

A Data-Driven Approach to Characterize the Impact of Connected and Autonomous Vehicles on Traffic Flow

Amir Bahador Parsa ^{a,1}, Ramin Shabanpour ^a, Abolfazl (Kouros) Mohammadian ^a, Joshua Auld ^b,
Thomas Stephens ^b

^a*University of Illinois at Chicago, 842 W Taylor Street, Chicago, IL 60607, USA*

^b*Argonne National Laboratory, 9700 S Cass Ave Lemont, IL 60439, USA*

Abstract

The current study aims to present a model to characterize changes in network traffic flows as a result of implementing connected and autonomous vehicle (CAV) technology based on traffic network and built-environment characteristics. To develop such a model, first, POLARIS agent-based modeling platform is used to predict changes in average daily traffic (ADT) under CAVs scenario in the road network of Chicago metropolitan area as the dependent variable of the model. Second, a comprehensive set of variables and indicators representing network characteristics and urban structure patterns are generated. Three machine learning models namely K-Nearest neighbors, Random Forest, and eXtreme Gradient Boosting are developed and validated to establish the relationship between network characteristics and changes in ADT under CAVs scenario. The estimated models are found to yield acceptable performance. In addition, SHapley Additive exPlanations (SHAP) analysis tool is employed to investigate the impact of important features on changes in ADT, which discloses the most important link properties, network features, and demographic information in predicting change in ADT under the analyzed CAVs scenario.

Keywords: Connected and Autonomous Vehicles, POLARIS, Traffic Flow, Machine Learning

1. Introduction

Emergence of the connected and autonomous vehicles (CAVs) is a controversial topic in transportation community since they are expected to revolutionize both human mobility and goods transport in near future. In the United States, it is predicted that penetration of CAVs in light-duty-vehicle fleet will be up to 24.8% by the year 2045, if the technology price annually drops by 5%, and Americans' willingness to pay (WTP) increases by 5% in each year (1). Accordingly, with 10% annual price reductions and WTP increases, penetration of CAVs can reach up to 87.2% (1). It is also reported in another study that, if CAVs prices decrease at rates of 15% or 20% per year, it is expected that their market share will be homogeneous and near 100% by the year 2050 (2).

¹ Corresponding author. Tel: +1 (312) 996-9840.

Email addresses: aparsa2@uic.edu (A. B. Parsa), rshaba4@uic.edu (R. Shabanpour), kouros@uic.edu (A. K. Mohammadian), jauld@anl.gov (J. Auld), tstephens@anl.gov (T. Stephens).

38 Aligned with the rapid technological advancements in this area, most of major car companies have
39 announced that they will release their fully autonomous vehicles in the next decade (3, 4).

40 Given that CAVs have the potential to substitute the current vehicle fleet (5), a growing
41 number of studies have focused on evaluating the impact of this new technology on different
42 aspects of transportation systems. One aspect that is expected to be substantially affected by
43 emergence this technology is travel demand (6). It is predicted that an AV fleet size of only one-
44 third the number of private vehicles would be enough to meet the demand generated today (7).
45 Potential change in peoples' preferences towards their vehicle ownership, mode of travel, and
46 timing and sequence of travels are some examples of impacts of CAVs on travel demand. Travel
47 safety is another dimension which is expected to be greatly affected as vehicle to vehicle
48 communication systems can reduce the chance of collision of vehicles (8). There are several
49 studies evaluating the effects of CAVs on travel safety (e.g., 6–8). Most of these studies have
50 reported a considerable enhancement of safety as a result of CAVs deployment (11–13). A report
51 published by KPMG indicates that about 90% of all types of vehicle accidents can be eliminated
52 when CAVs substitute current vehicle fleet (14). Quite a few studies have also investigated the
53 impact of connected vehicles or CAVs on energy consumption and emission and found
54 inconclusive results with respect to environmental impacts of the technology (15–18).

55 Network traffic condition is another major dimension of transportation system which is
56 anticipated to be affected by CAVs technology (19–22). Analysis of flow-density diagram has
57 shown that increasing partial penetration of CAVs can result in more stable traffic stream (15).
58 Stability of traffic is found to be higher under CAVs scenario since automation and connection
59 between vehicles can prevent shockwave formation (23). Regarding congestion, researchers found
60 out that CAVs can benefit travel time through smoothing the traffic (15, 24, 25). Analysis of
61 capacity under CAVs scenarios indicates that CAVs penetration rate of 75% increases the capacity
62 by 25-35% (14). In another study, impact of CAVs on heterogenous traffic flow is simulated under
63 different penetration rates. It is reported that by increasing CAVs penetration rate up to 30%,
64 capacity increases at a slow pace, and passing that penetration rate will result in faster capacity
65 enhancements (26). A 50% penetration rate of CAVs can increase vehicle miles traveled (VMT)
66 by 20%. Increasing penetration rate to 95% can result in 35% increase in VMT (14).

67 Studies focusing on the impact of CAVs on transportation network are suffering from dearth
68 of CAVs historical data, especially at the large scale, which can certainly affect reliability and
69 accuracy of their results. Recently, a number of transportation simulation platforms have started
70 to incorporate vehicle automation and connectivity features into their simulation process. For
71 instance, several researchers have used VISSIM to simulate the impact of AVs or CAVs on
72 highway capacity (14), car following behavior (27), emergency evacuation (28), etc. Zhang and
73 Cassandras combined MATLAB and VISSIM to simulate the impact of CAVs on performance of
74 a single urban intersection (16). To cope with limitation of microscopic simulation models,
75 Talebpour and Mahmassani, proposed a novel acceleration framework as an alternative (23).
76 Amoozadeh et al. employed VENTOS simulation framework to analyze impact of CAVs on
77 different aspects of transportation system (29). In addition, other researchers proposed new
78 simulation frameworks such as microscopic simulation framework of Rios-Torres and
79 Malikopoulos to understand interaction of CAVs and human driven vehicles (HVs) at on-ramp
80 merging area (15), and a java-based algorithm by Yang et al., to predict total flow and demand
81 ratio of CAVs at intersections (21).

