

Using GPS Data to Predict Accident Risk: Evidence from Beijing

Kanyao Han
University of Chicago

Summary

The popularization of telematics technology over the last decade have ushered in a new era for accident analysis since it can provide two critical sets of always-on information pertinent to accident analysis and prediction: individual-level moving behavior and sociodemographic characteristics. In this paper, I extract them from a car GPS dataset from a telematics company in Beijing, merge them with an insurance dataset from an insurance company, and test their role in accident analysis. I find the features extracted from the car GPS dataset can not only help us better analyze the probability of car accident, but improve the prediction of accident loss.

Objectives

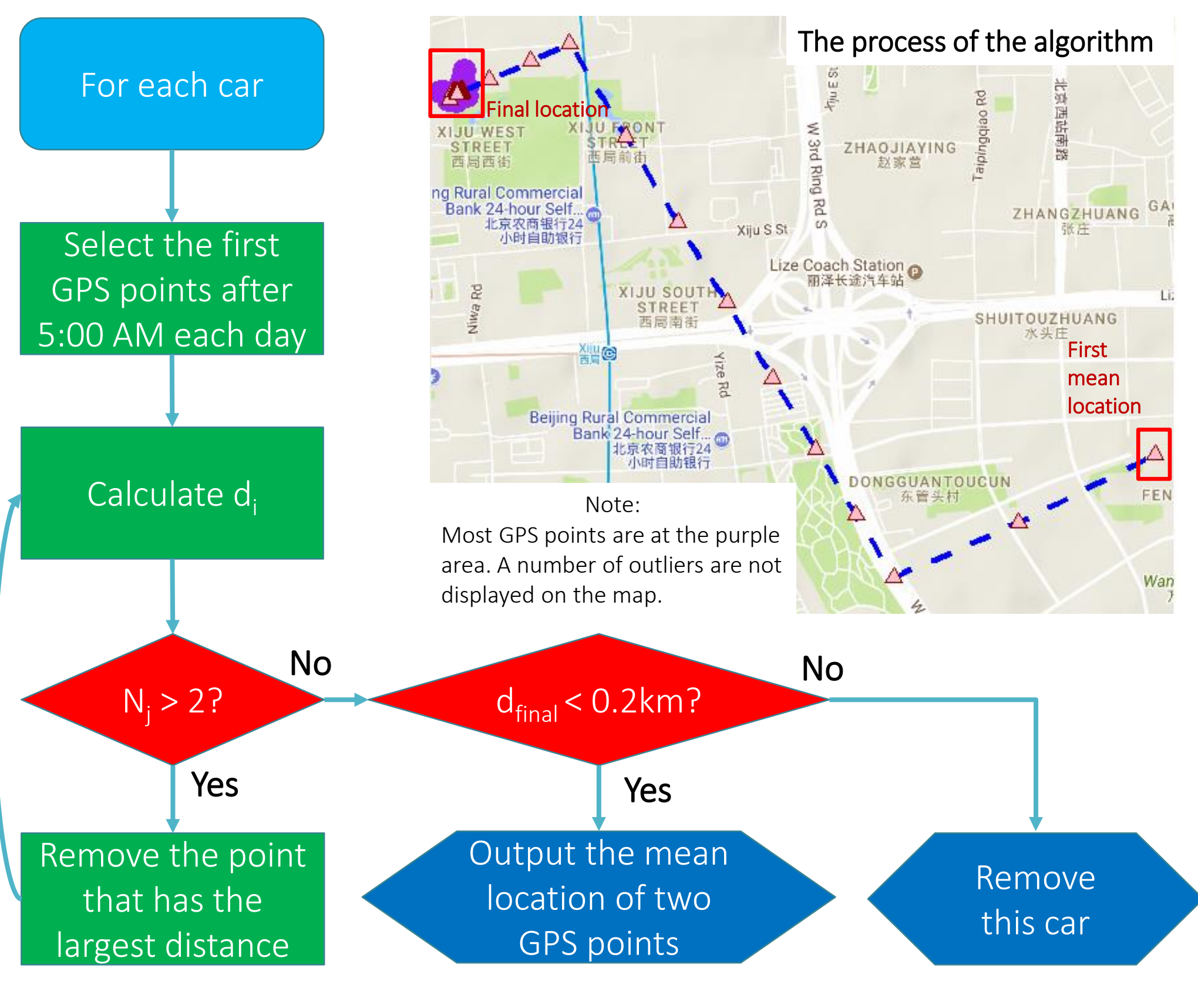
- Design a method to extract home address information from the GPS dataset.
- Further mine useful information from home address for prediction
- Extract driving behavior information from the GPS dataset
- Build predictive models for accident analysis and prediction

