Long-term prediction of vehicle speeds



Long-term prediction of vehicle speeds

by

Anne Siersema

Master Thesis Civil Engineering and Management University of Twente Faculty of Engineering Technology Transport Engineering and Management

Author

Educational institution External organisation Daily supervisor External supervisor Chair Date A. Siersema a.siersema@alumnus.utwente.nl University of Twente Simacan ir. O.A.L. Eikenbroek ir. J. Soesbergen prof.dr.ir. E.C. van Berkum February 10, 2020

Voorwoord

In je leven heb je doelen om te bereiken en tussen deze doelen in zitten geluksmomenten.

Dit is het eindverslag van een onderzoek waarmee ik een van deze grote doelen wil halen: het behalen van mijn studie en daarmee het eindigen van mijn studententijd.

Tijdens je studie, en ook tijdens dit onderzoek zaten geluksmomenten doordat je kleine resultaten boekte. Maar aan het eind van je studie, de tijd waarin je dit onderzoek schrijft, is het niet je studie die voor deze geluksmomenten zorgt.

Wel blijf je tijdens het schrijven een doel achtervolgen: *wat zou het toch mooi zijn, dat je door middel van dit verslag, aan anderen kan overbrengen wat voor moois je allemaal gedaan hebt*. Echter zorgt het ook voor frustraties: een afstudeerverslag is immers nooit perfect, maar het doel om duidelijk aan mensen te kunnen overbrengen wat je gedaan hebt kan zeker wel nagestreefd worden. Ik wil daarom ook graag iedereen bedanken die het mogelijk heeft gemaakt dit grote doel te gaan bereiken door mijn studie af te ronden.

De geluksmomenten zitten in alle kleine dingen eromheen. Voor sommigen is dit misschien de pakketbezorger die exact op hetzelfde tijdstip aankomt als dat voorspelt is of het navigatiesysteem die je informeert om een andere route te nemen omdat er een ongeluk is gebeurd. Voor mij is het misschien meer het sneller kunnen fietsen dan de auto's in de file naast je op de weg.

Maar mijn voornaamste geluksmomenten zijn de trainingen 's avonds bij Aloha, het sporten of afspreken met vrienden, het bezoeken van familie en het leven ontdekken met je vriend. Daarnaast zijn ook collega's en studievrienden altijd een goede reden voor mij geweest om elke dag weer vol enthousiasme aan de slag te gaan. Ik ben iedereen altijd heel dankbaar voor alle geluksmomenten die ze me geven en ik zal dit ook altijd blijven.

Anne

Summary

Research objective In this research, a method is designed and evaluated for predicting vehicle speeds. The results of this study are intended to make a better prediction of the arrival times of parcel deliverers. That is why this designed prediction method is suitable for the next three aspects:

- For a prediction horizon up to 3 hours;
- During regular and irregular situations;
- On freeways and urban roads for multiple regions.

Additionally, this study has the restriction that predictions are calculated real-time without storing speed predictions for every location, time and scenario before-hand.

To my knowledge, there are no studies that make a prediction that (i) meets these three requirements (ii) meets this restriction, (iii) make a prediction based on only speed measurements and (iv) improves the prediction by adding traffic information.

Method A short- and long-term prediction is made using roadway speed data based on the idea that traffic conditions have repetitive patterns; these patterns are also analyzed short-term and long-term. The long-term prediction is based on clustered historical data. In order to predict short-term, this prediction is updated with real-time measurements at the location itself and the road segments downstream of the location. The Kalman filter is used to remove the noise from the measurements. In addition, the prediction is adjusted if an accident occurs by analyzing the impact of accidents on speeds.

Discussion A challenge of this study is that the used speed data is rightcensored. This means that speeds above a certain boundary value, mostly during freeflow conditions, are unknown. This results in that predictions will not be accurate around the freeflow conditions and that this accuracy could not be evaluated around freeflow conditions. Moreover, it is hard to conclude whether the designed method makes good predictions, because it depends a lot on the type of data.

Conclusion and recommendations Although the method is only evaluated on freeways, the method can also be applied on urban roads. Besides, the method works for all circumstances by only use of vehicle speed data. However, the prediction accuracy can be improved if it is known that irregular situations occurred.

In this study, the prediction is only adapted if it is known that an accident happened. It is advised to investigate the impact of other irregular situations on the traffic conditions for better prediction accuracies.

The long-term prediction method gives more accurate predictions for all time horizons longer than 30 minutes compared to other prediction methods. Adding realtime measurements of the location itself improves the accuracy of the prediction method short-term. However, real-time measurements of road segments downstream the location, do not improve the accuracy. The lack of success is most probably caused by the fact that congestion propagation is dependent on the location. Therefore it is advised to investigate better the impact of location characteristics on the congestion propagation to improve the accuracy of the short-term prediction.

Samenvatting

Onderzoeksdoel In dit onderzoek wordt een methode ontworpen en geëvalueerd voor het voorspellen van snelheden van auto's. De uitkomsten van dit onderzoek zijn bedoelt om de aankomsttijden van pakketbezorgers beter voorspellen. Daarom is deze methode vooral gericht op drie speerpunten:

- Voor een voorspellingshorizon tot 3 uur;
- Onder reguliere en irreguliere omstandigheden;
- Voor snelwegen en gebiedsontsluitingswegen.

Daarnaast heeft dit onderzoek de restrictie dat de methode een voorspelling real-time kan maken zonder dat er data van tevoren wordt opgeslagen voor elk wegsegment, tijd en scenario.

Zover ik weet, zijn er geen onderzoeken bekend die een voorspelling maken die (i) zich focust op deze drie speerpunten, (ii) aan de restrictie voldoet, (iii) waarbij de voorspelling is gemaakt op basis van snelheidsmetingen en (iv) waarbij de voorspelling kan worden voorbeterd door verkeersinformatie toe te voegen.

Methode Een lange en korte termijn voorspelling is gemaakt aan de hand van verkeerssnelheden gebasseerd op het idee dat verkeersconditions hebben herhalende patronen; deze patronen zijn geanalyseerd voor zowel korte als lange termijn. De lange termijn voorspelling is gemaakt op basis van gegroepeerde historische data. Om op korte termijn te voorspellen, is deze voorspelling geupdate met actuele metingen op de locatie zelf en op de wegsegmenten stroomafwaards van de locatie. Het Kalman filter is gebruikt om de ruis uit de metingen te verwijderen. Aanvullend wordt de voorspelling aangepast als er een ongeluk plaatsvindt. Dit is gedaan door impact van ongelukken op snelheden te analyseren.

Discussie Een uitdaging van dit onderzoek is dat de gebruikte snelheidsdata censurering bevat. Dit houdt in dat snelheden boven een bepaalde waarde, rondom vrije doorstroming, onbekend zijn. Dit zorgt ervoor dat tijdens vrije doorstroming een voorspelling minder precies de werkelijkheid voorspelt en bovendien ook moeilijk te evalueren is. Daarnaast is het moeilijk te concluderen of de ontworpen methode goede voorspellingen maakt, omdat het erg van het type data afhangt.

Conclusie en aanbevelingen Hoewel de voorspellingsmethode alleen geevalueerd is op snelwegen, kan de methode ook gebruikt worden voor andere wegen zoals gebiedsontsluitingswegen. Daarnaast werkt de methode onder alle omstandigheden en hoeft hiervoor enkel snelheidsdata gebruikt te worden. Echter kan de nauwkeurigheid van de voorspelling verbeterd worden als er bekend is of er niet-reguliere verkeerssituaties zijn. Nu houdt de ontworpen voorspellingsmethode alleen rekenening met ongelukken die hebben plaatsgevonden. Het wordt geadviseerd om andere niet-reguliere verkeerssituaties ook mee te nemen om de voorspelling nog nauwkeuriger te maken.

De lange termijn voorspellings methode geeft nauwkeurigere voorspellingen dan de andere vergeleken methoden voor tijdshorizons langer dan 30 minuten. Het toevoegen van real-time metingen van een te voorspellen locatie zorgt ervoor dat de voorspelling nauwkeuriger worden op korte termijn. Echter wordt de voorspelling niet nauwkeuriger door het toevoegen van real-time metingen van wegsegmenten stroomafwaards van de te voorspellen locatie. Dit wordt hoogstwaarschijnlijk veroorzaakt doordat de filevorming niet goed genoeg rekening houdt met wegkarakteristieken. Daarom wordt aanbevolen om beter de effecten van wegkarakteristieken op filevorming te onderzoeken zodat de nauwkeurigheid van de kortetermijn voorspelling kan worden verbeterd.

Terminology

Road types

highways	roads that are especially designed to drive on high- speed. They can have long stretches without any intersections, but can pass through urban areas and have traffic signals. (Dutch: snelwegen en N- wegen)
freeways	roads on which drivers are not hindered by inter- sections. The roads are seperated from other traf- fic directions and can only be accessed by ramps. (Dutch: snelwegen)
urban roads	roads smaller than highways, but having the func- tion of smooth flow and connecting neighborhoods (Dutch: gebiedsontsluitingswegen)

Prediction types and variation types

short-term prediction	in this research, this is the prediction for recurrent and non-recurrent conditions based on historical data and real-time data that will predict till 90 min- utes
long-term prediction	in this research, this is the prediction for recurrent conditions based on historical data that will predict for all time horizons but with less accuracy than the short-term prediction
baseline prediction	see long-term prediction
instantaneous predic-	in this research, this is the prediction using the
tion	current speed measurement as expected vehicle speed
temporal correlation	represents the relationship of speeds between ad- iacent times
spatial correlation	represents the relationship of speeds between ad- jacent locations
spatiotemporal correla- tion	represents the relationship of speeds between ad- jacent times and locations
noise	which are variations that do not give information about the next variations
systematic variations	variations that do give information about the next variations

residuals	errors of speed measurements compared to the baseline speed
Traffic conditions	
freeflow	the condition on a road when no congestion occurs
congestion	the condition on a road when it is too busy to drive the desired speed
regular conditions	traffic conditions that recur with a constant period at a location
irregular conditions	traffic conditions that are recurrent, but do not recur with with a constant period at a location
irregular situations	these are situations that recur incidentally; in this study these consist of accidents, roadworks, weather and holidays
incidents	see irregular situations
accident	disturbance on traffic conditions caused by vehi- cles; for example a car crash
holidays	these are days on which business, work or school are suspended or reduced. Most days are set by law or custom such as national celebration days or weeks that are meant for vacation. This results in days on which people do more recreation activities
Speeds	
freeflow speed	the speed that vehicles regularly drive in freeflow conditions. In this study, this speed is constant over time, but location dependent
boundary speed	this refers to the characteristics of the right cen- sored data. It is the maximum speed for which mea- surements are known. For the used TomTom data, this is constant over time, but location dependent.
freeflow	the condition on a road when no congestion occurs

Abbreviations and data sources

KNMI	Koninklijk	Nederlands	Meteorologisch	Instituut;
	weather in	formation is u	sed from this data	a source

NDW	National Data Warehouse for Traffic Information; accidents and roadwors are used from this data
	source
RWS-ADY	Rijkswaterstaat Acquisition Dynamic Data; matrix sign information is used from this data source
VILD	Traffic Information Location Database; a Dutch standard of static location referencing
openLR	Open Location Reference; a dynamic form of loca- tion referencing designed by TomTom which is open source
TomTom HDFlow	real-time road traffic information provided by Tom- Tom based on floating car data; speed are used from this data source
TomTom Speed Profiles	a realistic average roadway speed for all times of the day and for each day of the week provided by TomTom
FRC	Functional Road Class; describes the size and type of utility of the roads, such as freeways, local roads or bicycle paths
FOW	Form Of Way; describes the physical form of the road such as motorways, slip roads and round-abouts
alertC	codes that describe detailed traffic situations such as type of accidents, roadworks and weather

Contents

1	Intro 1.1 1.2 1.3 1.4 1.5	Description Simacan: company description Background Research Aim and research questions Outline Summer State	1 2 2 10 12
2	Data 2.1 2.2 2.3 2.4	a Speed measurements: TomTom HDFlow	14 14 16 18 19
3	App 3.1 3.2 3.3 3.4 3.5	roachLong-term prediction of vehicle speeds for regular conditionsAnalysis of short-term variations in vehicle speedsShort-term prediction of vehicle speedsEvaluation of the accuracy of the vehicle speed predictionsConclusion of the approach	22 27 30 36 39
4	Res 4.1 4.2 4.3 4.4	ultsLong-term prediction of vehicle speeds for regular conditionsAnalysis of short-term variations in vehicle speedsEvaluation of the accuracy of the vehicle speed predictionsConclusion of the results	41 42 46 54
4	Res 4.1 4.2 4.3 4.4 Disc 5.1 5.2 5.3 5.4	ults Long-term prediction of vehicle speeds for regular conditions Analysis of short-term variations in vehicle speeds Evaluation of the accuracy of the vehicle speed predictions Conclusion of the results Conclusion of the results Long-term prediction of vehicle speeds for regular conditions Construction Short-term prediction of vehicle speeds for regular conditions Long-term prediction of vehicle speeds Short-term prediction of vehicle speeds Short-term prediction of vehicle speeds Evaluation of the accuracy of the vehicle speed predictions	41 42 46 54 56 56 56 57 58
4 5 6	Res 4.1 4.2 4.3 4.4 Disc 5.1 5.2 5.3 5.4 Con 6.1 6.2 6.3	ults Long-term prediction of vehicle speeds for regular conditions Analysis of short-term variations in vehicle speeds Evaluation of the accuracy of the vehicle speed predictions Conclusion of the results Cong-term prediction of vehicle speeds for regular conditions Cong-term prediction of vehicle speeds for regular conditions Analysis of short-term variations in vehicle speeds Short-term prediction of vehicle speeds Evaluation of the accuracy of the vehicle speeds Conclusion Conclusion Evaluation of the accuracy of the vehicle speed Provide the accuracy of the vehicle speed predictions Conclusion Conclusion Recommendations for further research Recommendations for Simacan	41 41 42 46 54 56 56 56 56 57 58 59 61 63
4 5 6 Re	Res 4.1 4.2 4.3 4.4 Disc 5.1 5.2 5.3 5.4 Con 6.1 6.2 6.3 ferer	ults Long-term prediction of vehicle speeds for regular conditions Analysis of short-term variations in vehicle speeds Evaluation of the accuracy of the vehicle speed predictions Conclusion of the results Cong-term prediction of vehicle speeds for regular conditions cussion Long-term prediction of vehicle speeds for regular conditions Analysis of short-term variations in vehicle speeds Short-term prediction of vehicle speeds Evaluation of the accuracy of the vehicle speed predictions Conclusion Conclusion Conclusion Short-term prediction of vehicle speeds Evaluation of the accuracy of the vehicle speed predictions Conclusion Recommendations Conclusion Recommendations for further research Recommendations for Simacan	 41 42 46 54 56 56 57 58 59 61 63 65

A. Weather	66
B. Visualisation of the location of the accidents	72
C. Data Processing of speed measurements: TomTom HDFlow	74
D. AlertC codes indicating irregular situations	78
E. Extended foundation of the baseline prediction	82
F. State estimation of the Kalman Filter	84
G. The estimation of the short-term prediction by combining several	
timeseries	89
H. Extended results of the evaluation	92

1 Introduction

Logistic operations of companies providing home deliveries are inefficient with respect to route plans. An important cause is that the actual travel times differ from the expected travel times. These variations in travel times are difficult to predict. This especially holds for predictions of hours to days ahead which typically are used to make the initial route planning for home deliverers.

The estimation of speeds during the trip (online) of a home deliverer is important. More accurate speed predictions result in more accurate travel time predictions. More accurate travel time predictions ensure better online route adaptations and arrival time estimations, which on its turn, result in more efficient logistic operations.

Variations in speeds occur for several reasons. A part of the variations is explained by regular conditions. These variations are often known in advance so they can be predicted offline. However, another part of the variations is explained by irregular conditions such as incidents. Of course, it is unknown beforehand whether these conditions occur. But online, if a specific condition is known, logistic operators quickly want to know the effects of these incidents on travel speeds and possibly redirect the driver.

1.1 Problem context

As stated in the introduction, companies providing home deliveries want online and quickly accurate estimations of speeds. Moreover, they want predictions for many origin-destination paths so the predictions must be possible for urban roads and freeways. Also, home deliverers want predictions several hours ahead and under all traffic situations, because their trips last several hours and the situations are unknown beforehand.