82 POLARIS, as an advanced transportation simulation framework, is another simulation
83 platform which is recently equipped with new modules to incorporate the simulation of CAVs

84 (30). The framework is developed by Argonne National Laboratory for Chicago and Detroit
85 regions. The POLARIS modeling suite is an open-source agent-based modeling platform
86 specifically designed for simulating large-scale transportation systems. The platform has been used
87 to successfully simulate ITS interactions with regional demand, statewide long-distance passenger
88 travel, and evacuations (31). Polaris is designed as a continuously integrated activity-based model
89 and network supply model, where individuals plan and schedule their activities dynamically,
90 engage in simulated travel, and re-plan activities on the fly due to changing traffic conditions, new
91 information or external control. The dynamic, integrated nature of POLARIS means that it is well
92 suited for simulating vehicle connectivity and automation and the impact on individual travel
93 behavior.

94 Using CAVs simulation results in POLARIS (32), the current study aims to present a data-
95 driven model to relate changes in network traffic flows as a result of implementing connected and
96 autonomous vehicle (CAV) technology to traffic network and built-environment. It is worthwhile
97 to note that the objective of this study is not assessing the impacts of CAVs, rather it aims to
98 provide a data-driven model to model traffic flow changes as a function of network characteristics
99 and built-environment factors. The proposed model can be applied in other geographical contexts
100 where a CAV-based network simulator is not available. To develop such a model, we used results
101 of simulations of CAVs in the Chicago metropolitan area, which were generated by POLARIS
102 agent-based platform (32), specifically taking the changes in average daily traffic (ADT) as the
103 dependent variable of the model. We have also integrated several other data sources along with
104 feature engineering through link-based analysis to train three powerful machine learning models,
105 K-Nearest Neighbors (KNN), Random Forest (RF), and eXtreme Gradient Boosting (XGBoost) to
106 find out and analyze significant features and their level of importance in predicting changes in
107 ADT of links under fully CAVs scenario. It is worth noting that although frequently-used machine
108 learning techniques have performed very well in transportation studies (33), more advanced
109 techniques such as deep learning and XGBoost along with a powerful analysis tool, SHAP, are
110 recently employed and resulted in more robust and great performance (34–36).

111 The remainder of this paper is organized as follows. First, different sources of data and
112 feature generation are described in detail. Second, machine learning techniques employed in this
113 study are explained in the methodology section. Then, in the results section, final models are
114 presented and performance of them are compared. Finally, conclusion and limitations of this study
115 are discussed.

116 **2. Data**

117 **2.1. POLARIS Simulation Output**

118 As previously pointed out, one of the major challenges in conducting research on CAVs
119 implications is the lack of historical data. Result of traffic simulation platform under CAVs
120 scenarios can be an acceptable alternative to the historical data. In this study, we use results of
121 simulations (29) of CAVs that were generated using the POLARIS platform that is developed by
122 Argonne National Laboratory for Chicago and Detroit regions is employed (30).

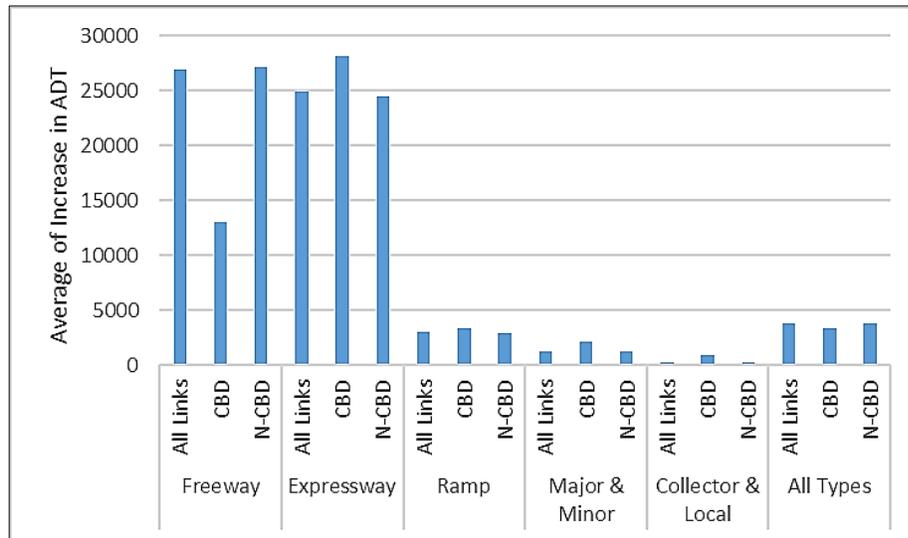
123 In (29) CAVs were represented in POLARIS by modifying several aspects of the simulation
124 to account for the expected impacts. For example, travelers who use an AV are assumed to have a
125 lower value of travel time, due to the reduced burden of driving, so the travel time utility
126 parameters for all choices (mode, timing, destination, etc.) within the demand models were
127 reduced for AV drivers based on literature review, with the value of time ranging from 100% down

128 to 50% of the current value. Traffic flow impacts were represented using empirically estimated
 129 link capacity changes with cooperative adaptive cruise control (CACC) penetration from (37). This
 130 function relates increases in link capacity to the penetration of CACC-equipped vehicles in a
 131 vehicle stream. For more information about the implementation of CAV scenario in the POLARIS
 132 framework, the reader is referred to (32).

133 In this study, we used POLARIS simulation results for two extreme scenarios: 0%
 134 penetration of CAVs (base scenario) and 100% penetration of CAVs (CAVs scenario) and
 135 calculated the difference of traffic flow for links of the Chicago network between these two
 136 scenarios. Traffic flow of links is simulated for a duration of 24 hours under the two scenarios. The
 137 daily traffic flows are referred to as ADT in this study, assuming that the POLARIS simulation
 138 results for a whole day period is a representative of the traffic condition during the year.
 139 Accordingly, the target variable is calculated through **Equation 1**.

$$141 \Delta ADT = ADT_{CAVs \text{ scenario}} - ADT_{base \text{ scenario}} \quad (1)$$

142
 143 Total number of 22,465 links from Chicago traffic network are considered in the POLARIS
 144 platform to generate traffic flow. **Figure 1** shows the average value of change in ADT across
 145 different five road types: Freeway (8.5% of roads), Expressway (1.1% of roads), Ramp (9.3% of
 146 roads), Major & Minor (71.9% of roads), Collector & Local (9.2% of roads). Based on this figure,
 147 average change in ADT is significantly higher for freeway and expressway. Furthermore,
 148 classifying links into central business district (CBD) and non-CBD groups displays a tangible
 149 difference between the average change in ADT in links located inside and outside the CBD.



