

1 Traffic flow estimation using higher-order speed statistics

2 Submission date: 07/31/2012

3 Edouard Bulteau¹
4 Master's student, Systems Engineering
5 Department of Civil and Environmental Engineering
6 University of California, Berkeley
7 Berkeley, CA 94720, USA
8 edouard.bulteau@polytechnique.org

9 Romain Leblanc
10 Master's student, Systems Engineering
11 Department of Civil and Environmental Engineering
12 University of California, Berkeley
13 Berkeley, CA 94720, USA
14 romain.leblanc@polytechnique.org

15 Sebastien Blandin
16 Research Scientist
17 IBM Research Collaboratory, Singapore
18 sblandin@sg.ibm.com

19 Alex Bayen
20 Associate Professor, Systems Engineering
21 Department of Electrical Engineering and Computer Sciences
22 Department of Civil and Environmental Engineering
23 University of California, Berkeley
24 642 Sutardja Dai Hall, Berkeley, CA 94720-1764, USA
25 bayen@berkeley.edu

26 Word count:
27 Number of words (Body): 4,125
28 Number of words (Abstract): 190
29 Number of figures: 7 (250 each)
30 Number of tables: 3 (250 each)
31 Total: 6,815

¹Corresponding author

Abstract

In this article, we consider the problem of estimating traffic flow on a multi-lane road using a set of point speeds, either crowd-sourced or collected from the fixed infrastructure. We specifically investigate the relation between higher-order speed moments and the expected value of traffic flow. The algorithm proposed is based on the selection of optimal covariates constructed as speed moments, for a class of conditional mean predictors. The second contribution of this article consists in the analysis of specific components of the speed moments with significant correlation with flow values. In particular, we show that for more than 75% of the fixed sensing devices considered, the correlation coefficient between the inter-lanes speed variance and the aggregate flow is more than 0.75. Additionally, for more than 70% of these fixed sensing devices the lane speed variance increases with flow. The third contribution of this article consists of identifying the explanatory features for the high correlation between speed moments and flow values. The algorithms presented in this article are trained and tested on a large dataset from the Mobile Millennium system, collected in the Bay Area from August 2009 to October 2009.

1 Introduction

2 1.1 Motivation

3 Macroscopic traffic flow modeling is at the core of real-time traffic monitoring. The main assumption of
4 macroscopic traffic modeling is that vehicles traveling on a road behave similarly to a fluid in a pipe, hence
5 the state of traffic can be completely captured by three macroscopic variables: the *density* ρ of vehicles, the
6 *flow* q of traffic and the *space-mean speed* v of vehicles (Cassidy and Coifman (1), Knoop et al. (2)).

7 The *Lighthill-Whitham-Richards* (LWR) equation, expressing the conservation of mass, is classically used
8 to model traffic flow. By definition, ρ , q and v also satisfy the relation $q = \rho v$. However, a third relation
9 between these three quantities is needed in order to have a well-defined system.

10 The so-called *fundamental diagram* (Newell (3), Daganzo (4)) of traffic flow expresses an empirical relation
11 between density, flow and speed defined for any two couples of these three quantities. Lighthill and Whitham
12 (5) describe in their classical work the fundamental diagram as the fact that "at any point of the road the
13 flow (vehicles per hour) is a function of the concentration (vehicles per mile)". Historically, research work
14 has been focused on modeling the relation for the couple density-flow.

15 The recent explosion of mobile devices (Mobile Millennium (6)) measuring speeds requires innovative
16 approaches allowing the estimation of traffic flow from reported speeds only. This problem is a major open
17 problem in traffic flow modeling. For the triangular model, it is not possible to infer flow from speed as flow
18 is a multi-valued function of speed, in particular for high speeds.

19 Recent work from Blandin et al. (7) has shown a promising venue for estimating traffic flow from speed via
20 the computation of individual speed variances. The goal of this work is to investigate innovative approaches
21 to further push the research frontier in this topic. Using speed and flow measurements recorded by 116 fixed
22 radar sensors in the Bay Area, in this article we investigate the design of a speed-flow mapping based on
23 statistical analysis.

24 1.2 General background

25 Research on fundamental relations between macroscopic traffic quantities date back to the work of Green-
26 shields (8), who proposed a parabolic flow-density relation corresponding to a linear speed-density relation.
27 Over the past decades, traffic scientists developed a wide range of mathematical models for uninterrupted
28 traffic flow. In particular, the *triangular fundamental diagram* (Daganzo (4)) models the relation between
29 traffic flow and density by a triangular relation.

30 In the triangular fundamental diagram, the uncongested phase is characterized by an increasing linear
31 relation between flow and density. The congested phase is characterized by a decreasing affine relationship
32 between flow and density. In the uncongested phase, for this model, the speed is constant equal to the
33 free-flow speed, whereas in the congested phase, the speed is decreasing as the density increases.

34 In the uncongested phase, the flow is a multi-valued function of speed. This models the fact that up
35 to a certain density, the spacing between vehicles traveling on a roadway segment is sufficiently large for
36 drivers not to be affected by other vehicles. Thus, while the speed-to-flow conversion is straightforward in
37 the congested phase, this very conversion is impossible in the uncongested phase as speed is theoretically
38 constant at the free-flow speed for the whole phase.

39 1.3 Related work

40 While a lot of efforts focus on estimating vehicle density (Alvarez-Icaza et al. (9), Atrimy (10), Coifman
41 (11), Gazis and Knapp (12), Gazis and Szeto (13), Jerbi et al. (14), Panichpapiboon and Pattara-atikom
42 (15)), analyzing the relation between traffic flow and individual traffic speed is much less documented. The
43 analysis of the relation between traffic flow and individual traffic speeds and the ability to estimate flow
44 using a set of point speeds, either crowd-sourced or collected from the fixed infrastructure, currently
45 constitutes a major research question for traffic modeling theory.

