## Development of origin–destination matrices using mobile phone call data

1

2

# Development of Origin-Destination Matrices Using Mobile Phone Call Data: A Simulation Based Approach

5

6          Md. Shahadat Iqbal
7          Department of Civil Engineering.
8   Bangladesh University of Engineering and Technology, Dhaka 1000, Bangladesh
9          Shahadat.buet05@gmail.com
10
11
12          Charisma F. Choudhury*
13          Institute for Transport Studies
14          University of Leeds, Leeds LS2 9BJ, UK
15          cfc@alum.mit.edu

16
17          Pu Wang
18   School of Traffic and Transportation Engineering
19   Central South University, Hunan 410000, P.R. China
20
21
22          Marta C. Gonza´lez
23   Department of Civil and Environmental Engineering,
24   Massachusetts Institute of Technology, Cambridge, MA 02139, USA
25
26
27
28
29
30
31

| Word Count | Tables and Figures | 13 x 250 = 3250 |
|---|---|---|
| | Word Count | 3814 |
| | Total | 7064 |

35
36
37
38
39
40
41
42   *Corresponding Author

43  **Abstract**

44  In this research, we propose a methodology to develop OD matrices using mobile phone Call
45  Detail Records (CDR), which consist of time stamped tower locations with caller IDs, and
46  limited traffic counts. CDR from 2.87 million users from Dhaka, Bangladesh over a month and
47  traffic counts from 13 key locations of the city over 3 days of the same period are used in this
48  regard. The individual movement patterns within certain time windows are extracted first from
49  CDR to generate tower-to-tower *transient* OD matrices. These are then associated with
50  corresponding nodes of the traffic network and used as seed-OD matrices in a microscopic traffic
51  simulator. An optimization based approach, which aims to minimize the differences between
52  observed and simulated traffic counts at selected locations, is deployed to determine scaling
53  factors and the actual OD matrix is derived. The applicability of the methodology is supported by
54  a validation study.

55

## 1. Background

Reliable Origin-Destination (OD) matrices are critical inputs for analyzing transportation initiatives. Traditional approaches of developing OD matrices rely on roadside and household surveys, and/or traffic counts. The roadside and household surveys for origin destination involve expensive data collection and thereby have limited sample sizes and lower update frequencies. Moreover, they are prone to sampling biases and reporting errors (e.g.1,2,3). Estimation of reliable OD matrices from traffic link count data on the other hand is extremely challenging since very often the data is limited in extent and can lead to multiple plausible non-unique OD matrices (4,5). A number of Bayesian methods (e.g.6,7,8), Generalized Least Squares approaches (e.g.9,*10*), Maximum Likelihood Approaches (*11*), and Correlation Methods (e.g.*12,13,14*) have been used to tackle the indeterminacy problem. These approaches typically use *target* matrices based on prior information for generating the plausible route flows and are very sensitive to this prior information as well as to the chosen methodology (*15*). More recent approaches for OD estimation include automated registration plate scanners (*16*) and mobile traffic sensors such as portable GPS devices (e.g.*17,18,19*) . The practical successes of these approaches have however been limited due to high installation costs of the license plate readers and the low penetration rates of GPS devices (especially in developing countries).

Mobile phone users on the other hand also leave footprints of their approximate locations whenever they make a call or send an SMS. Over the last decade, mobile phone penetration rates have increased manifold both in developed and developing countries: the current penetration rates being 128% and 89% in developed and developing countries respectively (*20*). Subsequently, mobile phone data has emerged as a very promising source of data for transportation researchers. In recent years, mobile phone data have been used for human travel pattern visualization (e.g. *21,22,23*), mobility pattern extraction (e.g. *24,25,26,27,28,29*), route choice modeling (e.g. *30,31*), traffic model calibration (e.g. *32*), traffic flow estimation (*33*) to name a few. There have been several limited scale researches to explore the feasibility of application of mobile phone data for OD estimation as well. Wang et al. (*34*) for instance use a correlation based approach to dynamically update a prior OD matrix using time difference of phone signal receipt times of base stations and Caceras et al. (*35*) use a GSM network simulator to simulate the detailed movements of phones that are turned on. But both of these feasibility studies are based on synthetic data in small networks and the practical application is challenging given the need to collect and process detailed location data (which are currently processed by the mobile phone companies for load management purposes but are not stored). The potential estimate OD matrices using mobile phone Call Detail Records (CDR) (which are stored by operators for billing purposes and hence more readily available) have also been explored (e.g. *36,37,38*). Mellegård et al. (*36*) have developed an algorithm to assign mobile phone towers extracted from CDR to traffic nodes and Calabrese et al. (*37*) have proposed a methodology to reduce the noise in the CDR data but both studies have focused more on computation issues and the relationship between the mobile phone OD and the traffic OD have not been explored in

96   detail. Wang et al. (*38*) have used an analytical model to scale up the ODs derived from CDR by
97   using the population, mode choice probabilities and vehicle occupancy and usage ratios and have
98   validated it using probe vehicle data. The methodology however relies heavily on availability of
99   traffic and demographic data in high spatial resolution which may not be always available,
100  particularly in developing countries.

101  In this research, we propose a methodology to develop OD matrices using mobile phone CDR
102  and limited traffic counts. CDR from 2.87 million users from Dhaka, Bangladesh over a month
103  are used to generate the OD patterns on different time periods and traffic counts from 13 key
104  locations of the city over a limited time are used to scale it up to derive the actual ODs using a
105  microscopic traffic simulator. The methodology is particularly useful in situations when there is
106  limited availability of high resolution traffic and demographic data. The ODs are validated by
107  comparing the simulated and observed traffic counts of a different location (which has not been
108  used for calibration).

109  The rest of the paper is organized as follows. First we describe the data followed by the
110  methodology used for development of the OD matrix. The estimation and validation results are
111  presented next. We conclude with the summary of findings and directions for future research.

