

ORIGINAL PAPER

Open Access

Modelling travel time uncertainty in urban networks based on floating taxi data



Dietmar Bauer^{1*} , Mirsad Tulic² and Wolfgang Scherrer³

Abstract

The prediction of the uncertainty of route travel time predictions for all possible routes in an urban road network is of importance for example for logistics. Such predictions need to take the essential features of the data set as well as the underlying traffic dynamics into account.

In this paper a large floating taxi data set is used in order to derive predictions of route travel time uncertainty based on link travel time uncertainty predictions. Prediction errors, that is actual travel times minus predicted travel times, are differentiated from model errors, that is measured travel times minus predicted travel times. These two errors are related, but not identical, as model errors contain measurement noise while the prediction errors do not. Detailed models for the variance of the link travel time prediction errors as well as the correlation between the model errors for different links are derived. The models are validated in depth using two different validation data sets.

Estimates for the variance of prediction errors are obtained. The standardized model error distributions show a remarkable stability, such that modelling the variance appears to be sufficient for quantifying the uncertainty of the model errors.

Furthermore we show that the model errors for adjacent links are highly correlated but correlations fade with increasing distance. Additionally usage of the road network plays a role with high correlation for links along common routes and low correlations for links along seldom used routes. We assume identical features for the prediction errors which is partly validated based on additional data.

The paper provides a way to estimate the complete distribution of route travel time prediction errors for any given route in the street network.

Keywords: Taxi floating car, Travel time uncertainty, Travel time prediction

1 Introduction

The quantification of route travel time uncertainty is of importance for logistics applications as well as for publicly available routing services (see [3] for a survey of studies dealing with valuing reliability; compare also [14]). Some authors even found that travel time reliability is valued more highly than travel time itself [8]. For logistic applications planning usually involves pre trip decisions, in many cases several hours or even days before the trips are executed. Both private individuals as well as logistics companies typically have asymmetric costs with being late implying higher penalties than being early. Accordingly for cost optimal decisions not only the predicted travel

time but also the uncertainty involved in the prediction is of interest.

Typically routing services are based on a directed mathematical graph (a set of nodes connected by links) representing the street network such that the shortest path between two points can be obtained using the Dijkstra algorithm. Here 'shortest' is to be understood in a broad sense and could involve predictions of link travel times for a particular departure time. For such predictions a huge amount of different methods based on a large number of different data sets have been obtained, first for highways (compare the papers contained in the compendium [2]), subsequently for general road networks. Excellent surveys of the many contributions can be found in [16, 17]. Usually in these approaches the predicted route travel time is obtained as the sum of the predicted link travel times.

*Correspondence: Dietmar.Bauer@uni-bielefeld.de

¹Bielefeld University, Universitätsstrasse 25, 33619 Bielefeld, Germany
Full list of author information is available at the end of the article

With respect to travel time reliability a number of different measures could be used, compare the survey in [6]. It is important to note that these measures not only depend on the respective link quantities but also on the relationship between the various link travel times. The most basic uncertainty measure is constituted by the variance. This measure for highways has been criticized as not telling the whole story [10]. Consequently approaches such as quantile regression for quantifying the whole distribution of travel time prediction errors have been developed (see the contributions [7, 11] and the references contained therein). Often confidence intervals for prediction uncertainty based on standard deviations (that is, square roots of variances) are constructed assuming Gaussian distribution of the errors. A more elaborate approach would imply a constant distribution scaled by standard deviations. Quantile regression methods replace this simple model by a detailed model for a number of quantiles depending on influencing factors. Such methods provide better quantifications of uncertainty compared to the scaling approach if the shape of the distribution changes a lot depending on influencing factors such as the time-of-the-day for example.

The variance of a sum of random variables (such as the sum of link travel times) equals the sum of the variances plus the sum of all possible covariances between pairs of random variables. In [9] the covariances of prediction errors are neglected and several different measures for route travel time variances are compared without reaching a compelling conclusion. It is clear that the omission of covariances is unjustified if the contribution of the covariances is substantial. It is unclear, however, if this is the case.

Travel time uncertainty is related to different levels of congestion. It may be argued that congestion is spread only along routes driven by many cars while crossing traffic might be unaffected by congestion in the orthogonal direction. This conjecture will be investigated in this paper by using a variable called *trip count ratio* indicating for each pair of links the proportion of trips along one link also traversing the second link.

As routing applications potentially build routes including all possible combinations of links, the estimation of the variance of an arbitrary route travel time prediction error requires the estimation of all covariances between the link travel time prediction errors for all pairs of links. For a large map this is impossible and hence a model is needed that provides an estimate of the correlation of the travel time prediction errors for any two given links based on some influential factors such as the distance as well as the trip count ratio.

An ideal source for modelling is constituted by taxi floating car data as a large fleet can cover the whole street network and be active throughout the day. Taxis

show the advantage – compared to other fleets – to be operated continuously. Consequently this paper investigates the properties of the uncertainties of link and route travel time prediction errors based on models developed for a large floating taxi data set (FCD) in Vienna [15]. In the paper we distinguish prediction errors, that is actual travel time minus predicted travel time, from model errors which additionally include measurement errors. Both errors depend on the traffic state and contain inter-driver and intra-driver variation. In the paper we discuss the relation between the various components of the two errors and their impact on the estimation in detail.

Note that the related paper [13] deals with a slightly different problem by assuming low covering of the floating car data. This is countered by imposing much structure (in the form of regression equations) on the relation between measured variables while our data set is large enough in order to achieve a good coverage of the network (see below). However, for a map with several thousands of links, estimation of the covariances of link travel time prediction errors for all pairs of links still is not feasible.

The main contribution of this paper thus is the thorough investigation of link and route travel time prediction and model errors. First, it is shown that for our data set the model error variances for the link travel time models show a strong dependence on the measurement conditions. In particular the number of single taxi observations for one link and one time interval as well as the current traffic conditions influence the variance of the model errors. Second, the distribution of the standardized model errors (such that the conditional mean is zero and the conditional variance equal to one) is remarkably stable such that quantifying the variance is sufficient in order to obtain the whole model error distribution. Third, the correlations of link travel time model errors for different pairs of links are investigated in detail, identifying two influencing factors: the driving distance between the two links and the trip count ratio. Fourth, all models are thoroughly validated by means of two separate validation data sets: the first consists of a second period of the floating taxi measurements which is not used for modelling. The second comprises an even tougher test by comparing the estimated route travel time uncertainty obtained from the models to the uncertainty of the measured route travel time based on trajectory data for single taxi observations for a collection of routes.

The paper is organized as follows: In the next section the data set is described while Section 3 provides the methodology for the estimation of the route travel time uncertainty. The empirical results are described in Section 4. Finally Section 5 concludes the paper.

2 Data set

In Vienna the movement of some taxis is observed using low frequent (with a sampling frequency of about 1 minute) GPS sensing since 2004 with a fleet of approximately 3500 taxis in total of which around 2000 are on the road at any one time. These raw data are used in different ways as discussed in the following two subsections.

2.1 Floating taxi data

The floating taxi dataset used in this paper covers the time period from July 1st 2008 until July 31st 2010, a total of 761 days. For each GPS observation information on the status of the taxis is available which allows to filter out only those movements that are made when carrying a passenger.

The raw data set is map matched (using trajectory to route map matching; a commercial map is used to encode the street network and hence the location and length of the links are given exogenously; note that the links in the map are directed such that two way segments of a street are represented by two links in opposite direction in the map) and interpolated between observations in order to obtain an estimated route as a continuous sequence of links in combination with an estimated entry and exit time for each link. Additionally the routes found are assessed and unreliable routes (implying very high speeds) are excluded from further analysis. Details on data collection and preprocessing can be found in [15].

The corresponding obtained route data is aggregated to obtain link specific average travel times within a given time interval by using arithmetic averaging. The time intervals have been chosen heuristically to equal 15 minutes leading to 96 time intervals per day.

This procedure results in the generation of two separate data sets:

1. One data set contains estimated taxi routes including the estimated timing of link entry and exit events providing direct measurements of route travel times.
2. The second one consists of average travel time measurements $\check{\Pi}_{d,i}^l$ for each link l and each fifteen minute interval i on any given day d .

2.2 Link travel time data sets

In this paper the averaged travel time measurements $\check{\Pi}_{d,i}^l$ for four different locations are used: (the location is presented in Fig. 1):¹

- (a) Hietzing (H): 191 links around the main arterial in the West of the city leading past the tourist attraction Schönbrunn castle.

- (b) Westbahnhof (WBH): 122 links in the area of the Westbahnhof rail station. This area is adjacent to the shopping street Mariahilferstrasse.
- (c) Ring (R): 79 links in the innermost city with lots of tourist attractions.
- (d) Südosttangente (SOT): 58 links on the largest (in terms of traffic) inner city highway including a number of feeder links.

The four sites are selected as a compromise of including many different traffic environments such as city highway, main arterial as well as inner city regions on the one hand and respecting time restrictions for analysis with the soft- and hardware tools available to the authors.

