

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

1 **Can longitudinal electronic health record data identify patients at**
2 **higher risk of developing long COVID?**

3 Priya Shanmugam¹, Molly Bair¹, Emma Pendl-Robinson¹, Xindi C. Hu^{1*}, on behalf of N3C consortium

4 ¹Mathematica, Inc., Princeton, NJ 08543

5 *Corresponding author:

6 Xindi C. Hu, ScD

7 Mathematica, Inc., Princeton, NJ 08543

8 chu@mathematica-mpr.com

9

10 **Table of contents**

11	Abstract.....	1
12	1. Introduction.....	2
13	2. Our approach to predicting long COVID.....	5
14	3. Characteristics of Sample.....	6
15	4. Summary statistics for Cognitive, Fatigue, and Respiratory variations of Long	
16	COVID.....	9
17	5. Predictive model for long COVID outcomes.....	10
18	6. Apply the fair machine learning framework to the long COVID predictive models.	12
19	Conclusion.....	13
20	Acknowledgements.....	15
21	References.....	19
22	Tables.....	23
23	Table 1: Analytic Sample.....	23
24	Table 2: Model Performance Results.....	25
25	Table 3: Coefficient Values for the Binary Logistic Regression Models.....	26
26	Table 4: Impurity-Based Variable Importance Scores for the Random Forest Models.	27

27

28 **Abstract**

29 With hundreds of millions of COVID-19 infections to date, a considerable portion of the population has
30 developed or will develop long COVID. Understanding the prevalence, risk factors, and healthcare costs

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

31 of long COVID can be of significant societal importance. To investigate the utility of large-scale
32 electronic health record (EHR) data in identifying and predicting long COVID, we analyzed a sample of
33 1.23 million participants from the National COVID Cohort Collaborative (N3C), a longitudinal EHR data
34 repository from 80 sites in the US with over 8 million COVID-19 patients. We characterized the
35 prevalence of long COVID using a few different types of definitions to illustrate their relative strengths
36 and weaknesses. Then we developed machine learning models to predict the risk of developing long
37 COVID using demographic factors and comorbidity in the EHR. The risk factors for long COVID include
38 patient age; sex; smoking status; and comorbidities characterized by the Charlson Comorbidity Index
39 (CCI). We were able to predict three types of long COVID with low to moderate levels of accuracy (AUC
40 0.599 – 0.734). We found that age and CCI were most predictive of long COVID diagnosis. Ongoing
41 work includes applying the fair machine learning framework to the long COVID predictive models. We
42 are implementing fairness and bias mitigation methods to model fitting through the following steps,
43 selecting fairness metrics, preparing data and model, evaluating fairness metrics, applying bias mitigation
44 methods to the dataset, and comparing model results and fairness metrics before and after the mitigation.
45 The objective is to achieve equalized odds, a statistical notion that ensures classification algorithms do not
46 discriminate against protected groups (such as sex and race/ethnicity). Results from the fairness-based
47 machine learning will be included in the conference presentation.

48

49 1. Introduction

50 By November 2022, 94% of the US population was estimated to have been infected by SARS-CoV-2 at
51 least once (Klaassen et al., 2022). This proportion continues to rise as COVID-19 remains as an important
52 public health threat. Some people recover quickly after the initial infection, while others (between 10% to
53 30%) continue to experience symptoms even months after the initial infection (Herman et al., 2022). Long
54 COVID is defined by the World Health Organization (WHO) as persistent symptoms and/or long-term
55 complications following a probable or confirmed SARS-CoV-2 infection, usually 3 months from acute

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

56 infection and lasting longer than 2 months with no probable alternative diagnosis (Nittas et al., 2022;
57 World Health Organization, 2021). However, many challenges exist around diagnosing and managing
58 patients with long COVID and identifying individuals who may have elevated risks of developing long
59 COVID. Partially because the long-term effects of COVID-19 may exhibit differently for different
60 people, or because tracking symptoms that can wax and wane over a long period of time can be difficult.
61 Here we leverage a longitudinal electronic health record (EHR) data repository to test a few definitions of
62 long COVID and develop a risk prediction model to identify patients at higher risks of developing long
63 COVID. With hundreds of millions of COVID-19 infections to date, a considerable portion of the
64 population has developed or will develop long COVID. Understanding the prevalence, risk factors, and
65 healthcare costs of long COVID can be of significant societal importance. The healthcare sector needs to
66 better understand the demand for medical care due to long COVID and be prepared to deliver that. Public
67 health sector needs to be able to identify population segments at elevated risks for long COVID in order
68 to design and implement targeted intervention strategies.