Data

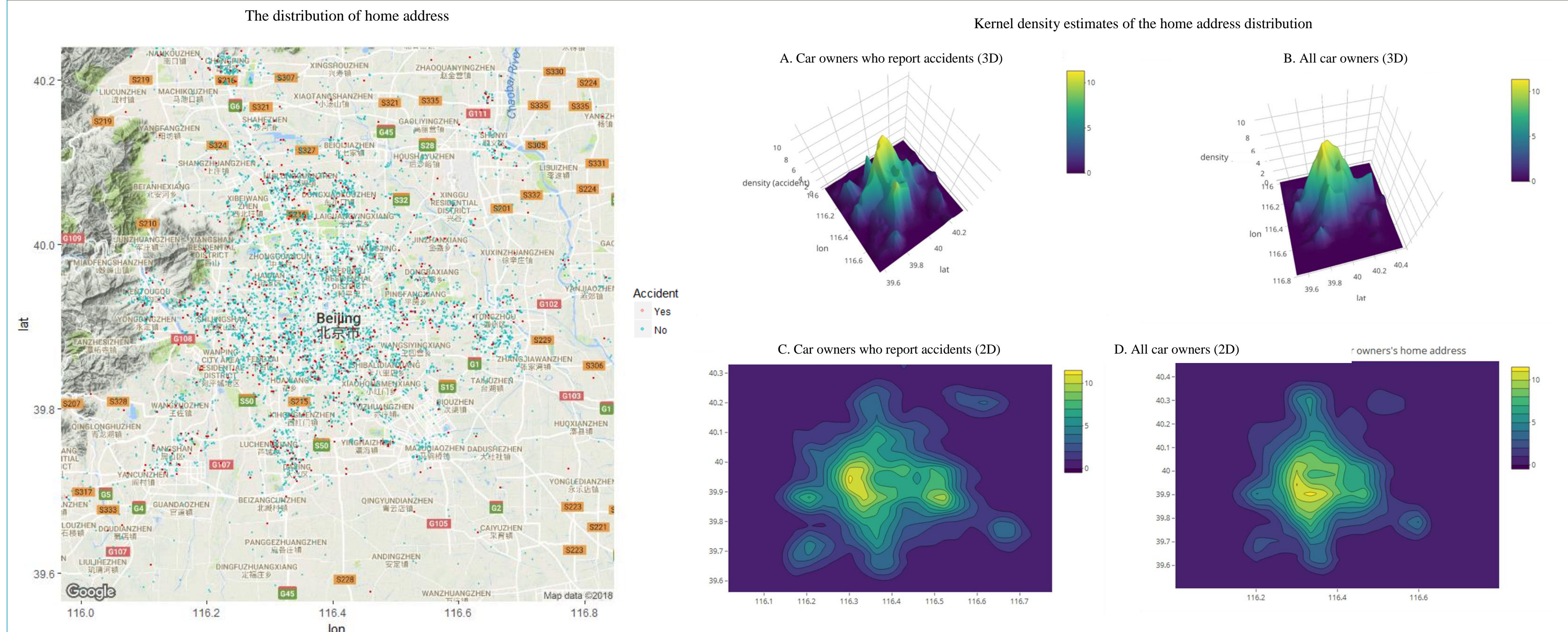
- **GPS DATA:** a three-month record of more than 10,000 cars, including GPS points and driving behavior records such as mileage and speed.
- **INSURANCE DATA:** claim history (used for measuring accident), insurance policy, a limited number of vehicle owners' demographic characteristics and some vehicle-related characteristics such as car price and type.