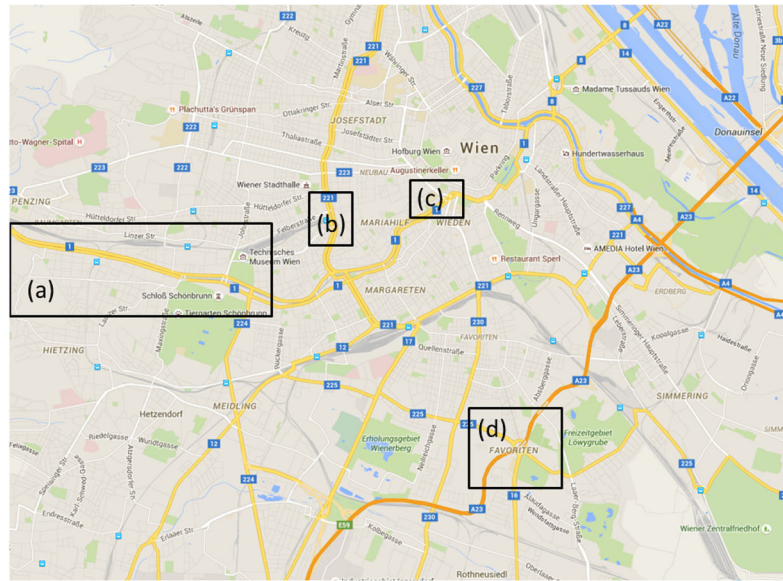
In total the dataset covers 38.9km of roads (approximately 19.4km in Hietzing, 6.1km at site WBH, 6.4km at R and 7.0km at SOT) containing approximately 32.9 million taxi observations. All four datasets include a varying number of missing observations in all dimensions. On 8.5% of days no observation exists at all due to coding errors either in the data collection or extraction from the database. The missing days occur on a continuous stretch of adjacent days, therefore there does not appear to be a systematic pattern of missing observation days. The fraction of missing observations per link (that is time intervals in which no taxi is observed on the corresponding link) varies from 10% to 80%. 19% of all measurements are based on only one taxi observation. For more details on the data set see [1].

Figure 2 provides information on the typical travel time measurements: (a) provides the empirical cumulative distribution function (ECDF) of all local travel time observations $\bar{\Pi}_{d,i}^l = \check{\Pi}_{d,i}^l / D_l$ (where D_l denotes the length in meter of the l -th link, d the day of observation, i the time-of-day-interval of the measured link travel time $\check{\Pi}_{d,i}^l$). Note that dealing with local travel time given in seconds per meter travelled helps in interpreting the results. As an example note that a travel speed of 50km/h corresponds to a local travel time of 0.072 seconds per meter.

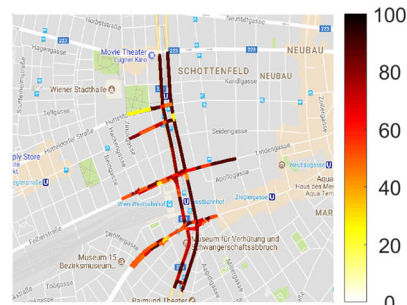
It can be seen that on the city highway Südosttangente in general smaller local travel times (corresponding to larger speeds) are observed while in the inner city larger local travel times are observed (R). Plot (b) provides a boxplot grouped across time-of-day-intervals for a link in WBH showing a number of characteristic features: Throughout the day congestion causes larger local travel times. Furthermore the standard deviation in general is large and varies a lot over the course of the day. During the afternoon peak period the standard deviation is a substantial fraction of the average local travel time.

In addition to $\check{\Pi}_{d,i}^l$ in the dataset also the empirical variance of local speed observations within one (link, day, time of day interval) combination before aggregation is provided for the Ring dataset. This information will be

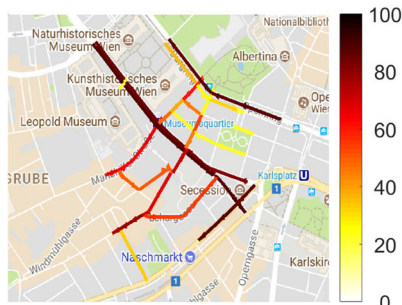
¹The same data set is also used in a companion paper [1] dealing with a different research question. Figure 1 is reprinted from this article.



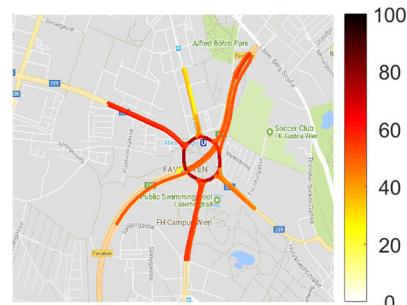
(a)



(b)



(c)



(d)

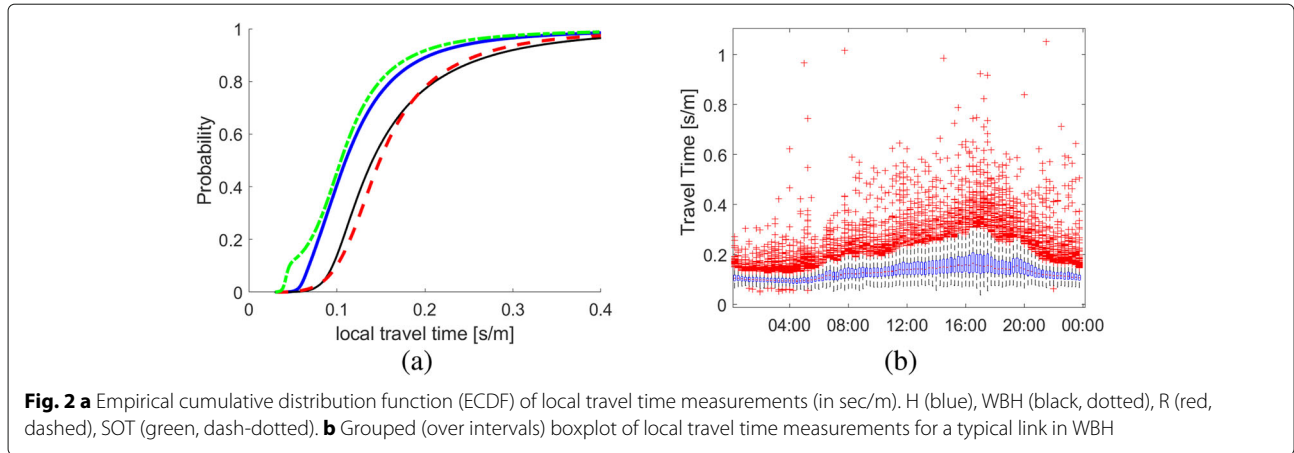
Fig. 1 Location of the test sites. **a** Hietzing, **b** Westbahnhof, **c** Ring, **d** Südosttangente. In the subfigures the links are colored according to percentage of available (non missing) observations (using the ‘hot’ colormap of MATLAB; white=0%, orange = 50%, black = 100%). Source: reprinted from [1]

used in order to estimate the measurement error variance in Section 3.2.

2.3 Link distance data

Driving distances (denoted with $D_{i,j}$) between the middle points of two links are obtained from the underlying

map. Intersection information such as turning restrictions are not used. Correspondingly the distance between two links is static over time. Two issues might arise in particular for links close to the boundary of the considered regions: As only the subgraphs in the four regions are used, there might exist shorter paths outside the



considered region connecting the two links. Second, the missing turning restrictions might reduce the driving distance. Both effects are considered minor problems for the chosen regions. In addition the number of links is relatively large such that such problems for a small number of pairs of links should be 'averaged out' for all results presented below.

2.4 Trip count ratios

In the analysis it will turn out that an influential factor for the size of spatial correlation of the link travel time model errors is given by a measure of connectedness between two links called `trip count ratio` (`tcr`) in the following. The calculation of `tcr` is based on trip data from 9 days² in 2012 (Sunday 1.1.-Wednesday 4.1., Sunday 8.1. - Tuesday 10.1., Wednesday 1.2., Tuesday 24.7.). Here the trip count ratio is defined as

$$\tau_{ij} = \frac{\max(\zeta_{ij}, \zeta_{ji})}{\zeta_i}$$

where ζ_i denotes the total number of trips via link i and ζ_{ij} the number of trips leading first via link i and then via link j for the whole nine days.

Note that this definition of the trip count ratio produces a symmetric measure in the sense that $\tau_{i,j} = \tau_{j,i}$. This measure will be used in order to model the correlation of model errors which inherently are symmetric. Typically high values of ζ_{ij} where link i lies upstream of link j imply low values of ζ_{ji} in the reverse direction. Alternatively in the model both the maximum and the minimum value of ζ_{ij} and ζ_{ji} could be used. This is left for future research.

Two instances of this measure for two links in the R and the H data set can be seen in Fig. 3. From these plots it is clearly visible that adjacent links on major routes reach τ values close to 1 while remote links show τ values of zero.

As detailed trip record data was not available for the same time span as the speed data, it is necessary to investigate the dependence of the τ values on the nine evaluation days. To this end a separate trip count ratio $\tau_{ij}^{(d)}$ is calculated for each of the 9 days. Figure 4 provides an ECDF for the absolute differences $|\tau_{ij}^{(d)} - \tau_{ij}|$ for all days and pairs of links. The WBH and H datasets show less variability with the 90% percentile of the deviation from the mean being equal to 0.05. For SOT (0.06) and R (0.08) we obtain a slightly larger – while still small – value.

Therefore using the overall trip count ratio τ_{ij} appears to be justified, where we have to be more careful with interpretation of the results for the R dataset.

3 Methods for uncertainty modelling

The main goal of this paper is to propose and validate a model for the distribution of the errors of the route travel time prediction along a given route $R = (L_j)_{j=1,\dots,J}$ (seen as a set of J links with indices L_j). The focus here is on long-term predictions, say at least one hour ahead, such that temporal correlation between deviations from 'usual' circumstances are no longer significantly different from zero. Therefore route travel time prediction and the corresponding uncertainty is modelled as a function of the time when embarking onto the route.