The prediction of speeds for home deliverers includes three main challenges. Vehicle speeds must be predicted:

- far ahead
- during all situations (regular vs. irregular)
- network-wide (urban vs. freeway)

1.2 Simacan: company description

This thesis is conducted by the University of Twente and executed at Simacan. Simacan has an open and supplier independent cloud platform for digital collaboration in transport logistics. This platform integrates real-time traffic information in the primary business processes of companies in fast-moving consumer goods or parcel transport execution, where timeliness and predictability are essential. The platform offers real-time trip information, including expected arrival times and last-mile guidance assistance.

At the moment, the arrival times are determined from the travel time of the fastest path to the destination. Travel times are calculated with the distance of a path and the predicted speed. For the first destination of a vehicle, the predicted speed is the current measured speed. TomTom Speed Profiles are used for the destinations after that. These are speed values based on historical measurements.

1.3 Background

This research focuses on the prediction of traffic speeds. This prediction implies that with gained knowledge, a statement is drafted about the future speed that vehicles can drive on the road. With this predicted speeds and a given origin and destination, travel times and arrival times can be estimated. These travel times are an important aspect for the optimization of the transport route plans of logistic operators.

This literature review discussed and compares methods for the prediction of speeds based on the following aspects:

- 1. Time horizon
- 2. Application on different conditions (regular vs. irregular)
- 3. Network-wide application (several road types and regions)

In this research, only studies are discussed that predict speeds with a minimum prediction horizon of 10 minutes, because the method to be designed will also predict horizons longer than 10 minutes. The methods that will be compared are divided into four different classes, which are also used by Oh et al. (2015):

Model-based Model-based methods use assumptions concerning the underlying flow dynamics (e.g. shockwave theory, queuing theory) to predict future traffic conditions. These models are often used for the prediction of speeds, because speeds measurements are not always available and reliable. An example of a study that uses the model-based method for a prediction of speeds for the next 20 minutes on freeways is done by Wang et al. (2006).

Linear regression and time series In this approach, the functional relationship between the measured values and the predicted values is known. The unknown parameters of the functions that represent these relationships will be estimated with previous measurements.

Three main different sub-types exist with this method which are linear regression, autoregressive integrated moving average (ARIMA) and variants of Kalman Filters.

Linear regression these models use linear functions to explain the relationship between the measured and predicted variables. An example of a study that uses this approach for speed predictions 1 hour in advance for freeways is Kamarianakis et al. (2012). They use autoregression and temporal clusters to determine the traffic state and to update speeds in real-time when measurements exceed a certain threshold. They optimize their model with the least absolute deviation of the residuals of the speed measurements.

AutoRegressive Integrated Moving Average (ARIMA) This approach is based on autoregression and moving average. The autoregression part represents a relationship between one or multiple variables which does not imply that the relationship is linear. The integrated part compensates for the non-stationarity in the data. The moving average part reduces the noise by taking linear combinations of the measurements. Min and Wynter (2011) use ARIMA with spatial-temporal regression for a 60-minute speed prediction on highways and urban roads.

Kalman Filter This method uses the theory that measurements are partly based on temporal patterns and partly based on variations that could not be explained (noise). It uses recent measurements and their joint probability distribution to estimate the impact of the last measurements on the prediction. Wang et al. (2006) use the Kalman Filter to predict the traffic state on freeways for a time horizon of 20 minutes. Xia et al. (2011) uses the method in combination with seasonal ARIMA, but the prediction accuracy on the freeways remains unknown.

Artificial Neural networks This method uses a large amount of data to learn underlying patterns from the measurements and output relationship without knowing the functional relationship in advance. The artificial neural networks (ANN) are based on multiple layers that together represent the total relationship. These relationships can typically also not be retrieved afterwards. Examples of studies that use neural networks to predict speeds on freeways are Hamad et al. (2009),

Innamaa (2009) and Li and Rose (2011). Hamad et al. (2009) decomposites the timeseries first to the basic components and then use an ANN to predict speeds up to 25 minutes in advance. Innamaa (2009) uses a two-layer network to predict 15 minutes. She made a distinction between the amount of wind, the visibility and road surface conditions. Li and Rose (2011) use an ANN with multiple layers to predict one hour during rush hours. Additionally, they classified their data beforehand on time of day, day of the week and the amount of rainfall.

Pattern recognition This method uses a large amount of data to discover regularities in timeseries. Often, pattern recognition is in combination with clustering whereby the regularities, or patterns, are labelled to a group. Timeseries can be predicted by matching the regularities of the timeserie to a cluster and use the characteristics of this cluster to predict. Heilmann et al. (2011) predicts speeds for heavy goods vehicles up to 2 hours in advance for some locations on freeways.

Additional methods The methods that are mentioned above can be performed in combination with classification and clustering.

Classification, also called supervised learning, is the process of classifying the data with the help of known class labels. Classes that are often defined, are: weather (e.g. rain, visibility) (Innamaa, 2009; Li and Rose, 2011; Min and Wynter, 2011), road surface conditions (Innamaa, 2009), road types (Pan et al., 2015), accidents (Min and Wynter, 2011) and roadworks (Min and Wynter, 2011).

Clustering, also called unsupervised learning, is similar to classification, but no predefined class labels are used. The classes are identified based upon found similarities in the data.

1.3.1 Time horizon

The time horizon contains the duration that speeds can be predicted.

Model-based methods typically make very accurate predictions short-term based upon several data types that, when combined, represent a very accurate traffic state. However, it is difficult to predict this traffic state far ahead. In its original form, the prediction depends on the adjacent road sections and barely on the boundary conditions. Nevertheless, with long-term predictions, the boundary conditions will become important. These boundary conditions do not always represent the reality well. Consequently, this will lower the accuracy of the model.

Basic variants of linear regression and time series methods do not use historical profiles, but only a combination of measurements of the recent past. Using real-time measurements is typically suitable for short-term predictions. However, regularities which can be found in historical data and classification methods can be added to the prediction. This method can result in accurate long-term predictions.

Li and Rose (2011) concluded that in combination with artificial neural networks and real-time data, predictions longer than 10 minutes are more accurate with the use of only historical patterns.

Patterns, which are regularities found in historical data, last often several hours. So it is assumed that this method can accurately predict speeds for several hours ahead. Heilmann et al. (2011) confirmed this statement, as they made accurate predictions on freeways for up to 2 hours. However, they discussed that their method must be perfected to make network-wide predictions.

Predictions are often based upon a combination of historical profiles and realtime measurements. Most studies did short-term predictions in which they tried to get very high accuracy, ranging from seconds to 1 minute (Clark, 2003; Rice and Van Zwet, 2004; Thakuriah, 1992). For predictions for a time horizon of 5 minutes, predictions will already be accurate using only real-time measurements (Hamad et al., 2009). In this study, the aim is to have long-term predictions up to a time horizon of 3 hours for which short-term and long-term accuracies are even relevant. For long-term predictions, the actual speeds on roads barely influence the speeds for several hours in advance. That is why historical profiles perform substantially better than real-time measurements for long-term predictions.

1.3.2 Different conditions

Conditions can be divided into regular and irregular conditions. Many traffic conditions are recurrent. Particular traffic conditions that repeat with a fixed period at the same location are called regular. Conditions that are recurrent but do not repeat with a fixed period at the same location are called irregular. Sometimes the cause of these irregular conditions is known (for example accidents, roadworks, weather, holidays).

Model-based methods can easily incorporate and account for the effects of disruptions if the size of these disruptions (i.e. change in volume) is known. However, as already mentioned, assumptions must be made of the dispersion of the vehicle volumes at intersections. Moreover, the exact size of disruptions is often unknown and must be assumed.

Linear time-series analysis performs well for prediction of the occurrence of congestion on freeways during regular conditions (Davis et al., 1990). However, for non-regular conditions, models that use methods like ARIMA or Kalman Filters, perform better than methods only based on linear regression in predicting extreme values for congestion.

Neural networks can only explain and predict regular conditions, because they do not generalize outside the training set. During off-peak hours, neural networks perform well, but the accuracy is almost equal to real-time profiles that also predict accurately during this time-frame (Hamad et al., 2009). For long-term prediction during peak hours, neural networks perform better than real-time profiles.

Pattern recognition gives accurate predictions during regular and irregular conditions (Heilmann et al., 2011). Furthermore, it is assumed that pattern recognition mainly performs well during irregular conditions compared to other methods, because the method uses other (most probably irregular) conditions with comparable patterns to predict speeds.

If it is known that irregular situations will occur in the future or are occurring at the moment, models can accommodate their prediction to get more accuracy. This addition of irregular situations is possible with methods that use neural networks or methods that use time series. Furthermore, it is useful to apply pattern recognition for more accurate predictions during irregular conditions. Additionally, if irregular situations are known, the data can be classified in these irregular situations.

1.3.3 Network-wide application

The network-wide application includes the usability of the prediction method on different type of roads and different regions. The different road types that are distinguished are freeways and urban roads.

Model-based methods typically predict speeds on a relatively small scale, like stretches of freeways, and are influenced by the assumptions at the boundary of the model. These methods are difficult to apply on urban networks for several reasons. An important reason is that urban networks have many intersections; for all these intersections, assumptions must be made of the dispersion at an intersection. Furthermore, model-based predictions often use more types of data (speeds, flow, occupancy) for speed prediction than other prediction methods, which make them difficult to use in urban areas (Wang et al., 2006).

ARIMA models, but also other time series and regression models, are proper methods for network-wide predictions. ARIMA does work with only speed or travel time measurements which are often available on a network-wide scale (Min and Wynter, 2011). In general, neural networks have a good prediction accuracy if much data from several data types are available. It is shown that accurate speeds on freeways can be predicted with neural networks, because in most cases, much data is available for this road type. Innamaa (2009) states that the performance of neural networks in a complex environment, which are urban networks, depends on many input variables. She assumes that the prediction will only perform well if many of these input variables are known. However, we are not aware of any study that uses neural networks for the prediction of urban roads.

The advantage of the pattern recognition method is that traffic speeds can be predicted with only speed measurement data. Hamad et al. (2009) stated that if only speed measurements are available, self-learning methods are a good choice. However, the best of my knowledge, no studies are known that predicts speeds with a time horizon longer than 10 minutes based on pattern recognition. It is assumed that pattern recognition is difficult in these areas. The complex network structure causes less clearly visible patterns.

Most studies focus on freeways, because more data is available and freeways have a less complicated network structure due to the fewer intersections. A complex network structure causes many interactions within the system. The more interactions, the more difficult to predict the traffic speeds. All methods will have less accurate predictions on urban roads than on highways. The model-based will only predict accurately if the interactions in the network are known. For the other methods, the interactions can be unknown, but can cause less accurate predictions.

1.3.4 Other important aspects

It is assumed that in general at this moment and in the future, enough data will be available for any data type. In the past, it often occurred that limited data was available which resulted in large data aggregations (i.e. some studies aggregate data to 10 until 30 minutes intervals, which results in less accurate predictions (Abdulhai et al., 1999; Vythoulkas, 1993)). Also, if too many specific classes or clusters will be formed, each class will have limited data, so it must be avoided to have too few classes.

A significant disadvantage of neural networks and pattern recognition methods is that these methods have long calculation times or require much database storage. By usage of classification of (for example the weather or road surface) conditions and spatial and temporal clustering, the amount of database storage can be reduced while remaining good predictions (Innamaa, 2009). Classification is a good additional method if data is available to classify on, and if the amount of data per class is large enough for the method that will be used.

1.3.5 Conclusion

Vlahogianni et al. (2004) and Vlahogianni et al. (2014) stated that the quality of the prediction is difficult to compare, because the quality of the prediction not only depends on the used method, but also depends on:

- 1. spatial characteristics of the test set (highway or urban ways)
- 2. time horizon of the prediction (short-term or long-term)
- 3. data quality: delay in real-time data
- 4. data quantity: data resolution, availability of measurement types of traffic conditions and availability of spatial and temporal incident data
- 5. evaluation: compare with which methods and with which evaluation method

It depends on the availability and quality of data (variety in data types and quantity and quality of a data type) which method will give the most accurate predictions. This dependency makes it difficult to compare these methods, because every study uses other datasets. As a result of that, studies have different conclusions about the most accurate prediction method. That is why this literature study does not use a quantitative but purely a qualitative comparison. Independent of the type of method, it can be concluded that if descriptive data is available, interactions between different data types will give a better understanding, which on its turn could get better prediction accuracies.

Model-based methods are less applicable for this research because they typically predict speeds for a relatively small scale. It is possible to do it for larger scales by predicting the model input variables, but many input variables must be predicted. To predict these input variables accurately, it will take a relative intensive study for large scale models compared to other models. However, without accurate input variables, the prediction will not be accurate.

The basic linear regression and time series method does not make use of historical profiles. However, these methods make a prediction based on a combination of measurements of the recent past. This results in accurate short-term predictions. However, by adding historical profiles based on classification, high accuracies can be obtained for long-term predictions.

On the other hand, neural networks seem to have good long-term prediction accuracies. However, the relationships that contribute to the prediction are unknown,

which will make these predictions less believable. Especially for irregular conditions, people want to know the cause of the irregularity of the prediction.

Pattern recognition has good prediction results in the long-term if only speed measurements are available. However, if additional information is available about irregular conditions and the impact of these irregular conditions on the prediction results can be determined, pattern recognition does not give more accurate results than other methods.

1.3.6 Research gap

Most of the studies that predict traffic speeds made short-term predictions for freeways. No studies are known that focuses on long-term and network-wide predictions that deal with all conditions. This can be probably explained by the fact that urban areas, which are a complex environment, depends on many input variables which results in lower prediction accuracies. We are only aware of studies that focus on only long-term, network-wide or during regular and irregular conditions.

Min and Wynter (2011) predicted network-wide traffic speeds with good accuracies until 60 minutes. They use a multivariate autoregressive model with spatiotemporal correlations. However, they did not report the prediction accuracies during irregular conditions.

Studies that made accurate predictions for traffic speeds during irregular conditions did not predict on a network-wide scale, but only for freeways. Pan et al. (2015) especially predict irregular conditions caused by accidents and road works. Their model gives good accuracies for incident predictions for a time horizon of 30 minutes by optimizing their baseline prediction using autoregression. However, they cluster their data based upon occupancy and intensity measurements which give promising results, but are often not available.

Other studies only show methods that have accurate predictions of traffic speeds during peak hours on freeways, like the model-based (Wang et al., 2006), linear regression and time series (Heilmann et al., 2011; Kamarianakis et al., 2012), neural networks (Li and Rose, 2011) and pattern recognition (Hamad et al., 2009) methods. However, they did not mention the effects of irregular situations which can result in other method preferences.

From the literature study, it is recommended to design a method based upon a regression method with clustering to achieve accurate network-wide long-term predictions during regular and irregular conditions. This is encouraged because:

- Predictions could give relatively accurate predictions on a network-wide scale with only the use of speed measurements during regular conditions;
- Relationships between different parameters, like incidents, road types and daily variations become explicit;
- The prediction will be fast without using much database storage;
- Clustering mostly improves the prediction accuracy if reliable measurements for different parameters are available, especially during irregular conditions.

Thus, this study will research the variations of traffic speeds. A prediction will be performed with a linear regression and time series method. Moreover, the dynamic behaviour of traffic congestion as a result of incidents will be quantified, using the underlying method of Pan et al. (2015).

1.4 Research Aim and research questions

In this section, the research aim will be introduced. Next, the limits and boundaries of this research will be presented. This both results in the last paragraph: the research questions.

1.4.1 Research aim

As already mentioned in the introduction, the research aim is:

To design a method that predicts accurate vehicle speeds and is suitable for several aspects: for a prediction horizon up to 3 hours; during regular and irregular conditions; and on freeways and urban roads.

1.4.2 Limits and Boundaries

Concluding from the literature study and the available data, this research will have some limits and boundaries with respect to the challenges. Besides that, it also has limits and boundaries on the calculation time and needed data storage of the method.

Different conditions For highways, predictions will be made for regular and irregular conditions. Accidents, roadworks, holidays and extreme weather situations that apply to large areas are irregular situations that, furthermore, will explain most of the variations in traffic speeds. For these irregular situations, the

relationships between these variables will be explained. The prediction will also handle other irregular conditions if they are occurring.

