

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Conceptualizing bias in EHR data: A case study in performance disparities by demographic subgroups for a pediatric obesity incidence classifier

Short Title: EHR Data Bias in Pediatric Obesity Classification

Elizabeth A. Campbell, PhD*^{1,2,3}, Saurav Bose, MS^{1,4}, Aaron J. Masino, PhD^{1,5}

1. Department of Biomedical and Health Informatics, Children’s Hospital of Philadelphia, Philadelphia, PA, United States of America

2. Department of Biomedical Informatics, Columbia University Medical Center, New York, NY, United States of America

3. Department of Information Science, College of Computing & Informatics, Drexel University, Philadelphia, PA, United States of America

4. Foursquare Labs Inc., New York, NY, United States of America

5. University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, United States of America

*Corresponding Author:

Email: ec3696@cumc.columbia.edu (EC)

17
18

19 **Abstract**

20
21 Electronic Health Records (EHRs) are increasingly used to develop machine learning models in
22 predictive medicine. There has been limited research on utilizing machine learning methods to
23 predict childhood obesity and related disparities in classifier performance among vulnerable
24 patient subpopulations. In this work, classification models are developed to recognize pediatric
25 obesity using temporal condition patterns obtained from patient EHR data. We trained four
26 machine learning algorithms (Logistic Regression, Random Forest, XGBoost, and Neural
27 Networks) to classify cases and controls as obesity positive or negative, and optimized
28 hyperparameter settings through a bootstrapping methodology. To assess the classifiers for bias,
29 we studied model performance by population subgroups then used permutation analysis to
30 identify the most predictive features for each model and the demographic characteristics of
31 patients with these features. Mean AUC-ROC values were consistent across classifiers, ranging
32 from 0.72-0.80. Some evidence of bias was identified, although this was through the models
33 performing better for minority subgroups (African Americans and patients enrolled in Medicaid).
34 Permutation analysis revealed that patients from vulnerable population subgroups were over-
35 represented among patients with the most predictive diagnostic patterns. We hypothesize that our
36 models performed better on under-represented groups because the features more strongly
37 associated with obesity were more commonly observed among minority patients. These findings
38 highlight the complex ways that bias may arise in machine learning models and can be
39 incorporated into future research to develop a thorough analytical approach to identify and
40 mitigate bias that may arise from features and within EHR datasets when developing more
41 equitable models.

42

43 **Author Summary**

44

45 Childhood obesity is a pressing health issue. Machine learning methods are useful tools to study
46 and predict the condition. Electronic Health Record (EHR) data may be used in clinical research
47 to develop solutions and improve outcomes for pressing health issues such as pediatric obesity.
48 However, EHR data may contain biases that impact how machine learning models perform for
49 marginalized patient subgroups. In this paper, we present a comprehensive framework of how
50 bias may be present within EHR data and external sources of bias in the model development
51 process. Our pediatric obesity case study describes a detailed exploration of a real-world
52 machine learning model to contextualize how concepts related to EHR data and machine
53 learning model bias occur in an applied setting. We describe how we evaluated our models for
54 bias, and considered how these results are representative of health disparity issues related to
55 pediatric obesity. Our paper adds to the limited body of literature on the use of machine learning
56 methods to study pediatric obesity and investigates the potential pitfalls in using a machine
57 learning approach when studying social significant health issues.

58

59

60

61 **Introduction**

62
63
64 Throughout most disciplines, massive amounts of data are being digitally generated, collected,
65 and stored at a rapidly expanding pace. Additionally, advances in computational methods enable
66 extraction of information from such datasets that produce useful insights and knowledge. (1)

67
68 In healthcare, there is increasing use of large and variable data sources that include medical
69 imaging, wearable devices, genome sequencing, and payer records among others; electronic
70 health records (EHRs) are one particularly robust source of healthcare data. This data is available
71 in an exceptionally high volume, spans the healthcare sector's digital ethos, and is extremely
72 variable in its structure, semantics, and information content. (2,3)

73
74 Advanced data mining and analytical methods are necessary to obtain, transform, and analyze
75 EHR data for secondary uses such as clinical and health policy research; machine learning
76 methods are key to addressing challenges in secondary EHR uses. However, although EHR data
77 and the models that may be trained with this data, hold tremendous potential to transform clinical
78 care and research, caution must be exercised in how this data is utilized and interpreted
79 analytically. (4,5) Bias is an inherent property to statistical models and within data collection,
80 and can also be introduced algorithmically or found within the data used to train and test
81 machine learning models.(6)

82
83 Issues in using machine learning methods to analyze EHR data often arise when letting data
84 speak for itself. Algorithms may be subject to biases that are present in EHR datasets from
85 several sources including study population characteristics, systemic errors in how EHR data is
86 collected, missing data, misclassification, and sample size. (7) Spurious correlations and other
87 dataset deficiencies such as multicollinear, correlated predictors may lead to algorithms
88 overfitting predictions to biased data and producing unstable estimates. In turn, this affects the
89 models' performance and generalizability, potentially causing or perpetuating health system
90 disparities. Machine learning models may be subject to new biases not typically seen in more
91 traditional observational studies or statistical methods, such as adjusting away healthcare quality

92 differences between patients or misinterpreting treatment outcomes when making therapy
93 recommendations.(8,9)

94

95 **Bias Definitions**

96

97 Bias that is present in EHR data may result from numerous sources including measurement
98 errors or selection bias in populations that are represented in EHR data versus the communities
99 that they represent. (10) Biased data may reflect existing prejudices or disparities of the contexts
100 from which data are collected. For example, inadequate access to insurance or under-diagnosis of
101 certain conditions may lead to a misrepresentation of a condition's prevalence among vulnerable
102 populations. In EHR data, these pernicious biases may also manifest from inequities in usage and
103 access to care or in the care that vulnerable subgroups may receive in healthcare settings. Bias
104 may also be introduced algorithmically or in the model design processes, which makes
105 measuring bias when evaluating machine learning models an important area for promoting
106 equity. (11,12) In Figure 1, we conceptualize these EHR data bias sources and how they
107 contribute to developing biased machine learning models in clinical research. In this study, we
108 focus on contextualizing pernicious bias in a particular dataset and how such biases may be
109 characterized.