112  **2. Data**

113  *2.1 Study Area*

114  The central part of the Dhaka city has been selected as the study area and the major roads in the
115  network has been coded.  This consists of 67 nodes and 215 links covering an area of about
116  300km$^2$ with a population of about 10.7million (*39*). The average trip production rate is 2.74 per
117  person per day with significant portions of walking (19.8%) and non-motorized transport trips
118  (38.3%) (*39*).The traffic is subjected to severe congestion in most parts of the day, the average
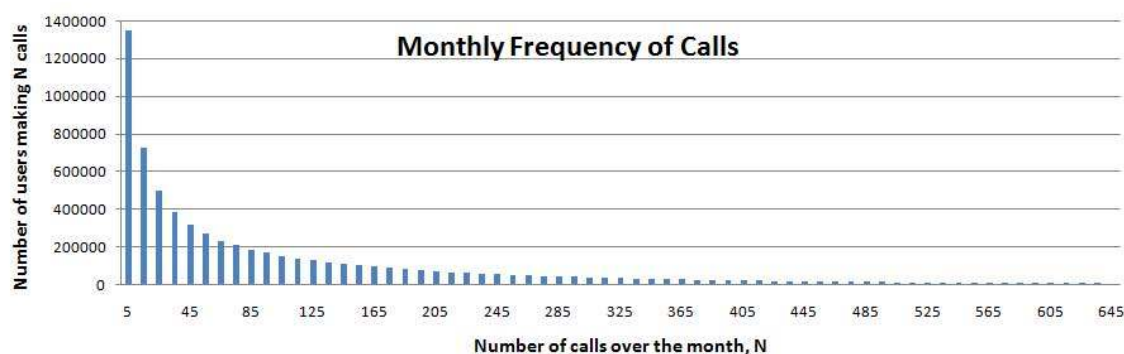119  speed being only 17km/hr[1].

120  The mobile phone penetration rate is approximated to be more than 90% in Dhaka (66.36%
121  being the national average) and Grameenphone Ltd. has the highest market share with 42.7m
122  mobile phone subscribers nationwide (*40*).

123  *2.2 CDR Data*

124  The CDR data, collected from Grameenphone Ltd, consists of calls from 6.9 million users
125  (which are more than 65% of the population of the study area) over a month. This comprises of
126  971.33 million anonymized call records in total made in between June 19, 2012 and July 18,
127  2012. The majority of the users (63%) have made 100 calls or less over the month. The
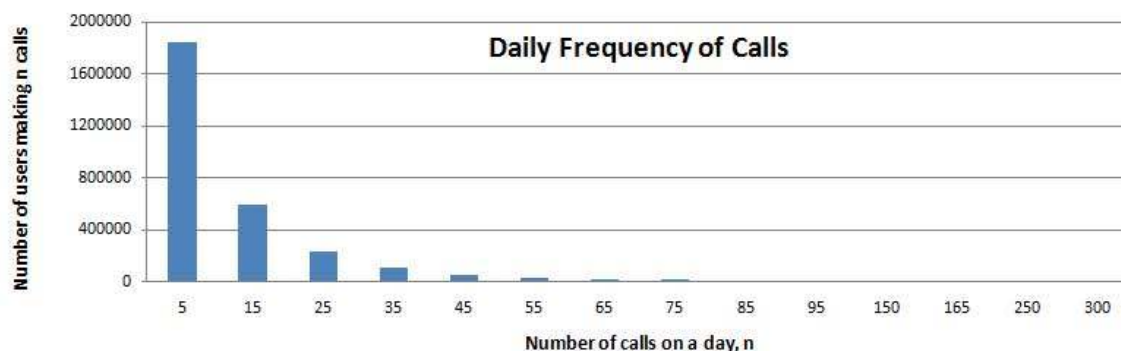128  frequencies of users making certain number of calls over the month and on a randomly selected

---

[1] Excluding the non-motorized vehicles which are restricted from entering the major roads

129  day (15[th] July, 2012) are presented in Figure 1. It may be noted that no demographic data related
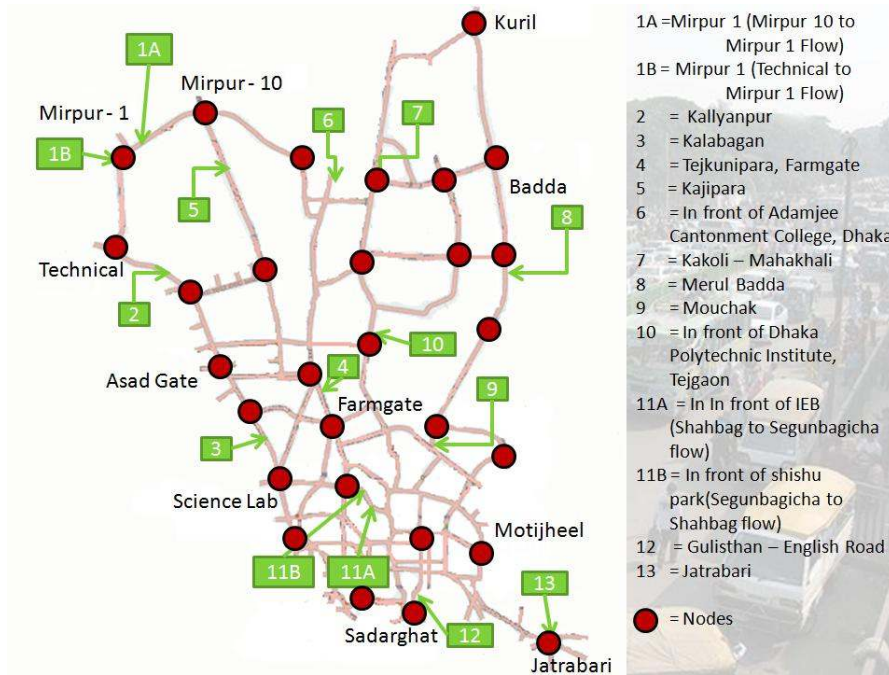130  to the phone users are available.

131



132



133
134  **Figure 1:**  Frequency of calls per user

135  *2.3 Traffic Count Data*

136  Video data, collected from 13 key locations of Dhaka city network over 3 days (12[th], 15[th] and
137  17[th] July 2012) have been used in this study to extract the traffic counts[2]. The locations (shown
138  in Figure 2) have been selected such that they cover the major roads (links) of Dhaka city with
139  flows from major generators and governed by the availability of foot over bridges for mounting
140  video cameras. Since MITSIMLab is developed for lane-based motorized traffic, care has been
141  taken to avoid roads that have high percentages of non-motorized transport and where lane-
142  discipline is not strictly followed. The data has been collected for 8 hrs (8.00 am to 12.00 noon
143  and 3.00 pm to 7.00pm) and analyzed using the software TRAZER (*41*) to generate classified
144  vehicle counts. Due to inclement weather and poor visibility some portion of the data is non-
145  usable though. Moreover, TRAZER (which is the only commercial software that can deal with
146  mixed traffic streams with *'weak'* lane discipline) has high misspecification rates in presence of
147  high congestion levels and in those cases, manual counting has been performed instead.
148

