Using Mobile Phones to Forecast Arterial Traffic through Statistical Learning

Ryan Herring^{*} Aude Hofleitner[†] Saurabh Amin[‡] Tania Abou Nasr[§]

Amin Abdel Khalek[¶] Pieter Abbeel^{\parallel} Alexandre Bayen^{**}

Submitted For Publication 89th Annual Meeting of the Transportation Research Board August 1, 2009

Word Count:

Number of words:	6245
Number of figures:	7 (250 words each)
Number of tables:	0 (250 words each)
Total:	7995

^{*}Corresponding Author, Department of Industrial Engineering and Operations Research, University of California, Berkeley, 2105 Bancroft Way, Suite 300, Berkeley, CA 94720, (510) 642-5667, ryanherring@berkeley.edu

[†]Department of Electrical Engineering and Computer Science, University of California, Berkeley, 2105 Bancroft Way, Suite 300, Berkeley, CA 94720, aude.hofleitner@polytechnique.edu, and Ecole Doctorale Ville, Transport et Territoire, Universit Paris-Est, Marne-La-Valle, France

[‡]Department of Civil and Environmental Engineering, Systems Engineering, University of California, Berkeley, 760 Davis Hall, Berkeley, CA 94720, amins@berkeley.edu

[§]California Center for Innovative Transportation, 2105 Bancroft Way, Suite 300, Berkeley, CA 94720, tania.abou-nasr@polytechnique.edu

[¶]California Center for Innovative Transportation, 2105 Bancroft Way, Suite 300, Berkeley, CA 94720, ana36@aub.edu.lb

^{||}Department of Electrical Engineering and Computer Science, University of California, Berkeley, 746 Sutardja Dai Hall #1758, Berkeley, CA 94720, pabbeel@cs.berkeley.edu

^{**}Department of Civil and Environmental Engineering, Systems Engineering, University of California, Berkeley, 642 Sutardja Dai Hall, Berkeley, CA 94720, bayen@berkeley.edu

Abstract

This article introduces the new component of *Mobile Millennium* dedicated to arterial traffic. *Mobile Millennium* is a pilot system for collecting, processing and broadcasting real-time traffic conditions through the use of GPS equipped smartphones. Two algorithms that use data from GPS equipped smartphones to estimate arterial traffic conditions are presented, analyzed and compared. The algorithms are based on *Logistic Regression* and *Spatio-Temporal Auto Regressive Moving Average* (STARMA), respectively. Each algorithm contains a *learning* component, which produces estimates of spatio-temporal parameters for describing interactions between the states of arterial links in the network. Additionally, each algorithm contains an *inference* component, which gives the procedure for processing real-time data into short-term forecasts using these parameters. The algorithms are tested with simulation data obtained from Paramics software, and from a field test in New York. Both methods provide encouraging results in forecasting arterial traffic conditions using sparse GPS data.

$_{1}$ 1 Introduction

In the United States and numerous other parts of the world, traffic is an unavoidable part of economic activity. The 2007 Urban Mobility Report [3] states that traffic congestion causes 4.2 billion hours of extra travel in the United States every year, which accounts for 2.9 billion extra gallons of fuel, which cost taxpayers an additional \$78 billion.

Numerous measures can be taken to address problems due to traffic congestion. An essential step is to create the ability to forecast traffic conditions with significant 8 accuracy and reliability. Numerous challenges stand in the way of this type of effort. 9 A significant portion of the transportation network has little or no dedicated infras-10 tructure for collecting traffic data. Areas equipped with this infrastructure generally 11 only cover highways and have high installation and maintenance costs in addition to 12 providing data of variable reliability. An alternative to using dedicated communication 13 infrastructure is to leverage an existing system such as the cellular phone network. 14 The Mobile Millennium project [2] was conceived as a response to these challenges, to 15 explore the capability of using cellular phones to provide traffic data. 16

The mobile internet is the underlying technology enabling the existence of the 17 Mobile Millennium system. User-generated content (in the present case, smartphone 18 measured traffic data) is sent to a central system, which provides information back 19 to the cell phone owner for personal use. This "web 2.0" application framework is 20 commonly referred to as "participatory sensing," which refers to the ad hoc process 21 of voluntarily providing sensing data to a system. In general, there are a number 22 of challenges to overcome with any nomadic sensing technology, including unknown 23 location of upcoming measurements, sparsity of the data, and unpredictability of the 24 frequency of data collection. 25

The Mobile Millennium system was officially launched on November 10, 2008 when 26 the team released a software client for GPS enabled smartphones to the public, avail-27 able for download (see figure 1(a)). Traffic conditions are broadcast back to drivers' 28 mobile phones, enabling commuters to make more informed route and trip decisions. 29 Additionally, traffic data can be analyzed in the *Mobile Millennium* live traffic visual-30 izer, shown in figure 1(b), which is currently on display in the CITRIS [1] Tech Museum 31 on the UC Berkeley campus. The deployment area of the pilot system is focused on 32 commuters in Northern California, including the San Francisco Bay Area and Sacra-33 mento. The project is a follow up to the *Mobile Century* experiment, in which 165 UC 34 Berkeley graduate students were hired to drive a 10-mile stretch of I880 in California 35 for a day, demonstrating the feasibility of a real-time highway traffic estimation service 36 using only GPS enabled devices [29]. 37

This article focuses on estimating and forecasting *arterial* traffic conditions using only GPS enabled devices using two statistic models. Section 2 reviews previous traffic models and examines the need for a new (statistical) approach. Section 3 is a formal presentation of the problem with details of the data types and models used. Sections 4 and 5 present the details of the logistic regression and STARMA models, respectively. Results are presented on simulation and field experiment data in section 6. Further analysis and future directions are presented in the conclusion (section 7).



(a) Traffic client

(b) Web interface

Figure 1: *Mobile Millennium* traffic information services. (a) The *Mobile Millennium* traffic client. (b) The *Mobile Millennium* interface for highway and arterial traffic visualization.