Identifying Home Address

For each car:
 d_i is the distance from GPS point i to the mean center of all points
 N_j is the number of GPS points in period j
 d_{final} is the distance between two GPS points when $N_j = 2$

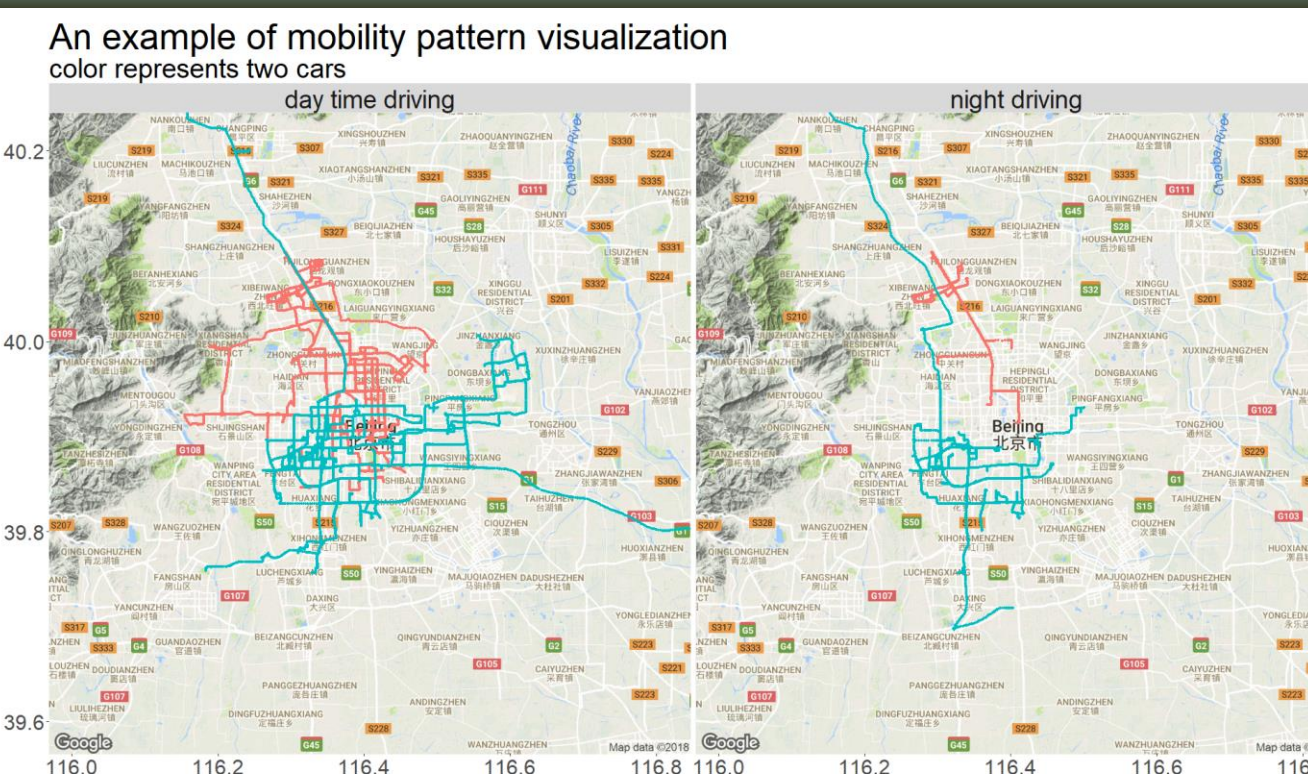


The Distribution of Car Owners' Home Address



The distributions of all car owners' home address and that of those who report accident are similar since more residents signifies more accidents in terms of quantity instead of proportion. In other words, the differences between them may imply different probability of accident. According to the maps and the socioeconomic layout in Beijing, I speculate that Haidian District, the south part of Beijing and northern resident area might provide some important information for accident prediction.

