint the solution as follows:  
 22

$$23 \quad \begin{aligned} & \text{minimize } \|\mathcal{A}(x) - b\|_2 \\ & \text{s.t. } \|x\|_0 \leq k \\ & \quad x \geq 0 \end{aligned} \quad (12)$$

24 We use the optimization program (8) with  $l_1$  norm for solving the transit route OD estimation  
 25 problem. The problem is convex and can be solved easily using a standard convex optimization  
 26 solver such as CVX (27). We could also employ an iterative algorithm proposed in (19) to evaluate  
 27 a sparse solution but the algorithm does not gurantee convergence to a unique solution. In the next  
 28 section, we present numerical examples to show the application of the porposed method.  
 29  
 30  
 31

## 32 **RESULTS**

33 In this section, we present two numerical examples of OD estimation using the proposed  
 34 methodology. First, simulation is used to assess the consistency and accuracy of the estimation  
 35 method. Second, the OD estimation of a bus route in Twin Cities, MN is presented.  
 36  
 37

### 38 **OD estimation using simulation**

39 We use the APC data provided by Metro Transit, which is a primary transit service provider in the  
 40 Twin Cities, MN area offering an integrated network of buses, light rail, bus rapid transit, and a  
 41 commuter train. To prepare a synthetic OD matrix, we make some assumptions on the probability  
 42 distribution of arrival of the passengers on different stops. We consider 10 stops along a transit  
 43 route to facilitate the presentation of results. The passenger arrival at the stop is assumed to follow  
 44 a Poisson distribution.  
 45

$$46 \quad b_i \sim \text{Poisson}(k) \forall i \quad (13)$$



1  
2 where  $k$  is the mean arrival rate at the stop and  $b_i$  is the number of boarding at stop  $i$ . We  
3 recommend fitting a Poisson distribution to the real data to calculate the value of  $k$ . We calculated  
4 the mean value of arrival rate of the passengers on the A line, an arterial BRT route in Twin Cities,  
5 MN. The mean arrival rate was found to be equal to 0.86 during peak hours, which is quite low. To  
6 assess the significant errors produced by the estimation, we assumed the mean value equal to 15  
7 passengers. Then we set the sparsity level for the O-D matrix. The sparsity level will make the  
8 value of probability of flow from one stop to another stop zero if this probability value is less than  
9 the threshold sparsity level. This is done to create sparsity in the matrix and to test whether the  
10 method works more efficiently when the sparsity is high. We then assign the flow from one stop to  
11 others by assuming a multinomial distribution i.e.

$$12 \quad x_{ij} \sim MNL(b_i, p_{i1}, p_{i2}, \dots, p_{i|N|}) \quad (14)$$

13  
14 Where  $p_{ij}$  is the probability of movement from stop  $i$  to stop  $j$ . We then make the diagonal and  
15 lower triangle of the matrix zero because of the constraints (4-5). To calculate the boarding and  
16 alighting flows for O-D estimation, we sum the rows and columns of the simulated matrix. After  
17 that, we set up an optimization model using the python API of CVX (27). To avoid choosing the  
18 value of  $\mu$  in optimization program (10), we solved the program (8) with  $l_1$  norm. However, we  
19 recommend using the optimization program (10) when the sum of boarding and alighting count do  
20 not match in the APC data, which happens because of the errors in data collection. We generated  
21 200 Monte-Carlo samples of OD matrices and calculated the  $l_2$  error between the actual OD  $x$  and  
22 estimated OD vector  $x_{est}$  as:  
23  
24

$$25 \quad \|x - x_{est}\|_2 = \sqrt{\sum_i (x_i - x_{est,i})^2} \quad (15)$$

26  
27  
28  
29  
30  
31  
32 Figure 3 shows the histogram of  $l_2$  error in the estimation for each sample. We can observe that  
33 the mean value of the error is 4.99 and with a standard deviation of 1.32. The 95% confidence  
34 interval of the  $l_2$  error was found to be equal to (4.81, 5.19). This shows that the results obtained  
35 from this estimation method are consistent and small. To see how the method performed in  
36 predicting the individual origin-destination pair flow value, we created a box plot for the  
37 estimation error (Figure 4). The proposed method predicted the actual value of the non-zero entries  
38 41.5 % of the time. In case of errors, the method seems to overpredict the values except some of the  
39 O-D pairs such as 0-4, 1-2 and 5-6.  
40  
41  
42  
43  
44  
45  
46

1  
2  
3  
4 Figure 5(a) shows the average load profile of the passengers on the transit route. The width of the  
5 95% confidence interval is small which shows that the method is reliable in estimating demand and  
6 therefore in deciding the adequate frequency to handle the load of the passengers. We can also  
7 observe that the errors in estimating the load of the passengers is also quite small (Figure 5(b)).  
8  
9  
10  
11  
12  
13  
14