1 Jørgensen (16) analyzes a large amount of combined speed and flow observations collected for the Ring
2 3 Motorway of the greater Copenhagen area over a 2 month period. The measurements are analyzed with
3 the aim of extracting information which can be used to realistically and reasonably model speed and flow
4 in a road traffic assignment model. It is observed that the US Bureau of Public Road (BPR) speed-flow
5 curve bpr (17) is an adequate average description of the non queuing conditions in static assignment models.

6 Blandin et al. (7) show that the conventional speed-flow mapping produces significantly more accurate
7 results than a speed variance-flow approach in the case of an increasing uncongested branch in (q, v) co-
8 ordinates. However, a classical assumption of the free-flow phase is a flat uncongested branch in (q, v)
9 coordinates. In that case, (7) show that the speed variance regression is more accurate than the classical
10 approach.

11 Wang et al. (18) focus essentially on the speed variance as a function of the density. They show that
12 the empirical variance takes a parabolic shape which first increases to a local maximum and then decreases
13 as traffic density increases. Although some flow estimation schemes are proposed, most of them are not
14 designed to be used in the uncongested phase.

15 The remainder of this article is organized as follows. Section 2 describes the dataset used in the subsequent
16 analysis. Section 3 presents the mathematical formulation of the estimation problem and introduces the
17 covariates considered. Numerical and experimental results for the flow estimation algorithm are provided in
18 Section 4. The features of the covariate that gives the best mapping with the flow are detailed in Section 5.
19 Finally, Section 6 concludes the paper.

20 2 Large scale Bay Area dataset

21 2.1 Sensing devices

22 The dataset used for the analysis presented in this paper consists of speeds and flows of vehicles traveling
23 on highway and freeway segments, as recorded by 116 fixed sensing devices r_k ($k = \{1, \dots, 116\}$) in the Bay
24 Area, California between August and October 2009.

25 2.2 Measurements

26 The sensing devices output speed (miles per hours) and flow (vehicles per minute) measurements per lane of
27 roadway at a sampling period $\Delta t = 1$ min. These measurements are made available in the *Mobile Millennium*
28 system hosted at the Center of California for Innovative Transportation. For every sampling time $m\Delta t$ (with
29 $m \in \mathbb{N}$ the time index), we have a direct access to $v_i(m\Delta t)$ the average speed on lane i during the time
30 interval $[(m-1)\Delta t, m\Delta t]$, with $n_i(m\Delta t) = q_i(m\Delta t)\Delta t$ the number of cars measured on lane i during the
31 previous period Δt . Additional static information at the location of the sensing device is also available, for
32 instance the number of lanes L and their characteristics (HOV or not), the direction the sensing device faces
33 (N, NE, E, SE, S, SW, W or NW) and the speed limit.

34 3 Mathematical formulation

35 3.1 Problem Statement

In this section we formulate the problem of finding a mapping between flow values and traffic variables
available in an infrastructure-free environment. For a given dataset \mathcal{D} of joint speed flow measurements,
we build a set $X(\mathcal{D})$ of traffic quantities derived from point speeds. The estimation problem, posed as an
optimization problem, reads as:

$$x_{opt} = \operatorname{argmin}_{x \in X(\mathcal{D})} \|\hat{q}(x) - q\| \quad (1)$$

36 where:

- 1 • q denotes the measured flow,
- 2 • $\hat{q}(\cdot)$ denotes a flow estimator,
- 3 • $\|\cdot\|$ denotes a norm chosen to measure the error.

4 3.2 Covariates considered

5 Finding a mapping between the flow and speed data for the uncongested state requires to make sure that
 6 the considered covariates are built using speed data from the uncongested state. The critical speed \bar{v}_c is
 7 defined as the speed at which the congested and uncongested branches of the speed/flow diagram intersect
 8 and partitions the set of speed data V into two subsets $C = [0, \bar{v}_c[$ and $U = [\bar{v}_c, +\infty[$ is the set of speed
 9 data in the free-flow state. The critical velocity \bar{v}_c is computed from observed data for each sensing device
 10 and the covariates considered in the remainder of this article are built at each time of interest for which the
 11 average speed belongs to U .

12 In this section we present the most relevant covariate of each category, whose performance in terms of
 13 flow estimation are compared in Section 4.2. The two covariates we present are the flow-weighted speed \bar{v}_q
 14 and the speed variance over time of the individual speeds Var_T^{is} .

15 The flow-weighted speed at time index m reads:

$$\bar{v}_q(m\Delta t) = \frac{\sum_{i=1}^L n_i v_i(m\Delta t)}{\sum_{i=1}^L n_i(m\Delta t)} \quad (2)$$

16 where:

- 17 • L is the number of lanes of the roadway segment at the sensing location,
- 18 • $v_i(m\Delta t)$ is the average speed on lane i at time index m ,
- 19 • $n_i(m\Delta t)$ is the number of cars on lane i at time index m .

20 The second covariate considered, suggested by Blandin et al. (7) is the variance over time of the aggregate
 21 speeds. It is the variance of the speeds measured on time interval. In this work, the size T of the time intervals
 22 is taken to be 11 min to be consistent with Blandin et al. (7). This choice is a trade-off between a short
 23 interval size required for consistency of the assumption of stationary traffic state, and a large number of
 24 speed observations required for accurate variance computation. In this article we focus our work on finding
 25 a mapping between the flow and variables built from individual speeds.

26 The variance over time of the individual speeds at time index m over a period T reads:

$$Var_T^{is}(m\Delta t) = \frac{1}{N-1} \sum_{i=1}^L \sum_{j=m-k}^{m+k} n_{ij} (v_{ij} - \bar{v}(m\Delta t))^2 \quad (3)$$

27 where:

- 28 • $[(m-k)\Delta t, (m+k)\Delta t]$ is the time interval of interest around the time index m . Note that $(2k+1)\Delta t =$
 29 T and for our analysis we take $k = 5$ to have $T = 11$ min.
- 30 • L is the number of lanes of the roadway segment at the sensing location,
- 31 • v_{ij} is the average speed of cars on lane i at time index j ,
- 32 • n_{ij} is the number of cars with speed v_{ij} ,
- 33 • $N = \sum_{i=1}^L \sum_{j=m-k}^{m+k} n_{ij}$ is the total number of cars observed in the time interval of interest,

- 1 • $\bar{v}(m\Delta t)$ is the flow weighted speed in the time interval of interest: $\frac{1}{N} \sum_{i=1}^L \sum_{j=m-k}^{m+k} n_{ij} v_{ij}$

The superscript *is* and subscript *T* are used to remind the reader that this variance aims to evaluate the variance over the period *T* of the individual speed (*is*). Note that this expression of the variance assumes that all the cars on lane *i* at time index *j* drive at the same speed v_{ij} , the average speed over this set of car. This assumption is discussed in 5.2.