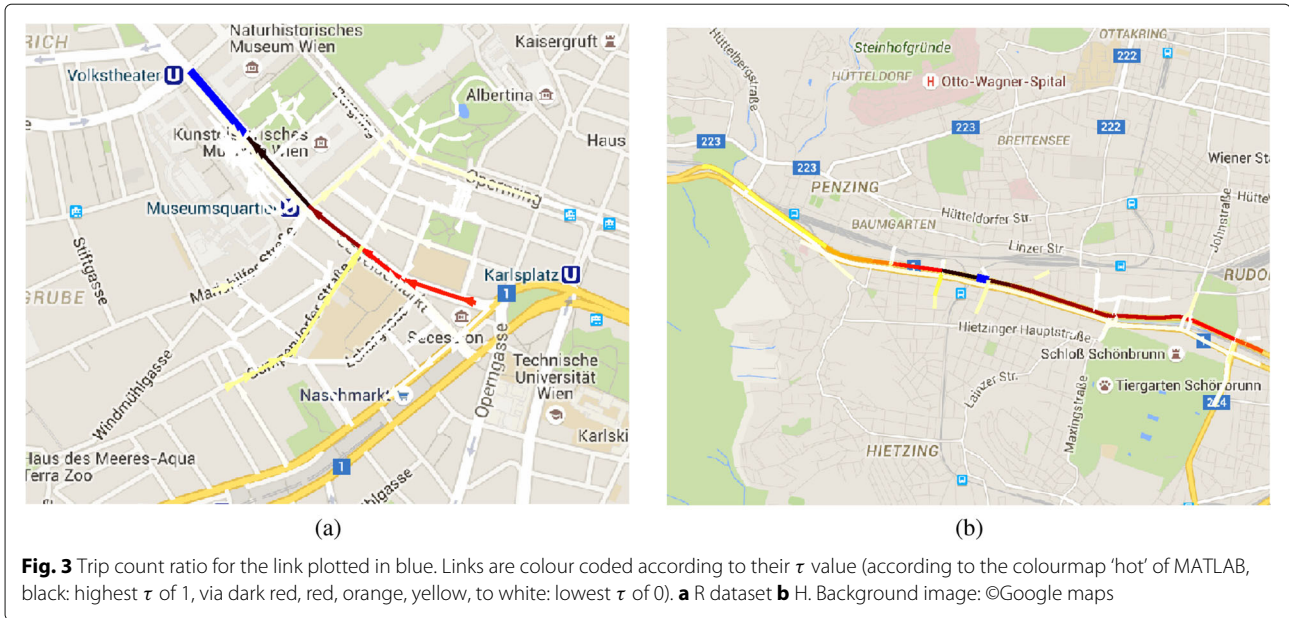
One of the difficulties related to our data set is that we don't have observations of actual link travel times $\Pi_{d,i}^l$ of a (single) taxi on given link l on day d and time-of-the-day interval i . We only have given an aggregate³, $\check{\Pi}_{d,i}^l$ say, of noisy measurements of (single) taxi travel times for a (random) number, $N_{d,i}^l$ say, of taxis.

As a prediction for the link travel times we simply use an estimate $\hat{\mu}_{d,i}^l$ for the (link and time dependent) expectation $\mu_{d,i}^l = E\check{\Pi}_{d,i}^l$. We call

$$\check{u}_{d,i}^l = \check{\Pi}_{d,i}^l - \hat{\mu}_{d,i}^l \quad (1)$$

²The selection of days was due to a lack of data availability.

³The aggregate is computed via the geometric mean of the respective link speeds.



the *model error* and

$$u_{d,i}^l = \Pi_{d,i}^l - \hat{\mu}_{d,i}^l \quad (2)$$

the *prediction error*. Note that $u_{d,i}^l$ and $\check{u}_{d,i}^l$ are closely related but they are not identical.

Based on the link travel time predictions the predicted route travel time on day d in time interval i equals the sum of the predicted link travel times:

$$\hat{\Pi}_{d,i}^R = \sum_{j=1}^J \hat{\mu}_{d,i}^{L_j}$$

The actual travel time will be denoted as $\Pi_{d,i}^R = \sum_{j=1}^J \Pi_{d,i}^{L_j}$ and the route travel time prediction error is $u_{d,i}^R = \Pi_{d,i}^R - \hat{\Pi}_{d,i}^R = \sum_{j=1}^J u_{d,i}^{L_j}$.

The variance of the route travel time prediction errors equals the sum of all variances and covariances of link travel time prediction errors:

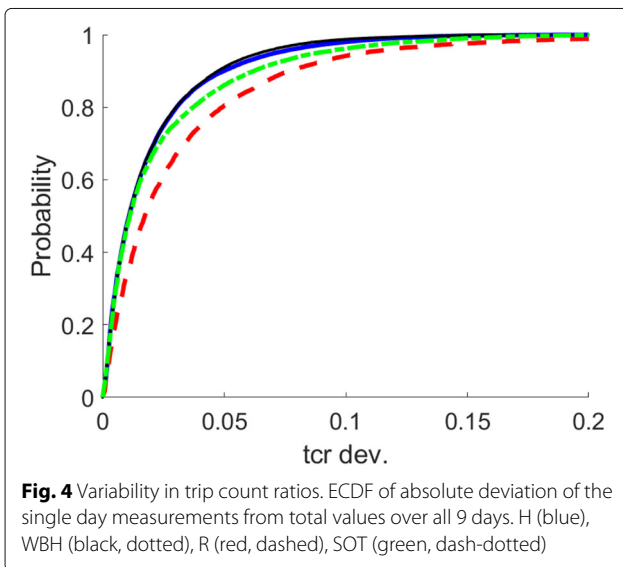
$$V(u_{d,i}^R) = \sum_{a=1}^J \sum_{b=1}^J \text{Cov}(u_{d,i}^{L_a}, u_{d,i}^{L_b}) \quad (3)$$

$$= \sum_{a=1}^J \sum_{b=1}^J \text{Corr}(u_{d,i}^{L_a}, u_{d,i}^{L_b}) \sqrt{V(u_{d,i}^{L_a})V(u_{d,i}^{L_b})}. \quad (4)$$

As has been noted above we do not have direct observations of the link travel times and hence it is not possible to directly estimate the variances and correlations of the link travel time prediction errors $u_{d,i}^l$. Instead we propose estimates of these quantities which are based on the model errors $\check{u}_{d,i}^l$.

The uncertainty in the measured travel time of an individual taxi include three components (cf. [12]):

- **inter driver variability:** under free flow conditions every driver sets his/her free speed which differs between drivers.
- **varying traffic conditions:** congestion is not identical on different days leading to deviations (that is random variables with zero mean) from the expected travel times. Additionally weather conditions and further noise factors might lead to deviations from usual traffic conditions.
- **measurement errors:** as link travel times are only measured based on low frequent GPS signals there is



a measurement error. We assume that these measurement errors are independent of the traffic state and of the time-of-day.

These three factors are all mixed in the observations. It is hard to separate them based only on the (aggregated) link travel time observations. Due to the aggregation the first and the third factor (inter driver variability and measurement errors) diminish with increasing number of observations ($N_{d,i}^l$) per time-of-day-interval while the second (varying traffic conditions) does not.

For the prediction uncertainty of a single (taxi-) driver travel time the first two components are relevant, while the third is not. Therefore in the following we will need to provide a detailed model for the variance of link travel time model errors $\check{u}_{d,i}^l$ as a function of time, traffic conditions as well as the number $N_{d,i}^l$ of observed taxis on day d in time interval i on link l :

$$V(\check{u}_{d,i}^l) = \check{\sigma}_{l,d,i}^2(\mu_{d,i}^l, N_{d,i}^l). \tag{5}$$

If ω_l denotes the variance of the measurement errors (which is assumed to be independent of the time) then the variance of the travel time prediction errors is

$$V(u_{d,i}^l) = V(\check{u}_{d,i}^l) - \omega_l = \check{\sigma}_{l,d,i}^2(\mu_{d,i}^l, 1) - \omega_l. \tag{6}$$

The correlations of the prediction errors are estimated via the correlations of the normalized (with the inverse of the standard deviation $\sqrt{V(\check{u}_{d,i}^l)}$) model errors. These are seen as proxies for the correlations of the prediction errors to which we do not have to access.

In the following we will explain these modelling steps in more detail:

- in Section 3.1 we discuss the modelling of the mean travel time $\mu_{d,i}^l$.
- in Section 3.2 we present a model for the variances of the measured link travel times and discuss how to construct estimates for the variance of the link travel time predictions from this model and suitable estimates of the measurement error variance.
- in Section 3.3 a model for the (spatial) correlations $\text{Corr}(\check{u}_{d,i}^a, \check{u}_{d,i}^b)$ for all pairs of links a, b is presented.
- finally Section 3.4 shows how these pieces are put together to get an estimate of the variance of the route travel time prediction error.

For each model we will investigate the dependence on links, days and time-of-day-intervals.

3.1 Model for the expected link travel time

Following [15] we will use the model

$$\begin{aligned} \check{\Pi}_{d,i}^l &= \mu_{d,i}^l + e_{d,i}^l, \\ \mu_{d,i}^l &= x_d^l \beta_{l,i} \end{aligned}$$

for the measured link travel time $\check{\Pi}_{d,i}^l$, where the regressor vector x_d contains the constant, dummies for the day category, school holidays and additional cyclical terms to model potential yearly effects ($\cos(\omega_j d), \sin(\omega_j d), \omega_j = 2\pi j/365, j = 1, \dots, 5$).