15 To understand the effect of sparsity, we solved the problem for several levels of sparsity and  
16 calculated the root mean square error (RMSE) between the estimated and actual OD matrix. Figure  
17 6 shows the RMSE value with respect to the sparsity in the matrix. We can observe that the RMSE  
18 value is reduced with increased sparsity. For example, when the OD matrix has only 10% non-zero  
19 values, the corresponding RMSE value was found to be less than 0.35, which is quite impressive.  
20 This shows that the accuracy of the method is improved when there is more sparsity. Comparing  
21 the results to common least squares solution (Figure 6), the proposed method is able to recover  
22 solutions with lower RMSE value. It can also be observed that when the sparsity is low, the  
23 proposed method is more efficient than least squares as the gap between two lines is high but when  
24 there are a greater number of non-zero entries, the RMSE gap between these two methods reduces.  
25

26  
27  
28  
29  
30  
31 To see how different demand patterns affect the OD estimation, we performed a similar simulation  
32 for several mean arrival rates( $k$ ) of passengers at stops. Figure 7 shows normalized RMSE values  
33 with respect to sparsity in the random matrix for different mean arrival rate. The normalization is  
34 done by simply dividing RMSE by mean arrival rate. The results are presented in separate panels.  
35 We can see that the normalized RMSE decreases with an increase in demand. At  $k = 2$ , the  
36 normalized RMSE value was found to be almost equal to 1, which is still quite low. At lower  
37 demand, the matrix is already sparse, so we see less effect of sparsity parameter.  
38  
39  
40  
41  
42  
43

#### 44 **OD estimation of A Line BRT route in Twin Cities**

45 We use the APC data from Twin Cities, MN to calculate the origin-destination flow of a route. The  
46 Automatic Passenger Count (APC) data used for this research contains transit trip information,  
47 such as date and time of the operation, routeID, stopID, departure and arrival time, number of  
48 boarding and alighting on each stop, and the geographical coordinates of the stops. We select A  
49 Line, which is a bus rapid transit (BRT) route in Twin Cities for this analysis. It serves 20 stations

1 along Snelling Av and 46<sup>th</sup> St. We select a trip from the data during peak hour. The number of  
 2 boarding and alighting at different stops in the northbound direction is shown in Figure 8. We can  
 3 observe the popular boarding locations such as 46<sup>th</sup> street station, 46<sup>th</sup> & Minnehaha station, and  
 4 Snelling & Highland station and alighting stops such as Rosedale transit center, Snelling &  
 5 Highland station and Snelling & Clair st. station. We use the optimization program (10) to solve  
 6 the given problem with a value of  $\mu = 0.2$ . Some recommendations for choosing the value of  $\mu$  is  
 7 given in (19).

8  
 9  
 10  
 11  
 12  
 13  
 14 The total ridership of the trip is 16. Because of low ridership, flow along most of the O-D pairs  
 15 should be equal to zero. We apply the proposed method to the given data and calculate the  
 16 origin-destination flows. Figure 9 shows the origin-destination flows between different O-D pairs.  
 17 We can see that the flow occurred only between 11 O-D pairs out of 400 pairs (2.75%). The highest  
 18 flow was observed between Snelling & Highland Av and Rosedale Transit Center, which is the last  
 19 station along this route. Other popular OD pairs are 46<sup>th</sup> St and Snelling & St. Clair, Snelling &  
 20 Minnehaha and Snelling & Highland Av. Because of the low ridership, the sparse matrix recovery  
 21 seems to perform well.

## 22 23 24 25 26 27 28 **LIMITATIONS AND DISCUSSION**

29 In this section, we discuss the limitations of the proposed method and provide some  
 30 recommendations for future research to address these limitations. Various studies use a prior  
 31 matrix as a regularizer which can also be included in the proposed framework as follows:

$$32 \quad \text{minimize}_{x \geq 0} \quad \|\mathcal{A}(x) - b\|_2 + \mu_1 \|x\|_0 + \mu_2 \|x - x^{prior}\|_2 \quad (16)$$

33  
 34  
 35 The choice of parameters  $\mu_1$  and  $\mu_2$  will control the weight of different objectives which can be  
 36 obtained by observing the error rate for different values of these parameters. Due to lack of the  
 37 suitable prior matrix, we have not included any results using this program in this study. Also, no  
 38 unique choice of  $\mu_1$  and  $\mu_2$  can make this method unattractive to practitioners.

39  
 40 The problem of OD estimation has been well studied in the literature both in the context of road  
 41 and transit network. This is an interesting problem with solution methods using both optimization  
 42 and statistical techniques. Depending on the available data, the problem can be formulated as an  
 43 underdetermined or overdetermined system of equations. The classic techniques such as entropy  
 44 maximization, least squares, etc. produce some good results but may not evaluate the correct  
 45 solution as there can be infinitely many or no solutions, which again depends on the quality of data  
 46 and the set of equations obtained from the setup. Recent work in the field of compressed sensing  
 47 has established that under suitable conditions, we can evaluate a unique sparse solution out of

1 given ill-posed system of linear equations. We tried to solve the problem using this approach and  
2 found impressive results but not an exact solution. We also did not prove that the linear mapping  
3 produced by a set of linear equations in this case satisfies the restricted isometry property, which  
4 may be the source of error in our results. The condition is hard to prove and can be a future  
5 research topic.

6  
7 The compressed sensing technique can also be used to design a linear map  $\mathcal{A}$  so that we can  
8 guarantee an exact solution to this problem. The future work in this regard should be focused on  
9 how to engineer a system in order to create a linear map  $\mathcal{A}$  so that the recovery of an exact solution  
10 can be guaranteed. This can be done by finding optimal sensor locations on highway and transit  
11 network or integrating different data sources to produce an appropriate  $\mathcal{A}$ .