69 More than 80 symptoms have been identified in the literature to be potentially associated with long
70 COVID (Nasserie et al., 2021). Like the initial SARS-CoV-2 infection, long COVID could affect multiple
71 organs and body systems. The most common symptoms include fatigue, dyspnea, cardiac abnormalities,
72 cognitive impairment, sleep disturbances, symptoms of post-traumatic stress disorder, muscle pain,
73 concentration problems, and headache (Crook et al., 2021). In addition to a long list of symptoms,
74 information on symptom severity, duration, and co-occurrence is generally lacking in the literature.
75 Several approaches exist among published studies to characterize long COVID: some uses a very liberal
76 definition of having at least one symptom which result in high percentage of patients estimated to have
77 long COVID but is likely an overestimate (CDC, 2023), some focuses on the most salient symptoms like
78 fatigue, headache, dyspnea and anosmia (Sudre et al., 2021), some focuses on symptom clusters to
79 capture the overlapping symptoms (Global Burden of Disease Long COVID Collaborators, 2022).

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

80 A wide range of prevalence estimates for long COVID may be a result of non-standardized definition,
81 virulence of different SARS-CoV-2 variants, and population immunity. Prevalence of as high as 80% was
82 observed in early studies, using a lenient definition of at least one symptom, among patients who were
83 previously hospitalized (Nasserie et al., 2021). The Centers of Disease Control and Prevention (CDC)
84 estimated that nearly one in five American adults who have had COVID-19 still have long COVID (CDC,
85 2022). More recently, the Global Burden of Disease Long COVID Collaborators estimated that 6.2%
86 individuals have at least one of the three select-reported long COVID symptom clusters (Global Burden
87 of Disease Long COVID Collaborators, 2022). As a result of the uncertainty around the prevalence of
88 long COVID, the economic cost of long COVID is hard to fully grasp, but safe to assume it would be an
89 enormous number. Cutler and Summers estimated the economic costs of long COVID to be \$2.6 trillion
90 in 2020, and later updated it to be \$3.7 trillion as a result of higher prevalence of long COVID than
91 previously assumed (D. Cutler, 2022; D. M. Cutler & Summers, 2020).

92 The enormity of the economic costs and societal impact of long COVID highlights the necessity to
93 develop better methods for detecting, treating, and preventing long COVID. In this study, we analyzed
94 data from the National COVID Cohort Collaborative (N3C), a longitudinal EHR data repository from 80
95 sites in the US with over 8 million COVID-19 patients (Haendel et al., 2021). We characterize the
96 prevalence of long COVID using a few different types of definitions to illustrate their relative strengths
97 and weaknesses. Then we develop machine learning models to predict the risk of developing long COVID
98 using demographic factors and comorbidity in the EHR. We discuss model performance with a focus on
99 features with strong predictive power that can be utilized to design early detection or targeted intervention
100 strategies. We the discuss how EHR data derived risk prediction model can be used to enhance an
101 individual COVID-19 risk calculator, 19andMe, to support individual and clinical decision-making.
102 Ongoing work also includes assessing algorithmic bias and mitigating them using techniques including
103 resampling and reweighting.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

104 2. Our approach to predicting long COVID

105 Several features of long COVID make identifying positive cases using claims and EHR data challenging.
106 First, some symptoms of long COVID may or may not be detectable, depending on which diagnostic test
107 is used (Mancini et al., 2021). Second, due to early skepticism within the medical profession about the
108 existence of COVID-19, there have been significant informal, patient-led efforts to report and track the
109 various symptoms of long COVID (Patient Led Research Collaborative, 2023). Third, symptoms
110 occurring during the long COVID episode must be separated from the initial episode of COVID itself,
111 which may vary in length and severity. Fourth, like other medical conditions, a claims-based definition
112 must be clinically relevant with respect to disease severity, and account for health conditions a patient
113 may have had prior to their COVID and long COVID episodes.

114 Given these challenges, we tested several approaches to defining long COVID using data from the N3C
115 Enclave. We first defined a six-month lookback & look-forward window around each patients' index
116 COVID-19 episode and limited our sample to patients who have healthcare utilization (defined as having
117 a condition recorded in the N3C database) prior to, and after, these windows. This allowed us to
118 accurately classify a lack of healthcare utilization during the lookback or look-forward windows as true
119 non-utilization, as opposed to a patient death, relocation, or otherwise missing data. We then identified all
120 OMOP concept ID condition codes that appear on a patient's record between 90 and 180 days following
121 their index COVID-19 diagnosis. Excluding conditions recorded within 90 days after an initial COVID-
122 19 diagnosis ensured that we do not flag symptoms of COVID-19 as symptoms of long COVID. From the
123 resulting set of symptoms, we removed any symptoms that appeared on a patient's record during their
124 lookback window. A drawback of this approach is that it reduces the likelihood that we detect long
125 COVID cases among people with chronic conditions or comorbidities, who are actually most likely to
126 contract long COVID. We additionally removed any symptoms during the look-forward period that do
127 not appear on the list of 76 conditions in the N3C Long COVID OMOP condition code set. The
128 remaining symptoms were considered as the patient's long COVID symptom burden.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

129 We considered that the N3C long COVID OMOP condition code set contains several potential
130 redundancies, and that these conditions might be over-reported in claims and EHR data due to upcoding.
131 For example, the following pairs of symptoms were similar to one another, yet each appeared as an
132 independent item in the OMOP code set: sleep disturbances and insomnia, nausea and vomiting, body
133 aches and muscle pain, and fatigue and exercise intolerance. We consolidated these redundancies in a
134 patient's look-forward period prior to calculating each patient's long COVID symptom score. This
135 reduced the number of long covid symptoms from 76 to 72.