Network-wide The focus of this study will be on freeways and not on urban roads. As a result of that, the prediction on freeways will be based upon an intensive study on the relationships between spatial, temporal and incidental variables on freeways. This intensive study will not be performed on urban roads, because it is assumed that the relationships on urban roads are more complicated and more difficult to visualize.

In addition, the prediction will not be evaluated on every freeway in the Netherlands. Only two datasets will be used: one for the analysis used for the designed method and one for the evaluation of the results of the designed method. Although it is not the whole sample, each dataset is intended to represent congestion propagation on general freeways in the Netherlands.

Used data This study will make use of floating car data, because floating car data has a high coverage world-wide for freeways and most urban roads.

In this study, the used floating car data consists of only vehicle speeds. Therefore, flow and occupancy data will not be used. Although it is assumed that more accurate predictions can be made with the inclusion of this data, the data is not always available. So, the aim is to still get accurate predictions without using flow and occupancy data. Because many variations in flow are explained by time of day, it is assumed that nearly the same performance can be achieved by using temporal and spatial clusters as by using flow or occupancy.

This study will predict the size of the congestion, given that it is known if there is congestion or not at the moment. It will not predict the chance that congestion arises from the recent measurements, which is not possible with the used data in this research.

Fast method without many data storage The method to be designed must give accurate predictions real-time. Therefore, the predictions must be able to be calculated fast. Moreover, it is not possible to store speed predictions for every location, time and scenario, because of limited data storage.

1.4.3 Research Questions

The research aim and limits and boundaries lead to the main research question:

To what extent is the designed method for vehicle speed predictions suitable for several aspects: for a prediction horizon up to 3 hours; during regular and irregular conditions; and on freeways and urban roads?

The sub-questions are:

- 1. To what extent is the designed method for vehicle speed predictions suitable network-wide?
 - (a) How can a method be designed so that it can predict vehicle speeds on urban roads and freeways?
 - (b) How can a method be designed so that it can predict vehicle speeds for different regions whereby not all road characteristics of the regions must be separately investigated?
- 2. To what extent is the designed method for vehicle speed predictions suitable for a prediction horizon up to 3 hours?
 - (a) Which long-term recurrent patterns in vehicle speeds occur?
 - (b) What is a suitable method to predict vehicle speeds using historical data that gives accurate predictions long-term?
 - (c) To what time horizon do short-term temporal and spatiotemporal systematic variations in vehicle speeds occur?
 - (d) How can the designed method deal with noise; which are variations that do not give information about the next variations?
 - (e) What is a suitable method to predict vehicle speeds using historical and real-time data that give accurate predictions short-term?
- 3. To what extent is the designed method for vehicle speed predictions suitable during regular and irregular conditions?
 - (a) What are the effects of irregular situations (accidents, roadworks, weather and holidays) on the speeds and the variations of these speeds?
 - (b) How can the designed method be adjusted if irregular conditions occur?

1.5 Outline

In Chapter 2, the used data is described. The main data source is vehicle speeds from floating car data which is chosen because it contains vehicle speeds for freeways and urban roads world-wide.

Next, in Chapter 3, the designed prediction method for vehicle speeds and the method for evaluation of the accuracy of this prediction method is described. In Chapter 4, the results are shown of the accuracy of the prediction method and the time horizon that accurate prediction could be made.

In Chapter 5, the used data and the prediction method is discussed. Finally, the conclusion, the recommendations for further research and the recommendations for Simacan can be found in Chapter 6.

2 Data

This chapter describes the data sources that are used in this study. First, the roadway speed measurements are described on which the prediction method is based.

The speed measurements are categorized into measurements during *regular* and *irregular traffic conditions*. With other data sources, the traffic conditions are described.

Section 2.3, shows the case study (region and period) that is used for the design and the evaluation of the prediction method.

2.1 Speed measurements: TomTom HDFlow

TomTom HDFlow is a data source from TomTom that provides real-time road traffic information based on floating car data. The road traffic information is provided every minute for many road segments. These road segments have a high coverage of urban roads and freeways world-wide. In this study, only the roadway speeds are used from this source. This data is aggregated to intervals of 5 minutes before it is analyzed and used in this study.

TomTom uses road segment lengths of approximately 200 meters to 3000 meters on highways in The Netherlands. These road segments are defined in *openLR*. This is an open source dynamic form of location referencing. Moreover, every segment contains a *Functional Road Class (FRC)* and a *Form Of Way (FOW)* value. These values are used to classify segments. FRC represents the size and type of utility of the roads, such as freeways, local roads or bicycle paths. FOW describes the physical form of the road such as motorways, slip roads and roundabouts.

This research focuses on freeways (FRC 0 and 1), because these type of roads contains the most . However, also other major roads (FRC 2, 3 and 4) show temporal correlations between ascending measurements errors compared to the average.

The roadway speeds are determined every minute by TomTom using real-time aggregated vehicle travel times. TomTom do not mention how they calculates the roadway speeds. Moreover they do not quantify the reliability of the data source. According to Aarts et al. (2015), the proportion of drivers that contribute to the TomTom dataset is approximately 4% on freeways and 1% on highways. Although this proportion seems limited, they concluded that the given speeds per segment represent the actual speeds well.

2.1.1 Data exploration: Right-censored measurements

The speed data is right censored which means that values above some *boundary speed* value, are not registered. For the completeness, Simacan replaced this unknown values with the *freeflow speed*. The freeflow speed is the speed that vehicles regularly drive in *freeflow conditions* on a specific segment. This freeflow speed is constant over time, but location dependent. An example of these measurements is shown in Figure 1. This is a timeseries of a segment at Barneveld in East direction.



Figure 1: A timeserie of speed measurements at Tuesday September 4th, 2018 at the A1 near Barneveld in East direction. The dotted line represents the boundary speed of this segment

The boundary speed is dependent per segment. For freeways, this boundary speed is always lower or equal to 80 km/hour in the used data samples (see Figure 2). In this study, these unknown measurements, which Simacan replaces with freeflow speed, are dealt as $x \in \mathbb{R}^+$: $v_{boundary} < x \le v_{freeflow}$ where $v_{boundary}$ is the boundary speed and $v_{freeflow}$ is the freeflow speed.

2.1.2 Data preparation: making 5 minute aggregates

The measurements are aggregated to 5-minute intervals, because more detailed measurements are not needed for the aim of this study. An additional advantage is that aggregates make sure that noise is reduced. If some measurements exist below the boundary speed in the interval, only these measurements are used in the aggregation, because of the uncertainty of the freeflow speed measurements (see equation 1). This results in that the calculated speed aggregate might be only based on one measurement if the other measurement values are above the



Figure 2: The distribution of the speed measurements of the data set in The Hague

boundary speed. It is assumed that this will not cause high variations, because the measurements are already edited by TomTom.

$$\hat{v} = \begin{cases} \sum_{t=0}^{T} \frac{v_t}{T} \ \forall \ t \in T : v_t \neq v_{freeflow}, \ \text{if} \ \exists \ v_t \neq v_{freeflow} \\ \sum_{t=0}^{T} \frac{v_t}{T} \ \forall \ t \in T, \ \text{otherwise} \end{cases}$$
(1)

With \hat{v} is the speed aggregate per 5 minutes, v_t is the measured speed of time t, T is the set of all measurements in the 5-minute interval of the floating car data and $v_{freeflow}$ is the speed during freeflow.

However, Simacan stored the data in compressed files. Each file contains the data of all locations of one minute. Therefore, the needed data for making a timeseries of one segment, must be searched and extracted from many files. This results in low performance times which makes the analysis impossible. Therefore, the data is processed and stored differently. Details about this processing can be found in Appendix C: Data Processing of Speed measurements: TomTom HDFlow.

2.2 Irregular situations

Traffic speeds are influenced by traffic conditions. In this research, two conditions are distinguished: regular traffic conditions and irregular traffic conditions. These

traffic conditions consist of repetitive patterns. Regular traffic conditions do occur every week at the same time at the same location. Irregular traffic conditions are conditions that are caused by situations. These situations are repetitive, but do not occur at the same moment and at the same location. Therefore, they are called *irregular situations*. In this research, the irregular situations consist of:

- accidents
- roadworks
- weather: snow, icing and extreme weather alerts
- holidays

The locations and moments that are impacted by these irregular situations are not further considered by the calculation of the long-term prediction during regular conditions, because these conditions causes uncommon speed values which can influence the results.

2.2.1 Accidents and roadworks: NDW situations and RWS-ADY

For the accidents and roadworks on highways, information about the roadworks and accidents from the *National Data Warehouse for Traffic Information (NDW)* is used. This includes all accidents on highways that are reported or noticed by the Dutch road authority. Also, the planned and active roadworks on highways are available in this dataset.

Many information is available about the accidents and roadworks. However, not all information is available for every accident and roadwork, and if it is available, the information is not always correct. The relevant types of information are converted to *AlertC* codes by Simacan. The AlertC codes that exist in the dataset are allocated to three classes: accidents, roadworks, other situations that affect the traffic situation and non-relevant AlertC codes (see Appendix D: AlertC codes indicating irregular situations).

The locations of the NDW situations always contains a coordinate. However, these coordinates are difficult to map to the segments. Almost all situations on highways also contain another Dutch location referencing called *Traffic Information Location Database (VILD)*. These codes can be converted so that it can be matched to the OpenLR location reference of the speed measurements.

A disadvantage of these datasets is that the reliability of the number of closed lanes during accidents is low. To check if these numbers are correct, matrix sign information on highways retrieved by *Rijskwaterstaat Acquisition Dynamic Data*

(*RWS-ADY*) is used. For the analysis of the impact of accidents, all accident information is removed for which the number of lanes is not identical to the number of matrix signs that display if lanes are closed. All highways in the used datasets have matrix signs. If these will not be available, the intensities from the detection loops could be used to investigate if lanes are closed.

2.2.2 Weather

Weather information is collected from the *Koninklijk Nederlands Meteorologisch Instituut (KNMI)*. These include measurements when icing and snow occur. For region Amersfoort, information from weather station De Bilt is used. For region The Hague, information from weather station Rotterdam is used.

It is investigated which weather information influences traffic conditions on a national level so that the weather information can be used for several regions. Weather conditions, specifically road surface conditions and visibility, can influence the travel speed.

Days with snow, icing and extreme weather alerts (code orange and red) are categorized as irregular situations. These weather circumstances influences traffic conditions by changed road surface conditions. However, visibility and heavy rain are not categorized as irregular situations, because it can vary from road to road and the coverage of the weather data is not high enough to say something on road level. An extensive analysis of the weather data can be found in Appendix A: Weather impact.

2.2.3 Holidays

Holidays for which officials do not have to work and pupils are free are categorized as irregular conditions.

2.3 Case study

Two datasets are used in this study. One set is used for the analysis of the variations and the design of the prediction method. The other set is used for the evaluation of the prediction method. The used regions and time periods are shown in Table 6.

Purpose	Region	Period
Baseline prediction	-	1 July 2018 until 30 April 2019
Analysis and design of short-term prediction	Amersfoort	1 September 2018 until 30 Decem- ber 2019
Evaluation	The Hague	1 January until 30 April 2019

Table 6: The regions and periods for each used set

2.3.1 Regions

It is expected that the used locations are representative for a large part in the Netherland that are relatively congested. The prediction method only can be well evaluated if many measurements are known. Therefore, the used datasets must contain congested conditions so that it contain many speed measurements below the *boundary speed*. As an example, the number of known measurements, expressed in hours per day are visualized for the region Amersfoort (see Appendix B: Visualisation of location characteristics).

The region Amersfoort is used for the analysis of the variations and the design of the prediction method (see Figure 3). The region The Hague is used for the evaluation of the prediction method (see Figure 4).

Fortunately, some segments exists for which relatively many measurements are known. These segments are used for analysis of the variations in speeds. Moreover, it can better demonstrate the influences of measurements on the predicted speeds. The segment used in this research is the road section on the A1 near Barneveld in the west direction during the morning peak (see Figure 5).

2.4 Conclusion of the data

TomTom HDFlow speed measurements are used in this research for the prediction of speeds. This data source uses aggregates of vehicle travel times to determine roadway speeds.

TomTom HDFlow is used, because it has a high coverage of urban roads and freeways world-wide. Moreover, it is expected that these roadway speeds represents the actual speeds well. However, a disadvantage of this data source is that the measurements are right-censored: all roadway speeds around freeflow conditions are unknown.

Above that, data about irregular situations are used to seperate the speed measurements in regular and irregular conditions. The data source of the irregular



Figure 3: Used roads for region Amersfoort. The red lines are freeways which are used for the calibration and validation and black are other roads for which data is available.

situations are: information about accidents and roadworks (NDW situations), lane closures from matrix signs (RWS-ADY), weather information (KNMI) and holidays.

The case study focuses on highways in Amersfoort (The Netherlands) for the analysis of the variations and the design of the prediction method. Next, highways in The Hague are used for the evaluation of the prediction method.



Figure 4: Used roads for region The Hague. The red lines are freeways which are used for testing the final prediction and black are other roads for which data is available.



Figure 5: Used road section for the examples of the correlation results and the prediction method and results

3 Approach

This chapter describes the method used to predict vehicle speeds. Moreover, it describes the method for the evaluation of the accuracy of the prediction.

Two prediction methods are designed: one for the *long-term prediction*, called the *baseline prediction*, and one for the short-term prediction. These methods give two results: an average expected speed, and an interval in which the speed lies for a particular chance.

The long-term prediction is based on the idea that traffic conditions have repetitive patterns; the conditions for a particular location, time of the day, day of the week are almost identical during regular circumstances. Moreover, it is assumed that seasonal variations in these conditions occur for the long-term prediction.

Vehicle speeds are a measure of the traffic conditions. The current variation in speed gives an indication of the variation of the future speed. However, still variations in speeds occur while the traffic conditions are constant. These variations are not predictable and are called *noise*. Short-term variations in vehicle speeds are analyzed whether they say something about the next speeds. Moreover, the effects of irregular situations on vehicle speeds are analyzed.

With the knowledge of the behaviour of these variations, a *short-term prediction* method is designed. This method can give more accurate predictions than the long-term prediction by using the actual speed measurements, and additionally, information about the current situation.

Finally, a method is presented that shows the accuracy of the prediction method.

3.1 Long-term prediction of vehicle speeds for regular conditions

Used symbols for the long-term prediction:

$p_{s,t}$	the probability that the censored measurements are known, depen-
	dent of segment s and time t
$v_{b,s}$	the boundary speed, speed measurements higher than this value
	are unknown, dependent on segment s
$v_{ff,s}$	the speed that is driven under freeflow conditions, dependent on
	segment s
V_{low}	the speed assuming that all unknown measurements are equal to
	the boundary speed

V_{high}	the speed assuming that all unknown measurements are equal to the freeflow speed
$E(V_{s,t})$	the average of the expectation of the speed, dependent on segment s and time t
$\mu_{low,s,t}$	the average of the speed assuming that all unknown measurements are equal to the boundary speed, dependent on segment s and time t
$\mu_{high,s,t}$	the average of the speed assuming that all unknown measurements are equal to the freeflow speed, dependent on segment s and time t
$var(V_{s,t})$	the variance of the expectation of the speed, dependent of segment \boldsymbol{s} and time \boldsymbol{t}
$\sigma_{s,t}$	the standard deviation of the expectation of the speed, dependent on segment s and time t
$\sigma_{low,s,t}$	the standard deviation of the speed assuming that all unknown measurements are equal to the boundary speed, dependent on segment s and time t
pdf cdf	the probability density function of the normal distribution the cumulative density function of the normal distribution

The baseline prediction is the long-term prediction of vehicle speeds for regular conditions. It is based upon historical data, and it can give a prediction for any location at any time independent of the size of the time horizon.

The baseline prediction does not make use of real-time measurements, because the current condition does not impact the vehicle speeds for a large time horizon ahead.

This data is classified based on two characteristics:

- Per 5-minute time interval per day of the week (temporary)
- Per road segment (spatially)

For all these clusters, the average speed, μ , and the size of the variation in speeds, σ^2 , are determined. From these values, the speed interval is determined in which the speed lies for a particular chance:

$$interval = \mu \pm 2\sigma$$

These values only change if the network changes in such a way that it impacts the speeds. Therefore, these values only have to be calculated once and can be saved in a database.
However, many values are unknown due to right-censored measurements (see Section 2.1.1 for explanation), which makes it challenging to determine the real average speed and size of the variation.

