

1 Article

2 Predicting Freeway Travelling Time Using Multiple- 3 Source Data

4 Kejun Long¹, Wukai Yao², Jian Gu^{1*} and Wei Wu¹

5 ¹ Hunan Provincial Key Laboratory of Smart Roadway and Cooperative Vehicle-Infrastructure Systems,
6 Changsha University of Science & Technology, Changsha 410004, China; longkejun@csust.edu.cn;
7 jiaotongweiwu@csust.edu.cn; gujian87@126.com

8 ² School of Traffic and Transportation Engineering, Changsha University of Science & Technology, Changsha
9 410004, China; 185803450@qq.com

10 * Correspondence : gujian87@126.com; Tel.: +86-731-8525-8575

11

12 **Abstract:** Freeway travelling time is affected by many factors including traffic volume, adverse
13 weather, accident, traffic control and so on. We employ the multiple source data-mining method
14 to analyze freeway travelling time. We collected toll data, weather data, traffic accident disposal
15 logs and other historical data of freeway G5513 in Hunan province, China. Using Support Vector
16 Machine (SVM), we proposed the travelling time model based on these databases. The new SVM
17 model can simulate the nonlinear relationship between travelling time and those factors. In order
18 to improve the precision of the SVM model, we applied Artificial Fish Swarm algorithm to
19 optimize the SVM model parameters, which include the kernel parameter σ , non-sensitive loss
20 function parameter ε , and penalty parameter C . We compared the new optimized SVM model
21 with Back Propagation (BP) neural network and common SVM model, using the historical data
22 collected from freeway G5513. The results show that the accuracy of the optimized SVM model is
23 17.27% and 16.44% higher than those of the BP neural network model and the common SVM model
24 respectively.

25 **Keywords:** Support Vector Machine; Travelling time; Intelligent Transportation System; Artificial
26 Fish Swarm algorithm; Big data.

27

28 1. Introduction

29 Travel time is one of the main indexes that reflect the traffic operation level of a freeway, and it
30 is also the basis for Advanced Traveler Information System (ATIS), Traffic Guidance System(TGS),
31 and Traffic Control System(TCS). The challenges and difficulties of travel time prediction are
32 identified below.

- 33 • Diverse influencing factors such as weather, holidays, traffic accidents, out of sample
34 prediction, and mechanisms contributing to congestion. It's difficult to describe and predict the
35 influence mechanism by using traditional conventional mathematical models.
- 36 • The complexity and incompleteness of basic data. Although there are many flow detectors and
37 video detection equipment on the freeway, captured data are incompatible, redundant, and
38 include error or loss. To avoid these, techniques which utilize multi-source data to improve the
39 accuracy of travel time prediction is extremely important.

40 In China, practical application of travel time prediction focuses mainly on the following two
41 aspects.

42 The first aspect is the prediction of travel time by map navigation providers using their
43 personalized GPS data. Many map service providers employ their personalized data for travel time
44 forecast services and commercial products. For instance, Bai-du, Gao-de, and other Chinese map
45 providers collect real-time GPS data from users while providing map navigation services. Then, a

46 correlation algorithm is proposed to obtain the travel time prediction result at road sections, which
47 depends on the market share of the map navigation service. The higher frequency of people using
48 the navigation service, the more complete of GPS data and the higher the prediction accuracy will
49 be. However, according to the Chinese market report, the market share of Bai-du and Gao-de
50 services are presently 29.3% and 32.6% respectively. Therefore, the accuracy of results should be
51 further improved by increasing market share.

52 The second aspect is the prediction based on the traffic detection data of urban traffic managers
53 and historical data. In recent years, numerous fixed detector devices have been installed in most of
54 urban roads and rural freeways for the prediction of travel time, including inductive loops, video
55 recorder, microwave, and laser detection. However, unavoidable damage to flow detection
56 equipment and transmission error of partial data make traffic detection data incomplete redundant
57 or error. In addition, different detectors have different data formats and data accuracy. With the
58 rough use of mistake data for precise travel-time prediction, the Traveler Information Service system
59 cannot recommend optimal travel routes or warn of potential traffic congestion and users cannot
60 determine the optimal departure time or estimate their expected arrival time based on predicted
61 travel times.

62 Theoretical research on freeway travel time prediction can be divided into two categories based
63 on single source data and multi-source data.

64 *1.1 Overview of Prediction Method of Single Source Data*

65 A single data source was earlier method used for predict travel time. Many researches
66 prediction results were obtained upon a single data source.

67 Gipps, P. G. [1] used the occupancy and arrival time to predict the travel time in a road in 1997.
68 Mehmet Y, Nikolas G [2] set statistical predictive algorithms to predict the future travel time. Shen,
69 L., & Hadi, M. [3] employed data obtained from detector in freeway. Kyung et al. [4,5] used inductive
70 loop detectors to obtain the front position and capture the interactions between trucks and
71 non-trucks. But fixed detector devices is easily affected by external environment and cannot directly
72 access some important parameters, such as travel time, etc.

73 In addition, many researches consider using GPS data to predict travel time. Ramezani et al. [6]
74 and Zhang et al. [7,8] considered the diversity of GPS data and investigated the application of
75 Markov chain to travel time estimation and implemented good prediction accuracy. Woodard,
76 Nogin & Paul et al. [9] used the GPS data of the current highest volume GPS data source, and applied
77 the TRIP method to predict the travel time. Based on GPS data sets, Bahuleyan, H., & Vanajakshi, L.
78 D. [10] proposed a prediction method for urban trunk lines which was only suitable for traffic
79 conditions in India. But the GPS data only gets the speed and real-time position information and
80 uncertainty due to the route of the vehicle, it affects the coverage and accuracy of detection data.

81 Above methods indeed are innovations and improvements in travel time prediction, and
82 results are more accurate. However, many predictions that use a single data source do not consider
83 the impact of other unexpected events or the result was not accurate enough because single data
84 source cannot reflect traffic state of road network exactly. It will result in certain errors between
85 prediction results and true values.