110 **Figure 1.** A Framework to Understand Bias in EHR Data and Machine Learning Models. We
111 identify three sources of bias (pernicious, measurement, and sampling biases) that may occur in
112 raw data that translates to biased datasets. Biased data along with bias introduced from design
113 processes and algorithmically may lead to differential machine learning model performance,
114 interpretation, and implementation, which in turn may perpetuate health system disparities.

115 In this work, we investigate potential biases (which we define here to mean systemic errors or
116 misrepresentations embedded within datasets) that may exist in machine learning models
117 developed from EHR data in the context of a childhood obesity incidence case study. Within
118 the United States, approximately one third of children are overweight (age- and sex-specific
119 body mass index (BMI) greater than or equal to the 85th percentile per Centers for Disease
120 Control and Prevention (CDC) growth charts) or obese (age- and sex-specific BMI greater than
121 or equal to the 95th percentile per CDC growth charts). (13,14) Obesity is linked with numerous
122 comorbidities, including an increased risk of developing asthma, diabetes, hypertension, and

123 psychological conditions during childhood and later in life. (15,16) Pediatric obesity is a socially
124 significant health issue that disproportionately impacts American Indian, African American, and
125 Latino children, compared to non-Hispanic whites. Obesity prevalence is also higher among low-
126 income, rural, or less-educated population subgroups. (13,17)

127
128 We developed a classification model to predict childhood obesity incidence using our previously
129 published temporal condition patterns surrounding pediatric obesity that were derived from EHR
130 data and patient demographic data. (18) Our study aims to address the following research
131 questions:

- 132
- 133 1. Can a machine learning classification model using temporal condition patterns and
134 demographic information from EHR data accurately predict future childhood obesity
135 incidence?
 - 136 2. How can machine learning models developed using EHR data be analyzed for bias in
137 model performance amongst population subgroups?
 - 138 3. How can biased model performance be understood in the context of the individual
139 condition that a researcher is working to address?

140 141 **Materials and Methods**

142 143 **Setting**

144
145 EHR data was obtained from the Pediatric Big Data (PBD) resource at the Children's Hospital of
146 Philadelphia (CHOP). Patients in this study were from a retrospective cohort of newly obese
147 patients and matched control patients with a healthy BMI identified in a previous study. (19) The
148 PBD resource was created for secondary research use, and contains health-related information,
149 including demographic, encounter, medication, procedure, and measurement (e.g. vital signs,
150 laboratory results) elements for a large, unselected population of children.

151
152 Ethics statement: Non-study personnel extracted all data from the EHR and removed protected
153 health information (PHI) identifiers, except for dates, prior to transfer to the study database. Date

154 information was removed from the analysis dataset used in this study. The CHOP Institutional
155 Review Board approved this study and waived the requirement for consent.

156

157 **Temporal Condition Pattern Mining Methodology**

158

159

160 In a previous study, (18) we applied a sequential pattern mining algorithm to a dataset from a
161 large retrospective cohort of newly obese pediatric patients (n = 49 694) at CHOP from 2009-
162 2017. Patients were identified using the CDC definition of childhood obesity (BMI z-score at or
163 above the 95th percentile for age and sex). (13,14) The BMI z-scores were centrally calculated in
164 this analysis. The same definition of obesity was used across study sites for the entire study
165 period. Patients had at least one obesity measurement during a CHOP primary care visit and at
166 least one visit prior to the first obesity measurement where an obese BMI was not recorded. The
167 analysis aimed to identify common temporal condition patterns derived from visits immediately
168 before (pre-index) and after (post-index) the first visit with an obese BMI (index). We found 163
169 condition patterns present in at least 1% of the obese patients, of which 80 were more
170 significantly more common than in matched controls. Campbell, et al includes a full study
171 diagram detailing the inclusion criteria implementation for obtaining the study population and
172 derivation of the temporal condition patterns. (18)

173

174 **Study Population**

175

176 To obtain the study population for the machine learning case study presented here, we began
177 with 49,694 pairs of matched cases and controls from the prior study. Patients in our final study
178 population must have had both a BMI measurement in the pre-index and index visit (for control
179 patients the index visit was the date for the visit that they were matched on with their
180 corresponding case patient). For case patients, this meant that they needed to have a non-obese
181 BMI measurement in the pre-index visit, and an obese BMI measurement in the index visit;
182 15,522 case patients met these criteria. Control patients needed a healthy BMI measurement in
183 both their pre- and index visits; 31,366 control patients met this criterion. Finally, only patients
184 from case-control pairs where both patients met the BMI inclusion criteria were kept in the study

185 population; 4,843 case-control pairs met the criteria and 44,851 did not. A total of 9,686 patients
186 met the BMI criterion for inclusion.

187
188 Patients and their corresponding matched case or control were eliminated if they did not have
189 insurance information within 2 years of the matched index visit. For controls, 45 were missing
190 insurance information from within two years or altogether; these 45 controls and their matched
191 cases were eliminated from the study population. Seven cases were missing this information and
192 were eliminated from the study population (along with their matched controls). The final study
193 population contained 4,777 matched pairs, and 9,554 total patients. The study population and
194 data acquisition process are summarized in Figure 2. Table 1 presents the demographic
195 characteristics of the total study population, as well as the case and control populations
196 respectively.

197
198 **Figure 2.** A flow chart illustrating how patients in the final study population were obtained after
199 filtering through the study's inclusion criteria. M represents the total number of patients, and N
200 represents the total number of matched case-control pairs.