---

[2] There are no loop detectors or any other automatic traffic counters in Dhaka
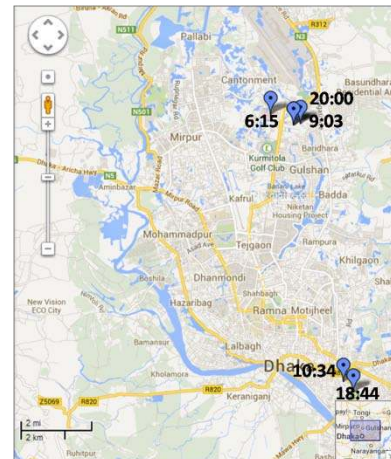
149



**Figure 2:** Locations of video data collection and position of OD generating nodes

## 3. Methodology
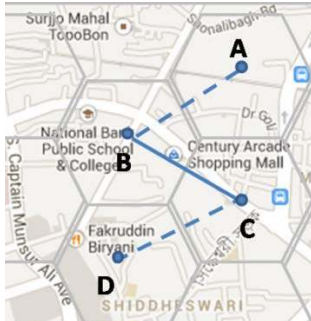
Each entry in the CDR contains unique caller id (anonymized), the date and time of the call, call duration and latitude and longitude of the Base Transceiver Station (BTS). A snapshot of the data is presented in Figure 1. As seen in the figure, if a person traverses within the city boundary and uses his/her phone from different locations that is captured in the CDR. CDR can thus provide an abstraction of his/her physical displacements over time (Figure 3).

| ID | Call Date | Call Time | Duration | Latitude | Longitude |
|---|---|---|---|---|---|
| AH03JAC8AAAbXtAId | 20120701 | 09:34:19 | 18 | 23.8153 | 90.4181 |
| AAH03JABiAAJKnPAa5 | 20120707 | 06:15:20 | 109 | 23.8139 | 90.3986 |
| AAH03JABiAAJKnPAa5 | 20120707 | 09:03:06 | 109 | 23.7042 | 90.4297 |
| AAH03JABiAAJKnPAa5 | 20120707 | 10:34:19 | 16 | 23.6989 | 90.4353 |
| AAH03JABiAAJKnPAa5 | 20120707 | 18:44:53 | 154 | 23.6989 | 90.4353 |
| AAH03JABiAAJKnPAa5 | 20120707 | 20:00:08 | 154 | 23.8092 | 90.4089 |
| AAH03JAC5AAAdAYAE | 20120701 | 09:15:05 | 62 | 23.7428 | 90.4164 |
| AAH03JAC+AAAcVKAC | 20120707 | 08:56:34 | 242 | 23.7908 | 90.3753 |
| AAH03JAC+AAAcVKAC | 20120701 | 18:03:06 | 36 | 23.9300 | 90.2794 |
| AAH03JAC5AAAdAYAA | 20120701 | 11:15:55 | 12 | 23.7428 | 90.4164 |

158

**Figure 3:** An excerpt from CDR data (entries of the same user are highlighted) and locations of a random user "AAH03JABiAAJKnPAa5" throughout the day as observed in data

161  However, in the CDR data, a user's location information is lost when he/she does not use his/her
162  phone. As shown in Figure 4, according to the CDR, a user may be observed to move from zone
163  B to zone C, but his/her initial origin (O) and final destination (D) may actually be located in
164  zone A and zone D. In such cases, a segment of the trip information is unobserved in the CDR.
165  However, the mobile phone call records enable us to capture the *transient* origins and
166  destinations which still retain a large portion of the actual ODs. Thus, we use the concept of
167  transient origin destination (*t*-OD) matrix (as used by Wang et al. (*38*)), which uses the mobile
168  phone data to efficiently and economically capture the pattern of travel demand.
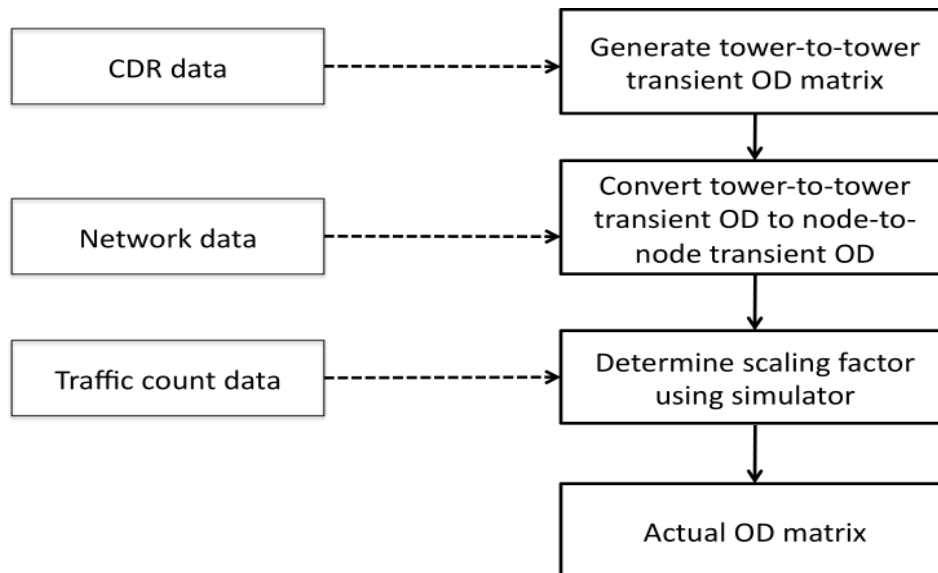
169


170  **Figure 4:**  Actual vs. Transient OD

171  The second source of data used in this research is classified traffic counts extracted from video
172  recordings collected from 13 key locations of Dhaka. These counts represent the *ground truth*
173  but are more expensive to collect[3] and limited in extent (only 3 days). This limited point source
174  data therefore cannot be used as a stand-alone source to reliably capture the OD pattern.

175  In this research, we therefore plan to combine the two data sources. The OD pattern is generated
176  using the CDR data and scaled up to match the traffic counts. The scaling factors are determined
177  using a microscopic traffic simulator platform MITSIMLab (*42*) using an optimization based
178  approach which aims to minimize the differences between observed and simulated traffic counts
179  at the points where the traffic counts are available.

