

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

**An ecological examination of early adolescent e-cigarette use: A machine learning approach to understanding a health epidemic**

Alejandro L. Vázquez <sup>1\*</sup>, Cynthia M. Navarro Flores <sup>1</sup>, Byron H. Garcia <sup>2</sup>, Tyson S. Barrett <sup>3</sup>,  
Melanie M. Domenech Rodríguez <sup>4</sup>

<sup>1</sup> Department of Psychology, University of Tennessee, Knoxville, Knoxville, Tennessee, United States

<sup>2</sup> Department of Psychology, Arizona State University, Tempe, Arizona, United States

<sup>3</sup> Highmark Health, Pittsburg, Pennsylvania, United States

<sup>4</sup> Department of Psychology, Utah State University, Logan, Utah, United States

\* Correspondence author

Email: [avazquez4004@gmail.com](mailto:avazquez4004@gmail.com)

## 22 **Abstract**

23 E-cigarette use among adolescents is a national health epidemic spreading faster than researchers  
24 can amass evidence for risk and protective factors and long-term consequences associated with  
25 use. New technologies, such as machine learning, may assist prevention programs in identifying  
26 at-risk youth and potential targets for intervention before adolescents enter developmental  
27 periods where e-cigarette use escalates. The current study utilized machine learning algorithms  
28 to explore a wide array of individual and socioecological variables in relation to patterns of  
29 lifetime e-cigarette use during early adolescence (i.e., exclusive, or with tobacco). Extant data  
30 was used from 14,346 students middle school students ( $M_{\text{age}} = 12.5$ ,  $SD = 1.1$ ; 6<sup>th</sup> and 8<sup>th</sup>  
31 grades) who participated in the Utah Prevention Needs Assessment survey. Students self-  
32 reported their substance use behaviors and related risk and protective factors. Machine learning  
33 algorithms examined 112 individual and socioecological factors as potential classifiers of  
34 lifetime e-cigarette use outcomes. The elastic net algorithm achieved outstanding classification  
35 for lifetime exclusive ( $AUC = .926$ ) and dual use ( $AUC = .944$ ) on a validation test set. Six high  
36 value classifiers were identified that varied in importance by outcome: Lifetime alcohol or  
37 marijuana use, perception of e-cigarette availability and risk, school suspension(s), and perceived  
38 risk of smoking marijuana regularly. Specific classifiers were important for lifetime exclusive  
39 (parent attitudes regarding student vaping, best friend[s] tried alcohol or marijuana) and dual use  
40 (best friend[s] smoked cigarettes, lifetime inhalant use). Our findings provide specific targets for  
41 the adaptation of existing substance use prevention programs to address early adolescent e-  
42 cigarette use.

43 *Keywords:* e-cigarette use, vaping, machine learning, early adolescence, ecological

44

## 45 **Introduction**

46 E-cigarette use among adolescents is a national epidemic [1]. The popularity of these  
47 devices has spread faster than health researchers could amass evidence for the potential  
48 deleterious effects of e-cigarette use. As the prevalence of traditional cigarette smoking among  
49 U.S. adolescents has declined, e-cigarette use, or vaping, has become the most commonly used  
50 form of nicotine uptake among youth in the U.S. [2]. Adolescents' decisions to engage in e-  
51 cigarette use may be understood through an ecological framework that accounts for complex  
52 interactions between spheres of influence [3]. Research is underway to identify individual and  
53 socioecological risk-factors associated with e-cigarette use [4–11]. However, this literature has  
54 prominently focused on high school samples resulting in a dearth of knowledge regarding e-  
55 cigarette risk-factors during early adolescence. Identifying factors associated with the emergence  
56 of e-cigarette use during early adolescence may facilitate intervention prior to developmental  
57 periods where use escalates (i.e., middle to late adolescence [12]). These efforts may be bolstered  
58 by new methodologies that allow researchers to efficiently explore the importance of a wide  
59 range of variables in relation to e-cigarette use [13]. In this study we use machine learning  
60 algorithms to simultaneously consider a large number of individual and socioecological factors  
61 in relation to patterns of e-cigarette usage among middle school students [7].

62 The use of the e-cigarettes has been touted as a healthier alternative to tobacco cigarettes,  
63 despite their delivery of nicotine and other potentially harmful chemicals [14]. A major concern  
64 of nicotine consumption during early adolescence is the possible alteration of function in the  
65 brain's reward systems at a sensitive developmental period, in ways that can increase risk for  
66 other substance use, mood disorders, and difficulties with concentration and learning [14]. In  
67 addition to nicotine-related risks, other carcinogenic agents found in chemicals in the e-liquid as

68 well as those produced in the vaporizing product or even associated with the e-cigarette materials  
69 (i.e., nickel, chromium, cadmium; [14]). Chemicals used to flavor e-liquid have also been found  
70 to have sufficiently high toxicity to warrant medical concerns [15] or even cause death [16]. The  
71 possible harms of e-cigarettes go well beyond exposure to nicotine.

72         Researchers have documented complex relationships between individual (e.g., academic  
73 performance, substance use, perceptions of use) and socioecological (e.g., access,  
74 advertisements, peer and parental factors) influences implicated in e-cigarette use during middle  
75 to late adolescence [7–9]. Less is known about e-cigarette use risk-factors during early  
76 adolescence. The early adolescent e-cigarette literature has predominately focused on the  
77 prevalence and reasons for use, or factors associated with susceptibility rather than initiation  
78 [12,17–20]. Studies examining adolescent e-cigarette use have also had a narrow focus when  
79 considering potential individual and socioecological influences on adolescent e-cigarette use. For  
80 example, studies commonly focus on specific adolescent attitudes (e.g., perceived danger of e-  
81 cigarettes and tobacco), substance use behaviors (e.g., alcohol, tobacco, marijuana), aspects of  
82 the environment (e.g., access, advertisements), and social influences (e.g., peer and parental e-  
83 cigarette or cigarette use [7–9,12,19]). Research is needed to examine these risk-factors in  
84 conjunction with a broader array of influences traditionally associated with early substance use  
85 (e.g., anti-social behavior; parenting practices; school involvement, performance, environment;  
86 community attachment, norms, drug use, delinquency; [21–25]).