3.1.1 Determination of the variation in speed

The variance of the baseline prediction can not directly be calculated because the measurements are right-censored. Therefore, it is approximated from the assumption that the speeds in a cluster are normally distributed (see Figure 6).

For the determination of the variance, two sample standard deviations are calculated and then combined with a weighted average. The first one, σ_{low} , is the sample standard deviation of the speeds whereby all unknown speed measurements are replaced by the boundary speed.

The second sample standard deviation is approximated by calculating the interval, which contains 95% of the measurements. This equals 4σ . It is assumed that the interval is between the 2.5th percentile of V_{low} and the speed during freeflow conditions. This is illustrated in Figure 6.



Figure 6: Calculation of the variance of the speed assuming that all measurements are equal to the speed in freeflow condition

Finally, these two sample standard deviations are combined with a weighted average, which results in Equation 2. The weights are defined as the number of known measurements. This is chosen because, if many speed measurements are unknown, it is assumed that most of these are freeflow speed. If few speed measurements are unknown, most of these are boundary speed.

$$var(V) = \left((1-p) \cdot \sigma_{high} + p \cdot \sigma_{low} \right)^2$$

$$= \left((1-p) \frac{v_{ff} - (\mu_{low} - 3\sigma_{low})}{4} + p \cdot \sigma_{low} \right)^2$$
(2)

3.1.2 Determination of the average speed

Two different methods are investigated for the determination of the average speed per cluster: one using the assumption that the speeds follow a normal distribution under regular conditions per cluster and one for which the speeds do not follow a normal distribution.

It is assumed that the speeds can follow a normal distribution if the average, $E(V_{s,t})$, is within some interval. The lower limit of the interval is defined as the average of the speeds for which all unknown measurements are replaced by the boundary speed ($\mu_{low,s,t}$). The average of the speeds for which they are replaced by the freeflow speed represents the upper limit ($\mu_{high,s,t}$). This is shown in Equation 3.

$$\mu_{low,s,t} < E(V_{s,t}) < \mu_{high,s,t} \,\forall \, s, t \tag{3}$$

This method is chosen over a statistical test, because most clusters only contain a few known measurements, which do not give significant results with statistical tests.

Assumption that speeds follow a normal distribution Using the assumption that the variations of the measurements follow a normal distribution, the baseline equations (the average and the variation per segment per time) that are constructed are shown in equation 5 and 2. In these equations, all variables can be determined from the historical measurements. The extended foundation of these equations can be found in Appendix E: Extended foundation of the baseline prediction.

$$\sigma_{s,t} = \sqrt{Var(V_{s,t})} \tag{4}$$

$$E(V_{s,t}) = p_{s,t} \cdot \left(\frac{\mu_{low,s,t}}{p_{s,t}} + \sigma_{s,t} \cdot \frac{pdf(cdf(p_{s,t}))}{cdf(p_{s,t})}\right) - \frac{(1-p_{s,t})}{p_{s,t}}v_{b,s}\right) + (1-p_{s,t})v_{ff,s}$$
(5)

The disadvantage of the method for the calculation of the average is that it could not well construct a distribution if fewer measurements are known. Therefore, it is assumed that the average of the baseline prediction is the freeflow speed if less than 15% of the measurements are known.

assumption that speeds do not follow one normal distribution It this method it is assumed that the clusters follow one of the two normal distributions. However, from the data, it could not be concluded which measurements follow the distribution for freeflow conditions and which follow the distribution for congested conditions.

Therefore, the average is not calculated by the estimation of the normal distributions. To calculate the baseline prediction, first, a weighted average is made with:

- Average of V_{low}
- Average of V_{high}

The weights are proportional to the number of known values. This is with the same reason as with the calculation as the variation in speeds in Section 3.1.1.

This results in Equation 6 for which the size of the weights are equal to p. The variation is the same as for hypothesis A, which is shown in Equation 2.

$$E(V_{s,t}) = p_{s,t} \cdot E(V_{low,s,t}) + (1 - p_{s,t}) \cdot E(V_{high,s,t})$$
(6)

3.1.3 Analysis of long-term variations in vehicle speeds

It is assumed that repetitive patterns exist in the variations of speed measurements compared to the baseline speed. In this study, variations of measurements compared to the baseline are also called *residuals*. If these repetitive patterns exist, the baseline prediction can be updated.

Repetitive patterns, so *systematic variation* occurs in the residuals, can be proven with auto-correlation. If the residuals are correlated, this indicates that the residuals do have systematic variation. The following factors are investigated on correlation:

- Residuals of successive days
- Residuals of successive weeks

Many of these residuals are explained by short-term systematic variations. To account for that, hourly averages of these residuals are made. The auto-correlation of the residuals is calculated for a fixed hour of the day. For example, 8 till 9 o'clock on Monday is compared with the same time interval on Tuesday and Wednesday. Above that, unknown values caused by the censored measurements must be replaced to calculate the auto-correlation. The average value of the residual replaces the unknown values.

The systematic variations in the residuals can be explained by increasing or decreasing values (trend) or by repeating short-term cycles (seasonality). It is investigated which type of systematic variations occur, by plotting the timeseries of the residuals of successive days and weeks. Two different timeseries are visualized. In the first timeserie all unknown speeds are removed. However, the disadvantage of removing all unknown speeds is that only congested conditions are compared. Therefore, it does not show the amount of congestion compared to freeflow conditions. Therefore, another timeserie is visualized for which all unknown speeds are replaced by the freeflow speed.

It is also possible that correlations between successive years occur. However, we do not have data available to find these correlations. It is decided to investigate the correlations based on hourly or daily averages, because these averages are less influenced by noise.

3.2 Analysis of short-term variations in vehicle speeds

The short-term variations of the residuals are analyzed for freeways in the region Amersfoort. The residuals are the measurements compared to the baseline prediction. These variations are investigated, because it is assumed that they consist of repetitive patterns. If repetitive patterns exist, the prediction can be improved for a specific time horizon with the expected variations in the residuals. These expected variations are determined from the knowledge of the current traffic situation and the current variations in the residuals. Two types of variations are distinguished: variations during regular conditions and variations during irregular conditions.

Sometimes variations in the residuals occur for which the cause of these variations is known (for example accidents, roadworks, weather, holidays). In this study, these conditions are called irregular. However, mostly, variations of the residuals do occur for which the cause is unknown. These are called regular. It is assumed that both conditions consist of repetitive patterns. However, they are investigated separately, because it is assumed that the repetitive patterns are different in terms of pattern duration and the recurring period.

3.2.1 Analysis of variations under regular conditions

Always variations in the residuals occur. However, the cause of the variations is often unknown. Nevertheless, it is still possible to find recurrent patterns in the residuals. The existence of the *temporal* and *spatiotemporal* recurrent patterns in these variations have been investigated in this study. These recurrent patterns have been demonstrated by the calculation of the autocorrelation of the residuals.

Temporal variations Temporal variations are the variations in adjacent times at the same location. It is investigated if repetitive patterns in the variations of the residuals occur up to a time horizon of 3 hours at a specific location.

Only the measurements during the morning rush hour are investigated for the calculation of the correlation. In this period the most measurements are known, because the most congestion occurs. In this case, the morning rush hour is between 6:30 and 9:30 hour. If also the other moments would be considered, the unknown measurements must be guessed, which results in untrustworthy correlations.

Spatiotemporal variations Spatiotemporal variations are the variations in adjacent locations and next times upstream these locations. This analysis is performed, because many variations on freeways are a consequence of the propagation of disruptions (so-called traffic jams). The propagation of these disruptions travels upstream the road with a certain speed. If the road is fully occupied, these disruptions have an almost constant speed (Cassidy and Windover, 1995; Kerner and Rehborn, 1996). That is why the variations of the residuals are analyzed at a constant speed in the upstream direction.

The constant speed value is determined with the lag, which has the highest crosscorrelation.

With the same reason as with the temporal variation, the correlations are only calculated for the morning rush hour. Also, it depends on the location how many times congestion occurs during the morning rush hour. Therefore, for this analysis, a segment is used that has much congestion.

Moreover, it is location-dependent if recurrent patterns occur. Congestion propagates dissimilar for different locations in the network such as roads near intersections. A big chance exists that if congestion occurs at a segment, the congestion propagates upstream to the adjacent segment. However, if an intersection is on the adjacent segment, it is possible that the congestion stops, because the adjacent segment is an intersection which influences the traffic conditions. Therefore, the correlation of the known residuals between adjacent segments and next times is calculated.

3.2.2 Analysis of variations under known irregular conditions

Variations in the residuals can occur because of a irregular situation. Irregular conditions are conditions for which the situation do not occur every week at the same time at the same location (see Section 2.2 for more information). However, it is assumed that the variations, caused by these situations, also show repetitive patterns; independent of the time and location.

The duration of the patterns is analyzed by determining if variations are occurring. In this study, it is assumed that variations are occurring caused by the situation if the chance the measured speed value occurs is less than 2.5% in regular conditions. The baseline prediction defines the regular conditions.

The variations caused by accidents are the only irregular conditions that are analyzed in this study, because only for these situations, enough data is available.

Accidents It is assumed that accidents caused variations in the residuals, by way of congestion. This congestion propagates in time and space. Therefore, these variations are investigated for two characteristics:

- The duration that the variations in the residuals occur at the location of the accident (time);
- The distance that the variations in the residuals occur (space).

It is expected that the size of the variations in the residuals mainly depends on the road occupancy and the size of the accident. However, the information about the occupancy of the road and the size of the accident is not directly available. Therefore, the effects are grouped dependent on two different characteristics:

- The relative number of closed lanes;
- The speed of the baseline prediction compared to the freeflow speed.

The relative number of closed lanes indicates the size of the capacity reduction. This reduction influences the size of the impact of the situation. The relative number of closed lanes can be found in the accident information provided by NDW. With the matrix sign information from RWS-ADY, this number of closed lanes is checked, because this is not always correct. All accidents for which the accident information of NDW is not equal to the matrix sign information are not further considered.

The relative error of the baseline prediction compared to the freeflow speed indicates the regular road occupancy. The occupancy is high if the baseline prediction is significantly lower than the freeflow speed. High occupancy is more likely to cause congestion than a low occupancy.

3.3 Short-term prediction of vehicle speeds

The short-term prediction is a method that predicts vehicle speeds by updating the long-term prediction with the real-time speed measurements and by the knowledge of the current traffic situation. With the knowledge of the recurrent patterns of the variations in the residuals from Section 3.2, the long-term prediction can be improved to give more accurate predictions for regular and irregular conditions.

In this study, two challenges are tackled that follow from predicting vehicle speeds using real-time measurements:

- The influence of noise in the real-time measurements;
- The influence of the real-time measurements and historical measurements on the prediction of multiple time steps ahead.

The Kalman filter is used to determine the noise of the measurements. However, the original filter can only predict one timestep because it uses the previous measurement. For the prediction of multiple time steps, the previous measurement must be presumed, because it is unknown. These future measurements are presumed from a combination of the baseline prediction and the previous known real-time measurements of the current location and the locations upstream.

The method could be used for different regions because it takes characteristics between adjacent locations into account. Moreover, in the method is added that it predicts congestion for a particular time if it is known that an accident occurred. The method can also be used for other irregular situations, but it always returns to the baseline prediction after a particular time.

Used symbols for the short-term prediction:

h	the maximum steps ahead, expressed per 5 minutes
$\sigma_{v,base}$	standard deviation of the baseline speed, dependent on time and
	segment
$\sigma_{v,pred}$	standard deviation of the predicted speed, dependent on time and
	segment
$\mu_{v,base}$	average of the baseline speed, dependent on time and segment
$\mu_{v,pred}$	average of the predicted speed, dependent on time and segment

$\mu_{v,type}$	average of the speed during irregular conditions, dependent on type
	of cause, known consequence and road class
v_b	boundary speed, speed measurements higher than this value are
	unknown, dependent on segment
v_{ff}	speed during freeflow conditions, dependent on segment
v_m	speed measurement, expressed per 5 minutes, dependent on time
	and segment
$r_{k k-1}$	correlation between segment k and segment k-1, dependent on
I	segment
$\hat{x}_{k k-1}$	a priori state for time interval k
$x_{k k}$	a posteriori state for time interval k
F	the state transition value
w_k	the process noise or disturbance for time interval k
y_k	the measurement value for time interval k
P_k	the uncertainty of the state for time interval k
K_k	Kalman gain for time interval k
Q_k	Variance of the model noise for time interval k
R_k	Variance of the measurement noise for time interval k
$R_{spat,k}$	Covariance of the measurement noise, dependent on segment k
• /	and segment $k-1$
Ι	identity matrix
C_{R1} ,	Constants that must be calibrated
C_{R2} ,	
C_{var1} ,	
C_{var2} ,	
$C_{ ho}$	

3.3.1 Original Kalman filter

The purpose of a Kalman filter is to remove the noise of a timeseries with a relatively fast method. A requirement is that this noise is normally distributed. The Kalman filter is based on a state estimate and a state uncertainty. In every iteration k, this state is updated with a measurement value and measurement uncertainty. This new state represents the state (in this study: vehicle speed) without the noise. The Kalman filter exists of five main equations.

Every iteration over k, the state, x, will be estimated from the previous estimate, the state transition and the process noise.

$$x_{k|k-1} = F \cdot x_{k-1|k-1} + w_k \tag{7}$$

Moreover, the uncertainty of the state, P, will be estimated from the previous

estimate, the state transition and the process noise. F^T is the transposed matrix of F.

$$P_{k|k-1} = F \cdot P_{k-1|k-1} \cdot F^T + Q \tag{8}$$

The Kalman gain represents the weight of the measurement and the weight of the previous estimate to make a new estimate. Next, the Kalman gain will be estimated from the state uncertainty and the variance of the measurement noise.

$$K_{k} = \frac{P_{k|k-1}}{P_{k|k-1} + R_{k}}$$
(9)

The old estimate will be updated with the new measurement and the Kalman gain.

$$x_{k|k} = x_{k|k-1} + K_k \cdot (y_k - x_{k|k-1}) \tag{10}$$

Finally, the uncertainty of the state will be updated with the Kalman gain.

$$P_{k|k} = P_{k|k-1} \cdot (I - K_k) \tag{11}$$

3.3.2 Initialization of the Kalman filter

The state and the variance of the state must be initialized. In this study, the state x is the relative difference of the speed estimate compared to the baseline prediction for a particular segment and time.

For iteration 0, the state equals the relative difference of the speed measurement compared to the baseline prediction for segment 0 and time 0. Due to the censored measurements, many speeds above a certain boundary, are unknown. These unknown speeds are replaced by the baseline speed if the baseline speed is higher than the boundary speed. If the baseline speed is lower than the boundary speed, these unknown speeds are replaced by the boundary speed.

$$x_{0,temp} = x_{0,spat} = \begin{cases} v_{m,0}/\mu_{v,base,0}, & \text{if } v_{m,0} \le v_b \\ 1, & \text{if } v_{m,0} > v_b \cup \mu_{v,base,0} > v_b \\ v_b/\mu_{v,base,0}, & \text{if } v_{m,0} > v_b \cup \mu_{v,base,0} \le v_b \end{cases}$$
(12)

The variance of the state is assumed as the square root of the speed of the baseline prediction.

$$P_{0,temp} = P_{0,spat} = \sqrt{\mu_{v,base,0}} \tag{13}$$

3.3.3 Adaptation of Kalman filter for speed prediction

The Kalman filter is adapted to predict speeds for all locations, including multiple time steps in the future, and is able to predict for all situations.

In the used Kalman filter, two states are defined. One represents the state caused by the temporal variation (*temp*), and one represents the state caused by the spatiotemporal variation (*spat*).

