

1 **Investigation on the Driver-Victim Pairs in Pedestrian and Bicyclist Crashes by Latent Class**
2 **Clustering and Random Forest Algorithm**

3

4 **Chunwu Zhu***

5 Texas A&M Transportation Institute (TTI), Texas A&M University, Texas, USA
6 Department of Landscape Architecture and Urban Planning, School of Architecture, Texas A&M
7 University, Texas, USA
8 Email: chunwu.zhu@tamu.edu

9

10 **Charles T. Brown**

11 Equitable Cities LLC, New Jersey, USA
12 Email: charlesbrown@equitablecities.com

13 **Bahar Dadashova***

14 Texas A&M Transportation Institute (TTI), Texas A&M University, Texas, USA
15 Email: b-dadashova@tti.tamu.edu

16

17 **Xinyue Ye**

18 Department of Landscape Architecture and Urban Planning, School of Architecture, Texas A&M
19 University, Texas, USA
20 Email: xinyue.ye@tamu.edu

21

22 **Soheil Sohrabi**

23 Safe Transportation Research and Education Center (SafeTREC), University of California, Berkeley
24 Email: sohrabi@berkeley.edu

25

26 **Ingrid Potts**

27 Texas A&M Transportation Institute (TTI), Texas A&M University, Texas, USA
28 Email: i-potts@tti.tamu.edu

29

30

31

32

33 **Acknowledgment**

34

35 The work conducted in this paper was funded by the University Transportation Center Safety through
36 Disruption (Safe-D) Project ID 06-001. The authors would like to thank the three anonymous reviewers
37 whose feedback and comments helped to improve the quality of the manuscript significantly.

38

39 ***Corresponding Author:** *Bahar Dadashova, Ph.D. b-dadashova@tti.tamu.edu*

40

41 **ABSTRACT**

42
43 Pedestrians and bicyclists from marginalized and underserved populations experienced disproportionate
44 fatalities and injury rates due to traffic crashes in the US. This disparity among road users of different
45 races and the increasing trend of traffic risk for underserved racial groups called for an urgent agenda for
46 transportation policy making and research to ensure equity in roadway safety. Pedestrian and bicyclist
47 crashes involved drivers and pedestrians/bicyclists; the latter were usually victims. Traditional safety
48 studies did not account for the interaction between the two parties and assumed that they were
49 independent from each other. In this study we paired the driver and pedestrian/bicyclist involved in the
50 same crash to understand the socioeconomic and demographic make-up of the two parties involved in
51 crashes and assessed the geographic distribution of these crashes and crash-contributing factors. For this
52 purpose, we applied the latent class clustering analysis (LCA) to classify different crash types and analyze
53 the patterns of the crashes based on the income and ethnicity of both drivers and victims involved in
54 pedestrian and bicyclist crashes. We then used random forest algorithms and partial dependence plots
55 (PDPs) to model and interpreted the contributing factors of the clusters in both pedestrian and bicyclist
56 models. The clustering results showed a pattern of social segregation in pedestrian and bicyclist crashes
57 that drivers and victims with similar socioeconomic characteristics tend to be involved in one crash.
58 Pedestrian/bicyclist exposure, driver's age, victim's age, year of the car in use, annual average daily
59 traffic (AADT), speed limit, roadbed width, and lane width were the most influential factors contributing
60 to this pattern. Crashes that involved drivers and victims with lower income and non-white ethnicity
61 tended to happen in the location with higher pedestrian/bicyclist exposure, higher speed limit, and wider
62 road. The findings of this research can help to inform the decision-making process for improving safety to
63 ensure equitable and sustainable safety for all road users and communities.

64 **Keywords:** driver-victim pairs, pedestrian crashes, bicyclist crashes, latent class clustering, random forest

65

66 **1. INTRODUCTION**

67 Crash statistics in the US showed that vulnerable road users (VRUs) from marginalized and underserved
68 populations experienced disproportionate fatalities and injury rates due to traffic crashes. According to a
69 report from Governors Highway Safety Association (GHSA), Black, Indigenous, and People of Color
70 (BIPOC) experienced disproportionate traffic crash fatalities in the US from 2015-2019. The nationwide
71 total traffic deaths were 145.6 and 68.5 per 100,000 population for American Indian/Alaska Native and
72 Black, respectively, higher than 58.1 per 100,000 for total population (GHSA, 2021). For pedestrian
73 crashes, a report from National Highway Traffic Safety Administration (NHTSA) showed that the
74 pedestrian fatality rate for the white population is 1.5/100,000 in 2018, while the pedestrian fatality rate
75 for the Black population is 2.94/100,000, which is twice than the white pedestrian fatality rate
76 (Glassbrenner et al., 2022). Meanwhile, the motor vehicle traffic fatality for the Black population has
77 increased by 23 percent from 2019 to 2020, while for total population, it only increased by 7 percent
78 (NHTSA, 2021). These disparities in the distribution of traffic crashes among VRUs of different races
79 and the increasing trend of traffic risk for underserved racial groups suggested an urgent agenda for
80 transportation policy and research to ensure equity in roadway safety.

81 Traditional approaches to roadway safety, such as predictive and systemic tools safety analysis, usually
82 studied various road users and roadway infrastructure characteristics to predict the crash frequency and
83 severity and develop implementable solutions for preventing crashes. At the individual level, roadway
84 safety research investigated the influential factors such as the demographic and economic and behavioral
85 features of both parties involved in the crash (Hasheminejad et al., 2018; Balakrishnan et al., 2019;
86 Mokhtarimousavi et al., 2020), roadway environment such as speed limit, number of lanes, and traffic
87 control of the road segment where the crash happened (Sivasankaran & Balasubramanian, 2020; Xiao et
88 al., 2022), and other circumstances of the crash like weather condition and surface condition (Weiss et al.,
89 2014; Li et al., 2018). Human factors play an important role in a traffic crash for both parties. Typical
90 human factors like belligerent driving behavior and violations of traffic rules are deeply rooted in the road
91 users' socioeconomic backgrounds, which shape the different levels of vulnerability for road users with
92 different socioeconomic characteristics. Age, gender, income, and ethnicity were found to be major
93 demographic and socioeconomic features in the disparity of crash vulnerabilities (Boufous et al., 2011;
94 Zhao et al., 2013; Lombardi et al., 2017; Barajas, 2018; Billah et al., 2022). The difference in income and
95 ethnicity for both parties not only have a potential influence on the road user's driving, walking and
96 cycling behavior, but also result in an environmental difference in roadway infrastructure of a traffic crash
97 due to residential segregation of road users. For example, disadvantaged communities with more minority
98 populations and populations of lower socioeconomic status were found to have less access to bike lanes
99 across 22 large US cities (Braun et al., 2019). This disparity in crash risks among income and ethnic
100 groups was one of the major concerns for scholars and practitioners who want to ensure the principle of
101 environmental justice by mitigating the crash risk for low-income and minority groups through improving
102 the roadway infrastructure for them (Kravetz and Noland, 2012; Rebentisch et al., 2019).

103 VRU crashes usually involve two parties: drivers and VRUs like pedestrians or bicyclists. Drivers are
104 typically reported as the party at-fault in pedestrian/bicyclist-involved crashes, and pedestrians/bicyclists
105 are the victims. Previous research has investigated both parties' demographic and behavioral factors in
106 disaggregated analysis (Hasheminejad et al., 2018; Salon and McIntyre, 2018; Balakrishnan et al., 2019;).
107 Although these studies do include the characteristics of both drivers and VRUs in the analysis, they
108 usually treated the characteristics of drivers and victims as unrelated independent variables in their
109 theoretical assumptions and modeling process, which might overlook the potential interaction between
110 two parties. The close-to-home effect in roadway crashes suggested that the drivers and VRUs involved in
111 the same crash might live near each other and might share similar socioeconomic and demographic

112 characteristics (Burdett et al., 2017; Ulak et al., 2019). Under this assumption, the characteristics of
113 drivers and victims might be correlated, and understanding the occurrence of a crash should consider the
114 similarity of drivers and victims. This raised research questions about the socioeconomic patterns of
115 drivers and victims involved in one crash: To what extent the drivers and victims involved in one crash
116 share similar demographic and economic features? Are there potential crash patterns that can be found
117 based on their demographic and economic features? How are the different crash patterns distributed
118 geographically? And what factors shape the distribution of these patterns of crashes?

119 The remainder of this paper is organized as follows. In the next section we provide a literature review. In
120 section 3 we describe the data which is then followed by the methodological approach and modeling
121 techniques. In section 5 we describe the study findings and provide discussions. The paper ends with the
122 conclusions, references and Appendix.

123 **2. LITERATURE REVIEW**

124 **2.1 Vulnerable Road Users in Traffic Crashes**

125 Crashes and their consequences are not created equally for all road users. VRUs, such as pedestrians and
126 bicyclists, are more likely to be injured than drivers since they are less protected. There are also “implicit”
127 VRUs of certain demographic and economic groups who are usually found to have a higher chance of
128 getting involved in a crash or receiving more severe consequences. For example, children and the elderly
129 were considered more vulnerable than adult pedestrians and bicyclists (Braver, 2004; Ivan et al., 2019;
130 Ding et al., 2020). Behavioral and environmental differences were two major reasons contributing to the
131 vulnerability of implicit VRUs. Behavioral difference refers to the particular groups of VRUs who
132 showed riskier behavior when driving, walking, or biking. For example, younger drivers were more likely
133 to intentionally engage in risky driving behaviors such as mobile phone use (Scott-Parker and Oviedo-
134 Trespalacios, 2017; Oviedo-Trespalacios and Scott-Parker, 2018; Eren and Gauld. 2022). Environmental
135 difference refers to specific groups of VRUs who might live and travel in places with higher traffic
136 exposure and more unsafe roadway infrastructure. For example, Rothman et al. (2020) compared the road
137 infrastructure for low-income and high-income communities and found that fewer speed humps and lower
138 road classification might result in higher rates of child pedestrian crashes in low-income communities in
139 Toronto, Canada.

140 **2.2 Vulnerability of Pedestrians and Bicyclists**

141 Pedestrians/bicyclists are usually considered as VRUs in road safety literature, but certain groups of
142 pedestrians/bicyclists are more vulnerable according to their age (Boufous et al., 2011; Koopmans et al.;
143 2015, Boele-Vos et al.; 2017, Das et al., 2019), gender (Zhao et al., 2013; Toran Pour et al., 2018;
144 Algurén and Rizzi, 2022), income (Siddiqui et al. 2012; Barajas, 2018), ethnicity (Kravetz and Noland;
145 2012, Steinbach et al. 2016, Barajas, 2018), among others. Nearly one-third of pedestrian crashes and
146 two-thirds of bicyclist crashes involved school-aged children, according to police-reported crash data in
147 26 states in the US (Wheeler-Martin et al., 2020). Significant higher crash risks have also been found in
148 bicyclists younger than 30 years and older than 65 years of age when controlling for exposure in Spain
149 from 1993 to 2019 (Martínez-Ruiz et al., 2014). Though there was no solid evidence showing that male
150 pedestrians or bicyclists have higher crash risks than their counterparts, a few studies found male
151 pedestrians and bicyclists have less rule compliance and lower risk perception than females (Tom and
152 Granié, 2011; Prati et al., 2019). The behavioral differences among age and gender groups play a major
153 role in the disparity of roadway crashes, while the environmental differences better explained the
154 disparity among income and ethnic groups. Research has found that low-income and minority groups
155 were exposed to higher crash risk in pedestrian and bicyclist crashes in regions and cities of the United
156 States (Kravetz and Noland, 2012; Barajas, 2018). Scholars have also linked the disparity between low-

157 income and minority communities and high-income and majority communities with traffic exposure and
158 quality of roadway infrastructure and provided a potential explanation for this disparity from the
159 environmental difference (Fuller and Winters, 2017, Wang and Lindsey, 2017, Braun et al., 2019). For
160 example, Ferencak and Marshall (2021) investigated the installation of bicycling facilities across 29 US
161 cities and found a lower rate of bicycling facility installation in the block groups with more people of
162 color. Recent research in Oregon has found that lower median income and a higher proportion of the
163 BIPOC population are associated with more pedestrian crashes at the census tract level considering
164 factors from roadway infrastructure, land use, and socioeconomic background (Roll and McNeil, 2022).
165 The disproportionate share of low-income and minority groups in traffic crashes has called for equity and
166 environmental justice considerations in transportation planning and policy (Kravetz and Noland, 2012;
167 Rebentisch et al., 2019). Besides, the population with lower educational attainment and limited English
168 speaking has also been found to have higher crash risks at the aggregated level. (Barajas, 2018; Saha et
169 al., 2018).