87         It is also important to consider patterns of use when identifying correlates of early  
88 adolescent vaping. Research suggests risk-factors vary between youth who have utilized e-  
89 cigarettes exclusively and those who have used them in combination with tobacco, with dual use  
90 being associated with greater behavioral problems (i.e., lifetime use;  $M = 14.6$  years,  $SD = 0.7$ )

91 and substance use (i.e., lifetime alcohol, marijuana, drug use prescription drug misuse; 9<sup>th</sup> and  
92 12<sup>th</sup> graders; [7,26]). Exclusive use may represent adolescent using e-cigarettes as a “safer”  
93 alternative to traditional tobacco cigarettes [19]. Dual use may be associated with tobacco  
94 cessation or recreational use in conjunction with other substances [19,26]. Research has yet to  
95 determine whether differences in risk-factors for exclusive or dual e-cigarette use exist during  
96 early adolescence.

97         Methodological challenges may explain the limited number of studies examining a broad  
98 array of correlates of e-cigarette use during early adolescence. For example, lifetime e-cigarette  
99 is a low base rate behavior during early adolescence relative to later developmental periods [26].  
100 Prevalence rates are even lower when researchers examine exclusive and dual e-cigarette use  
101 relative to general lifetime use [7]. Furthermore, limitations associated with traditional statistical  
102 methodologies may pose a barrier to examining the broad array of potential factors implicated in  
103 early adolescent e-cigarette use (e.g., statistical power issues, multicellularity; familywise error  
104 rate[13]). These limitations can be addressed with large datasets, however meeting statistical  
105 assumptions for multicollinearity and reducing family-wise error may limit the number of  
106 potential risk-factors that can be simultaneously considered in relation to early adolescent e-  
107 cigarette use.

108         Machine learning may facilitate the examination of factors associated with early  
109 adolescent e-cigarette use. Machine learning provides an efficient method of simultaneously  
110 examining large numbers of variables representing youth individual and socioecological factors  
111 to determine their *importance* in classifying substance use [13]. Within the context of machine  
112 learning, variable importance refers to the relative ability for variables to reduce the error in  
113 models’ predictions of group membership (e.g., exclusive e-cigarette user or non-user) compared

114 to other covariates in the model [27]. Elastic net, random forest, k-nearest neighbors, and neural  
115 networks are examples of common algorithms that are capable of provided superior accuracy in  
116 classifying lifetime substance use relative to traditional logistic regression [13,28]. Each of these  
117 algorithms approach classification with contrasting linear (i.e., elastic net) and nonlinear (i.e.,  
118 random forest, k-nearest neighbors, neural networks) methods, providing the opportunity to  
119 identify the algorithm that best performs the classification task for each outcome [27].

120 Identifying high value correlates of e-cigarette initiation during early adolescence could  
121 improving our ability to identify at-risk youth prior to developmental periods were the  
122 prevalence and frequency of vaping escalates (i.e., middle to late adolescence; [12]). While  
123 machine learning may provide an additional tool for informing substance use prevention effort  
124 [29], few studies have utilized machine learning to identify factors associated with patterns of  
125 early e-cigarette use (i.e., lifetime exclusive use, dual use with tobacco) or determined which  
126 method provides the best classification accuracy. Prior applications of machine learning have  
127 predominately focused on unstructured data (i.e., pictures, text) to classify e-cigarette use  
128 [30,31]. While a recent study trained machine learning algorithms on survey data collected on  
129 older teens ( $M_{\text{age}} = 15.36$  years old;  $SD = 1.85$ ), this research had a narrow focus on tobacco  
130 related substance use predictors that precludes a broader understanding of factors associated with  
131 early vaping initiation (i.e., LASSO and Random Forest; [32]). Thus, our aim was to (a) explore  
132 a wide array of factors using machine learning to identify important classifiers of lifetime  
133 exclusive e-cigarette and dual use within a sample of middle school students, (b) and identify the  
134 algorithm that best performs the classification task. These analyses can help quantify the relative  
135 importance of predictors and established the extent to which e-cigarette use can be classified by  
136 individual and socioecological factors.

## 137 **Materials and method**

138           The current study utilized data from the Utah Student Health and Risk Prevention  
139 (SHARP) survey project, which has been collecting and disseminating information on substance  
140 use prevalence and related behaviors since 2007 (Utah Department of Human Services [UDHS];  
141 [33]). SHARP was developed as a collaboration between multiple state agencies with the  
142 purpose of assessing risk and protective factors for problem behaviors among Utah middle and  
143 high school students. Students complete the Utah Prevention Needs Assessment (PNA) survey  
144 biannually, during the spring of odd numbered years, as a part of the SHARP survey project. The  
145 PNA survey gathers statewide data on substance use and individual/socioecological factors that  
146 influence the use of alcohol, tobacco, and other drugs. PNA surveys are implemented and used to  
147 inform statewide prevention policy and programming across the United States. Surveys are  
148 completed in schools and are self-administered using paper and pencil. The present study used  
149 data collected during the Spring of 2017 (i.e., March-June 2017) as part of a PNA survey in  
150 Utah. Parents provided written consent for their child to participate in the survey. Parents of  
151 youths that did not consent to their child participating in the PNA were not administered the  
152 survey. Student also provided verbal assent prior to participating in the PNA survey.  
153 Participation in the survey was voluntary and students could opt to participate in an alternative  
154 activity or discontinue at any time. The Utah State University Institutional Review Board  
155 approved secondary analyses of the 2017 Utah PNA survey data as non-human subjects research  
156 as participants could not be re-identified (protocol #10108). Previous research has utilized  
157 similar statewide school-based samples to identify factors associated with e-cigarette use among  
158 adolescents across the U.S. (e.g., Hawaii, Texas, Connecticut, New Jersey; [6–9,20,34]).