12  
13 This research can also be expanded in multiple directions. The method can be used to estimate a  
14 full transit network OD matrix. The problem can be formulated as a bi-level program with sparse  
15 recovery optimization at the upper level and transit assignment (28, 29) at the lower level to  
16 capture route choice behavior in the model. We believe that the network level OD will also be  
17 sparse because it is unlikely that passengers boarding at one stop can alight at all other stops in the  
18 network. The concept can also be extended to matrix sensing which will be helpful in estimating a  
19 time-dependent transit OD matrix. As the boardings and alightings follow a regular pattern during  
20 various hours of the day, data from several days can be used to learn this pattern. This means the  
21 high dimensional data for several days can be used to minimize the rank of the matrix to extract a  
22 regular pattern. This can be done by minimizing the nuclear norm of the matrix, which is a convex  
23 surrogate for the rank of the matrix. The problem is computationally challenging and needs further  
24 attention.

## 25 26 CONCLUSIONS

27 In this research, we proposed a method for estimating an origin-destination OD matrix for a transit  
28 route along one direction. The problem can be formulated as an undetermined system of linear  
29 equations. The adopted strategy was to estimate a sparse O-D matrix, using  $l_0$  norm. Using its  
30 convex surrogate  $l_1$  regularizer, the problem can be solved efficiently. The sparsity in the matrix is  
31 generated because there are only a few popular O-D pairs along a transit route where the flow  
32 occurs. We also discussed the complexity of solving the proposed optimization program. The  
33 constraints and sparsity try to force the solution to an actual value. We tested the efficiency of the  
34 estimator using simulation. The errors were found to be bound within a small range. With an  
35 increased level of sparsity in the matrix, the method was able to recover more accurate results. We  
36 also found small errors even for higher demand. For example, the normalized RMSE between  
37 estimated and actual matrix value was found to be at most 0.1. We also presented a numerical  
38 example for A-line BRT route in Twin Cities, MN. Finally, we discussed various limitations and  
39 directions for future research in section 4. Further studies are required to show under which  
40 constraints, the OD linear map satisfy the RIP property. Other statistical methods are also required  
41 to assess the accuracy of the estimation.

## 42 43 44 ACKNOWLEDGEMENT

45 This research was conducted at the University of Minnesota Transit Lab  
46 (<http://umntransit.weebly.com/>), currently supported by the following, but not limited to, projects:

- 47 - National Science Foundation, award CMMI-1637548

- 1 - Minnesota Department of Transportation, Contract No. 1003325 Work Order No. 15
- 2 - Minnesota Department of Transportation, Contract No. 1003325 Work Order No. 44
- 3 - Transitways Research Impact Program (TIRP), Contract No. A100460 Work Order No.
- 4 UM2917

5 The authors are grateful to Metro Transit for sharing the data. We are also grateful to the  
 6 anonymous referees for their constructive input to improve the quality of this article. Any  
 7 limitation of this study remains the responsibility of the authors.