180  The methodology is summarized in Figure 5 and described in the subsequent sections.

---

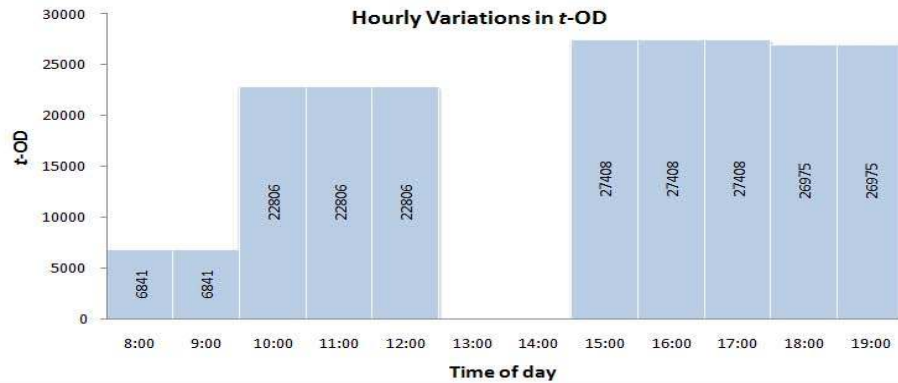[3] There are no detectors or any other traffic count mechanisms in Dhaka

**Figure 5:** Framework for developing OD Matrix

*3.1 Generation of tower-to-tower transient OD matrix*
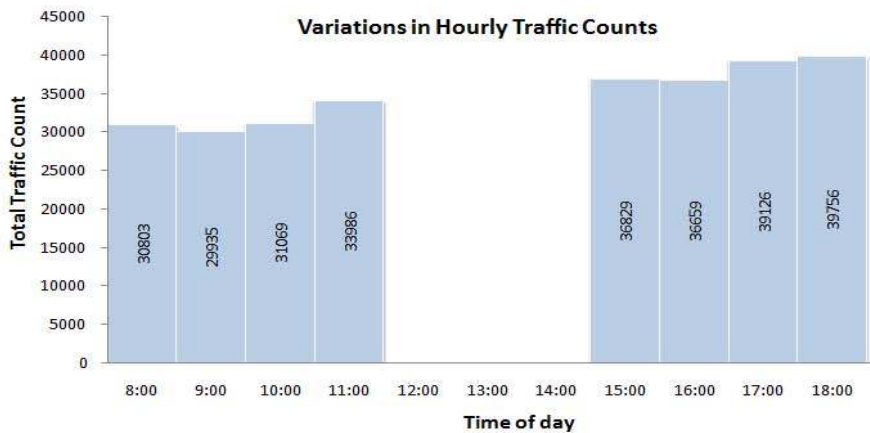
The time-stamped BTS tower locations of each user are first extracted from the mobile phone CDR data and used for generating tower-to-tower transient OD matrix. The CDR however only contains sparse and irregular records (*28*), in which user displacements (consecutive non-identical locations) are usually observed with long travel intervals i.e. the first location may be observed at 8:56 and next location may be observed at 18:03 with no information about intermediate locations (if any) or the time when the trip in between these two locations have been made.

Another limitation is the CDR data often records changes in towers in spite of no actual displacement (as the operator balances call traffic among adjacent towers). To better identify timing and origin-destinations of specific trips and reduce the number of *false displacements*, we therefore extract displacements that have occurred within a specific *time window*. A lower bound in the time window (10 minutes) is imposed to reduce the number of *false displacements* without affecting the number of physical displacements occurring within short intervals. An upper bound in the time window (1 hr) is imposed to ensure that meaningful numbers of trips are retained. Therefore, a person trip is recorded if in the CDR, subsequent entries of the same user indicate a displacement (change in tower) with a time difference of more than 10 minutes but less than 1 hour.

Further, both call volumes (from CDR data) and traffic volumes (from traffic counts) had significant variations throughout the day. Based on correlation analysis of total mobile call volumes and total traffic counts (Figure 6), four time periods (7:00-9:00, 9:00-12:00, 15:00-17:00 and 17:00-19:00), have been chosen for analysis.

205



206

**Figure 6:** Hourly variations a. traffic count b. transient ODs from mobile call records

*3.2 Conversion of tower-to-tower t-OD to node-to-node t-OD*

For application of the *t*-ODs in traffic analyses, the origin and destination towers need to be associated with corresponding nodes of the traffic network. The typical tower coverage area can be represented as a combination of three hyperbolas (Figure 7), the size varying depending on tower height, terrain, locations of adjacent towers and number of users active in the proximity (which can vary dynamically).



214

**Figure 7:** Typical coverage area of a tower (http://www.truteq.co.za/tips_gsm/)

216   The population density in the chosen study area is very high (more than 8111 inhabitants/sq. km
217   (*44*) and the tower locations are very close to each other (1 km on average). Because of the high
218   user density, it can be assumed that the area between two towers is equally split among the two
219   towers (Figure 8) that is, each tower $t$ has a coverage area ($A_t$) approximately defined by a circle
220   of radius $0.5l$, where $l$ is the tower-to-tower distance.

221



222
223   Tower 6 and Node 3 need to be added to Figure

224

**a. Tower-to-tower OD**

| ID | Call Date | Call Time | Origin Tower | Destination Tower |
|---|---|---|---|---|
| AAH03JA | 20120718 | 15:54 | 6 | 1 |
| AAH03JA | 20120718 | 16:13 | 1 | 2 |
| AAH03JA | 20120718 | 16:15 | 2 | 1 |
| AAH03JA | 20120718 | 18:53 | 1 | 6 |
| AAH03JA | 20120718 | 20:49 | 6 | 1 |
| AAH03JA | 20120718 | 23:41 | 1 | 6 |

**b. Intermediate OD with candidate nodes**

| ID | Call Time | Origin Tower | Origin Candidate Node | Destination Tower | Destination Candidate Node |
|---|---|---|---|---|---|
| AAH03JA | 14:54 | 6 | 3 | 1 | 1 |
| AAH03JA | 16:13 | 1 | 1 | 2 | 2 Or 1 |
| AAH03JA | 16:15 | 2 | 2 Or 1 | 1 | 1 |
| AAH03JA | 18:53 | 1 | 1 | 6 | 3 |
| AAH03JA | 20:49 | 6 | 3 | 1 | 1 |
| AAH03JA | 23:41 | 1 | 1 | 6 | 3 |

**c. Node-to-node OD**

| ID | Call Time | Origin Node | Destination Node |
|---|---|---|---|
| AAH03JA | 14:54 | 3 | 1 |
| AAH03JA | 16:13 | 1 | 1 |
| AAH03JA | 16:15 | 1 | 1 |
| AAH03JA | 18:53 | 1 | 3 |
| AAH03JA | 20:49 | 3 | 1 |
| AAH03JA | 23:41 | 1 | 3 |