The regression coefficients $\beta_{l,i}$ are specific to the time-of-day-interval i and the link l .

According to [15] the models are estimated using stepwise regression techniques and excessive model selection in order to identify the most relevant regressors. A model for the variance of $e_{d,i}^l$ (see the next section) as a function of the underlying number of observations as well as the average mean speed reduces the influence of heteroskedasticity. For details see [15].

As a result we obtain estimates $\hat{\mu}_{d,i}^l$ which serve as predictions for the actual link travel times as explained above.

3.2 Model for the variance of link travel times

We will use estimates for the variance of $e_{d,i}^l = \check{\Pi}_{d,i}^l - \mu_{d,i}^l$ as estimates for the variance of the model errors $\check{u}_{d,i}^l = \check{\Pi}_{d,i}^l - \hat{\mu}_{d,i}^l = e_{d,i}^l + (\mu_{d,i}^l - \hat{\mu}_{d,i}^l)$. This simplification is justified since, due to the large data set used for the estimation of $\mu_{d,i}^l$, the estimation error $(\hat{\mu}_{d,i}^l - \mu_{d,i}^l)$ is "small" compared to the noise $e_{d,i}^l$.

The link travel time measurements show a pronounced heteroskedasticity, that depends on the number of measurements as well as the mean travel time. Again following [15] we model the variance of $e_{d,i}^l$ as a function of the number of observations as well as the predicted travel time:

$$\begin{aligned} V(e_{d,i}^l) &= \sigma_{l,d,i}^2(\mu_{d,i}^l, N_{d,i}^l) \\ &= \exp\left(\alpha_{l,i} + \gamma_{l,i} \mu_{d,i}^l + \phi_{l,i} (N_{d,i}^l)^{-1/2} + \delta_{l,i} \bar{N}_{d,i}^l\right) \end{aligned}$$

where $\bar{N}_{d,i}^l$ denotes the dummy variable indicating that the corresponding measurement is only based on one taxi observation. This makes the model for one taxi measurement insensitive to misspecifications of the dependence on $N_{d,i}^l$.

As in [15] this can be estimated in logarithms using $\log((\check{\Pi}_{d,i}^l - \hat{\mu}_{d,i}^l)^2)$ as the dependent variable. Here the coefficients $\delta_{l,i} \geq 0, \phi_{l,i} \geq 0$ are restricted to be positive, since we expect that averaging individual taxi observations decreases the variance.

Additionally a penalization is introduced in order to obtain smooth (over time-of-day-intervals) variation of coefficients. For details on the penalization approach used see [5]. In the following let $\hat{\sigma}_{l,d,i}^2$ denote the estimate for $\sigma_{l,d,i}^2$.

Note that this variance contains all three components of the link travel time uncertainty. The inter driver variability

and the varying traffic conditions act as influencing factors. The uncertainty related to a single trip from a single driver is obtained by setting $N_{d,i}^l = 1$.

With regard to the third component, the variance ω_l of the travel time measurement error for link l is assumed to be independent of time while the other components of the variance of the travel times vary with the traffic state: in conditions of synchronized traffic, inter driver variability vanishes. Varying traffic conditions lead to varying levels of the variance of travel time measurements. Therefore the measurement error variance can be bounded by the minimum of all observed variances. Assuming that in all cases the long data set contains worst case scenarios we will use the minimum of all observed variances as the measurement error variance.

In this respect note that for the Ring dataset the empirical variances $V_{d,i}^l$ of the speed measurements for each (link, day, time-of-day) observation is contained in the data set available for this study. For each link these empirical variances are modelled as

$$V_{d,i}^l = \eta_i^l + v_{d,i}^l$$

where $v_{d,i}^l$ denotes the error terms deviating from the expectation η_i^l not depending on the day of measurement. Estimating this model using penalization in order to obtain smooth curves over time-of-the-day

$$\hat{\omega}_l^V = \min_i \hat{\eta}_i^l$$

estimates the measurement error variance for the speed measurements specific to link l .

In order to transfer this result obtained for the Ring dataset to all datasets a nonlinear regression is used for explaining the measurement error variance $\hat{\omega}_l^V$ for all links l using regressors built from the following variables: (i) typical speed k_l according to the map,⁴ (ii) the empirical variance \hat{v}_l^2 of the aggregated speed measurements over all days and time-of-day intervals, (iii) the variance \hat{n}_l^2 of the number of observations and (iv) the variance \hat{p}_l^2 of the predicted speeds. With these variables the following model is estimated:

$$\hat{\omega}_l^V = \beta_0 + \beta_v \hat{v}_l + \beta_k \frac{1}{k_l} + \beta_2 \hat{v}_l^2 + \beta_3 \hat{v}_l \hat{p}_l + \beta_p \hat{p}_l + \beta_n \frac{1}{\hat{n}_l} + v_l$$

The corresponding R^2 equals 0.88 indicating a very good fit.

In the final step the Delta method is used to transfer the estimated variances $\hat{\omega}_l^V$ for the speed measurements to the corresponding variance $\hat{\omega}_l$ of the travel time measurements.

Therefore the estimates of the link travel time prediction error variances are obtained as

$$\hat{V}(u_{d,i}^l) = \hat{\sigma}_{l,d,i}^2 (\hat{\mu}_{d,i}^l, 1) - \hat{\omega}_l. \tag{7}$$

We note that this estimate uses a number of assumptions that are not obvious. Therefore a detailed verification of the assumptions using thorough validation on different data sets will be presented below.

3.3 Modelling spatial correlations

The pairwise correlation ρ_{ij} between the model errors $\check{u}_{d,i}^l$ of the link travel time for two links i, j is modelled as a function of the driving distance D_{ij} between the two links as well as the trip count ratio τ_{ij} .

The driving distance depends on the sequence of the two links (as is immediate from two adjacent links, where one follows the other in the direction of driving while getting from the downstream link to the upstream link requires either two U-turns or driving around a block). The trip count ratio has been defined already taking a symmetrization into account such that $\tau_{ij} = \tau_{ji}$. This is necessary since we are modelling correlations ρ_{ij} which by definition are symmetric such that $\rho_{ij} = \rho_{ji}$. Therefore also for the distance a symmetrization is needed:

$$d_{i,j} = \min(D_{i,j}, D_{j,i}).$$

Alternative specifications (such as using the average) have been tested but did not result in superior models.

As the dependent variable the empirical correlation of the normalized model errors is considered accounting for the observed heteroskedasticities:

$$\rho_{ij} = f(d_{i,j}, \tau_{ij}) + u_{i,j} = \check{\rho}_{ij} + u_{i,j}. \tag{8}$$

Note that we model the correlation of the normalized model errors, whereas for the estimation of the uncertainty of the route travel times the correlations of the normalized prediction errors would be needed.

Since we do not have access to the latter, we assume here that the two correlations show similar features such that the correlations of model errors are indicative of the correlations of the prediction errors. Limitations in our data set do not allow a detailed investigation.

3.4 Estimation of the route travel time variance

The travel time variance is estimated as

$$\hat{V}(u_{d,i}^R) = \sum_{a=1}^J \sum_{b=1}^J \text{Corr}(\check{u}_{d,i}^{L_a}, \check{u}_{d,i}^{L_b}) \sqrt{\hat{V}(u_{d,i}^{L_a}) \hat{V}(u_{d,i}^{L_b})} \tag{9}$$

with $\hat{V}(u_{d,i}^{L_a})$ computed according to (7).

With respect to the estimation of the correlations five different versions are tested:

⁴The Teleatlas map used contains a classification of roads into classes characterized by typical speed levels.

- NC: no correlation, $\text{Corr}(\hat{u}_{d,i}^{L_a}, \hat{u}_{d,i}^{L_b}) = 0$ for $L_a \neq L_b$. It is expected that this underestimates variability by neglecting typically positive correlations.
- EM: empirical correlation matrix including all pairwise estimates. It is to be expected that the corresponding estimates are noisy in particular for pairs of links with only a few joint observations.
- ET: using the information in Fig. 9 below empirical correlations are used for all pairs of links with distance smaller than 1km and correlations are set to zero for higher distances.
- MI: the correlations are estimated with an individual model according to Eq. 8 for each of the four datasets
- MJ: correlations are predicted from a model according to Eq. 8, estimated based on all datasets jointly
- SD: at the opposite end of the spectrum lies the case of perfect correlation: $\text{Corr}(\hat{u}_{d,i}^{L_a}, \hat{u}_{d,i}^{L_b}) = 1$ for all L_a, L_b which provides an upper bound of uncertainty.

3.5 Validation procedures

All components of the model are validated carefully in the most appropriate context. The predictions for the observed (aggregated) link travel times as well as the corresponding variance estimates $\hat{\sigma}_{l,d,i}^2(\hat{\mu}_{d,i}^l, N_{d,i}^l)$ are evaluated on a validation data by splitting the dataset into the first 701 days as the estimation data and the last 60 days as the validation data.

Secondly, we achieve a detailed validation of the estimation procedure from comparing the route travel time predictions as well as the variance estimates $\hat{V}(u_{d,i}^R)$ to estimation of single trip observations based on the additional trip data.

4 Results and discussion

4.1 Mean and variance model

Figure 5 presents the mean (over links) root mean square errors (in sample, over all observed days) of (scaled) model errors $(\check{\mu}_{d,i}^l - \hat{\mu}_{d,i}^l)/D_l$ for the four datasets over all time intervals. The strong dependence on the time of the day is clearly visible with RMSE varying between 0.04 s/m and 0.11 s/m.