Mobility Pattern and Driving Behavior



The mobility pattern and driving behavior can be easily obtained because the GPS dataset also directly contains driving parameters. Generally, they can be divided into three groups:

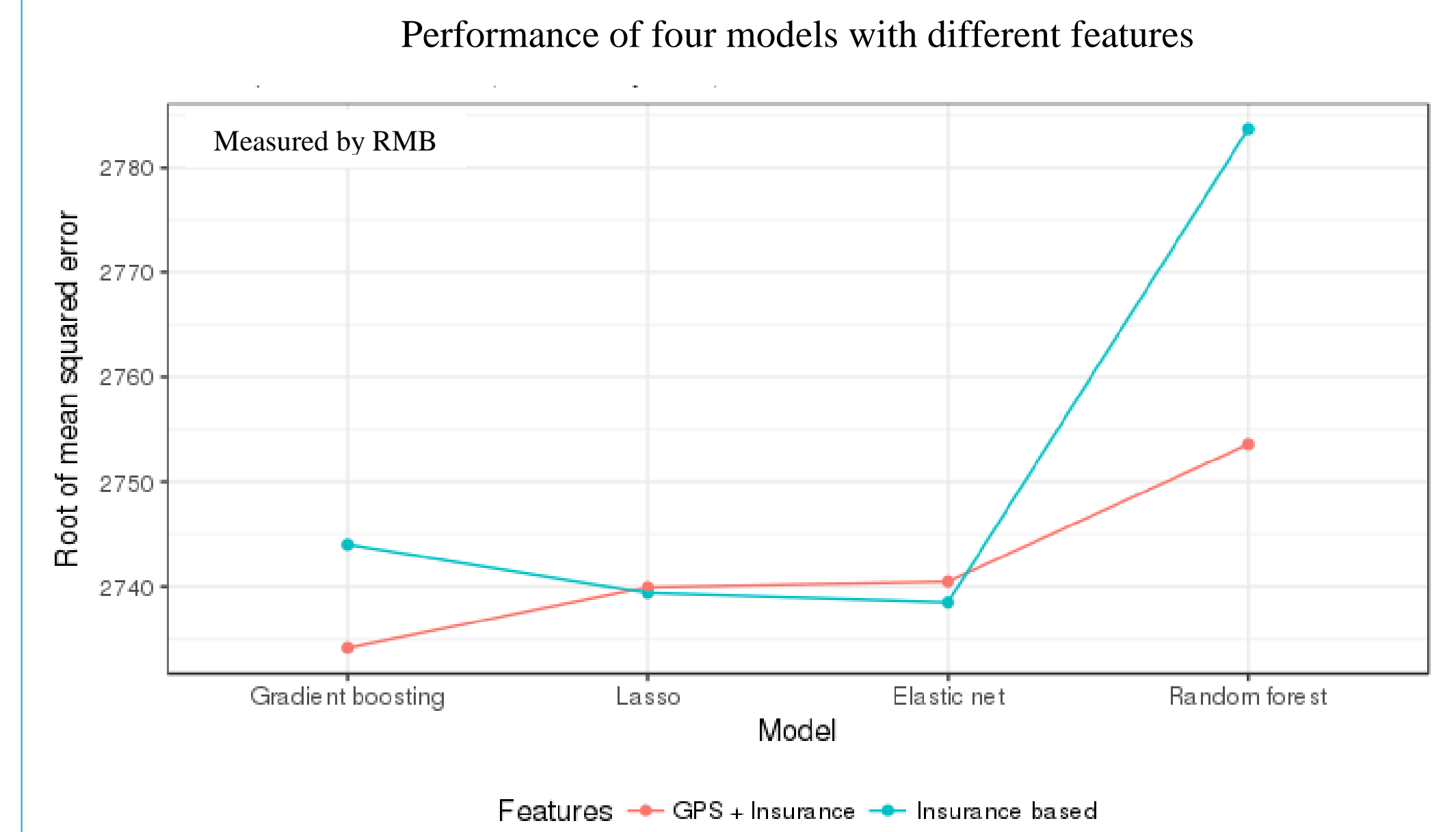
Pure driving behavior	Social driving behavior
Acceleration and brake behavior, driving speed and familiar road driving	Driving record based on time and place such as night driving and urban driving