This variance can be decomposed into components involving explicitly the variance over time of the individual speeds on each unique lane. Indeed for each lane *i*, we can introduce $\bar{v}_i(m\Delta t)$ the average speed on lane *i* during the time interval of interest and n_i the number of cars passing on lane *i* during the same interval:

$$\bar{v}_i(m\Delta t) = \frac{1}{n_i} \sum_{j=m-k}^{m+k} n_{ij} v_{ij}$$

$$n_i = \sum_{j=m-k}^{m+k} n_{ij}$$

- 2 Using (3) we can write $Var_T^{is}(m\Delta t)$ as:

$$\begin{aligned} Var_T^{is}(m\Delta t) &= \frac{1}{N-1} \sum_{i=1}^L \sum_{j=m-k}^{m+k} n_{ij} [(v_{ij} - \bar{v}_i(m\Delta t)) + (\bar{v}_i(m\Delta t) - \bar{v}(m\Delta t))]^2 \\ &= \frac{1}{N-1} \sum_{i=1}^L \left[\sum_{j=m-k}^{m+k} n_{ij} (v_{ij} - \bar{v}_i(m\Delta t))^2 + \sum_{j=m-k}^{m+k} n_{ij} (\bar{v}_i(m\Delta t) - \bar{v}(m\Delta t))^2 \right] \\ &\quad + \frac{2}{N-1} \sum_{i=1}^L [(\bar{v}_i(m\Delta t) - \bar{v}(m\Delta t)) \sum_{j=m-k}^{m+k} n_{ij} (v_{ij} - \bar{v}_i(m\Delta t))] \end{aligned} \quad (4)$$

We can introduce the individual speed variances over time for a single lane *i*:

$$Var_{T,i}^{is}(m\Delta t) = \frac{1}{n_i - 1} \sum_{j=m-k}^{m+k} n_{ij} (v_{ij} - \bar{v}_i(m\Delta t))^2$$

In (4) the first term can be expressed with the lane variances, and the last term can be simplified. The total variance can thus be written as:

$$Var_T^{is}(m\Delta t) = \frac{1}{N-1} \sum_{i=1}^L (n_i - 1) Var_{T,i}^{is}(m\Delta t) + \frac{1}{N-1} \sum_{i=1}^L n_i (\bar{v}_i(m\Delta t) - \bar{v}(m\Delta t))^2 \quad (5)$$

- 3 Equation (5) shows that the total individual speed variance over time is composed of the weighted sum of the
 4 variances on each lane, and of the weighted sum of the lanes average speed deviation from the overall speed
 5 average. This last component can be seen as a variance of the average speed over the lanes. The respective
 6 weights of each of these two components in the total individual speed variance is investigated in 5.2.

7 3.3 Flow estimation method

Our aim is to find a mapping between the variables described above and the flow, to then be able to produce flow estimations from raw speed data. The accuracy of the method used is determined using the root mean square error. For a given predictor variable *x* (in our case either flow-weighted speed or speed variance), we construct the flow predictor as:

$$\hat{q}(x_i) = E(q | [x] = x_i) \quad (6)$$

1 where $\lfloor x \rfloor$ is the floor of the variable x . The conditional mean estimator is justified by Sherman (19) in the
2 following lemma:

3 **Lemma 1** *If the probability density function associated with a variable q is symmetric around the mean and*
4 *unimodal then $E(q)$ is the optimal estimator of q in the sense of the quadratic error.*

In our case for each category x_i we assume that the probability density of q verifies these properties, hence
the conditional mean $E(q|x = x_i)$ is the optimal estimator of q when $x = x_i$. Using equation (6), equation (1)
becomes:

$$x_{opt} = \underset{x \in X(\mathcal{D})}{\operatorname{argmin}} \|E(q|\lfloor X \rfloor = x) - q\| \quad (7)$$

5 where:

- 6 • q denotes the measured flow,
- 7 • $\hat{q}(\cdot)$ denotes the estimator described in equation (6),
- 8 • $\|\cdot\|$ denotes a norm chosen to measure the error.

9 3.4 Algorithm

10 For each of the two covariates proposed we use the method illustrated in Figure 1 for assessing the estimator
11 performance for flow estimation. We use a cross-validation technique, where the cross-validation subsets
12 consist of the months of August 2009 and September 2009 - October 2009.

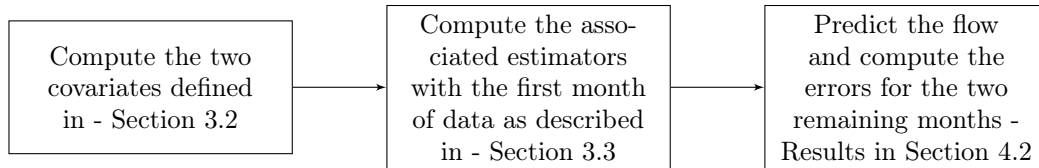


FIGURE 1 Methodology for best mapping selection

13 4 Flow estimation results

14 4.1 Properties of conditional flow distributions

15 In this section, we empirically validate the assumption made in the previous section on the symmetry around
16 the mean, and unimodality, of the conditional probability of the flow. Figure 2 represents the quartiles and
17 the mean of the conditional flow distribution as a function of the speed variance (Figure 2a) and as a
18 function of the flow-weighted speed (Figure 2b) for different sensing devices. For the two covariates the
19 symmetry of the quartiles around the median (which stays close to the mean) partially validates the use of
20 Lemma 1. In order to assess the unimodality of the conditional distributions, we propose to represent in
21 Figure 3 the conditional flow distributions and their mean for different values of speed variance (Figure 3a)
22 and flow-weighted speed (Figure 3b).