## 9 REFERENCES

- 10 1. Attanucci, J. & Wilson, N. H. M. Bus Transit Monitoring Manual: Volume 1: Data Collection  
 11 Program Design. *US Department of Transportation*, Vol. 1, No. August, 1981.
- 12 2. Kumar, P., A. Khani, and Q. He. A Robust Method for Estimating Transit Passenger Trajectories  
 13 Using Automated Data. *Transportation Research Part C: Emerging Technologies*, Vol. 95, 2018.  
 14 <https://doi.org/10.1016/j.trc.2018.08.006>.
- 15 3. Nassir, N., A. Khani, S. Lee, H. Noh, and M. Hickman. Transit Stop-Level Origin-Destination  
 16 Estimation Through Use of Transit Schedule and Automated Data Collection System.  
 17 *Transportation Research Record: Journal of the Transportation Research Board*, Vol. 2263, 2011,  
 18 pp. 140–150. <https://doi.org/10.3141/2263-16>.
- 19 4. Trépanier, M., N. Tranchant, and R. Chapleau. Individual Trip Destination Estimation in a Transit  
 20 Smart Card Automated Fare Collection System. *Journal of Intelligent Transportation Systems*, Vol.  
 21 11, No. 1, 2007, pp. 1–14. <https://doi.org/10.1080/15472450601122256>.
- 22 5. Zhao, J., a Rahbee, and N. H. . Wilson. Estimating a Rail Passenger Trip Origin-Destination Using  
 23 Automatic Data Collection Systems. *Computer-Aided Civil and Infrastructure Engineering*, Vol. 22,  
 24 No. 5, 2007, pp. 376–387.
- 25 6. Ben-Akiva, E. al. Alternative Methods to Estimate Route-Level Trip Tables and Expand On-Board  
 26 Surveys. *Transportation Research Record 1037, TRB, National Research Council, Washington,*  
 27 *D.C.*, p. 1–11.
- 28 7. Cascetta, E., and S. Nguyen. A Unified Framework for Estimating or Updating Origin/Destination  
 29 Matrices from Traffic Counts. *Transportation Research Part B*, Vol. 22, No. 6, 1988, pp. 437–455.  
 30 [https://doi.org/10.1016/0191-2615\(88\)90024-0](https://doi.org/10.1016/0191-2615(88)90024-0).
- 31 8. Mishalani, R., Y. Ji, and M. McCord. Effect of Onboard Survey Sample Size on Estimation of  
 32 Transit Bus Route Passenger Origin-Destination Flow Matrix Using Automatic Passenger Counter  
 33 Data. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. 2246,  
 34 2011, pp. 64–73. <https://doi.org/10.3141/2246-09>.
- 35 9. Maher, M. J. Inferences on Trip Matrices from Observations on Link Volumes: A Bayesian  
 36 Statistical Approach. *Transportation Research Part B: Methodological*, Vol. 17, No. 6, 1983, pp.  
 37 435–447.
- 38 10. Gary A . Davis. Estimating Freeway Demand Patterns and Impact of Uncertainty on Ramps  
 39 Controls. *ASCE Journal of Transportation Engineering*, Vol. 119, No. 4, 1993, pp. 489–503.
- 40 11. Li, B. Markov Models for Bayesian Analysis about Transit Route Origin-Destination Matrices.  
 41 *Transportation Research Part B: Methodological*, Vol. 43, No. 3, 2009, pp. 301–310.  
 42 <https://doi.org/10.1016/j.trb.2008.07.001>.
- 43 12. Hazelton, M. L. Statistical Inference for Transit System Origin-Destination Matrices.  
 44 *Technometrics*, Vol. 52, No. 2, 2010, pp. 221–230. <https://doi.org/10.1198/TECH.2010.09021>.
- 45 13. Van Zuylen, H. J., and L. G. Willumsen. The Most Likely Trip Matrix Estimated from Traffic  
 46 Counts. *Transportation Research Part B: Methodological*, Vol. 14, No. 3, 1980, pp. 281–293.  
 47 [https://doi.org/10.1016/0191-2615\(80\)90008-9](https://doi.org/10.1016/0191-2615(80)90008-9).
- 48 14. Nihan, N. L., and G. A. Davis. Application of Prediction-Error Minimization and Maximum  
 49 Likelihood to Estimate Intersection O-D Matrices from Traffic Counts. *Transportation Science*. 2.  
 50 Volume 23, 77.

- 1 <http://ezproxy.unicartagena.edu.co:2048/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=a9h&AN=6798710&lang=es&site=ehost-live>.
- 2
- 3 15. Spiess, H. A Maximum Likelihood Model For Estimating Origin-Destination Matrices.
- 4 *Transportation Research Board*, Vol. 21B, No. 5, 1987, pp. 395–412.
- 5 [https://doi.org/10.1016/0191-2615\(87\)90037-3](https://doi.org/10.1016/0191-2615(87)90037-3).
- 6 16. Vardi, Y. Network Tomography: Estimating Source-Destination Traffic Intensities from Link Data.
- 7 *Journal of the American Statistical Association*, Vol. 91, No. 433, 1996, pp. 365–377.
- 8 <https://doi.org/10.1080/01621459.1996.10476697>.
- 9 17. Menon, A. K., C. Cai, W. Wang, T. Wen, and F. Chen. Fine-Grained OD Estimation with Automated
- 10 Zoning and Sparsity Regularisation. *Transportation Research Part B: Methodological*, Vol. 80,
- 11 2015, pp. 150–172. <https://doi.org/10.1016/j.trb.2015.07.003>.
- 12 18. Candes, E. J., and M. B. Wakin. An Introduction To Compressive Sampling. *IEEE Signal*
- 13 *Processing Magazine*, Vol. 25, No. 2, 2008, pp. 21–30. <https://doi.org/10.1109/MSP.2007.914731>.
- 14 19. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical*
- 15 *Society. Series B (Methodological)*, Vol. 58, No. 1, 1996, pp. 267–288.
- 16 20. Kikuchi, S., and N. Kronprasert. Constructing Transit Origin-Destination Tables from Fragmented
- 17 Data. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. 2196,
- 18 2010, pp. 34–44. <https://doi.org/10.3141/2196-04>.
- 19 21. Sanandaji, B. M., and P. Varaiya. Compressive Origin-Destination Estimation. *Transportation*
- 20 *Letters*, Vol. 8, No. 3, 2016, pp. 148–157. <https://doi.org/10.1179/1942787515Y.0000000018>.
- 21 22. Mardani, M., and G. B. Giannakis. FEF. No. 00, pp. 1–5.
- 22 23. Chawla, S., Y. Zheng, and J. Hu. Inferring the Root Cause in Road Traffic Anomalies. 2012.
- 23 <https://doi.org/10.1109/ICDM.2012.104>.
- 24 24. Zhang, Y., Z. Ge, A. Greenberg, M. Roughan, and F. Park. Network Anomography.
- 25 25. Candès, E., and T. Tao. Decoding by Linear Programming Emmanuel Candès†. *IEEE Trans. Inf.*
- 26 *Theory*, Vol. 51, No. 12, 2005, pp. 4203–4215.
- 27 26. Park, T., and G. Casella. The Bayesian Lasso. Vol. 103, No. 482, 2008, pp. 681–686.
- 28 <https://doi.org/10.1198/016214508000000337>.
- 29 27. Michael Grant and Stephen Boyd. CVX: Matlab Software for Disciplined Convex Programming,
- 30 Version 2.1. <http://cvxr.com/cvx>.
- 31 28. Alireza Khani, Elizabeth Sall, Lisa Zorn, M. H. Integration of the FAST-TripS Person-Based
- 32 Dynamic Transit Assignment Model, the SF-CHAMP Regional, Activity-Based Travel Demand
- 33 Model, and San Francisco’s Citywide Dynamic Traffic Assignment Model. 2013.
- 34 29. Khani, A., M. Hickman, and H. Noh. Trip-Based Path Algorithms Using the Transit Network
- 35 Hierarchy. *Networks and Spatial Economics*, Vol. 15, No. 3, 2014, pp. 635–653.
- 36 <https://doi.org/10.1007/s11067-014-9249-3>.
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52