201
202

203 **Feature Selection, Data Acquisition and Preprocessing**

204
205 The features selected for the machine learning case study included temporal condition patterns
206 uncovered in Campbell, et al. (18) For obesity incidence prediction, only temporal condition
207 patterns from the pre-index and index visits were considered in this model. Of the original 80
208 patterns identified in the previous study, 70 temporal condition patterns were selected for
209 inclusion in this analysis. Each temporal condition pattern is considered separately as a feature
210 for this study. Patient EHR data were analyzed for the presence of each temporal condition
211 pattern, and patients were assigned a binary value of 0 (indicating a patient did not have a record
212 of the temporal condition pattern) or a 1 (indicating that a patient did have a record of the
213 temporal condition pattern) for each variable. Diagnoses in the temporal condition patterns are
214 coded using the Expanded Diagnostic Clusters (EDCs) from the Adjusted Clinical Group (ACG)
215 System. (20,21) The temporal diagnoses that comprised the condition patterns used in this study
216 may be found in Table 1 in the Supporting Information section.

217

218 Person-level characteristics including race, sex, ethnicity, age at index visit, and insurance were
219 also extracted from EHR data within the PBD database and included in the final dataset. The
220 demographic variables considered were sex assigned at birth, race, Medicaid enrollment (a proxy
221 for socioeconomic status at the time of obesity incidence), (22,23) and age at index visit (coded
222 as a categorical variable, for patients who were 2-5 years, 6-11 years, and 12-18 years). Patients
223 were classified as Hispanic if their self-identified ethnicity was specified as Hispanic or Latino;
224 otherwise, they were categorized by the value of their self-identified race from the EHR. Patients
225 with missing race and ethnicity information were classified as unknown. Patients were classified
226 as being enrolled in Medicaid if they used multiple insurance plans and one of those was
227 Medicaid or Children’s Health Insurance Program (CHIP), Pennsylvania’s state program to
228 provide health insurance to uninsured children and teens who are ineligible or not enrolled in
229 Medicaid. (24) For patients who did not have insurance information recorded for their index
230 visit, all insurance information for their visits within a year of their index visit was obtained from
231 the PBD database and analyzed. If patients had a record of Medicaid/CHIP enrollment within a
232 year of their index visit, they were classified in the Medicaid/CHIP enrollment category.

233

234 **Machine Learning Analysis**

235

236 We trained four machine learning models (Logistic Regression, Random Forest, XGBoost, and
237 Neural Networks) to classify cases and controls as obesity positive or negative, and optimized
238 hyperparameter settings through a bootstrapping methodology. We randomly shuffled the data
239 and split it into training and validation folds in a stratified fashion relative to the 50:50 class
240 balance. We trained each model with all hyperparameter settings on the training fold and
241 evaluated its Area Under the Receiver Operating Curve (AUC-ROC) on the validation fold. We
242 repeated the process 200 times to obtain 200 validation AUC-ROCs for each hyperparameter
243 setting for each model, then selected the hyperparameter combination with the highest mean
244 validation AUC-ROC for a given model class. We implemented all algorithms using the Scikit-
245 learn library in Python 3. (25) We calculated the mean and standard deviation (SD) AUC-ROC
246 values for the total study population and demographic subgroups for each algorithm.

247

248 **Permutation Analysis**

249

250 We performed a permutation feature importance analysis on the data split with median validation
 251 AUC-ROC for the model with hyperparameters corresponding to the highest mean validation
 252 AUC-ROC. The feature importance is computed by measuring the change in the AUC-ROC on
 253 the validation set when the values in the dataset for a given feature are randomly shuffled among
 254 samples. Feature importance is reflected by a decrease in AUC-ROC as compared to when the
 255 feature is not permuted, with higher importance indicated by a larger decrease.

256
 257

258 Results

259

260 Study Population

261
 262

263 Table 1 presents the demographic characteristics of the total study population, as well as the case
 264 and control populations respectively.

265

266 **Table 1.** Demographic Characteristics of Obesity Incidence Study Case and Control Populations

267

	Total Study Population (n= 9 554)	Case Population (n= 4 777)	Control Population (n= 4 777)
	n (%)	n (%)	
<i>Sex</i>			
Male	5 294 (55.4%)	2 647 (55.4%)	2 647 (55.4%)
Female	4 260 (44.6%)	2 130 (44.6%)	2 130 (44.6%)
<i>Race/ethnicity</i>			
Non-Hispanic Asian	285 (2.9%)	114 (2.4%)	171 (3.6%)
Non-Hispanic Black/African American	2 417 (25.3%)	1 512 (31.7%)	905 (18.9%)
Non-Hispanic White	5 744 (60.1%)	2 575 (53.9%)	3 168 (66.3%)
Hispanic	334 (3.5%)	208 (4.4%)	126 (2.6%)
Non-Hispanic Multiple Race	124 (1.3%)	68 (1.4%)	56 (1.2%)
Non-Hispanic Heterogeneous Other	9 (<1%)	7 (<1%)	2 (<1%)
Unknown	641 (6.7%)	292 (6.1%)	349 (7.3%)
<i>Medicaid Enrollment</i>			

Medicaid/CHIP	3 067 (32.1%)	1 822 (38.1%)	1 245 (26.1%)
---------------	---------------	---------------	---------------

Age at index visit

	<i>Logistic Regression</i>	<i>Random Forest</i>	<i>XGBoost</i>	<i>Neural Net</i>
--	----------------------------	----------------------	----------------	-------------------

2-4 years	2 426 (25.4%)	1 213 (25.4%)	1 213 (25.4%)
5-11 years	3 791 (39.7%)	1 896 (39.7%)	1 895 (39.7%)
12-18 years	3 337 (34.9%)	1 667 (34.9%)	1 669 (34.9%)

268
 269 The study population is majority male (55.4%) and majority White (60.1%). African Americans
 270 are the second largest racial /ethnic group (25.3%). Approximately 1/3 of patients (32.1%) were
 271 enrolled in Medicaid at the time of their index visit. The case population is majority male
 272 (55.4%) and majority White (53.9%) but is comprised of a higher proportion of African
 273 Americans (31.7% vs. 25.3%) and Hispanic patients (4.4% vs. 3.5%) patients compared to the
 274 entire study population. Additionally, a greater proportion of case patients (38.1%) were enrolled
 275 in Medicaid compared to the overall study population (32.1%). The control population has a
 276 higher proportion of White patients compared to the case population (66.3% versus 53.9%) and a
 277 lower proportion of racial minorities. The control population also had a lower rate of Medicaid
 278 enrollment than the case population (26.1% versus 38.1%).