23 The majority of the 116 sensing devices exhibit unimodal conditional distributions with a good symmetry
24 around the mean in the case of the speed variance. The conditional distributions for the flow-weighted speed
25 are much less well-behaved (Figure 3b).

26

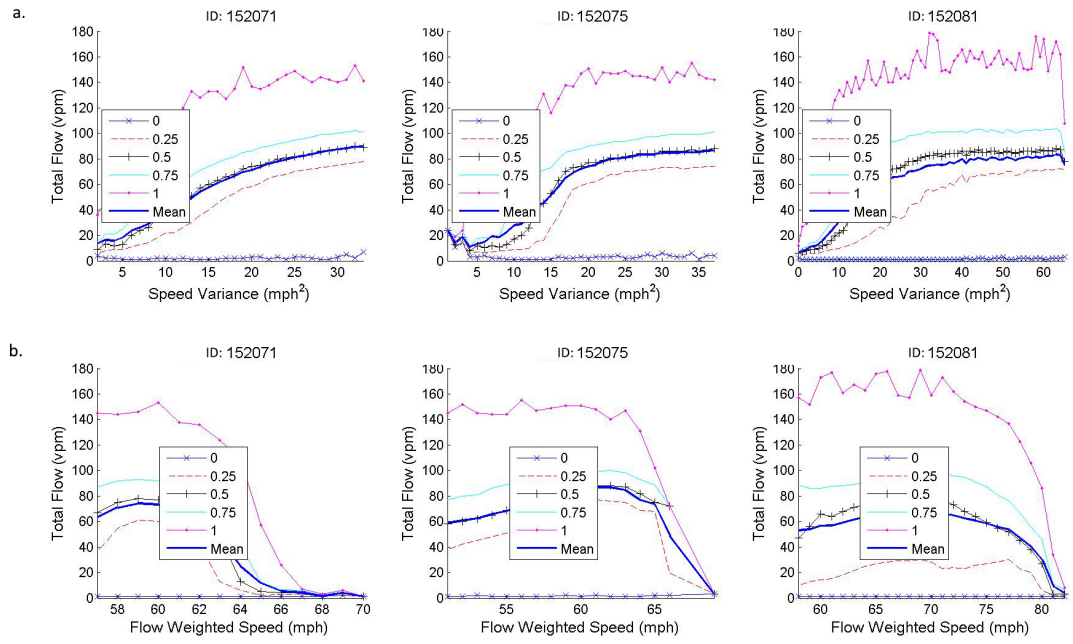


FIGURE 2 Quartiles and mean of the flow distribution conditioned by: a. The speed variance b. The flow-weighted speed. Results are given for three sensing devices (identifiers are displayed at the top of each graph)

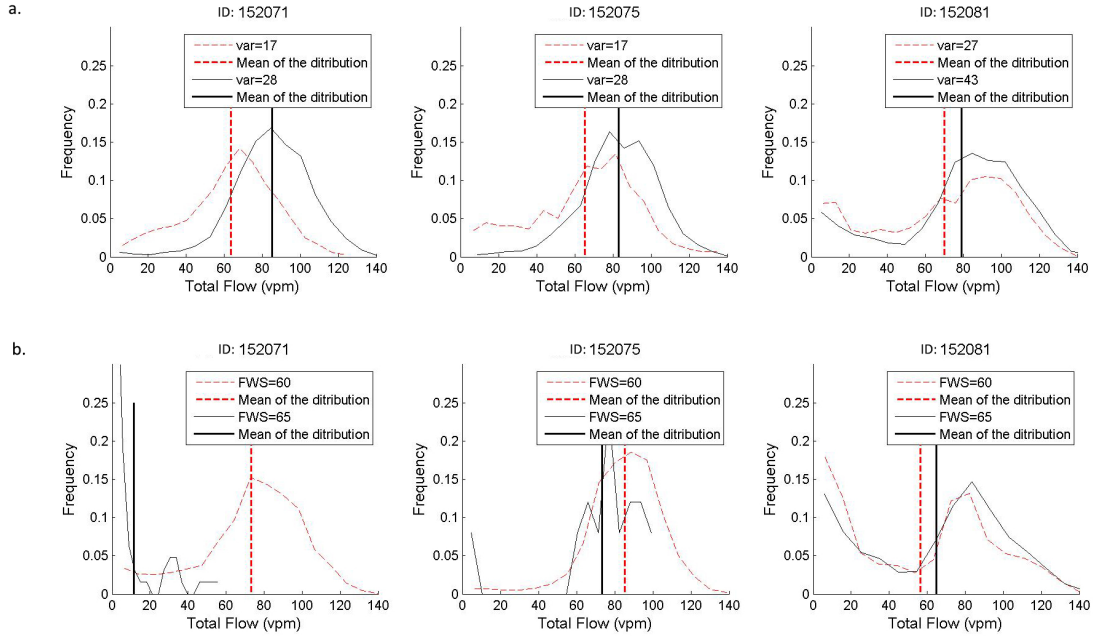


FIGURE 3 Flow distribution for given: a. Speed variances b. Flow-weighted speeds

1 These results show that:

- 2 • the first and third quartiles of the conditional flow distribution are closer to the mean and median in
 3 the case of the speed variance,
 4 • the distribution of the flow conditioned by the speed variance satisfy the assumption of unimodality and
 5 symmetry around the mean which does not seem to be always the case for the distribution conditioned
 6 by the flow-weighted speed.

