

# Car-following behavior modeling driven by small data sets based on mnemonic extreme gradient boosting framework

Baichuan LOU<sup>1</sup>, Yufang LI<sup>1,2\*</sup>, Xiaoding LU<sup>1</sup> & Zhe XU<sup>2</sup>

<sup>1</sup>Department of Vehicle Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China;

<sup>2</sup>Key Laboratory of Advanced Manufacture Technology for Automobile Parts (Chongqing University of Technology), Ministry of Education, Chongqing 400000, China

Received 17 April 2020/Revised 29 June 2020/Accepted 1 August 2020/Published online 30 March 2021

**Citation** Lou B C, Li Y F, Lu X D, et al. Car-following behavior modeling driven by small data sets based on mnemonic extreme gradient boosting framework. *Sci China Inf Sci*, 2022, 65(6): 169203, <https://doi.org/10.1007/s11432-020-3044-6>

Dear editor,

In recent years, data-driven car-following models have been developed based on their ability to drill down to information in driving data and their flexibility. According to a study by Fleming et al. [1], the main limitation of existing methods is an insufficient amount of natural driving data. In addition, there are many situations during actual driving processes, such as those under extreme conditions [2] and in the early stages of car-following behavior modeling. In these cases, sparse data learning algorithms are indispensable.

As deep learning has developed, data-driven car-following models based on recurrent neural network and long short-term memory (LSTM) neural network algorithms have become key points of interest in current studies. Chen et al. [3] has summarized the latest deep learning models in the area of driver behavior modeling and highlighted the great potential of LSTM applied to driver behavior modeling.

Recent data-based car-following behavior modeling studies have contributed much to developments in this field. However, this research is mostly supported by large data sets. Although a smaller number of samples generally limit the performance of learning algorithms, there are some algorithms that are specifically designed to learn from sparse samples. Chen et al. [4] proposed a gradient boosting decision tree based on the extreme gradient boosting (XGBoost) learning framework, which is a well-known method for solving complex prediction issues with sparse data.

In this study, further research is reported on sparse data-driven car-following model based on the memory units of LSTM and the XGBoost framework. A test car with intelligent sensors is shown in Figure 1(a). The sampling frequency is 1 Hz. The Cell structure of LSTM and the XGBoost framework are integrated as mnemonic XGBoost (M-XGBoost), which is accurate for the car-following behavior model trained by small data sets and benefits from the capability of XGBoost to process sparse samples and the memory

ability of LSTM. Furthermore, personalized dynamic constraints are designed to identify and counteract iterative errors of memory units (Cells) to ensure driving safety, driving comfort, and optimal vehicle energy consumption. Hence, representations of different car-following styles are achieved.

*XGBoost car-following model.* The main structure of XGBoost is composed by a series of classification and regression trees, as shown in Figure 1(b). The inputs of this model include the relative speed  $\Delta v$ , relative distance  $\Delta s$ , and velocity  $v$ , and the output is acceleration  $a$ .

The objective function with regularization coefficients is expressed by (1), where  $K$  is the number of subtrees (weak learners),  $T$  is the number of leaf nodes,  $\hat{y}_i$  is model output,  $l(y_i, \hat{y}_i)$  is the residual of the current model,  $\omega_j$  is the weight of the corresponding leaf node, and  $\gamma$  and  $\lambda$  are regularization coefficients:

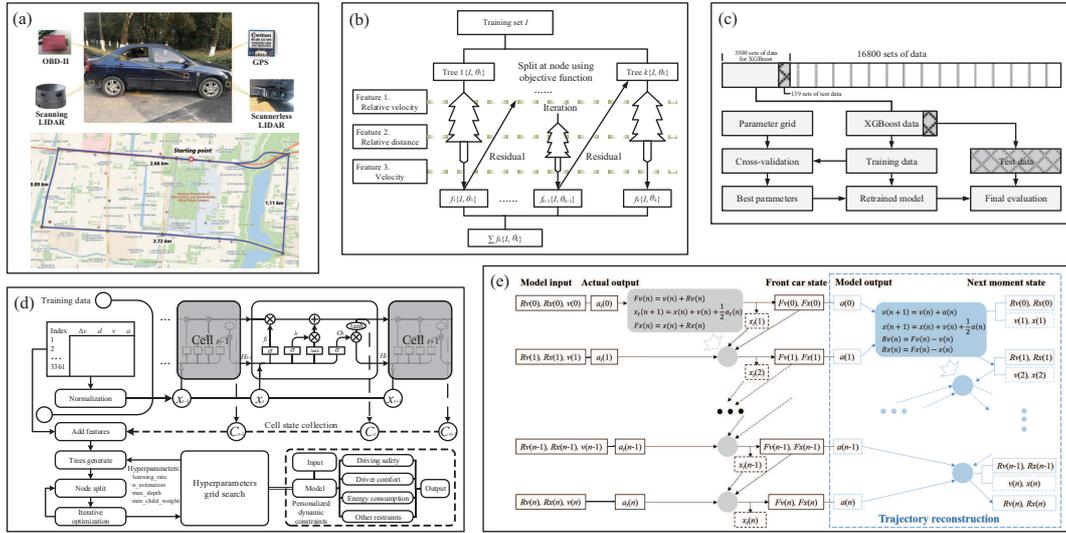
$$obj(\theta) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \left( \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T \omega_j^2 \right) + \text{constant}. \quad (1)$$

The processes of model training and parameter optimization are shown in Figure 1(c). Compared with other popular algorithms, the XGBoost is considerably better at sparse data learning, and the comparison is given in Appendix A. There are also some defects: the prediction process of car-following is a time series issue but XGBoost does not consider autocorrelation and misalignment of time series, which restricts the upper limit of this algorithm.