225
226   a. Tower-to-tower OD          b. Intermediate OD with candidate nodes  c. Node-to-node OD
227   **Figure 8:**  Example of tower to node allocation

228   If a unique traffic node $i$ overlaps with $A_t$, the calls handled by $t$ are associated with node $i$ (as in
229   the case of Tower 1 in Figure 6). However, if $A_t$ has two (or more) candidate nodes for
230   association, then the candidate nodes are ranked based on the proportion of $A_t$ feeding to each
231   node. That is, the node serving greatest portion of $A_t$ is ranked 1, the node serving second highest
232   portion of $A_t$ is ranked 2, etc. For example, in Figure 6, network connectivity (feeder roads) and
233   topography (presence of a canal with no crossing facility in the vicinity) denote that Node 1 and
234   Node 2 are candidate nodes for association with Tower 2. As the major portion of $A_t$ is connected
235   to Node 2 and the remaining portion is connected to Node 1, they are ranked 1 and 2 respectively
236   for Tower 2. The data format after this step is presented in Figure 7b. As seen in the figure, this
237   typically consists of call records associated with unique nodes and some calls associated with
238   *multiple candidate nodes*. The calls are then sorted and ranked based on the frequency of the
239   unique nodes used by each user.  The frequency of occurrence of the candidate nodes are

240    compared and used as the basis of replacement. For example, frequency analysis of User
241    "AAH03JA" indicates a higher frequency of Node 1. Therefore, in cases where there are
242    ambiguities between Nodes 2 and 1, Node 1 is used (for this particular user).
243    The same process is used for all users and node-to-node $t$-OD matrices for each time period of
244    each day are derived.
245
246    *3.3 Finding the scaling factor and determining the actual OD matrix*
247    As discussed, the node-to-node $t$-OD matrix ($t\text{-}OD_{ij}$) provides the trip patterns for developing the
248    actual OD matrix ($OD_{ij}$). However, in order to determine the actual OD matrix, the $t$-OD needs to
249    be scaled to match the real traffic flows. A scaling factor $\beta_{ij}$ is used in this regard:

$$OD_{ij} = \sum_{ij} (t\text{-}OD_{ij}) * \beta_{ij}$$

250    It may be noted that $\beta_{ij}$ takes into account the market penetration rates (i.e. not every user has a
251    mobile phone or uses the specific service provider), the mobile phone non-usage issue (i.e.
252    mobile phone calls are not made from every location traversed by the user), the vehicle usage
253    issue (i.e. users may not use cars for every trip). The potential error introduced due to *false*
254    *displacement* (described in Section 2.1) is also accounted for in the scaling factors.
255    The scaling factors are determined using the open-sourced microscopic traffic simulator platform
256    MITSIMLab (*42*) by applying an optimization based approach. The movements of vehicles in
257    MITSIMLab are dictated by driving behavior models based on decision theories and estimated
258    with detailed trajectory data using econometric approaches. Route choices of drivers are based
259    on a discrete choice based probabilistic model where the utilities of selecting and re-evaluating
260    routes are functions of path attributes, such as path travel times and freeway bias (see *43* for
261    details). The inputs of the simulator include network data, driving behavior parameters and OD
262    matrix. The generated outputs include traffic flow at specified locations in the network.

263    The node-to-node OD matrix derived from the mobile phone data are provided as the initial or
264    seed-OD in this case. The simulated traffic flows are compared with the actual traffic flows
265    extracted from video recordings. The objective function seeks to minimize the difference
266    between the actual and simulated traffic flows in each location by changing the scaling factors.
267    The optimization problem can be represented as follows:
268
269    $minimize, Z = \sum_{k=1}^{K}(V_{actual}^{k} - V_{simulated}^{k})^2$                  (1)
270               $Such\ that, OD_{ij,t} = \sum_{i,j=1}^{N} t\text{-}OD_{ij,t} * \beta_{ij,t}$
271    Where,
272    $V_{simulated}^{k}$= Traffic flow of link $k$ of the road network from simulation
273    $OD_{ij,t}$    = Actual OD between nodes $i$ and $j$ in time period $t$
274    $t\text{-}OD_{ij,t}$   = Transient OD between nodes $i$ and $j$ in time period $t$
275    $\beta_{ij,t}$         = Scaling factor associated with the node pair $i$ and $j$ and time period $t$

276     $K$         = Total number of links for which traffic flow data is available

277     $N$         = Total number of nodes in the network

278

279     However, to make the optimization problem more tractable, group-wise scaling factors are used

280     rather than an individual scaling factor for each OD pair. The grouping is based on the analyses

281     of the CDR data. This simplifies the problem as follows:

282

283     $minimize, Z = \sum_{k=1}^{K}(V_{actual}^{k} - V_{simulated}^{k})^{2}$                   (2)

284             Such that, $OD_{ij,t} = \sum_{m=1}^{M} t\text{-}OD_{ij,t}^{m} * \beta_{t}^{m}$

285     Where,

286     $t\text{-}OD_{ij,t}^{m}$    = Transient OD between node pair $i$ and $j$ in time period $t$ where the node pair $i,j$

287             belong to group $m$

288     $\beta_{t}^{m}$         = Scaling factor for group $m$ and time period $t$

289     M          = Total number of groups of OD-pairs

290

291     **4. Results**

292     The mobile phone network within the study area comprises of 1360 towers which have been

293     assigned to 29 OD generating nodes (812 OD pairs). Out of the one month CDR data, the

294     weekend data have been discarded. For each day, the calls of each user originating from two

295     different towers in each of the time period have been extracted. After application of the transient

296     trip definitions (displacements occurring more than 10mins but less than 1hr apart) and the tower

297     to node conversion rules (elaborated in Section 3.2), the node-to-node $t$-ODs are derived. The

298     total number of node-to-node $t$-ODs are presented in Table 1.