7 In the following section we present the performance results for these two estimators.

8 4.2 Flow estimation errors

For each sensing device r the root mean squared error (20) of the estimator $\hat{q}(\cdot)$ is given by:

$$\varepsilon_r = \sqrt{\frac{\sum_{i=1}^N (\hat{q}(x_i) - q_i)^2}{N - 1}}$$

9 where:

- 10 • $\hat{q}(x_i)$ is the estimator of the flow at time of measurement i
 11 • q_i is the measured flow at time i
 12 • N is the total number of flow measurements in uncongested state

1 We represent in Figure 4 the cumulative distribution of the error over all sensing devices for the two methods.
2 The blue plain line corresponds to the speed variance method, whereas the red dashed line corresponds to
3 the flow weighted speed method. We notice that the plain line is always above the dashed line. It indicates
4 that for a given error ε the proportion of sensing devices which have an error inferior to ε is always higher
5 for the speed variance method than for the flow weighted speed method.

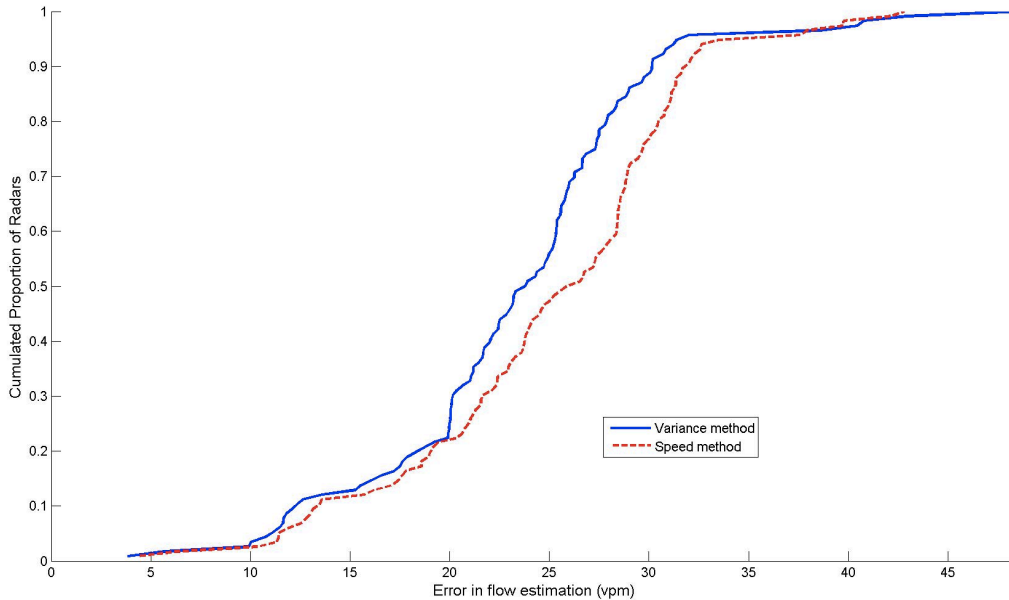


FIGURE 4 Cumulative distribution of the error

6 In the rest of the article, we investigate more thoroughly the relation between flow and speed variance.

7 **5 Features of the speed variance over time**

8 **5.1 Correlation with the flow**

9 In this section we analyze the relation between flow and speed variance by computing the distribution of the
10 correlations between these two quantities for each sensing device. Figure 5 represents the flow estimator as
11 a function of the speed variance. For all the presented sensing devices, a clear positive relation between the
12 two variables is displayed.

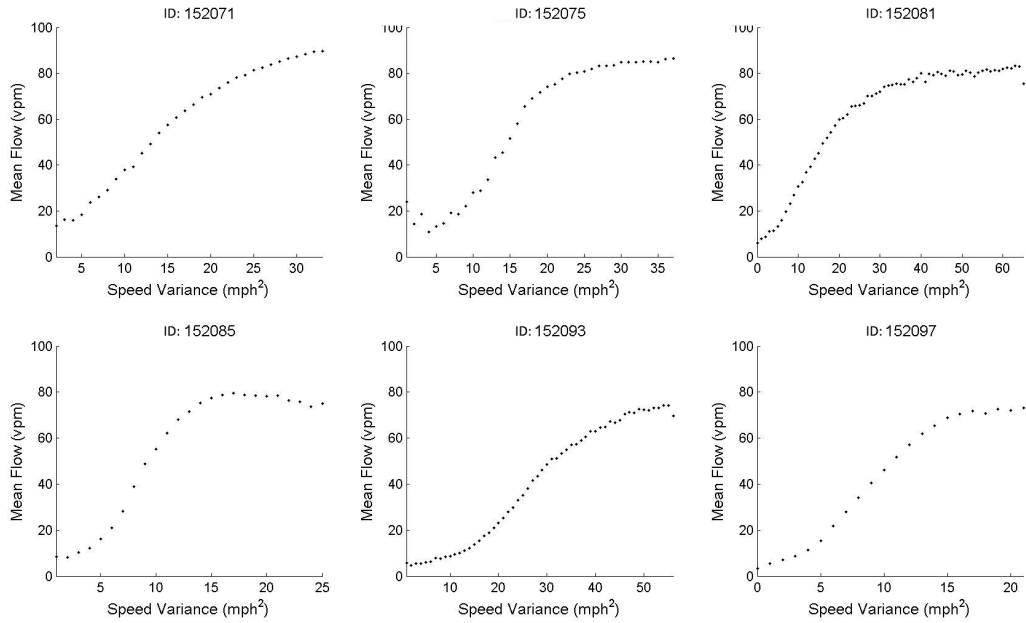


FIGURE 5 Estimator of the flow vs. speed variance over time

- 1 In order to get a better insight into the relation between speed variance and flow, we compute the Pearson's
- 2 correlation coefficients between the speed variance and the mean flow conditioned by the speed variance. We
- 3 represent the distribution of these coefficients over the 116 sensing devices in Figure 6.