159           The current study focused on 14,346 middle school students (i.e., 6<sup>th</sup> and 8th grade) that  
160 participated in the 2017 Utah PNA survey. Participants were approximately 12 years old on  
161 average ( $M = 12.5$ ;  $SD = 1.1$ ), were relatively balanced on sex (girls;  $n = 7,532$ , 52.5%) and 6<sup>th</sup>  
162 grade ( $n = 7,473$ , 52.1%), and were predominantly White (9,491, 71%). Nearly a third of  
163 students attended school within Salt Lake County ( $n = 4,173$ , 29.1%) in Utah. Youths in this  
164 sample reported a 9.4% ( $n = 1,343$ ) prevalence of lifetime e-cigarette use and 5.4% ( $n = 784$ )  
165 tobacco use. Students largely reported abstaining from both tobacco and e-cigarette use (91%;  $n$   
166 = 13,003) and reported greater lifetime exclusive e-cigarette use (5.5%;  $n = 791$ ) relative to  
167 exclusive tobacco use (1.6%;  $n = 232$ ). Within the sample, 3.8% ( $n = 552$ ) of students reported  
168 dual lifetime use of tobacco and e-cigarettes. See Table 1 for sample demographic information  
169 by outcomes.

## 170 **Measures**

### 171 **Individual and socioecological variables**

172           The measures utilized in the current study have been traditionally used and reported by  
173 the SHARP survey project as individual items [33]. Variables examined in the current study have  
174 been identified as being theoretically and/or empirically important factors in the substance use  
175 literature. We decided to examine individual items to provide a nuanced understanding of e-  
176 cigarette use risk-factors [13]. Variables included a wide array of factors representing *individual*  
177 (i.e., antisocial behaviors/attitudes, rebelliousness, academic performance, perceived risk of drug  
178 use, intentions for adulthood substance use, lifetime substance use), *community* (i.e., attachment,  
179 prosocial involvement and reward, drug use consequences and antisocial behavior, perceived  
180 availability of substances), *school* (i.e., learning environment perceptions, enjoyment,  
181 commitment, benefits of learning, truancy), *home* (i.e., parenting practices, family history of



182 substance abuse, rewards for prosocial behavior, parental attitudes regarding antisocial behavior  
183 and drug use, relationship quality with parents), and *social* (i.e., best friends engaged in  
184 antisocial behavior, tried alcohol or drugs, exhibited prosocial behavior; social rewards for  
185 antisocial and prosocial behaviors) influences. See supplementary Table S1 for all items  
186 examined in the current study.

## 187 **Outcome**

188 Students reported whether they ever tried electronic cigarettes or e-cigarettes (i.e., *yes* or  
189 *no*). They also reported whether they had ever tried tobacco cigarettes, even just a puff (i.e., *yes*  
190 or *no*). Two dichotomous outcome variables were created from these items to represent lifetime  
191 exclusive e-cigarette use and dual use (i.e., tobacco and e-cigarette). The comparison group for  
192 each outcome were students who did not use either substance.

## 193 **Analytic plan**

194 In our sample, 47% ( $n = 6,744$ ) of participants were missing at least one covariate. Prior  
195 to imputation, data was randomly resampled into training (70%;  $n = 9,657$  e-cigarette;  $n = 9,490$   
196 dual) and testing sets (30%;  $n = 4,237$  e-cigarette;  $n = 4,065$  dual). We then used mode  
197 imputation, wherein missing values were replaced with the mode for each variable to address  
198 missingness independently for training and testing sets. Mode imputation is commonly utilized  
199 within the context of machine learning for classification task [28]. As algorithms can struggle to  
200 predict low base rate outcomes, a method known as down sampling was used to randomly  
201 resample and reduce the negative class (i.e., those that did not use e-cigarettes or tobacco) until it  
202 was equal to the positive class within the training set [27]. Thus, rates of lifetime use and non-  
203 use were equal for each outcome within the resampled training sets. The training sets were  $n =$   
204 1,108 for exclusive e-cigarette use and  $n = 774$  for dual use.

205 Five dissimilar machine learning algorithms—elastic net, random forest, neural networks,  
206 k-nearest neighbors, and logistic regression—were then fitted to the training set to create  
207 classification models for each outcome [35]. Each classification algorithm drew information  
208 from 112 variables representing student individual and socioecological factors. 5-fold cross-  
209 validation was used to identify variables that improved classification accuracy across random  
210 subsets of data within the training set [28]. Model performance was assessed on a test set using  
211 the Area Under (AUC) of the Receiving Operator Characteristic (ROC) curve, which represents  
212 the ability of a model to classify outcomes across all possible cut points [27]. The top performing  
213 classification algorithm on the test set was selected for each outcome (i.e., AUC; sensitivity,  
214 specificity; [27]). Variable importance figures reflect results from the best performing algorithms  
215 for each outcome. High value classifiers were then identified through visual inspection of the  
216 relative importance figures. Variables that demonstrate large increase in relative importance over  
217 subsequent covariates were said to be high value classifiers [28]. High value classifiers were  
218 examined using a crosstabulation visualization to determine the nature of the relationship  
219 between each variable and the corresponding outcome [13].

## 220 **Results**

221 Patterns of lifetime e-cigarette use differed by demographic variables within the current  
222 sample. Chi-square test of independence suggest e-cigarette usage was significantly ( $p < .001$ )  
223 associated with student gender, grade, and race/ethnicity. Boys, 8<sup>th</sup> graders, Latinxs, Native  
224 Americans, and mixed-race students reported the greatest proportion of use across outcomes.  
225 Exclusive and dual use were generally associated with a greater proportion of lifetime use across  
226 substances relative to no-users. See Table 1 for demographic variables by outcomes.

### 227 **Exclusive use**

228           Algorithmic performance on the exclusive e-cigarette use classification task ranged from  
229 good to outstanding (AUC = .787 - .926) on the test set. See supplemental Fig S1 for ROCs for  
230 classification algorithms. Elastic net was the best performing algorithm in classifying exclusive  
231 e-cigarette use (AUC = .926, sensitivity = .857, specificity = .848). In contrast, logistic  
232 regression was the worst performing algorithm in classifying lifetime e-cigarette use (AUC =  
233 .787, sensitivity = .768, specificity = .806). Elastic net identified perceived availability of e-  
234 cigarettes, lifetime alcohol use, parents' attitudes regarding their use of vape products, school  
235 suspension, perceived risk of e-cigarette use, lifetime marijuana use, best friend(s) tried alcohol,  
236 best friend(s) used marijuana, and perceived risk of smoking marijuana regularly as the best  
237 discriminators between lifetime exclusive e-cigarette users and non-users. See Fig 1 for variable  
238 importance. Visual inspection of cross-tabulation mosaics suggests that perceived availability of  
239 e-cigarettes (i.e., *sort of hard*, *very easy*, *sort of easy*), lifetime substance use (i.e., alcohol,  
240 marijuana), school suspensions (i.e., 1 or more), lower levels of perceived risk associated with e-  
241 cigarette use (i.e., *none to moderate*), best friend(s) tried alcohol or used marijuana (i.e., 1 or  
242 more), and less perceived risk associated with smoking marijuana regularly were all associated  
243 with a greater proportion of lifetime e-cigarette use. Students who reported that their parents  
244 would view their use of vape products as "*very wrong*" had the lowest proportion of use relative  
245 to other levels of approval (i.e., *wrong to not wrong at all*). See supplemental Fig S3-11 for  
246 cross-tabulation visualizations.