Accident classification in binary logistic models

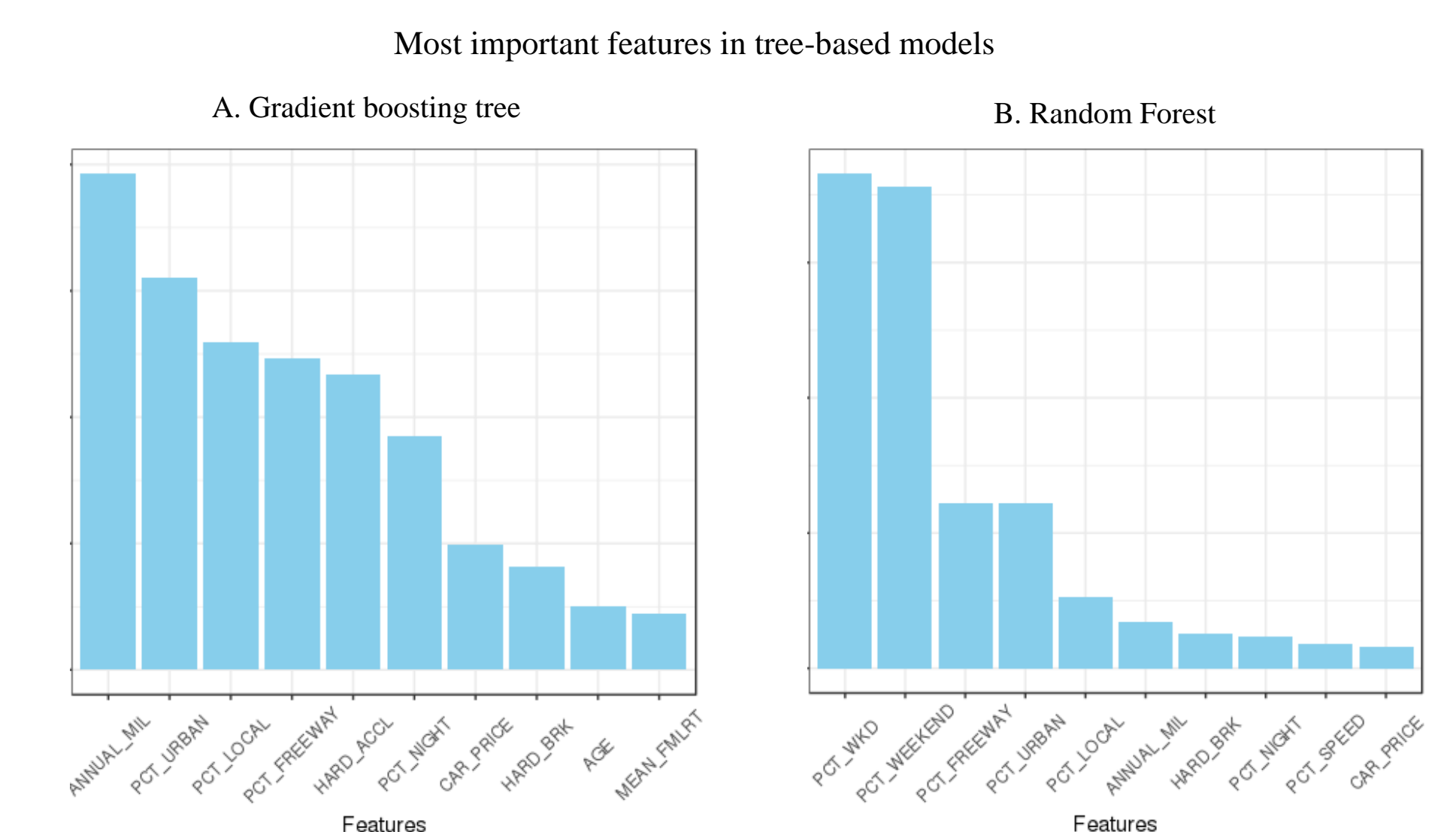
Variables	Definition	Models		
		Insurance	GPS	Stepwise
Driver-related characteristics in the insurance dataset				
FEMALE	Drivers are female	0.008 (0.004)	0.030 (0.005)	
AGE	Drivers' age	-0.001 (0.004)	0.001 (0.107)	
STATE_JOB	Working for the state	-0.016 (0.090)	-0.008 (0.091)	
INTERNET_SALE	Buying insurance via internet	-0.323 *** (0.125)	-0.339 *** (0.126)	-0.324 *** (0.120)
ANNUAL_MIL	Annual driving mileage	0.160 *** (0.025)	0.130 *** (0.028)	0.133 *** (0.028)
Vehicle-related characteristics in the insurance dataset				
CAR_PRICE	Car price	1.578e-07 (4.3e-07)	-1.629e-09 (4.37e-07)	
NEW_CAR	Car was purchased in recent three years	0.503 *** (0.098)	0.498 *** (0.099)	0.501 *** (0.098)
BIG_CAR	Car belongs to Minivan/SUV	-0.190 (0.188)	-0.194 (0.189)	
AIRBAG	Number of airbags	0.003 (0.028)	0.005 (0.028)	
ALARM	Equipped with a safe belt alarm	-0.0143 (0.349)	-0.132 (0.354)	
Home address extracted from GPS data				
HAIDIAN	Living at Haidian District	-0.409 (0.279)	-0.417 (0.275)	
SOUTH	Living at the south of the main urban area	0.278 *** (0.107)	0.262 ** (0.103)	
NORTH	Living in the northern resident area	0.255 * (0.139)	0.269 ** (0.133)	
Driving behaviors and patterns extracted from GPS data				
HARD_ACCL	Average number of hard accelerations in one hour	-0.025 (0.032)		
HARD_BRK	Average number of hard brakes in one hour	0.187 *** (0.054)	0.159 *** (0.042)	
PCT_SPEED	The fraction of mileage of driving with speed below 90 km/h	-0.028 (0.687)		
PCT_URBAN	The fraction of mileage of driving in the urban area	0.205 (0.256)		
PCT_FREEWAY	The fraction of mileage of driving on the freeways	-0.148 (0.535)		