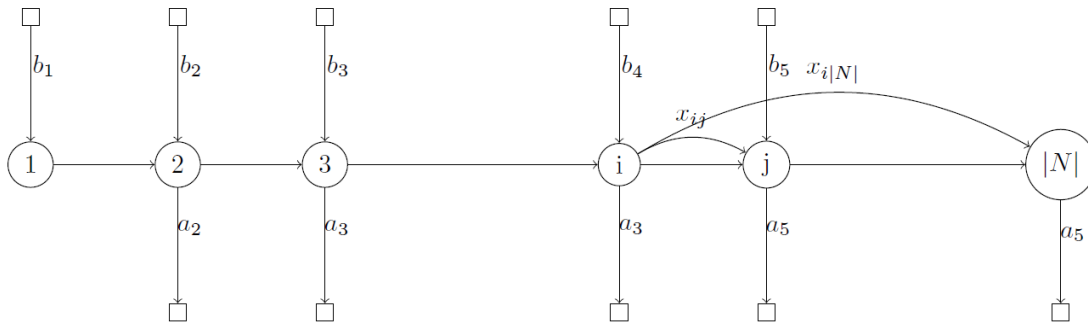
1  
2  
3  
4  
5  
6  
7  
8  
9

10 **TABLE 1. Notations**

<i>Variable</i>	<i>Definition</i>
N	Set of stops/stations along a transit route
i	Index for transit stop
$b_i$	Number of passengers boarding at stop i
$a_i$	Number of passengers alighting at stop i
X	Origin destination flow matrix
x	Vector form of OD matrix X
$\ x\ _0$	$l_0$ norm of vector x, $\ x\ _0 = \lim_{p \rightarrow 0} \sum_i  x_i ^p$
$\ x\ _1$	$l_1$ norm of vector x, $\ x\ _1 = \sum_i  x_i $
A	Linear map on a vector
Z	Set of integers

11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

1  
2  
3  
4  
5  
6  
7  
8  
9  
10



11  
12  
13  
14

**FIGURE 1. Transit route origin-destination (OD) flow**

Alighting

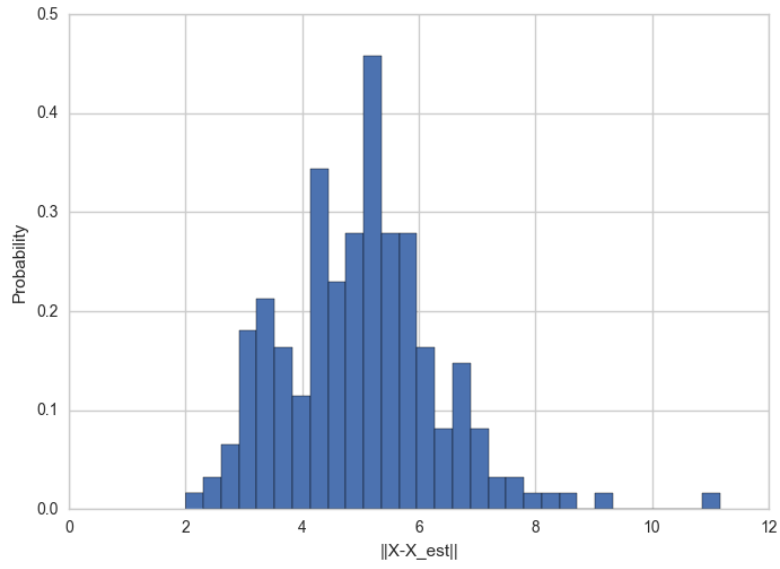
	1	2	.	.	.	n	
1	0						$b_1$
2	0	0					.
.	0	0	0				.
.	0	0	0	0			.
.	0	0	0	0	0		.
n	0	0	0	0	0	0	$b_n$
	$a_1$	.	.			$a_n$	$T$