170 **2.3 Vulnerability of Drivers**

171 Drivers' socioeconomic features, attitudes toward driving, and driving behavior are primary contributing
172 factors to the occurrence of roadway crashes (Adanu et al., 2017; Kemnitzer et al., 2019). Like
173 pedestrians and bicyclists, certain groups of drivers are more vulnerable to roadway crashes, primarily
174 due to differences in behavior and environmental factors. In the safety literature, these groups of drivers
175 were divided mainly by their socioeconomic and demographic characteristics in the literature, like age
176 (Lombardi et al., 2017; Gong and Fan, 2017; Liang and Yang, 2022) and gender (Russo et al., 2014;
177 Pulido et al., 2016; Billah et al., 2022). Regev et al. (2018) found that crash risk is highest for drivers
178 aged 21 to 29 in single-vehicle and multi-vehicle crashes from 2002 to 2012 in Great Britain when
179 controlling an exposure measurement considering the driver's trip number and population size. Billah et
180 al. (2022) found a more significant association between male drivers and the likelihood of crashes mainly
181 due to riskier driving behaviors of male drivers compared to their counterparts, such as speeding, driving
182 under the influence, and lane departure. Since drivers' income level and ethnicity were usually not
183 publicly available in police-reported crash data, research represents the economic status of drivers using
184 aggregated census data of drivers' residential ZIP code (Lee et al., 2021; Sagar et al., 2021). Though it
185 was not without bias, this surrogate measurement provides a feasible way to investigate the driver's
186 economic status in police-reported crashes. In the region where drivers' ethnic information was
187 unavailable, some researchers have also developed alternative approaches to estimate the drivers' race
188 and ethnicity. For example, Sartin et al. (2021) employed a Bayesian Improved Surname Geocoding
189 (BISG) method to estimate the population-level ethnic information for drivers in New Jersey.

190 **2.4 Linking Drivers and Victims in Crash Analysis**

191 Demographic and economic characteristics of drivers and victims should be considered in the crash
192 analysis since specific demographic and economic groups of drivers and victims are more vulnerable than
193 others. Existing literature considered demographic and economic features from both parties (Salon and
194 McIntyre, 2018; Balakrishnan et al., 2019). For example, Behnood and Mannering (2017) incorporated
195 both bicyclists' characteristics (gender, age, ethnicity, etc.) and drivers' characteristics (gender, age,
196 ethnicity, etc.) in their crash severity model of bicyclist crashes and found bicyclists' and drivers' race
197 and gender are the most important determinants of injury severity. However, these studies treated the
198 characteristics of drivers and victims as unrelated variables independent in their quantitative analysis.
199 This assumption might be problematic since potential spatial association might exist between drivers and
200 victims, which might lead to the similarity of social characteristics between drivers and victims. A series
201 of research investigating the proximity of crashes to the residential location of drivers/victims found a
202 close-to-home effect in crashes in which most of the crashes happened near the residence of both drivers
203 and victims (Burdett et al. 2017; Ulak et al., 2019). This close-to-home effect indicated that drivers and
204 victims involved in a crash might share the same neighborhood and similar socioeconomic and

205 demographic characteristics. Treating the socioeconomic characteristics of drivers and victims as
206 uncorrelated variables might ignore the spatial similarity of both parties and may lead to potential bias in
207 estimation. Thus, it is vital to consider the similarity of their characteristics in crash analysis.

208 Linking the characteristic of drivers and victims as driver-victim pairs and finding the hidden crash
209 patterns within driver-victim pairs can reveal the similarity between the drivers and victims in the same
210 crash. Clustering approaches have been usually employed to classify crashes by maximizing similarity
211 and minimizing dissimilarity among clusters to find the potential patterns in roadway crashes. Latent class
212 clustering analysis is one of the most popular approaches for revealing different crash patterns recently.
213 Sun et al. (2019) employed a latent class clustering method to classify pedestrian crashes in Louisiana and
214 found five clusters based on the factors from pedestrians' demographic features, crash-related factors, and
215 environmental factors. Samerei et al. (2021) also used latent class clustering analysis to classify bicyclist
216 crashes in Australia and found two clusters of crashes with different characteristics of bicyclists, road
217 environment, traffic control, and crash circumstance.

218 **3. DATA PREPARATION**

219 Data used in this study includes counts of pedestrian and bicyclist crashes, crash specific information,
220 socioeconomic characteristics of drivers and victims, roadway infrastructure characteristics, and traffic
221 exposure. Descriptive information of the variables is shown in Table 1.

222 **3.1 Crashes and Crash Specific Information**

223 This research aimed to investigate the spatial distribution and contributing factors for driver-victim pairs
224 in pedestrian/bicyclist crashes in Harries County, Texas, whose county seat is Houston. To collect the
225 crashes and related information, we obtained four-year (2017-2020) records of pedestrian and bicyclist
226 crashes from the Crash Records Information System (CRIS) of the Texas Department of Transportation
227 (TxDOT). We identified pedestrian/bicyclist crashes based on the type of primary victim (pedestrian or
228 bicyclist) involved in the crash. After removing redundant information and crash cases with missing
229 critical information, we kept only one driver and one primary victim for each crash event. As a result,
230 2,822 pedestrian crashes and 1,123 bicyclist crashes were identified with both the driver's and victim's
231 economic and demographic information available. There were 1,659 (58.8%) male and 1,163 (41.2%)
232 female victims in pedestrian crashes, with an average age of 39.3. For bicyclist crashes, there were 924
233 (82.3%) male and 199 (17.7%) female victims with an average age of 37.5. Eight factors in crash specific
234 information were retrieved from the CRIS database, including time of the day (*CR_TimeDay*), whether
235 the crash happened on a workday (*CR_Workday*), season (*CR_Season*), weather condition (*CR_Weather*),
236 surface condition (*CR_Surface*), whether crash happened on construction zone (*CR_Construct*), whether
237 the crash occurred at the intersection (*CR_Intersec*), and years of the car was in use (*CR_CarUsedYr*).

238 **3.2 Economic and Demographic Characteristics**

239 The drivers' and victims' economic and demographic characteristics were usually missing in a publicly
240 accessible crash database for privacy and liability concerns. From the CRIS database, we retrieved the
241 driver's ethnicity (*DR_Ethincity*), age (*DR_Age*), and gender (*DR_Gender*), and victim's ethnicity
242 (*VT_Ethincity*), age (*VT_Age*), and gender (*VT_Gender*). However, the CRIS database did not include the
243 income information of drivers and victims. Thus, we estimated the driver's and victim's income
244 information by the income level of their residential census tract based on median household income in
245 2019 American Community Survey (ACS) 5-year estimates. To obtain the driver's census tract, we
246 matched the ZIP code of drivers with their census tract in ArcGIS Pro. The corresponding census tract of
247 the driver was where the centroid of the ZIP code is situated. The victim's residential census tract was
248 hypothesized to be the same as the census tract where the crash happened. There were 786 census tracts in
249 Harris County with an average area of 5.9 km², which is within the range of acceptable walking and
250 biking distance (1,750-2,122 meters) (Rahul and Verma, 2014). Besides, most pedestrian and bicyclist
251 crashes happened near the victim's home (Steinbach et al., 2013; Ulak et al., 2019). Therefore, we

252 assumed the crash location was the same as the victim's residential census tract. Finally, we recoded the
253 driver's income (*DR_Income*) and victim's income (*VT_Income*) to ordinal variables in five levels: low
254 income, lower to medium income, medium income, medium to high income, and high income, according
255 to the five quintiles of their residential census tracts in the research area. In our dataset, there are 73
256 driver-victim pairs in which the driver and victim live in the same census tract, which accounts for the
257 3.14% of the total number of pedestrian crashes; for bicyclist crashes, there are 50 driver-victim pairs in
258 which the driver and victim live in the same census tract, which accounts for 5.18% of the total bicyclist
259 crashes. This small proportion of driver-victim pairs shows that assigning the driver-victim pair to the
260 same income level (despite the potential differences) does not introduce significant error to overall model
261 performance.

262 3.3 Roadway Infrastructure Characteristics

263 Characteristics of roadway infrastructure were collected from the Roadway Inventory of TxDOT, which
264 was a GIS-based road network database storing roadway information in Texas. Data for the roadway
265 inventory was updated annually, and we used the 2020 version, which conformed with the time span of
266 our crash events. We selected 11 characteristics of the roadway infrastructure where the crash happened,
267 including road functional classification (*RD_FuncCls*), speed limit (*RD_SpdLmt*), whether the crash
268 occurred in an urban area (*RD_Urban*), roadbed width which comprises shoulder width and surface
269 width (*RD_RdWth*), the number of lanes (*RD_LnNum*), lane width (*RD_LnWth*), median width
270 (*RD_MedWth*), inside shoulder width (*RD_SWthIn*), outside shoulder width (*RD_SWthOut*), existence of
271 left curb (*RD_CurbL*), and existence of right curb (*RD_CurbR*).

272 3.4 Exposure

273 Vehicular, pedestrian, and bicyclist exposure variables were also taken into consideration in this study.
274 Vehicular exposure of the road segments where the crash happened was measured as the Annual Average
275 Daily Traffic (AADT) from the TxDOT Roadway Inventory database for each year of the crash events.
276 However, the scarcity of pedestrian and bicyclist exposure data was one of the primary limitations in
277 crash modeling. Scholars have used emerging crowdsourced data to estimate pedestrian and bicyclist
278 exposure, such as bicycle count data from Strava (Dadashova and Griffin, 2020; Dadashova et al., 2020).
279 In this study, we used a scaling approach to estimate the bicyclist and pedestrian exposure leveraging
280 observed pedestrian and bicyclist count data available from Texas Bicycle and Pedestrian Data Exchange
281 (BP|CX) (<https://mobility.tamu.edu/bikepeddata/>) and crowd-sourced pedestrian and bicyclist count data
282 from Strava. The estimated pedestrian and bicyclist counts were averaged daily and calculated annually
283 using the scaling approach for Strava data by Dadashova et al. (2020). Using this approach, we calculated
284 the scaling factors (See Table S1) for each year by dividing the observed pedestrian/bicyclist counts with
285 the crowd-sourced pedestrian/bicyclist counts and calculated the estimated pedestrian/bicyclist counts by
286 multiplying the scaling factors with the crowd-sourced pedestrian/bicyclist counts on other road
287 segments.