Second, the models for the variances as given in (7) are estimated and specified. The results can be seen in Fig. 6 that provides the average of estimated standard deviation $\hat{\sigma}_{l,d,i}(\hat{\mu}_{d,i}^l, N_{d,i}^l)/D_l$ over days and links in the four datasets. In order to put the scale into perspective note that 50 km/h correspond to 0.072 seconds per meter. It can be seen that for the urban highway the estimated uncertainty is high in the two rush hour periods while it is comparatively low during the day. For the highway also an increase of uncertainty during night-time is visible due to a drop in the number of observations.

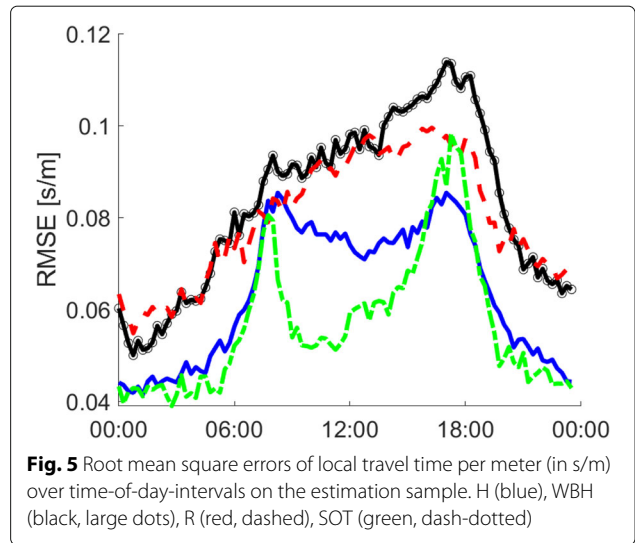


Fig. 5 Root mean square errors of local travel time per meter (in s/m) over time-of-day-intervals on the estimation sample. H (blue), WBH (black, large dots), R (red, dashed), SOT (green, dash-dotted)

Figure 7 demonstrates that the model for the variance is appropriate for all datasets. Plot 7(a) provides the boxplot of percentile values (computed over days and time of the days intervals) for normalized model errors $\hat{u}_{d,i}^l/\hat{\sigma}_{l,d,i}(\hat{\mu}_{d,i}^l, N_{d,i}^l)$ for all links on the four test sites for the estimation data (left columns) and the validation data (right columns). In all cases a skew distribution is visible. Except for the H data, estimation and validation data are in reasonable agreement with more variability over links on the shorter validation data as expected. For H a larger difference occurs due to a few outlying observations being present in the validation data.

Plot (b) of Fig. 7 provides the (2.5%, 50%, 97.5%) percentiles (computed over days and links) grouped into time-of-day-intervals (validation data is plotted in bold, estimation data in thin lines; the results on the validation data set being almost identical to the ones on the

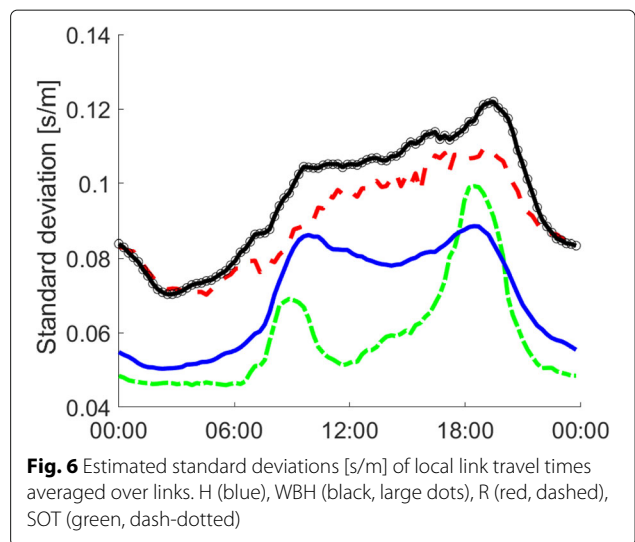
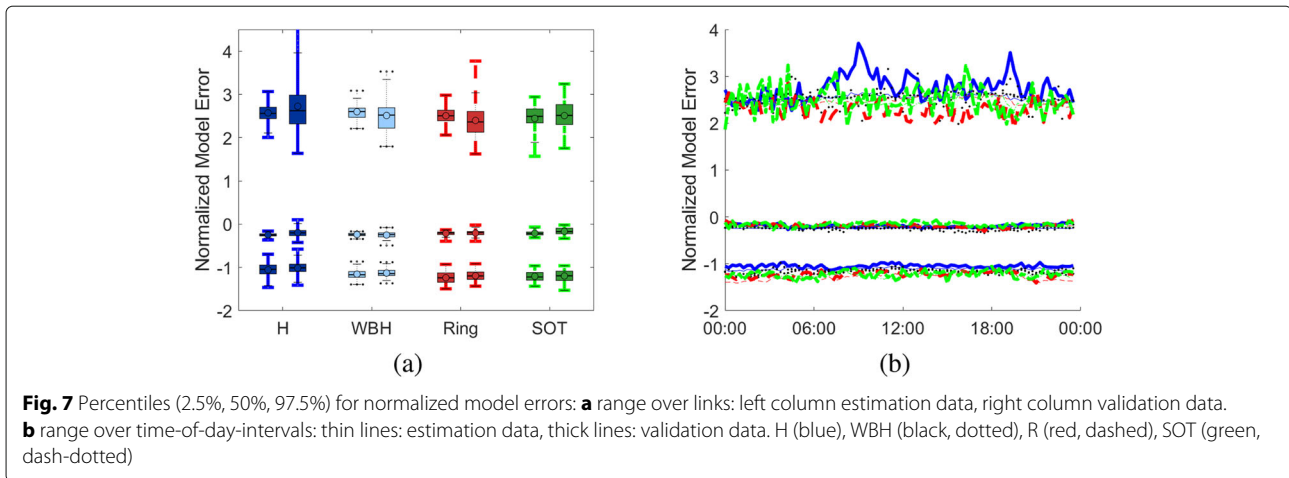


Fig. 6 Estimated standard deviations [s/m] of local link travel times averaged over links. H (blue), WBH (black, large dots), R (red, dashed), SOT (green, dash-dotted)



estimation data which are hence almost invisible). Only the H data shows deviations over time-of-day-intervals for the morning and the evening peak in the validation data set. The three percentiles are located at approximately -1 (2.5%), 0 (50%) and 2.5 (97.5%). The location of the given percentiles of the normalized residuals is very stable across different datasets, links and time-of-day-intervals. This indicates that the distribution of the normalized model errors is identical in all cases such that the whole distribution can be characterized by the normalized distribution times the scaling using the estimated standard deviation.

Note, however, that the prediction errors equal model errors minus measurement errors. Since measurement errors cannot be measured directly, it is unclear whether prediction errors also equal standard deviation times a random variable with distribution not depending on the factors influencing the standard deviation. This is left for future research.

4.2 Spatial correlation models

In this section the spatial correlation of the normalized model errors

$$\tilde{u}_{d,i}^l = \check{u}_{d,i}^l / \hat{\sigma}_{l,d,i}(\hat{\Pi}_{d,i}, N_{d,i}^l)$$

is investigated. Empirical correlations for all link pairs for each of the 95 time-of-day-intervals are calculated on the estimation data set. In this way a time series of 95 observations of the correlation matrices is obtained which is subjected to a Giraitis, Kokoszka, Leipus and Teyssièrè test [4] for time constancy. In the majority of the cases the test does not find evidence (at significance level $\alpha = 0.01$) for variation over time: for the H data set evidence is found only in 15%, 51% for R, 36% for SOT and 43% for WBH. Thus for the inner city settings there is some evidence of changes in correlation structure for a sizeable fraction of

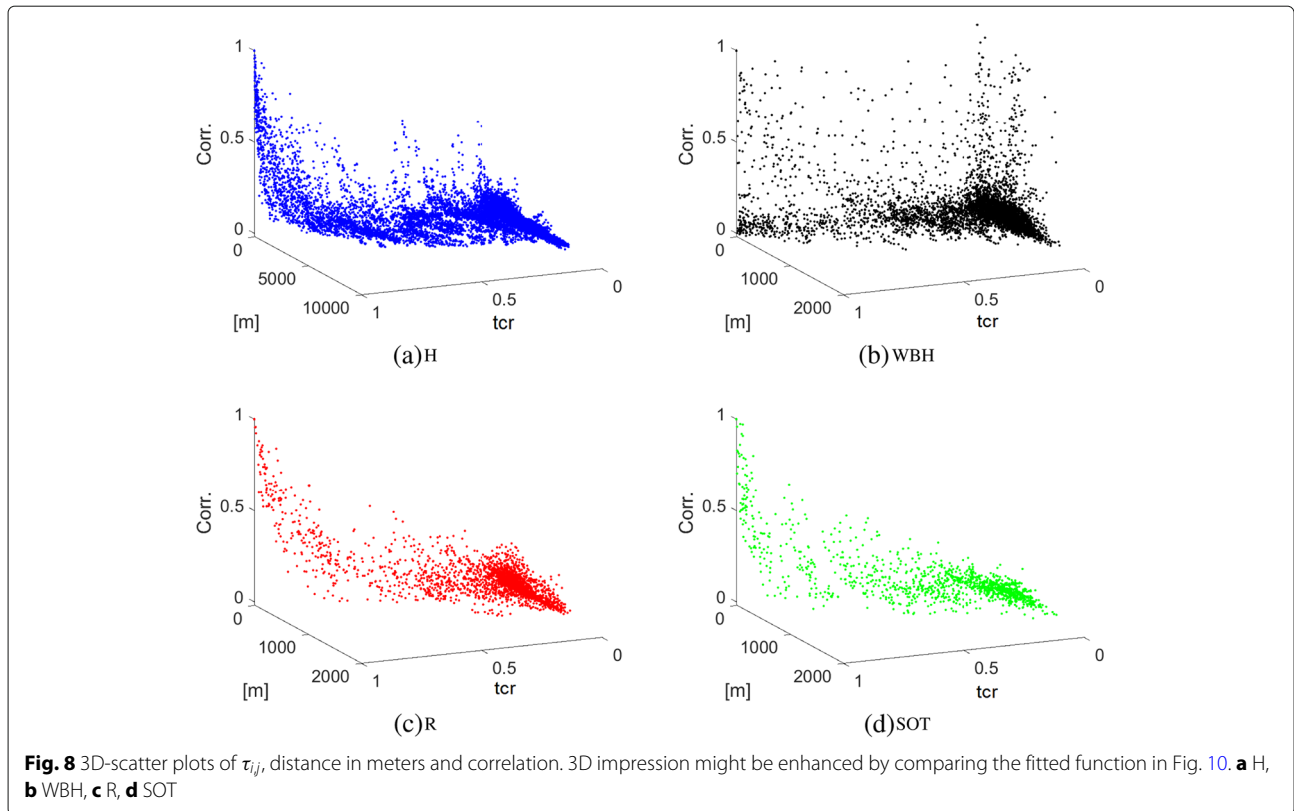
link pairs. Changing temporal aggregation to hourly values reduces the rejection rates to 4% (H), 25% (WBH), 20% (R) and 19% (SOT). Going to aggregation for two hours, however, eliminates all rejections. Based on these results, in the following we will use the assumption of constant correlation over time-of-day-intervals.