279
 280 **Machine Learning Results**

281 **Table 2.** Mean(SD) AUC-ROC for Study Population and Demographic Subgroups by
 282 Classification Algorithm

<i>Total Study Population</i>	0.78 (0.01)	0.77 (0.01)	0.78 (0.01)	0.76 (0.01)
<i>Sex</i>				
Male	0.78 (0.01)	0.77 (0.01)	0.78 (0.01)	0.76 (0.01)
Female	0.78 (0.01)	0.77 (0.01)	0.78 (0.01)	0.77 (0.01)
<i>Race</i>				
Asian	0.76 (0.05)	0.76 (0.05)	0.77 (0.05)	0.74 (0.06)
Black/African American	0.79 (0.01)	0.79 (0.01)	0.79 (0.02)	0.78 (0.02)
White	0.75 (0.04)	0.75 (0.04)	0.76 (0.01)	0.74 (0.01)
Hispanic	0.75 (0.04)	0.75 (0.04)	0.77 (0.04)	0.74 (0.04)
Multiple Race	0.73 (0.08)	0.73 (0.08)	0.76 (0.08)	0.75 (0.07)
Unknown	0.73 (0.03)	0.73 (0.03)	0.73 (0.04)	0.72 (0.04)
<i>Medicaid Enrollment</i>				
Medicaid/CHIP	0.80 (0.01)	0.79 (0.01)	0.80 (0.01)	0.79 (0.01)
Not Enrolled in Medicaid/CHIP	0.76 (0.01)	0.75 (0.01)	0.76 (0.01)	0.74 (0.01)
<i>Age at index visit</i>				
2-4 years	0.76 (0.02)	0.75 (0.01)	0.76 (0.02)	0.75 (0.01)
5-11 years	0.80 (0.01)	0.79 (0.01)	0.80 (0.01)	0.79 (0.01)
12-18 years	0.75 (0.02)	0.75 (0.02)	0.76 (0.01)	0.75 (0.02)

284
 285 Mean AUC-ROC values were consistent across algorithms, ranging from 0.72-0.80. On the full
 286 study population, Neural Net had a mean AUC-ROC value of 0.76, and mean AUC-ROC values
 287 ranged from 0.70-0.79 across demographic subgroups. On the full study population, Random
 288 Forest had a mean AUC-ROC value of 0.77, and mean AUC-ROC values ranged from 0.73-0.79
 289 across demographic subgroups. On the full study population, Logistic Regression had a mean
 290 AUC-ROC value of 0.77, and mean AUC-ROC values ranged from 0.73-0.80 across
 291 demographic subgroups. On the full study population, XGBoost had a mean AUC-ROC value of
 292 0.78, and mean AUC-ROC values ranged from 0.73-0.80 across demographic subgroups.
 293 XGBoost and Logistic regression tended to perform the best on the full study population and
 294 when evaluated by demographic subgroups. Some evidence of bias was identified, although
 295 surprisingly this was through the models performing better for minority subgroups. The highest
 296 mean AUC-ROC value (0.80) was observed among African American patients, patients enrolled
 297 in Medicaid, and patients ages 5-11 years.

298 **Permutation Analysis Findings**

299

300 A permutation analysis was undertaken to investigate why models tended to perform slightly
 301 better for under-represented groups in the study population. We hypothesized that the features
 302 that are most predictive of obesity may be more common among marginalized subpopulations.
 303 Thus, we undertook a permutation feature analysis to identify which features were most
 304 important in classifying patients as obese or not obese for each algorithm.

305
 306
 307

Table 3. Most Predictive Variables by Classifier (Value (2 sigma))

	Classifier			
Variable Predictive Rank	XGBoost	Random Forest	Neural Net	Logistic Regression
1	1-ALL04 2.04% (.37%)	1-ALL04 2.10% (.60%)	1-ALL04 1.13 % (.57%)	race_3 1.52 % (.55%)
2	2-MUS01 1.36% (.32%)	2-ALL03 1.41% (.65%)	2-MUS01 0.93% (.30%)	2-MUS01 1.26% (.28%)
3	2-ALL03 1.15% (.25%)	2-NUR19 1.31% (.62%)	1-ALL03, 1-ALL04 0.85% (.35%)	1-ALL04 1.11% (.41%)
4	2-SKN04 0.83% (.29%)	race_2 1.05% (.41%)	1-SKN02 0.79% (.34%)	2-MUS04 0.99% (.20%)
5	1-SKN02 0.79% (.33%)	1-SKN02 0.82% (.40%)	2-ALL03 0.79% (.40%)	2-NUR19 0.97% (.41%)
6	2-END05 0.75% (.18%)	2-MUS01 0.82% (.32%)	race_3 0.70% (.62%)	1-ALL03, 1-ALL04 0.87% (.20%)
7	race_2 0.72% (.26%)	2-RES01 0.77% (.37%)	2-SKN04 0.61% (.22%)	2-RES01 0.69% (.41%)
8	2-ALL04 0.72% (.41%)	1-GAS03 0.70% (.54%)	2-ALL04 0.60% (.86%)	2-ALL03 0.67% (.50%)
9	2-MUS17 0.70% (.11%)	2-ALL04 0.65% (.47%)	1-GAS03 0.60% (.22%)	2-SKN04 0.66% (.39%)
10	2-MUS04 0.68% (.31%)	2-SKN04 0.59% (.27%)	2-RES01 0.54% (.45%)	medicaid 0.55% (.29%)

308 **Table 3** The top ten most predictive sequences for each classification algorithm. The gray highlighted cells represent
 309 sequences that were most predictive across all four classifiers. The orange highlighted cells indicate race variables
 310 that were most predictive.
 311

312 Table 3 presents the top ten most predictive sequences for each classification algorithm. Four
 313 temporal condition patterns were among the top ten most predictive features across all four
 314 algorithms: 1-ALL04 (a diagnosis of asthma in the pre-index visit), 2-MUS01 (a diagnosis of
 315 Musculoskeletal signs and symptoms in the index visit), 2-ALL03 (a diagnosis of allergic rhinitis
 316 in the index visit), and 2-SKN04 (a diagnosis of acne in the index visit). A diagnosis of asthma
 317 in the pre-index visit (1-ALL04) was the most predictive feature for the XGBoost, Neural
 318 Network, and Random Forest algorithms, and was the third most predictive for Logistic
 319 Regression. To better understand the disparate machine learning model performance and
 320 permutation analysis findings, the prevalence of these four temporal condition patterns were
 321 assessed among demographic subgroups in the study population, Table 4.