299         **Table 1:** Node-to-node $t$-OD

300

| Time Period | Time | t-OD | |
|---|---|---|---|
| | | Total Over the Month[4] | Weekday Average |
| 1 | 7:00-9:00 | 397355 | 13681.86 |
| 2 | 9:00-12:00 | 1915417 | 68418.48 |
| 3 | 15:00-17:00 | 2255859 | 82226.05 |
| 4 | 17:00-19:00 | 1549109 | 53950.57 |

301

302

303

---

[4] Includes weekends

304  Analyses of the node-to-node transient flows indicate that the flows between adjacent nodes are
305  substantially higher than those between non-adjacent nodes (Figure 9).  This is reasonable since
306  given the low travel speed in Dhaka, a traveler may not be able to move very far in the 50min
307  time window and the *t*-ODs mostly capture segments of a longer trip. However, part of it may
308  also be due to the *false displacement* problem discussed in section 3.1. Therefore, the OD-pairs
309  have been divided into two groups (adjacent and non-adjacent nodes) and the objective function
310  to determine scaling factors has been formulated as follows:

311  $minimize, Z = \sum_{k=1}^{K}(V_{actual}^{k} - V_{simulated}^{k})^2$ (3)

312  Such that, $OD_{ij,t} = \sum_{adj} t\text{-}OD_{ij,t}^{adj} * \beta_t^{adj} + \sum_{non\text{-}adj} t\text{-}OD_{ij,t}^{non\text{-}adj} * \beta_t^{non\text{-}adj}$
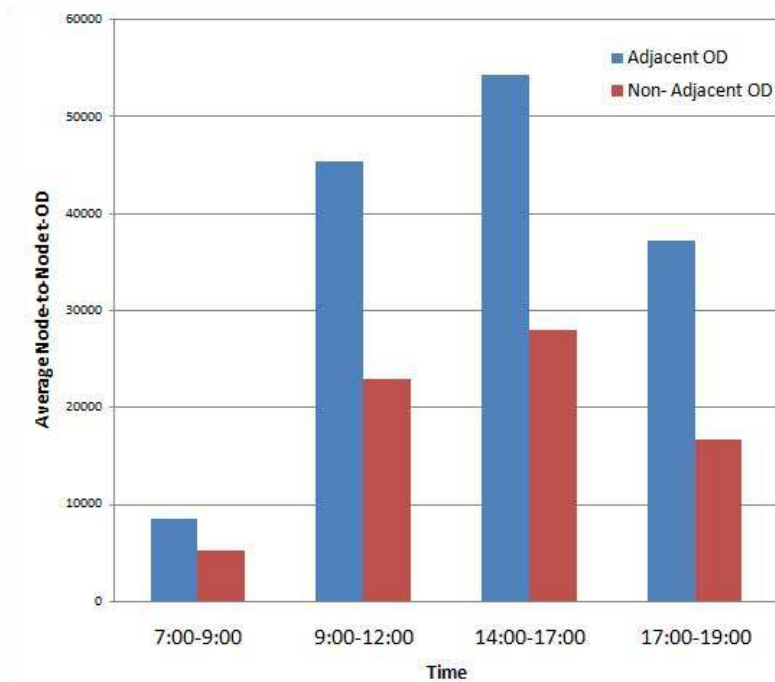
313  Where,

314  $t\text{-}OD_{ij}^{adj}$   = Transient OD between node pair *i* and *j* in time period *t* where the node pair *i,j*
315       are adjacent nodes

316  $t\text{-}OD_{ij}^{non\text{-}adj}$   = Transient OD between node pair *i* and *j* in time period *t* where the node pair *i,j*
317       are non-adjacent nodes

318  $\beta_t^{adj}, \beta_t^{non\text{-}adj}$ = Scaling factors for time period *t* and adjacent and non-adjacent nodes
319       respectively

320

321



322  **Figure 9:** Comparison of  *t*-ODs between adjacent and non-adjacent nodes
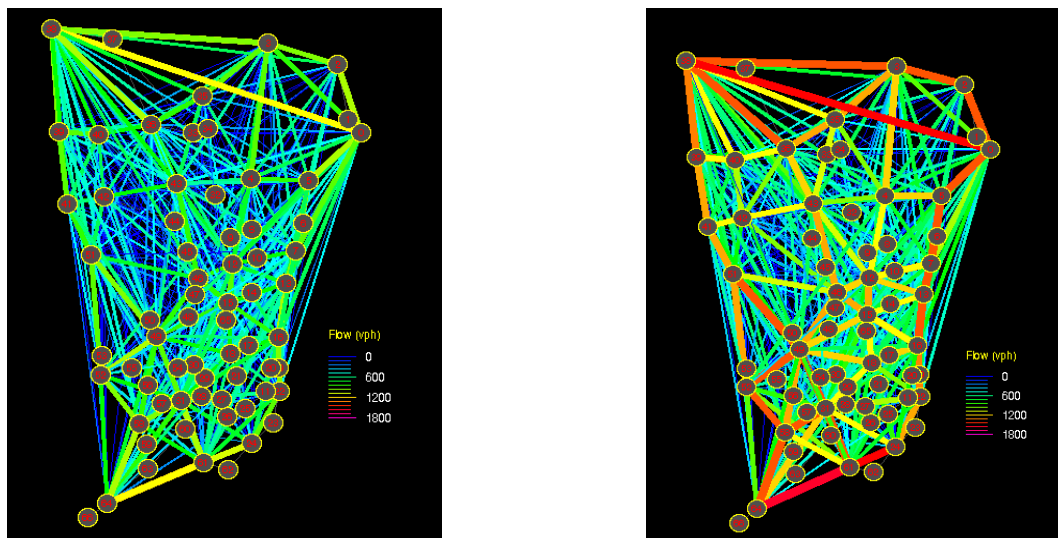
323     This yielded eight scaling factors in total that needed to be estimated from the simulation runs of
324     MITSIMLab. Running the optimization process in MATLAB (that invokes MITSIMLab) and
325     using a BOX algorithm (*45*), the following values of scaling factors have been derived.