288

Table 1. Descriptive Statistics of Variables

| Categorical Variables | Pedestrian crashes | Bicyclist crashes | |
|------------------------------|------------------------------|--------------------------|-------------|
| <i>CR_TimeDay</i> | 1 = 0:00-6:00 | 254(9.0%) | 51(4.5%) |
| | 2 = 6:00-12:00 | 749(26.5%) | 303(27.0%) |
| | 3 = 12:00-18:00 | 839(29.7%) | 449(40.0%) |
| | 4 = 18:00-24:00 | 980(34.7%) | 320(28.5%) |
| <i>CR_Workday</i> | 1 = Monday to Friday | 2199(77.9%) | 856(76.2%) |
| | 2 = Saturday to Sunday | 623(22.1%) | 267(23.8%) |
| <i>CR_Season</i> | 1 = Spring | 746(26.4%) | 243(21.6%) |
| | 2 = Summer | 679(24.1%) | 300(26.7%) |
| | 3 = Autumn | 606(21.5%) | 293(26.1%) |
| | 4 = Winter | 791(28.0%) | 287(25.6%) |
| <i>CR_Weather</i> | 1 = Clear | 2093(74.2%) | 873(77.7%) |
| | 2 = Others | 729(25.8%) | 250(22.3%) |
| <i>CR_Surface</i> | 1 = Dry | 2495(88.4%) | 1039(92.5%) |
| | 2 = Others | 327(11.6%) | 84(7.5%) |
| <i>CR_Construction</i> | 1 = At construction zone | 51(1.8%) | 5(0.4%) |
| | 2 = Not at construction zone | 2771(98.2%) | 1118(99.6%) |
| <i>CR_Intersec</i> | 1 = At intersection | 1034(36.6%) | 661(58.9%) |
| | 2 = Not at intersection | 1788(63.4%) | 462(41.1%) |
| <i>DR_Income</i> | 1 = low income | 454(16.1%) | 193(17.2%) |
| | 2 = low to medium income | 614(21.8%) | 229(20.4%) |
| | 3 = medium income | 814(28.8%) | 280(24.9%) |
| | 4 = medium to high income | 368(13.0%) | 161(14.3%) |
| | 5 = high income | 572(20.3%) | 260(23.2%) |
| <i>DR_Ethnicity</i> | 1 = White | 888(31.5%) | 384(34.2%) |
| | 2 = Hispanic | 896(31.8%) | 350(31.2%) |
| | 3 = Black | 800(28.3%) | 297(26.4%) |
| | 4 = Asian | 185(6.6%) | 70(6.2%) |
| | 5 = Others | 47(1.7%) | 20(1.8%) |
| <i>DR_Gender</i> | 1 = Male | 1632(57.8%) | 626(55.7%) |
| | 2 = Female | 1190(42.2%) | 497(44.3%) |
| <i>VT_Income</i> | 1 = low income | 600(21.3%) | 211(18.8%) |
| | 2 = low to medium income | 762(27.0%) | 286(25.5%) |
| | 3 = medium income | 559(19.8%) | 210(18.7%) |
| | 4 = medium to high income | 389(13.8%) | 156(13.9%) |
| | 5 = high income | 512(18.1%) | 260(23.2%) |
| <i>VT_Ethnicity</i> | 1 = White | 935(33.1%) | 485(43.2%) |
| | 2 = Hispanic | 852(30.2%) | 274(24.4%) |

| Categorical Variables | Pedestrian crashes | | | | Bicyclist crashes | | | | |
|-----------------------------------|---------------------------|-------------|---------|--------|--------------------------|-------------|-------|--------|--|
| | 3 = Black | 843(29.9%) | | | | 299(26.6%) | | | |
| | 4 = Asian | 131(4.6%) | | | | 52(4.6%) | | | |
| | 5 = Others | 54(1.9%) | | | | 11(1.0%) | | | |
| <i>VT_Gender</i> | 1 = Male | 1659(58.8%) | | | | 924(82.3%) | | | |
| | 2 = Female | 1163(41.2%) | | | | 199(17.7%) | | | |
| <i>RD_FuncCls</i> | 1 = Collectors | 782(27.7%) | | | | 268(23.9%) | | | |
| | 2 = Local roads | 2044(72.4%) | | | | 855(76.1%) | | | |
| <i>RD_Urban</i> | 1 = Urban area | 2818(99.9%) | | | | 1117(99.5%) | | | |
| | 2 = Rural area | 4(0.1%) | | | | 6(0.5%) | | | |
| <i>RD_CurbL</i> | 1 = Left curb exists | 2689(95.3%) | | | | 1092(97.2%) | | | |
| | 2 = No left curb | 133(4.7%) | | | | 31(2.8%) | | | |
| <i>RD_CurbR</i> | 1 = Right curb exists | 2690(95.3%) | | | | 1093(97.3%) | | | |
| | 2 = No right curb | 132(4.7%) | | | | 30(2.7%) | | | |
| <i>RD_LnWth (feet)</i> | 1 = Less than 10 | 1527(54.1%) | | | | 187(16.7%) | | | |
| | 2 = 10 to 12 | 884(31.3%) | | | | 534(47.6%) | | | |
| | 3 = 12 to 14 | 139(4.9%) | | | | 313(27.0%) | | | |
| | 4 = Greater than 14 | 272(9.6%) | | | | 89(7.9%) | | | |
| Continuous Variables | Mean | Min | Max | SD | Mean | Min | Max | SD | |
| <i>CR_CarUsedYr</i> | 8.3 | 0.0 | 43.0 | 5.8 | 8.2 | 0 | 43 | 6.1 | |
| <i>DR_Age</i> | 41.6 | 15.0 | 118.0 | 16.6 | 43.4 | 8 | 118 | 17.3 | |
| <i>VT_Age</i> | 39.3 | 1.0 | 100.0 | 19.5 | 37.5 | 3 | 100 | 19 | |
| <i>RD_SpdLmt (miles per hour)</i> | 36.7 | 20.0 | 65.0 | 11.2 | 37.4 | 20 | 65 | 12 | |
| <i>RD_RdWth (feet)</i> | 33.3 | 14.0 | 106.0 | 14.1 | 30.7 | 16 | 106 | 12.8 | |
| <i>RD_LnNum</i> | 2.9 | 1.0 | 6.0 | 1.1 | 2.7 | 2 | 6 | 1 | |
| <i>RD_LnWth</i> | 11.4 | 5.0 | 27.0 | 2.8 | 11 | 5 | 27 | 2.3 | |
| <i>RD_MedWth (feet)</i> | 0.3 | 0.0 | 138.0 | 3.6 | 1.1 | 0 | 138 | 9.4 | |
| <i>RD_SWthIn (feet)</i> | 0.1 | 0.0 | 10.0 | 0.6 | 0.1 | 0 | 10 | 0.6 | |
| <i>RD_SWthOut (feet)</i> | 0.1 | 0.0 | 10.0 | 0.8 | 0.1 | 0 | 10 | 1 | |
| <i>EX_Ped</i> | 36.3 | 0.3 | 340.7 | 104.0 | Not Applicable | | | | |
| <i>EX_Cyc</i> | Not Applicable | | | | 14.1 | 0.1 | 340.7 | 30.8 | |
| <i>AADT</i> | 9864.7 | 50.0 | 49968.0 | 9954.5 | 8601 | 69 | 49968 | 9565.4 | |

290 4. METHODOLOGY

291 In this study, we applied a latent class clustering analysis (LCA) to identify the patterns in driver-victim
292 pairs according to the driver's and victim's income and ethnicity in pedestrian and bicyclist crashes. We
293 also mapped the crash patterns in the study area to reveal their spatial distribution. Then, we used random
294 forest algorithm to investigate the relative contribution of factors to the crash patterns from crash specific
295 information, economic and demographic characteristics of drivers and victims, roadway infrastructure,
296 and exposure. Finally, we drew partial dependence plots (PDPs) for the most important factors to interpret
297 their influences on certain crash patterns.

298 4.1 Latent Class Clustering Analysis

299 To investigate the possible patterns in driver-victim pairs, we applied LCA to divide the pedestrian and
300 bicyclist crashes according to the victim's and driver's socioeconomic characteristics. Clustering analysis
301 is an unsupervised machine learning method that can separate the crashes into homogenous subgroups,
302 which have the largest similarities within each subgroup and the largest dissimilarity between each
303 subgroup (Sivasankaran and Balasubramanian, 2020). We used a probability-based clustering approach
304 (i.e., Latent Class Clustering), which has recently been applied in several roadway safety research studies
305 (Sun et al., 2019, Samerei et al., 2021). The Latent Class Clustering approach has several advantages over
306 other clustering approaches (e.g., K-means) in that it 1) can calculate the probability of a crash of being in
307 a certain cluster by maximum likelihood method; 2) does not necessarily need to standardize the variables
308 beforehand; 3) does not need to specify the number of clusters before performing the clustering; 4) can
309 generate statistical criteria afterward to select the best model with a certain number of clusters
310 (Sasidharan et al., 2015, Sun et al., 2019). The mathematical formula of the LCA approach is shown
311 below (Samerei et al., 2021):

$$312 \quad P(Y_i = y) = \sum_{k=1}^{K_c} \rho \prod_{m=1}^M \prod_{n=1}^{r_m} \theta_{mn|l}^{l(y_m=n)}$$

313 Where $Y_i = (Y_{i1}, \dots, Y_{iM})$ is the observation (crash) i 's responses in M category, where the possible
314 values of Y_{iM} are $1, \dots, r_m$; r_m represents the crash i 's r th attribute in m category; K_c represents the
315 number of latent classes to be estimated; $l(y_m = n)$ is the indicator function to be 1 if y equals n and to
316 be 0 when y is not 1; ρ is the probability of latent class membership probability and θ is the conditional
317 probabilities of responses on latent class membership. The number of clusters can influence the goodness-
318 of-fit of the latent class clustering model. We employ Bayesian Information Criteria (BIC) to select the
319 appropriate number of clusters. LCA modeling and BIC calculation were conducted by package *polPCA*
320 in R.

321 4.2 Random Forest Algorithm and Partial Dependence Plot

322 Random forest algorithm is known as a tree-based ensemble machine learning technique. It is built upon a
323 multitude of weak decision tree models to form a strong "forest" by averaging the predictions from all the
324 individual regression trees or by taking the majority vote from the classification tree. It can be applied in
325 both classification and regression, and we use the random forest algorithm for classification in this task.
326 The random forest algorithm employs a bagging technique to repeatedly select a random sample from the
327 training dataset and use the sample to fit a decision tree. Let feature set X be $\{x_1, x_2, \dots, x_n\}$, target set Y
328 be $\{y_1, y_2, \dots, y_n\}$, and $i = 1, 2, \dots, I$, the process of random forest can be represented as:

- 329 1) Select a random sample set from $\{X, Y\}$, which is denoted as $\{x_i, y_i\}$;
- 330 2) Train a decision tree f_i on the sample set $\{x_i, y_i\}$;
- 331 3) Repeat procedures 1 and 2 for I times to get I decision trees $\{f_1, f_2, \dots, f_I\}$;
- 332 4) Aggregating the prediction results for any random sample \hat{x} to get function \hat{f} for the random
333 forest. For classification, it takes the majority vote of the target from all individual decision trees,
334 denoted as $\hat{f}(\hat{x}) = \max_{i=1,2,\dots,I} f_i(x_i)$.

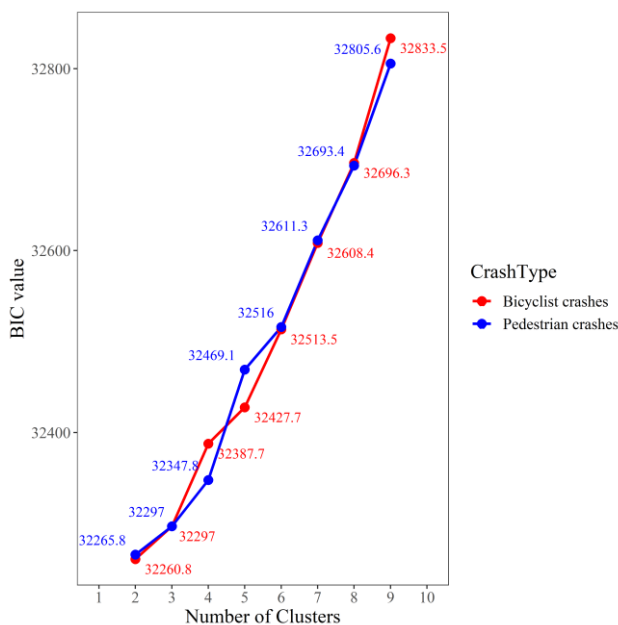
335 Several parameters can affect the performance of the model: for example, the number of decision trees
336 (I). To optimize performance, we employed a random search method for optimal parameters with
337 successive halving to automatically find the best combination of parameters (Scikit-Learn, 2022). To
338 investigate the impact of variables in clusters of driver-victim pairs, we calculated the feature importance
339 for each variable to assess the relative contribution of all the variables (Masís, 2021). Furthermore, we
340 used the PDPs, which is one of the model-agnostic interpretable machine learning approaches to reveal
341 the marginal effect of a feature in machine learning models (Masís, 2021). Random forest algorithm was
342 implemented by *Scikit-learn*, and PDPs are generated by *pdpbox* in Python.