It might be suspected that there are a few factors driving the deviation from normal conditions such as an unexpectedly high level of congestion uniformly on the whole urban region (corresponding for example to weather incidents such as snowfall, heavy rain etc.). The data in all four regions do not support this hypothesis. In all four regions approximately 60% of the factors are needed to reach a cumulative explanation of 90% of the variance in a factor model.

Spatial correlations are low with an average of 0.10 to 0.05 (at a standard deviation of 0.12 to 0.18). For high values of trip count ratio $\tau_{i,j}$, however, substantial correlation exists in particular for the urban highway and the arterial in Hietzing. For $\tau_{i,j} > 0.8$ we obtain mean values of 0.41 (H, standard deviation: 0.29), 0.35 (WBH, std: 0.35), 0.61 (R, std: 0.26) and 0.56 (SOT, std: 0.29). This demonstrates that the $\tau_{i,j}$ values influence the correlation.

Figure 8 provides four 3D scatter plots showing the dependence structure of the correlation on the $\tau_{i,j}$ values and the driving distances for the four datasets. In all cases almost only positive correlations are observed. The minimum value of all correlations over all four datasets amounts to -0.04. Most of the correlations are small with some cases reaching almost perfect correlation values equal to 1.

The H dataset shows the expected behaviour with high correlations occurring exclusively for links with small distances and high $\tau_{i,j}$ values. Note that this dataset shows mainly two arterials in and out of the city with few alternative routes. A similar behaviour can be seen for the urban highway dataset SOT where, however, very few



pairs of links with larger distance and high τ_{ij} values are contained. The same behaviour also occurs in the R dataset.

The urban WBH dataset shows distinctly different patterns with a cloud of few points corresponding to small τ_{ij} values and high correlations, the remaining points possessing very small correlations. Such data points are occasionally seen also in the H dataset for small distances. Many instances of such pairs of links occur on segments of streets in opposite direction which indicates influences of common disturbances affecting both directions. On the urban highway SOT and the arterials in the H data set such scenarios do not occur. Interestingly this also does not occur in the R dataset.

We fit feedforward neural networks with the logistic function as the activation function with two hidden layers with two nodes each and additionally a constant as input variable to the data (compare (8))

$$\rho_{ij} = f(d_{ij}, \tau_{ij}; \theta) + u_{ij}$$

for all pairs i, j of links. A separate model is fitted for each data set. As for all but the H data set only data for distances less than 2km are contained, the property that correlations tend to zero for increasing distance is explicitly imposed in these data sets by data augmentation methods.

The results can be seen in Fig. 9. As expected, high correlations are obtained only for small distances and high values of τ . For τ of 0.8 correlations already are smaller than 0.2 in all four datasets except for extremely small distances. Also for distances of 0.5km correlations are smaller than 0.4.

Again the WBH dataset is special: The pseudo- R^2 value for H (0.93), R (0.92) and SOT (0.91) are quite high while for WBH we achieve only $R^2 = 0.31$.

For the comparisons in the next section a joint model for data from all four datasets is computed and presented in Fig. 10. The R^2 values for the four datasets decrease slightly to 0.91 (H), 0.88 (R), 0.89 (SOT) and 0.21 (WBH).

4.3 Comparison to single trips

The previous discussion led to the development of a number of models for the variance of travel time prediction errors which are validated in this section using trip data of single taxis obtained out of sample after the modeling took place. Here nine days (Sunday 1.1.-Wednesday 4.1., Sunday 8.1.-Tuesday 10.1., Wednesday 1.2. and Tuesday 24.7) in 2012 of trip data are used. These days contain weekdays and weekends, holiday periods (1.1.-4.1., 24.7.) and school periods.

On these days for a total of 8 heavily used routes in the four data sets single trip start and end points are estimated. Details on the routes are given in Table 2 of [1], the location of the routes is presented in Fig. 7 in [1].

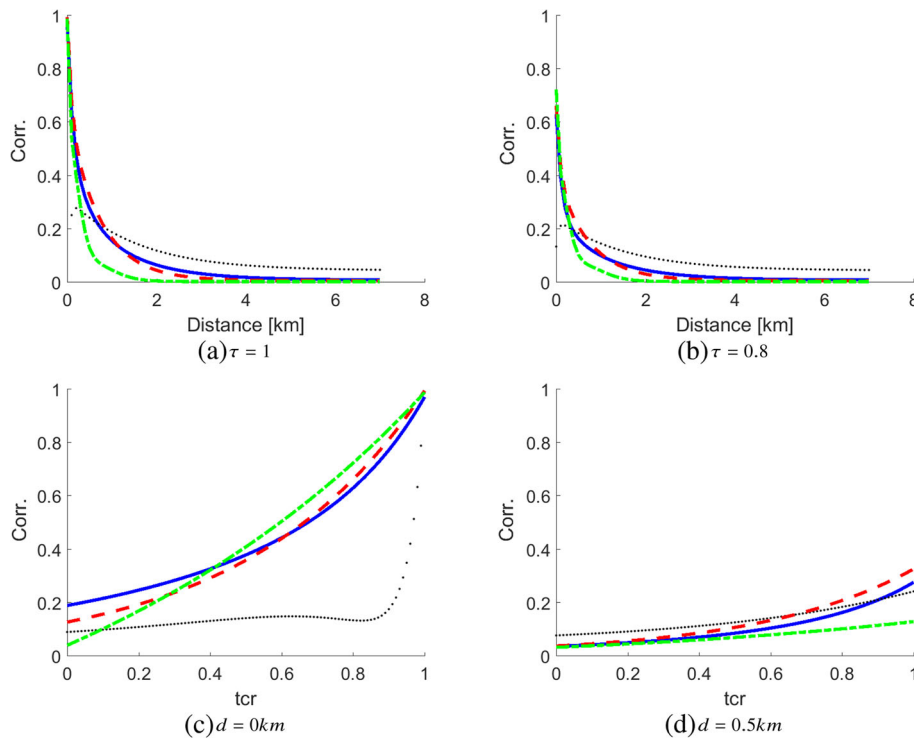


Fig. 9 Model for the correlation as a function of distance and τ values for the four datasets. H (blue), WBH (black, dotted), R (red, dashed), SOT (green, dash-dotted). **a** $\tau = 1$, **b** $\tau = 0.8$, **c** $d = 0km$, **d** $d = 0.5km$

The thus measured travel times are classified into time-of-day-intervals using the estimated time of entry to the route. Subsequently the corresponding predicted route travel time is subtracted in order to obtain deviations from predictions. One example for each setting is given in Fig. 11. It is visible that the route travel times are not estimated unbiased in all cases: For the urban settings the predictions appear to be slightly larger than the measured

travel times, for the urban highway SOT the contrary is visible.

The average bias as a percent of mean measured travel times amounts to 12 and 13% for H, 32 and 23% for WBH, 10 and 12% for R. For the SOT we underestimate travel time on average by 8 and 19%. This holds although on the validation sample no bias in the predictions has been detected (see Fig. 7). Note, however, that the validation period is limited to a few days in January 2012 where the weather conditions might interfere with predictions.

Corresponding to the various routes the travel time variance is estimated according to Eq. 9. For each single trip we calculate the deviation between the measured and the predicted route travel time and divide by $\sqrt{\hat{V}(u_{d,i}^R)}$. If the variance is correctly estimated then the corresponding sample should show unit variance. If the variance is underestimated then the normalized prediction errors have empirical variance larger than one; if the variance is overestimated the normalized prediction errors show empirical variance smaller than unity. Naturally the route travel time measurements are also subject to measurement errors.

Table 1 presents standard errors of the normalized prediction errors. It can be seen that assuming zero correlations underestimates the uncertainty in almost all cases.