136 **3 symptom cluster approach**

137 Lastly, we built off of the recent literature by calculating a cluster-specific score for each of three
138 symptom clusters (Cognitive, Respiratory, and Fatigue) (Global Burden of Disease Long COVID
139 Collaborators, 2022). We created three cluster-specific, binary long COVID definitions based on a cutoff
140 of at least 3 distinct long COVID symptoms within the relevant cluster. This approach provided a few
141 advantages: it built off the clinical literature identifying these three symptom clusters as distinct, allowing
142 us to classify each patient into each subtype without concern for whether a patient exhibits symptoms
143 from multiple clusters. Second, it allowed us to capture variation in the relationship between patient
144 demographics and various long COVID subtypes. Third, it allowed us to draw distinctions between more
145 and less severe symptom clusters of long COVID, and potentially identify predictive targets that capture
146 high-morbidity or high-cost events.

147 **3. Characteristics of Sample**

148 We used the N3C de-identified "tier 2 access" data set.¹ A limitation of our analysis is that the Patient
149 Severity and Scores dataset, condition occurrence, and drug exposure datasets are frequently updated

¹ Tire 2 access is the patient level EHR data where 17 patient identifier variables are removed and longitudinal data is data-shifted to safeguard privacy (Haendel et al., 2021). We selected patients for the analysis from the Patient Severity and Scores dataset (Release-v70-2022-03-19). The dataset contains patient level information from all sites. N3C identifies positive COVID-19 cases through their COVID-19 Phenotype Inclusion Criteria (Pfaff, 2022). Though this did not factor into our analysis, another feature of the N3C data is a 1:2 case to control ratio of patients with identifies lab-confirmed, suspected, and possible cases of COVID-19 to patients who have been screen and tested negative for COVID-19. In addition, we linked patient ids from the Patient Severity and Scores dataset with

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

150 within the N3C data enclave and new records being added. Our analysis used the Patient Severity and
151 Scores dataset version 70 which was released on March 19, 2022. Thus, all patients in our analysis were
152 diagnosed on or before March 19, 2022. Therefore, our results did not capture recent changes in the
153 COVID-19 pandemic including the rise in the prevalence of the omicron variant, vaccination, and vaccine
154 booster doses.

155 For our analysis, we started with 13,151,716 patients from the Patient Severity and Scores database and
156 filtered down to the patients with lab confirmed positive cases (3,724,542 patients). Finally, we selected
157 the patients who had conditions occurrences within the 6 months preceding and proceeding diagnosis of
158 COVID-19,² whose COVID-19 visit start date was between September 1, 2020, and September 1, 2021,
159 and who did not die from COVID-19. This provided us with our final analytic sample of 1,234,119
160 patients.

161 We selected covariates based on our literature review of covariates other researchers found were related
162 to long COVID symptoms (Tsampasian et al., 2023). We identified the following covariates using the
163 Patient Severity and Scores dataset: age, sex, race, ethnicity, and current or former smoking status. We
164 categorized the patient's age at diagnosis into 9 categories (0-9, 10-19, 20-29, 30-39, 40-49, 50-59, 60-69,
165 70-79, and 80+ years old) for patients missing age, we imputed age as the median age 41.2 which
166 translates to the age category 40-49. We categorized sex as male, female, or other. Race was defined as
167 White, Black or African American, Asian, Native Hawaiian or Other Pacific Islander, Other, and
168 Missing/Unknown. Ethnicity was defined as Hispanic or Latino, Not Hispanic or Latino, and Ethnicity
169 Missing or Unknown.

two additional datasets provided in N3C, the condition occurrence dataset (Release-v70-2022-03-19) and the drug exposures dataset (Release-v70-2022-03-19) to provide additional longitudinal information about the patients in the database. These condition occurrence dataset records diagnosis, signs, and symptoms of conditions and drug exposure dataset records introduction of drugs into the body of the patient overtime (Observational Health Data Sciences and Informatics, 2018).

² For selecting the patients how has conditions occurrences within 6 months preceding and proceeding diagnosis, we used the condition occurrence dataset (Release-v70-2022-03-19)

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

170 In addition, we identified the following covariates using the OMOP code set ids in the condition
171 occurrence and drug exposure datasets for the dates prior to the covid diagnosis: COVID-19 vaccination,
172 pregnancy, hypertension, obesity, immunocompromised, and Charlson Comorbidity Index (CCI) score
173 (Glasheen et al., 2019). For each patient, if the OMOP codes for the comorbidity's hypertension, obesity,
174 immunocompromised, or CCI score appear in the condition occurrence data before the COVID-19
175 diagnosis, then we categorized the patient as having the respective comorbidity. We considered a patient
176 vaccinated if a patient received one or more doses of a Pfizer, Moderna, or Johnson and Johnson vaccine
177 before diagnosis of COVID-19 as recorded in the drug exposure dataset. We considered a patient
178 pregnant if the condition start date was within a year of COVID-19 diagnosis. See Appendix A1 for full
179 list of the OMOP codes we used to create the covariates.