150
 151 **Figure 1.** Average of increase in ADT under CAVs scenario

152 **Figure 2** displays the change in ADT across the study area. According to this figure, freeways
 153 and expressways connecting downtown to suburban areas are impacted more than the other road
 154 types by CAVs scenario. That means, by implementing CAVs technology, changes in ADT of these
 155 expressways and freeways can exceed 32000 vehicles per day.

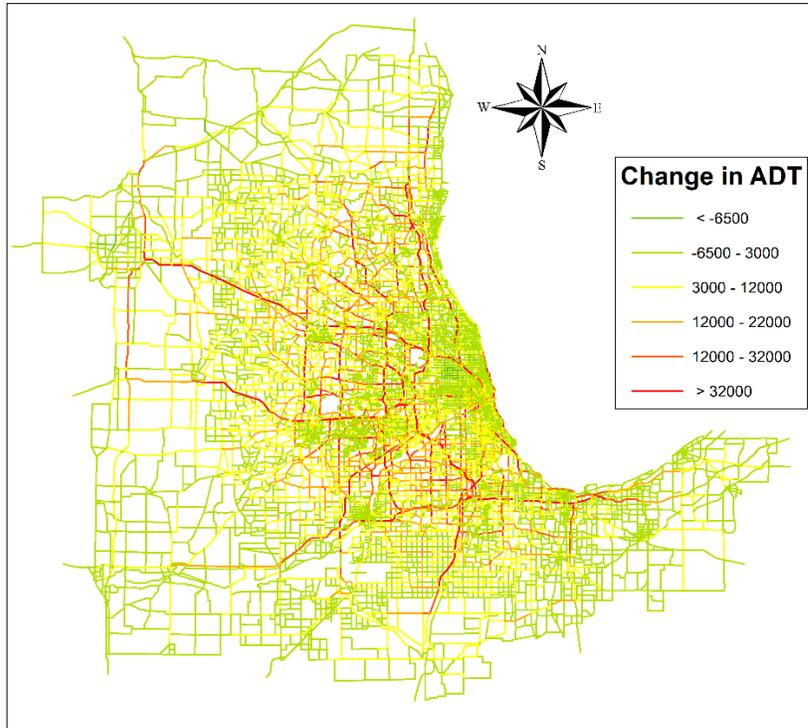
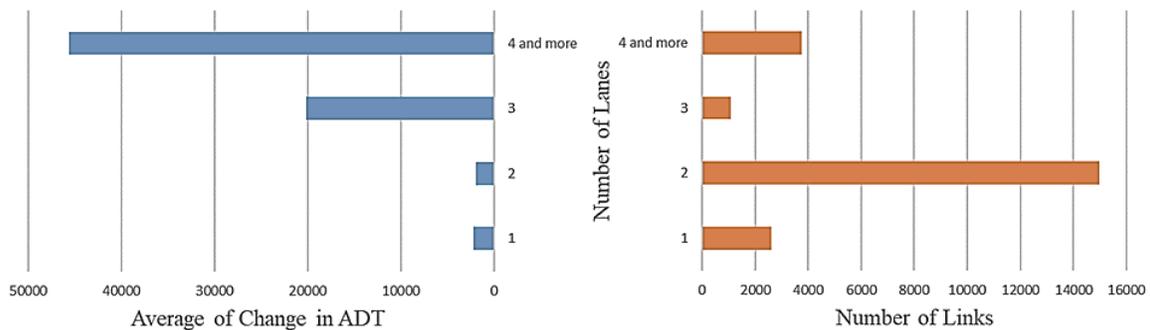


Figure 2. Change in ADT across the study area

156
157

2.2. Link Properties Data

158 Since our target variable, change in ADT, is a link-based parameter, all the features should be
 159 generated at the link level as well. Link properties include type of road, slope, length, and number
 160 of lanes. Intuitively, road type and other features such as connectivity are expected to have
 161 meaningful impacts on traffic flow. Finally, number of lanes is another important attribute of links
 162 which is used in training the models. **Figure 3** shows the proportion of links with different number
 163 of lanes and the average of target variable for links with different number of lanes. Based on this
 164 figure, most of links in the data have 2 lanes. Accordingly, increasing number of lanes can increase
 165 the change in ADT.
 166



167
168

Figure 3. Number of lanes and average change in ADT

169 **2.3. Network Data**

170 Features generated from traffic network of Chicago include connectivity index, distance to CBD,
171 road density, and intersection density. Connectivity index is a feature which is generated for each
172 link to represent the role of the link in creating linkage in the traffic network. It is expected that
173 increasing number of links connected to the start node and end node of a link increases connectivity
174 of that link. Therefore, connectivity index is defined through **Equation 2** for each link.

175
$$\text{Connectivity index} = (\# \text{ of links connected to the start node} + \# \text{ of links connected to the end node}) / \text{Link length} \quad (2)$$

176
177 Distance to CBD is another traffic network feature which is generated based on the distance
178 between each link and the CBD of Chicago. To do so, first centroid of Chicago CBD and all the
179 links are specified in the ArcMap and then the distance between the centroid of each link and
180 centroid of CBD is calculated and assigned to the links.

181 Finally, road density and intersection density are two traffic network features that are
182 generated using the Environmental Protection Agency's (EPA) Smart Location Database (38). EPA
183 is a comprehensive source of data which includes demographic, employment, and built
184 environment information for every census block group in the US. As mentioned earlier, all
185 variables in this study should be prepared at the link level. Therefore, attributes of each census
186 block group are assigned to the links which are passing through that block group. However, since
187 there are many links which pass through multiple block groups, weighted average of attributes of
188 those block groups is assigned to the link passes through them.

189 **2.4. Demographic Data**

190 Demographic data is another category that includes some features such as population, vehicle
191 ownership, job per household, job density within 45-minute drive, etc. Similar to road and
192 intersection density features, these features are generated from the EPA Smart Location Database.
193 Hence, they are assigned from block groups to the links, as was done with the density of roads and
194 intersections.

195 **2.5. Transportation Data**

196 Another type of data used in this study is transportation related data. Although there are several
197 transportation related features in the EPA Smart Location Database, trip equilibrium index is
198 selected to be used in model training. Trip equilibrium index is generated by calculating trip
199 productions and trip attractions of block groups in such a way that values closer to one indicate
200 that trip making at block group level is more balanced. Although we tested several transportation-
201 related variables in the EPA Smart Location Database, only trip equilibrium index is found to be
202 significant in the models.