86 *1.2 Overview of Prediction Method of Multiple Source Data*

87 Nowadays, the development of traffic big data environment has progressed rapidly. With the
88 support of a large amount of data, it is possible to clearly visualize traffic flow changes under the
89 joint action of different factors, that is, the traffic state presented, which is more favorable. The
90 construction of the predictive model improves the adaptability and accuracy of the model, [11] and if
91 the same state occurs, it can be predicted based on historical results. The more populated the
92 database is, the higher the quality and the higher the likelihood of finding commonalities and
93 predicting accurate results will be. This concept can be applied by searching for common traffic
94 states for prediction.

Owing to progress in the dynamic traffic information acquisition system, various traffic data can be collected more easily. And data fusion is finished in the dynamic traffic information acquisition system, which is jointly determined by the advantages of multi-source data and the characteristics of traffic conditions.

And using multi-source data for prediction can overcome the limitations of the single data source. In other words single data source has low quality and is not comprehensive. The traffic state is described from different angles and directions to improve the accuracy of prediction and reduce disturbance from unexpected factors.

At present, many studies have been conducted on travel time prediction, especially studies based on the historical data travel time of multi-source data.

The common predicting methods and their characters are summarized in the table below.

Table 1. Comparison of common prediction methods

Prediction method	Author	Data source
Kalman filter	Lin J W C V.2006[12]	Travel time data
	Zhou, J. 2014[13]	Floating vehicle and fixed detector data
	Tang-Hsien Chang.2016[14]	Electronic Toll Collection (ETC) and traditional Vehicle Detector data
Bayesian estimation	Fei, X.2011[15]	The real loop detector data of an I-66 segment in Northern Virginia
	Zhan, X. 2016[16]	A large-scale taxi trip dataset from New York City
Statistical decision theory	Wosyka,J.2012[17]	Two detectors data in Prague and also in the Czech Republic.
Neural network	Innamaa, S.2005[18]	travel time data
	Lin J W C V.2005[19]	Travel data and Some missing or corrupt travel data

Compared to the single source data, the multi-source data method can extract deep information within data, significantly reduce the cost of data acquisition, and make up for lack of information and packet loss of single source data. At present, big data application technology are widely used in the traffic field, many studies have been conducted in the field of freeway travel time prediction based on big data analytics. However, there are several deficiencies including:

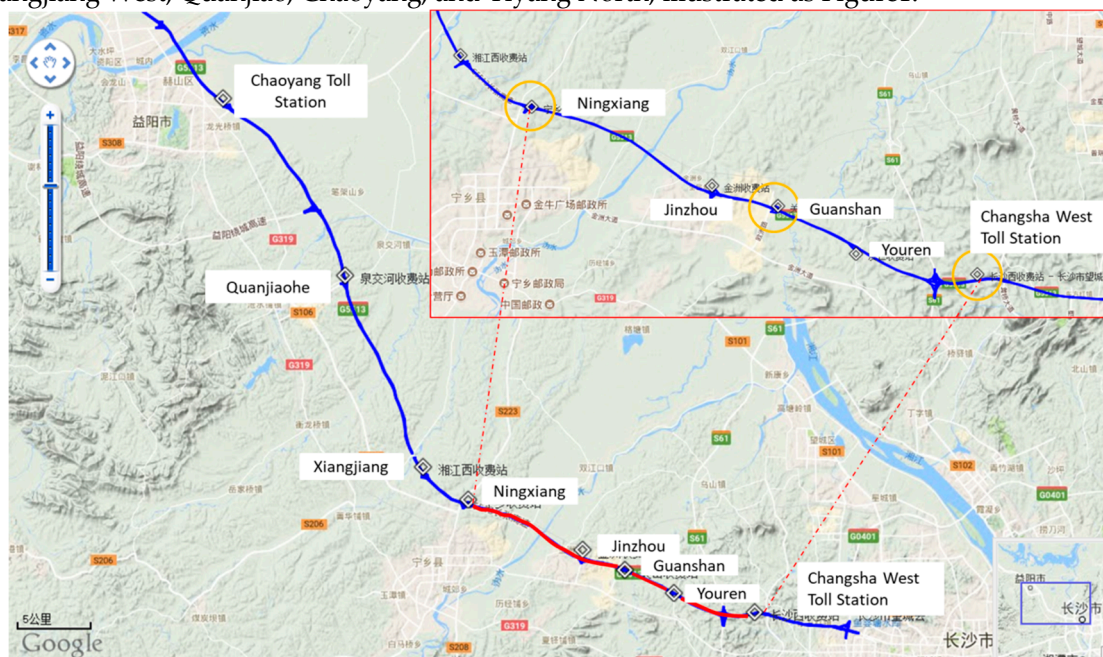
- The method pays attention to machine learning algorithms and lacks the mastery of the characteristics of the traffic flow, resulting in the uncoordinated and unsuitable correspondence between the data and the traffic flow.
- With the continuous updating of big data, it provides conditions for traffic travel time prediction, but some advantages and characteristics of these data are not noticed, resulting in many useful data not being used and mined
- Some model parameter calibration is too subjective, which largely depend on researchers' experience;
- Some model mostly aimed at a specific example, and cannot be easily adapted to other situations.

Therefore, in this study, historical data of a freeway toll station were collected, and were categorized using the support vector machine (SVM) algorithm. Although the predecessors have done some work: Wu, Chun-Hsin[20]used the method of support vector regression to predict the time; Vanajakshi[21] obtained the support vector for short-term prediction of travel time by algorithm. Machine technology; Mendes-Moreira [22]obtained a regression method for comparing long-term travel time prediction through intelligent data analysis. But their analysis is based on machine learning algorithms and does not better understand or improve the transportation system. So this paper uses SVM model based on historical data to predict the common traffic state and the method used for model construction was simplified. The practicability of the prediction algorithm was enhanced to overcome assumptions and uncertainties in the existing traffic flow theory.

132 2. Data Collection and Preprocessing

133 2.1 Data Description

134 Data for this study was collected in the FreewayG5513 from Changsha to Yiyang, Hunan
 135 Province, starting from the Changsha toll station and ending at Yiyang toll station. This freeway
 136 G5513 is a standard freeway with a two-way four lane of 100 km/h and a roadbed width of 26 m. The
 137 total length was approximately 63 km, the daily average flow reached 58,000 vehicles, and the peak
 138 flow during long vacation was up to 96,600 vehicles. Because of heavy traffic, the freeway has been
 139 rated as one of the six most congested sections in the Hunan Province. There are nine toll stations
 140 along the road, from east to west, that is Changsha West, Youren, Guanshan, Jinzhou, Ningxiang,
 141 Xiangjiang West, Quanjiao, Chaoyang, and Yiyang North, illustrated as Figure1.