343 **5. RESULTS**

344 **5.1 Results of Latent Class Clustering Analysis**

345 *Identifying the Number of Clusters*

346 The LCA models were performed on the economic and demographic characteristics of the driver-victim
347 pairs to find the patterns within the driver-victim pairs regarding their income level and ethnicity. Figure
348 1 shows the graphs of BIC value with a different number of clusters for pedestrian crashes and bicyclist
349 crashes. As indicated by the graph, the BIC values in the LCA of pedestrian crashes increase along with
350 the increasing number of clusters, and the minimum BIC value is generated when the number of clusters
351 is set as two. The trend of BIC values in the LCA of bicyclist crashes is similar to pedestrian crashes,
352 achieving its lowest value when there are two clusters. Thus, we report the results of the LCA for
353 pedestrian crashes and bicyclist crashes when there are two clusters in each model.

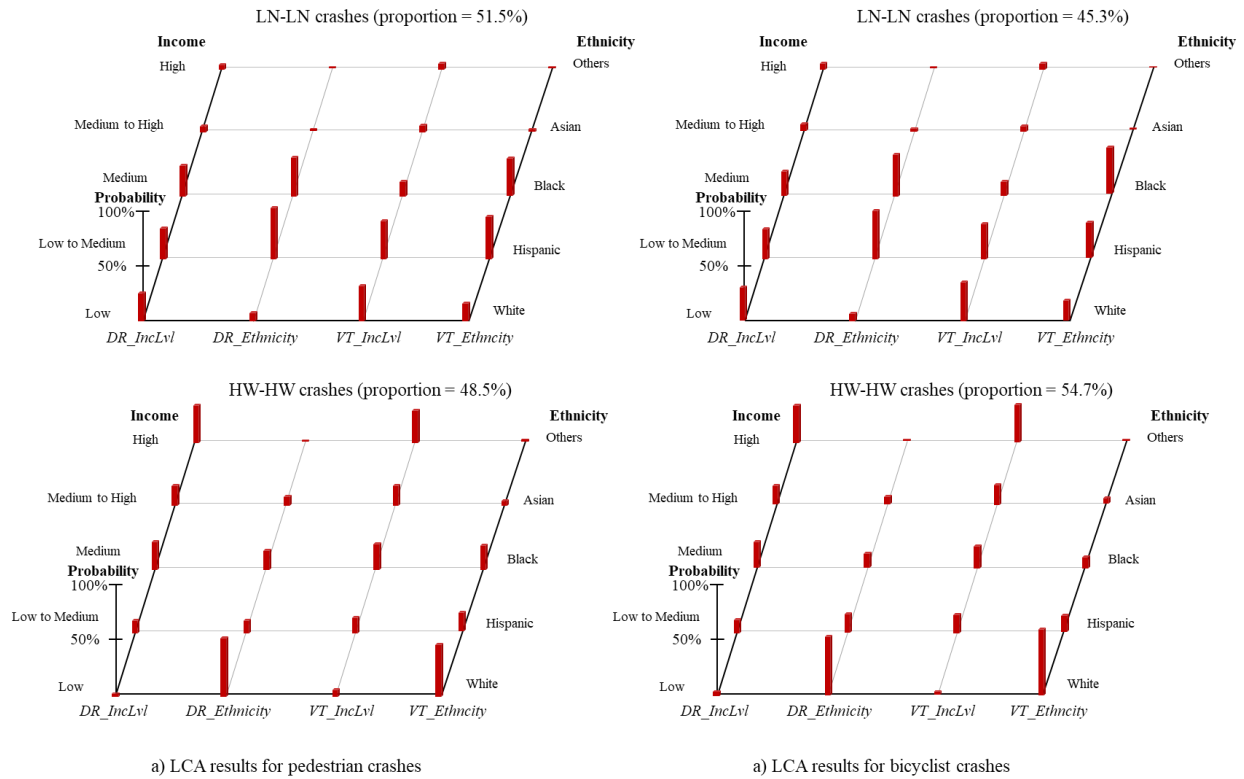


354
355 **Figure 1. Number of Clusters and Their Respective BIC Value in Pedestrian and Bicyclist Crashes**

356 *Clustering Crashes by LCA*

357 Figure 2 shows the distribution of the driver’s and victim’s income level and ethnicity in each cluster in
358 both models. Detailed information about the clustering results can be found in Table S2. For
359 pedestrian/bicyclist crashes, the driver-victim pairs are clustered into crashes involving “lower income
360 non-white driver and lower income non-white victim” (LN-LN crashes) and crashes involving “higher
361 income white driver and higher income white victim” (HW-HW crashes), respectively. In the pedestrian
362 crash model, two clusters are almost evenly divided (51.5% for LN-LN crashes and 48.5% for HW-HW
363 crashes). Figure 2.a shows the two clusters and the corresponding distribution of drivers’ and victims’
364 income levels and ethnicity in the pedestrian crash model. For driver’s characteristics, white drivers make
365 9.2% of LN-LN crashes, while its probability is 55% in HW-HW crashes. The income level of drivers in
366 LN-LN crashes concentrates in the low income to medium income categories. In contrast, the income
367 level of drivers in HW-HW crashes is distributed in medium income to high income categories. Victims
368 in LN-LN crashes have a higher probability of being non-whites (81.9%), while victims in HW-HW
369 crashes have the highest probability of being white (49.1%). Victims’ income level is also distributed on
370 low income to medium income in LN-LN crashes and medium income to high income in HW-HW
371 crashes. Clustering results in bicyclist crashes appear to have similar patterns of economic and
372 demographic characteristics for drivers and victims with pedestrian crashes. The bicyclist LN-LN crashes

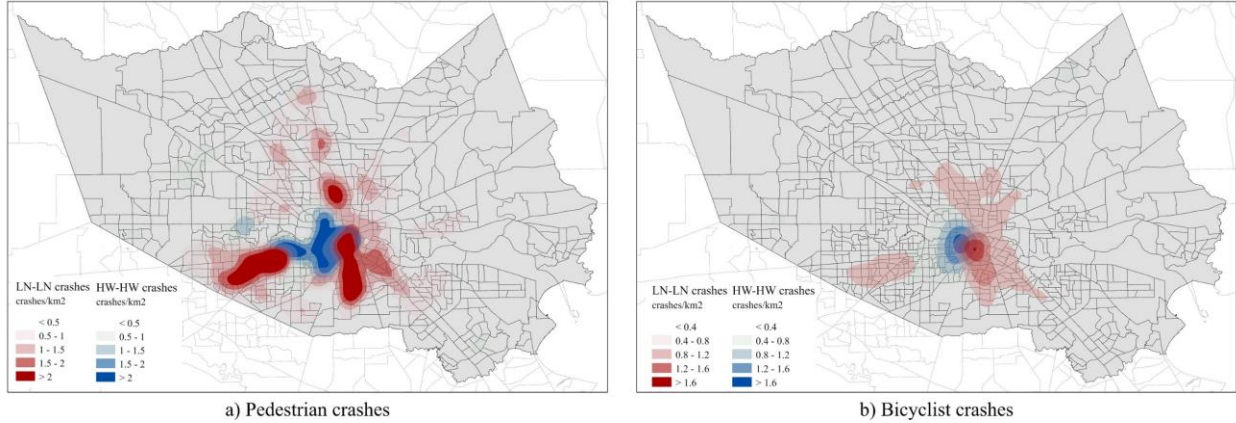
373 have a higher probability of involving non-white drivers (90.6%), drivers from lower income levels and
 374 non-white victims (79.7%), and victims from lower income levels. In comparison, bicyclist HW-HW
 375 crashes have a higher chance of involving white drivers (54.7%), drivers from higher income levels, white
 376 victims (62.1%), and victims from higher income levels. Notable socioeconomic patterns of driver-victim
 377 pairs have been revealed in these results, which show the social segregation of pedestrian and bicyclist
 378 crashes. This social segregation of crashes demonstrates that the driver and victim involved in a crash are
 379 likely to be similar regarding their income and ethnicity. Non-white and low-income drivers and non-
 380 white and low-income victims are more likely to be involved in one crash, while white and high-income
 381 victims are more likely to get into crashes by white and high-income drivers.



382 a) LCA results for pedestrian crashes
 383 a) LCA results for bicyclist crashes

383 **Figure 2. Clustering Results for Pedestrian Crashes and Bicyclist Crashes**

384 As discussed before, there are potential spatial patterns due to the spatial proximity of crashes, so we plot
 385 the density maps of pedestrian and bicyclist crashes based on crash location to observe the spatial
 386 distribution of LN-LN and HW-HW crashes (See Figure 3). Figure 3.a and Figure 3.b show the density of
 387 pedestrian LN-LN crashes and HW-HW crashes of pedestrian crashes and bicyclist crashes, respectively.
 388 The concentrated area of both LN-LN and HW-HW crashes overlay in the downtown area, which
 389 suggests that downtown is the nucleus for crashes of all kinds. Except for downtown Houston, two types
 390 of crashes show spatial segregation in which the LN-LN crashes concentrate on three major areas,
 391 including southern, northern, and further southwest areas near downtown. In contrast, HW-HW crashes
 392 happened more in a closer west region near downtown. Bicyclist crashes show a more segregated pattern
 393 for which LN-LN crashes are denser in the western area near downtown, and HW-HW crashes happened
 394 more in the eastern, southern, northern, and further southwest areas near downtown. Compared with
 395 bicyclist crashes, pedestrian crashes have a denser distribution within the research area.



396

397

Figure 3. Density Map of LN-LN and HW-HW Crashes for Pedestrian and Bicyclist Crashes

398

399

400

401

402

403

404

405

406

407

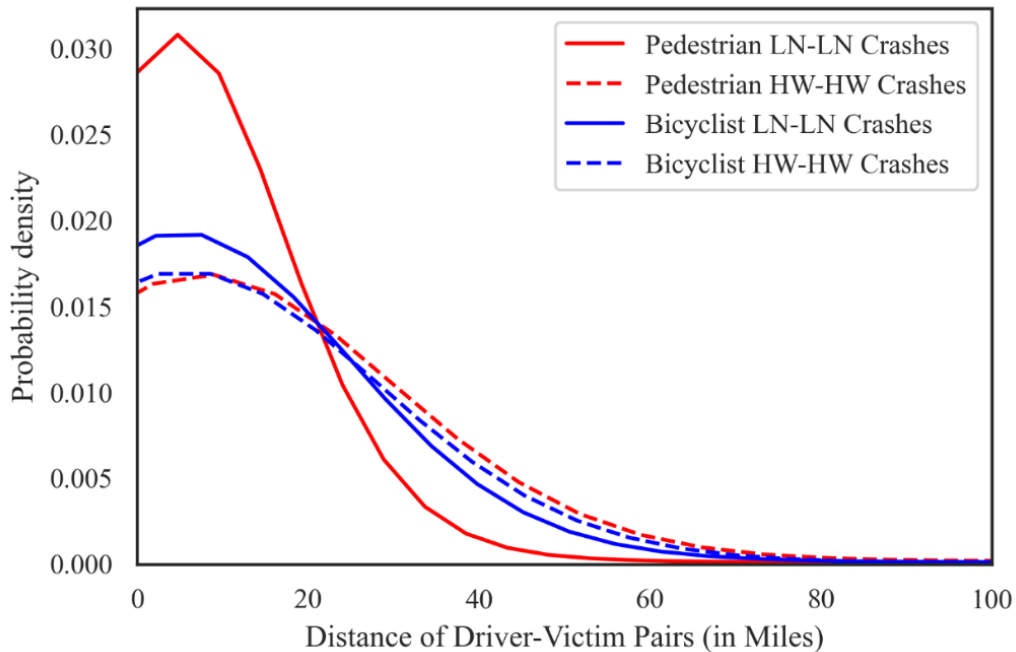
408

409

410

411

Distance between the crash and residence of both parties is important in understanding LN-LN crashes and HW-HW crashes due to the difference in travel behaviors and activity space between drivers and victims from different demographic and socioeconomic backgrounds. We map the trajectory of driver-victim pairs based on the crash location and the centroid of the driver's ZIP code to investigate the spatial characteristics of driver-victim pairs for each cluster (See Figure S1). Downtown Houston is the nucleus of all four types of crashes according to the distribution of driver-victim pairs. Compared with LN-LN crashes, HW-HW crashes trajectory are sparser for both pedestrian and bicyclist crashes. Figure 4 plots the probability density for the distance of driver-victim pairs, showing that LN-LN crashes have a more positively skewed distribution than their counterparts. The geographical and probability distribution of driver-victims pairs indicate that LN-LN crashes are more likely to involve a crash location when a driver lives nearby, while drivers in HW-HW crashes might live farther away from the crash location. This might be because drivers from higher-income and majority-white communities have more resources and capability to travel further away, while drivers in lower-income and minority communities are limited in their activity space hence getting into a crash within their community.