180 Our sample was 56% (n = 692,547) female, 68% (n = 1,031,221) White, 2% (n = 25,083) received at least
181 one dose of a COVID-19 vaccine, 22% (n = 27,3255) were 60 years old or older, see Table 1 for
182 additional attributes of the analytic sample.

183 Several limitations exist in our data. First, the observed conditions dataset only recorded diagnoses, signs,
184 or symptoms of a condition either observed by a provider or reported by the patient (Blacketer, 2021).
185 This limitation might have affected the makeup of the sample through: 1) underreporting conditions
186 because patients might have had conditions treated at another facility not sharing data with N3C, 2)
187 excluding the healthiest patients because they did not have any reported conditions in the six months
188 before and after the COVID-19 diagnosis. The vaccination status had a similar limitation where
189 vaccination data was reported by clinics, pharmacies, or patients (Blacketer, 2021). Due to this limitation,
190 our sample may have underestimated the vaccination rates if patients received vaccinations at a facility
191 not sharing data with N3C. The low vaccination rates in our data might have led to a weaker relationship
192 between vaccination and long COVID symptoms in our findings. Due to our concerns with the
193 vaccination variable, we choose not to include vaccination in our modeling.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

194 **4. Summary statistics for Cognitive, Fatigue, and Respiratory variations of Long COVID**

195 We examined long COVID by three symptom clusters -- cognitive, fatigue, and respiratory -- as defined
196 by the Global Burden of Disease Long COVID Collaborators (Global Burden of Disease Long COVID
197 Collaborators, 2022). We classified patients as being positive for the individual symptom clusters if the
198 patients had a qualifying symptom present between 3 to 6 months after having a lab confirmed case of
199 COVID-19 and that symptom was not reported in the 6 months preceding their lab confirmed case of
200 COVID-19. For a full list of the symptoms and OMOP codes we used as qualifying cognitive, fatigue,
201 and respiratory long COVID symptoms, see Appendix A2, A3, and A4 respectively. Then we focused our
202 analysis on the three separate symptom clusters and predicted the risk of a patient having had each
203 symptom cluster of long COVID systems.

204 The fatigue long COVID symptom cluster was the most common in our sample with 4.7% (n = 57,483) of
205 all patients, followed by respiratory cluster with 2.5% (n = 30,668) of our sample, and the cognitive
206 symptom cluster with only 0.2% (n = 2,937) of the sample (Table 1). In general, for all three symptom
207 clusters, there was an age gradient where patients in younger age groups have lower rates of symptom
208 clusters than patients in older groups. Female patients had higher rates (5.7%) of the fatigue symptom
209 cluster than male patients (3.3%) and there was no difference in rates between the sexes for the cognitive
210 symptom clusters (female = 0.3%; male = 0.2%) and respiratory (female = 2.6%; male = 2.3%). Black
211 and White patients had higher prevalence than other race groups for all three symptom clusters. There
212 was not a meaningful difference in prevalence by ethnicity. Patients who had cognitive, fatigue, and
213 respiratory symptom clusters had higher average CCI scores than those who did not experience said
214 symptom clusters. The highest average CCI scores were patients who had cognitive cluster symptoms
215 (mean = 2.7), compared to fatigue (mean = 1.4) and respiratory (mean = 1.4) symptom clusters. For other
216 variables we did not observe large differences between groups.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

217 5. Predictive model for long COVID outcomes

218 **Model training methods**

219 We trained two types of machine learning models (a binary logistic regression model and a binary random
220 forest model) on all three long COVID clusters separately: respiratory, fatigue, and cognitive. We used a
221 70-30 train-test split, resulting in 863,883 training observations and 370,236 test observations. To
222 accommodate for long COVID being a rare outcome in all three of our clusters, we trained all models on
223 down sampled data and tested them on the original, unbalanced data. To create the down sampled training
224 data, we took five bootstrapped samples with replacement. Each bootstrapped sample had an equal
225 number of positive and negative cases, which was also equal to the total number of positive cases in the
226 original training data.

227 **Model performance**

228 The random forest model outperformed the logistic regression model at predicting the respiratory and
229 cognitive clusters, but the logistic regression model slightly outperformed the random forest model at
230 predicting the fatigue outcome (Table 2). In all six models, the AUC remained fairly consistent between
231 the train and test data, but the model precision (a measure of how often a model's positive predictions are
232 correct) dropped quite dramatically between the train and test data. This could be because long COVID
233 was a rare outcome, or that there is low signal in the data.