203 **2.6. Land-Use Data**

204 Another source of data used in this study is land use data provided by Chicago Metropolitan
205 Agency for Planning (CMAP) which includes very detailed land-use information for the Chicago
206 metropolitan area. CMAP land-use types can be divided into eight groups of 1) Residential, 2)
207 Commercial, 3) Institutional, 4) Industrial, 5) Transportation, communication, utilities, and waste
208 land uses, 6) Agriculture, 7) Open space, and 8) vacant/under construction. In order to assign land
209 use variables to the links, a comprehensive GIS-based analysis has been conducted and different
210 sizes of buffer area are created and tested around links. Having analyzed different sizes of buffer

211 area, an area which covers 150 meters around a link is selected as the preferred size of buffer area
 212 for this study. Accordingly, for each link’s buffer area, percentage of area which is covered by each
 213 land use type is calculated and assigned to that link. Further, **Table 1** shows the final set of
 214 explanatory variables used in the next step to train the models.

215 **Table 1.** Description of explanatory variables

Variable	Description	Mean
Link Properties		
Freeway	1: if link type is freeway; 0: otherwise	0.08
Expressway	1: if link type is Expressway; 0: otherwise	0.01
lanes	Number of lanes	2.32
Network		
Dist_CBD	Distance from centroid of link to the centroid of CBD	32673.4
Connectivity	Role of link in making connection between links of network	0.027
Road_den	Total road network density	18.24
Intersect_den	Street intersection density	84.91
Demographic		
Job	Jobs within 45 minutes auto travel time	264676
HH_1veh	Number of households in block group that own 1 auto	260.87
HH_2veh_	Number of households in block group that own 2 or more auto	355.44
Pop_work_aged	Percent of population that is working aged	0.77
G_pop_den	Gross population density (people/acre) on unprotected land	14.53
Jobs_HH	Jobs per household	63.52
Entertain_job	Entertainment jobs within a 5-tier employment classification scheme	233.70
Pop	Population of block group	1673.87
Transportation		
Trip_equ_ind	Trip equilibrium index	0.41
Land-use		
Residential	Area of buffer zone around link covered by residential land use	0.024
Commercial	Area of buffer zone around link covered by commercial land use	0.017
Transport	Area of buffer zone around link covered by transportation related land use	0.011

216 3. Methodology

217 Three Machine Learning techniques, namely, K-Nearest Neighbors (KNN), Random Forest (RF)
 218 and eXtreme Gradient Boosting (XGBoost) are employed in this study due to their high estimation
 219 accuracy compared to the other ML and statistical models. A brief introduction to these models is
 220 provided in the following sub-sections.

221 3.1. K-Nearest Neighbors

222 One of the most popular supervised machine learning techniques, which is widely used for
 223 classification and regression, is the K-Nearest Neighbors technique. In this study, the KNN
 224 regression algorithm is used in which the output is a continuous value (change in ADT).
 225 Based on the training data points, which are described by multiple attributes, a feature space is
 226 formed, and each record is positioned in this space. Then, each unknown record (i.e., a data point
 227 from test data) is located in the feature space based on the value of its attributes and the KNN

228 technique looks for the k nearest neighbors for this record in the training data points. Thus, the
229 value of target variable for this record is predicted based on the arithmetic average of the value of
230 target variable of those data points which are selected as the k nearest neighbors.

231 To measure the distance in order to find the closest (i.e., most similar) data points, different
232 distance metrics could be used such as Euclidean distance which is one of the famous ones.
233 **Equation 3** represents the Euclidean distance between two points of $X_1 = (x_{11}, x_{12}, \dots, x_{1n})$ and
234 $X_2 = (x_{21}, x_{22}, \dots, x_{2n})$ with n attributes.

$$235 \text{Euclidean distance } (X_1, X_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2} \quad (3)$$

236
237
238 Finally, performance of the model is evaluated by comparing true value of test data points to
239 the values which are predicted for these test data points by the model.

240 **3.2. Random Forest**

241 Random Forest (RF) is a Machine Learning technique which utilizes combination of several
242 random Decision Trees (DTs). In DT technique, during the training process a feature selection
243 method is used in order to choose the best attribute to be used at each node of the tree; this heuristic
244 procedure also determines how to best split the node to two or more branches. Among different
245 functions such as Mean Squared Error (MSE), Friedman_MSE, and Mean Absolute Error (MAE)
246 to measure the quality of a split, the MSE technique is used in DT regressor model of this study
247 which is equal to variance reduction as the feature selection criterion. **Equation 4** presents the
248 MSE function:

$$249 \text{MSE} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad (4)$$

250
251
252 In this equation, μ is the average of x_i when i goes from 1 to n. Splitting on nodes is
253 accomplished through reduction of variance in such a way that the weighted variance of lower
254 level nodes should be less than the variance of upper level node.

255 RF is capable of working with categorical and numerical data. One disadvantage of DT is
256 that they are sensitive to the data on which they are trained. Hence, changing the training data can
257 significantly impact the resulting DT. To this end, aggregating several trees can result in higher
258 accuracy and decrease the probability of overfitting which might happen in an individual tree.

259 In the RF, a technique called Bootstrap Aggregation which is also known as bagging is used to
260 combine DTs. Bagging is a powerful method which is used to combine machine learning
261 techniques in order to achieve higher accuracy than the individual machine learning technique.
262 That is, different DTs are trained in parallel on different samples, selected randomly with
263 replacement from the data, and the aggregation of these trees would be the output prediction of the
264 RF model.

265 **3.3. eXtreme Gradient Boosting (XGBoost)**

266 Although RF usually performs well by combining a large number of DTs and taking average of
267 their outputs, DTs are generated independently in this technique. On the other hand, a more
268 advanced model called XGBoost which is created from gradient boosted decision trees can
269 improve the model performance through combining DTs in such a way that each new tree is
270 impacted by previously trained trees, and this can help to reduce errors. In this ensemble learning
271 technique, there are more parameters which need to be tuned to maximize model performance.
272 Proper parameter tuning is essential for XGBoost to avoid overfitting or being too complex. It is
273 also worth noting that RFs combine the results at the end of modeling procedure while XGBoost
274 does it along the process.