PCT_LOCAL	The fraction of mileage of driving on the local roads	0.360 (0.645)		
PCT_WEEKEND	The fraction of mileage of driving during weekends	-0.334 (0.382)		
PCT_NIGHT	The fraction of mileage exposure of driving during nights	0.611 * (0.361)	0.621 * (0.349)	
MEAN_FMLRT	The frequency in the same roads	-0.123 *** (0.036)	-0.121 *** (0.035)	
Intercept		-2.403 *** (0.392)	-2.225 *** (0.921)	-2.390 *** (0.152)
Log-likelihood		-2039.3	-2015.5	-2017.7
AIC		4100.6	4077	4055
BIC		4175	4232	4123

standard errors are in parentheses. ***p≤0.01, **p≤0.05, *p≤0.1

Accident Loss prediction



Feature Influence / Importance



The gradient boosting model with features from both GPS and insurance data heavily relies on driving behavior and mobility pattern.

Conclusion and Discussion

- It is possible to extract useful socioeconomic information from GPS data for predictive purpose. However, the difficulty in this step is not only to design an appropriate method/algorithm, but to interpret the GPS points and assign the social meaning to it.
- I identify car owners' home address and speculate some areas that might provide information for prediction. The logistic regression results show that this method is useful since the variables are significant or quasi-significant. .
- The combination of socioeconomic and driving behavior features extracted from GPS data can improve the prediction of both accident and its loss, since they increase the log likelihood in the logistic model and lower the mean squared error via gradient boosting model.
- Driving behaviors sometimes seem more important than socioeconomic characteristics in accident prediction.