From the original Kalman filter, the state transition value F is assumed to be 1 and the process noise w is assumed to be 0. This results in that these parameters could be removed from the equations. Q is a constant that is guessed in this study based on the intuition gathered from the analysis of the variations in speeds. This results in:

$$\hat{x}_{k|k-1,temp} = \hat{x}_{k-1|k-1,temp}$$
 (14)

$$\hat{x}_{k|k-1,spat} = \hat{x}_{k-1|k-1,spat}$$
 (15)

$$P_{temp_{k|k-1}} = P_{k-1|k-1} + Q_{temp}$$
(16)

$$P_{spat_{k|k-1}} = P_{k-1|k-1} + Q_{spat}$$
(17)

The variation of the measurement noise R, denotes the chance the variation is caused by noise. How larger R, how fewer the next state will be updated towards the measurement. It is assumed that the variation in the measurement noise is the *n*th root of the absolute speed at the model state. The index *n* is defined by some constant C_{R1} . The absolute speed at the model state can be calculated by multiplying the average of the baseline prediction μ_{base} by the model state *x*.

$$R_{k,temp} = \sqrt[C_{R_1}]{\hat{x}_{k|k-1,temp} \cdot \mu_{v,k,base}}$$
(18)

Location-dependent predictions As already mentioned in Section 3.2.1, it depends on the location, if the variation in the residuals of a specific location influences the adjacent location upstream. Therefore, the spatiotemporal correlation ρ is added to the measurement noise *R*. The size of the correlation indicates the chance the variation of a road influences the adjacent road. The equation is adjusted in such a way so that for small correlations, the measurement noise becomes large and vice versa.

$$R_{k,spat} = \frac{\sqrt[C_{R2}]{\hat{x}_{k|k-1,spat} \cdot \mu_{v,k,base}}}{\rho_{k|k-1}^{C_{\rho}}}$$
(19)

Estimating the current state In the original Kalman filter, the state is updated with the new measurement. Due to censored measurements, not all measurements are known. In this method, these unknown measurements are assumed as freeflow speed.

$$y_{k,temp} = \begin{cases} v_{m,k}/\mu_{v,base,k}, & \text{if } v_{m,k} \le v_b \\ v_{ff}/\mu_{v,base,k}, & \text{if } v_{m,k} > v_b \end{cases}$$
(20)

$$y_{k,spat} = \begin{cases} v_{m,k}/\mu_{v,base,k}, & \text{if } v_{m,k} \le v_b \\ v_{ff,k}/\mu_{v,base,k}, & \text{if } v_{m,k} > v_b \end{cases}$$
(21)

The visualization of the behaviour of the Kalman filter to estimate the new state, can be found in Appendix F: State estimation of the Kalman Filter.

Estimating the state multiple time steps ahead The original Kalman filter could only predict for one timestep. Therefore, the Kalman filter is adapted so that it can be used for prediction purposes. Therefore, the measurements are estimated from the combination of several timeseries.

During prediction, the speed measurements are unknown. Every iteration, these measurements are estimated from a combination of the baseline prediction, the state estimate of the temporal variations and the state estimate of the spatiotemporal variations.

Combining several timeseries The prediction of the measurement is based on the inverse-variance weighted mean for normal distributions (see Equation 23). An assumption that is made is that the baseline prediction is independent on the temporal prediction and the spatiotemporal prediction. Also, the covariance between the temporal prediction and the spatiotemporal prediction is independent on time and location.

$$V_{base}, V_{temp}, V_{spat} \sim N(\mu, \sigma^2)$$
 (22)

$$E(V_k) = \frac{1/\sigma_{v,base,k}^2 \cdot \mu_{v,base,k} + 1/\sigma_{v,temp,k}^2 \cdot \mu_{v,temp,k} + 1/\sigma_{v,spat,k}^2 \cdot \mu_{v,spat,k}}{1/\sigma_{v,base,k}^2 + 1/\sigma_{v,temp,k}^2 + 1/\sigma_{v,spat,k}^2}$$
(23)

$$Var(V_k) = \frac{1}{1/\sigma_{v,base,k}^2 + 1/\sigma_{v,temp,k}^2 + 1/\sigma_{v,spat,k}^2}$$
(24)

The estimated measurement becomes:

$$\mu_{v,k,pred} = \frac{\left(\frac{1}{\sigma_{v,k,base}^2} + \frac{\hat{x}_{k|k-1,temp}}{P_{k|k-1,temp} + k \cdot C_{var1}} + \frac{\hat{x}_{k|k-1,spat}}{P_{k|k-1,spat} + k \cdot C_{var2}}\right) \cdot \mu_{v,k,base}}{\frac{1}{\sigma_{v,k,base}^2} + \frac{1}{P_{k|k-1,temp} + k \cdot C_{var1} + \frac{1}{P_{k|k-1,spat} + k \cdot C_{var2}}}$$
(25)

The visualization of the behaviour of this method to combine the predictions are shown in Appendix G: The estimation of the short-term prediction by combining several timeseries for several timesteps.

Irregular conditions The long-term prediction method determines the variation under regular conditions. For irregular conditions, it is assumed that the variance could be higher. Therefore, the minimum variance of the baseline prediction is $\sigma_{v.k.base}^2 = v_{ff}$.

In the method, only accidents are considered. However, other irregular situations could be added. If it is known that an accident happened, the estimated measurements are a constant value depending on the type of the situation $\mu_{v,type}$. The value is the average speed that is driven during that type of situation and followed from the analysis of variations under known irregular conditions in Section 3.2.2.

The relative difference of the speed measurement compared to the baseline prediction, becomes:

$$y_{k,temp} = \begin{cases} \mu_{v,k,pred}/\mu_{v,k,base}, & \text{if } v_{m,k} \text{ is in the future during regular conditions} \\ \mu_{v,type}/\mu_{v,k,base}, & \text{if } v_{m,k} \text{ is in the future during irregular conditions} \end{cases}$$
(26)

$$y_{k,spat} = \begin{cases} \mu_{v,k,pred}/\mu_{v,k,base}, & \text{if } v_{m,k} \text{ is in the future during regular conditions} \\ \mu_{v,type}/\mu_{v,k,base}, & \text{if } v_{m,k} \text{ is in the future during irregular conditions} \end{cases}$$
(27)

The Kalman gain, the update of the estimate of the state and the update of the uncertainty of the state remains unchanged:

$$K_{k,temp} = \frac{P_{k|k-1,temp}}{P_{k|k-1,temp} + R_{k,temp}}$$
(28)

$$K_{k,spat} = \frac{P_{k|k-1,spat}}{P_{k|k-1,spat} + R_{k,spat}}$$
(29)

$$x_{k|k,temp} = \hat{x}_{k|k-1,temp} + K_{k,temp} \cdot (y_{k,temp} - \hat{x}_{k|k-1,temp})$$
(30)

$$x_{k|k,spat} = \hat{x}_{k|k-1,spat} + K_{k,spat} \cdot (y_{k,spat} - \hat{x}_{k|k-1,spat})$$
(31)

$$P_{k|k,temp} = P_{k|k-1,temp} \cdot (I - K_{k,temp})$$
(32)

$$P_{k|k,spat} = P_{k|k-1,spat} \cdot (I - K_{k,spat})$$
(33)

3.3.4 Used values for the constants

In the Kalman filter, several constants must be determined in order to give accurate predictions. In this study, these constants are approximated by the knowledge from the analysis of the variation in the residuals. Then, they are manually adapted so that the variations in the residuals of the prediction visually behaves like expected.

The used constants in the Kalman filter are:

$Q_{k,spat}$	$200 \ km^2/h^2$
$Q_{k,temp}$	$100 \ km^2/h^2$
C_{R1}	2
C_{R2}	1
C_{var1}	5
C_{var2}	15
$C_{ ho}$	2

3.3.5 Final predicted speed

The iteration stops if k equals the set prediction horizon. The average predicted speed and the variation of the prediction speed is:

$$E(v) = x_{k|k-1} \tag{34}$$

$$Var(v) = \frac{1}{\frac{1}{\sigma_{v,k,base}^2 + \frac{1}{(P_{k|k-1,temp} + k \cdot C_{var1}) + \frac{1}{(P_{k|k-1,spat} + k \cdot C_{var2})}}}$$
(35)

3.4 Evaluation of the accuracy of the vehicle speed predictions

The short-term and long-term prediction is tested on their accuracy of the predicted speed averages and the predicted speed intervals. The performance is evaluated with two different measures: the Mean Absolute Percentage Error and the likelihood that measurements lie in the predicted interval.

To manage that the presented prediction accuracies could be achieved independent of the region, the evaluation is performed on another region and time than the analysis. These datasets are described in Section 2.3.

To manage that the prediction accuracies hold for all conditions and for all prediction horizons, the accuracy of the predictions is investigated for different categories: prediction horizons, time periods and traffic conditions (see Section 3.4.2).

3.4.1 Measures

Two different measures are used to evaluate the performance of the designed prediction methods. The predicted average is compared with other prediction methods using the Mean Absolute Percentage Error (MAPE). Furthermore, the likelihood that the measurement lies in the predicted interval is calculated.

MAPE of the predicted average For companies providing home deliveries, the punctuality of the arrival times is important. This punctuality can be achieved by minimizing the relative error of the average speed prediction compared to the actual speed measurements. Therefore, the MAPE is used, which is a measure for the determination of the absolute relative error. The MAPE is calculated with equation 36, where v_m is the speed measurement and E(v) is the predicted speed.

$$MAPE = \frac{1}{n} \sum_{t=1}^{n} \left| \frac{v_m - E(v)}{v_m} \right| \cdot 100\%$$
(36)

With the MAPE, the average of the baseline prediction and the short-term prediction is compared with other prediction methods to see if the designed prediction method is an improvement. The other predictions for which the predicted average speed is compared are:

- TomTom Speed Profiles
- Instantaneous prediction

The *TomTom Speed Profiles* give, like the baseline prediction, a long-term speed prediction for all locations, all times of the day and for each day of the week, independent of the current traffic conditions. Further details are explained in

Chapter 2. The *instantaneous prediction* denotes that the speed measurement occurring at the moment is used as the expected speed.

The evaluation of the accuracy of the prediction during freeflow conditions is difficult, because the measurement during freeflow conditions are unknown. Therefore, the prediction is evaluated using different scenario's, whereby the unknown speeds are replaced with: freeflow speed, the baseline prediction, and the Tom-Tom Speed Profiles.

Likelihood of the predicted interval The uncertainty of the prediction can be shown by predicting an interval in which the speeds lies. With the likelihood, the percentage is calculated that the predicted interval is correct.

For known measurements, the interval is correct if the measurement is inside the interval. For unknown measurements, only an interval is known in which the measurements lie. In this case, the interval is correct if the measurement may be inside the interval, so the intervals intersect.

3.4.2 Categories for evaluation

Time horizon The prediction averages and prediction intervals are evaluated up to a time horizon of 90 minutes in steps of 5 minutes. For longer time horizons, it is assumed that under regular conditions, the actual condition does not influence the prediction, so the short-term prediction does not give more accurate predictions than the baseline prediction. The accuracy of the baseline prediction is constant over the predicted time horizon, because the baseline prediction is independent of the actual situation.

Traffic conditions Traffic conditions influence traffic speeds. Mostly, the conditions depend on the location and are recurrent in time (for example, every Monday morning, a rush-hour exists). Sometimes, traffic conditions are influenced by irregular situations such as accidents.

The baseline prediction gives an indication of the recurrent traffic conditions. The average speeds of this prediction represent if the conditions for a particular location and time are regular congested or are regular freeflow. The congested conditions are if the speeds of the baseline prediction are below the boundary speed. If they are above the boundary speed, it is assumed that freeflow conditions exist.

Time periods Traffic conditions variate over the day. Especially during rush hours, traffic is more congested than outside the rush hours. Therefore, the accuracies of the predictions are compared for different periods. These periods are:

- Morning rush hour (between 7:00 and 9:00 hour)
- Evening rush hour (between 16:30 and 18:30 hour)
- Between rush hours (between 9:00 and 16:30 hour)
- Weekends

3.5 Conclusion of the approach

The method for the prediction of vehicle speeds contains two methods: the baseline prediction meant for the long-term speed predictions during regular conditions and a short-term prediction that could be used for the regular and irregular conditions.

The baseline prediction contains an average speed and a speed interval per cluster: per segment, per day of the week and time of the day. The determination of the speeds is based on clustered historical speed data during regular conditions.

For the short-term, the baseline prediction does not accurately predict vehicle speeds. The prediction can be optimized by the knowledge of the temporal and spatiotemporal recurrent patterns in the speed residuals. The duration of these recurrent patterns is analyzed by calculating the autocorrelation in the residuals. Above that, the variations in vehicle speeds due to accidents are analyzed.

A method is designed for the short-term prediction of vehicle speeds, whereby the baseline prediction is updated with real-time data of vehicle speeds. A Kalman filter is used to account for the noise in the real-time speed data and to make an estimate of the next speed. The Kalman filter is adapted to predict multiple time steps, by combining every time step the temporal and spatiotemporal variations of the residuals with an inverse variance weighting method. Moreover, the prediction includes the expected duration of the congestion caused by the accident if an accident occurs. To manage that accurate predictions are achieved independent of the location that must be predicted, the spatiotemporal correlation between successive locations is taken into account.

In the end, the predicted average speed and the predicted speed interval of the short-term and long-term predictions are evaluated. The predicted average is evaluated in two ways. The prediction is compared with other prediction methods with the Mean Absolute Percentage Error (*MAPE*). Furthermore, it is evaluated

if the predicted average is optimal by the determination correlation exists in the residuals. Finally, the likelihood that the actual measurement lies in the predicted interval is calculated.

4 Results

In this chapter, the results from the analyses are showed. These analyses show the importance of the use of real-time vehicle speed measurements, use of historical vehicle speeds, and knowing irregular situations in the prediction of vehicle speeds short- and long-term.

Moreover, the results of the used prediction method for the long-term prediction under regular conditions and the short-term prediction for regular and irregular conditions are evaluated. In the evaluation, the short-term prediction is divided in a temporal prediction for which the temporal variations are taken into account in the prediction method; the spatiotemporal prediction for which the spatiotemporal variations are taken into account; and the short-term prediction for which both of these variations are taken into account.

To manage that the prediction accuracy is not dependent on location, the prediction results used for this evaluation are carried out on another area than for which the prediction method is calibrated.

4.1 Long-term prediction of vehicle speeds for regular conditions

The baseline prediction is the method for the long-term prediction of vehicle speeds for regular conditions. The average of the baseline prediction is determined with two different methods. It is investigated which method will give the most reliable average speed predictions.

Moreover, it is investigated if variations exist in the vehicle speeds so that these variations can be added to the baseline prediction.

4.1.1 Determination of the average speed

Two methods for the determination of the average speed of one cluster are investigated: one using the assumption that the speeds follow a normal distribution within a cluster and one under the assumption that speeds do not follow a normal distribution.

The speeds within a cluster will not follow a normal distribution, because 10% of the calculated averages are outside the interval for which it is stated that the average will lie (see Equation 3). Therefore, Equation 6 will be used for the calculation of the average speed for the baseline prediction.

4.1.2 Analysis of long-term variations in vehicle speeds

The baseline prediction is independent of the week number or day number within the year, because no significant correlations are found in the variations in the average hourly residuals for successive days or weeks ($\alpha = 0.05$).

By plotting a timeseries of the relative residuals, also no systematic variations are found (see Figure 7). In this figure, all unknown speeds are replaced by the freeflow speed. However, it is unknown if, in reality, these unknown speeds are freeflow. Therefore another timeseries is made for which the unknown speeds are removed, because it is unknown if this replacement is correct and therefore it could cause noise. Nevertheless, both figures show the same pattern.



Figure 7: The relative residuals of successive days between 8 and 9 'o clock. Unknown speeds are replaced by the freeflow speed.

4.2 Analysis of short-term variations in vehicle speeds

Repetitive patterns occur in the residuals short-term. A distinguish will be made by: repetitive patterns during regular conditions, and repetitive patterns caused by irregular situations.

4.2.1 Analysis of variations under regular conditions

Repetitive patterns under regular conditions are shown by the calculation of autocorrelation of the residuals. Herein, both temporal and spatiotemporal autocorrelations are investigated.

Temporal variations Temporal variations in the residuals are recommended to implement in the short-term prediction, because the current variations have an influence on the next variations up to 95 minutes.

It is demonstrated that the residuals are significantly correlated up to 95 minutes in the morning rush hour at the A1 in Barneveld during regular conditions for $\alpha = 0.05$ (see Figure 8). Moreover, it is assumed that these correlations also occur for other locations and times of the day.