## 247 **Dual use**

248           Algorithmic performance on the dual tobacco and e-cigarette use classification task also  
249 ranged from excellent to outstanding (AUC = .725 - .944) on the test set. See supplemental Fig  
250 S2 for ROCs for classification algorithms. Elastic net and random forest had the same AUC

251 score (.944). However, elastic net (sensitivity = .824, specificity = .947) outperformed random  
252 forest (sensitivity = .818, specificity = .939) based on sensitivity and specificity. Logistic  
253 regression was the worst performing algorithm in classifying lifetime dual use (AUC = .725,  
254 sensitivity = .630, specificity = .779). Elastic net identified lifetime alcohol use, lifetime  
255 marijuana use, perceived availability of e-cigarettes, best friend(s) cigarette use, perceived risk of  
256 e-cigarette use, lifetime inhalants use, school suspension, and perceived risk of smoking  
257 marijuana regularly as the best discriminators between lifetime dual users and non-users. See Fig  
258 2 for variable importance. Visual inspection of cross-tabulation mosaics suggests that lifetime  
259 substance use (i.e., alcohol, marijuana, inhalants), higher levels of perceived availability of e-  
260 cigarettes (i.e., *very easy, sort of easy*), best friend(s) that have smoked cigarettes (i.e., 1 or  
261 more), school suspensions (i.e., 1 or more), lower levels of perceived risk associated with e-  
262 cigarette use and using marijuana regularly (i.e., *none to moderate*) were all associated with a  
263 greater proportion of lifetime dual use. See supplemental Fig S12-19 for cross-tabulation  
264 visualizations.

## 265 **Discussion**

266 The current study expands the literature through the simultaneous exploration of  
267 established correlates of e-cigarette initiation and traditional factors associated with substance  
268 use in relation to early adolescent vaping. Algorithms utilizing information regarding student  
269 individual characteristics and socioecological context demonstrated high levels of classification  
270 accuracy for both lifetime exclusive and dual e-cigarette use. Elastic net generally outperformed  
271 other algorithm in classification accuracy. While the order of importance of classifiers differed  
272 by outcome, elastic net consistently identified six high value classifiers across usage groups:  
273 lifetime alcohol or marijuana use, perception of e-cigarette availability and risk, school

274 suspension(s), and perceived risk of smoking marijuana regularly. Several high value classifiers  
275 differed between youth who reported lifetime exclusive (i.e., parent's attitudes regarding their  
276 use of vaping products, best friend[s] tried alcohol, best friend[s] used marijuana) and dual e-  
277 cigarette use (i.e., best friend[s] smoked cigarettes, lifetime inhalants use). These findings  
278 highlight important commonalities and difference in risk profiles between lifetime exclusive and  
279 dual e-cigarette users.

280         Research using high school samples have documented higher rates of life substance use  
281 among dual versus exclusive e-cigarette users [7]. Within our sample, rates of substance use  
282 were generally higher among exclusive and dual users relative to those who abstained from both.  
283 Consistent with prior research, the greatest portions of substance use were found among youth  
284 who had reported dual use [26]. However, only lifetime alcohol and tobacco were found to be  
285 important classifiers of both e-cigarette use outcomes among middle school students, which is  
286 consistent with prior findings in high school samples [7]. Our findings also extend prior work  
287 through the identification of inhalant use as a novel risk factor specifically related to lifetime  
288 dual use during early adolescence. It is possible that dual users may access a wide variety of  
289 substances recreationally and may utilize inhalants as they are easy to access within the home  
290 [26,36]. Further research is needed to understand the relationship between lifetime inhalant and  
291 dual use.

292         Our findings confirm that availability of e-cigarettes is an important influence on early  
293 adolescent vaping [19], despite a Utah state law that restrict the sale of these products to  
294 individuals under the age of 20. Accessibility was especially important to exclusive e-cigarette  
295 use, which is concerning as this may translate to future traditional cigarette use among  
296 adolescents who may have otherwise abstained from tobacco use [8]. Consistent with prior

297 research, students who reported lower perceived danger of using e-cigarettes reported a greater  
298 proportion of lifetime use [12]. Our findings suggest a need to consider perceptions regarding the  
299 danger of other substances, such as marijuana, when assessing risk for early adolescent e-  
300 cigarette use. While recent findings have highlighted the importance of school-based factors in  
301 assessing risk for e-cigarette using among high school samples (i.e., truancy and poor academic  
302 performance; [10]), our findings suggest that student suspensions were the most relevant aspect  
303 of school in relation to early adolescent e-cigarette use. It is possible school suspensions may be  
304 associated with an increased risk for e-cigarette use as a potential proxy for rule breaking  
305 behaviors or through greater unsupervised time outside of school [7,37]. Further research is  
306 needed to elucidate the relationship between school suspensions and early adolescent e-cigarettes  
307 use.

308 It is important to mention that factors traditionally associated with substance use such as  
309 adolescent and peer delinquency, community substance use norms, and school involvement were  
310 not important predictors of lifetime exclusive and dual e-cigarette use within the current sample.  
311 These findings may signal potential differences between factors underlying e-cigarette and other  
312 forms of substance use. Additionally, factors identified by prior research as relevant predictors of  
313 vaping (e.g., parenting practices, perceived risk of smoking tobacco) were not relevant correlates  
314 of patterns of e-cigarette use within the current sample [7,12]. It is possible that when competing  
315 against other variables within a machine learning approach, these important predictors are truly  
316 of lesser importance relative to high value classifiers identified in the current study.