322
 323 **Table 4.** Demographic Characteristics of Patient Subgroups with most predictive sequences
 324

	Total Study Population (n= 9 554)	1-ALL04 (n= 851)	2-MUS01 (n= 190)	2-ALL03 (n= 459)	2-SKN04 (n= 140)
	n (%)	n (%)	n (%)	n (%)	n (%)
<i>Sex</i>					
Male	5 294 (55.4%)	533 (62.6%)	99 (52.1%)	251 (54.7%)	61 (43.5%)
Female	4 260 (44.6%)	318 (37.4%)	91 (47.9%)	208 (45.3%)	79 (56.4%)
<i>Race</i>					
Asian	285 (2.9%)	18 (2.1%)	1 (<1%)	13 (2.8%)	5 (3.6%)
Black/AA	2 417 (25.3%)	421 (49.4%)	37 (19.5%)	212 (46.2%)	57 (40.7%)
White	5 744 (60.1%)	342 (40.1%)	132 (69.5%)	189 (41.2%)	61 (43.6%)
Hispanic	334 (3.5%)	30 (3.5%)	6 (3.2%)	15 (3.3%)	3 (2.1%)
Multiple Race	124 (1.3%)	29 (1.3%)	4 (2.1%)	3 (<1%)	1 (<1%)
Heterogeneous	9 (<1%)	0 (0%)	1 (<1%)	0 (0%)	0 (0%)
Unknown	641 (6.7%)	29 (3.4%)	9 (4.7%)	27 (5.9%)	13 (9.3%)
<i>Medicaid Enrollment</i>					
Medicaid/CHIP	3 067 (32.1%)	404 (47.5%)	55 (28.9%)	184 (40.1%)	44 (31.4%)
<i>Age at index visit</i>					
2-4 years	2 426 (25.4%)	240 (28.2%)	13 (6.8%)	78 (17.0%)	2 (1.4%)
5-11 years	3 791 (39.7%)	375 (44.1%)	67 (35.3%)	233 (50.8%)	17 (12.1%)
12-18 years	3 337 (34.9%)	240 (27.7%)	110 (57.9%)	148 (32.2%)	121 (86.4%)

325

326 Two features that were most common among patients were also a most predictive sequence for
327 all four classifiers: 1-ALL04 (asthma in the pre-index visit), of which 851 patients had a record
328 in their EHR data, and 2-ALL03 (allergic rhinitis in the index visit) of which 459 patients had a
329 record (compared to 190 patients with the 2-MUS01 sequence and 140 patients with the 2-
330 SKN04 diagnosis). African American patients and patients enrolled in Medicaid are over-
331 represented among patients who have these diagnoses. Almost half of patients (49.4%) with a 1-
332 ALL04 diagnosis were African American, even though African American patients make up only
333 25.3% of the study population, and 47.5% of patients with this diagnosis were enrolled in
334 Medicaid compared to only 32.1% of the total study population. Similarly, 46.2% of patients
335 with the 2-ALL03 diagnosis were African American and 40.1% were enrolled in Medicaid.

336

337 **Discussion**

338

339 In this study, four supervised machine learning algorithms were trained to identify pediatric
340 patients as obese or not obese using demographic variables and temporal condition patterns
341 previously found to be associated with obesity incidence. Model performance was evaluated for
342 the total population and by demographic subgroups using mean AUC-ROC values. Mean AUC-
343 ROC values were consistent across algorithms, ranging from 0.72-0.80. XGBoost and Logistic
344 regression tended to perform the best on the full study population and when evaluated by
345 demographic subgroups. Algorithms tended to perform relatively consistently compared to one
346 another and when each classifier's performance was analyzed by demographic subgroups. This
347 result is consistent with existing work demonstrating that complex machine learning models
348 often perform no better than logistic regression when using EHR input data. (26) Our conjecture
349 is that this is likely true when inputs consist of highly informative structured data such as that
350 found in the EHR and similarly, the temporal patterns used in our work.

351

352 A permutation analysis was conducted on the classifiers developed in our case study. The ten
353 variables that were most predictive for each of the four classifiers developed in the obesity
354 incidence prediction research were identified and their average impact for each model's AUC-
355 ROC was computed. Two variables, 1-ALL04 (asthma in the pre-index visit) and 2-ALL03
356 (allergic rhinitis in the index visit), were among the most predictive variables for all four

357 classifiers and were the most prevalent condition patterns among the study population. African
358 American patients and patients enrolled in Medicaid were over-represented among patients who
359 had this temporal condition pattern. Our findings align with prior research on the association
360 between asthma and pediatric obesity which provide insight into why these variables were so
361 informative. Pediatric obesity and asthma are strongly associated, and early-life asthma
362 contributes to the onset of pediatric obesity.(27) Although the relationship between allergic
363 rhinitis and pediatric obesity is unclear, (28,29) allergic rhinitis has been shown to be comorbid
364 with pediatric asthma.(30) Low-income, urban, and racial minority children are
365 disproportionately impacted by both pediatric obesity and asthma, (17,31,32) which explains
366 their over-representation amongst patients with the most predictive features and the slight
367 classifier performance bias in their favor.