*M-XGBoost.* To remedy the limitations of XGBoost on the memory of time series, the memory units (Cells) of LSTM [5] are added to the XGBoost framework, and constraints are considered according to real driving experience. The following research is based on this improved XGBoost framework, M-XGBoost.

The basic framework of the M-XGBoost is shown in Figure 1(d). The training process of M-XGBoost is as follows:

\* Corresponding author (email: lyf2007@nuaa.edu.cn)



**Figure 1** (Color online) (a) The test car and the experiment route; (b) the structure of XGBoost car-following model; (c) model training and parameter optimization; (d) M-XGBoost framework; (e) trajectory reconstruction.

- (a) Normalize the original training data;
- (b) Train the Cells of LSTM using normalized data;
- (c) Get the memory factor ( $C_n$ ) passed between Cells;
- (d) Synchronize  $C_n$  with the original data;
- (e) Train the XGBoost model;
- (f) Do GridSearch and cross-validation.

The prediction process of M-XGBoost is as follows:

- (a) Normalize the test data;
- (b) Put the normalized data into the Cells, run the model, and obtain the  $C_n$  of test data;
- (c) Synchronize  $C_n$  with the test data as a new data set;
- (d) Run the XGBoost model;
- (e) Use personalized dynamic constraints to adjust the prediction results as the final output.

**Personalized dynamic constraints.** Personalized dynamic constraints contain the driving safety, energy consumption, and driving comfort factors, as shown in Figure 1(d), which help to confine the prediction results and ensure the driving safety and the diversity of car-following characteristics. The specific constraints can be found in Appendix B.

**M-XGBoost verification.** Trajectory reconstruction is used to verify the follow-up ability of the predicted results. The overall process is to iteratively calculate the state of the next position based on the state of the current position as shown in Figure 1(e). The prediction results of the XGBoost, LSTM, the M-XGBoost and constrained M-XGBoost models are compared with the actual values. It is assumed that the driver wants to follow the car ahead comfortably and safely, and maintain a relative distance of 60 m, as personalized dynamic constraints.

Compared with the original M-XGBoost, although the constrained results have some offset from the actual values, the constrained M-XGBoost improves the accuracy of trajectory tracking. Further, these results reflect a different driving style (following distance 60 m) from the actual value (following distance 30 m), which achieves switching of driving styles based on the M-XGBoost model trained by a small data set. The energy consumption values of the models are equivalent to the superposition of instantaneous power consumption from the tires. Specifically, the follow-up performance of the constrained M-XGBoost model is better than

that of the original M-XGBoost in the case of consumption of the same amounts of energy. The specific processes can be found in Appendix C.

**Conclusion.** In this study, to solve the problem of small data-based car-following behavior modeling, the M-XGBoost framework is proposed, which involves the memory units of LSTM and the XGBoost framework. It is verified that the M-XGBoost framework can achieve sparse data-based car-following model training accurately. On the basis of these results, personalized dynamic constraints are implemented to counteract the iteration errors generated by M-XGBoost and ensure driving safety, comfort and fuel efficiency. Hence, this also yields prediction results for different driving characteristics.

**Acknowledgements** This work was supported by Opening Foundation of Key Laboratory of Advanced Manufacture Technology for Automobile Parts, Ministry of Education, China (Grant No. 2019KLMT05) and Natural Science Foundation of Chongqing (Grant No. cstc2019jcyj-msxmX0119).

**Supporting information** Appendixes A–C. The supporting information is available online at [info.scichina.com](http://info.scichina.com) and [link.springer.com](http://link.springer.com). The supporting materials are published as submitted, without typesetting or editing. The responsibility for scientific accuracy and content remains entirely with the authors.

**References**

- 1 Fleming J M, Allison C K, Yan X D, et al. Adaptive driver modelling in ADAS to improve user acceptance: a study using naturalistic data. *Saf Sci*, 2019, 119: 76–83
- 2 Bärngman J, Boda C N, Dozza M. Counterfactual simulations applied to SHRP2 crashes: the effect of driver behavior models on safety benefit estimations of intelligent safety systems. *Accid Anal Prev*, 2017, 102: 165–180
- 3 Chen S T, Jian Z Q, Huang Y H, et al. Autonomous driving: cognitive construction and situation understanding. *Sci China Inf Sci*, 2019, 62: 081101
- 4 Chen T Q, Guestrin C, Yan X. XGBoost: a scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, 2016. 785–794
- 5 Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput*, 1997, 9: 1735–1780