142
 143 **Figure 1.** Layout of the freeway and toll station

144 The main data set collected in this study includes:

- 145 • Toll data of the whole toll stations along G5513 in February 2018 (vehicles entering and leaving
 146 toll station), with a total of 561,081 data items, including the name of the toll station, the time of
 147 vehicle entering and leaving the toll station, vehicle type and weight.
- 148 • Weather information surveying station located near the freeway, which was collected from the
 149 Chinese Weather Network in February 2018, with a total of 672 data items, including 24 hour
 150 daily weather, temperature, relative humidity, precipitation, and wind direction.
- 151 • Freeway blockage record statistics, which was obtained from the freeways management
 152 department, a total of 260 freeway blockage information reports were collected in February,
 153 March, April, and May, including blockage location, reasons for the blockage, blockage start
 154 time, and blockage end time.
- 155 • Freeway traffic control measures report, which was obtained from the Traffic Police
 156 Department, with a total of 7 data items collected on April 5 Qingming traditional
 157 national Festival, May 1 International Labor Day, and other holiday control information.

158 2.2 Data Preprocessing

159 Many abnormal data items were found in the data, which need to be preprocessed before use.

- 160 • Data sharing the same entry and exit toll

161 On the freeway, some drivers turn around in the service area or other sections to avoid the
 162 charges and even exchange the toll tags, which is likely to make the entry and exit of the vehicles at
 163 toll stations consistent.

164 Therefore, it is necessary to determine whether a data item is consistent with the toll gates and
 165 eliminate invalid data items.

166 • Abnormal time record data

167 Owing to the failure of the time system associated with toll station to synchronize or system
 168 failure, the time of accessing the toll station can be earlier than the time of exiting the toll station. In
 169 addition, there are other factors that can lead to long travel time, such as the breakdown of vehicle
 170 on road, accidents, and the situation where driver may have a long rest in service area. All of these
 171 situations will result in unusual time record data.

172 In the process of data preprocessing, abnormal time data record can be eliminated by screening.

173 • Missing data

174 There are two main reasons for missing data: on the one hand, it's mainly from equipment
 175 problems or road environment including the unstable scanning frequency of detector, faulty of
 176 transmission equipment, and traffic jam. On the other hand, eliminating wrong data items will also
 177 lead to partial data missing. Lack of data will cause the road real traffic conditions to change directly
 178 or indirectly. Therefore, it is essential to make up for the missing data for historical data. Because of
 179 the strong continuity of the traffic flow travel time parameter, the trend in the change of the traffic
 180 flow travel time parameter with time is consistent, although its fluctuation will change as the
 181 collection period changes. Therefore, the following data fill formula is obtained, as given in equation
 182 (1).

$$183 \quad data(t) = \frac{3}{6} \times data(t-1) + \frac{2}{6} \times data(t-2) + \frac{1}{6} \times data(t-3) \quad (1)$$

184 Where, $data(t)$ represents the current missing data, $data(t-1)$, $data(t-2)$ and $data(t-3)$ are the traffic
 185 flow travel time data of the past period, two cycles, and three cycles are respectively represented.

186 3. Support Vector Machine Model

187 3.1 Problem Description of Freeway Travel Time Prediction

188 Travel time of a freeway has strong continuity in a certain time range. That is, there are some
 189 complex functional relationships between the current travel time and the past travel time. By
 190 analyzing the changes in travel time, we can obtain rules and establish a real-time prediction model
 191 updated every 5 min. Then, the accuracy and reliability of the predicted results can be improved by
 192 using an optimization algorithm to find the optimal solution of the model.

193 The change in the freeway travel time in different time periods is not a simple linear
 194 relationship, and it will neither increase indefinitely, nor decrease indefinitely. But it will only
 195 change continuously within a floating interval. Therefore, using a simple least squares regression
 196 prediction or similar methods is not sufficient to predict the travel time. The SVM nonlinear
 197 regression theory can be employed to solve this problem.

198 SVM uses nonlinear transformation to map the original variables to a high-dimensional feature
 199 space, so that the problem of nonlinear separability in the original sample space is transformed into
 200 high-dimensional feature space. The linear separable problem and the application of expansion
 201 theorem of the kernel function in the calculation process do not require the explicit expression of
 202 nonlinear mapping. In addition, since the linear learning machine is established in the
 203 high-dimensional feature space, it can be compared to the linear model. The comparison not only
 204 increases the complexity of the calculation, but also solves the problem of "dimensional disaster".

205 Owing to changes in the traffic environment, sudden traffic accidents, weather, and other
 206 special events, the SVM algorithm will eventually be transformed into a quadratic programming
 207 problem. In theory, a global optimal solution can be obtained, thus solving the traditional neural
 208 network. The network can't avoid the local optimal problem, and should adequately accommodate

209 the influencing factors due to these sudden changes to improve the accuracy of the travel time
210 prediction result.

211 Therefore, this study used the nonlinear support vector machine regression theory[23].

212 3.2 Model Overview

213 The SVM is a machine learning method, which is based on the statistical learning theory
214 developed by Vapnik. The theory has been further extended to diversified application algorithms,
215 including the linear SVM classification algorithm, the nonlinear SVM classification algorithm, the
216 linear SVM regression algorithm, and the nonlinear SVM regression algorithm[24]. These SVM
217 algorithms have been widely used in many fields owing to their simple structure and high
218 computational efficiency.

219 Consider a training sample set of l training samples, $S = \{(x_i, y_i) | x_i \in R, y_i \in R\}_{i=1}^l$, which is
220 non-linear, where x_i is the input column vector of the i training sample,
221 $x_i = [x_i^1, x_i^2, \dots, x_i^d]^T$, $y_i \in R$ is the corresponding output of the kernel function
222 $K(x_i, y_i) = \varphi(x_i)^T * \varphi(y_i)$.