326     **Table 2:** Scaling Factors

| Time Period | OD Type | Scaling Factor |
|---|---|---|
| 7:00-9:00 | Adjacent | 6.787 |
| | Non-adjacent | 1.712 |
| 9:00-12:00 | Adjacent | 0.971 |
| | Non-adjacent | 0.345 |
| 15:00-17:00 | Adjacent | 1.647 |
| | Non-adjacent | 3.407 |
| 17:00-19:00 | Adjacent | 9.404 |
| | Non-adjacent | 6.779 |

328

329     It is interesting to note that the scaling factors for adjacent nodes are higher than those of non-
330     adjacent in all time periods other than 15:00-17:00. This does not however indicate that most of
331     the actual trips are to the adjacent nodes (since a full trip may consist of several segments each
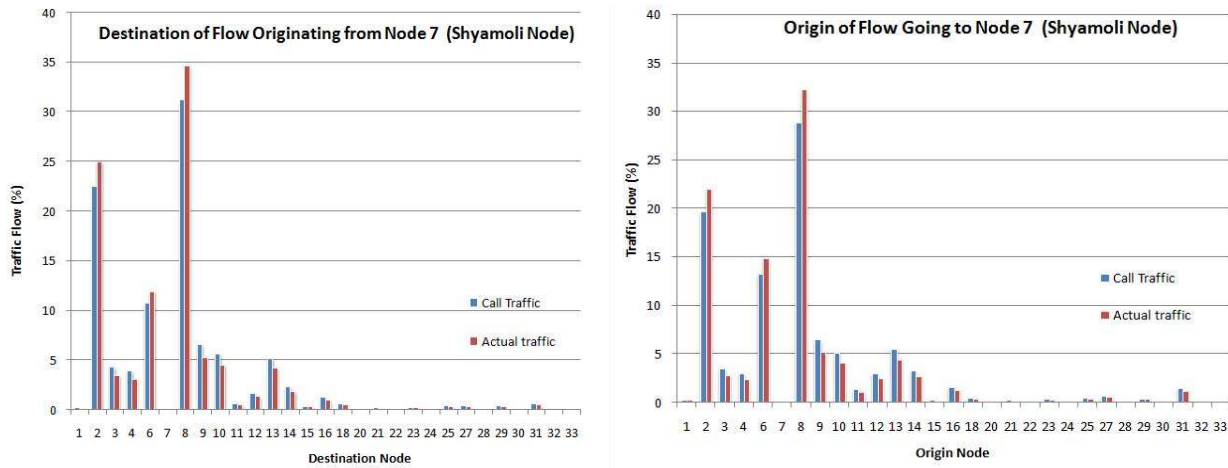332     represented by a separate *t*-OD).

333     The graphical representation of the *t*-ODs and actual ODs across the network for one of the time
334     periods and the variations for an example node are presented in Figures 10 and 11 respectively.



335     a. t-OD                                  b. actual OD

336     **Figure 10:** *t*-ODs and actual ODs across the network for 7:00-9:00

337



**Figure 11:** Example of Transient and Actual Traffic Flows To and From a Node (Shyamoli) between 7:00-9:00.

**5. Validation**

In addition to the aggregate data used for calibration, traffic counts are collected from four additional locations on a different day. For validation purposes, the scaled up ODs have been applied to simulate the traffic between 9:00-12:00 in MITSIMLab and the simulated traffic counts are compared against the observed counts from these locations. In order to quantify the prediction error, Root Mean Square Error and Root Mean Square Percent Errors have been calculated and are found to be 335.09 and 13.59% respectively.

**6. Conclusion**

The main outcome of this research is the methodology for development of the OD matrix using mobile phone CDR and limited traffic count data. The strengths of both data sources are utilized in this approach: the trip patterns are extracted from mobile phones and the ground truth traffic scenario are derived from the counts. The methodology is demonstrated using data collected from Dhaka.

There are several limitations of the current research though. Firstly, in this research a simplified objective function with grouped scaling factors has been used. This overlooks the heterogeneity in call rates from different locations (e.g., more calls may be generated to and from railway stations compared to and from offices with land telephone lines, etc.). A more detailed classification of scaling factor can be used to overcome this bias and may yield better results. Moreover, in this particular context, detailed network data and extensive calibration data were not available which may have increased the simulation errors and affected the validation results. However, initial validation results indicate promising success in real life application by transport planners and managers.

Since CDR is already recorded by mobile phone companies for billing purposes, the approach is more economic than the traditional approaches which rely on expensive household surveys and/or extensive traffic counts. It is also convenient for periodic update of the OD matrix and extendable for dynamic OD estimation. This method is particularly effective for generating complex OD matrix where land use pattern is heterogeneous and asymmetry in travelling pattern prevails throughout the day but there is a limitation of traditional data sources.

**Acknowledgment**

**References**

1. Hajek, J. J. (1977). *Optimal sample size of roadside-interview origin-destination surveys* (No. RR 208).
2. Kuwahara, M., and Sullivan, E. C. (1987). Estimating origin-destination matrices from roadside survey data. *Transportation Research Part B*,*21*(3), 233-248.
3. Groves, R. M. (2006). Nonresponse rates and nonresponse bias in household surveys. Public Opinion Quarterly, 70(5), 646-675.
4. Lo, H. P., Zhang, N., and Lam, W. H. (1996). Estimation of an origin-destination matrix with random link choice proportions: a statistical approach.*Transportation Research Part B*, *30*(4), 309-324.
5. Van Zuylen, H. J., and Willumsen, L. G. (1980). The most likely trip matrix estimated from traffic counts. *Transportation Research Part B*,*14*(3), 281-293.
6. Maher, M. (1983). Inferences on trip matrices from observations on link volumes: a Bayesian statistical approach. *Transportation Research Part B, 20 (6)*, 435–447.
7. Tebaldi, C., West, M. (1998). Bayesian inference on network traffic using link count data (with discussion). *Journal of the American Statistical Association,93,* 557–576.
8. Li, B. (2005). Bayesian inference for origin–destination matrices of transport networks using the EM algorithm. *Technometrics 47 (4)*, 399–408.
9. Cascetta, E. (1984). Estimation of trip matrices from traffic counts and survey data: a generalized least squares estimator. *Transportation Research Part B, 18(4–5)*, 289–299.
10. Bell, M. (1991). The estimation of origin–destination matrices by constrained generalized least squares. *Transportation Research Part B, 25 (1)*, 13–22.
11. Spiess, H. (1987). A maximum likelihood model for estimating origin-destination matrices, *Transportation Research Part B,21(5)*, 395-412.
12. Vardi, Y. (1996). Network tomography: estimating source-destination traffic intensities from link data. *Journal of the American Statistical Association,91,* 365–377.
13. Hazelton, M.L. (2000). Estimation of Origin–Destination matrices from link flows on uncongested networks. *Transportation Research Part B, 34 (7)*, 549–566.

14. Hazelton, M.L. (2003). Some comments on origin–destination matrix estimation. *Transportation Research Part A, 37 (10)*, 811–822.

15. Hazelton, M.L., 2001b. Inference for origin–destination matrices: estimation, reconstruction and prediction. *Transportation Research Part B, 35 (7),* 667–676.

16. Castillo, E., Menéndez, J., Jiménez, P. (2008). Trip matrix and path flow reconstruction and estimation based on plate scanning and link observations. *Transportation Research Part B ,42 (5),* 455–481.