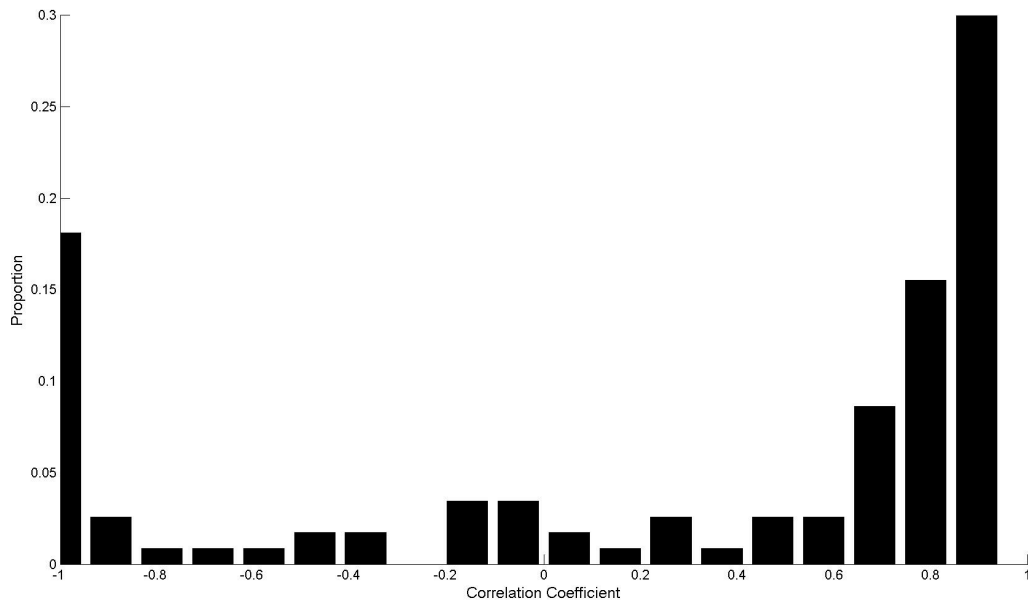


FIGURE 6 Distribution of the correlation coefficients

1 Figure 6 illustrates that for more than 75% of the sensing devices, the absolute correlation coefficient between
 2 the two described variables is more than 0.75. The histogram in Figure 6 illustrates both the strong relation
 3 between speed variance and the mean of the conditional distribution of the flow and the relevance of the
 4 mapping introduced in equation (1).

5 It is also important to notice that for most of these sensing devices (70%) individual speed variance
 6 increases with the flow. To understand this phenomenon and what could be the physical reasons of this
 7 relation, it is necessary to consider specific components of the speed variance.

8 5.2 Physical interpretation

9 We showed in (5) that the total variance results from the addition of two components: the first one is a
 10 weighted sum of the variances over time of the speed on each lane (we call it component 1) and the second
 11 one is a weighted variance over the lane of the lane average speed (called component 2). The first component
 12 quantifies the variability of the cars speed on the lanes over a given time interval. The second component
 13 represents how different lanes differ in terms of the average lane speed of the cars.

14 In order to explore the impact of each of these components in the total variance, we compute separately
 15 each of them and look at the ratios $r_1 = \text{component1}/\text{totalvariance}$ and $r_2 = \text{component2}/\text{totalvariance}$.
 16 We compute for each sensing device the mean of these ratios for all the uncongested states and analyze the
 17 distribution of these means over all sensing devices. We find that the median of the ratio r_2 is 60%. It shows
 18 that the second component explains slightly more of the total variance than the first component.

19 A more significant result describes the dependence of r_2 on the value of the total variance. Indeed Figure 7
 20 represents the mean of r_2 conditioned by the value of the total variance. We find that for all the sensing
 21 device, an increase of the total variance is associated with an increase of r_2 . This result seems to indicate
 22 that the increase of the total variance is mainly driven by an increase of the speed variance between the
 23 lanes. Thus the physical process which leads to an increase of the variance is less due to the increase of the
 24 speed variance within one same lane than to the difference of average speeds between the lanes.

1 The fact that the speed variance over time is not the main driver of the total variance justifies the
 2 assumption made in the definition of the total variance in equation (3) that along each time period Δt all
 3 the cars on a same lane have identical speed equal to their average speed.

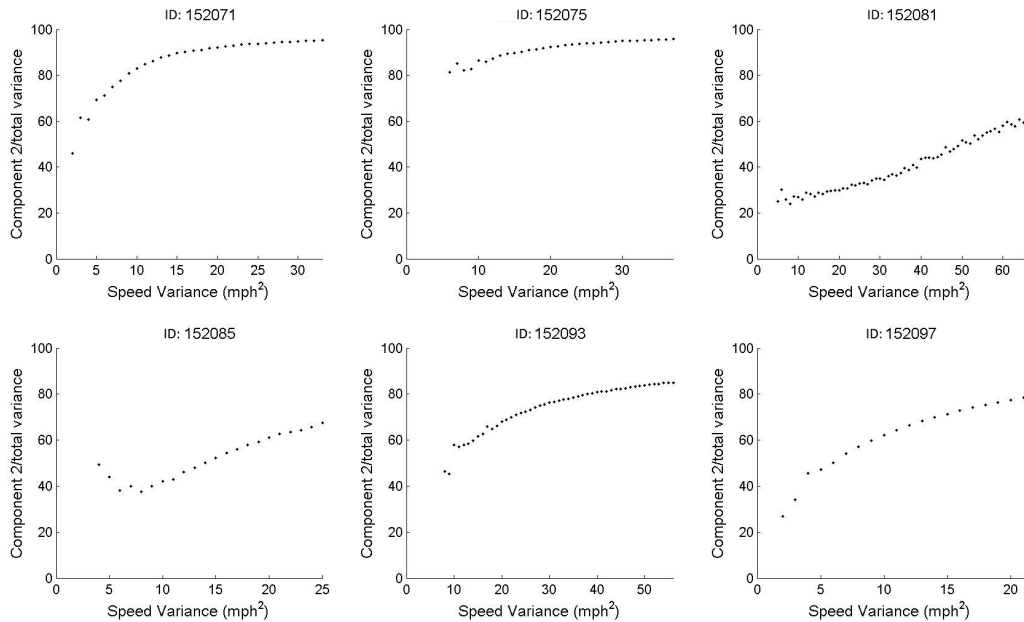


FIGURE 7 Weight of the inter-lane variance component in the total variance depending on the total variance

4 These results can be used to physically explain the phenomenon observed in the previous part (5.1)
 5 displaying for the majority of the sensing devices an increase of the total speed variance associated with an
 6 increase of the flow. Indeed an increase of the flow (note that we are still in free flow situation) could lead
 7 to a strong lane heterogeneity of the cars and thus to an increase of the speed variance due to an increase
 8 of the speed difference between the lanes. On the contrary, at low flow values, the distribution is fairly
 9 homogeneous between lanes because the drivers do not feel the need to stay into a specific lane and the
 10 overall variance is low.