## 317 **Implications**

318 Our findings support addressing early adolescent vaping through prevention programs  
319 aiming to address substance use prevention broadly (i.e., alcohol, tobacco, marijuana, inhalants,

320 e-cigarette). There is ample evidence for the efficacy of prevention programs whether focused on  
321 a single substance or multiple ones [38]. Our findings provide specific structural and behavioral  
322 targets that may inform the adaptation of programs seeking to prevent different patterns of early  
323 adolescent e-cigarette use. Substance use prevention programs may benefit from adding  
324 components that equip parents to effectively communicate disapproval regarding their child's  
325 use of e-cigarette and teach youth skills to resist peer substance use influences [39,40]. Targeted  
326 prevention programs are already supported by research documenting perception of risk  
327 associated with vaping as a consistent predictor of e-cigarette use [7,8] and our findings highlight  
328 the importance of also considering perceptions of danger regarding marijuana use during early  
329 adolescence. Specific to programming, youth's perceptions of risk do not align with research  
330 evidence providing an important point of content for preventive interventions [41]. Altering  
331 youth's perception of e-cigarette accessibility, however, may require intervention at a broader  
332 social level (e.g., public media campaigns). Alternatively, decreasing accessibility to e-cigarettes  
333 may be achieved by actions external to youths such as strong enforcement of laws regarding  
334 possession and/or consumption for underage users and/or those selling e-cigarette products to  
335 them, or by way of increasing prices for goods associated with e-cigarette use.

336 Machine learning appears to be a promising screening tool for the identification of risk  
337 factors that can accelerate the development of the e-cigarette knowledge base needed to curb the  
338 rapid spread of vaping among adolescents nationally. Algorithms were able to efficiently explore  
339 a wide range of factors in association with early adolescent e-cigarette use, which confirmed  
340 findings from later developmental stages and identified several novel risk factors (i.e., inhalant  
341 use, perceived risk of marijuana, school suspensions). An important consideration in this  
342 research is that the tools utilized to identify e-cigarette use classifiers are publicly available. R

343 offers open access statistical packages for machine learning. Additionally, there are substantial  
344 training materials available for free online. These tools can provide an accessible and replicable  
345 method of generating and disseminating scientific knowledge regarding e-cigarette use classifiers  
346 nationally. We encourage researchers to apply machine learning algorithms to their data to draw  
347 new insight regarding factors contributing to a variety of e-cigarette use outcomes among  
348 adolescents. Examining cross-sectional markers of e-cigarette use could also identify important  
349 variables that can be examined longitudinally in prospective research. Machine learning  
350 algorithms have many exciting applications when applied to longitudinal data, including  
351 identifying context specific predictors of service use, specific targets for substance use  
352 prevention programs, and ensure that important factors are not excluded from causal models  
353 examining mechanism underpinning early vaping initiation.

## 354 **Limitations**

355 Results from algorithms used in the current study do not necessarily imply causal  
356 mechanisms explaining patterns of lifetime e-cigarette use but rather identify factors that are  
357 strong correlates of group membership (i.e., use or non-use). Longitudinal research is needed to  
358 establish causal links. The current study examined lifetime substance use that may range from  
359 experimentation to habitual use. Future research may consider using machine learning as a  
360 method of identifying youth at-risk for habitual e-cigarette use. Although a large number of the  
361 Utah adolescent population was captured, the PNA survey does not include students in private  
362 schools, correctional facilities, or treatment centers. Additionally, students who were not in  
363 attendance, declined participation, or did not return parental consent forms are not represented.  
364 Furthermore, findings may not generalize to students in other states. Further research is needed  
365 to replicate our findings in different contexts and developmental periods.



## 366 **Conclusions**

367           The current study utilized a machine learning approach to efficiently explore and identify  
368 high value correlates of early adolescent lifetime e-cigarette use. This approach identified several  
369 shared risk factors for exclusive and dual e-cigarette use such as lifetime use of specific  
370 substances (i.e., alcohol, marijuana), perception of e-cigarette availability and risk, school  
371 suspension(s), and perceived risk of smoking marijuana regularly. Several differences were also  
372 identified between youth who reported lifetime exclusive (i.e., parent's attitudes regarding their  
373 use of vaping products, best friend[s] tried alcohol or used marijuana) and dual use (i.e., best  
374 friend[s] smoked cigarettes, lifetime inhalants use) relative to non-users. This information  
375 provides a first step towards identifying youth at-risk for e-cigarette use during early  
376 adolescence. Further research is needed to examine high value classifiers identified by the  
377 current study using explanatory models and longitudinal data to understand mechanism  
378 underlying their importance in accounting for differences in risk profiles between e-cigarette  
379 usage groups during early adolescence.

## 380 **Reporting**

### 381 **Funding**

382           This research was funded by a collaboration between multiple state agencies in Utah (i.e.,  
383 Department of Health, Department of Human Services, and the State Board of Education).

### 384 **Disclosure Statement**

385           The authors have no conflicts of interest associated with the publication of this  
386 manuscript.

### 387 **Acknowledgements**

388            This research was funded by a collaboration between multiple state agencies in Utah (i.e.,  
389 Department of Health, Department of Human Services, and the State Board of Education).

### 390 **Data Availability**

391            Data used in the current study can be requested from the Utah Department of Human and  
392 Health Services.

### 393 **Data Deposition**

394            We are not authorized to share the Utah Prevention Needs Assessment data used in the  
395 current study as it is owned and managed by the Utah Department of Human and Health Services.