17. Parry, K., & Hazelton, M. L. (2012). Estimation of origin–destination matrices from link counts and sporadic routing data. *Transportation Research Part B, 46(1),* 175-188.

18. Morimura, T., and Kato, S. (2012). Statistical origin-destination generation with multiple sources. 21st International Conference on In Pattern Recognition (ICPR), November 11-15, 2012. Tsukuba, Japan.

19. Herrera, J., Work D. B., Herring R., Ban X., Jacobson Q., Bayen A. (2010). Evaluation of traffic data obtained via GPS-enabled mobile phones: The Mobile Century field experiment, *Transportation Research Part C: Emerging Technologies, 18(4)*, 568-583.

20. http://www.itu.int/en/ITU-D/Statistics/Documents/facts/ICTFactsFigures2013.pdf [accessed 20July, 2013]

21. Phithakkitnukoon, S., Horanont, T., Di Lorenzo, G., Shibasaki, R., and Ratti, C. (2010). Activity - Aware Map: Identifying human daily activity pattern using mobile phone data, *Human Behavior Understanding, 6219(3),* 14-25,Springer Berlin / Heidelberg.

22. Phithakkitnukoon, S., and Ratti, C., (2011), Inferring Asymmetry of Inhabitant Flow using Call Detail Records, *Journal of Advances in Information Technology*, 2 (4), 239-249.

23. Reades, J., Calabrese, F., and Ratti, C. (2009). Eigenplaces: analyzing cities using the space-time structure of the mobile phone network, *Environment and Planning B: Planning and Design, 36(5),* pp. 824-836.

24. Wang, P., Hunter, T., Bayen, A. M., Schechtner, K., and González, M. C. (2012). Understanding Road Usage Patterns in Urban Areas. *Scientific reports, 2*.

25. G onzález, M. C., Hidalgo, C. A., and Barabási, A. L.(2008).Understanding individual human mobility patterns, *Nature, 453,* 779–782.

26. Song, C, Koren, T, Wang, P, and Barabási, A. L. (2010). Modelling the scaling properties of human mobility, *Nature Physics, 6*, 818–823.

27. Simini, F., Gonza´lez, M. C., Maritan, A., and Baraba´si, A. L.(2012). A universal model for mobility and migration patterns, *Nature, 484*, 96–100.

28. Candia, J., González, M. C., Wang, P., Schoenharl, T., Madey, G., and Barabási, A. L. (2008). Uncovering individual and collective human dynamics from mobile phone records. *Journal of Physics A: Mathematical and Theoretical, 41(22)*, 224015.

29. Sevtsuk, A., and Ratti, C. (2010). Does Urban Mobility Have a Daily Routine? Learning from Aggregate Data of Mobile Networks, *Journal of Urban Technology, 17 (1)*, 41-60.

30. Schlaich, J., Otterstätter, T., Friedrich, M., 2010, Generating Trajectories from Mobile Phone Data, TRB 89th Annual Meeting Compendium of Papers, Transportation Research Board of the National Academies, Washington, D.C., USA.

31. Becker, R.A., Caceres, R., Hanson, K., Loh, J.M., Urbanek, S., Varshavsky, A., Volinsky, C., Ave, P., Park, F., 2011. Route classification using cellular handoff patterns. In: Proceedings of the 13th International Conference on Ubiquitous Computing. ACM, Beijing, China.

447     32. Bolla, R., Davoli, F., and Giordano, A. (2000). Estimating road traffic parameters from mobile
448         communications. In *Proceedings 7th World Congress on ITS, Turin, Italy*.
449     33. Demissie, M. G., de Almeida Correia, G. H., and Bento, C. (2013). Intelligent road traffic status
450         detection system through cellular networks handover information: An exploratory study.
451         Transportation *Research Part C: Emerging Technologies, 32*, 76-88.
452     34. Wang J., Wang D. Song X. Sun Di. (2011). Dynamic OD Expansion Method Based on Mobile
453         Phone Location, Fourth International Conference on Intelligent Computation Technology and
454         Automation, Shenzhen, China.
455     35. Caceres, N., Wideberg, J. P., and Benitez, F. G. (2007). Deriving origin destination data from a
456         mobile phone network. *Intelligent Transport Systems, IET*, *1*(1), 15-26.
457     36. Mellegard, E., Moritz, S., and Zahoor, M. (2011, December). Origin/Destination-estimation using
458         cellular network data. In *Data Mining Workshops (ICDMW), 2011 IEEE 11th International
459         Conference on* (pp. 891-896). IEEE.
460     37. Calabrese F., Lorenzo G. D., Liu L. and Ratti C. (2011). Estimating Origin-Destination Flows
461         using Mobile phone Location Data. IEEE Pervasive Computing, vol. XX, no. XX, 200XX, pp.
462         36–43.
463     38. Wang, P., Hunter, T., Bayen, A. M., Schechtner, K., and González, M. C. (2012). Understanding
464         Road Usage Patterns in Urban Areas. *Scientific reports*, *2*.
465     39. DHUTS. (2010). Dhaka Urban Transport Network Development Study, Draft Final Report.
466         Prepared by Katahira and Engineers International, Oriental Consultants Co. Ltd., and Mitsubishi
467         Research Institute, Inc.
468     40. Grameenphone Ltd. Bangladesh. http://grameenphone.com, accessed on 15.12.2012
469     41. Kritikal Solutions Ltd., India. http://www.kritikalsolutions.com/products/traffic-analyzer.html,
470         accessed on 15.12.2012
471     42. Yang Q. and Koutsopoulos, H. N., (1996). A microscopic traffic simulator for evaluation of
472         dynamic traffic management systems, *Transportation Research C, 4(3)*,113-129
473     43. Ben-Akiva M., Koutsopoulos H. N., Toledo T., Yang Q., Choudhury C. F., Antoniou C., and
474         Balakrishna R. (2010). Traffic simulation with MITSIMLab, in Fundamentals of Traffic
475         Simulation, 1st ed., ser. International Series in Operations Research and Management Science, J.
476         Barceló, Ed. Springer, 233-268.
477     44. Population and Housing Census: Preliminary Results (2011), Bangladesh Bureau of Statistics,
478         Statistics Division, Ministry of Planning, Government of the People's Republic of Bangladesh
479     45. Box M. J. (1965), A new method of constrained optimization and a comparison with other
480         methods, *Computer Journal, 8(1)*,42-52.

481