275 The parameters which should be tuned for XGBoost are as follows. First, number of
276 iterations which is the number of trees fitted in the model. Second, maximum depth of the tree
277 which is maximum number of splits and increasing this parameter can cause overfitting. Third,
278 subsample which is the fraction of observations randomly selected for the training instances and
279 can prevents overfitting. Forth is the learning rate used to shrink the weights and change the impact
280 of each individual tree at each step which results in a more robust model. Next parameter is
281 colsample_bytree which is subsampling the columns and can help prevent overfitting. The last two
282 parameters are lambda and alpha that are L2 and L1 regularization terms on weights, respectively,
283 and increasing their value makes the model more conservative. In this study parameters are tuned,
284 and their values are as follows. The optimal XGBoost hyper-parameters values after cross-
285 validation process are: Number of iterations: 700, Max Depth: 7, Subsample: 0.8, Colsample
286 bytree: 0.4, Lambda: 1.5, Alpha: 0.2, Learning Rate: 0.02.

287 **3.4. SHapley Additive exPlanations (SHAP)**

288 Interpreting output of machine learning techniques is often challenging. However, SHapley
289 Additive exPlanations (SHAP) is a powerful tool for this which was proposed by Lundberg and
290 Lee (39). SHAP is based on game theory rules (40) and local explanations (41), and it can provide
291 a means for estimating the contribution of each feature to the output of the model. Given an
292 XGBoost model with a set of N features is used to predict an output $v(N)$, SHAP values are
293 determined using several axioms to allocate the contribution of each feature through **Equation 5**.

$$294 \phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n-|S|-1)!}{n!} [v(S \cup \{i\}) - v(S)] \quad (5)$$

295
296 Where ϕ_i is contribution of feature i in the model output, and it is allocated based on their
297 marginal contribution (42). A linear function of binary features g is defined based on an additive
298 feature attribution method shown in **Equation 6** where M is the number of input features and
299 $z' \in \{0, 1\}^M$, equals to 1 when a feature is observed, otherwise it equals to 0 (39).

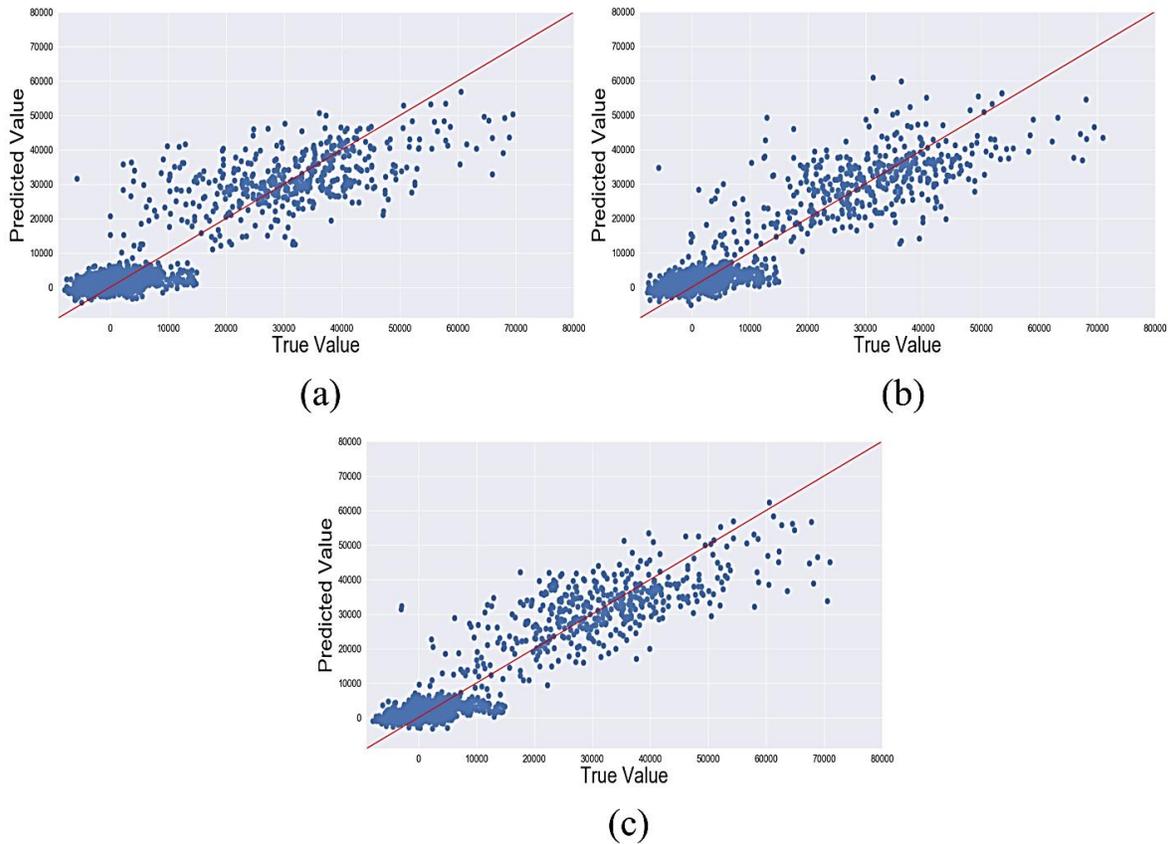
$$300 g(z') = \phi_i + \sum_{i=1}^M \phi_i z'_i \quad (6)$$

303 **4. Results**

304 To train the models, 70% of the data is randomly selected for training and the remaining 30% is
305 used to validate the models. In addition, a 5-fold cross-validation procedure is applied on the
306 training data. Therefore, at first the training data is divided to five subsamples randomly, and then

307 four subsamples are used to train the models while the remaining subsample is used as the
308 validation data. We repeated this procedure 5 times so that each subsample is used exactly once as
309 the validation data. This procedure helps us to measure whether a model is performing well
310 consistently.

311 Validation of three models shows that the KNN model, for which the optimal number of
312 neighbors is found to be six, results in the accuracy of 83.5%, the RF model achieves accuracy of
313 87.1%, and XGBoost yields the accuracy of 89.7%. Thus, XGBoost outperforms the other two
314 models in terms of accuracy. In the **Figure 4**, true values and predicted values of test dataset are
315 plotted for KNN, RF, and XGBoost techniques.
316



317
318 **Figure 4.** Predicted values against true values: (a) KNN, (b) RF, (c) XGBoost

319 After training the models, SHAP values of every feature are plotted in the **Figure 5** to show which
320 features are most important for the model as well as how these features can impact the XGBoost
321 model. In this figure, first 11 important features are sorted by the sum of SHAP value magnitudes,
322 then distribution of the impacts each feature has on the model output are displayed using SHAP
323 values. The color spectrum from blue to red represent the magnitude of feature values from low to
324 high, respectively.