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

## 411 **References**

- 412 1. Office of the Surgeon General. Surgeon General’s Advisory on E-cigarette Use Among  
413 Youth. 2019. Available: [https://www.cdc.gov/tobacco/basic\\_information/e-](https://www.cdc.gov/tobacco/basic_information/e-cigarettes/surgeon-general-advisory/index.html)  
414 [cigarettes/surgeon-general-advisory/index.html](https://www.cdc.gov/tobacco/basic_information/e-cigarettes/surgeon-general-advisory/index.html)
- 415 2. Centers for Disease Control and Prevention. Tobacco product use among middle and high  
416 school students—United States, 2011–2017. *Morb Mortal Wkly Rep.* 2018;63: 629–633.  
417 doi:[http://dx.doi.org/ 10.15585/mmwr.mm6722a3](http://dx.doi.org/10.15585/mmwr.mm6722a3)
- 418 3. Bronfenbrenner U. Toward an experimental ecology of human development. *Am Psychol.*  
419 1977;32: 513–531. doi:<http://doi.org/10.1037/0003-066X.32.7.513>
- 420 4. Barrington-Trimis JL, Berhane K, Unger JB, Cruz TB, Urman R, Chou P, et al. The E-  
421 cigarette Social Environment, E-cigarette Use, and Susceptibility to Cigarette Smoking.  
422 2016 [cited 26 Jun 2019]. doi:10.1016/j.jadohealth.2016.03.019
- 423 5. Bold KW, Morean ME, Kong G, Simon P, Camenga DR, Cavallo DA, et al. Early age of  
424 e-cigarette use onset mediates the association between impulsivity and e-cigarette use  
425 frequency in youth. *Drug Alcohol Depend.* 2017;181: 146–151.  
426 doi:10.1016/j.drugalcdep.2017.09.025
- 427 6. Giovenco DP, Casseus M, Duncan DT, Coups EJ, Lewis MJ, Delnevo CD. Association  
428 Between Electronic Cigarette Marketing Near Schools and E-cigarette Use Among Youth.  
429 *J Adolesc Heal.* 2016;59: 627–634. doi:10.1016/j.jadohealth.2016.08.007
- 430 7. Wills TA, Knight R, Williams RJ, Pagano I, Sargent JD. Risk factors for exclusive e-  
431 cigarette use and dual e-cigarette use and tobacco use in adolescents. *Pediatrics.* 2015;135:  
432 e43–e51. doi:10.1542/peds.2014-0760
- 433 8. Wills TA, Sargent JD, Knight R, Pagano I, Gibbons FX. E-cigarette Use and Willingness

- 434 to Smoke in a Sample of Adolescent Nonsmokers. 2016;2: 1–16.
- 435 doi:10.14440/jbm.2015.54.A
- 436 9. Simon P, Camenga DR, Morean ME, Kong G, Bold KW, Cavallo DA, et al.
- 437 Socioeconomic status and adolescent e-cigarette use: The mediating role of e-cigarette
- 438 advertisement exposure. *Prev Med (Baltim)*. 2018;112: 193–198.
- 439 doi:10.1016/j.ypmed.2018.04.019
- 440 10. McCabe SE, West BT, Veliz P, Boyd CJ. E-cigarette Use, Cigarette Smoking, Dual Use,
- 441 and Problem Behaviors Among U.S. Adolescents: Results From a National Survey. *J*
- 442 *Adolesc Heal*. 2017;61: 155–162. doi:10.1016/j.jadohealth.2017.02.004
- 443 11. Rocheleau GC, Vito AG, Intravia J. Peers, Perceptions, and E-Cigarettes: A Social
- 444 Learning Approach to Explaining E-Cigarette Use Among Youth. *J Drug Issues*. 2020;50:
- 445 472–489. doi:10.1177/0022042620921351
- 446 12. Kwon E, Seo DC, Lin HC, Chen Z. Predictors of youth e-cigarette use susceptibility in a
- 447 U.S. nationally representative sample. *Addict Behav*. 2018;82: 79–85.
- 448 doi:10.1016/j.addbeh.2018.02.026
- 449 13. Vázquez AL, Domenech Rodríguez MM, Barrett TS, Schwartz S, Amador Buenabad NG,
- 450 Bustos Gamiño MN, et al. Innovative Identification of Substance Use Predictors: Machine
- 451 Learning in a National Sample of Mexican Children. *Prev Sci*. 2020;21: 171–181.
- 452 doi:10.1007/s11121-020-01089-4
- 453 14. National Institute on Drug Abuse. What are electronic cigarettes? 2018.
- 454 15. Tierney PA, Karpinski CD, Brown JE, Luo W, Pankow JF. Flavour chemicals in
- 455 electronic cigarette fluids. *Tob Control*. 2016;25: e10–e15. doi:10.1136/tobaccocontrol-
- 456 2014-052175

- 457 16. Centers for Disease Control and Prevention. Outbreak of lung injury associated with the  
458 use of e-cigarette , or vaping, products. 2020.
- 459 17. Carey FR, Rogers SM, Cohn EA, Harrell □ MB, Wilkinson A V, Perry CL.  
460 Understanding susceptibility to e-cigarettes: A comprehensive model of risk factors that  
461 influence the transition from non-susceptible to susceptible among e-cigarette naïve  
462 adolescents. *Addict Behav.* 2019;91: 68–74. doi:10.1016/j.addbeh.2018.09.002
- 463 18. Cullen KA, Gentzke AS, Sawdey MD, Chang JT, Anic GM, Wang TW, et al. e-Cigarette  
464 Use among Youth in the United States, 2019. *JAMA - J Am Med Assoc.* 2019;322: 2095–  
465 2103. doi:10.1001/jama.2019.18387
- 466 19. Kong G, Morean ME, Cavallo DA, Camenga DR, Krishnan-Sarin S. Reasons for  
467 electronic cigarette experimentation and discontinuation among adolescents and young  
468 adults. *Nicotine Tob Res.* 2015;17: 847–854. doi:10.1093/ntr/ntu257
- 469 20. Krishnan-Sarin S, Morean ME, Camenga DR, Cavallo DA, Kong G. E-cigarette use  
470 among high school and middle school adolescents in Connecticut. *Nicotine Tob Res.*  
471 2015;17: 810–818. doi:10.1093/ntr/ntu243
- 472 21. Eitle DJ, Eitle TM. School and county characteristics as predictors of school rates of drug,  
473 alcohol, and tobacco offenses. *J Health Soc Behav.* 2004;45: 408–421.  
474 doi:10.1177/002214650404500404
- 475 22. Evans WD, Powers A, Hersey J, Renaud J. The influence of social environment and social  
476 image on adolescent smoking. *Heal Psychol.* 2006;25: 26–33.  
477 doi:<https://doi.org/10.1037/0278-6133.25.1.26>
- 478 23. Mason MJ, Zaharakis NM, Rusby JC, Westling E, Light JM, Mennis J, et al. A  
479 longitudinal study predicting adolescent tobacco, alcohol, and cannabis use by behavioral