Boarding

15  
16  
17  
18

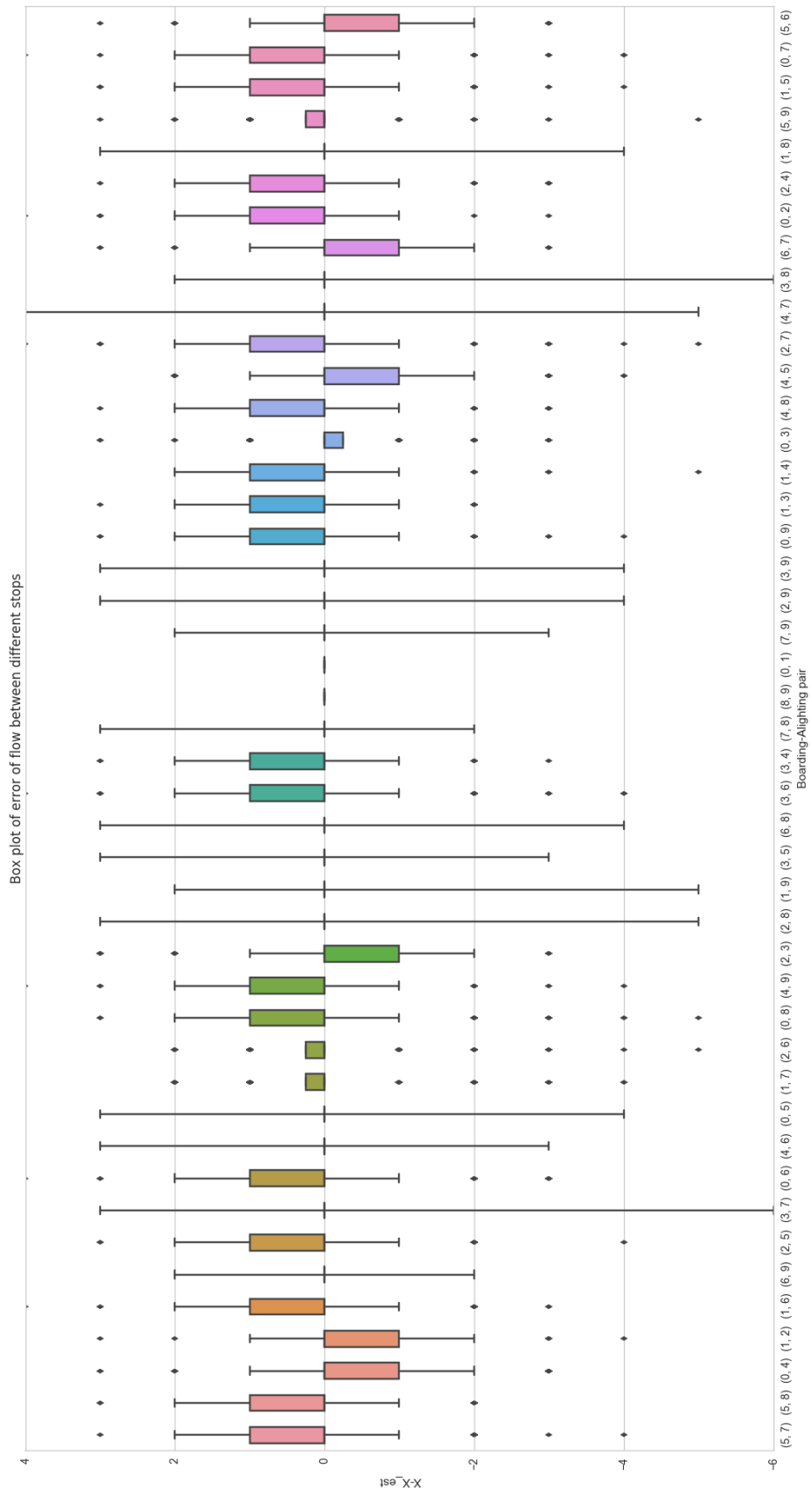
**FIGURE 2. OD matrix for a route in a single direction**





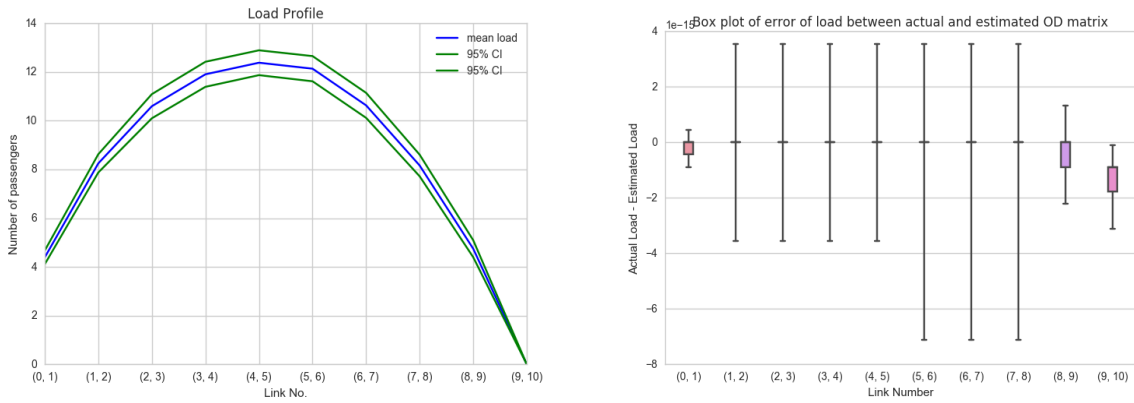
**FIGURE 3.  $l_2$  error between the actual and estimated OD matrix**