412

413

Figure 4. Probability Density Plot for Distance of Driver-Victim Pairs

414 **5.2 Random Forest Results**

415 *Relative importance of selected variables*

416 Table 3 shows the variables' feature importance of the random forest algorithm to classify whether a
 417 crash belongs to LN-LN crashes or HW-HW crashes for pedestrian and bicyclist crashes, respectively.

418

419 **Table 3. Feature Importance of Random Forest Model for Pedestrian and Bicyclist Crashes**

| Pedestrian Crash model | | | Bicyclist Crash model | | |
|------------------------|--------------------|------|-----------------------|--------------------|------|
| Variables | Feature Importance | Rank | Variables | Feature Importance | Rank |
| <i>EX_Ped</i> | 0.260 | 1 | <i>EX_Cyc</i> | 0.227 | 1 |
| <i>DR_Age</i> | 0.132 | 2 | <i>VT_Age</i> | 0.114 | 2 |
| <i>VT_Age</i> | 0.117 | 3 | <i>DR_Age</i> | 0.103 | 3 |
| <i>CR_CarUsedYr</i> | 0.100 | 4 | <i>EX_AADT</i> | 0.091 | 4 |
| <i>EX_AADT</i> | 0.098 | 5 | <i>CR_CarUsedYr</i> | 0.089 | 5 |
| <i>RD_SpdLmt</i> | 0.049 | 6 | <i>RD_SpdLmt</i> | 0.066 | 6 |
| <i>RD_RdWth</i> | 0.045 | 7 | <i>RD_RdWth</i> | 0.056 | 7 |
| <i>CR_Season</i> | 0.033 | 8 | <i>CR_TimeDay</i> | 0.037 | 8 |
| <i>CR_TimeDay</i> | 0.031 | 9 | <i>CR_Season</i> | 0.034 | 9 |
| <i>RD_LnWth</i> | 0.024 | 10 | <i>RD_LnWth</i> | 0.031 | 10 |
| <i>VT_Gender</i> | 0.014 | 11 | <i>VT_Gender</i> | 0.025 | 11 |
| <i>CR_Intersec</i> | 0.013 | 12 | <i>RD_LnNum</i> | 0.019 | 12 |
| <i>RD_LnNum</i> | 0.013 | 13 | <i>CR_Workday</i> | 0.015 | 13 |
| <i>DR_Gender</i> | 0.012 | 14 | <i>CR_Weather</i> | 0.015 | 14 |
| <i>CR_Surface</i> | 0.012 | 15 | <i>CR_Intersec</i> | 0.015 | 15 |
| <i>CR_Workday</i> | 0.011 | 16 | <i>DR_Gender</i> | 0.014 | 16 |
| <i>CR_Weather</i> | 0.010 | 17 | <i>RD_FuncCls</i> | 0.014 | 17 |
| <i>RD_FuncCls</i> | 0.010 | 18 | <i>CR_Surface</i> | 0.011 | 18 |
| <i>CR_Construt</i> | 0.005 | 19 | <i>RD_MedWth</i> | 0.006 | 19 |
| <i>RD_CurbR</i> | 0.004 | 20 | <i>RD_SWthIn</i> | 0.005 | 20 |
| <i>RD_CurbL</i> | 0.004 | 21 | <i>RD_SWthOut</i> | 0.005 | 21 |
| <i>RD_SWthIn</i> | 0.002 | 22 | <i>RD_CurbL</i> | 0.004 | 22 |
| <i>RD_SWthOut</i> | 0.001 | 23 | <i>RD_CurbR</i> | 0.003 | 23 |
| <i>RD_MedWth</i> | 0.001 | 24 | <i>RD_Urban</i> | <0.001 | 24 |
| <i>RD_Urban</i> | <0.001 | 25 | <i>CR_Construt</i> | <0.001 | 25 |

420

421 Since we only got two clusters in each model, the LN-LN crashes are taken as the reference group. Thus,
 422 the higher value of feature importance a variable has, the larger contribution the variables will make in
 423 determining whether a crash belongs to the LN-LN crash. The ranks of feature importance imply the
 424 relative contribution of a feature in the random forest model. Exposures are the most relevant factors in
 425 determining crash clusters. The estimated pedestrian exposure and estimated pedestrian exposure rank
 426 first in their respective model. AADT ranks fifth in pedestrian crashes and ranks fourth in bicyclist
 427 crashes. The high rank of exposure variables indicates a strong association between the traffic volume of
 428 both vehicles and pedestrians/bicyclists with patterns of driver-victim pair. The driver's age and victim's

429 age are also among the most important variables, while their gender is less influential. Driver's age and
430 victim's age rank second and third in pedestrian crashes, and driver's age and victim's age rank third and
431 second in bicyclist crashes. For crash specific information, the year of the car in use ranks fourth in
432 pedestrian crashes, and fifth in bicyclist crashes, indicating the vehicles involved in LN-LN and HW-HW
433 crashes might have different used years. Time of the day and season ranks eighth in both models,
434 indicating a relatively sizeable temporal variation of the crash pattern. For road infrastructure
435 characteristics, speed limit and roadbed width rank sixth and seventh in both pedestrian and bicyclist
436 crashes, showing the considerable influence of roadway infrastructure characteristics in determining the
437 crash clusters. However, their feature importance is relatively low compared to previous factors.

438

439 **5.3 Partial Dependence Plots**

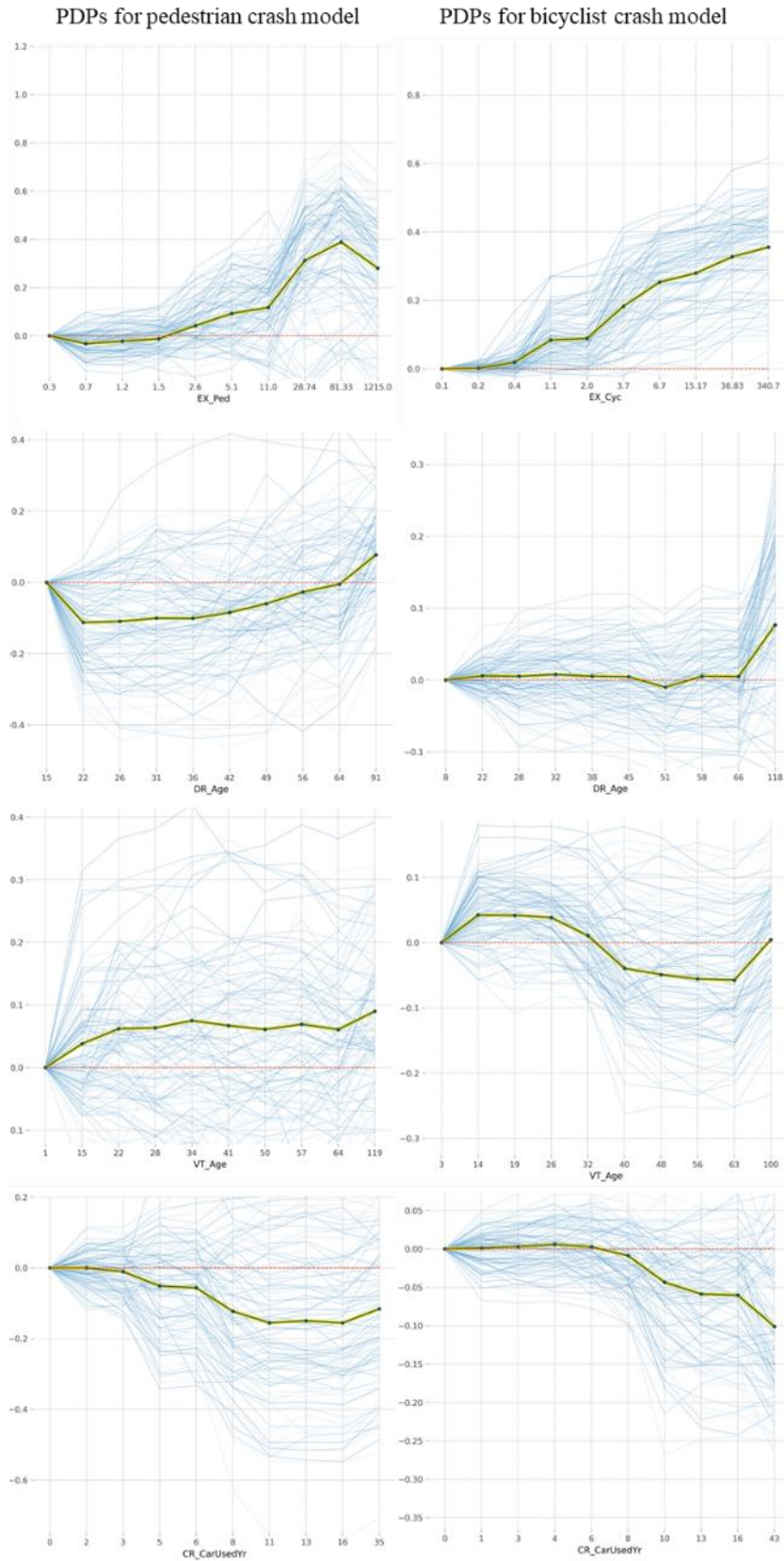
440 To investigate variables' impact on driver-victim pairs' patterns, we draw the PDPs for the top eight
441 variables in feature importance for both pedestrian and bicyclist models (Figure 5). For exposure
442 variables, when annual daily pedestrian volumes are less than 2.6, pedestrian exposure is not influential.
443 When it is larger than 2.6, it becomes positively associated with the probability of a crash being a LN-LN
444 crash. This indicates that LN-LN crashes will likely happen on the road with larger pedestrian exposure.
445 HW-HW crashes will be less likely to occur on the road with larger pedestrian exposure. In bicyclist
446 crashes, the positive marginal effect of bicyclist exposure on the probability of a crash being a LN-LN
447 crash will increase when the bicyclist exposure becomes larger. This indicates that LN-LN crashes will be
448 more likely to happen on the road with larger bicyclist exposure, and the larger the bicyclist exposure, the
449 higher the probability of LN-LN crashes. One of the potential explanations for this could be missing
450 active transportation-friendly infrastructure in low income and minority communities, which may force
451 the bicyclists to share the road with oncoming traffic, increasing their crash probability. However, this
452 speculation needs further explored and proven by accounting for the bicyclist infrastructure in data
453 analysis. For vehicular exposure, the pedestrian and bicyclist crashes have similar patterns, which shows
454 lower AADT does not have significant influence on the probability of a crash being LN-LN crash. While
455 within the highest quantile of the AADT, it will have a larger positive association for both pedestrian and
456 bicyclist crashes. This means both pedestrian and bicyclist LN-LN crashes tend to occur on the road with
457 a larger vehicular volume.