368
369 Some evidence of model bias relative to population subtypes was detected, although
370 unexpectedly the bias manifested as the classifiers tending to perform better on vulnerable
371 subgroups including African American patients and patients enrolled in Medicaid (a proxy for
372 lower socioeconomic status) than the entire population. These findings illustrate that there are
373 many complex ways that bias may emerge within EHR data, and that the context in which data is
374 collected and population impacted by a condition must be carefully considered when assessing
375 EHR data for bias. Prior work has shown that bias in machine learning typically results in lower
376 model performance for minorities due to an under-representation in data. (33,34) However,
377 possible biases need to be examined carefully. We hypothesize that the features most predictive
378 of obesity are more represented among patients in under-represented subgroups in this study,
379 which lends to the classification algorithms performing generally equitably if not slightly better
380 for African American patients and patients enrolled in Medicaid (a bias that is in favor of
381 vulnerable subpopulations). While these are simply associations and causality cannot be inferred,
382 our results support the idea that causes of bias in datasets and the models trained from them are
383 much more nuanced than initially thought.

384

385 **Limitations**

386

387 While this study serves to illustrate challenges and nuances associated with bias in machine
388 learning models developed using EHR, it does have limitations. First, the temporal condition
389 patterns mined from EHR data that were utilized as features for the machine learning models
390 only show associations. Findings are descriptive and the discovered temporal patterns and
391 comorbidities should be viewed in this light. No causality can be attributed to the associations
392 uncovered in this study. Similarly, the comparisons made in classifier performance differences
393 between population subgroups and the prevalence of temporal condition pattern prevalence are
394 also descriptive. No tests of statistical significance were performed, as this was a descriptive case
395 study that represents a first step in further research into bias in machine learning models
396 developed from EHR data. Finally, we acknowledge that when considering the outcome of
397 obesity in the machine learning prediction problem, minority patients comprised a greater
398 proportion of the case population compared to controls. This may have contributed to the model
399 performance bias in favor of vulnerable subgroups.

400

401 **Conclusion**

402

403 Our paper presents a comprehensive framework of how bias may be present within EHR data
404 and external sources of bias in the model development process, which in turn impacts machine
405 learning model development and clinical applications. Our pediatric obesity case study describes
406 a detailed exploration of a real-world machine learning model to contextualize how concepts
407 related to EHR data and machine learning model bias occur in an applied setting. We describe
408 how we evaluated our models for bias, and considered how these results are representative of
409 health disparity issues related to pediatric obesity. Finally, our paper presents a novel application
410 of data-driven temporal condition patterns that surround pediatric obesity incidence into a
411 predictive machine learning model. This adds to the limited body of literature on the use of
412 machine learning methods to study pediatric obesity and investigates the potential pitfalls in
413 using a machine learning approach when studying social significant health issues.

414

415 Bias is a complex and multi-faceted issue that is present in society and translates into data
416 collected in applied settings. We expect that our study may be used to define the types of bias
417 that researchers working with EHR data to develop machine learning models may look for, and

418 to understand that bias may manifest in machine learning models in unexpected ways. Our
419 approach to evaluating a machine learning model for bias and contextualizing our model
420 evaluation alongside clinical and psychosocial knowledge surrounding pediatric obesity provides
421 a useful blueprint for researchers developing and evaluating machine learning models with EHR
422 data in the obesity space and beyond. Finally, our findings support more equitable model
423 development, and may be used to guide researchers and clinicians in the precision medicine
424 space to consider the types of bias that may be present in machine learning models and how to
425 implement these models in clinical settings in a way that helps to address and not advance
426 existing systemic disparities.

427

428 Contributors

429
430 EC, SB, and AM contributed to the study's conception and design. EC conducted the study's
431 literature review and drafted the paper. EC and SB contributed to data acquisition, dataset
432 development, and data analysis. EC, SB, and AM contributed to interpretation of findings. AM
433 obtained funding to support this study and supervised study implementation. EC, SB, and AM
434 revised the paper. All authors approved the final version of the manuscript. All authors had full
435 access to the data in the study and were involved in data interpretation and writing of the report.
436 All authors had final responsibility for the decision to submit for publication.

437

438 Competing Interests

439
440 The funder of the study had no role in study design, data collection, data analysis, data
441 interpretation, or writing of the report. The authors declare that they have no conflict of interest.

442

443 Data sharing

444 Data cannot be shared publicly because of HIPAA requirements. Please contact
445 pedsnet@chop.edu with questions regarding data availability.

446

447 Acknowledgements

448
449 This work was supported by a grant from the Commonwealth Universal Research Enhancement
450 (C.U.R.E.) program funded by the Pennsylvania Department of Health—2015 Formula award—
451 SAP #4100072543. This work was also supported by funding from The Children's Hospital of
452 Philadelphia (CHOP)-Drexel Research Fellowship Program: Informatics and Analytics
453 Collaborative Research. We would like to thank the investigators of the Pediatric Big Health
454 Data initiative for their contributions. These individuals include: Christopher B. Forrest, MD,
455 PhD; L. Charles Bailey, MD, PhD; Shweta P. Chavan, MSEE; Rahul A. Darwar, MPH; Daniel
456 Forsyth; Chén C. Kenyon, MD, MSHP; Ritu Khare, PhD; Mitchell G. Maltenfort, PhD; Xueqin
457 Pang, PhD; Hanieh Razzaghi, MPH; Justine Shults, PhD; Levon H. Utidjian, MD, MBI from the
458 Children's Hospital of Philadelphia; Ana Diez Roux, MD, PhD, MPH; Amy H. Auchincloss,
459 PhD, MPH; Kimberly Daniels, MS; Anneclaire J. De Roos, PhD, MPH; J. Felipe Garcia-Espana,
460 MS, PhD; Irene Headen, PhD, MS; Félice Lê-Scherban, PhD, MPH; Steven Melly, MS, MA;
461 Yvonne L. Michael, ScD, SM; Kari Moore, MS; Abigail E. Mudd, MPH; Leah Schinasi, PhD,
462 MSPH from Drexel University and, Yong Chen, PhD; John H. Holmes, PhD; Rebecca A.
463 Hubbard, PhD; A. Russell Localio, JD, MPH, PhD from the University of Pennsylvania.