234 **Covariates**

235 To measure feature importance, we considered coefficient values for the binary logistic regression models
236 and impurity-based variable importance scores for the random forest models (Tables 3 and 4). Among our
237 six models, there were some variations in which covariates were most predictive. The CCI score was the
238 most important feature in all three clusters for the random forest models. However, for the binary logistic
239 regression models, all three clusters had different most important features (Table 3). For the respiratory
240 cluster, the two most important features are belonging to the 0 – 9 age group and being a current or
241 former smoker. For the fatigue cluster, the two most important features are belonging to the 0 – 9 age

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

242 group and being male. Finally, for the cognitive cluster, the two most important features are being 80
243 years old or older and belonging to the 70 – 79 age group.

244 **Comparisons with other long COVID predictive models**

245 One challenge in comparing models is that there is not a consistent definition of long COVID that is
246 universally accepted. The two papers described below offer two alternative definitions of long COVID.

247 One study also uses data from N3C to predict outcomes of long COVID. The authors consider visiting a
248 long COVID clinic as an indicator that a patient has long COVID (Pfaff, 2022). This is a narrower
249 definition of long COVID than our definition in this study. Their XGBoost model results in an AUC of
250 0.92 for all patients in their sample, 0.90 for hospitalized patients, and 0.85 for non-hospitalized patients.

251 The simplicity of their measure could explain why their results are better than the results from our
252 models. However, there are also drawbacks to defining long COVID as a visit to a long COVID clinic. a
253 long COVID clinic is a very specific type of care and someone’s long COVID symptoms have to be
254 pretty severe to be referred to a long COVID clinic. Thus, the Pfaff definition undercounts the number of
255 patients with long COVID especially patients with more mild symptoms of long COVID.

256 Another study used age, sex, and the number of symptoms a patient experienced in the first week of
257 infection to classify duration of COVID-19 as either short (less than 10 days) or long (28 days or more,
258 which is shorter than the WHO recommendation) using k-means clustering (Sudre et al., 2021). The
259 authors analyzed self-reported data from 4,182 cases of COVID-19 through the COVID Symptom Study
260 app and used a random forest model and a logistic regression model to predict long COVID, as defined by
261 having symptoms for 28 days or longer. The random forest model, which included first week symptoms
262 and other comorbidities, performed moderately well with an AUC of 0.768. The logistic regression model
263 was much simpler and included only age, sex, and number of symptoms in the first week, with an AUC of
264 0.767. One key difference between Sudre et al and our model is that Sudre et al included the number of
265 symptoms in the first week which is likely a proxy for identify which patients had more severe initial
266 COVID infections, whereas our model had limited variables representing the severity of initial infection.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

267 6. Apply the fair machine learning framework to the long COVID predictive models

268 Timely risk assessment of long COVID outcomes can improve patient care and healthcare resource
269 allocation. However, disparities across different sex, race/ethnicity, and social economic status in
270 COVID-19 patient outcomes have been well-documented(Kharroubi & Diab-El-Harake, 2022; Webb
271 Hooper et al., 2020). For example, while males have high risks for severe COVID outcomes such as death
272 and ICU admission, females have been reported in the literature to have higher risks for long COVID
273 (Cohen & van der Meulen Rodgers, 2023). In addition, compared to white patients, patients from
274 racial/ethnic minority groups had significantly different odds of developing long COVID symptoms and
275 conditions (Khullar et al., 2023). It is important to carefully examine model fairness across different
276 population subgroups to achieve optimal and fair clinical decision-making.

277 We plan to confirm our models achieve similar performances across different sex/race/ethnicity groups by
278 calculating the performance score using AUROC, F1 score, precision, recall, sensitivity, and specificity
279 for males/females, different race groups including White, Black or African American, Asian, Native
280 Hawaiian or Pacific Islander, and Other or Missing/Unknown, and different ethnicity (Hispanic and non-
281 Hispanic) groups. This will allow us to detect if there is any bias in model performance. In addition to
282 model performance, other fairness metrics will also be evaluated since the outcome (diagnoses codes in
283 EHR) is not considered “ground truth”. We will also evaluate predictive equality, equal opportunity, and
284 statistical parity of the model predictions.

285 If we detect disparity in fairness metrics, we plan to address it using various techniques including
286 disparate impact remover (a preprocessing technique that edits input values to increase fairness between
287 groups), reweighting (producing weights for each subgroup for each outcome class to achieve statistical
288 parity), and resampling (oversampling or under-sampling to achieve statistical parity). We will compare
289 the effects of different mitigation techniques and discuss their impact on the balance between model
290 performance and fairness metrics. The method maximizing model predictive performance may be
291 different from the method achieving the highest fairness metrics. We will discuss the trade-off between

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

292 model performance and fairness in the context of using real-world EHR data to develop predictive models
293 for long COVID.