458 The driver's and victim's age are among the most influential factors in socioeconomic characteristics for
459 both crash types. In the pedestrian crash model, when the driver's age is less than 64, the probability of a
460 crash being a LN-LN crash will decrease. When the driver's age is larger than 64, the probability of a
461 crash being a LN-LN crash will increase. This means younger drivers are less likely to be involved in a
462 pedestrian LN-LN crash, while older drivers are more likely to be involved in a pedestrian LN-LN crash.
463 The PDP shows that as the victim's age increases, the marginal effect of the probability of being in a LN-
464 LN crash will rise, indicating that older victims are more likely to be involved in a pedestrian LN-LN
465 crash. In the bicyclist crash model, the driver's age does not have much influence on the probability of a
466 LN-LN crash in its lower quintiles. It only has a positive marginal effect when the driver's age exceeds 66
467 years old, indicating that older drivers will be more likely to be involved in a bicyclist LN-LN crash. For
468 the victim's age, a victim aged 32 or younger will increase the probability of a crash being a LN-LN
469 crash, while a victim aged 33 or higher will decrease the probability of a crash being a LN-LN crash. This
470 means bicyclist LN-LN crashes are more likely to involve older drivers and younger bicyclists.

471 For crash specific information, the year of the car in use, time of the day, and season rank among the most
472 influential variables. When the year of the car in use is less than six, it has negligible influence on the
473 probability of a pedestrian LN-LN crash for both crash types. As the year of the car in use increases in
474 pedestrian crashes, its marginal effect will become larger in a negative direction for both crash types. This
475 indicates older cars are less likely to be involved in a LN-LN crash and more likely to be involved in an

476 HW-HW crash for both pedestrian and bicyclist crashes. The influence of time of the day on pedestrian
477 LN-LN crashes is positive in summer and autumn and negative in winter, but the effect of the influence is
478 minimal. For bicyclist crashes, from 6:00 am to 12:00 pm and from 12:00 pm to 6:00 pm, there will be a
479 higher chance of bicyclist crashes. Lower-income and non-white groups might choose biking as their
480 mean of transportation to commute during the daytime more frequently than their counterparts due to
481 economic affordability or behavioral difference, which forms a higher bicyclist crash probability.

482 For road infrastructure characteristics, the road speed limit has the same patterns in its influence on
483 pedestrian and bicyclist crashes. When the road speed limit is less than 35 miles per hour, its impact on
484 the crash clusters is negligible. When the road speed limit exceeds 45 miles per hour, the probability of a
485 crash being a LN-LN crash will increase in both pedestrian and bicyclist crashes. This indicates that LN-
486 LN crashes for pedestrians and bicyclists are more likely to happen on the road with a higher speed limit.
487 Roadbed width has little effect when less than 40 feet and only has a positive marginal effect on the
488 highest quantile, indicating that LN-LN pedestrian crashes are more likely to happen on wider roads. In
489 the bicyclist crash model, the effect of roadbed width is not influential when it is less than 24 feet but
490 becomes negative when it is larger than 24 feet, suggesting that LN-LN bicyclist crashes are less likely to
491 happen on wider roads.

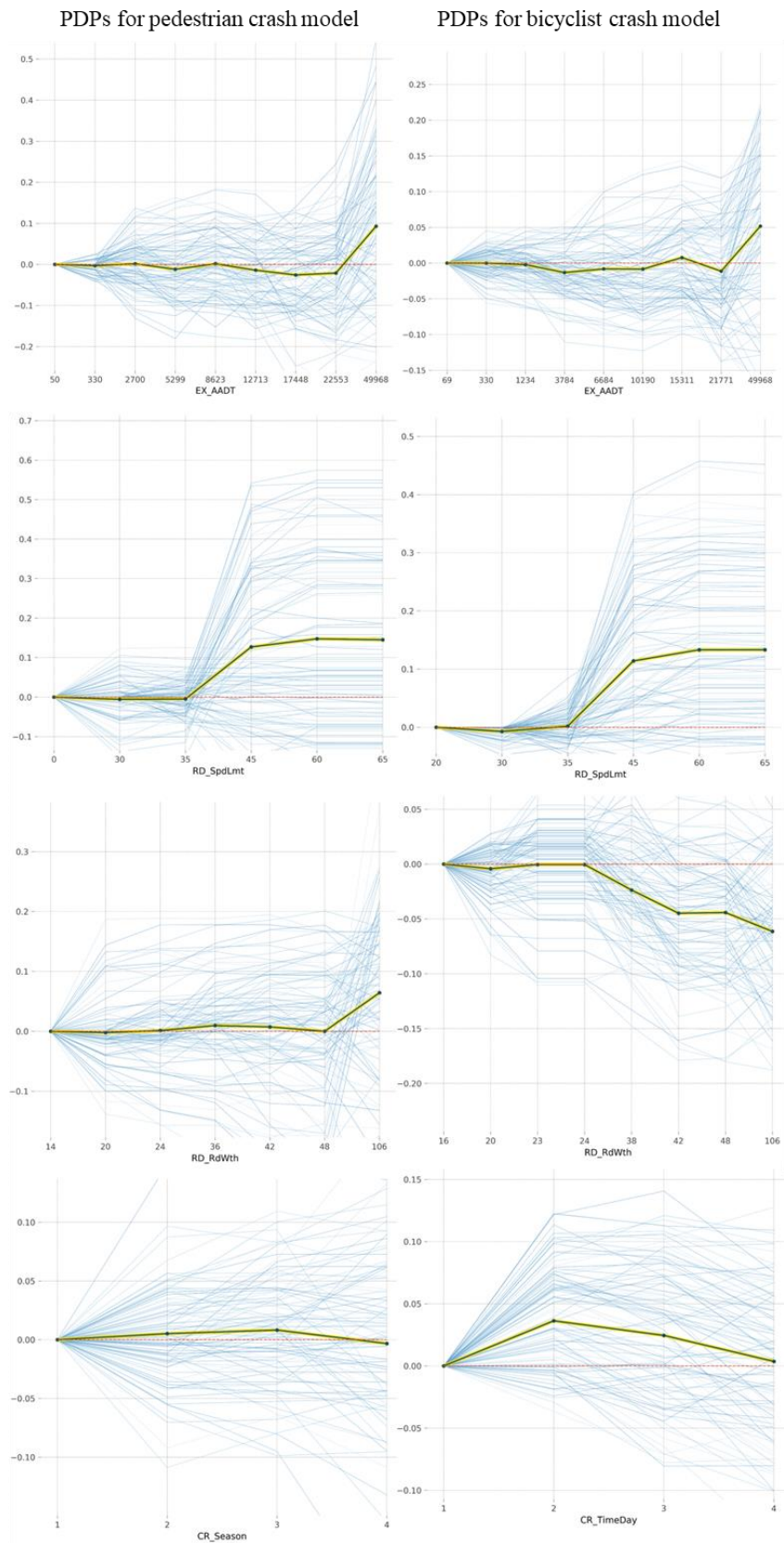


492

493

Figure 5. PDPs for Variables in Pedestrian and Bicyclist Crash Model

(Figure 5 continued)



496 6. CONCLUSION

497 In this study, we used driver-victim pairs to reveal the crash patterns based on clustering drivers' and
498 victims' ethnicity and income level. Using crash data from Harris County, we applied a probability-based
499 latent class clustering analysis to classify pedestrian and bicyclist crashes. The clustering results showed
500 that lower income and non-white drivers tend to be involved in crashes with lower income and non-white
501 victims (LN-LN crashes). While higher income and white drivers tend to be involved in crashes with
502 higher income and white victims (HW-HW crashes). This result showed a certain degree of social
503 segregation in pedestrian and bicyclist crashes, indicating that drivers and victims of similar
504 socioeconomic characteristics are more likely to be involved in the same crash, while those from different
505 socioeconomic backgrounds are not. We further analyzed the trajectories of driver-victim pairs and found
506 all crash types tend to concentrate in downtown Houston. The trajectories of HW-HW crashes are sparser
507 in their geographic distribution, which suggests higher income and white drivers are driving a long
508 distance and getting involved in a crash in farther geographic areas than their counterparts.

509 To explore how the LN-LN and HW-HW crash patterns were shaped, we employed a random forest
510 algorithm and partial dependence plots to model and interpret the clustering outcomes from LCA models.
511 Contributing factors for the crash patterns were selected from crash specific information, drivers' and
512 victims' age and gender, roadway infrastructure, and traffic exposure. Pedestrian/bicyclist exposure,
513 driver's age, victim's age, year of the car in use, AADT, speed limit, roadbed width, time of the day, and
514 season are the most influential variables in pedestrian and bicyclist models. We drew partial dependence
515 plots for the most influential variables to interpret how the variables are associated with crash patterns.
516 The results showed that LN-LN crashes tend to happen on the road with larger traffic exposure of
517 pedestrians/bicyclists and vehicle, which is contradictory to safety in number theory, indicating that the
518 European model of bicycling/walking is not always implementable for underserved communities in the
519 US (Elvik and Bjørnskau, 2017). Older drivers and older pedestrians are more likely to be in the same
520 LN-LN crash, while older drivers and younger bicyclists are more likely to be in the same LN-LN crash.
521 Longer years of the car in use will increase the probability of HW-HW crashes. Higher speed limits and
522 wider roads are associated with a higher probability of LN-LN crashes for both pedestrian and bicyclist
523 crashes. The results indicated the coexistence of LN-LN crashes and road conditions of higher traffic
524 exposure, higher speed limit, and wider roads. The communities where low-income and ethnic minorities
525 are concentrated might have higher traffic exposure and less safe road environments, which shapes the
526 distribution of LN-LN crashes.

527 This study contributes to the existing body of literature in several ways. First, from a planning and
528 engineering perspective, this study confirms long-believed hypotheses that there is a clear
529 sociodemographic and economic segregation of crashes. We also find that the crash-contributing factors
530 are not usually the same across different communities. These results can help safety practitioners in both
531 engineering and planning fields to develop and implement practices that will target the main concerns of
532 each community instead of developing one size fits all strategies. Safe systems approach can be one of the
533 potential strategies to accomplish this goal. Another significant contribution of this study concerns the
534 methodological approach. We innovatively use machine learning techniques to address a largely
535 unexplored research question where the driver's and victim's characteristics are analyzed simultaneously.

536 Despite these contributions, the study does have limitations. In this study, we used the police-reported
537 crash data, which have been considered to underestimate the actual number of crashes. Besides, the
538 police-reported crash data also lacks other economic and demographic information for drivers and
539 victims, such as educational level and occupation. The detailed income level is also not reported by the
540 police agents. On the other hand, collecting individual level income data is not be feasible in an

541 observational study and may require additional data collection efforts by implementing experimental
542 design studies. The success of such experimental design study however is not guaranteed given that many
543 drivers may be reluctant to share their crash history due to potential liabilities. We therefore use the
544 surrogate measurement of income level based on drivers' and victims' residential census tract. This
545 approach may be biased, but it is an acceptable alternative in the absence of readily available data on
546 income measurement. Another limitation of the study is related to the exposure data. Although we
547 account for bicycle and pedestrian exposure by developing scaling factors, the measurement of exposure
548 can be improved by implementing more rigorous models.

549 In this study we used Harris County as the pilot site, which might not be robust, but the analytical
550 methods can be generalized to other cities and regions with the availability of data. We also do not
551 account for the bicycle and pedestrian infrastructure such as the quality of sidewalk or bike lane, which
552 can help to explain some of the findings of this research. Future studies will try to address these
553 limitations by implementing rigorous statistical models and image analysis tools to obtain the
554 infrastructure information.