⁴⁵ 2 Challenges and Motivation for the Statistical ⁴⁶ Learning Approach

The development of traffic theory has led to numerous modeling contributions since 47 the pioneering work of Lighthill, Whitham and Richards, [15, 21], which relied on hy-48 drodynamic theory. By nature, arterial traffic has very high variability, which make 49 it challenging to use flow models for arterial networks. Studies have typically focused 50 on modeling single intersections [22, 28, 4, 23] using dedicated traffic sensors. This 51 modeling approach is difficult to adapt to a general traffic information system on a 52 dense arterial network because it requires a high density of traffic sensors (which are 53 prohibitively expensive at the scale of the arterial network). In particular, one of the 54 major challenges of *Mobile Millennium* is that the system does not have access to flow 55 counts. A statistical approach is suitable because sensing every vehicle is impracti-56 cal and because this allows for the incorporation of other information types (such as 57 human mobility patterns [10]). The present article article suggests to develop moni-58 toring capabilities for arterial traffic in two directions: (1) using alternate data sources 59 such as privacy aware cell phone information; (2) developing new arterial models based 60 on statistical learning which overcome some major issues faced by analytical flow or 61 queuing-based models. The motivations for these two items are described in the re-62 mainder of this section. 63

⁶⁴ 2.1 Using Cell Phones as Traffic Probes

Experimental research on cell phone based traffic monitoring [5, 30, 24, 31] has investigated the ability to locate the position of users using trilateration- or triangulationbased methods. It has shown limited success for estimation of travel times due to the position measurement inaccuracy, particulary on short distances and dense networks [16, 12]. The complexity of traffic patterns in the arterial networks gives the use of GPS-based traffic information enormous growth potential.

⁷¹ 2.2 Application of Machine Learning to Arterial Traffic

Machine learning techniques have been used to estimate and produce short term traffic 72 predictions for both freeway [8, 27, 7, 6, 17] and arterial networks [26, 28, 14, 19, 25, 9]. 73 These studies present encouraging results. One of the limitations of these approaches 74 for our problem is that they present results for specific traffic variables (in particular 75 for flow/density). The present article focuses on estimating and predicting congestion 76 states and travel times. Congestion states, also referred to in the literature as Level of 77 Service (LoS), represent traffic conditions on the road segment as experienced by the 78 network user. They also represent the level of service offered by the network manager. 79 They can be interpreted as a discrete representation of traffic states. Traffic states 80 (for example travel time) have their own statistical distribution depending on traffic 81 conditions. 82

In the remainder of this article, we present regression techniques corresponding to two different approaches:

- Logistic Regression model. An example of a neural network used as a clustering algorithm between discrete traffic congestion states.
- Spatio-Temporal Auto-Regressive Moving Average (STARMA) model. An example of a time series model in which the traffic variable studied (travel time) depends on the previous values of this variable.

To our knowledge, logistic regression has not been used in arterial traffic estima-90 tion or prediction. We compare its results to a more widely used model, the STARMA 91 model. This comparison is of significant interest for the transportation community 92 since it is between a discrete output (congestion states from the clustering of the lo-93 gistic regression) and a continuous output (travel time estimations from the linear 94 regression of the STARMA model). Furthermore, the system and results presented 95 here are one of the first instantiations of real-time arterial monitoring using machine 96 learning with streaming data collected from smartphone. Mobile Millennium is cur-97 rently implemented and operational in all of Northern California [2]. 98

⁹⁹ 3 Problem Formulation

83

84

85

86

87

88

89

This section formally presents the problem formulation, namely estimating LoS indicators which are the aggregate travel times and congestion states for an arterial road network. First, we introduce our sampling paradigm, the *Virtual Trip Line* in section 3.1. This leads to the problem of sensing on a graph (section 3.2) and the formal definitions of LoS indicators (section 3.3). The problem description of estimating the LoS indicators based on STARMA and logistic regression is presented in section 3.4.

¹⁰⁶ 3.1 Virtual Trip Line Sensing Infrastructure

A GPS-enabled smartphone is capable of recording its GPS location every few seconds. Over time, this vehicle trajectory information produces a rich history of the vehicle and the velocity field through which it evolves [11]. While this level of detail can be useful for traffic estimation, it can be privacy invasive, since the device is ultimately carried by a single user. Even if personally identifiable information from the data is replaced with a randomly chosen ID through a process known as pseudo-anonymization, it is still possible to re-identify individuals from trajectory data [13].

Virtual Trip Lines (VTLs) [12] are spatial triggers for phones to collect measurements and send updates. Each VTL consists of two GPS coordinates which make a virtual line drawn across a roadway of interest. Instead of time-based periodic sampling, VTLs trigger disclosure of speed and location updates by sampling in space, creating updates at predefined geographic locations on roadways of interest. Additionally, the travel time between pairs of VTLs can be extracted and this type of travel time data will be considered the primary data source used in this article.

3.2 Graph Model of the Road Network

114

115

116

117

118

119

120

121

142

145

Consider an arterial network with a total of N pairs of VTLs deployed. Each pair has 122 a unique identification number $i \in \{1, \ldots, N\}$. The set of all VTL pairs is denoted by 123 $\mathcal{V} = \{1, \ldots, N\}$. Each VTL pair has a segment of road in between with a possibility 124 of one or more road features such as an intersection (with or without traffic lights), 125 pedestrian walkways, stop/slow signs etc. The characteristics of these road features 126 can be static (such as presence of a stop sign) or dynamic (such as phase of a signalized 127 intersection) with respect to time. The travel time experienced by a vehicle traveling 128 through a VTL pair depends on the characteristics of the road features as well as the 129 demand-capacity restrictions imposed by the dynamics of traffic flow. We also assume 130 that each VTL pair is associated with unidirectional traffic flow. For arterial links 131 consisting for bidirectional traffic, we associate a VTL pair corresponding to each flow 132 direction. 133

We say that the upstream (resp. downstream) VTL for the pair i is the VTL at 134 which the traffic enters (resp. leaves) the corresponding stretch of road. For pair i, let 135 the upstream and downstream VTLs be denoted by i_u and i_d respectively. Then the 136 VTL sensor network can be represented as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the 137 set of all VTL pair as defined earlier and \mathcal{E} is the set of all edges. Two VTL pairs i 138 and j form an edge directed from pair i to pair j, denoted e_{ij} , if i_d and j_u correspond 139 to same VTL. Then i (resp. j) is called the upstream (resp. downstream) node of edge 140 141 e_{ij} .

We define the set of first order neighbors for VTL pair j as

$$\mathcal{N}^1(j) = \{j\} \cup \{i \in \mathcal{V} : e_{ij} \in \mathcal{E}\} \cup \{k \in \mathcal{V} : e_{jk} \in \mathcal{E}\}$$

which is simply the set of all the upstream and downstream VTL pairs for the pair j(in which we include pair j itself).

We can extend the above definition to define n^{th} $(n \ge 1)$ order neighbors as:

$$\begin{cases} \mathcal{N}^{0}(j) = \{j\} \\ \mathcal{N}^{n}(j) = \mathcal{N}^{n-1}(j) \cup \left(\bigcup_{l \in \mathcal{N}^{n-1}(j)} \{i \in \mathcal{V} : e_{il} \in \mathcal{E}\} \cup \{k \in \mathcal{V} : e_{lk} \in \mathcal{E}\}\right) \end{cases}$$
(1)

¹⁴⁶ 3.3 Traffic Level of Service Indicators

We assume that for any VTL pair $i \in \mathcal{V}$, the travel time data is available at times $0 \leq t_1 \leq t_2 \leq \ldots$ As an alternative to travel time data, we can also compute the pace (travel time divided by the length of road for the VTL pair). We denote the data obtained at time t_1 for VTL pair i as $X_{t_1,i}$. Since the acceptable values of $X_{t_1,i}$ generally lie between a minimum and maximum value, we reject data that do not fall in this range. For VTL pair i, let us denote this range as $[\underline{X}_i, \overline{X}_i]$.