294 Conclusion

295 Overarching takeaways regarding model performance

296 In this study, we explored the feasibility to train a predictive model to provide insights into the risk
297 factors for long COVID. We leveraged a diverse and harmonized electronic health record data source and
298 applied minimum filters to build a model with high potential to be applied to general population. After
299 using techniques to handle unbalanced sample, our predictive models achieved moderate performance
300 accuracy (AUC 0.599 – 0.730 for the random forest models, AUC 0.602 – 0.734 for the logistic
301 regression models). It was interesting to observe that logistic regression models have comparable
302 performance to random forest models, potentially because we included relatively limited number of
303 predictors that are easy to acquire. Other predictive models developed in the past had stronger
304 performances but they either leveraged clinical factors that can be hard to acquire such as blood oxygen,
305 blood pH (Bennett et al., 2021), or they limited the population to a small and relatively homogenous
306 groups of people, such as users of the COVID Symptom Study app or patients visiting a long COVID
307 clinic (Pfaff, 2022; Sudre et al., 2021).

308 Comparison with 19&Me severe COVID risk model

309 Previously we leveraged the same data source (N3C) and machine learning models to build a predictive
310 model for severe COVID-19 outcomes. Severe COVID-19 outcome is defined as having a score of 6 or
311 above using the Clinical Progression Scale (CPS) established by the World Health Organization for
312 COVID-19 clinical research (Marshall et al., 2020). In a sample of 864,8080 COVID-19 positive
313 patients, 15,401 (1.75%) has this outcome. Many risk factors are related to elevated COVID-19 risk, and
314 others have published machine learning models that achieve high predictive performance (AUROC =
315 0.87) using a long list of predictors including demographics (age/sex/race), health conditions
316 (comorbidities) and clinical characteristics (blood pH, respiratory rate, oxygen saturation...). We limit

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

317 the predictors to those easily accessible in the 19andMe app (age/sex/race/comorbidities/smoking
318 status/vaccination status). Since the outcome is highly unbalanced, we used down sampling and up
319 weighting to handle the sample unbalance. We experimented with different modeling techniques such as
320 random forest and logistic regression, with and without the downsampling techniques. The best
321 performing model was the logistic regression with downsampled data (Table 2).

322 **Implications for future work: learnings on N3C data**

323 The N3C is a nationally representative, harmonized data resource that we can leverage for COVID-19
324 research, among other data sources that our team already have experience with. In addition to the data
325 tables that we explored in the current work, several other data tables such as visitation occurrence can be
326 potentially useful for future similar work. In our current work, we can only create a flag variable for
327 whether or not the patient has certain symptoms, without the ability to rate its severity. Information in the
328 visitation occurrence table can be helpful because it has visit start and end date, which can then be used to
329 characterize the severity of the condition. This can be a future area of modeling enhancement.

330 Machine learning models have utilities in building predictive models for both severe COVID outcomes
331 and long COVID. The severe COVID model had better performance, partially because of a clear
332 definition of severe COVID and relative follow-up period. The model building process could benefit
333 from close collaborations between data scientists, health service researchers and clinicians. The models
334 developed in this study may be valuable to both payers to improve their understanding of the risks of their
335 insured population, and to public health officials to better plan for the resources needed to improve
336 population health in the aftermaths of COVID-19 pandemic.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

337 Acknowledgements

338 This work is supported by Mathematica.

339 N3C Attribution

340 The analyses described in this manuscript were conducted with data or tools accessed through the
341 NCATS N3C Data Enclave <https://covid.cd2h.org> and N3C Attribution & Publication Policy v 1.2-2020-
342 08-25b supported by NCATS Contract No. 75N95023D00001, Axle Informatics Subcontract: NCATS-
343 P00438-B, and [insert additional funding agencies or sources and reference numbers as declared by the
344 contributors in their form response above]. This research was possible because of the patients whose
345 information is included within the data and the organizations ([https://ncats.nih.gov/n3c/resources/data-
346 contribution/data-transfer-agreement-signatories](https://ncats.nih.gov/n3c/resources/data-contribution/data-transfer-agreement-signatories)) and scientists who have contributed to the on-going
347 development of this community resource [<https://doi.org/10.1093/jamia/ocaa196>].

348

349 Disclaimer The N3C Publication committee confirmed that this manuscript msid:1865.948 is in
350 accordance with N3C data use and attribution policies; however, this content is solely the responsibility of
351 the authors and does not necessarily represent the official views of the National Institutes of Health or the
352 N3C program.

353

354 IRB

355 The N3C data transfer to NCATS is performed under a Johns Hopkins University Reliance Protocol #
356 IRB00249128 or individual site agreements with NIH. The N3C Data Enclave is managed under the
357 authority of the NIH; information can be found at <https://ncats.nih.gov/n3c/resources>.