464

465

466 **References**

467

- 468 1. Murdoch TB, Detsky AS. The Inevitable Application of Big Data to Health Care. *JAMA*.
469 2013 Apr 3;309(13):1351–2.
- 470 2. J. Andreu-Perez, C. C. Y. Poon, R. D. Merrifield, S. T. C. Wong, G. -Z. Yang. Big Data for
471 Health. *IEEE Journal of Biomedical and Health Informatics*. 2015 Jul;19(4):1193–208.
- 472 3. Ferrao JC, Oliveira MD, Janela F, Martins HM. Preprocessing structured clinical data for
473 predictive modeling and decision support. *Applied clinical informatics*. 2016;7(04):1135–53.
- 474 4. Char DS, Shah NH, Magnus D. Implementing machine learning in health care—addressing
475 ethical challenges. *The New England journal of medicine*. 2018;378(11):981.
- 476 5. Carter SM, Rogers W, Win KT, Frazer H, Richards B, Houssami N. The ethical, legal and
477 social implications of using artificial intelligence systems in breast cancer care. *The Breast*.
478 2020;49:25–32.
- 479 6. Shah DS, Schwartz HA, Hovy D. Predictive Biases in Natural Language Processing Models:
480 A Conceptual Framework and Overview. In 2020. p. 5248–64.
- 481 7. Yadav P, Steinbach M, Kumar V, Simon G. Mining electronic health records (EHRs) A
482 survey. *ACM Computing Surveys (CSUR)*. 2018;50(6):1–40.
- 483 8. Gianfrancesco MA, Tamang S, Yazdany J, Schmajuk G. Potential biases in machine learning
484 algorithms using electronic health record data. *JAMA internal medicine*.
485 2018;178(11):1544–7.
- 486 9. Auerbach A, Fihn SD. Discovery, Learning, and Experimentation With Artificial
487 Intelligence–Based Tools at the Point of Care—Perils and Opportunity. *JAMA Network*
488 *Open*. 2021;4(3):e211474–e211474.
- 489 10. Cutillo CM, Sharma KR, Foschini L, Kundu S, Mackintosh M, Mandl KD. Machine
490 intelligence in healthcare—perspectives on trustworthiness, explainability, usability, and
491 transparency. *NPJ digital medicine*. 2020;3(1):1–5.
- 492 11. McCradden MD, Joshi S, Anderson JA, Mazwi M, Goldenberg A, Zlotnik Shaul R. Patient
493 safety and quality improvement: Ethical principles for a regulatory approach to bias in
494 healthcare machine learning. *Journal of the American Medical Informatics Association*.
495 2020;27(12):2024–7.
- 496 12. Chen IY, Szolovits P, Ghassemi M. Can AI help reduce disparities in general medical and
497 mental health care? *AMA journal of ethics*. 2019;21(2):167–79.
- 498 13. Skinner AC, Ravanbakht SN, Skelton JA, Perrin EM, Armstrong SC. Prevalence of obesity
499 and severe obesity in US children, 1999–2016. *Pediatrics*. 2018;141(3).

- 500 14. Kuczmariski RJ. 2000 CDC Growth Charts for the United States: methods and development.
501 Department of Health and Human Services, Centers for Disease Control and ...; 2002.
- 502 15. Karnik S, Kanekar A. Childhood obesity: a global public health crisis. *Int J Prev Med*.
503 2012;3(1):1–7.
- 504 16. Pulgarón ER. Childhood obesity: a review of increased risk for physical and psychological
505 comorbidities. *Clinical therapeutics*. 2013;35(1):A18–32.
- 506 17. Ogden CL, Carroll MD, Fakhouri TH, Hales CM, Fryar CD, Li X, et al. Prevalence of
507 obesity among youths by household income and education level of head of household—
508 United States 2011–2014. *Morbidity and mortality weekly report*. 2018;67(6):186.
- 509 18. Campbell EA, Qian T, Miller JM, Bass EJ, Masino AJ. Identification of temporal condition
510 patterns associated with pediatric obesity incidence using sequence mining and big data.
511 *International Journal of Obesity*. 2020;44(8):1753–65.
- 512 19. Campbell EA, Bass EJ, Masino AJ. Temporal condition pattern mining in large, sparse
513 electronic health record data: A case study in characterizing pediatric asthma. *Journal of the*
514 *American Medical Informatics Association*. 2020;27(4):558–66.
- 515 20. Bailey LC, Milov DE, Kelleher K, Kahn MG, Del Beccaro M, Yu F, et al. Multi-institutional
516 sharing of electronic health record data to assess childhood obesity. *PloS one*.
517 2013;8(6):e66192.
- 518 21. Weiner J, Abrams C. The Johns Hopkins ACG System Technical Reference Guide, Version
519 10.0. John Hopkins Bloomberg School of Public Health. 2011;
- 520 22. Arpey NC, Gaglioti AH, Rosenbaum ME. How socioeconomic status affects patient
521 perceptions of health care: a qualitative study. *Journal of primary care & community health*.
522 2017;8(3):169–75.
- 523 23. Schechter MS, Shelton BJ, Margolis PA, FitzSimmons SC. The association of
524 socioeconomic status with outcomes in cystic fibrosis patients in the United States.
525 *American journal of respiratory and critical care medicine*. 2001;163(6):1331–7.
- 526 24. About CHIP: Commonwealth of Pennsylvania [Internet]. 2019. Available from:
527 <https://www.chipcoverspakids.com/AboutCHIP/Pages/default.aspx>
- 528 25. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
529 Machine learning in Python. *the Journal of machine Learning research*. 2011;12:2825–30.
- 530 26. Christodoulou E, Ma J, Collins GS, Steyerberg EW, Verbakel JY, Van Calster B. A
531 systematic review shows no performance benefit of machine learning over logistic regression
532 for clinical prediction models. *Journal of clinical epidemiology*. 2019;110:12–22.