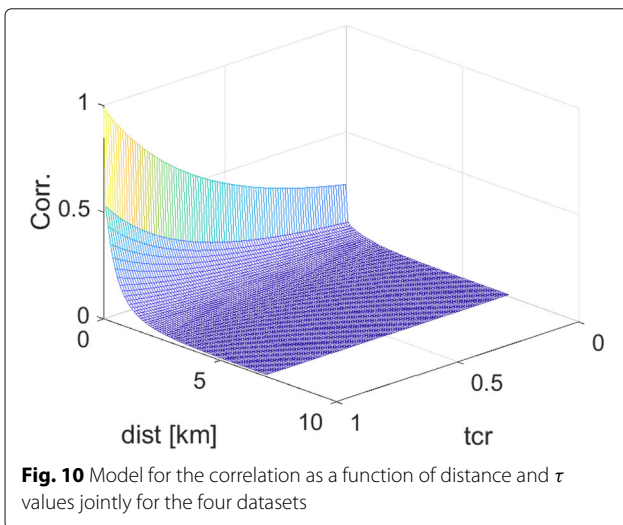
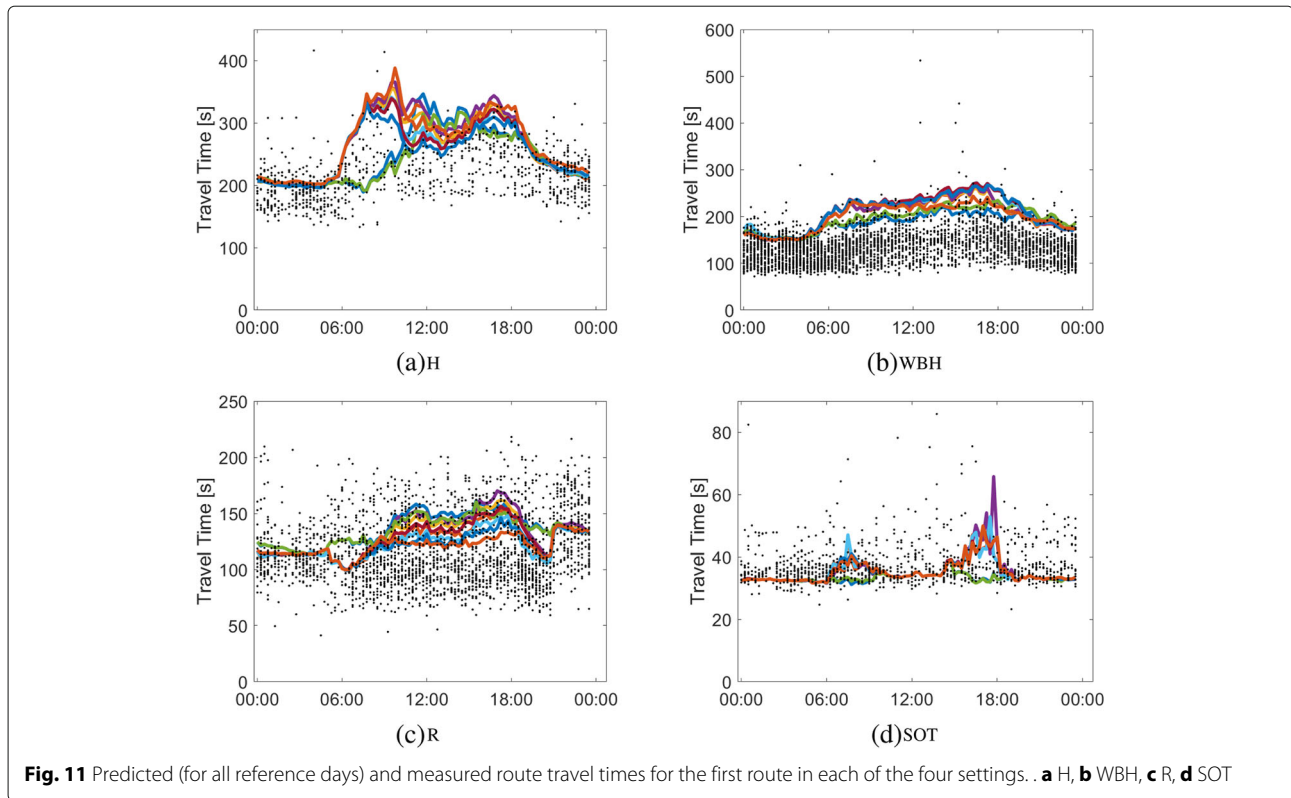


Fig. 10 Model for the correlation as a function of distance and τ values jointly for the four datasets



On the other end of the spectrum assuming total correlation leads to the smallest variances in almost all cases. The normalizations assuming perfect correlation over all links in most cases overestimate the variances, in many cases substantially so. All other four methods lie in between these two extremes. The differences between EM and ET are negligible. The individual models and the joint model (MI and MJ) perform similar.

In more detail one observes that the normalization works better during daytime and underestimates variability at night-time to a higher degree, see Fig. 12: in (a) the normalized (using model MI) residuals and all routes are plotted. It can be seen that during daytime (between

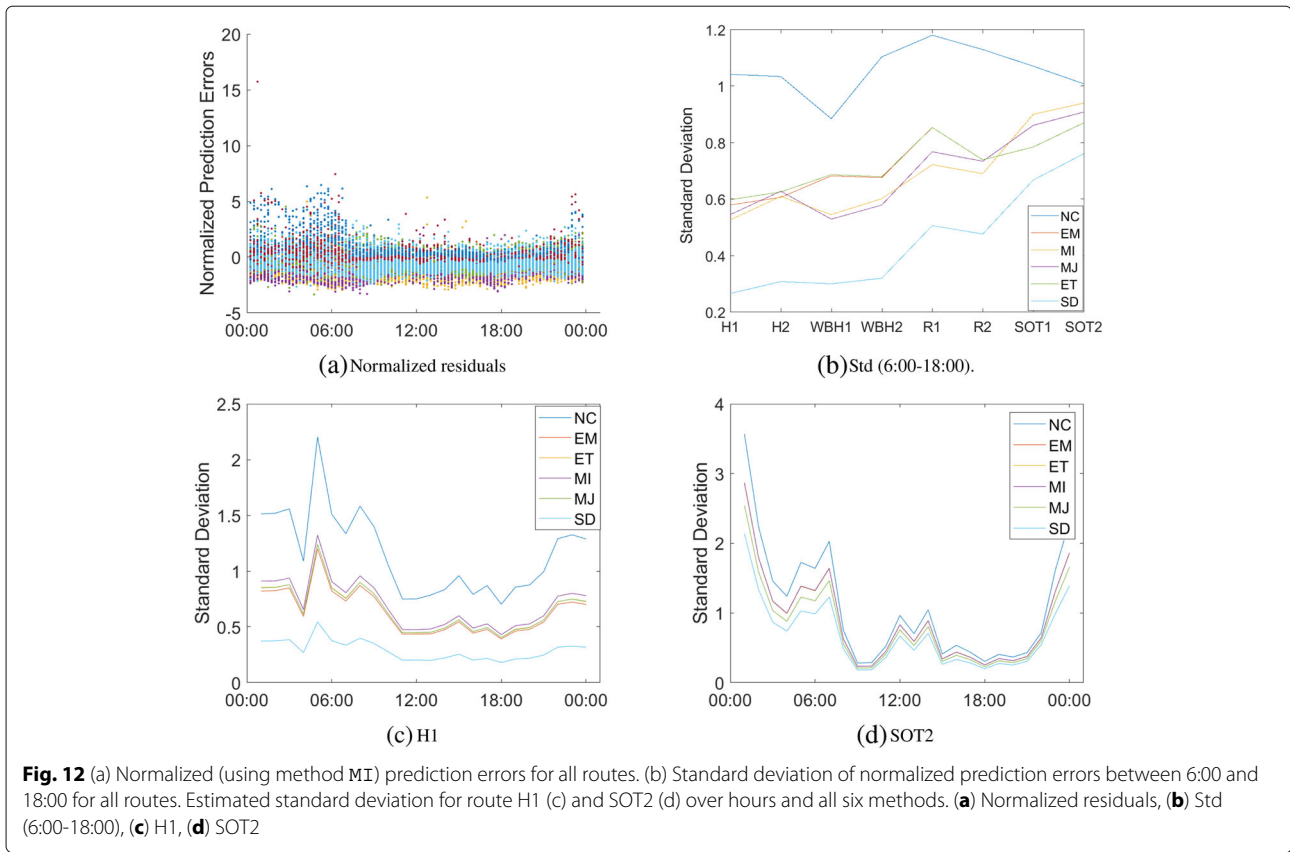
6:00 and 18:00) the mass of the residuals approximately lies in an interval of length 6 (corresponding roughly to a three sigma confidence interval), while for night-time the interval is larger. Plot (b) confirms this by providing the standard deviation in all eight routes only during daytime. Here the lower bound given by SD is substantially smaller than 1 for H, WBH and R, while the four methods EM, ET, MI, MJ all show standard deviations between 0.6 and 1.

Note that Fig. 12a indicates that the normalized prediction errors for all links and all time instants show stable distributions over time. Figure 13 investigates this in more detail by providing kernel density estimates of the normalized route travel time prediction errors which have been centered by subtracting the mean for better comparability. The normalized distributions are similar for one route in Hietzing and Westbahnhof, slightly different for the Ring dataset and completely different for SOT. These plots add to the evidence that the heteroskedasticity can be adjusted for by modelling the standard deviation while the shape of the distribution appears to be less affected. However, it also shows that this is not necessarily the case as some distributions deviate.

This also provides some indication that the assumption of correlations between model errors being similar to correlations between prediction errors is realistic. However, more research based on more appropriate data sets is needed.