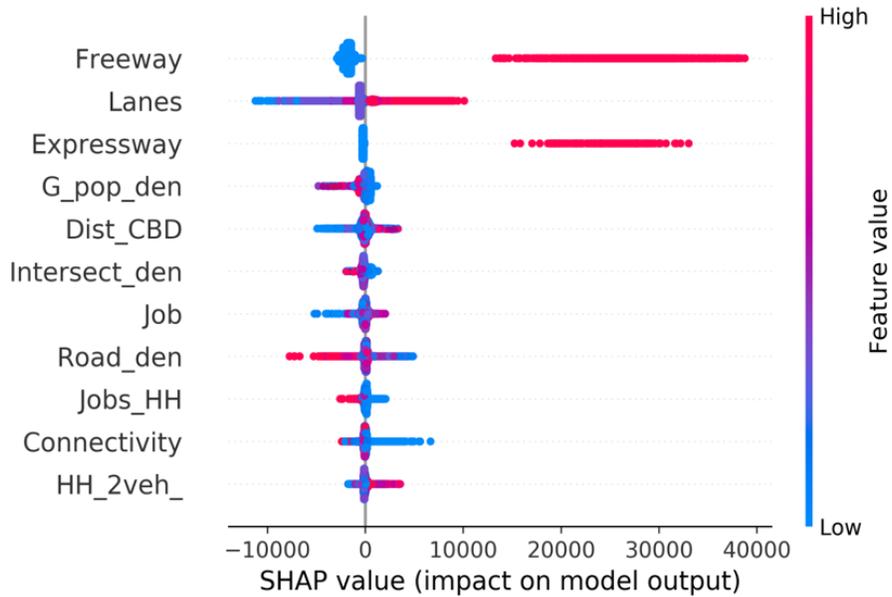


Figure 5. SHAP summary plot

325
326

327 Based on **Figure 5**, link properties including type of the roadway and number of lanes have
 328 the highest impact on the target variable (i.e., change in ADT) in such a way that increasing number
 329 of lanes and changing road type to Freeway and Expressway can increase likelihood of higher
 330 ADT in the CAVs scenario. Interestingly, next most important feature is gross population density
 331 so that for the roads passing through zones with denser gross population the change in ADT
 332 between base and CAVs scenario can decrease. The next important feature is distance to CBD, and
 333 it has a direct impact on the target variable meaning for the roads close to the CBD the change in
 334 ADT is less than that of roads far from the CBD. It could stem from that traffic of roadways which
 335 are close to the CBD are already higher than other roadways so that the impact of CAVs in
 336 increasing the ADT is less for these roads. **Figure 1** can also show the increase in ADT under
 337 CAVs scenario is slightly less for CBD roadways.

338 Intersection density of the zones through which roadways are passing is the next important
 339 features. However, based on **Figure 5**, when intersection density is lower, ADT would increase
 340 slightly more. Next important feature is number of jobs near the roadway and this feature has a
 341 direct impact on the target variable meaning that when there are more job opportunities around a
 342 road, the change in ADT would be higher. Road density is the next important feature which has a
 343 similar impact to intersection density. The next two features are jobs per household and
 344 connectivity, respectively, which have indirect impact on the target variable. That is, for lower
 345 values of these features impact of CAVs on change in ADT increases. Finally, according to **Figure**
 346 **5**, when number of households with two or more vehicles increases in a block group, ADT of
 347 roadways passing through this block group is expected to increase more under CAVs scenario.

348 It is worth noting that the reasons provided in this section are not definite and we tried to
 349 analyze features based on the observed data and our understanding about it. In addition, although
 350 some features might seem correlated, it doesn't impact the performance of the models, especially
 351 in tree-based models.

352 **5. Conclusion**

353 This study presented a data-driven model to relate changes in network traffic flows as a result of
354 implementing CAV technology to characteristics of the traffic network and built environment. To
355 develop such a model, we used changes in ADT under CAVs scenario in traffic network of Chicago
356 metropolitan area, which is generated by POLARIS agent-based platform. Using other sources of
357 data and feature engineering techniques, three machine learning models, KNN, RF and XGBoost,
358 are trained to predict impact of CAVs on traffic flow based on link-based features. Changes in
359 daily traffic flows of traffic network links is an indicator considered in this study and using data-
360 driven methods, it was modeled at the regional level and cross-validated in the same context. This
361 study demonstrates approaches that are useful for identifying the most important factors that
362 influence the changes in traffic flow attributable to widespread adoption of CAVs and for
363 quantifying the importance of each of these factors. We demonstrated these methods using results
364 of previous simulations of a CAVs scenario in the POLARIS (from (32)), and we took advantage
365 of different sources of data and powerful machine learning techniques to model the impacts of
366 CAVs on ADT.

367 It is found that traffic flows will most likely increase in most of the road types in case of fully
368 CAVs scenario. SHAP feature analysis also shows that properties of links have the highest impact
369 on target variable. Gross population density is the next important feature which has an indirect
370 impact on ADT. Next, distance of links from the CBD as well as other network features are the
371 second most important, and finally, attributes of block groups around the links such as
372 demographic, transportation and land uses are, respectively, less important, but still significant
373 features in predicting traffic flow in the CAV scenario analyzed.

374 Results of this study offer powerful methods that we validated for the Chicago metropolitan
375 area. Future work should test and hopefully validate these methods in analyzing simulations under
376 other conditions, such as different levels of CAVs penetration, or in other geographical contexts,
377 or perhaps for transferring simulation results from one geographical area to others. This is
378 important, since agent-based transportation demand models that model appropriate behaviors and
379 choices for metropolitan areas are difficult to develop and validate, and methods to analyze and
380 generalize results from existing models would be very valuable.