555 REFERENCES

- 556 Adanu, E. K., Smith, R., Powell, L. & Jones, S. 2017. Multilevel analysis of the role of human factors in
557 regional disparities in crash outcomes. *Accident Analysis & Prevention*, 109, 10-17.
- 558 Algurén, B., & Rizzi, M. 2022. In-depth understanding of single bicycle crashes in Sweden—Crash
559 characteristics, injury types and health outcomes differentiated by gender and age-groups. *Journal*
560 *of Transport & Health*, 24,
- 561 Balakrishnan, S., Moridpour, S., & Tay, R. 2019. Sociodemographic influences on injury severity in
562 truck-vulnerable road user crashes. *ASCE-ASME Journal of Risk and Uncertainty in Engineering*
563 *Systems, Part A: Civil Engineering*, 5(4), 04019015.6.
- 564 Barajas, J. M. 2018. Not all crashes are created equal: Associations between the built environment and
565 disparities in bicycle collisions. *Journal of Transport and Land Use*, 11(1), 865-882.
- 566 Behnood, A., & Mannering, F. 2017. Determinants of bicyclist injury severities in bicycle-vehicle
567 crashes: A random parameters approach with heterogeneity in means and variances. *Analytic*
568 *Methods in Accident Research*, 16, 35–47.
- 569 Billah, K., Sharif, H. O. & Dessouky, S. 2022. How gender affects motor vehicle crashes: A case study
570 from San Antonio, Texas. *Sustainability*, 14(12), 7023.
- 571 Boele-Vos, M. J., Van Duijvenvoorde, K., Doumen, M. J. A., Duivenvoorden, C. W. A. E., Louwse, W.
572 J. R., & Davidse, R. J. 2017. Crashes involving cyclists aged 50 and over in the Netherlands: An
573 in-depth study. *Accident Analysis & Prevention*, 105, 4–10.
- 574 Boufous, S., Rome, L. D., Senserrick, T., & Ivers, R. 2011. Cycling Crashes in Children, Adolescents,
575 and Adults—A Comparative Analysis. *Traffic Injury Prevention*, 12(3), 244–250.
- 576 Braun, L. M., Rodriguez, D. A., & Gordon-Larsen, P. (2019). Social (in)equality in access to cycling
577 infrastructure: Cross-sectional associations between bike lanes and area-level sociodemographic
578 characteristics in 22 large U.S. cities. *Journal of Transport Geography*, 80, 102544.

- 579 Braver, E. R. 2004. Are older drivers actually at higher risk of involvement in collisions resulting in
580 deaths or non-fatal injuries among their passengers and other road users? *Injury Prevention*,
581 10(1), 27–32.
- 582 Burdett, B. R. D., Starkey, N. J., & Charlton, S. G. 2017. The close to home effect in road crashes. *Safety*
583 *Science*, 98, 1–8.
- 584 Dadashova, B., & Griffin, G. P. 2020. Random parameter models for estimating statewide daily bicycle
585 counts using crowdsourced data. *Transportation Research Part D: Transport and Environment*, 84,
586 102368.
- 587 Dadashova, B., Griffin, G. P., Das, S., Turner, S., & Sherman, B. 2020. Estimation of Average Annual
588 Daily Bicycle Counts using Crowdsourced Strava Data. *Transportation Research Record: Journal*
589 *of the Transportation Research Board*, 2674(11), 390–402.
- 590 Das, S., Bibeka, A., Sun, X., Zhou, H. “Tracy,” & Jalayer, M. 2019. Elderly pedestrian fatal crash-related
591 contributing factors: applying empirical Bayes geometric mean method. *Transportation Research*
592 *Record: Journal of the Transportation Research Board*, 2673(8), 254–263.
- 593 Ding, H., Sze, N. N., Li, H., & Guo, Y. 2020. Roles of infrastructure and land use in bicycle crash
594 exposure and frequency: A case study using Greater London bike sharing data. *Accident Analysis*
595 *& Prevention*, 144, 105652.
- 596 Elvik, R. and Bjørnskau, T., 2017. Safety-in-numbers: a systematic review and meta-analysis of evidence.
597 *Safety science*, 92, 274-282.
- 598 Eren, H., & Gauld, C. 2022. Smartphone use among young drivers: Applying an extended Theory of
599 Planned Behaviour to predict young drivers’ intention and engagement in concealed responding.
600 *Accident Analysis & Prevention*, 164, 106474.
- 601 Ferenchak, N. N., & Marshall, W. E. 2021. Bicycling facility inequalities and the causality dilemma with
602 socioeconomic/sociodemographic change. *Transportation Research Part D: Transport and*
603 *Environment*, 97, 102920.
- 604 Fuller, D., & Winters, M. 2017. Income inequalities in Bike Score and bicycling to work in Canada.
605 *Journal of Transport & Health*, 7, 264–268.
- 606 Glassbrenner, D., Herbert, G., Reish, L., Webb, C., & Lindsey, T. 2022. Evaluating disparities in traffic
607 fatalities by race, ethnicity, and income (Report No. DOT HS 813 188).
608 <https://crashstats.nhtsa.dot.gov/Api/Public/ViewPublication/813188>. Accessed October 5th, 2022.
- 609 Gong, L. & Fan, W. D. 2017. Modeling single-vehicle run-off-road crash severity in rural areas:
610 Accounting for unobserved heterogeneity and age difference. *Accident Analysis & Prevention*,
611 101, 124-134.
- 612 Governors Highway Safety Association (GHSA). 2021. An Analysis of Traffic Fatalities by Race and
613 Ethnicity. [https://www.ghsa.org/resources/Analysis-of-Traffic-Fatalities-by-Race-and-](https://www.ghsa.org/resources/Analysis-of-Traffic-Fatalities-by-Race-and-Ethnicity21)
614 [Ethnicity21](https://www.ghsa.org/resources/Analysis-of-Traffic-Fatalities-by-Race-and-Ethnicity21). Accessed October 5th, 2022.

- 615 Hasheminejad, S. H.-A., Zahedi, M., & Hasheminejad, S. M. H. (2018). A hybrid clustering and
616 classification approach for predicting crash injury severity on rural roads. *International Journal of*
617 *Injury Control and Safety Promotion*, 25(1), 85–101.
- 618 Ivan, K., Benedek, J., & Ciobanu, S. 2019. School-aged pedestrian–vehicle crash vulnerability.
619 *Sustainability*, 11(4), 1214.
- 620 Kemnitzer, C. R., Pope, C. N., Nwosu, A., Zhao, S., Wei, L. & Zhu, M. 2019. An investigation of driver,
621 pedestrian, and environmental characteristics and resulting pedestrian injury. *Traffic Injury*
622 *Prevention*, 20, 510-514.
- 623 Koopmans, J. M., Friedman, L., Kwon, S., & Sheehan, K. 2015. Urban crash-related child pedestrian
624 injury incidence and characteristics associated with injury severity. *Accident Analysis &*
625 *Prevention*, 77, 127–136.
- 626 Kravetz, D. & Noland, R. B. 2012. Spatial analysis of income disparities in pedestrian safety in northern
627 New Jersey: is there an environmental justice issue? *Transportation Research Record: Journal of*
628 *the Transportation Research Board*, 2320, 10-17.
- 629 Lee, J., Li, X., Mao, S., Fu, W. & Moridpour, S. 2021. Investigation of contributing factors to traffic
630 crashes and violations: A random parameter multinomial logit approach. *Journal of Advanced*
631 *Transportation*, 2021, 1-11.
- 632 Li, Z., Chen, C., Ci, Y., Zhang, G., Wu, Q., Liu, C., & Qian, Z. (Sean). (2018). Examining driver injury
633 severity in intersection-related crashes using cluster analysis and hierarchical Bayesian models.
634 *Accident Analysis & Prevention*, 120, 139–151.
- 635 Liang, O. S. & Yang, C. C. 2022. How are different sources of distraction associated with at-fault crashes
636 among drivers of different age gender groups? *Accident Analysis & Prevention*, 165, 106505.
- 637 Lombardi, D. A., Horrey, W. J., & Courtney, T. K. 2017. Age-related differences in fatal intersection
638 crashes in the United States. *Accident Analysis & Prevention*, 99, 20–29.
- 639 Martínez-Ruiz, V., Jiménez-Mejías, E., Luna-del-Castillo, J. de D., García-Martín, M., Jiménez-Moleón,
640 J. J., & Lardelli-Claret, P. 2014. Association of cyclists’ age and sex with risk of involvement in a
641 crash before and after adjustment for cycling exposure. *Accident Analysis & Prevention*, 62, 259–
642 267.
- 643 Masís, S. 2021. *Interpretable machine learning with Python: Learn to build interpretable high-*
644 *performance models with hands-on real-world examples*, Packt Publishing Ltd.
- 645 Mokhtarimousavi, S., Anderson, J. C., Azizinamini, A., & Hadi, M. (2020). Factors affecting injury
646 severity in vehicle-pedestrian crashes: A day-of-week analysis using random parameter ordered
647 response models and Artificial Neural Networks. *International Journal of Transportation Science*
648 *and Technology*, 9(2), 100–115.
- 649 National Highway Traffic Safety Administration (NHTSA). 2021. Early Estimates of Motor Vehicle
650 Traffic Fatalities and Fatality Rate by Sub-Categories in 2020.
651 <https://crashstats.nhtsa.dot.gov/Api/Public/ViewPublication/813118>. Accessed October 5th, 2022.

652 Oviedo-Trespalacios, O., & Scott-Parker, B. (2018). Young drivers and their cars: Safe and sound or the
653 perfect storm? *Accident Analysis & Prevention*, 110, 18–28.

654 Prati, G., Fraboni, F., De Angelis, M., Pietrantoni, L., Johnson, D., & Shires, J. 2019. Gender differences
655 in cycling patterns and attitudes towards cycling in a sample of European regular cyclists. *Journal*
656 *of Transport Geography*, 78, 1–7.

657 Pulido, J., Barrio, G., Hoyos, J., Jimenez-Mejias, E., Martin-Rodriguez Mdel, M., Houwing, S. &
658 Lardelli-Claret, P. 2016. The role of exposure on differences in driver death rates by gender and
659 age: Results of a quasi-induced method on crash data in Spain. *Accident Analysis Prevention*, 94,
660 162-7.

661 Rahul, T. M., & Verma, A. 2014. A study of acceptable trip distances using walking and cycling in
662 Bangalore. *Journal of Transport Geography*, 38, 106–113.

663 Rebentisch, H., Wasfi, R., Piatkowski, D. P. & Manaugh, K. 2019. Safe streets for all? Analyzing
664 infrastructural response to pedestrian and cyclist crashes in New York City, 2009–2018.
665 *Transportation Research Record: Journal of the Transportation Research Board*, 2673, 672-685.

666 Regev, S., Rolison, J. J. & Moutari, S. 2018. Crash risk by driver age, gender, and time of day using a
667 new exposure methodology. *Journal of Safety Research*, 66, 131-140.

668 Roll, J., & McNeil, N. 2022. Race and income disparities in pedestrian injuries: factors influencing
669 pedestrian safety inequity. *Transportation Research Part D: Transport and Environment*, 107,
670 103294.

671 Rothman, L., Cloutier, M.-S., Manaugh, K., Howard, A. W., Macpherson, A. K., & Macarthur, C. 2020.
672 Spatial distribution of roadway environment features related to child pedestrian safety by census
673 tract income in Toronto, Canada. *Injury Prevention*, 26(3), 229–233.

674 Russo, F., Biancardo, S. A. & Dell'acqua, G. 2014. Road safety from the perspective of driver gender and
675 age as related to the injury crash frequency and road scenario. *Traffic Injury Prevention*, 15, 25-
676 33.

677 Sagar, S., Stamatiadis, N. & Stromberg, A. 2021. Effect of socioeconomic and demographic factors on
678 crash occurrence. *Transportation Research Record: Journal of the Transportation Research Board*,
679 2675, 80-91.

680 Saha, D., Alluri, P., Gan, A. & Wu, W. 2018. Spatial analysis of macro-level bicycle crashes using the
681 class of conditional autoregressive models. *Accident Analysis & Prevention*, 118, 166-177.

682 Salon, D., & McIntyre, A. 2018. Determinants of pedestrian and bicyclist crash severity by party at fault
683 in San Francisco, CA. *Accident Analysis & Prevention*, 110, 149–160.

684 Samerei, S. A., Aghabayk, K., Shiwakoti, N. & Mohammadi, A. 2021. Using latent class clustering and
685 binary logistic regression to model Australian cyclist injury severity in motor vehicle-bicycle
686 crashes. *Journal of Safety Research*, 79, 246-256.