223 Equation (1) in the following is the linear regression equation established in the high
224 dimensional space, and ε is introduced as a linear insensitive loss function:

$$225 \quad f(x) = \omega \varphi(x) + b \quad (2)$$

$$226 \quad L(f(x), y, \varepsilon) = \begin{cases} 0, & |f(x) - y| \leq \varepsilon \\ |f(x) - y| - \varepsilon, & |f(x) - y| > \varepsilon \end{cases} \quad (3)$$

227 Where, $\varphi(x)$ is the nonlinear mapping function, $f(x)$ is the prediction function, which returns
228 the predicted value, and y is the corresponding real value.

229 Under the above constraints, we can find the optimal classification hyper-plane, that is, find the
230 solution to the following optimization problem.

$$231 \quad \begin{cases} \min \frac{\|\omega\|^2}{2} \\ s.t. \|\omega^T \varphi(x_i) + b - y_i\| \leq \varepsilon, i = 1, 2, \dots, l \end{cases} \quad (4)$$

232 This problem can be solved by solving the saddle point of the Lagrange function, and its dual
233 theory can be applied to solve the dual problem.

$$234 \quad \begin{cases} \min \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) K(x_i, x_j) + \varepsilon \sum_{i=1}^l (\alpha_i - \alpha_i^*) - \sum_{i=1}^l y_i (\alpha_i - \alpha_i^*) \\ s.t. \sum_{i=1}^l (\alpha_i - \alpha_i^*) = 0, \alpha_i^*, \alpha_i \geq 0, i = 1, 2, \dots, l \end{cases} \quad (5)$$

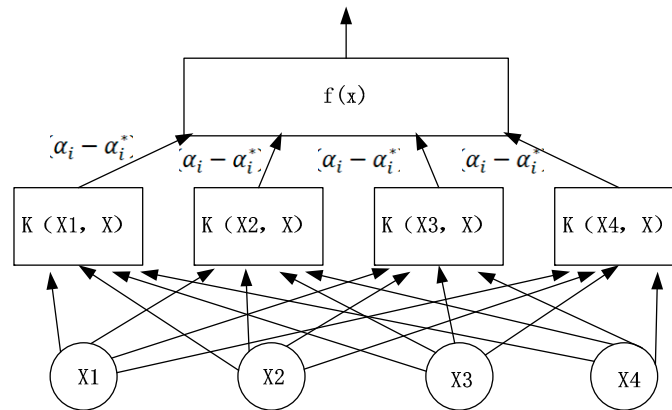
235 To solve the dual problem, a relaxation factor can be set for each data point. After introducing
236 these two relaxation factors, ξ_i, ξ_i^* ($\xi_i, \xi_i^* \geq 0, i = 1, 2, \dots, l$), the function can be optimized as:

$$237 \quad \begin{cases} \min \frac{\|\omega\|^2}{2} + C \sum_{i=1}^l (\xi_i - \xi_i^*) \\ s.t. \omega^T \varphi(x_i) + b - y_i \leq \xi_i + \varepsilon, i = 1, 2, \dots, l \\ s.t. y_i - \omega^T \varphi(x_i) - b \leq \xi_i^* + \varepsilon, i = 1, 2, \dots, l \end{cases} \quad (6)$$

238 In the above equation, C is the penalty factor; the smaller the value, the smaller the penalty to
 239 the error data.

240 Next, the Lagrangian multiplier method can be used to solve the optimization algorithm, and
 241 the nonlinear regression function can be further used to solve the double optimization problem.

$$242 \quad f(x) = \omega \partial(x) + b = \sum_{i=1}^l (\alpha_i - \alpha_i^*) \partial(x_i)^T * \partial(x_j) + b = \sum_{i=1}^l (\alpha_i - \alpha_i^*) K(x_i, y_i) + b \quad (7)$$



243

244

Figure 2. Structure of SVM model

245 For the freeways, there is essentially no significant error in the travel time between the low peak
 246 period and even the flat peak period. However, the impact of different traffic conditions on the
 247 travel time is inevitable during peak periods. Therefore, the two cases should be discussed
 248 separately. Furthermore, whether or not the working day has a different influence on the travel time
 249 of vehicles traveling on the highway. The commuting time, travel purpose, and travel mode will be
 250 also different. Therefore, these two points should be viewed separately. In addition, road weather
 251 conditions and traffic control will have certain influence on the prediction results and should be
 252 considered.

253 Based on the SVM model, input workdays, non-working days, the morning and evening peak
 254 periods, and off-peak hours as original values, finally can get six time periods: the morning peak
 255 hours on workdays, the evening peak hours on workdays, off-peak hours on workdays and so on.
 256 Weather and traffic control factors for the four scenarios can also be analyzed. However, the
 257 difference in highway traffic conditions between the working day and the non-working day, the
 258 morning and evening peak periods and flat peak period is not considered due to the limitation of the
 259 length of the article.

260 This study used the travel time prediction of evening peak hours in the classified working
 261 days and non-working day peak hours as an example by comparing the traffic data of multiple
 262 working days and non-working days.

263 3.3 Model Construction

264 The freeway travel time prediction model is based on the SVM algorithm and is constructed
 265 based on the relationship between the current travel time of the road segment and the past travel
 266 time of the road segment, the current weather, and the possibility of traffic control.

267 In this study, data from two toll stations with different distances from east to west of G5513
 268 were selected for analysis. Moreover, as the first-class passenger car (7 passenger car) accounts for
 269 the vast majority of the data, the travel time of the first-class passenger car was taken as the
 270 prediction object, and the analysis time interval is 5 minutes.

271 The characteristic of the toll station is presented in table 2.

272

273

Table 2. Analysis objects of freeway travel time prediction

Toll station	Start point	End point	Distance (kilometer)
1	Changsha West	Guanshan	10.6
2	Changsha West	Ningxiang	23.2

274 The structure of the SVM is similar to the neural network. The output is a linear combination of
 275 intermediate nodes, and each intermediate node corresponds to a support vector. To determine the
 276 optimal classification function, this study takes the four travel times of the time period before the
 277 prediction time as the input, namely $t_{k-1}, t_{k-2}, t_{k-3}, t_{k-4}$.

$$278 \quad t_k = g(t_{k-1}, t_{k-2}, t_{k-3}, t_{k-4}) \quad (8)$$

279 k is the current time period, and t_k represents the average travel time of all vehicles in the
 280 current predicted time period.