- 480 characteristics of close friends. *Psychol Addict Behav.* 2017;31: 712–720.
- 481 doi:<https://doi.org/10.1037/adb0000299>
- 482 24. Novak SP, Clayton RR. The influence of school environment and self-regulation on  
483 transitions between stages of cigarette smoking: A multilevel analysis. *Heal Psychol.*  
484 2001;20: 196–207. doi:10.1037/0278-6133.20.3.196
- 485 25. Voelkl KE, Frone MR. Predictors of substance use at school among high school students.  
486 *J Educ Psychol.* 2000;92: 583–592. doi:10.1037/0022-0663.92.3.583
- 487 26. Curran KA, Burk T, Pitt PD, Middleman AB. Trends and Substance Use Associations  
488 With E-Cigarette Use in US Adolescents. *Clin Pediatr (Phila).* 2018;57: 1191–1198.  
489 doi:10.1177/0009922818769405
- 490 27. Kuhn M, Johnson K. *Applied Predictive Modeling with Applications in R.* 2013.  
491 Available:  
492 [http://appliedpredictivemodeling.com/s/Applied\\_Predictive\\_Modeling\\_in\\_R.pdf](http://appliedpredictivemodeling.com/s/Applied_Predictive_Modeling_in_R.pdf)
- 493 28. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning.* Second. New  
494 York: Springer; 2013. doi:10.1007/978-1-4419-9863-7\_941
- 495 29. Barenholtz E, Fitzgerald ND, Hahn WE. Machine-learning approaches to substance-abuse  
496 research: emerging trends and their implications. *Curr Opin Psychiatry.* 2020;33: 334–  
497 342. doi:10.1097/YCO.0000000000000611
- 498 30. Ketonen V, Malik A. Characterizing vaping posts on instagram by using unsupervised  
499 machine learning. *Int J Med Inform.* 2020;141: 104223.  
500 doi:10.1016/j.ijmedinf.2020.104223
- 501 31. Malik A, Khan MI, Karbasian H, Nieminen M, Ammad-Ud-Din M, Khan SA. Modeling  
502 Public Sentiments About JUUL Flavors on Twitter Through Machine Learning. *Nicotine*

- 503 Tob Res. 2021; 1–11. doi:10.1093/ntr/ntab098
- 504 32. Choi J, Jung H-T, Ferrell A, Woo S, Haddad L. Machine Learning-Based Nicotine  
505 Addiction Prediction Models for Youth E-Cigarette and Waterpipe (Hookah) Users. *J Clin*  
506 *Med.* 2021;10: 972. doi:10.3390/jcm10050972
- 507 33. Utah Department of Human Services. Student Health and Risk Prevention (SHARP)  
508 survey reports. 2019.
- 509 34. Temple JR, Shorey RC, Lu Y, Torres E, Stuart GL, Le VD. E-cigarette use of young  
510 adults motivations and associations with combustible cigarette alcohol, marijuana, and  
511 other illicit drugs. *Am J Addict.* 2017;26: 343–348. doi:10.1111/ajad.12530
- 512 35. James G, Witten D, Hastie T, Tibshirani R. *An Introduction to Statistical Learning.*  
513 *Current Medicinal Chemistry.* Springer; 2013.
- 514 36. Medina-Mora M´a E, Real T. Epidemiology of inhalant use. *Curr Opin Psychiatry.*  
515 2008;21: 247–251. doi:<https://doi.org/10.1097/yco.0b013e3282fc9875>
- 516 37. Lee KTH, Vandell DL. Out-of-School Time and Adolescent Substance Use. *J Adolesc*  
517 *Heal.* 2015;57: 523–529. doi:10.1016/j.jadohealth.2015.07.003
- 518 38. Das JK, Salam RA, Arshad A, Finkelstein Y, Bhutta ZA. Interventions for Adolescent  
519 Substance Abuse: An Overview of Systematic Reviews. *J Adolesc Heal.* 2016;59: S61–  
520 S75. doi:10.1016/j.jadohealth.2016.06.021
- 521 39. Marsiglia FF, Kulis S, Yabiku ST, Nieri TA, Coleman E. When to Intervene: Elementary  
522 School, Middle School or Both? Effects of keepin’ It REAL on Substance Use  
523 Trajectories of Mexican Heritage Youth. *Prev Sci.* 2011;12: 48–62. doi:10.1007/s11121-  
524 010-0189-y
- 525 40. Patterson G. The next generation of PMTO models. *Behav Ther.* 2005;28: 25–32.

526 41. Cheeta S, Halil A, Kenny M, Sheehan E, Zamyadi R, Williams AL, et al. Does perception  
527 of drug-related harm change with age? A cross-sectional online survey of young and older  
528 people. *BMJ Open*. 2018;8. doi:10.1136/bmjopen-2017-021109