Figure 8: Correlation of the measurement error in the morning rush hour compared to the average for lags of 5 minutes

For urban roads, the interval for which correlations occur is smaller than for freeways. For example, on the 50 km/h road near the Hospital in Amersfoort, correlation occurs up to 10 minutes. This small interval for which correlation occurs, means that only a slightly more accurate prediction can be achieved using realtime measurements. Therefore, urban roads will not be considered in this study. **Spatiotemporal variations** Spatiotemporal variations of the residuals do also take an important part in the accuracy of the short-term prediction. However, the current measurements influence the next variations only up to 35 minutes.

Correlations are found for a distance up to 10 km (see Figure 9), which corresponds with a time horizon up to 35 minutes for a propagation velocity of 17 km/hour upstream the roads. This velocity is used, because it gives the highest cross-correlation during regular conditions for this case study.

These correlations confirm that real-time measurements upstream the to be predicted segment can be used in the prediction scheme to achieve more accurate predictions. Correlations for a longer time horizon could not be investigated because the length of the roads in the case study is limited.



Figure 9: Correlation of the relative residuals in the morning rush hour with a speed of 17 km/h upstream the roads upstream road in The Hague. The lags are expressed per km distance, corresponding with 3.5 minutes travel time

4.2.2 Analysis of variations under known irregular conditions: accidents

For the irregular conditions, only the variations caused by accidents are investigated. The number of accidents per segment used for this analysis, is visualized in Appendix B: Visualisation of location characteristics. The distance and the duration that variations in the residuals occur are investigated. **Distance** No conclusions could be drawn about the effects of the accidents on the distance that variations are visible. It is assumed that this distance mainly depends on the characteristics of the location.

Duration However, some small conclusions can be drawn about the duration that variations are visible at the location of the accident. Herein, two characteristics are distinguished: the effects of the relative number of closed lanes and the relative speed of the baseline prediction are investigated on the duration of these variations.

The distribution of the duration of the congestion at the position of the accident is shown in Figure 10. In this figure, all accidents are investigated, despite the characteristics of the accidents. The duration lies in 90% of the cases between 5 and 160 minutes. Sometimes no congestion occurs, while an accident is reported in the information provided by NDW. On the other hand, some high peaks exist in the duration. All in all, in most cases, the congestion at the position of the accident takes 30 minutes. Therefore, in the short-term prediction, it is assumed that up to 30 minutes, the speed is reduced at the place of the accident.



Figure 10: The duration of the congestion at the position of the accident in Amersfoort between June, 1st 2018 and May, 1st 2019

No difference is visible between the number of closed lanes and the duration of the congestion at the position of the accident. Moreover, the results are too noisy to conclude something about the effects of the regular speed on the congestion duration. Therefore, in the designed short-term prediction method, the predicted vehicle speed is independent of the number of closed lanes and the regular speed.

4.3 Evaluation of the accuracy of the vehicle speed predictions

This section presents the evaluation of the accuracy of the vehicle speed predictions. This evaluation includes the evaluation of the accuracy of the predicted average and the likelihood that the measurement will be in the predicted interval.

The accuracy of the predicted average is evaluated with the Mean Absolute Percentage Error (MAPE). This evaluation method could give a distorted view about the reliability of a good prediction, because high deviations of the prediction errors occur during congestion. However, to give an idea about the reliability, with an example, the predicted speeds are shown.

In this chapter, only the figures of the predictions during congested conditions and the figures during freeflow conditions are showed, because these give the two extreme results. Predictions are a combination of freeflow and congested conditions, so the results will be somewhere in between. The results of different time periods will be described to substantiate the accuracy of the predictions. However, the figures of the results of these different time periods are shown in Appendix Appendix H: Extended results of the evaluation.

4.3.1 Example of the prediction: Barneveld

The behaviour of the predicted speed average and the predicted speed interval are shown with an example. For this example, the location Barneveld is chosen because this location has much congestion which results in many known measurements. The used road section is shown on a map in Figure 5. A time window is chosen on which congestion occurs. It is expected that this example is representative of a large part of the situations.

Predicted average The average of the short-term and long-term prediction is compared with the TomTom Speed Profile and the instantaneous prediction in Figure 11. It shows that the short-term prediction mainly makes use of temporal and/or spatiotemporal variations in the first 30 minutes. After 30 minutes, the short-term prediction will follow the baseline prediction. The instantaneous prediction deviates much from the real measurements for time horizons larger than

10 minutes. These large deviations are the results of the fact that the instantaneous prediction does not account for that the start and end of a rush-hour are recurrent in time and location. Especially at the beginning or at the end of a rush hour, these large deviations are visible.



Figure 11: A prediction based on real-measurements until 9:30 compared to the actual measurements and other prediction methods at the intersection at Barn-eveld.

Moreover, the baseline prediction is less smooth and predicts lower speeds than the TomTom Speed Profile. However, it is assumed that these variations in speeds of the baseline prediction will not influence the prediction accuracy, because the variation is small. The baseline prediction is not smooth, because a small sample is used for the prediction, and this is not compensated by smoothening the speeds in time or space. The TomTom Speed Profile is a function over time, always result in smoother speeds over time.

It is shown in Figure 12 that the short-term prediction method can predict congestion propagation in time and space by using both temporal and spatiotemporal variations in the prediction method. The predicted speeds, starting from the blue vertical line, follow the same pattern as is shown by the actual measurements in Figure 13. In this figure, the dotted horizontal line at distance zero corresponds to the location Barneveld for which the timeseries is shown.



Figure 12: TD chart of the prediction on the intersection at Barneveld at the start of 9:30, knowing measurements till that moment.



Figure 13: TD chart of the actual measurements on the intersection at Barneveld.



Figure 14: A prediction based on real-time measurements until 9:30 compared to the actual measurements at the intersection at Barneveld. The light coloured band shows the standard deviation of the prediction and measurements.

Predicted speed interval The short-term and long-term prediction are shown in Figure 14. The light coloured band indicates the predicted intervals. Assuming that the variations are normally distributed, the interval of the baseline prediction covers 96% of the measurements. In this example, the interval of the baseline prediction is broad, because the speeds have large variations within a cluster. Although the short-term prediction also has a large interval, the interval is much smaller than the baseline prediction. Moreover, this figure indicates that the interval of the short-term prediction covers almost all measurements. It is investigated if this expectation is correct with the likelihood in Section 4.3.3.

4.3.2 MAPE of the predicted average

The prediction of the average is compared with other prediction methods with the Mean Absolute Percentage Error (MAPE). In Figure 15 shows the MAPE of different predictions for which the regular conditions are congested.

The short-term prediction that is designed in this study does not give the most accurate predictions. Up to a time horizon of approximately 10 minutes, the instantaneous prediction gives more accurate predictions than the short-term prediction. Moreover, from a time horizon of approximately 30 minutes, the baseline



Figure 15: MAPE for several prediction methods for which the regular conditions are congested. Unknown measurements are replaced by the baseline prediciton

prediction is more accurate than the final prediction.

However, other results are obtained for which the regular traffic conditions are freeflow (see Figure 16). In this figure, a scenario is used that vehicles drive freeflow speed if the measurements are unknown. Again it is shown that the baseline prediction has fewer errors compared to the measurements than the TomTom Speed Profile. The instantaneous prediction also has the assumption that the vehicles drive freeflow speed if the speed measurements are unknown. Because these assumptions are the same while other prediction methods make use of other assumptions, the MAPE of the instantaneous prediction is lower than the other prediction methods.

However, it is unknown if the scenario is correct that vehicles drive freeflow speed if the measurements are unknown. Two other scenarios are evaluated whereby vehicles drive the speed defined by the TomTom Speed Profiles and the speed defined by the baseline prediction. These scenarios show the same interaction between the prediction methods compared to the height of the MAPE. Figures of the MAPE for these scenarios for several prediction horizons can be found in Appendix H: Extended results of the evaluation.

The figures show that the baseline prediction is a suitable prediction method if no real-time measurements are available or for long-term predictions. However, the



Figure 16: MAPE for several prediction during freeflow conditions. Unknown measurements are replaced by the freeflow speed.

accuracy of the short-term prediction is strongly dependent on the availability of known measurements. The prediction reduces in accuracy if measurements of the near past are unknown at the location and downstream the location. Moreover, the relative errors of the baseline prediction are lower than for the TomTom speed profiles.

The accuracy of the prediction is also evaluated for different periods of the day. Pictures can be found in Appendix H: Extended results of the evaluation. These results show that the prediction accuracy is a combination of the prediction accuracy during regular congested and regular freeflow conditions.

4.3.3 Likelihood of the predicted interval

Moreover, the overall likelihood that the measurement is in the predicted interval is shown for several conditions. Only the likelihood of the baseline prediction and the short-term prediction are shown, because other prediction methods that are used for comparison does not predict speed intervals.

It is investigated if the overall likelihood that the predicted interval is correct depends on the period and the traffic situation. Above that, the aim is to have an as small as possible prediction interval. Therefore, the average interval will be shown.

The likelihood that the predicted interval will be correct during regular congested locations and hours is shown in Figure 17. Figure 18 shows the likelihood that the predicted interval will be correct for all conditions. These figures show that the predicted intervals are most of the times correct. During rush hours, the prediction interval is fewer times correct than during freeflow conditions. The accuracy of the prediction is also evaluated for different periods of the day. This evaluation can be found in Appendix H: Extended results of the evaluation.

Above that, the aim is to have an as small as possible prediction interval. The average size of the predicted intervals in these two cases are shown in Figure 19 and 20. Especially, the interval of the baseline prediction is quite large. If this interval is halved, a likelihood of 75% that the predicted interval is correct is still achieved.



Figure 17: Likelihood if the prediction interval is true during congested conditions



Figure 18: Likelihood if the prediction interval is true during freeflow conditions



Figure 19: The average size of the prediction interval during congested conditions



Figure 20: The average size of the prediction interval during freeflow conditions

4.4 Conclusion of the results

The baseline prediction method calculates an average speed and a speed interval from historical clustered data. The measurements are right-censored; the unknown values are approximated, whereby is assumed that the data is not normally distributed within a cluster. This calculated average speed is used as the predicted average.

Autocorrelation in the temporal and spatiotemporal variations of the residuals shows that the baseline prediction can be updated with real-time measurements to get more accurate predictions for a time horizon until 90 minutes. This updated baseline prediction method is called the short-term prediction. However, variations from day to day or week to week could not be found.

It is difficult to explain the effects of accidents on the change in vehicle speed in time and space, because these effects are strongly dependent on many characteristics. However, it is concluded that in most cases, the duration of the congestion is 30 minutes at the position of the accident. Therefore, the short-term prediction expects the same duration if an accident occurs.

From the evaluation of the results of the prediction methods, several things are concluded. Long-term the average of the baseline prediction gives more accurate predictions than the other prediction methods that are compared. The likelihood

that the speed will lie in the predicted interval is large, because the interval is large.

Up to a time horizon of approximately 10 minutes, the average of the instantaneous prediction is the most accurate. Between a time horizon of 10 to 30 minutes, the short-term prediction is the most accurate if only temporal variations and no spatiotemporal variations are used in the prediction method. Although the example shows that the spatiotemporal variations are essential for good prediction accuracies, it is assumed that the performance of these variations on the result is to location dependent.

Finally, the interval of the short-term prediction method is much smaller than the interval of the long-term prediction. The likelihood that the actual speed is in the predicted interval seems low if the regular traffic conditions are congested. However, the regular traffic conditions are not always congested, which results that for all periods of the day, the likelihood is high that the actual speeds lie in the interval.

5 Discussion

In this chapter, the long-term and short-term analysis of the variations in speeds are discussed. Moreover, the designed short-term and long-term prediction method and the evaluation method are discussed.

5.1 Long-term prediction of vehicle speeds for regular conditions

It is expected that the baseline prediction of a particular segment is quite a stable prediction unless the road network is drastically changed around that segment.

Results show that the average of the baseline prediction is accurate. However, it could be that the predicted average does not represent the actual average well. One reason is that it is unknown if the method to estimate the unknown speed values give accurate results. Moreover, the baseline prediction is based on a limited number of measurements.

The assumption that the variations are normally distributed is incorrect for the prediction of the average speed. However, the predicted interval of the baseline prediction assumes that the variation in speeds within a cluster is normally distributed. Therefore, a standard deviation is used around the average. It is expected that this assumption can be made, because the results show that the likelihood that the predicted interval is correct corresponds with the likelihood for which the predicted interval is designed.

5.2 Analysis of short-term variations in vehicle speeds

The baseline prediction does not distinguish long-term variations, because longterm variations could not be found. However, it is expected that long-term variations do occur. If they are occurring, they influence the short-term variations which result that longer recurrent patterns are found.

The used HDFlow data is not raw data, but processing is already done on this data. This processing can influence the correlation and so the prediction results. However, it is expected that, despite this processing, the data represents the actual speeds well. Therefore, it is expected that if the prediction method is used in reality, the same results will be obtained.

A short inspection indicates that the duration of recurrent temporal patterns is dependent on the type of roads such as freeways and urban roads. Recurrent patterns have a longer duration on freeways than on urban roads. Therefore, it is expected that the short-term prediction improves the baseline prediction for a much smaller time horizon.

In this study, it is analyzed if repetitive patterns exist in accidents. However, no distinction is made between several types of accidents, because no differences could be found. It is expected that no differences are found, because the impact of accidents depends on many factors which are not investigated, which causes variations. Examples of the factors are the location characteristics, the traffic intensities and the duration that the accident barricades a lane.

Other irregular situations than accidents are not analyzed because of limited data. It is expected that prediction accuracy is smaller than in the case of knowing the irregular situation.

5.3 Short-term prediction of vehicle speeds

The duration of recurrent spatiotemporal patterns is strongly dependent on the location within the network. The short-term prediction method uses correlations between adjacent segments to make a distinction between locations. However, this addition in the prediction method results in less accurate predictions on freeways than without the addition. Several reasons could cause this less accuracy.

The first reason is that the correlations are not well determined, because estimations of the unknown speed measurements influence the correlations. Furthermore, it is an option that this location dependency on the propagation of congestion could not be measured with correlations. Moreover, the less accuracy could be due to incorrect values of the parameters in the prediction method.

The last reason is that the propagation velocity of 17 km/h does not represent the case study used for the evaluation. In reality, this shockwave speed varies between 17 km/hour till 19 km/hour for regular situations for the used regions. If wrong propagation velocities are used, wrong congestion propagation is predicted, which result in less accurate predictions.

How long the flow of the traffic is obstructed at the location of the accident, is included in the short-term prediction method. However, the duration of the obstruction is assumed as a constant duration independent of the type of accident, because no distinction between different types of accidents could be found. It is expected that this assumption is incorrect and will cause inaccurate speed predictions if accidents occurred.

5.4 Evaluation of the accuracy of the vehicle speed predictions

Two different evaluation methods are used: the likelihood the predicted interval is correct and the Mean Absolute Percentage Error (MAPE) of the predicted average.

The likelihood evaluates if the actual measurement is within the predicted interval. In this evaluation, the interval is also correct if the measurement may be inside the interval if the measurement is unknown. In reality, it could be that the real vehicle speed is higher than the predicted speed interval. Therefore, the results of the likelihood in this study are higher than in reality. However, it is expected that this difference is not significant.

The MAPE evaluates the absolute relative error of the predicted average compared to the reality. However, if the measurement is unknown because of the censoring, the measurement is guessed using different scenarios. Moreover, in reality, a vehicle is driving faster one moment than predicted and slower next. This results in that the punctuality of the arrival time is expected to be better than the MAPE is showing, because the MAPE use absolute values. That is why the MAPE could show which prediction method performs best, but it is a worse evaluation method to show the accuracy of the prediction method.

6 Conclusion and recommendations

In this chapter, first, the conclusion of the research is presented. It shows the suitability of the designed method according to the multiple challenges described in the research aim. Then, the recommendations are listed for further research. Lastly, the recommendations for Simacan are described.

6.1 Conclusion

The research aim is to design a prediction method that is suitable for the next three aspects:

- On freeways and urban roads for multiple regions;
- During regular and irregular conditions;
- For a prediction horizon up to 3 hours.

Moreover, this designed method must be fast enough so that real-time vehicle speed predictions can be made.

Two prediction methods are designed: the long-term prediction method makes use of recurrent historical patterns in the speed variations in time and space. For all locations, all times of the day and for each day of the week, an average speed and a speed interval is determined.

The short-term prediction method uses real-time measurements to adapt the long-term prediction for a more accurate speed prediction short-term. This improvement is possible with the knowledge of short-term recurrent patterns in the variations of speeds.