Since the data obtained is event-based, it cannot be directly used for training statistical models that needs regular sampling rates. We aggregate the travel time data in t second windows to obtain a time series of observations at times k = 0, t, 2t, ...Here t is the aggregation interval. Henceforth, we will use k to denote the time interval [(k-1)t, kt). The set of available observations during the time period k for any VTL pair i is denoted as $A_{k,i}$, that is,

$$A_{k,i} = \{ X_{t_m,i} \mid (k-1)t \le t_m < kt \}$$

The penetration rate for VTL pair *i* during time *k*, denoted $p_{k,i}$, is the fraction of available observations out of the total number of vehicles $D_{k,i}$ traveling through pair *i* during time *k*:

$$p_{k,i} = \frac{A_{k,i}}{D_{k,i}}.$$

We define the spatial aggregation function for VTL pair $i, h_i(\cdot) : A_{k,i} \mapsto [\underline{X}_i, \overline{X}_i]$, as the function that aggregates the set of observations $A_{k,i}$ in to an aggregate representative quantity, denoted $Z_{k,i}$ with values in the range $[\underline{X}_i, \overline{X}_i]$. In the remainder of this article, $Z_{k,i}$ is an aggregated travel time (seconds). Thus, the aggregate travel time for VTL *i* during interval *k* is

$$Z_{k,i} = h_i(\{X_{t_m,i} \mid (k-1)t \le t_m < kt\}).$$

We define the *mode* of a VTL pair as the categorical variable indicative of the extent 167 of delay experienced in navigating through the VTL pair. For example, a binary mode 168 classification can be uncongested or congested. Thus, the mode of a VTL pair can also 169 interpreted as a *congestion state*. Let the mode of VTL pair i during time interval k be 170 denoted as $Q_{k,i}$. In order to convert the total number of observations available at VTL 171 i during time interval k into a mode of the VTL pair, we define a congestion indicator 172 function $g_i(\cdot): A_{k,i} \mapsto \{1, \ldots, M\}$ where M is the desired number of modes the VTL 173 pairs should be classified to. Thus, 174

$$Q_{k,i} = g_i(\{X_{t_m,i} | (k-1)t \le t_m < kt\}).$$

From a statistical modeling perspective, both the aggregate speed or travel time, $Z_{k,i}$, and the congestion state, $Q_{k,i}$, for $i \in \mathcal{V}$ and $k \in \{0, 1, ...\}$ can be considered as random processes generated by space-time varying traffic flow phenomena on the arterial network. Both $Q_{k,i}$ and $Z_{k,i}$ can be regarded as LoS indicators.

¹⁷⁹ 3.4 Estimating Level of Service Indicators

If we had data from all the vehicles for all the VTL pairs over the entire time horizon of interest, the penetration rate $p_{k,i}$ would satisfy $p_{k,i} = 1$ for all $i \in \mathcal{V}$ and $k \in \{0, 1, \ldots\}$. We could then compute the entire probability distribution of $Z_{k,i}$ and $Q_{k,i}$. However, the challenge of arterial traffic state estimation and forecast is that the typical penetration rates are very low. Our focus in this article is to develop reliable estimation and forecasting methods for such situations.

We now describe the problem formulation for estimation or $nowcast^1$. We typically 186 only have data from a small percentage of the total number of vehicles $(p_{k,i} \sim 0.02 -$ 187 0.05). Thus, the choice of the aggregation function $h_i(\cdot)$ (resp. the congestion indicator 188 function $g_i(\cdot)$ becomes critical to obtain reliable estimates of $Z_{k,i}$ (resp. $Q_{k,i}$). For 189 a given choice of $h_i(\cdot)$, the best estimate of the aggregate travel time or speed for 190 VTL pair i during interval k is given by the conditional expectation of $Z_{k,i}$ given the 191 aggregate travel times up-to (and excluding) the current time interval: 192

$$\hat{Z}_{k,i} = \mathbb{E}[h_i(A_{k,i}) | h_j(A_{j,v}), j < k, v \in \mathcal{V}] = \mathbb{E}_{h_i}[Z_{k,i} | Z_{j,v}, j < k, v \in \mathcal{V}],$$

where notation $\mathbb{E}_{h_i}[\cdot]$ is used to indicate the dependence of the expectation on the 193 aggregation function h_i . We now introduce the following conditional independence 194 assumption: $Z_{k,i}$ is conditionally independent of all other data conditioned on the data 195 from the past r time intervals for VTL pairs in the set $\mathcal{N}^{s}(i)$. Under this assumption, 196 we can write 197

$$\hat{Z}_{k,i} \approx \mathbb{E}_{h_i}[Z_{k,i} | Z_{j,v}, k - r \le j < k, v \in \mathcal{N}^s(i)]$$
⁽²⁾

Thus, $Z_{k,i}$ only depends on data with r temporal dependencies in the past and s 198 spatial dependencies from the neighbors. Similarly, for given choices of the aggregation 199 function $h_i(\cdot)$ and the congestion indicator function $q_i(\cdot)$, we can write the conditional 200 expectation of $Q_{k,i}$ given all the aggregate travel times up-to (and excluding) the 201 current time interval as^2 202

$$\hat{Q}_{k,i} = \mathbb{E}[g_i(A_{k,i})|h_j(A_{j,v}), j < k, v \in \mathcal{V}]$$

$$\approx \mathbb{E}_{h_i, q_i}[Q_{k,i}|Z_{j,v}, k - r \le j < k, v \in \mathcal{N}^s(i)], \qquad (3)$$

In the statistics terminology, the quantities $Z_{k,i}$ and $Q_{k,i}$ in (2) and (3) are known as 203 the response variables; the conditioned variables $Z_{j,v}$ and $Q_{j,v}$ are called the dependent 204 variables or covariates. The present article compares the two estimators which we now 205 introduce. 206

The first estimator is based on expressing (2) as a *linear regression problem*. For 207 a temporal and spatial dependence of orders r and s respectively, we assume a linear 208 dependence of response $Z_{k,i}$ on the covariates $Z_{j,v}$: 209

$$\hat{Z}_{k,i} = \beta_i^0 + \sum_{v \in \mathcal{N}^s(i)} \left(\sum_{j=k-r}^{k-1} \beta_i^{j,v} Z_{j,v} \right).$$

$$\tag{4}$$

210

In order to make the notation concise, let $\mathbf{Z}_{k,i}^{r,s}$ be the $r \times \mathcal{N}^{s}(i)$ vector of covariates or dependent variables obtained by stacking the aggregate travel times $Z_{i,v}$ for $k-r \leq 1$ 211

¹Depending on the convention used, this can also be treated as one-step ahead forecast. In this article, we do not distinguish between one-step forecast and estimation.