11 Nevertheless it is important to note that for 30% of the sensing devices the relation between flow and
 12 speed variance is negative. In the following part we study how specific parameters impact the sign of the
 13 relation between flow and speed variance.

14 5.3 Analysis of factors influencing the correlation

15 Our goal is to determine if some features of the sensing device and its location could explain why some
 16 sensing devices show an increasing relation between flow and speed variance whereas other sensing devices
 17 exhibit a decreasing relation. Our analysis is based on two techniques: the Chi-test for independence is
 18 used to consider categorical parameters and a logistic regression model is employed to consider more type of
 19 parameters.

1 **5.3.1 Chi-test of independence and Cramer’s V**

2 Cramer’s V, denoted ϕ_c , represents the intercorrelation of two discrete variables (Cramér (21)). It varies
 3 from 0 to 1 and can reach 1 only when the two variables are equal. Our objective is to understand the
 4 relation between the fact that a sensing device displays a strong positive correlation (i.e. $r \geq 0.75$) between
 5 the speed variance and the flow, and some other characteristics of the sensing device and its location -such
 6 as direction, speed limit, number of lanes or type of sensing device (two different types of devices, based on
 7 two different measurement process, are used). The formula for the ϕ_c coefficient is:

$$\phi_c = \sqrt{\frac{\chi^2}{N(k-1)}} \tag{8}$$

8 where:

- 9 • χ^2 is derived from Pearson’s chi-squared test,
- 10 • N is the total number of observations,
- 11 • k is the number of rows or columns, whichever is less.

12 Below are the results of the Cramer’s coefficients and Chi-squared test of independence between the listed
 13 possible explanatory variables and the property of the sensing device of exhibiting an increasing or decreasing
 14 relation between flow and speed variance:

Explanatory variable	χ^2	ϕ_c	p-value
Direction	3.013	0.188	0.390 (*)
Speed limit	1.672	0.140	0.643 (*)
Sensing device type	48.749	0.757	<0.001
Number of lanes	5.496	0.254	0.358 (*)

TABLE 1 Cramer’s coefficients

15 Standards for interpreting Cramer’s V were proposed in (22):

$0.10 \leq \phi_c \leq 0.30$	Small effect
$0.30 \leq \phi_c \leq 0.50$	Medium effect
$\phi_c \geq 0.50$	Strong effect

TABLE 2 Standards for interpreting Cramer’s V as proposed by (22)

16 Using the above standard and the p-values of the chi-squared test for independence we notice that in our
 17 case only one variable seems to have an important impact on the sign of the relation between speed variance
 18 and flow. This variable is the sensing device type (p -value $\leq 10^{-4}$ and $\phi_c = 0.757$). The fact that a sensing
 19 device is of type 1 is a strong predictor of a negative relation between individual speed variance and flow.

20 **5.3.2 Logistic regression model**

21 The statistical analysis performed in 5.3.1 enables us to test independence between categorical variables.
 22 To complete our analysis and to take into account continuous and ordinal predictors, we perform a binary
 23 logistic regression (Hosmer and Lemeshow (23)). The output of this model is for a sensing device the sign
 24 of the relation between speed variance and flow, and the explanatory variables are the speed limit, the
 25 type of sensing device, the number of lanes and the distance from the closest on-ramp/off-ramp. Based
 26 on the maximum likelihood procedure, this model returns for one category of the output (‘positive relation

1 between speed variance and flow’) an intercept and regression coefficients for each predictor. For the second
 2 output category (‘negative relation between flow and speed variance’) the coefficients are taken to be zero.
 3 These coefficients can then be used to determine for a given set of explanatory variables observations, the
 4 probabilities of being in a positive or negative case. We present in Table 3 the results of this analysis
 5 conducted using all the sensing devices.

Variable	Coefficient	t-value	p-value
Intercept	4.702	0.533	0.594
Speed limit	0.020	0.142	0.887
Type	-1.163	-3.623	<0.001
Number of lanes	1.530	2.037	0.042
Distance from an exit	-0.833	-0.969	0.333

TABLE 3 Coefficients of the logistic regression model

6 These results confirm the significant impact of the sensing device type on the sign of the correlation between
 7 speed variance and flow. It also reveals the impact of the number of lanes on this outcome. Indeed the
 8 coefficient associated with the number of lanes has a low p-value ($p\text{-value} \leq 0.042$). Because this coefficient
 9 is positive, it implies that the higher the number of lanes is, the higher the probability for the sensing device
 10 of displaying a positive relation, is.

11 This result supports the interpretation proposed in 5.1: a higher number of lanes could lead to the
 12 categorization of the lanes into heterogeneous flow regimes. If we are in the case of a low number of lanes, a
 13 driver has to follow the flow and has limited freedom in the choice of his lane and of his speed: this increase
 14 of the flow leads to a decrease of the speed variance. Conversely, with a high number of lanes, a driver is able
 15 to choose a lane to keep a given speed. The difference between the lanes is more likely to increase compared
 16 to the case where we have a low number of lanes: an increase of the flow should lead to an increase of the
 17 average speed between the lane and thus drives the increase of the total speed variance.

18 6 Conclusion

19 In this article, we designed a novel algorithm for traffic flow estimation from fixed radars speed measure-
 20 ments. We showed that a fundamental relation exists between traffic flow and speed variance, which provides
 21 significant flow estimation improvements compared to standard regression techniques based on mean speed.