Network-wide application The designed method for the long-term and shortterm prediction is possible for urban roads and freeways. Moreover, it is demonstrated that the method predicts accurately on freeways. However, the accuracy of the prediction is not evaluated for urban roads. It is expected that also on urban roads, this method can predict speeds accurately. Nevertheless, the parameters of the short-term prediction method must be calibrated to get more accurate predictions. These parameters are also dependent on the type of road, because the recurrence in changes in traffic conditions are different for urban roads and freeways. Therefore it is expected that the real-time measurements could be used for a shorter prediction horizon on urban roads than for freeways.
The prediction methods make use of floating car data, because floating car data is available for many freeways and urban roads world-wide. Moreover, only historical speed data is needed for the prediction. Furthermore, for the long-term prediction method, no knowledge of road characteristics is needed.

The short-term prediction makes use of temporal and spatiotemporal recurrent patterns in the variations of speeds. The duration of recurrent spatiotemporal patterns are strongly dependent on the location within the network; straight roads show longer recurrent patterns than roads around intersections. The designed method uses correlations between adjacent segments to make a distinction between locations. However, this addition in the prediction method results in less accurate predictions on freeways than without the addition. It is expected that on urban roads, the spatiotemporal part will not improve the accuracy of the short-term prediction. Moreover, it is expected that with an improved version of the spatiotemporal part, accurate predictions can be obtained on freeways. However, it is doubted that the spatiotemporal part will give accurate predictions on urban roads, because these roads have many intersections which result in shorter spatiotemporal recurrent patterns.

In this study, the prediction method is established by the investigation of one region. The same short-term and long-term methods are applied to another region for the evaluation of the prediction. This shows that accurate predictions are obtained without investigation of the region. Therefore, it is expected that for all regions, the method gives accurate predictions.

All conditions The long-term prediction can only predict accurate speeds during regular conditions. However, the short-term prediction method does account for irregular conditions in two ways: knowing and unknowing that irregular situations are happening.

In this study, only accidents are analyzed and integrated into the method. Other irregular situations are not analyzed because of limited data. It must be further investigated how to integrate other irregular situations to improve the method.

Too much variation is visible on the duration and length of the congested caused by accidents to take account in the short-term prediction method. However, how long the flow of the traffic is obstructed at the location of the accident, is included in the short-term prediction method. If an accident has occurred, the short-term prediction method extends the current measured speed for a particular duration that follows from the analysis.

If irregular conditions occur, but the cause is unknown, the short-term prediction still accounts for irregular conditions by predicting the propagation of the current congestion. However, it will go back earlier to the baseline than knowing an accident happens. Therefore, the prediction accuracy is lower.

Far ahead It is expected that long-term variations in vehicle speeds occur during the year. However, no variations from day to day or week to week are found. Therefore, the long-term prediction does not account for long-term variations.

With autocorrelation, it is demonstrated that temporal and spatiotemporal recurrent patterns occur in the variations of speed: recurrent temporal patterns occur up to 90 minutes; recurrent spatiotemporal patterns occur up to 30 minutes. Therefore, the long-term prediction can be improved up to 90 minutes.

The short-term prediction method updates the long-term prediction if temporal or spatiotemporal deviations compared to the baseline are found in the real-time measurements. To be sure that the deviation is not caused by noise, a Kalman filter is used. The predicted speed of the short-term prediction is a weighted average of the expected speeds from the temporal prediction, spatiotemporal prediction and long-term prediction.

From the evaluation, it can be concluded that the long-term prediction method is more accurate than the other prediction methods that are compared long-term. Moreover, the short-term prediction method predicts more accurate speed than the long-term prediction method up to approximately 30 minutes. However, up to approximately 10 minutes, it is demonstrated that the short-term prediction method is not the most accurate method.

6.2 **Recommendations for further research**

From the discussion and conclusion follow several recommendations for further research.

Using non-censored measurements It is strongly recommended to use noncensored measurements for further research. By using non-censored measurements, it is expected that more accurate predictions could be achieved during freeflow conditions. Moreover, the prediction can be better compared to the actual situation within the evaluation of the accuracy of the designed prediction method during freeflow conditions.

Above that, it would be recommended to investigate whether congestion arises. If the arise of the congestion can be predicted, better prediction accuracies can be achieved. The arise of the congestion might be predicted by knowing the speeds during freeflow conditions. However, it is not known if it could be predicted with only speed measurements. The probability of succeeding to predict the arise of congestion is larger by the addition of traffic intensity data based on the fundamental diagram.

Calibrate short-term prediction method In this study, the parameters in the short-term prediction method are not calibrated. It is expected that better prediction accuracies can be achieved with a better estimation of the parameters. These parameters tune whether the errors in speeds can be explained by noise. Moreover, they tune the portions of the temporal, spatiotemporal and baseline speed variations on the short-term prediction.

It is not investigated if the temporal and spatiotemporal predicted intervals cover the expected amount of measurements. In this study, it is assumed that the portions to combine the temporal, spatiotemporal and baseline prediction for the short-term prediction corresponds to these intervals. It is expected that these predictions are dependent on each other. Therefore the portions to combine these predictions must be investigated.

Better analyse irregular situations In this study, no distinction is made between accident types. However, it is expected that the type of accidents plays an important role in the effects on the traffic conditions (for example the duration of a closed lane). Moreover, the location of the accident within the network will influence the effects of accidents on traffic conditions. Particularly accidents influence the velocity of the congestion propagation. It is recommended to use more accident data so that it is possible to classify it on spatial characteristics. Therefore, it is expected that differences between accident types and spatial characteristics on the traffic conditions are found.

Effects of road characteristics on speeds Road characteristics such as onand off-ramps influence the arise and propagation of congestion. These effects due to road characteristics, mainly occur at intersections. Moreover, an increase or decrease in the number of lanes influence the arise and propagation of congestion. The designed short-term prediction method uses spatiotemporal correlations of the traffic conditions between adjacent road segments. However, the spatiotemporal characteristics on the effects of this propagation are not investigated. It is recommended to analyze these effects so that the short-term prediction method can better predict the propagation of congestion.

6.3 **Recommendations for Simacan**

Use baseline prediction The results of the evaluation show that the baseline prediction is more accurate than the TomTom Speed Profiles. Therefore, it will be advised to use the baseline prediction instead of the TomTom Speed Profiles, if one year of measurements are available for a specific segment.

Moreover, except the fact that it is not proven in this study, it would be assumed that the effects of weather and holidays will result in better accuracies for the baseline prediction.

Adapt the instantaneous prediction The instantaneous prediction gives accurate predictions up to around 10 minutes. It is expected that accurate predictions can be obtained up to around 30 minutes by using a combination of the baseline prediction and the instantaneous prediction. In addition, the instantaneous prediction should be made with a weighted moving average, to account for noise. Furthermore, a weighted moving average should be made to estimate the speed from previous measurements during the transition of freeflow conditions to congested conditions or the other way around. Because, if the last measurement is unknown, it could be in the transition zone, while it indicates that the condition is freeflow.

Although the short-term prediction method is quite similar to the advised method described above, it is assumed that the improvements on the accuracy of the speeds of this short-term prediction method do not balance against the extra amount of calculation time and complexity of the method, especially the Kalman filter.

Do not use spatiotemporal variations The method of the spatiotemporal prediction has the potential to obtain accurate predictions; especially for the propagation of traffic jams. However, the results show that the accuracy is not good yet. Moreover, it is expected that the complexity of the algorithm increases much for getting good prediction accuracies, which results in higher calculation times. This complexity also will take a lot of research effort, especially because many unforeseen challenges may occur.

Use non-censored data It will be advised to consider if all vehicle speed data of TomTom HDFlow will be available instead of only the speeds in congested situations. Alternatively, other data sources that contain speed measurements are advised to consider if they become accessible for usage and have a sufficiently

update frequency and a high road coverage on urban roads and highways. Having non-censored data will make the application of the prediction of speed less complicated. Moreover, it will be assumed that the prediction accuracy will be better. Also, other algorithms for other purposes than the prediction of vehicle speeds will take advantage of this.

References

- Aarts, L. T., Bijleveld, F. D., and Stipdonk, H. L. (2015). Bruikbaarheid van snelheidsgegevens uit'floating car data'voor proactieve verkeersveiligheidsanalyses: analyse van tomtom-snelheidsgegevens en vergelijking met meetlusgegevens op het provinciale wegennet.
- Abdulhai, B., Porwal, H., and Recker, W. (1999). Short term freeway traffic flow prediction using genetically-optimized time-delay-based neural networks.
- Cassidy, M. J. and Windover, J. R. (1995). víethodology for assessing dynamics of freeway traffic flow.
- Catbagan, J. L. and Nakamura, H. (2008). Desired speed distributions on twolane highways under various conditions. *Transportation research record*, 2088(1):218–226.
- Clark, S. (2003). Traffic prediction using multivariate nonparametric regression. *Journal of transportation engineering*, 129(2):161–168.
- Davis, G. A., Nihan, N. L., Hamed, M. M., and Jacobson, L. N. (1990). Adaptive forecasting of freeway traffic congestion. *Transportation Research Record*, (1287).
- Hamad, K., Shourijeh, M. T., Lee, E., and Faghri, A. (2009). Near-term travel speed prediction utilizing hilbert–huang transform. *Computer-Aided Civil and Infrastructure Engineering*, 24(8):551–576.
- Hawkins, R. K. (1988). Motorway traffic behaviour in reduced visibility conditions. In Vision in Vehicles II. Second International Conference on Vision in VehiclesApplied Vision AssociationErgonomics SocietyAssociation of Optometrists.
- Heilmann, B., El Faouzi, N.-E., de Mouzon, O., Hainitz, N., Koller, H., Bauer, D., and Antoniou, C. (2011). Predicting motorway traffic performance by data fusion of local sensor data and electronic toll collection data. *Computer-Aided Civil and Infrastructure Engineering*, 26(6):451–463.
- Innamaa, S. (2009). Self-adapting traffic flow status forecasts using clustering. *IET Intelligent Transport Systems*, 3(1):67–76.
- Kamarianakis, Y., Shen, W., and Wynter, L. (2012). Real-time road traffic forecasting using regime-switching space-time models and adaptive lasso. *Applied stochastic models in business and industry*, 28(4):297–315.
- Kerner, B. S. and Rehborn, H. (1996). Experimental features and characteristics of traffic jams. *Physical Review E*, 53(2):R1297.
- Li, R. and Rose, G. (2011). Incorporating uncertainty into short-term travel time predictions. *Transportation Research Part C: Emerging Technologies*, 19(6):1006–1018.
- Min, W. and Wynter, L. (2011). Real-time road traffic prediction with spatiotemporal correlations. *Transportation Research Part C: Emerging Technolo-*

gies, 19(4):606–616.

- Oh, S., Byon, Y.-J., Jang, K., and Yeo, H. (2015). Short-term travel-time prediction on highway: a review of the data-driven approach. *Transport Reviews*, 35(1):4–32.
- Pan, B., Demiryurek, U., Gupta, C., and Shahabi, C. (2015). Forecasting spatiotemporal impact of traffic incidents for next-generation navigation systems. *Knowledge and Information Systems*, 45(1):75–104.
- Rice, J. and Van Zwet, E. (2004). A simple and effective method for predicting travel times on freeways. *IEEE Transactions on Intelligent Transportation Systems*, 5(3):200–207.
- Thakuriah, P. (1992). *Data needs for short term link travel time prediction*. Number 19. Urban Transportation Center, University of Illinois at Chicago.
- Vlahogianni, E. I., Golias, J. C., and Karlaftis, M. G. (2004). Short-term traffic forecasting: Overview of objectives and methods. *Transport reviews*, 24(5):533–557.
- Vlahogianni, E. I., Karlaftis, M. G., and Golias, J. C. (2014). Short-term traffic forecasting: Where we are and where we're going. *Transportation Research Part C: Emerging Technologies*, 43:3–19.
- Vythoulkas, P. (1993). Alternative approaches to short term traffic forecasting for use in driver information systems. *Transportation and traffic theory*, 12:485–506.
- Wang, Y., Papageorgiou, M., and Messmer, A. (2006). Renaissance–a unified macroscopic model-based approach to real-time freeway network traffic surveillance. *Transportation Research Part C: Emerging Technologies*, 14(3):190–212.
- Xia, J., Chen, M., and Huang, W. (2011). A multistep corridor travel-time prediction method using presence-type vehicle detector data. *Journal of Intelligent Transportation Systems*, 15(2):104–113.

Appendix A: Weather impact

Weather conditions are investigated if the same weather conditions are visible for the whole country, because then the effects of the irregular situations are identical for multiple road sections. Moreover, if these weather conditions are identical for the whole country, it is investigated which dates the irregular situations occur.

The weather conditions that are inspected are:

- minimum visibility length per day (km)
- fog occurrence per day (yes/no)
- snow occurrence per day (yes/no)
- icing occurrence per day (yes/no)

These conditions are chosen, because from literature, there can be concluded that the visibility and the road surface state have the most impact on speeds (Catbagan and Nakamura, 2008; Hawkins, 1988). Moreover, weather codes orange and red will be removed. The dates of these weather codes are shown in table 10.

Table 10: Dates with weather codes orange and red between September 2018 and March 2018

Date (dd-mm-yy)	Explanation
25-01-19 10-03-19	code red as result of icing code orange as result of heavy wind

Figure 21 shows that minimum visibility's per day varies much per area; for every weather station, totally different visibility's lengths are measured. However, there are no significant differences between the weather stations, as shown in table 11. Therefore, the confidence interval is calculated with $\mu = 2\sigma$ assuming that it is a normal distribution.

However, often when fog occurs, this occurs for several areas in The Netherlands (see figure 22). The definition for fog is that the visibility is less than 1 km. For speed decreases, the fog must be less than 300 meters (Hawkins, 1988). So many dates that are shown, for which theoretically no speed decrease will happen.

For snow and icing, many events happen on all weather stations (see figures 23 and 24). It is shown that in Twente, more dates with icing occur than in other re-



Figure 21: Minimum measured visibility per day for several weather stations

Table 11: Average difference in visibility between De Bilt and other weather stations with a 95% confidence interval

Weather station	average difference in km
Leeuwarden	$\textbf{0.3} \pm \textbf{21.0}$
Maastricht	$\textbf{0.3}\pm\textbf{20.4}$
Rotterdam	$\textbf{0.5} \pm \textbf{19.4}$
Twente	$\textbf{1.0} \pm \textbf{21.2}$

gions, probably because it has a more continental climate than the other stations, so it is colder during winters. However, overall, the dates with snow and icing for the Bilt seems quite representative for The Netherlands so these dates will be removed. These dates are shown in table 12, 13, 14 and 15.



Figure 22: Measured dates with fog for several weather stations



Figure 23: Measured dates with icing for several weather stations



Figure 24: Measured dates with snow for several weather stations

Table 12: Dates when icing occurs in De Bilt between May 1st 2018 and July 1st 2019

Date (yyyy-mm-dd)
2018-11-03
2018-12-16
2018-12-28
2019-01-18
2019-01-19
2019-01-24
2019-01-25
2019-01-31
2019-02-02
2019-02-03
2019-02-17
2019-03-24

Table 13: Dates when icing occurs in Rotterdam between May 1st 2018 and July 1st 2019

Date (yyyy-mm-dd)
2018-11-16
2018-12-12
2018-12-15
2018-12-16
2019-01-25
2019-01-31
2019-03-24

Table 14: Dates when snow occurs in De Bilt between May 1st 2018 and July 1st 2019

Date (yyyy-mm-dd)
2018-12-15
2018-12-16
2019-01-17
2019-01-22
2019-01-24
2019-01-30
2019-01-31
2019-02-01
2019-02-02

Table 15: Dates when snow occurs in Rotterdam between May 1st 2018 and July 1st 2019

Date (yyyy-mm-dd)
2018-12-15
2018-12-16
2019-01-22
2019-01-24
2019-01-25
2019-01-30
2019-01-31
2019-02-01
2019-02-04

Appendix B: Visualisation of location characteristics

In this appendix, first the number of known measurements per day for the region Amersfoort is visualized. Then, the number of accidents for the time period of this case study is visualized.

Number of known measurements

The speed measurements are right-censored, which means that many speed measurements are unknown during freeflow conditions.