²Alternatively, we can also condition $Q_{k,i}$ directly on congestion modes up-to (and excluding) the current time, that is, $\hat{Q}_{k,i} = \mathbb{E}_{g_i}[Q_{k,i}|Q_{j,v}, k-r \leq j < k, v \in \mathcal{N}^s(i)]$. However, we do not consider this type of estimator in this article.

 $j < k \text{ and } v \in \mathcal{N}^{s}(i), \beta_{i} \text{ be the corresponding } r \times \mathcal{N}^{s}(i) + 1 \text{ vector of parameters to be}$ estimated. Then the equation (4) can be re-written as

$$\hat{Z}_{k,i} = \beta_i^\top \mathbf{Z}_{k,i}^{r,s},$$

where \top stands for the transpose of a vector. As described later in Section 5, instead of a simple regression model (4), we consider a STARMA model.

Our second estimator is based on expressing (3) as a logistic regression problem which assumes a linear dependence of the *logit* or the log-odds ratio of conditional expectation $\hat{Q}_{k,i}$ on the response variables. That is, for a temporal and spatial dependence of orders r and s respectively, we have

$$\log\left(\frac{\hat{Q}_{k,i}}{1-\hat{Q}_{k,i}}\right) = \beta_i^\top \mathbf{Z}_{k,i}^{r,s}$$

220 We can express this equation as

$$\hat{Q}_{k,i} = f_{\beta_i}(\mathbf{Z}_{k,i}^{r,s}) := \frac{1}{1 + \exp\left(-\beta_i^\top \mathbf{Z}_{k,i}^{r,s}\right)},\tag{5}$$

where the subscript β_i in $f_{\beta_i}(\cdot)$ encodes the dependence on the β_i .

We detail the implementation of a logistic regression estimator in Section 4 and a STARMA-based estimator in Section 5. However, two important points need to be mentioned. First, the above formulation can be modified to include the case of multiple steps forecast. For example, an m-step forecast at time k for VTL pair i can be written as

$$\hat{Z}_{k+m,i} = \mathbb{E}_{h_i}[Z_{k,i}|Z_{j,v}, j < k, v \in \mathcal{V}], \tag{6}$$

where we consider data up to time k to predict traffic at time k + m.

Second, we note that for some VTL pairs and time intervals, we might not have any available data, that is, $A_{j,v} = \emptyset$ for some $j \in \{k - r, ..., k\}$ and $v \in \mathcal{N}^s(i)$. In this case, one has to employ a technique of *estimation with missing data*. We will briefly touch on the forecast problem for the STARMA model but will address the issue of missing data in later work.

4 Logistic Regression

We now discuss the estimator based on the logistic model (5) to estimate the congestion state $Q_{k,i}$ for a VTL pair *i* and time interval *k*. Suppose that $Q_{k,i}$ is binary-valued, that is $Q_{k,i} = \{0, 1\}$ and M = 2. When $Q_{k,i} = 1$ (resp. $Q_{k,i} = 0$), we say that the VTL pair *i* during interval *k* is in the congested mode (resp. uncongested mode). Then the estimator $\hat{Q}_{k,i}$ gives the conditional probability of the $Q_{k,i}$ given the dependent variables:

$$\hat{Q}_{k,i} = \mathbb{E}_{h_i,g_i}[Q_{k,i}|\mathbf{Z}_{k,i}^{r,s}] = 1 \cdot \mathbb{P}_{h_i,g_i}[Q_{k,i} = 1|\mathbf{Z}_{k,i}^{r,s}] + 0 \cdot \mathbb{P}_{h_i,g_i}[Q_{k,i} = 0|\mathbf{Z}_{k,i}^{r,s}] \\ = \mathbb{P}_{h_i,g_i}[Q_{k,i} = 1|\mathbf{Z}_{k,i}^{r,s}]$$

Now using (5), we can write the conditional probability of $Q_{k,i}$ given the aggregate travel time for r temporal and s spatial dependencies as

$$\mathbb{P}_{h_i,g_i}[Q_{k,i}|\mathbf{Z}_{k,i}^{r,s};\beta_i] = [f_{\beta_i}(\mathbf{Z}_{k,i}^{r,s})]^{Q_{k,i}}[1 - f_{\beta_i}(\mathbf{Z}_{k,i}^{r,s})]^{1 - Q_{k,i}}$$

We now assume that for a VTL pair *i*, the response process $\{Q_{k,i}\}$ and the covariate process $\{\mathbf{Z}_{k,i}^{r,s}\}$ is available for a number of time intervals $k = 0, \ldots, K$. Introducing the conditional independence assumption that the response variable $Q_{k,i}$ is independent of all other data given $\mathbf{Z}_{k,i}^{r,s}$. Then the joint conditional probability of $\{Q_{k,i}\}$ given $\{\mathbf{Z}_{k,i}^{r,s}\}$ (also known as the conditional likelihood) can be expressed as

$$\mathbb{P}_{h_i,g_i}[\{Q_{k,i}\}_{k=0}^K | \{\mathbf{Z}_{k,i}^{r,s}\}_{k=0}^K; \beta_i] = \prod_{k=0}^K [f_{\beta_i}(\mathbf{Z}_{k,i}^{r,s})]^{Q_{k,i}} [1 - f_{\beta_i}(\mathbf{Z}_{k,i}^{r,s})]^{1-Q_k},$$

For a given training data $\{Q_{k,i}\}_{k=0}^{K}$ and $\{\mathbf{Z}_{k,i}^{r,s}\}_{k=0}^{K}$, the *best* estimate of parameter β_i is obtained by maximizing the logarithm of the conditional likelihood which we state explicitly as follows:

$$\mathcal{L}(\beta_i; \{Q_{k,i}\}_{k=0}^K, \{\mathbf{Z}_{k,i}^{r,s}\}_{k=0}^K) = \sum_{k=0}^K \left(Q_{k,i} \cdot \beta_i^\top \mathbf{Z}_{k,i}^{r,s} - \log\left[1 + \exp\left(\beta_i^\top \mathbf{Z}_{k,i}^{r,s}\right) \right] \right)$$

The optimal estimate so obtained and denoted β_i^* , is called the *maximum likelihood estimate* (MLE). A number of standard iterative methods, all similar to Newton-Raphson method, can be used to obtain the MLE β_i^* . Examples of such method include Fisher scoring method, iterative reweighted least squares etc. Due to space limitations, we omit the details of the algorithm and refer the reader to [18].