281 In the prediction process, variables such as weather, traffic accident that affects the travel time,
 282 holiday or non-holiday, and day of the week, are evaluated as follows:

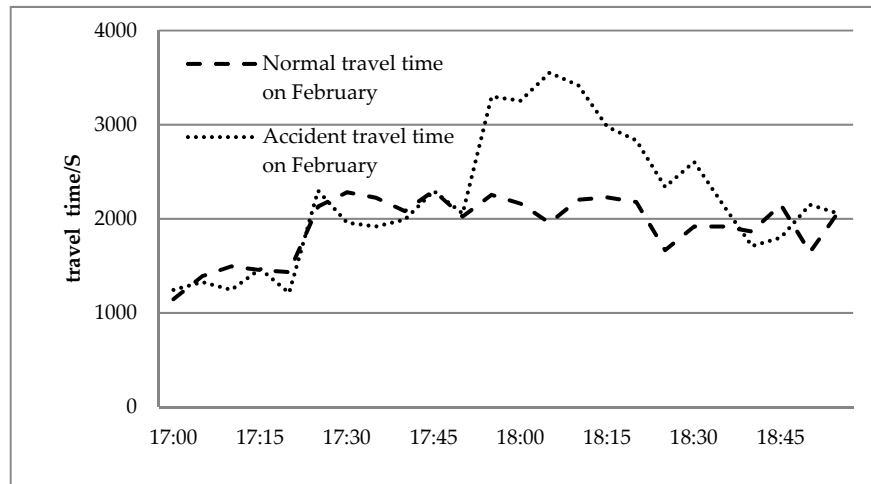
283

Table 3. Variables type and its meanings

Variable category	Variable name	Variable type	Variable value	Variable meaning
Meteorological	Weather status	Discrete	Clear; Cloudy; Fog;	Clear; Cloudy;
			Overcast; Light rain;	Fog; Overcast;
			mod rain; hvy rain	Light rain; mod rain; hvy rain
Time	Holiday	Discrete	0	N
			1	Y
	Weekly	Discrete	1	Monday
		
		7	Sunday	
Accident	Affecting travel time	Discrete	0	N
	traffic accidents		1	Y

284 Many traffic accidents occur on the freeway every day. To simplify the parameters, we divide
 285 traffic accidents into two categories: traffic accidents that affect the travel time and those that do not.
 286 In this study, traffic accidents that affect the travel time was regarded as invariant.

287 The following figure shows a comparison of travel time affected by an accident and normal
 288 travel time on the afternoon of February 17.



289

290 **Figure 3.** Comparison of accidents and normal travel time

291 Since SVM is a machine learning model, sample training is required before prediction. Multiple
 292 groups of any four consecutive travel times, daily weather, traffic accidents that affect the travel
 293 time, holidays, and weekday data were used as training samples to obtain a trained model. In the
 294 trained model, $t_{k-1}, t_{k-2}, t_{k-3}, t_{k-4}$, weather, accident, holiday, and week are used to predict the
 295 travel time of the next time period. When a certain number of training samples is achieved, real-time
 296 data input can be adopted to predict future results. Moreover, the model can be constantly modified
 297 based on the relation between the predicted data and the predicted data to prediction accuracy.

298 3.4 Parameter Calibration and Optimization

299 Parameter selection is very important to find the optimal hyper-plane in the SVM model used in
 300 this study. Existing studies mainly adopted the traditional grid search method, direct determination
 301 method, one-dimensional search method, and inverse ratio method to determine the insensitive loss
 302 function parameter ε and penalty parameter C . However, there are many shortcomings associated
 303 with these methods, and the resulting errors will significantly influence the accuracy of the
 304 prediction results.

305 Moreover, in the SVM model, kernel function selection is also an important factor that
 306 influences the performance of the SVM. The radial basis kernel function $K(x_i, y_i) = e^{-\frac{\|x-x_i\|^2}{\sigma^2}}$ (RBF) is
 307 an adaptive kernel function for low-dimensional space data and high-dimensional space, which
 308 have good convergence domains, and this function can be described as an ideal kernel function.
 309 Therefore, the RBF was selected as the classification prediction kernel function of the SVM, in which
 310 a kernel parameter σ needs to be optimized.

311 Therefore, three parameters need to be optimized, namely the core parameter σ , the
 312 non-sensitive loss function parameter ε , and the punishment parameter C . The kernel parameter σ is
 313 the distribution or range of the training sample data. The non-sensitive loss function parameter ε
 314 affects the number of support vectors. The larger the value of ε , the lower the regression precision,
 315 and the fewer the support vectors. The penalty parameter C is used to control the degree of
 316 punishment of samples beyond the allowable error range. The higher the value, the heavier the
 317 punishment of samples.

318 We used the artificial fish swarm algorithm to optimize the parameters of the regression model.
 319 The artificial fish swarm algorithm has unique advantages in parameter optimization and
 320 overcomes the blindness of traditional algorithms in parameter optimization and the defects of the
 321 linear model and neural network in parameter selection. It can be said that the parallel performance
 322 of the artificial fish swarm algorithm can ensure that the model parameters converge faster to the
 323 global optimization extreme[24,25].

324 The first step in the optimization process of the artificial fish swarm algorithm is to feed in the
 325 training value and the training target through the SVM model to calculate the fitness of the
 326 individual. The most adaptable individual is regarded as the optimal value of the current fish group
 327 and the corresponding parameters σ , ε , and C of the current optimal value are saved. In the
 328 subsequent iteration, σ , ε , and C corresponding to the maximum fitness value are taken as the final
 329 optimization results.

330 4. Case Study

331 4.1 Data Selection

332 The data used in this study was collected in February 2018 on G5513 (from Changsha West to
 333 Guanshan/Ningxiang Station)in Changsha, Hunan Province, China. The travel time was detected for
 334 all days and the detection interval is 5 minutes. Furthermore, 288 sequences are included in one day.