358

359 Individual Acknowledgements For Core Contributors

360 We gratefully acknowledge the following core contributors to N3C:

361

362 Adam B. Wilcox, Adam M. Lee, Alexis Graves, Alfred (Jerrod) Anzalone, Amin Manna, Amit Saha,
363 Amy Olex, Andrea Zhou, Andrew E. Williams, Andrew Southerland, Andrew T. Girvin, Anita Walden,
364 Anjali A. Sharathkumar, Benjamin Amor, Benjamin Bates, Brian Hendricks, Brijesh Patel, Caleb
365 Alexander, Carolyn Bramante, Cavin Ward-Caviness, Charisse Madlock-Brown, Christine Suver,
366 Christopher Chute, Christopher Dillon, Chunlei Wu, Clare Schmitt, Cliff Takemoto, Dan Housman,
367 Davera Gabriel, David A. Eichmann, Diego Mazzotti, Don Brown, Eilis Boudreau, Elaine Hill, Elizabeth
368 Zampino, Emily Carlson Marti, Emily R. Pfaff, Evan French, Farrukh M Koraisly, Federico Mariona,
369 Fred Prior, George Sokos, Greg Martin, Harold Lehmann, Heidi Spratt, Hemalkumar Mehta, Hongfang
370 Liu, Hythem Sidky, J.W. Awori Hayanga, Jami Pincavitch, Jaylyn Clark, Jeremy Richard Harper, Jessica
371 Islam, Jin Ge, Joel Gagnier, Joel H. Saltz, Joel Saltz, Johanna Loomba, John Buse, Jomol Mathew, Joni
372 L. Rutter, Julie A. McMurry, Justin Guinney, Justin Starren, Karen Crowley, Katie Rebecca Bradwell,
373 Kellie M. Walters, Ken Wilkins, Kenneth R. Gersing, Kenrick Dwain Cato, Kimberly Murray, Kristin
374 Kostka, Lavance Northington, Lee Allan Pyles, Leonie Misquitta, Lesley Cottrell, Lili Portilla, Mariam

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

375 Deacy, Mark M. Bissell, Marshall Clark, Mary Emmett, Mary Morrison Saltz, Matvey B. Palchuk,
376 Melissa A. Haendel, Meredith Adams, Meredith Temple-O'Connor, Michael G. Kurilla, Michele Morris,
377 Nabeel Qureshi, Nasia Safdar, Nicole Garbarini, Noha Sharafeldin, Ofer Sadan, Patricia A. Francis,
378 Penny Wung Burgoon, Peter Robinson, Philip R.O. Payne, Rafael Fuentes, Randeep Jawa, Rebecca
379 Erwin-Cohen, Rena Patel, Richard A. Moffitt, Richard L. Zhu, Rishi Kamaleswaran, Robert Hurley,
380 Robert T. Miller, Saiju Pyarajan, Sam G. Michael, Samuel Bozzette, Sandeep Mallipattu, Satyanarayana
381 Vedula, Scott Chapman, Shawn T. O'Neil, Soko Setoguchi, Stephanie S. Hong, Steve Johnson, Tellen D.
382 Bennett, Tiffany Callahan, Umit Topaloglu, Usman Sheikh, Valery Gordon, Vignesh Subbian, Warren A.
383 Kibbe, Wendy Hernandez, Will Beasley, Will Cooper, William Hillegass, Xiaohan Tanner Zhang.
384 Details of contributions available at covid.cd2h.org/core-contributors

385

386 Data Partners with Released Data

387 The following institutions whose data is released or pending:

388

389 Available: Advocate Health Care Network — UL1TR002389: The Institute for Translational Medicine
390 (ITM) • Aurora Health Care Inc — UL1TR002373: Wisconsin Network For Health Research • Boston
391 University Medical Campus — UL1TR001430: Boston University Clinical and Translational Science
392 Institute • Brown University — U54GM115677: Advance Clinical Translational Research (Advance-
393 CTR) • Carilion Clinic — UL1TR003015: iTHRIV Integrated Translational health Research Institute of
394 Virginia • Case Western Reserve University — UL1TR002548: The Clinical & Translational Science
395 Collaborative of Cleveland (CTSC) • Charleston Area Medical Center — U54GM104942: West Virginia
396 Clinical and Translational Science Institute (WVCTSI) • Children's Hospital Colorado —
397 UL1TR002535: Colorado Clinical and Translational Sciences Institute • Columbia University Irving
398 Medical Center — UL1TR001873: Irving Institute for Clinical and Translational Research • Dartmouth
399 College — None (Voluntary) Duke University — UL1TR002553: Duke Clinical and Translational
400 Science Institute • George Washington Children's Research Institute — UL1TR001876: Clinical and
401 Translational Science Institute at Children's National (CTSA-CN) • George Washington University —
402 UL1TR001876: Clinical and Translational Science Institute at Children's National (CTSA-CN) • Harvard
403 Medical School — UL1TR002541: Harvard Catalyst • Indiana University School of Medicine —
404 UL1TR002529: Indiana Clinical and Translational Science Institute • Johns Hopkins University —
405 UL1TR003098: Johns Hopkins Institute for Clinical and Translational Research • Louisiana Public
406 Health Institute — None (Voluntary) • Loyola Medicine — Loyola University Medical Center • Loyola
407 University Medical Center — UL1TR002389: The Institute for Translational Medicine (ITM) • Maine
408 Medical Center — U54GM115516: Northern New England Clinical & Translational Research (NNE-
409 CTR) Network • Mary Hitchcock Memorial Hospital & Dartmouth Hitchcock Clinic — None (Voluntary)
410 • Massachusetts General Brigham — UL1TR002541: Harvard Catalyst • Mayo Clinic Rochester —
411 UL1TR002377: Mayo Clinic Center for Clinical and Translational Science (CCaTS) • Medical University
412 of South Carolina — UL1TR001450: South Carolina Clinical & Translational Research Institute (SCTR)
413 • MITRE Corporation — None (Voluntary) • Montefiore Medical Center — UL1TR002556: Institute for
414 Clinical and Translational Research at Einstein and Montefiore • Nemours — U54GM104941: Delaware
415 CTR ACCEL Program • NorthShore University HealthSystem — UL1TR002389: The Institute for
416 Translational Medicine (ITM) • Northwestern University at Chicago — UL1TR001422: Northwestern
417 University Clinical and Translational Science Institute (NUCATS) • OCHIN — INV-018455: Bill and