22 We showed that for most of the 116 sensing devices considered, the main contribution to the positive
 23 correlation between speed variance and flow is due to a difference of speed between lanes.

24 The analysis presented in this article illustrates that speed measurements from fixed radars data provide
 25 fundamentally new insights into the field of macroscopic traffic monitoring, in particular regarding novel
 26 properties of the relation between flow and speed on multi-lane highways.

27 Extensions to this work include systematic calibration of the parameters used for data processing in this
 28 work (such as the time interval used to compute the speed variance over time), and a more thorough analysis
 29 of the stationarity property of this higher-order fundamental diagram. Also more work should be done to
 30 assess the relevance of the described method when using probe data from mobile devices. In particular the
 31 impact of speed population variance on the flow estimation algorithm could be explored.

1 **Acknowledgements**

2 We want to thank the California Center for Innovative Transportation (CCIT) team for their help with
3 this project, in particular Scott Myers for helping with data extraction, but also Joe Butler, Ali Mortazavi,
4 Saneesh Apte and Jonathan Felder. Finally, we want to thank the California DOT for their ongoing support
5 to UC Berkeley and their interest in data fusion and hybrid data. Particular thanks go to John Wolf, Joan
6 Sollenberger and Nicholas Compin for their guidance.

1 References

- 2 [1] Cassidy, M. and B. Coifman, Relation among average speed, flow, and density and analogous rela-
3 tion between density and occupancy. *Transportation Research Record: Journal of the Transportation*
4 *Research Board*, Vol. 1591, 1997, pp. 1–6.
- 5 [2] Knoop, V., S. Hoogendoorn, and H. Zuylen, Empirical differences between time mean speed and space
6 mean speed. *Traffic and Granular Flow 07*, 2009, pp. 351–356.
- 7 [3] Newell, G., A simplified theory of kinematic waves in highway traffic, part I: General theory. *Trans-*
8 *portation Research Part B: Methodological*, Vol. Volume 27, No. 4, 1993, pp. 281–287.
- 9 [4] Daganzo, C., The cell-transmission model, part II: network traffic. *Transportation 21 Research Part B:*
10 *Methodological*, Vol. 29, No. 2, 1995, pp. 79–83.
- 11 [5] Lighthill, M. and G. Whitham, On Kinematic Waves. II. A Theory of Traffic Flow on Long Crowded
12 Roads. *Proceedings of the Royal Society A Mathematical Physical and Engineering Sciences*, Vol. 229,
13 No. 1178, 1955, pp. 317–345.
- 14 [6] Mobile Millennium, *Mobile Millennium*. <http://traffic.berkeley.edu>, 2012.
- 15 [7] Blandin, S., A. Salam, and A. Bayen, Individual speed variance in traffic flow: analysis of Bay Area
16 radar measurements. *91th Annual Meeting of the Transportation Research Board*, 2011.
- 17 [8] Greenshields, B. D., A study of traffic capacity. *Highway Research Board Proceedings*, Vol. 14, 1935, pp.
18 448–477.
- 19 [9] Alvarez-Icaza, L., L. Munoz, X. Sun, and R. Horowitz, Adaptive observer for traffic density estimation.
20 Boston, MA, 2004, pp. 2705–2710.
- 21 [10] Atrimy, M., Local density estimation and dynamic transmission-range assignment in vehicular ad hoc
22 networks. *IEEE Trans. Intell Transp. Syst.*, Vol. 8, No. 3, 2007, pp. 400–412.
- 23 [11] Coifman, B., Estimating density and lane inflow on a freeway segment. *Transp. Res. Part A*, Vol. 37,
24 No. 8, 2003, pp. 689–701.
- 25 [12] Gazis, D. and C. Knapp, On-line estimation of traffic densities from timeseries of flow and speed data.
26 *Transp. Sci.*, Vol. 5, No. 3, 1971, pp. 283–301.
- 27 [13] Gazis, D. and M. Szeto, Design of density-measuring systems for roadways. *Transp. Res. Rec.*, Vol. 495,
28 1974, pp. 44–52.
- 29 [14] Jerbi, M., S. Senouci, T. Rasheed, and Y. Ghamri-Doudane, An infrastructure-free traffic information
30 system for vehicular networks. Baltimore, MD, 2007, pp. 2086–2090.
- 31 [15] Panichpapiboon, S. and W. Pattara-atikom, Exploiting Wireless Communication in Vehicle Density
32 Estimation. *IEEE Transaction on vehicular technology*, Vol. 60, No. 6, 2011, pp. 2742–2751.
- 33 [16] Jørgensen, O. N. R., Estimation of speed–flow and flow–density relations on the motorway network in
34 the greater Copenhagen region. *IET Intelligent Transport Systems*, 2008.
- 35 [17] *Traffic assignment manual*. US Dept. of Commerce, Urban Planning Division, Washington, DC, 1964.
- 36 [18] Wang, H., D. Ni, Q.-Y. Chen, and J. Li, Stochastic Modeling of the Equilibrium Speed-Density Rela-
37 tionship. *J. Adv. Transp.*, 2011.
- 38 [19] Sherman, S., Non-mean-square error criteria. *IRE Transactions on Information Theory*, Vol. 4, No. 3,
39 1958, pp. 125–126.

- 1 [20] Wasserman, L., *All of Statistics: A Concise Course in Statistical Inference*. Springer Texts in Statistics,
2 Springer, 2004.
- 3 [21] Cramér, H., *Mathematical methods of statistics*. Princeton mathematical series, Princeton University
4 Press, 1946.
- 5 [22] Cohen, J., *Statistical power analysis for the behavioral sciences*. Hillsdale, NJ: Lawrence Erlbaum As-
6 sociates, 2nd ed., 1988.
- 7 [23] Hosmer, D. and S. Lemeshow, *Applied Logistic Regression*. Wiley Series in Probability and Statistics:
8 Texts and References Section, Wiley, 2000.