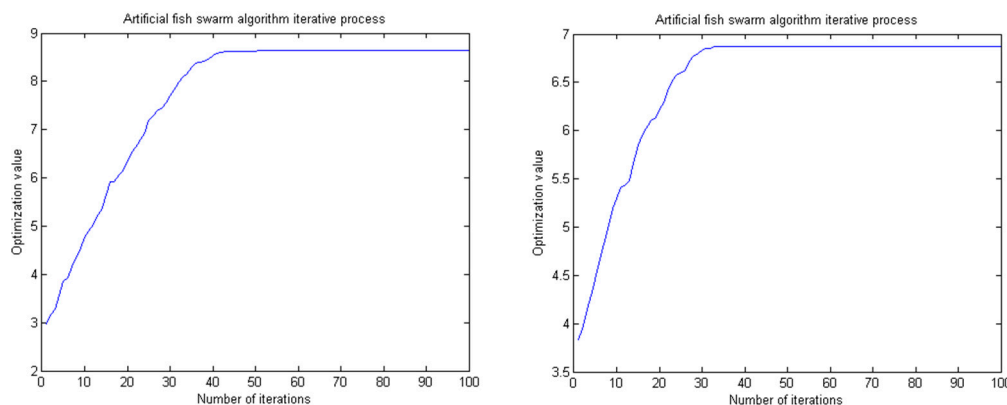
335 The daily evening peak (17:00-19:00) data of G5513 (Changsha West to Guanshan /Ningxiang
 336 Station) was selected as an example after comparing data of multiple working days and
 337 non-working days, which contains 204 to 228 items. Other variable data items also need to be filtered
 338 according to the above data. The dimension feature values are based on the time series and the data
 339 requirements according to the prediction model.

340 In this study, the regression SVM model is used to establish the model parameters, and the
 341 artificial fish swarm algorithm is used to establish the model parameter optimization algorithm. The
 342 optimization results are presented in Table 4. The optimization process for the optimal value of the
 343 penalty parameter C is shown in Figure 4.

344 **Table 4.** Optimization of parameter values

Section	Penalty parameter, C	Nuclear parameter, σ	Insensitive loss fun ction parameter, ε
Changsha West-Guanshan	6.8755	0.0064	0.3461
Changsha West-Ningxiang	8.6485	0.0034	0.6991

345



346

347 **Figure 4.** The parameter optimization curve

348 The parameters of the artificial fish swarm algorithm are set as follows: the maximum number
 349 of iterations of the artificial fish is 100; the population size is 5; the maximum number of trials is 5;
 350 the crowding factor δ is 0.618; the perceived distance is 0.5; and the moving step is 0.1.

351 4.2 Results and Comparative Analysis

352 The data adopted in this study were obtained during the Chinese Spring Festival from
 353 February 15 to 21, 2018. Therefore, the data was divided into working days and holidays. There were

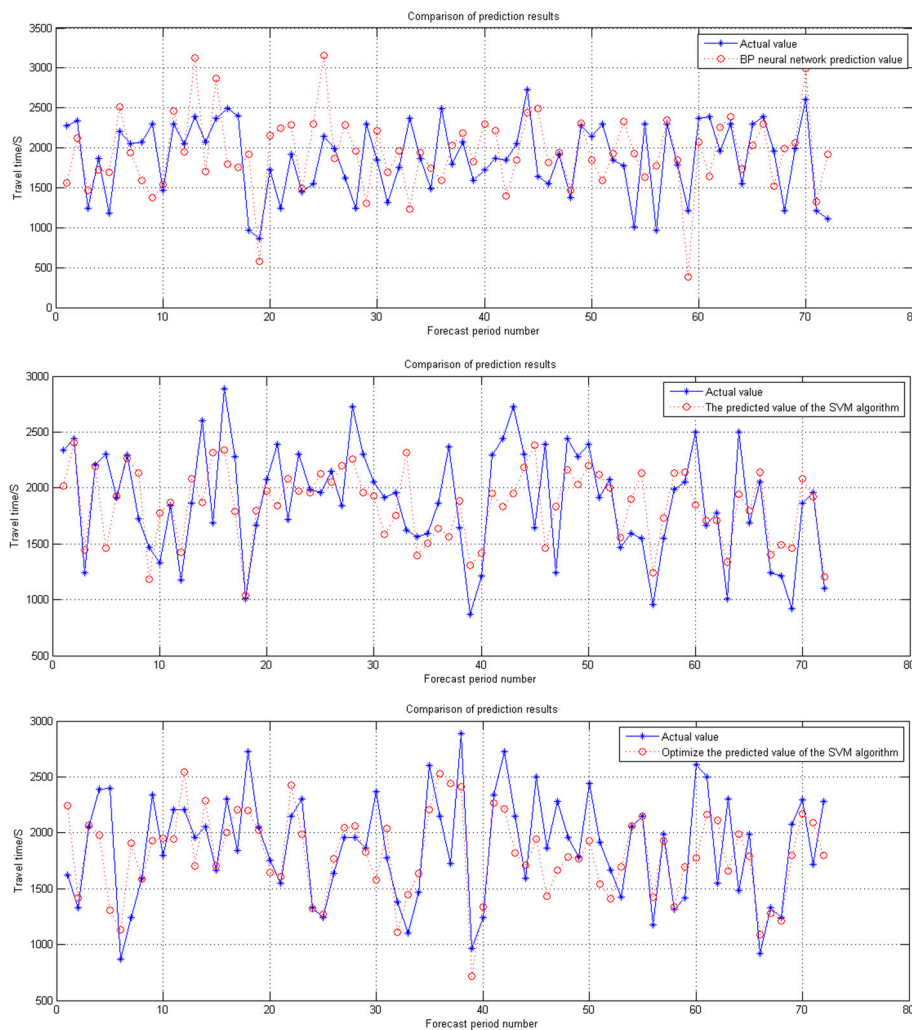
354 336 sets of data from February 1 to 14, 2018 (14 days data) in each group of toll stations; 264 groups (11
 355 days) were randomly selected as training data input, and 72 groups (3 days) were adopted as
 356 prediction numbers. There were 168 sets of data from February 15 to 21, 2018 (11 days data) in each
 357 group of toll stations; 96 groups (4 days) were randomly selected as training data input, and 72
 358 groups (3 days) were adopted as prediction numbers.

359 The detection time is from 17:00 to 19:00 P.M. and the value was taken as test data.

360 BP neural network, SVM, and optimized SVM were used for the prediction. The root mean
 361 square error (RMSE), the mean absolute percent error (MAPE) and the covariance protocol (CP)
 362 were selected as the error evaluation criteria in the prediction process [26].