Once the parameters are learned, *validation* can be done on a similar data set as the one used to obtain β_i^* . Validation is done to assess the ability of the learned model to correctly estimate the traffic status (congestion state in this case) on previously unseen data.

²⁵⁹ 5 STARMA

250

251

252

253

254

255

256

257

258

We now discuss the STARMA model which is a more efficient estimator than the simple linear regression model (4). The number of parameters to be estimated for (4), given by $r \times |\mathcal{N}^{s}(i)| + 1$ (|A| is the cardinality of A), can increase significantly as the spatial dependency s increases. In order to explain the model, we first present the *spatiotemporal autoregressive* (STAR) model and subsequently generalize to a full STARMA model.

Following (1), the set of n order neighbors $(0 \le n \le s)$ for a VTL pair i can be expressed as follows

$$\mathcal{N}^{s}(i) = \bigcup_{n=0}^{s} \mathcal{N}^{n}(i) \setminus \mathcal{N}^{n-1}(i).$$

Here we adopt the convention that $\mathcal{N}^{0}(i) \setminus \mathcal{N}^{-1}(i) = \{i\}$. Now, for the linear regression model (4), for any temporal order j, $(k - r \leq j < k)$ and spatial order n, $(0 \leq n \leq s)$, we introduce the assumption that

For all
$$v \in \mathcal{N}^n(i) \setminus \mathcal{N}^{n-1}(i), \quad \beta_i^{j,v} \equiv \beta_i^{j,n},$$
 (7)

and the definition of n-th order, spatially-weighted travel time as

$$\varphi_{i}^{(n)}(Z_{j}) = \frac{\sum_{l \in \mathcal{N}^{n}(i) \setminus \mathcal{N}^{n-1}(i)} w_{i,l}^{(n)} Z_{j,l}}{\sum_{l \in \mathcal{N}^{n}(i) \setminus \mathcal{N}^{n-1}(i)} w_{i,l}^{(n)}},$$
(8)

where $Z_j = (Z_{j,1}, \ldots, Z_{j,N})$ is the vector of aggregate travel times for all the N VTL pairs during time interval j and $w_{i,l}^{(n)}$ are the pre-defined *spatial weights of order n* for $Z_{j,l}$.

Under the assumption (7) and the definition (8), we can now write the STAR model of *autoregressive* (AR) temporal order r and spatial order s as

$$Z_{k,i} = \sum_{j=k-r}^{k-1} \sum_{n=0}^{s} \beta_i^{j,n} \varphi_i^{(n)}(Z_j) + \epsilon_{k,i}$$
(9)

where $\epsilon_{k,i}$ is the normally distributed error term with variance σ^2 with the properties that $\mathbb{E}[\epsilon_{k,i}] = 0$ for all k and $i \in \mathcal{V}$; and for all $i, j \in \mathcal{V}$

$$\mathbb{E}[\epsilon_{k,i}\epsilon_{k+s,j}] = \begin{cases} \sigma^2 & \text{if } s = 0\\ 0 & \text{otherwise.} \end{cases}$$

The number of parameters to be estimated for the STAR model (9), including σ^2 , is r(s+1)+1 which is (typically) much smaller than $r \times \mathcal{N}^s(i)+1$ for (4). The STAR model can now be generalized to STARMA model of autoregressive temporal order rand spatial order s, and moving average (MA) temporal order p and spatial order qas³

$$Z_{k,i} = \sum_{j=k-r}^{k-1} \sum_{n=0}^{s} \beta_i^{j,n} \varphi_i^{(n)}(Z_j) - \sum_{j=k-p}^{k-1} \sum_{n=0}^{q} \alpha_i^{j,n} \varphi_i^{(n)}(\epsilon_j) + \epsilon_{k,i},$$
(10)

where $\epsilon_j = (\epsilon_{j,1}, \dots, \epsilon_{j,N})^\top$.

Here $\alpha_i^{j,n}$ are the moving average parameters. The total number of parameters (including σ^2) to be estimated for the STARMA model (10), denoted as STARMA(r, s, p, q)are r(s+1) + p(q+1) + 1.

Following [20], we adopt the assumption in this article that that STARMA parameters are same for VTL pairs, that is, $\alpha_1^{j,n} = \ldots = \alpha_N^{j,n} \equiv \alpha_{j,n}$ and $\beta_1^{j,n} = \ldots = \beta_N^{j,n} \equiv \beta_{j,n}$. Then model (5) can be vectorized for all VTL pairs $i \in \mathcal{V}$ as

$$Z_k = \sum_{j=k-r}^{k-1} \sum_{n=0}^s \beta^{j,n} \Phi^{(n)}(Z_j) - \sum_{j=k-p}^{k-1} \sum_{n=0}^q \alpha^{j,n} \Phi^{(n)}(\epsilon_j) + \epsilon_k.$$
(11)

291

270

where $\Phi^{(n)}(\cdot) = (\varphi_1^{(n)}(\cdot), \dots, \varphi_N^{(n)}(\cdot))^\top$ and $\epsilon_k = (\epsilon_{k,1}, \dots, \epsilon_{k,N})^\top$.

³More generally, the AR spatial order s (resp. the MA spatial order q) can vary with the temporal order r (resp. p). However, we do not consider this generalization in this article.

For given training data $\{Z_k\}$, (k = 0, ..., K-1), the best estimate of the parameters $A := [\alpha^{j,n}]_{p \times (q+1)}$, $B := [\beta^{j,n}]_{r \times (s+1)}$ and σ^2 is given by maximizing the conditional likelihood expressed as

$$\mathbb{P}(\{Z_k\}_{k=0}^{K-1}; A, B, \sigma^2) = (2\pi)^{-\frac{KN}{2}} |\sigma^2 \mathbf{I}_{KN \times KN}|^{-\frac{1}{2}} \exp\left(-\frac{S(A, B)}{2\sigma^2}\right)$$
(12)

where $I_{KN \times KN}$ is the identity matrix, $S(A, B) := (\epsilon_0, \dots, \epsilon_{K-1})^\top (\epsilon_0, \dots, \epsilon_{K-1})$ and according to (11), we have

$$\epsilon_k = Z_k - \sum_{j=k-r}^{k-1} \sum_{n=0}^s \beta^{j,n} \Phi^{(n)}(Z_j) + \sum_{j=k-p}^{k-1} \sum_{n=0}^q \alpha^{j,n} \Phi^{(n)}(\epsilon_j).$$

The maximum likelihood estimate parameters, denoted A^*, B^* , are obtained by maximizing the logarithm of the conditional likelihood (12), and the corresponding σ^* is estimated by

$$\sigma^* = \sqrt{\frac{S(A^*,B^*)}{KN}}$$

For further details, we refer the reader to [20].