- 533 27. Azizpour Y, Delpisheh A, Montazeri Z, Sayehmiri K, Darabi B. Effect of childhood BMI on
534 asthma: a systematic review and meta-analysis of case-control studies. *BMC pediatrics*.
535 2018;18(1):1–13.
- 536 28. Sidell D, Shapiro NL, Bhattacharyya N. Obesity and the risk of chronic rhinosinusitis,
537 allergic rhinitis, and acute otitis media in school-age children. *The Laryngoscope*.
538 2013;123(10):2360–3.
- 539 29. Weinmayr G, Forastiere F, Büchele G, Jaensch A, Strachan DP, Nagel G, et al.
540 Overweight/obesity and respiratory and allergic disease in children: international study of
541 asthma and allergies in childhood (ISAAC) phase two. *PloS one*. 2014;9(12):e113996.
- 542 30. Stern J, Chen M, Fagnano M, Halterman JS. Allergic rhinitis co-morbidity on asthma
543 outcomes in city school children. *Journal of Asthma*. 2022;1–7.
- 544 31. Thakur N, Oh SS, Nguyen EA, Martin M, Roth LA, Galanter J, et al. Socioeconomic status
545 and childhood asthma in urban minority youths. The GALA II and SAGE II studies.
546 *American journal of respiratory and critical care medicine*. 2013;188(10):1202–9.
- 547 32. Persky VW, Slezak J, Contreras A, Becker L, Hernandez E, Ramakrishnan V, et al.
548 Relationships of race and socioeconomic status with prevalence, severity, and symptoms of
549 asthma in Chicago school children. *Annals of Allergy, Asthma & Immunology*.
550 1998;81(3):266–71.
- 551 33. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm
552 used to manage the health of populations. *Science*. 2019;366(6464):447–53.
- 553 34. Zhang H, Lu AX, Abdalla M, McDermott M, Ghassemi M. Hurtful words: quantifying
554 biases in clinical contextual word embeddings. In 2020. p. 110–20.
- 555
- 556

557 **Supporting Information**

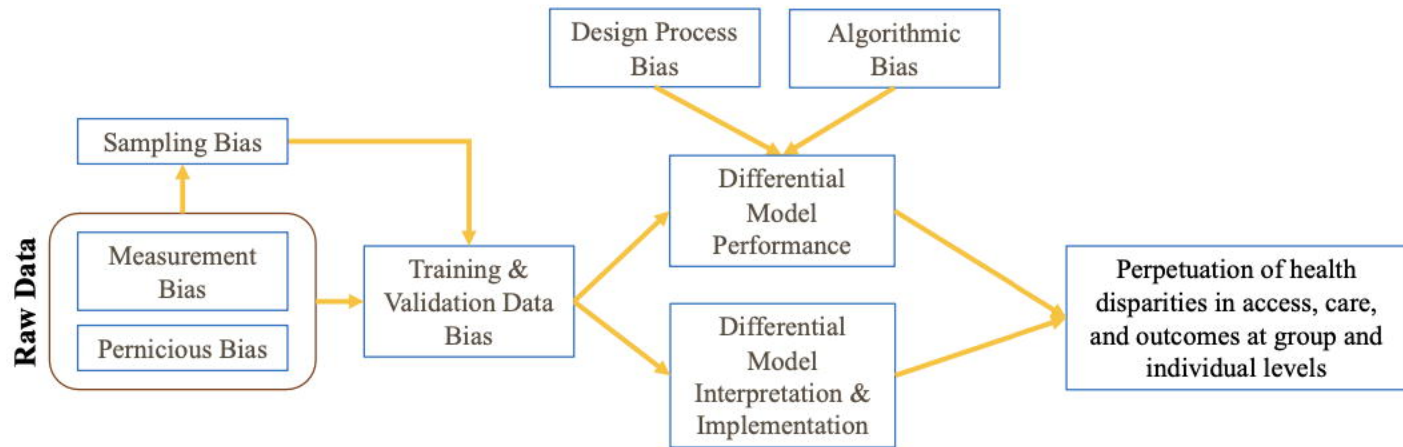
558

559 **S1 Table.** Temporal Diagnoses Included as Machine Learning Model Features

560

Diagnoses	561
<i>Pre-Index Visit Diagnoses</i>	
1-Chronic pharyngitis and tonsillitis	
1-Deafness, hearing loss	
1-Respiratory signs and symptoms	
1-Sleep Apnea	
1-Sleep Problems	
1-Dermatitis and eczema	
1-Seizure Disorder	
1-Asthma w/o Status Asthmaticus	
1-Constipation	
1-Urinary Symptoms	
1-Autism Spectrum Disorder	
1-Deafness, hearing loss	
1-Fever	
1-Gastroenteritis	
1-Headaches	
1-Nausea, vomiting	
<i>Index Visit Diagnoses</i>	
2-Chronic pharyngitis and tonsillitis	
2-Respiratory signs and symptoms	
2-Sleep Apnea	
2-Sleep Problems	
2-Allergic Rhinitis	
2-Dermatitis and eczema	
2-Developmental disorder	
2-Neurologic signs and symptoms	
2-Seizure Disorder	
2-Constipation	
2-Urinary Symptoms	
2-Allergic Rhinitis	
2-Asthma w/o Status Asthmaticus	
2-Dermatitis and eczema	
2-Autism Spectrum Disorder	
2-Developmental disorder	
2-Neurologic signs and symptoms	
2-Headaches	

Conceptual framework of Bias in ML



Initial Dataset:
M= 99,388
N = 49,694

Drop Cases without BMI
criterion
M = 34,172

Drop Controls without BMI
criterion
M = 18,328

M= 9,686
N = 4,843

Patients missing a Geo ID
for their address
M= 18
N= 9

M= 9,668
N = 4,834

Patients missing insurance
information within 2 years
of the index visit
M= 104
N= 52

Final Dataset
M= 9554
N= 4777