363 The RMSE is a comprehensive evaluation indicator of the prediction effect, the MAPE is the
 364 prediction relative error, while the CP is the error component analysis indicator.

365 The following figures show the forecasting effect diagram of the holiday evening peak.



366

367

Figure 5. Results of freeway travel time prediction

368

Table 5. Error evaluation of forecasting freeway travel time

Method	Changsha West- Guanshan Station			Changsha West-Ningxiang Station		
	RMSE	MAPE	CP	RMSE	MAPE	CP
Working day BP neural network	7.6544	5.0873	0.5489	10.5124	5.5104	0.5809
SVM	6.3405	4.8225	0.5810	9.0375	4.5422	0.6930

	Optimized SVM	6.0369	4.6536	0.5952	8.3069	4.4524	0.6404
	BP neural network	11.6308	9.1023	0.5990	16.0845	9.6148	0.6454
Holiday	SVM	11.5152	7.8645	0.6460	13.3215	8.1153	0.6905
	Optimized SVM	9.6218	6.2451	0.7621	12.2548	7.8651	0.7245

369 From Table 5, it can be observed that the working day was predicted by the RMSE, the MAPE,
 370 and the CP data between the two toll stations. It could be found that all three models can be used for
 371 predicting travel time.

372 Although the prediction error of the BP neural network may be larger than those of the SVM
 373 and the optimized SVM models, there is no deviation between the error of the SVM and the
 374 optimized SVM model.

375 However, when forecasting holidays with large traffic and long travel time, the RMSE of the
 376 optimized SVM model is significantly better than those of the BP neural network and the SVM
 377 model.

378 In the prediction of Changsha West-Guanshan, the accuracy of the optimized SVM model using
 379 artificial fish swarm algorithm is 17.27% higher than that of the BP neural network model and
 380 16.44% higher than the conventional SVM model.

381 In the prediction of Changsha West-Ningxiang, the accuracy of the optimized SVM model
 382 using artificial fish swarm algorithm is 23.80% and 8.01% higher than those of the BP neural network
 383 model and the conventional SVM model, respectively.

384 The optimized SVM model described in this paper has higher travel time prediction accuracy in
 385 the road segment, and the mapping law of the input and output are better represented by the
 386 optimized SVM model.

387 In terms of the relative prediction errors of the three prediction models, the MAPE of the
 388 optimized SVM model is lower than the prediction errors of the BP neural network and the
 389 conventional SVM model when using holiday and working day data. This indicates that the
 390 optimized SVM model described in this paper has certain advantages in terms of the travel time
 391 prediction model of the road segment, and the data requirements are lower.

392 4.3 Analysis of Influencing Factors of Travel Time

393 The freeway travel time is determined by various factors, such as weather, traffic accident,
 394 holiday, and week day. However, owing to limitations of sample data, only traffic accidents and
 395 holidays were considered.

396 4.3.1 Effect of traffic accidents on travel time

397 First, the optimized SVM model described in this paper is superior to the BP neural network
 398 and SVM model in terms of the CP.

399 In the prediction of Changsha West-Guanshan, the CP of the optimized SVM model is 21.40%
 400 higher than that of the BP neural network model, which is 7.28% higher than that of the conventional
 401 SVM model;

402 In the prediction of Changsha West-Ningxiang, the CP of the optimized SVM model is 10.9%
 403 and 6.53% higher than those of the BP neural network model and the conventional SVM model,
 404 respectively.

405 The results presented in this paper indicates that the optimized SVM model has better
 406 inclusiveness and stability when unexpected factors such as traffic accidents that affect the travel
 407 time are encountered, thereby avoiding the need for repeated trial and error to address network
 408 problems.

409 Second, freeway traffic accidents will cause the traffic capacity of certain sections of the road
410 network to decrease, and queues will be formed near the accident site, which increases the travel
411 time of the vehicle.

412 4.3.2 Effect of holidays on travel time

413 A comparison of the working day and holiday forecasting error evaluation criteria presented in
414 the previous section indicates that the BP neural network has a larger prediction error than the
415 working day prediction results of the other models probably due to the problem of construction of
416 the network structure. However, the conventional SVM and the optimized SVM models have similar
417 prediction error results. In the above analysis, there are large gaps in the holiday prediction results
418 of the three different models.

419 It can be observed that holidays have significant influence on the travel time. The effect of holidays
420 on the travel time, which was obtained from the analysis of the original toll station data, is that it
421 significantly increases the volume of traffic on the highway network.

422 5. Conclusions

423 This study performed an in-depth analysis of freeway travel time prediction to provide
424 high-quality travel experience for users and found that the freeway travel time is affected by travel
425 time, weather, and traffic. The effect of different factors was analyzed, such as accidents and
426 holidays.

427 Bad weather reduces the overall traffic rate, which increases the travel time. Traffic accidents
428 lead to reduced road traffic capacity, which affects the travel time. Free passage on highways during
429 holidays and the increased demand for travel result in increased vehicle flow, which also affects the
430 travel time.

431 In this study, basic data was analyzed, and the traffic state prediction method based on SVM
432 data mining technology was proposed to transform the problem into a quadratic programming
433 problem using artificial fish swarm algorithm, which reduces the computational and local optimal
434 problems of traditional neural networks. The parameters of the SVM were optimized using
435 traditional network optimization, and a global optimal solution was obtained.

436 Results show that the accuracy of the optimized SVM model is 17.27% and 16.44% higher than
437 those of the BP neural network model and the conventional SVM model, respectively. Accurate
438 prediction of the travel time on the freeway was realized, which can provide data support for
439 monitoring, early warning, and decision analysis for the freeway operation status.

440 In this study, influencing factors such as weather, traffic, accidents and holidays were included
441 in the optimization of the SVM prediction model. However, owing to limitations of the number of
442 samples, the model was not fully trained. Therefore, a certain error occurred in the prediction
443 results.

444 In the future, it is necessary to categorize traffic accidents, clarify the impact of each type of
445 accident on the travel time, categorize increase in the holiday traffic, and clarify the impact of each
446 level of accident on the travel time. Furthermore, the number of training samples, database capacity,
447 and prediction accuracy should be continuously increased. In this study, only data from freeway toll
448 stations was validated, and application to actual large-scale road networks should be further
449 explored in the future.