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

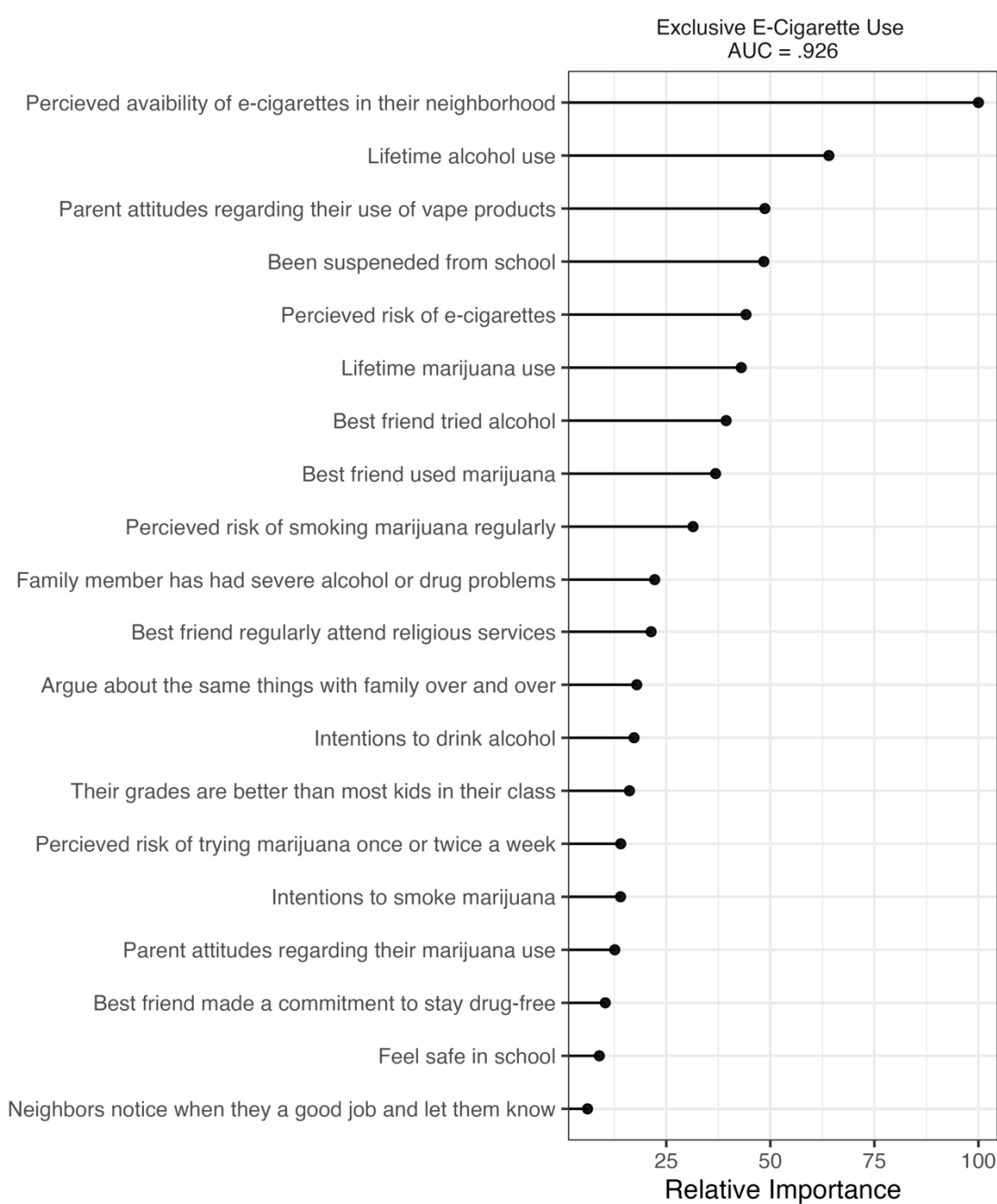
548



	<b>Total</b>	<b>Neither</b>	<b>Exclusive Use</b>	<b>Dual use</b>	<b><i>p</i> value<sup>a</sup></b>
<b>Sample</b>	<i>n</i> = 14346	<i>n</i> = 13003	<i>n</i> = 791	<i>n</i> = 552	
<b>Age <i>M</i>(<i>SD</i>)</b>	12.5 (1.1)	12.5 (1.1)	13.1 (1.0)	13.1 (1.0)	
<b>Sex</b>					< .001
<b>Boy</b>	6766 (47.2%)	6066 (89.7%)	427 (6.3%)	273 (4%)	
<b>Girl</b>	7532 (52.5%)	6894 (91.5%)	361 (4.8%)	277 (3.7%)	
<b>Grade</b>					< .001
<b>6th</b>	7473 (52.1%)	7110 (95.1%)	217 (2.9%)	146 (2%)	
<b>8th</b>	6873 (47.9%)	5893 (85.7%)	574 (8.4%)	406 (5.9%)	
<b>Race/ethnicity</b>					< .001
<b>White</b>	10191 (71%)	9491 (93.1%)	415 (4.1%)	285 (2.8%)	
<b>Native American</b>	289 (2%)	247 (85.5%)	20 (6.9%)	22 (7.6%)	
<b>Asian</b>	238 (1.7%)	221 (92.9%)	12 (5%)	5 (2.1%)	
<b>Black</b>	210 (1.5%)	186 (88.6%)	11 (5.2%)	13 (6.2%)	
<b>Latinx</b>	1888 (13.2%)	1554 (82.3%)	210 (11.1%)	124 (6.6%)	
<b>Pacific Islander</b>	210 (1.5%)	191 (91%)	13 (6.2%)	6 (2.9%)	
<b>Mixed Race</b>	1320 (9.2%)	1113 (84.3%)	110 (8.3%)	97 (7.3%)	
<b>Substance use<sup>b</sup></b>					
<b>Alcohol</b>					< .001
<b>Yes</b>	1329 (9.3%)	600 (45.1%)	353 (26.6%)	376 (28.3%)	
<b>No</b>	12982 (90.5%)	12371 (95.3%)	436 (3.4%)	175 (1.3%)	
<b>Marijuana</b>					< .001
<b>Yes</b>	602 (4.2%)	124 (20.6%)	201 (33.4%)	277 (46%)	
<b>No</b>	13671 (95.3%)	12812 (93.7%)	585 (4.3%)	274 (2%)	
<b>Inhalants</b>					< .001
<b>Yes</b>	669 (4.7%)	414 (61.9%)	105 (15.7%)	150 (22.4%)	
<b>No</b>	13579 (94.7%)	12499 (92%)	679 (5%)	401 (3%)	
<b>Prescription drugs</b>					< .001
<b>Yes</b>	546 (3.8%)	332 (60.8%)	81 (14.8%)	133 (24.4%)	
<b>No</b>	13562 (94.5%)	12468 (91.9%)	690 (5.1%)	404 (3%)	
<b>Hallucinogens</b>					< .001
<b>Yes</b>	99 (0.7%)	22 (22.2%)	25 (25.3%)	52 (52.5%)	
<b>No</b>	14128 (98.5%)	12871 (91.1%)	760 (5.4%)	497 (3.5%)	
<b>Synthetic marijuana</b>					< .001
<b>Yes</b>	86 (0.6%)	12 (14%)	21 (24.4%)	53 (61.6%)	
<b>No</b>	14199 (99%)	12940 (91.1%)	766 (5.4%)	493 (3.5%)	

549 Variable frequency is displayed by column for the total and row for usage groups. <sup>a</sup> Chi-square  
 550 test of independence. <sup>b</sup> Lifetime use.

551 **Fig 1. Top 20 Variables with the Highest Relative Importance in Classifying Lifetime E-**  
552 **Cigarette Use.** Results represent validation on a separate test dataset.



553  
554  
555  
556  
557

558 **Fig 2. Top 20 Variables with the Highest Relative Importance in Classifying Lifetime Dual**  
559 **Use.** Results represent validation on a separate test dataset.