301 6 Results

This section presents the results from logistic regression based classification and STARMAbased continuous linear regression. Each algorithm is implemented and tested on simulation and field experiment data, described in section 6.1. The framework for quantifying accuracy is described in section 6.2. Results are then presented for one-step forecast (section 6.3), followed by multi-step forecast for the STARMA model (section 6.5). Additionally, a study of the effect of the *penetration rate* on the forecast accuracy (section 6.4) is presented.

³⁰⁹ 6.1 Simulation and Field Experiment Data

There are two data sets used in this article. The first set was generated from Paramics 310 micro-simulation software. The road network modeled consists of 1,961 nodes, 4,426 311 links, 210 zones and is based on the SR41 corridor in Fresno, CA. We specifically 312 analyzed a sub-network that includes 9 arterial roads, 20 signals and 15 stop signs. 313 Paramics simulates every car in the network. From this simulation, we extract the po-314 sition of every vehicle at one-second time intervals. This provides detailed information 315 about speed and travel time through the network. The sub-network studied in this 316 article includes 380 different links, each one of which is characterized with a specific 317 length, a number of lanes, a direction, a speed limit and signal information. 99 VTLs 318 were placed on different links, which corresponds to 156 different pairs of VTLs, in 319 order to capture travel times along links and through intersections. 320

The second data set was obtained as part of the official *Mobile Millennium* launch demonstration in New York City at the *ITS World Congress*. Twenty drivers, each carrying a GPS equipped cell phone, drove for 3 hours (9:00am to 12:00pm) around a



Figure 2: **Experiment Design.** (a) Map of the Paramics network in Fresno, CA. (b) Experiment route for New York City field test used to collect the data (arrows represent the direction of traffic of probe vehicles). (c) Test vehicle used for the New York test.

2.4 mile loop of Manhattan (see figure 2). This number of drivers constituted approximately 2% of the total vehicle flow through the road of interest. The experiment was repeated 3 times in order to use two of the experiments as training data for the models and the other to validate the model results. The operational capabilities of the system were demonstrated at the *ITS World Congress* [2] on November 18, 2008, when live arterial traffic was displayed for conference attendees.

6.2 Validation Framework

In order to compute the accuracy of the model, one needs to define the "ground truth" 331 state of traffic. In this article, travel times are aggregated into a single value per time 332 interval (5 minutes for Paramics, 15 minutes for New York). This single value per time 333 interval is considered the true state for the interval. Determining ground truth for 334 the logistic regression method requires classifying each time interval as congested or 335 uncongested. The STARMA method uses the average travel time during each interval 336 as the ground truth value. Both of these methods correspond to choosing appropriate 337 $h_i(\cdot)$ and $g_i(\cdot)$ functions as described in section 3.3. 338

The aggregation function $h_i(\cdot)$ should capture the pattern of change in pace over 339 different intervals to provide an aggregate quantity that is sufficiently representative of 340 the congestion state, thus providing better accuracy in training the model and obtaining 341 the logistic regression parameters. Based on extensive testing and simulation, it is 342 observed that aggregating the travel times based on the entire data available in an 343 interval fails to capture the congestion state due to the high variance of travel times 344 when a link is congested. The probes most affected by congestion should thus have more 345 weight in the aggregation process. A simple yet fairly effective data-driven aggregation 346 method is as follows: given the set of observations for VTL pair i and interval k, $A_{k,i}$ 347



Figure 3: Average estimation accuracy vs. aggregation parameter w.

is sorted such that $t_{m_1} < t_{m_2} \implies X_{t_{m_1},i} > X_{t_{m_2},i}$, then take 348

$$Z_{k,i} = h_i(\{X_{t_m,i} \,|\, (k-1)t \le t_m < kt\}) := \frac{1}{w.M_{k,i}} \sum_{m=1}^{\lfloor w.M_{k,i} \rfloor} X_{t_m,i},$$

where $M_{k,i}$ is the number of observations in $A_{k,i}$ and $0 < w \leq 1$ is the fraction 349 of observations used for aggregation. The symbol |a| denotes the floor value of a. In 350 words, the aggregate pace is the mean of the $100 \times w\%$ observations with highest pace 351 or equivalently the worst observations. The simulation results for different values of w352 are shown in figure 3. From an application-driven point of view, we select the w that 353 maximizes estimation accuracy, in the present case w = 0.3. At this value, the travel 354 time envelope of the time series of observations is best captured. 355

The training phase of logistic regression requires as input a congestion threshold 356 along with the aggregate travel times $Z_{k,i}$. Since the congestion threshold should be 357 chosen to be consistent with the choice of aggregate travel times to provide meaningful 358 classification, we define the congestion threshold, T_i as the mean of the $100 \times w\%$ 359 observations in D_i with highest travel time where D_i is the set of available observations 360 in all intervals and w is essentially be the same value chosen for aggregation (w = 0.3361 in this section). This corresponds to choosing 362

$$Q_{k,i} = g_i(\{X_{t_m,i} | (k-1)t \le t_m < kt\}) = I(h_i(\{X_{t_m,i} | (k-1)t \le t_m < kt\}) > T_i),$$

where $I(\cdot)$ is the indicator function. The STARMA model does not use a g_i function 363 because it forecasts a continuous quantity. 364

The logistic regression algorithm produces a probability of congestion for each VTL 365 pair studied. If this probability is greater than .5, then the forecasted state is congested. 366 The accuracy of the logistic regression forecasts is defined as the percentage of correctly forecasted states over all intervals and VTL pairs studied. For the STARMA model, 368 the accuracy is defined as the percentage error between the forecasted travel time value 369 and the actual travel time value as defined by the h_i function described earlier. 370

367

6.3 Short-Term Forecast

Both regression methods are designed to do one-step (short-term) forecasts. For each 372 data set (as described in section 6.1), the performance of each model was evaluated 373 by dividing the data set into a training set and a validation set. For the Paramics 374 simulation data, the training set consisted of three simulation runs and the validation 375 set consisted of a separate, fourth simulation run. For the New York experiment data. 376 two days of data were used for training and the other day for validation. Through 377 a-priori experimentation, the temporal dependency for the logistic regression model 378 was set to r = 1 for the logistic regression, r = 2 for the STARMA model. The spatial 379 dependency is varied for comparison in the result figures described in the following 380 paragraph. 381