Figure 25 shows the average amount of known measurements to indicate the size of known measurements. Herein, for each time interval of 5 minutes it is investigated if at least one measurement is known.



Figure 25: Average amount of known measurements, expressed in hours per day

Accidents

The number of accidents per segment used for the analysis of the accidents is shown in Figure 26. Figure 27 shows the number of accidents per segment used for the evaluation of the prediction.



Figure 26: Times affected by accidents in 2018 per road section in Amersfoort



Figure 27: Times affected by accidents in 2018 per road section in Den Haag

Appendix C: Data Processing of Speed measurements: TomTom HDFlow

This appendix describes in short the data collection and data processing of the TomTom HDFlow data.

Collection

Real-time speed measurements are retrieved from tomtom and are stored in compressed files for which every file contains data of all road segments for a time interval of 30 seconds.

Variabele	Туре	Description
location	openLR	HDflow segment
region	2-letter country code	country
pub time src	Unix time in millisec- onds	publication time of source
received time	Unix time in millisec- onds	received time
speed	km/h	actual speed determined from the measured traveltime
freeflow	km/h	estimated speed during free flow traf- fic state
traveltime	Sec	actual measured traveltime on the segment
confidence	%	confidence of the measured travel- time
road blocked ¹	boolean	indicate if the road is blocked

The data contains the following important values:

Table 16: Available data from TomTom

The raw data contains of two files of the Netherlands per minute. These files have a size of approximately 2.3MB when they are compressed and a size of approximately 31MB when they are decompressed. Around 100.000 segments are available for which each file contains measurements.

Pre-processing

The data is processed to a classification to which every file contains all data for a specific segment. This processing is performed in AWS with an EC2 machine for the compute capacity and S3 for the data storage (see figure 28).



Figure 28: Overview of the processing of the data

The raw data is stored in files compressed with Lempel–Ziv–Markov algorithm (LZMA). Each file contains all data of a time interval of 30 seconds in a format with lines of JSON. After processing, the data is stored in files compressed with GNU zip. Each file contains all data of a segment in a CSV format. The overview of the process is shown in figure 29 and 30.

¹determined at Simacan from this data



Figure 29: Overview of the processing of the data part 1



Figure 30: Overview of the processing of the data part 2

Appendix D: AlertC codes indicating irregular situations

The tables in this appendix show the alertC codes that will be categorized as non-recurrent conditions in this study.

Accidents

Table 17: Used alertC codes representing accidents

alertC code	description
201	door een ongeluk
202	door een ernstig ongeluk
203	door een kettingbotsing
204	door een ongeluk met een vrachtwagen
205	door een ongeluk met gevaarlijke stoffen
206	door een ongeluk waarbij brandstof weglekt
207	door een ongeluk waarbij chemische stoffen weglekken
214	incident
335	door een ongeluk met een bus
336	ongeluk met olielekkage
337	door een gekantelde auto
338	door een gekantelde vrachtwagen
339	door een geschaarde auto met aanhanger
340	door een geschaarde caravan
341	door een geschaarde vrachtwagencombinatie
342	geslipte voertuigen op de rijbaan
345	secundair ongeluk
378	door een gekantelde auto

Roadworks

Table 18: Used alertC codes representing roadworks

alertC code	description
701	wegwerkzaamheden

Weather

alertC code	description
1001	gevaarlijke rij-omstandigheden
1003	gladde weg
1006	ijsvorming
1009	ijzel
1109	zware regenval
1209	windstoten
1301	dichte mist
1304	mist
1307	mistbanken
1318	zicht verminderd
1324	slecht zicht door opspattend water

Table 19: Used alertC codes representing weather

Other

alertC code	description
11	te hoog voertuig gesignaleerd, wordt afgehandeld
25	de tunnel is dicht
26	de brug is versperd
61	liggen voorwerpen op de weg
208	door een ongeluk op de andere rijbaan
209	door een ongeluk op de andere rijbaan
210	afgevallen lading
211	kapotte auto
212	kapotte vrachtauto
213	er staat een auto in brand
344	door een politieonderzoek
346	defecte bus
347	te hoog voertuig
397	bergingswerkzaamheden
401	dicht
402	dicht
406	de oprit is dicht
407	de afrit is dicht
408	de op- en afrit is dicht

479parallelrijbaan dicht480rechter parallelrijbaan dicht483de hoofdrijbaan is dicht490hoofdrijbaan versperd500rijbaan dicht501de rechterrijstrook is dicht503de linkerrijstrook is dicht506drie rijstroken zijn dicht507séén rijstroken zijn dicht518afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook is dicht648de spitsstrook is dicht649ge spitsstrook is dicht649ge spitsstrook is dicht641Eén rijstrook is dicht642gaslek920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1057grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
480rechter parallelrijbaan dicht483de hoofdrijbaan is dicht490hoofdrijbaan versperd500rijbaan dicht501de rechterrijstrook is dicht503de linkerrijstrook is dicht506drie rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht649gaslek920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
483de hoofdrijbaan is dicht490hoofdrijbaan versperd500rijbaan dicht501de rechterrijstrook is dicht503de linkerrijstrook is dicht506drie rijstroken zijn dicht507afwisselend slechts één rijstrook per richting beschikbaar518er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht648de spitsstrook is dicht649eer een hevige brand920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1043olie op het wegdek1044olie op het wegdek1045grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
490hoofdrijbaan versperd500rijbaan dicht501de rechterrijstrook is dicht503de linkerrijstrook is dicht505twee rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1044olie op het wegdek1044olie op het wegdek1044olie op het wegdek1044olie op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
500rijbaan dicht501de rechterrijstrook is dicht503de linkerrijstrook is dicht505twee rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht649gaslek900gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1057grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
501de rechterrijstrook is dicht503de linkerrijstrook is dicht505twee rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht649gaslek905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
503de linkerrijstrook is dicht505twee rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1044olie op het wegdek1057grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
505twee rijstroken zijn dicht506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1044olie op het wegdek1057grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
506drie rijstroken zijn dicht513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
513afwisselend slechts één rijstrook per richting beschikbaar514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
514er is één rijstrook open515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fictorer op de weg
515er zijn twee rijstroken open637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
637rijstrook voor hulpdiensten dicht641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
641Eén rijstrook is dicht648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
648de spitsstrook is dicht905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1043olie op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
905er zijn bomen omgewaaid920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
920gaslek921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
921door een hevige brand922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
922dieren op de weg924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
924door opruimwerkzaamheden977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg
977bermbrand1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fietsers op de weg
1037zeer gevaarlijke rij-omstandigheden1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fiotsors op do wog
1041water op het wegdek1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fiotsers op de weg
1042zand op het wegdek1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fiotsers op de weg
1044olie op het wegdek1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fiotsers op de weg
1067grote dieren op de weg1309slecht zicht door rookontwikkeling1472mensen op de weg1474fietsers op de weg
 1309 slecht zicht door rookontwikkeling 1472 mensen op de weg 1474 fietsers op de weg
1472 mensen op de weg
1/7/ fioteors on do woo
1474 Heisels op de weg
1476 veiligheidsincident
1501 door bezoekers aan een evenement
1515 veiligheidsmaatregelen
1583 langdurige vertraging door veiligheidsmaatregelen
1701 spookrijders
1804 verkeerslichten buiten werking
1805 verkeerslichten werken niet goed
3092 spoedreparatie
31/4 vakantieverkeer

Table 20: Used alertC codes representing other important non-recurrent situations

Appendix E: Extended foundation of the baseline prediction

This appendix describes the foundation of the equation of the baseline prediction by the assumption that the variation in speeds within a cluster are normally distributed.

From the historical data, the following parameters are known:

p	the fraction of the measurements that are known
x	a single measurement
v_b	the boundary speed, so the maximum speed that is measured
v_{ff}	the speed that will be driven under freeflow conditions
V_{low}	the speed assuming that all unknown measurments are equal to the
	boundary speed
V_{up}	the speed assuming that all unknown measurments are equal to the speed in freeflow condition
1£	the probability density function of the normal distribution
paj	the probability density function of the normal distribution
cdf	the cumulative density function of the normal distribution

The following parameters must be determined:

- var(V) the variance of the expected speed
- σ the standard deviation of the speed, assuming it follows a normal distribution
- μ the average speed, assuming it follows a normal distribution

It is known that the real $V \sim N(\mu_v, \sigma_v^2)$ exist of the two parts $E(V|V < v_b)$ and $E(V|V > v_b)$. Using the assumption that the speed measurements are normal distributed, the expected value for the two parts are:

$$E(V|V \le v_b) = \mu - \sigma \cdot \frac{pdf(b^*)}{cdf(b^*)}$$

$$E(V|V > v_b) = \mu + \sigma \cdot \frac{pdf(b^*)}{1 - cdf(b^*)}$$

$$b^* = \frac{v_b - \mu}{\sigma}$$
(37)

Where pdf and cdf are the probability density function and the cumulative density function. For every segment and time, two distributions will be created from the

historical data: $V_{low,s,t}$ is the distribution of all measurements for which all unknown measurements are replaced by the boundary speed. $V_{high,s,t}$ is the distribution of all measurements for which all unknown measurements are replaced by the freeflow speed.

The expected values of these two distributions are:

$$E(V_{low}) = cdf(b^{*}) \cdot E(V|V < v_{b}) + (1 - cdf(b^{*}))v_{b}$$

$$E(V_{low}) = cdf(b^{*}) \left(\mu - \sigma \cdot \frac{pdf(b^{*})}{cdf(b^{*})}\right) + (1 - cdf(b^{*}))v_{b}$$
(38)

$$E(V_{up}) = cdf(b^{*}) \cdot E(V|V < v_{b}) + (1 - cdf((v_{b} - \mu)/\sigma))v_{ff}$$

$$E(V_{up}) = cdf(b^{*})\Big(\mu + \sigma \cdot \frac{pdf(b^{*})}{1 - cdf(b^{*})}\Big) + (1 - cdf(b^{*}))v_{ff}$$
(39)

For the censored data, the only unknown values are μ en σ . $E(V_{low})$ is the sample mean and equation 40 denotes the number of measurements below the boundary speed v_b . Moreover, knowing the *cdf*, the *pdf* can be estimated with equation 41. This results in equation 42.

$$cdf((v_b - \mu)/\sigma)) \approx p$$
 (40)

$$pdf(b^*) \approx pdf(cdf^{-1}(p))$$
 (41)

$$E(V_{low}) \approx p \cdot E(V|V < v_b) + (1-p)v_b$$

$$\approx p \cdot \left(\mu - \sigma \cdot \frac{pdf(cdf^{-1}(p))}{p}\right) + (1-p)v_b$$
(42)

Rewriting for μ gives:

$$\mu \approx \frac{E(V_{low})}{p} + \sigma \cdot \frac{pdf(cdf^{-1}(p))}{p} - \frac{(1-p)}{p}v_b$$
(43)

Appendix F: State estimation of the Kalman Filter

This appendix shows, with an example of the location at Barneveld, the estimation of the speed, determined with the Kalman Filter. The speed is estimated from temporal speed variations to come to a temporal prediction. The spatiotemporal prediction is made with an estimated speed based upon the relative variation in measurement errors compared to the baseline.

It is chosen to use the relative variation of the errors compared to the baseline prediction, because the highest correlations can be found for this situation. On the other hand, the temporal prediction is made based upon the absolute variation because of the same reason that the variations are the highest.

Examples of the speed estimation for the temporal prediction, knowing the previous speed measurements (same segment previous times) are shown in figures 31 and 32. Figure 33 shows all the estimates and measurements if this is executed for all measurements of the recent past for the same segment.



Figure 31: Example of the estimation of the speed due to the temporal variations at 9:20, knowing the measurement, and the previous estimation at 9:15



Figure 32: Example of the estimation of the speed due to the temporal variations at 9:30, knowing the measurement, and the previous estimation at 9:25



Figure 33: Example of the estimation of the speed due to the temporal variations at 9:30, knowing the measurement, and the previous estimation at 9:25

Examples of the speed estimation for the spatiotemporal prediction, knowing the previous speed measurements (downstream segments, previous times) are shown in figures 34 and 35. Figure 36 shows all the estimates and measurements if this is executed for all measurements of the recent past for the same segment.



Figure 34: Example of the estimation of the speed due to the relative spatiotemporal variations in the errors at 9:00, knowing the measurement, and the previous estimation at 8:55



Figure 35: Example of the estimation of the speed due to the relative spatiotemporal variations in the errors at 9:30, knowing the measurement, and the previous estimation at 9:25



Figure 36: Example of the estimation of the speed due to the relative spatiotemporal variations in the errors at 9:30, knowing the measurement, and the previous estimation at 9:25

Appendix G: The estimation of the short-term prediction by combining several timeseries

In this appendix, the estimation of the short-term prediction from the baseline, temporal and spatiotemporal prediction is showed with examples. These predictions are combined with the inverse variance weighting method.

Figures 37, 38, 39 and 40 shows examples of the distributions of the baseline, temporal and spatiotemporal predictions. Also, it shows the final prediction, which corresponds to the short-term prediction and the afterwards known measurement. This measurement also has a distribution which shows the noise of the measurement, which is the variation that could not be explained by historical measurements or temporal or spatial variations.

It should be kept in mind that the Kalman Filter must be fed with new measurements to predict the next timestep. However, since the measurements are unknown, the predicted total speed of the previous time step is used as measurement.



Figure 37: Example of the short-term prediction for a time horizon of 5 minutes, compared to the three partly predictions and the afterwards known measurement.



Figure 38: Example of the short-term prediction for a time horizon of 30 minutes, compared to the three partly predictions and the afterwards known measurement.



Figure 39: Example of the short-term prediction for a time horizon of 60 minutes, compared to the three partly predictions and the afterwards known measurement.



Figure 40: Example of the short-term prediction for a time horizon of 90 minutes, compared to the three partly predictions and the afterwards known measurement.

Appendix H: Extended results of the evaluation

This appendix shows supplemental figures of the MAPE and the likelihood to support the results of the main report. Figures show the MAPE using different scenario is about the unknown speed measurements. Moreover, the evaluation of the prediction accuracy depending on the period of the day is shown.

Different scenarios about the unknown speed measurements

This section shows the MAPE in a figure for three different scenarios about the unknown speed measurements. The figures show almost the same relationships between the size of the errors of the prediction methods. Conclusions that could be drawn is that overall, the baseline prediction is more accurate than the TomTom Speed Profiles. Also, the temporal prediction is more accurate than the instantaneous prediction. In addition, the temporal prediction behaves better than the spatiotemporal prediction and the short-term prediction.

Point of interests is that the MAPE of the baseline prediction in Figure 41 is better than in reality, because many unknown values are replaced with the values of the baseline prediction. The same reasoning holds for the TomTom Speed Profiles in Figure 42 and the instantaneous prediction in Figure 43. The instantaneous prediction predicts freeflow if the speed measurement is unknown.



Figure 41: MAPE for several prediction methods during freeflow conditions. Unknown measurements are replaced by the baseline prediction



Figure 42: MAPE for several prediction methods during freeflow conditions. Unknown measurements are replaced by the TomTom Speed Profiles



Figure 43: MAPE for several prediction methods during freeflow conditions. Unknown measurements are replaced by the freeflow speed

Evaluation of the prediction accuracy depending on the time period of the day

These results show that for each time period the prediction accuracy is a combination of the prediction accuracy during regular congested and regular freeflow conditions.



Figure 44: Likelihood if the prediction interval is true during the morning rush hour



Figure 45: Likelihood if the prediction interval is true during the evening rush hour


Figure 46: Likelihood if the prediction interval is true between rush hour



Figure 47: Likelihood if the prediction interval is true during weekends



Figure 48: The average size of the prediction interval during the morning rush hour



Figure 49: The average size of the prediction interval during the evening rush hour



Figure 50: The average size of the prediction interval between the rush hours



Figure 51: The average size of the prediction interval during the weekends



Figure 52: MAPE for several prediction methods during the morning rush hour. Unknown measurements are replaced by the baseline prediction



Figure 53: MAPE for several prediction methods during the evening rush hour. Unknown measurements are replaced by the baseline prediction



Figure 54: MAPE for several prediction methods between the rush hours. Unknown measurements are replaced by the baseline prediction



Figure 55: MAPE for several prediction methods during the weekends. Unknown measurements are replaced by the baseline prediction