687 Sartin, E. B., Metzger, K. B., Pfeiffer, M. R., Myers, R. K. & Curry, A. E. 2021. Facilitating research on
688 racial and ethnic disparities and inequities in transportation: Application and evaluation of the

- 689 Bayesian Improved Surname Geocoding (BISG) algorithm. *Traffic Injury Prevention*, 22, S32-
690 S37.
- 691 Scikit-Learn. 2022. Tuning the hyper-parameters of an estimator [Online]. Available: [https://scikit-](https://scikit-learn.org/stable/modules/grid_search.html)
692 [learn.org/stable/modules/grid_search.html](https://scikit-learn.org/stable/modules/grid_search.html) [Accessed 07/30/2022].
- 693 Scott-Parker, B., & Oviedo-Trespalacios, O. 2017. Young driver risky behaviour and predictors of crash
694 risk in Australia, New Zealand and Colombia: Same but different? *Accident Analysis &*
695 *Prevention*, 99, 30–38.
- 696 Siddiqui, C., Abdel-Aty, M., & Choi, K. 2012. Macroscopic spatial analysis of pedestrian and bicycle
697 crashes. *Accident Analysis & Prevention*, 45, 382–391.
- 698 Sivasankaran, S. K., & Balasubramanian, V. (2020). Exploring the severity of bicycle–vehicle crashes
699 using latent class clustering approach in India. *Journal of Safety Research*, 72, 127–138.
- 700 Steinbach, R., Edwards, P. & Grundy, C. 2013. The road most travelled: The geographic distribution of
701 road traffic injuries in England. *International Journal of Health Geographics*, 12(1), 1-7.
- 702 Steinbach, R., Green, J., Kenward, M. G., & Edwards, P. 2016. Is ethnic density associated with risk of
703 child pedestrian injury? A comparison of inter-census changes in ethnic populations and injury
704 rates. *Ethnicity & Health*, 21(1), 1–19.
- 705 Sun, M., Sun, X. & Shan, D. 2019. Pedestrian crash analysis with latent class clustering method. *Accident*
706 *Analysis Prevention*, 124, 50-57.
- 707 Tom, A., & Granić, M.-A. 2011. Gender differences in pedestrian rule compliance and visual search at
708 signalized and unsignalized crossroads. *Accident Analysis & Prevention*, 43(5), 1794–1801.
- 709 Toran Pour, A., Moridpour, S., Tay, R., & Rajabifard, A. 2018. Influence of pedestrian age and gender on
710 spatial and temporal distribution of pedestrian crashes. *Traffic Injury Prevention*, 19(1), 81–87.
- 711 Ulak, M. B., Kocatepe, A., Ozguven, E. E., & Horner, M. W. 2019. How far from home do crashes
712 occur? A network based analysis. *Safety Science*, 118, 298–308.
- 713 Wang, J., & Lindsey, G. 2017. Equity of bikeway distribution in Minneapolis, Minnesota. *Transportation*
714 *Research Record: Journal of the Transportation Research Board*, 2605(1), 18–31.
- 715 Weiss, H. B., Kaplan, S., & Prato, C. G. (2014). Analysis of factors associated with injury severity in
716 crashes involving young New Zealand drivers. *Accident Analysis & Prevention*, 65, 142–155.
- 717 Wheeler-Martin, K. C., Curry, A. E., Metzger, K. B. & Dimaggio, C. J. 2020. Trends in school-age
718 pedestrian and pedalcyclist crashes in the USA: 26 states, 2000-2014. *Injury Prevention*, 26, 448-
719 455.
- 720 Xiao, D., Šarić, Ž., Xu, X., & Yuan, Q. (2022). Investigating injury severity of pedestrian–vehicle crashes
721 by integrating latent class cluster analysis and unbalanced panel mixed ordered probit model.
722 *Journal of Transportation Safety & Security*, 1–20.
723 <https://doi.org/10.1080/19439962.2022.2033900>

724 Zhao, H., Yang, G., Zhu, F., Jin, X., Begeman, P., Yin, Z., Yang, K. H., & Wang, Z. 2013. An
725 Investigation on the Head Injuries of Adult Pedestrians by Passenger Cars in China. *Traffic Injury*
726 *Prevention*, 14(7), 712–717.

727

728

729 APPENDIX

730

731 Table S1. Scaling Factors for Pedestrian and Bicyclist Crash from 2017 to 2019

| Crash type | Year | Scaling factor |
|--------------------|------|----------------|
| Pedestrian crashes | 17 | 9.99 |
| | 18 | 9.97 |
| | 19 | 10.29 |
| | 20 | 6.06 |
| Bicyclist crashes | 17 | 48.31 |
| | 18 | 47.7 |
| | 19 | 46.53 |
| | 20 | 21.32 |

732

733 Table S2. Latent Class Cluster Results for Pedestrian and Bicyclist Crashes

| Variables | Pedestrian Crashes | | | Bicyclist Crashes | | | |
|---------------------|--------------------|-------------|-------------|-------------------|------------|------------|------------|
| | Total | LN-LN | HW-HW | Total | LN-LN | HW-HW | |
| <i>DR_IncLvl</i> | Low | 456(16.2%) | 396(27.1%) | 60(4.4%) | 192(17.1%) | 160(31.7%) | 32(5.2%) |
| | Low to medium | 615(21.8%) | 439(30.0%) | 176(13.0%) | 229(20.4%) | 144(28.5%) | 85(13.7%) |
| | Medium | 815(28.9%) | 445(30.4%) | 370(27.2%) | 280(24.9%) | 123(24.3%) | 157(25.4%) |
| | Medium to high | 366(13.0%) | 98(6.7%) | 268(19.8%) | 161(14.3%) | 42(8.3%) | 119(19.3%) |
| | High | 569(20.2%) | 84(5.8%) | 485(35.7%) | 260(23.2%) | 36(7.2%) | 224(36.4%) |
| <i>DR_Ethnicity</i> | White | 883(31.3%) | 135(9.2%) | 748(55.0%) | 385(34.3%) | 47(9.4%) | 338(54.7%) |
| | Hispanic | 900(31.9%) | 722(49.3%) | 178(13.1%) | 349(31.1%) | 230(45.5%) | 119(19.3%) |
| | Black | 802(28.4%) | 539(36.8%) | 263(19.3%) | 296(26.4%) | 202(39.9%) | 94(15.3%) |
| | Asian | 184(6.5%) | 45(3.1%) | 139(10.2%) | 70(6.2%) | 18(3.5%) | 52(8.5%) |
| | Other | 53(1.9%) | 21(1.5%) | 32(2.3%) | 22(2.0%) | 9(1.7%) | 13(2.2%) |
| <i>VT_IncLvl</i> | Low | 602(21.3%) | 493(33.7%) | 109(8.0%) | 210(18.7%) | 188(37.2%) | 22(3.6%) |
| | Low to medium | 764(27.1%) | 545(37.3%) | 219(16.1%) | 285(25.4%) | 171(33.9%) | 114(18.5%) |
| | Medium | 558(19.8%) | 217(14.8%) | 341(25.1%) | 210(18.7%) | 75(14.9%) | 135(21.9%) |
| | Medium to high | 388(13.7%) | 112(7.7%) | 276(20.3%) | 156(13.9%) | 34(6.8%) | 122(19.8%) |
| | High | 510(18.1%) | 96(6.5%) | 414(30.5%) | 261(23.2%) | 37(7.3%) | 224(36.3%) |
| <i>VT_Ethnicity</i> | White | 932(33.0%) | 265(18.1%) | 667(49.1%) | 486(43.3%) | 103(20.3%) | 383(62.1%) |
| | Hispanic | 854(30.3%) | 600(41.0%) | 254(18.7%) | 273(24.3%) | 171(33.8%) | 102(16.6%) |
| | Black | 844(29.9%) | 525(35.9%) | 319(23.5%) | 298(26.5%) | 221(43.8%) | 77(12.4%) |
| | Asian | 130(4.6%) | 50(3.4%) | 80(5.9%) | 52(4.6%) | 9(1.8%) | 43(7.0%) |
| | Other | 61(2.2%) | 23(1.5%) | 38(2.8%) | 13(1.2%) | 2(0.3%) | 11(1.8%) |
| Total | 2822(100.0%) | 1463(51.5%) | 1359(48.5%) | 1123(100.0%) | 506(45.3%) | 617(54.7%) | |

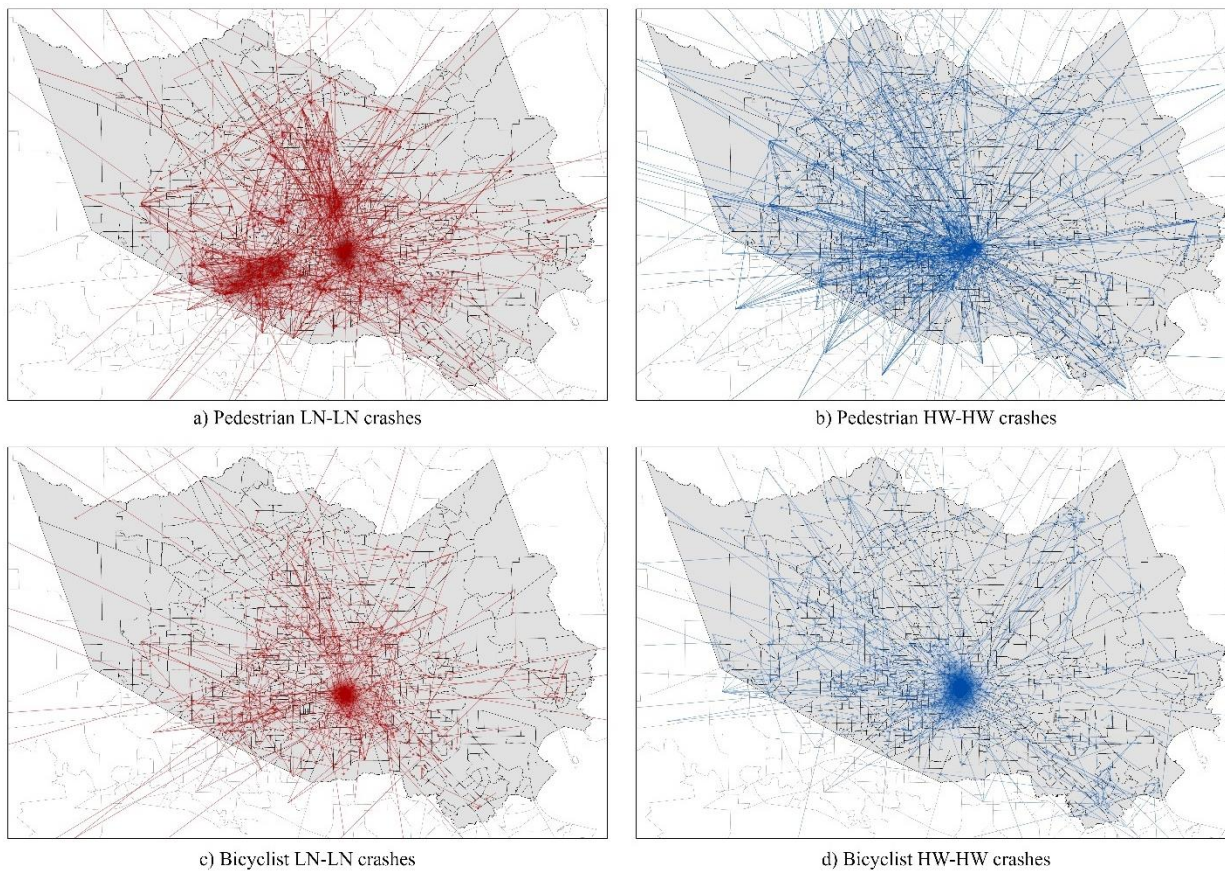
734

735

736

737

738



739

740

Figure S1. The Trajectory of Driver-Victim Pairs for LN-LN and HW-HW Crashes

741