The Paramics simulations give information about every vehicle. For testing the 382 methods, only a subset of the data is used for training and inference, corresponding 383 to the penetration rate. This was incorporated into the following analysis by requiring 384 each regression method to produce estimates for the validation data set using only a 385 small percentage of the available travel times. Figure 4 displays the one-step forecast 386 results of the logistic regression and STARMA methods on the Paramics validation set 387 respectively, using a penetration rate of 5%. Similarly, figure 5 displays the one-step 388 forecast results on the New York validation set. 389

³⁹⁰ 6.4 Penetration Rate Study

The value of 5% for the penetration rate used in section 6.3 was chosen based on the 391 prospects for future adoption of GPS equipped cell phones running traffic informa-392 tion software (such as that provided by *Mobile Millennium*). Therefore, a study of 393 the effect of the penetration rate on results is of interest to quantify the influence of 394 technology adoption on estimation and forecast accuracy. Figure 6 shows the one-step 395 forecast accuracy for the logistic regression and STARMA methods as a function of 396 the penetration rate. From these figures, one can infer that 2% penetration rate can 397 give reasonably good results, while 5% and higher give very accurate results. We also 398 note that using spatial neighbors of order 1 (direct neighbors) generally provides better 399 results. One can interpret this as indicating that second order neighbors lead to an 400 overfit model while no neighbors lead to an underfit model. 401

402 6.5 Multi-Step Forecast

The STARMA model is capable of producing forecasts of any number of steps by using the output of the model as input for the next time interval. It is not straightforward to do the same for the logistic regression model since it has an output that is fundamentally different from the input it requires. Therefore, the discrete output of the logistic regression model must be transformed back to a continuous value in order to do forecast in the same way. This avenue is not considered in this article and is left as further research.

In this section, the results of multi-step forecast for the STARMA model are presented. Figure 7 shows the forecast results for the New York data set. The best results for the first step forecast are obtained for an autoregressive temporal order of 1, a spatial order of 2, a moving average temporal order of 1 and a spatial of 1. The two



Figure 4: One-step forecast validation results on a given VTL pair of the Paramics simulation network (penetration rate: 5%). (a) Travel time data of the VTL pair and its aggregate value on 5 minutes time intervals. Both the data and the aggregate value are shown for the whole data set and for a 5 % penetration rate. (b) One-step forecast of the congestion state produced by the logistic regression algorithm. The bars represent the probability of congestion estimated by the models for different levels of spatial dependency. The real state of congestion is represented with circles. (c) One-step forecast of travel time produced by the STARMA algorithm.



Figure 5: One-step forecast validation results for logistic regression on one VTL pair of the New York network. (a) Travel time data of the VTL pair and its aggregate value on 15 minutes time intervals. (b) One-step forecast of the congestion state produced by the logistic regression algorithm. The bars represent the probability of congestion estimated by the models for different levels of spatial dependency. The ground truth state of congestion is represented with circles. (c) One-step forecast of travel time produced by the STARMA algorithm.



Figure 6: Average one-step forecast error vs. penetration rate for all VTL pairs in the Paramics dataset. (a) Logistic Regression Forecast Classification Error. (b) STARMA Travel Time Forecast Error.

plots for which the moving average temporal and spatial orders are both equal to 1 414 show the best result for the first step forecast, but the error becomes quickly significant 415 when the forecast step increases. On the other hand, the two other plots for which the 416 moving average orders are one temporally and two spatially show a worse result for the 417 first step forecast but considerably better results for more than one step. The choice of 418 the parameters is therefore a very important step and should take into consideration 419 the performance of the forecasting for more than one step ahead. Analysis of a larger 420 data set is necessary to come to a statistically significant conclusion about the best 421 way to chose the spatio-temporal parameters for the STARMA model. 422

7 423

432

433

434

435

Conclusion

This article presented two statistical learning algorithms for estimating and forecast-424 ing arterial traffic conditions on a network. A first implementation inside the *Mobile* 425 Millennium system demonstrates both algorithms' ability to successfully forecast arte-426 rial travel times when sufficient training data is available. In summary, this work has 427 achieved the following: 428

- 1. It established the validity of a new data collection paradigm on arterial road-429 ways, namely the inter-VTL travel time data collection method for travel time 430 estimation and forecast at low penetration rates. 431
 - 2. It created data aggregation methods for capturing trends in arterial travel times (functions $h_i(\cdot)$ and $g_i(\cdot)$).
 - 3. It applied logistic regression and STARMA methods for learning spatio-temporal parameters used for estimating arterial link travel times.
- 4. It validated both models using a training/validation partition of the data, includ-436 ing a Paramics simulation data set and the results from three field tests in New 437 York City. 438



Figure 7: Forecast accuracy. (a) and (c): Forecast error for a VTL pair in the Paramics and New York networks, respectively. (b) and (d): Average forecast error as a function of the number of forecast steps into the future, Paramics and New York networks, respectively. One step is 5 minutes. In (b) and (d), t represents the temporal dependency and s represents the spatial dependency.

- 5. It analyzed the effect of penetration rate on forecast accuracy. 439 6. It analyzed the forecast horizon and its effect on the forecast accuracy. 440 7. It implemented real-time algorithms inside the *Mobile Millennium* system, public 441 displaying arterial traffic information. 442 This work is foundational to future research using statistical learning techniques 443 to forecast travel times in urban networks. Extensions to other statistical learning 444 methods beyond logistic regression and STARMA are needed to assess which technique 445 is most appropriate and efficient to the case of arterial travel times. Additionally, 446 there are a number of analyses that could extend the current work. In particular, the 447 following questions are open future research topics: 448 • Given the segment by segment travel time forecasts, how can accurate forecasts 449 of route travel times be determined? 450 • The current work requires a full training set on which to operate and needs an 451 aggregated data value (h_i function). How can the data requirements be relaxed 452 while maintaining high accuracy? The goal here is to fill in "missing" data using 453 knowledge of typical traffic patterns. 454 • How can the results from the specific examples here be generalized to all roads by 455 using common features such as speed limit, number of lanes, number of signals, 456 number of stop signs, etc.? The goal here is to be able to estimate spatio-temporal 457 model parameters in locations where no validation data yet exists. 458 Those questions are part of the ongoing work in *Mobile Millennium*. Current efforts 459 are focused on giving more accurate forecasts of arterial travel times on a network-wide 460 scale. 461
- Acknowledgements 462

The authors are grateful to Jeff Ban and Peng Hao for the production of the simulation data 463 used in the present article with Paramics software. They thank Steve Andrews for the logistics 464 of the field tests in New York City used to collect data used in the present article. They also 465 wish to thank Quinn Jacobson, Toch Iwuchukwu and Ken Tracton for their work in setting up 466 the operational system to collect data in New York City. 467

468

469

470

471

472

473

474

475

476

References

- [1] CITRIS, Center for Information Technology Research in the Interest of Society. http://www.citris-uc.org/.
 - [2] The Mobile Millennium Project. http://traffic.berkeley.edu.
 - [3] TTI, Texas Transportation Institute: Urban Mobility Information: 2007 Annual Urban Mobility Report. http://mobility.tamu.edu/ums/.
- [4] X. Ban, R. Herring, P. Hao, and A. Bayen. Delay pattern estimation for signalized intersections using sampled travel times. In Proceedings of the 88th Annual Meeting of the Transportation Research Board, Washington, D.C., January 2009.