1  
2  
3

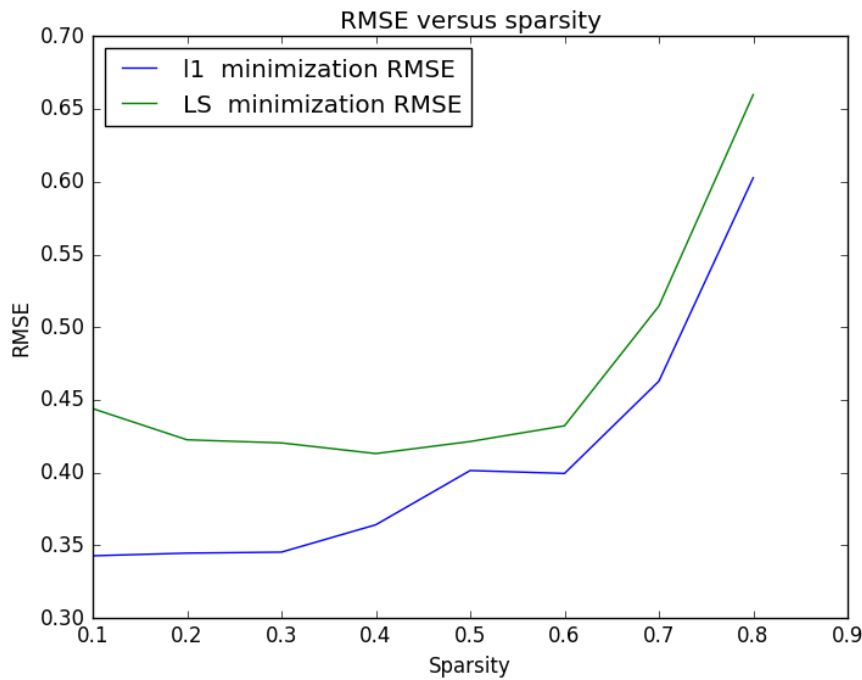


1  
2

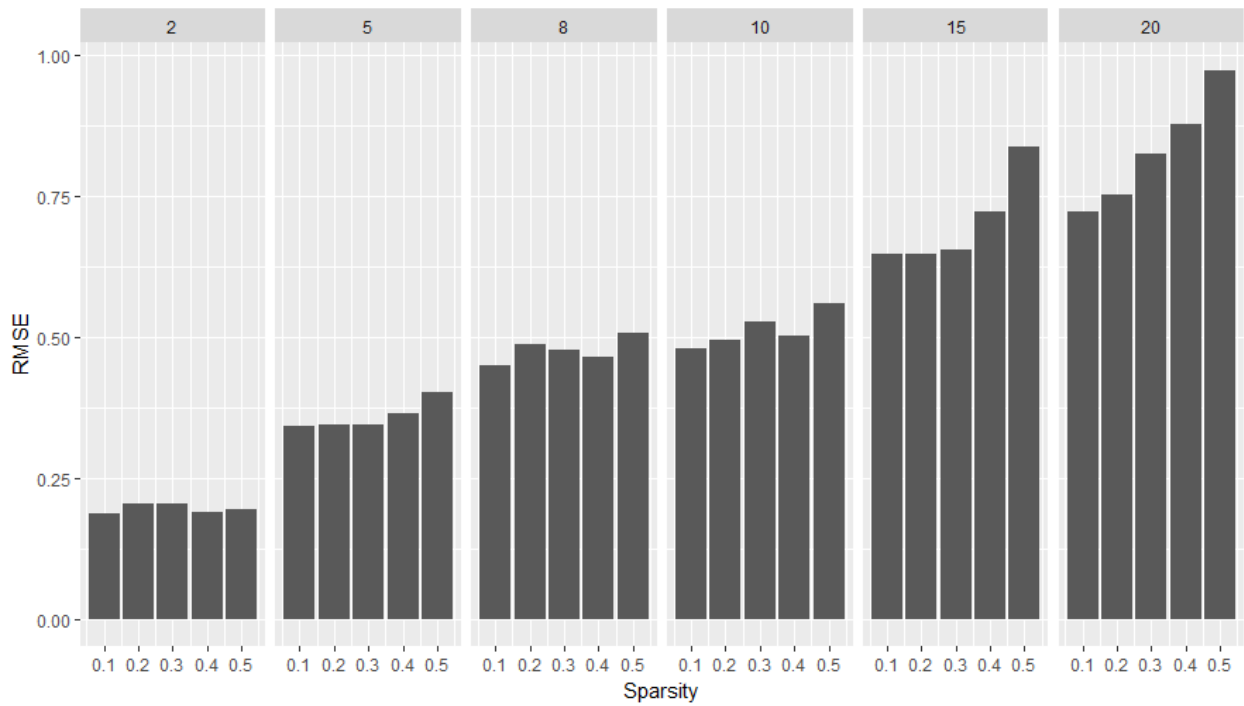
FIGURE 4. Box plot for the errors in estimation of O-D flows



1  
2 **FIGURE 5. (a) Average load profile of the transit route. (b) Box plot of error between actual**  
3 **load and estimated load**  
4