450

451 **Author Contributions:** This work was conducted by Kejun Long and Wei Wu with the help of graduate
452 student Wukai Yao. It was mainly drafted by Kejun Long and Wukai Yao, and checked and revised by Wei
453 Wu and Jian Gu. Kejun Long and Wukai Yao designed and analyzed the proposed model. Wukai Yao and Jian
454 Gu performed the simulation. Kejun Long is responsible for the English polish and proofreading of the work.

455 **Funding:** This research was funded by National Natural Science Foundation of China (NSFC), grant number
456 51678076, Hunan Provincial Key Laboratory of Smart Roadway and Cooperative Vehicle-Infrastructure
457 Systems, grant number 2017TP1016.

458 **Acknowledgments:** The authors express their thanks to all who participated in this research for their
459 cooperation. The authors would like to give great thank to the hard work by the peer reviewers and editor.

460 **Conflicts of Interest:** The authors declare no conflict of interest.

461 References

- 462 1. Gipps,P.G. The estimation of a measure of vehicle delay from detector output. Newcastle-Upon-Tyne
463 University, England1977, 18 p.
- 464 2. Yildirimoglu, M; Geroliminis, N. Experienced travel time prediction for congested freeways. *Transp. Res.*
465 *Part B-Method*, 2013, 53(4):45-63.
- 466 3. Shen,L; Hadi,M. Practical approach for travel time estimation from point traffic detector data. *J. Adv*
467 *Transportation*, 2013, 47(5):526-535.
- 468 4. Hyun ,K; Tok,A, Ritchie S G. Long distance truck tracking from advanced point detectors using a selective
469 weighted Bayesian model. *Transport. Res. Part C-Emerg. Technol*, 2016, 82:24-42.
- 470 5. Hyun, K. K.; Jeong, K. Assessing crash risk considering vehicle interactions with trucks using point
471 detector data. *Accid. Anal. Prevent*, 2018, 17p.
- 472 6. Ramezani, M.; Geroliminis, N. On the estimation of arterial route travel time distribution with Markov
473 chains. *Transp. Res. Part B*, 2012, 46(10): 1576-1590.
- 474 7. Zhang, J. T.; Zhou, J. An Arterial Travel Time Estimation Model Based on Discrete Time Markov Chains.
475 *Syst. Eng*, 2014,5:98-104.
- 476 8. Zhang, J. T.; Zhou, J. Travel time estimation model based on spatial Markov chains. *Syst. Eng*,2015,12
477 :72-77.
- 478 9. Woodard, D.; Nogin,G. Predicting travel time reliability using mobile phone GPS data. *Transport. Res.*
479 *Part C-Emerg. Technol*, 2017, 75:30-44.
- 480 10. Bahuleyan, H.;Vanajakshi, L. D. Arterial path-level travel-time estimation using machine-learning
481 techniques. *J. Comput. Civil. Eng*, 2016, 31(3), 04016070.
- 482 11. Tan, Y. Current Situation of Short-term Flow Forecasting and Discussion on Forecasting Method in Big
483 Data Environment. *ITS China*:2017:10.
- 484 12. Lin J W C V. Incremental and online learning through extended kalman filtering with constraint weights
485 for freeway travel time prediction. *ITSC. IEEE*, 2006:1041-1046.
- 486 13. Zhou, J.; Zhang, C.B. Travel Time Prediction Model for Urban Road Network based on Multi-source Data.
487 *Procedia - Social and Behavioral Sciences*, 2014,138.
- 488 14. Tang-Hsien Chang; Albert Y. Chen. Freeway Travel Time Prediction Based on Seamless Spatio-temporal
489 Data Fusion: Case Study of the Freeway in Taiwan. *Transportation Research Procedia*,2016,17.
- 490 15. Fei X; Lu C C. A Bayesian dynamic linear model approach for real-time short-term freeway travel time
491 prediction. *Transp. Res. Part C*, 2011, 19(6):1306-1318.
- 492 16. Zhan, X.; Ukkusuri, S. V. A Bayesian mixture model for short-term average link travel time estimation
493 using large-scale limited information trip-based data. *Autom.Constr*,2016, 72, 237-246.
- 494 17. Wosyka,J; Přebyl, P. Real-time travel time estimation on highways using loop detector data and license
495 plate recognition. *Elektro. IEEE*, 2012:391-394.
- 496 18. Innamaa, S. Short-Term Prediction of Travel Time using Neural Networks on an Interurban Highway.
497 *Transportation*, 2005, 32(6):649-669.
- 498 19. Lint, J.W.C.V; Hoogendoorn, S P. Accurate freeway travel time prediction with state-space neural
499 networks under missing data. *Transp. Res. Part C*, 2005, 13(5):347-369.
- 500 20. Wu, Chun-Hsin; Jan-Ming Ho. Travel-time prediction with support vector regression.
501 *IEEE.Trans.Intell.Transp.Systems*.5.4 ,2004: 276-281.
- 502 21. Vanajakshi, L; Rilett, L R. Support Vector Machine Technique for the Short Term Prediction of Travel
503 Time. *Intelligent Vehicles Symposium. IEEE*, 2007:600-605.
- 504 22. Mendes-Moreira, João. Comparing state-of-the-art regression methods for long term travel time
505 prediction. *Intell. Data. Anal.* 16.3 (2012): 427-449.
- 506 23. Wang, X.;Chen, X. H.; Yang, X. M. Short term prediction of expressway travel time based on k nearest
507 neighbor algorithm.*Chin.J.Highw.Transport*,2015,28(1), 102-111.
- 508 24. Li, S.; Yuan, Z. C.; Wang, C. Optimization of support vector machine parameters based on group
509 intelligence algorithm. *CAAI TIS*,2018,13(01):70-84.

- 510 25. Wang, Q.; Liu, Z.; Peng, Z. A PSO-SVM Model for short-term travel time prediction based on Bluetooth
511 Technology. *J. Harbin. Inst. Technol.*, 2015, 22(3), 7-14.
- 512 26. Yang, Z. S. Study on the Synthetic Link Travel Time Prediction Model of Key Theory of
513 ITS. *J. Traff. Transp. Eng.*, 2001, 01:65-67.