[5] H. Bar-Gera. Evaluation of a cellular phone-based system for measurements of 477 traffic speeds and travel times: A case study from Israel. Transportation Research 478 Part C, 15(6):380–391, December 2007. 479 [6] B.S. Chen, S.C. Peng, and K.C. Wang. Traffic modeling, prediction, and conges-480 tion control for high-speed networks: a fuzzy AR approach. *IEEE Transactions* 481 on Fuzzy Systems, 8(5), May 2000. 482 [7] H. Chen, S. Grant-Muller, L. Mussone, and F. Montgomery. A study of hybrid 483 neural network approaches and the effects of missing data on traffic forecasting. 484 Neural Computing & Applications, 10(3), April 2001. 485 [8] G. A. Davis and N. L. Nihan. Nonparametric regression and Short-Term freeway 486 traffic forecasting. Journal of Transportation Engineering, 117(2):178–188, March 487 1991. 488 [9] N. Geroliminis and A. Skabardonis. Prediction of arrival profiles and queue lengths 489 along signalized arterials by using a Markov decision process. Transportation 490 Research Record, 1934(1):116–124, May 2006. 491 [10] M. Gonzalez, C. Hidalgo, and A. Barabasi. Understanding individual human 492 mobility patterns. *Nature*, (453):779–782, June 2008. 493 [11] J.C. Herrera, D. Work, J. Ban, R. Herring, Q. Jacobson, and A. Bayen. Evalua-494 tion of traffic data obtained via GPS-enabled mobile phones: the mobile century 495 experiment. Submitted to Transportation Research Part C, December 2008. 496 [12] B. Hoh, M. Gruteser, R. Herring, J. Ban, D. Work, J. C. Herrera, and A. Bayen. 497 Virtual trip lines for distributed privacy-preserving traffic monitoring. In The Sixth 498 Annual International conference on Mobile Systems, Applications and Services 499 (MobiSys 2008), Breckenridge, U.S.A., June 2008. 500 [13] B. Hoh, M. Gruteser, H. Xiong, and A. Alrabady. Enhancing security and privacy 501 in traffic-monitoring systems. *IEEE Pervasive Computing*, 5(4):38-46, March 502 2006.503 [14] M. Kamarianakis and P. Prastacos. Forecasting traffic flow conditions in an urban 504 network: Comparison of multivariate and univariate approaches. Transportation 505 Research Record, 1857:74–84, May 2004. 506 [15] M. J. Lighthill and G. B. Whitham. On kinematic waves. II. a theory of traffic 507 flow on long crowded roads. Proceedings of the Royal Society of London. Series 508 A, Mathematical and Physical Sciences, 229(1178):317–345, May 1955. 509 [16] H. Liu, A. Danczyk, R. Brewer, and R. Starr. Evaluation of cell phone traffic data 510 in minnesota. Transportation Research Record, 2086:1–7, December 2008. 511 [17] I. Nagy, M. Karny, P. Nedoma, and S. Varacova. Bayesian estimation of traffic lane 512 state. International Journal of Adaptive Control and Signal Processing, 17(1):51-513 65, November 2003. 514 [18] A. Ng. Lecture notes. CS 229: Machine learning. Stanford University, 2003. 515 [19] T. Park and S. Lee. A Bayesian approach for estimating link travel time on urban 516 arterial road network. In Computational Science and Its Applications ICCSA 517 2004, pages 1017–1025. Perugia, Italy, May 2004. 518

519 520	[20]	P.E. Pfeifer and S.J. Deutsch. A three-stage iterative procedure for space-time modeling. <i>Technometrics</i> , 22(1):35–47, February 1980.
521 522	[21]	P. Richards. Shock waves on the highway. Operations Research, $4(1):42-51$, February 1956.
523 524 525 526	[22]	A. Skabardonis and N. Geroliminis. Real-time estimation of travel times on sig- nalized arterials. In <i>Proceedings of the 16th International Symposium on Trans-</i> <i>portation and Traffic Theory</i> , University of Maryland, College Park, MD, July 2005.
527 528 529	[23]	A. Skabardonis and N. Geroliminis. Real-Time monitoring and control on signal- ized arterials. <i>Journal of Intelligent Transportation Systems</i> , 12(2):6474, March 2008.
530 531	[24]	B. L. Smith and Smart Travel Laboratory. <i>Cellphone probes as an ATMS tool.</i> Center for ITS Implementation Research, University of Virginia, VA, June 2003.
532 533 534	[25]	A. Stathopoulos and M. Karlaftis. A multivariate state space approach for urban traffic flow modeling and prediction. Transportation Research Part C, $11(2)$:121–135, April 2003.
535 536 537	[26]	S. Sun, C. Zhang, and G. Yu. A Bayesian network approach to traffic flow fore- casting. <i>IEEE Transactions on Intelligent Transportation Systems</i> , 7(1), March 2006.
538 539 540	[27]	J.W.C. Van Lint, S. P. Hoogendoorn, and H.J. Van Zuylen. Accurate freeway travel time prediction with state-space neural networks under missing data. <i>Transportation Research Part C</i> , 13(5-6):347–369, August 2005.
541 542 543	[28]	E. I. Vlahogianni, N. Geroliminis, and A. Skabardonis. On traffic flow regimes and transitions in signalized arterials. In <i>Proceedings of the 86th TRB Annual Meeting, January, Washington, D.C.</i> , January 2007.
544 545 546 547	[29]	D. Work, O.P. Tossavainen, S. Blandin, A. Bayen, T. Iwuchukwu, and K. Tracton. An ensemble Kalman filtering approach to highway traffic estimation using GPS enabled mobile devices. In <i>Proceedings of the 47th IEEE Conference on Decision and Control</i> , pages 5062–5068, Cancun, Mexico, December 2008.
548 549 550	[30]	JL. Ygnace. Travel time estimation on the san francisco bay area network using cellular phones as probes. <i>California PATH Program, Institute of Transportation Studies, University of California at Berkeley</i> , September 2000.
551 552 553 554	[31]	Y. Yim and R. Cayford. Investigation of vehicles as probes using global position- ing system and cellular phone tracking: field operational test. <i>California PATH</i> <i>Program, Institute of Transportation Studies, University of California at Berkeley</i> , February 2001.