Table 1 Standard deviation of normalized travel time prediction errors

	NC	EM	ET	MI	MJ	SD
H1	1.33	0.73	0.76	0.67	0.69	0.33
H2	1.12	0.65	0.67	0.65	0.67	0.32
WBH1	0.96	0.74	0.75	0.59	0.57	0.32
WBH2	1.19	0.72	0.73	0.64	0.61	0.34
R1	1.22	0.88	0.88	0.75	0.79	0.52
R2	1.19	0.79	0.79	0.73	0.78	0.51
SOT1	1.50	1.08	1.08	1.25	1.19	0.91
SOT2	1.67	1.42	1.42	1.55	1.49	1.22



Somewhat of an outlier in these comparisons is the SOT data set. Here as can be seen in Fig. 12d the variability is overestimated heavily during daytime with all approaches while it is underestimated close to midnight. Fig. 11d shows that the observed route travel times do not contain observations of heavy congestion on the SOT during typical peak hours. Therefore the failure to match

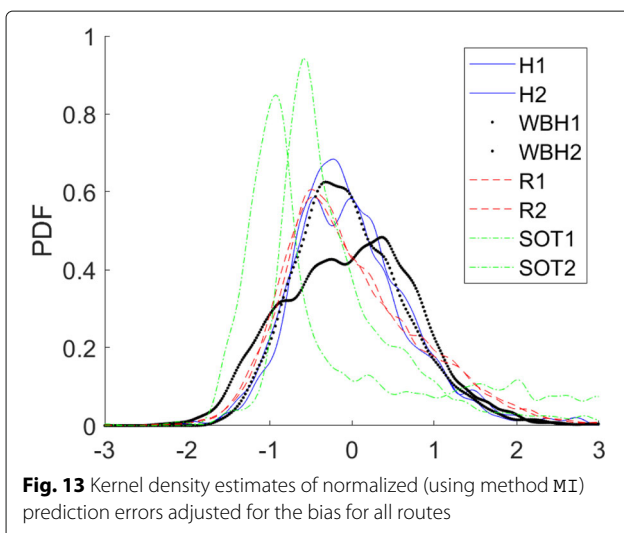
uncertainty might be an artefact of too few validation measurements.

5 Conclusion

In this paper, based on a large real world data set, models for the estimation of route travel times and the corresponding associated uncertainty have been obtained. Application of the models demonstrates that predictions of link travel times show considerable heteroskedasticity that needs to be taken into account for accurate estimation of route travel time uncertainties. We find that heteroskedasticity is related to the number of vehicles observed on each link in each time-of-day-interval but also to the traffic conditions. Explicit models for the dependency are derived.

Investigating the model errors in link travel time estimates further we found significant correlations between residuals on adjacent links that additionally have been shown to depend on the joint usage of roads. In this respect the trip count ratio is used as an indicator of joint usage and shown to have an impact on the spatial correlation.

Based on models for the correlation of link travel times as a function of the distance and the trip count ratio, formulas for the route travel time prediction error variance are suggested. Using directly measured route travel times



we find that including the correlation into the calculation of the route travel time uncertainty appears to result in partially superior results compared to simple assumptions of zero or perfect correlation. However, we also found that explicit modelling of the correlation only leads to minor performance enhancements compared to simple models using sample correlations for nearby links and setting the correlation to zero for distances larger than 1 km.

Concluding this leads to the suggestion to quantify route travel time uncertainty based on empirical spatial correlation estimates which can be confined to adjacent links and hence do not face the same data problems that empirical estimates in the whole network face. Using empirical correlation has the advantage of not requiring any other information (such as trip count ratios, location of traffic lights and so forth). Moreover this might also alleviate the restriction to correlations of model errors which needed to be imposed in this paper due to data availability. Alternatively prediction error correlations for adjacent links could be measured directly based on high frequent taxi FCD. The analysis in this paper justifies the usage of this simple method over more complex model based approaches.

Our model uses a simple scaling approach by modelling the model error being distributed according to a unique distribution scaled by a standard deviation depending on the current traffic conditions. This assumption has been verified empirically in the paper for the model errors. For the prediction errors a partial verification is contained in Fig. 13. However, this figure also contains substantial deviations that need to be investigated in more depth.

Summing up a model for the estimation of route travel time variability can be obtained based on the material in this paper which is also operational for a large street network without relying on excessive amounts of other data than the floating taxi measurements.

Acknowledgments

Part of the work has been done while the first author was with the AIT Austrian Institute of Technology GmbH.

We gratefully thank Taxi 31300 (taxi31300.at) and Taxi 40100 (taxi40100.at) for providing the taxi data used in this study and the AIT Austrian Institute of Technology (in particular Hannes Koller has been very helpful with the details) for processing the raw data and making the data available. We acknowledge support for the Article Processing Charge by the Deutsche Forschungsgemeinschaft and the Open Access Publication Fund of Bielefeld University.

Authors' contributions

DB: data analysis modelling, writing of the paper; supervising the work of MT. MT: data preparation, preliminary analysis, modelling, drafting results section. WS: structuring, writing and proof reading of the paper; supervising the work of MT. All authors read and approved the final manuscript.

Funding

The work did not receive any outside funding.

Availability of data and materials

The data set is proprietary.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Bielefeld University, Universitätsstrasse 25, 33619 Bielefeld, Germany. ²AIT Austrian Institute of Technology GmbH, Giefinggasse 2, A-1220 Wien, Austria. ³TU Wien, Wiedner Hauptstr. 8/105-2, A-1040 Wien, Austria.

Received: 27 March 2019 Accepted: 6 September 2019

Published online: 13 November 2019

References

- Bauer, D., & Tulic, M. (2018). Travel time predictions: should one model speeds or travel times? *European Transport Research Review*, 10(2), 46.
- Bovy, P.H., & Thijs, R. (2000). *Estimators of travel time for road networks: New developments, evaluation results, and applications*. Netherlands: Delft University Press.
- Carrion, C., & Levinson, D. (2012). Value of travel time reliability: A review of current evidence. *Transportation Research Part A: Policy and Practice*, 46(4), 720–741. <https://doi.org/10.1016/j.tra.2012.01.003>.
- Giraitis, L., Kokoszka, P., Leipus, R., Teyssi re, G. (2003). Rescaled variance and related tests for long memory in volatility and levels. *Journal of Econometrics*, 112(2), 265–294. [https://doi.org/10.1016/S0304-4076\(02\)00197-5](https://doi.org/10.1016/S0304-4076(02)00197-5).
- Heinze, C., Leodolter, M., Koller, H., Bauer, D. (2016). Transferring urban traveling speed model fits across cities. *European Transport Research Review*, 8(3), 19. <https://doi.org/10.1007/s12544-016-0206-8>.
- de Jong, G., & Bliemer, M. (2015). On including travel time reliability of road traffic in appraisal. *Transportation Research Part A: Policy and Practice*, 73, 80–95.
- Khosravi, A., Mazloumi, E., Nahavandi, S., Creighton, D., van Lint, J.W.C. (2011). Prediction intervals to account for uncertainties in travel time prediction. *IEEE Transactions on Intelligent Transportation Systems*, 12(2), 537–547. <https://doi.org/10.1109/TITS.2011.2106209>.
- Lam, T., & Small, K. (2001). The value of time and reliability: measurement from a value pricing experiment. *Transportation Research Part E: Logistics and Transportation Review*, 37, 231–251.
- Li, R., Chai, H., Tang, J. (2013). Empirical study of travel time estimation and reliability. *Mathematical Problems in Engineering*, 2013, 1–9. <https://doi.org/10.1155/2013/504579>.
- Lint, J.V., Zuyleen, H.V., Tu, H. (2008). Travel time unreliability on freeways: Why measures based on variance tell only half the story. *Research Part A: Policy and Practice*, 42(1), 258–277.
- O'Sullivan, A., Pereira, F.C., Zhao, J., Koutsopoulos, H.N. (2016). Uncertainty in bus arrival time predictions: Treating heteroscedasticity with a metamodel approach. *IEEE Transactions on Intelligent Transportation Systems*, 17(11), 3286–3296. <https://doi.org/10.1109/TITS.2016.2547184>.
- Pattanamekar, P., Park, D., Rilett, L., Lee, J. (2003). Dynamic and stochastic shortest path in transportation networks with two components of travel time uncertainty. *Transportation Research C*, 11(5), 331–354.
- Rahmani, M., Jenelius, E., Koutsopoulos, H.N. (2013). Route travel time estimation using low-frequency floating car data. In *16th International IEEE Conference on Intelligent Transportation Systems (ITSC 2013)*. <https://doi.org/10.1109/ITSC.2013.6728569> (pp. 2292–2297).
- Taylor, M. (2013). Travel through time: the story of research on travel time reliability. *Transportmetrica B: Transport Dynamics*, 1(3), 174–194. <https://doi.org/10.1080/21680566.2013.859107>.
- Tulic, M., Bauer, D., Scherrer, W. (2014). Link and Route Travel Time Prediction Including the Corresponding Reliability in an Urban Network Based on Taxi Floating Car Data. *Transportation Research Record: Journal of the Transportation Research Board*, 2442, 140–149.
- Vlahogianni, E.I., Karlaftis, M.G., Golias, J.C. (2014). Short-term traffic forecasting: Where we are and where we're going. *Transportation Research Part C: Emerging Technologies*, 43(May 2016), 3–19. <https://doi.org/10.1016/j.trc.2014.01.005>.
- Zheng, F. (2011). *Modelling urban travel times*. PhD thesis, Delft University.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.