381 **Acknowledgments**

382 The authors gratefully acknowledge the sponsorship of the Systems and Modeling for Accelerated
383 Research in Transportation (SMART) Mobility Laboratory Consortium, an initiative of the Energy Efficient
384 Mobility Systems (EEMS) Program, managed by David Anderson of the Vehicle Technologies Office of
385 the U.S. Department of Energy. This study was conducted by the University of Illinois at Chicago
386 and Argonne National Laboratory, a U.S. Department of Energy laboratory managed by UChicago
387 Argonne, LLC under Contract No. DE-AC02-06CH11357. The authors are solely responsible for
388 the findings of this research which do not necessarily represent the views of the U.S. Department of
389 Energy or the United States Government.

390

391 **References**

392 1. Bansal, P., and K. M. Kockelman. Forecasting Americans ' Long-Term Adoption of Connected and
393 Autonomous Vehicle Technologies. *Transportation Research Part A*, Vol. 95, 2017, pp. 49–63.
394 <https://doi.org/10.1016/j.tra.2016.10.013>.

395 2. Talebian, A., and S. Mishra. Predicting the Adoption of Connected Autonomous Vehicles : A New
396 Approach Based on the Theory of Diffusion of Innovations. *Transportation Research Part C*, Vol.
397 95, No. June, 2018, pp. 363–380. <https://doi.org/10.1016/j.trc.2018.06.005>.

398 3. Bimbraw, K. Autonomous Cars: Past, Present and Future a Review of the Developments in the Last
399 Century, the Present Scenario and the Expected Future of Autonomous Vehicle Technology. *12th*
400 *International Conference on Informatics in Control, Automation and Robotics (ICINCO)*, Vol. 01,
401 2015, pp. 191–198.

402 4. Taiebat, M., A. L. Brown, H. R. Safford, S. Qu, and X. Ming. A Review on Energy, Environmental,
403 and Sustainability Implications of Connected and Automated Vehicles. *Environmental science &*
404 *technology*, 2018, pp. 11449–11465.

405 5. Sheng, S., E. Pakdamanian, K. Han, B. Kim, P. Tiwari, I. Kim, and L. Feng. A Case Study of Trust
406 on Autonomous Driving. *arXiv: 1904.11007 [cs. HC]*, 2019.

407 6. Lee, Y.-J., and A. Nickkar. Optimal Automated Demand Responsive Feeder Transit Operation and
408 Its Impact. 2018.

409 7. Spieser, K., K. Treleaven, R. Zhang, E. Frazzoli, D. Morton, and M. Pavone. Toward a Systematic
410 Approach to the Design and Evaluation of Automated Mobility-on-Demand Systems: A Case Study
411 in Singapore. *Road vehicle automation*, 2014, pp. 229–245.

412 8. Mudalige, U. P. Fast Collision Detection Technique for Connected Autonomous and Manual
413 Vehicles. *U.S. Patent*, Vol. 2, No. 12, 2013.

414 9. Arvin, R., M. Kamrani, A. J. Khattak, and J. Rios-Torres. Safety Impacts of Automated Vehicles in
415 Mixed Traffic. *97st Annual Meeting of Transportation Research Board*, 2018.

416 10. Arvin, R., M. Kamrani, and A. J. Khattak. How Instantaneous Driving Behavior Contributes to
417 Crashes at Intersections: Extracting Useful Information from Connected Vehicle Message Data.
418 *Accident Analysis & Prevention*, Vol. 127, 2019, pp. 118–133.

419 11. Li, T., and K. M. Kockelman. Valuing the Safety Benefits of Connected and Automated Vehicle
420 Technologies. *Transportation Research Board 95th Annual Meeting*, 2016, pp. 1–22.

421 12. Li, S. , Y. Zheng, K. Li, Y. Wu, J. . Hedrick, F. Gao, and H. Zhang. Dynamical Modeling and
422 Distributed Control of Connected and Automated Vehicles: Challenges and Opportunities. *IEEE*
423 *Intelligent Transportation Systems Magazine*, Vol. 9, 2017, pp. 46–58.
424 <https://doi.org/10.1109/MITS.2017.2709781>.

425 13. Talebpour, A., H. S. Mahmassani, and F. E. Bustamante. Modeling Driver Behavior in a Connected
426 Environment and Mobile Wireless Telecommunication Systems. *Transportation Research Record*,
427 2016. <https://doi.org/10.3141/2560-09>.

428 14. Bierstedt, J., A. Gooze, C. Gray, J. Peterman, L. Raykin, and J. Walters. Effects of Next-Generation
429 Vehicles on Travel Demand and Highway Capacity. *FP Think Working Group*, 2014.

430 15. Rios-torres, J., A. A. Malikopoulos, and S. Member. Impact of Partial Penetrations of Connected
431 and Automated Vehicles on Fuel Consumption and Traffic Flow. *IEEE TRANSACTIONS ON*
432 *INTELLIGENT VEHICLES*, Vol. 3, No. 4, 2018, pp. 453–462.

433 16. Zhang, Y., and C. G. Cassandras. The Penetration Effect of Connected Automated Vehicles in Urban
434 Traffic : An Energy Impact Study. *IEEE Conference on Control Technology and Applications*, 2018,
435 pp. 1–6.

436 17. Greenblatt, J. B., and S. Saxena. Autonomous Taxis Could Greatly Reduce Greenhouse-Gas
437 Emissions of US Light-Duty Vehicles. *Nature Climate Change*, Vol. 5, No. September, 2015.
438 <https://doi.org/10.1038/NCLIMATE2685>.

439 18. Ahangari, S., Z. Rashidi Moghaddam, M. Jeihani, C. Chavis, H. Chen, H. Rakha, and K. Kang.
440 Investigating the Effectiveness of an Eco-Speed Control System in the Vicinity of Signalized

441 Intersections Using a Driving Simulator. 2019.

442 19. Kidando, E., R. Moses, M. Ghorbanzadeh, and E. E. Ozguven. Traffic Operation and Safety Analysis
443 on an Arterial Highway : Implications for Connected Vehicle Applications. *21st International
444 Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 2753–2758.

445 20. Mahmassani, H. S. M. 50th Anniversary Invited Article — Autonomous Vehicles and Connected
446 Vehicle Systems : Flow and Operations Considerations. *Transportation Science*, No. July 2019,
447 2016.

448 21. Yang, K., S. I. Guler, and M. Menendez. Isolated Intersection Control for Various Levels of Vehicle
449 Technology : Conventional , Connected , and Automated Vehicles. *Transportation Research Part C*,
450 Vol. 72, 2016, pp. 109–129. <https://doi.org/10.1016/j.trc.2016.08.009>.