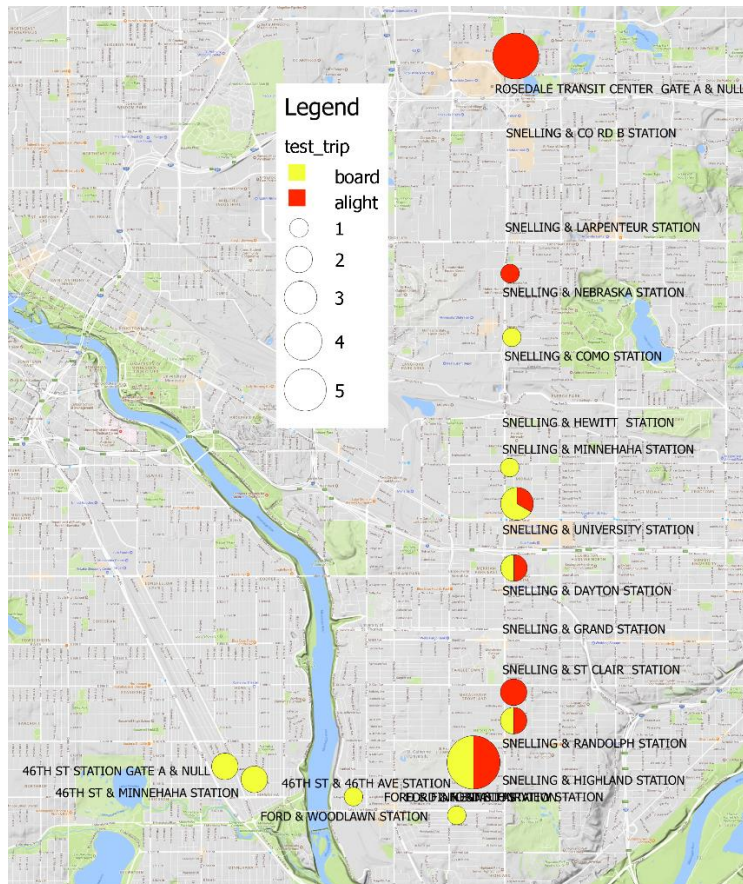


5  
6 **FIGURE 6. Root mean square error (RMSE) versus sparsity in OD estimation (Sparsity is**  
7 **in terms of proportion of non-zero values)**  
8



1  
2  
3  
4  
5

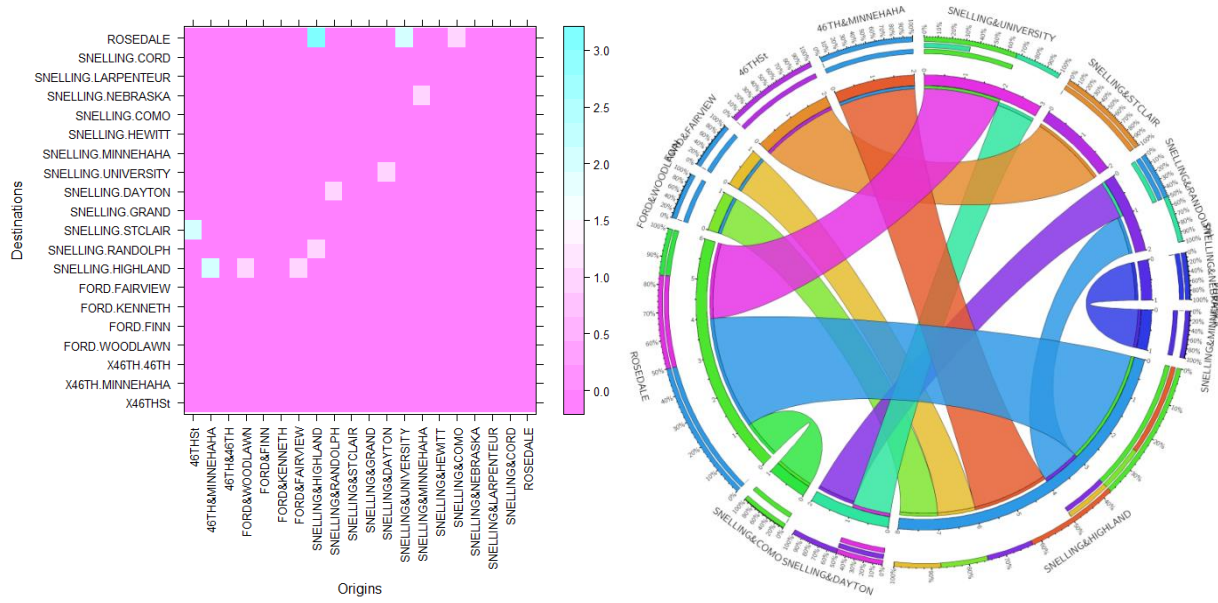
**FIGURE 7. Comparing RMSE with sparsity for different mean arrival rates (each panel represents different mean arrival rate of passengers)**



6

1  
2  
3

**FIGURE 8. Boarding and alighting count for A Line**



4  
5  
6

**FIGURE 9. Origin-Destination flow for A Line, Twin Cities, MN**