451 22. Azizi, L., M. S. Iqbal, and M. Hadi. Estimation of Freeway Platooning Measures Using Surrogate
452 Measures Based on Connected Vehicle Data. *97st Annual Meeting of Transportation Research
453 Board*, No. March, 2018.

454 23. Talebpour, A., and H. S. Mahmassani. Influence of Connected and Autonomous Vehicles on Traffic
455 Flow Stability and Throughput. *Transportation Research Part C*, Vol. 71, 2016, pp. 143–163.
456 <https://doi.org/10.1016/j.trc.2016.07.007>.

457 24. Approach, I. D. M., M. Zhou, X. Qu, S. Jin, A. M. Human, and D. Behavior. On the Impact of
458 Cooperative Autonomous Vehicles in Improving Freeway Merging : A Modified. *IEEE Transactions
459 on Intelligent Transportation Systems*, Vol. 18, No. 6, 2017, pp. 1422–1428.
460 <https://doi.org/10.1109/TITS.2016.2606492>.

461 25. Nezafat, R. V., E. Beheshtabar, M. Cetin, E. Williams, and G. F. List. Modeling and Evaluating
462 Traffic Flow at Sag Curves When Imposing Variable Speed Limits on Connected Vehicles.
463 *Transportation Research Record*, Vol. 2672, 2018, pp. 193–202.

464 26. Ye, L., and T. Yamamoto. Modeling Connected and Autonomous Vehicles in Heterogeneous Traffic
465 Flow. *Physica A*, Vol. 490, 2018, pp. 269–277. <https://doi.org/10.1016/j.physa.2017.08.015>.

466 27. Kockelman, K., P. Avery, P. Bansal, D. Stephen, P. Bujanovic, T. Choudhary, L. Clements, G.
467 Domnenko, D. Fagnant, J. Helsel, M. Levin, J. Li, T. Li, L. Loftus-, A. Nichols, M. Simoni, and R.
468 Hutchinson. Implications of Connected and Automated Vehicles on the Safety and Operations of
469 Roadway Networks : A Final Report. Vol. 7, 2016.

470 28. Al-Ahad, E., and R. Md Sharikur. Effects of Connected and Autonomous Vehicles on Contraflow
471 Operations for Emergency Evacuation: A Microsimulation Study. *Transportation Research Board
472 97th Annual Meeting*, 2018.

473 29. Amoozadeh, M., A. Raghuramu, C. Chuah, D. Ghosal, H. M. Zhang, and J. Rowe. Security
474 Vulnerabilities of Connected Vehicle Streams and Their Impact on Cooperative Driving. *IEEE
475 Communications Magazine*, Vol. 53, No. June, 2015, pp. 126–132.
476 <https://doi.org/10.1109/MCOM.2015.7120028>.

477 30. Auld, J., M. Hope, H. Ley, V. Sokolov, B. Xu, and K. Zhang. POLARIS: Agent-Based Modeling
478 Framework Development and Implementation for Integrated Travel Demand and Network and
479 Operations Simulations. *Transportation Research Part C: Emerging Technologies*, Vol. 64, 2016,
480 pp. 101–116. <https://doi.org/10.1016/j.trc.2015.07.017>.

481 31. Argonne National Laboratory. POLARIS Transportation System Simulation Tool.

482 32. Auld, J., O. Verbas, M. Javanmardi, and A. Rousseau. Impact of Privately-Owned Level 4 CAV
483 Technologies on Travel Demand and Energy. *Procedia Computer Science*, Vol. 130, 2018, pp. 914–
484 919. <https://doi.org/10.1016/J.PROCS.2018.04.089>.

485 33. Parsa, A. B., H. Taghipour, S. Derrible, and A. (Kouros) Mohammadian. Real-Time Accident
486 Detection : Coping with Imbalanced Data. *Accident Analysis and Prevention*, Vol. 129, No. January,
487 2019, pp. 202–210. <https://doi.org/10.1016/j.aap.2019.05.014>.

488 34. Parsa, A. B., A. Movahedi, H. Taghipour, S. Derrible, and A. K. Mohammadian. Toward Safer
489 Highways, Application of XGBoost and SHAP for Real-Time Accident Detection and Feature
490 Analysis. *Accident Analysis & Prevention*, 2020.

- 491 35. Parsa, A. B., R. S. Chauhan, H. Taghipour, S. Derrible, and A. Mohammadian. Applying Deep
492 Learning to Detect Traffic Accidents in Real Time Using Spatiotemporal Sequential Data. *arXiv*
493 *preprint*, Vol. arXiv:1912, 2019.
- 494 36. Movahedi, A., and S. Derrible. Interrelated Patterns of Electricity, Gas, and Water Consumption in
495 Large-Scale Buildings. (under Review). *engrXiv*, 2020.
496 <https://doi.org/https://doi.org/10.31224/osf.io/ahn3e>.
- 497 37. Shladover, S. ., C. Nowakowski, X. Lu, and R. Hoogendoorn. Using Cooperative Adaptive Cruise
498 Control (CACC) to Form High-Performance Vehicle Streams. Microscopic Traffic Modeling.
499 *eScholarship*. Available at: <http://escholarship.org/uc/item/3m89p611> .
- 500 38. Ramsey, K., and A. Bell. The Smart Location Database: A Nationwide Data Resource Characterizing
501 the Built Environment and Destination Accessibility at the Neighborhood Scale. *Cityscape*, 2014,
502 pp. 145–162.
- 503 39. Lundberg, S. M., and S.-I. Lee. A Unified Approach to Interpreting Model Predictions. *Advances in*
504 *Neural Information Processing Systems*, 2017, pp. 4765–4774.
- 505 40. Štrumbelj, E., and I. Kononenko. Explaining Prediction Models and Individual Predictions with
506 Feature Contributions. *Knowledge and information systems*, 2014, pp. 647–665.
- 507 41. Ribeiro, M. T., S. Singh, and C. Guestrin. Why Should i Trust You?: Explaining the Predictions of
508 Any Classifier. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge*
509 *discovery and data mining*, 2016.
- 510 42. Shapley, L. S. A Value for N-Person Games. *Contributions to the Theory of Games*, 1953, pp. 307–
511 317.
512