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

418 Melinda Gates Foundation grant to Sage Bionetworks • Oregon Health & Science University —
419 UL1TR002369: Oregon Clinical and Translational Research Institute • Penn State Health Milton S.
420 Hershey Medical Center — UL1TR002014: Penn State Clinical and Translational Science Institute •
421 Rush University Medical Center — UL1TR002389: The Institute for Translational Medicine (ITM) •
422 Rutgers, The State University of New Jersey — UL1TR003017: New Jersey Alliance for Clinical and
423 Translational Science • Stony Brook University — U24TR002306 • The Alliance at the University of
424 Puerto Rico, Medical Sciences Campus — U54GM133807: Hispanic Alliance for Clinical and
425 Translational Research (The Alliance) • The Ohio State University — UL1TR002733: Center for Clinical
426 and Translational Science • The State University of New York at Buffalo — UL1TR001412: Clinical and
427 Translational Science Institute • The University of Chicago — UL1TR002389: The Institute for
428 Translational Medicine (ITM) • The University of Iowa — UL1TR002537: Institute for Clinical and
429 Translational Science • The University of Miami Leonard M. Miller School of Medicine —
430 UL1TR002736: University of Miami Clinical and Translational Science Institute • The University of
431 Michigan at Ann Arbor — UL1TR002240: Michigan Institute for Clinical and Health Research • The
432 University of Texas Health Science Center at Houston — UL1TR003167: Center for Clinical and
433 Translational Sciences (CCTS) • The University of Texas Medical Branch at Galveston —
434 UL1TR001439: The Institute for Translational Sciences • The University of Utah — UL1TR002538:
435 Uhealth Center for Clinical and Translational Science • Tufts Medical Center — UL1TR002544: Tufts
436 Clinical and Translational Science Institute • Tulane University — UL1TR003096: Center for Clinical
437 and Translational Science • The Queens Medical Center — None (Voluntary) • University Medical
438 Center New Orleans — U54GM104940: Louisiana Clinical and Translational Science (LA CaTS) Center
439 • University of Alabama at Birmingham — UL1TR003096: Center for Clinical and Translational Science
440 • University of Arkansas for Medical Sciences — UL1TR003107: UAMS Translational Research
441 Institute • University of Cincinnati — UL1TR001425: Center for Clinical and Translational Science and
442 Training • University of Colorado Denver, Anschutz Medical Campus — UL1TR002535: Colorado
443 Clinical and Translational Sciences Institute • University of Illinois at Chicago — UL1TR002003: UIC
444 Center for Clinical and Translational Science • University of Kansas Medical Center — UL1TR002366:
445 Frontiers: University of Kansas Clinical and Translational Science Institute • University of Kentucky —
446 UL1TR001998: UK Center for Clinical and Translational Science • University of Massachusetts Medical
447 School Worcester — UL1TR001453: The UMass Center for Clinical and Translational Science
448 (UMCCTS) • University Medical Center of Southern Nevada — None (voluntary) • University of
449 Minnesota — UL1TR002494: Clinical and Translational Science Institute • University of Mississippi
450 Medical Center — U54GM115428: Mississippi Center for Clinical and Translational Research (CCTR) •
451 University of Nebraska Medical Center — U54GM115458: Great Plains IDeA-Clinical & Translational
452 Research • University of North Carolina at Chapel Hill — UL1TR002489: North Carolina Translational
453 and Clinical Science Institute • University of Oklahoma Health Sciences Center — U54GM104938:
454 Oklahoma Clinical and Translational Science Institute (OCTS) • University of Pittsburgh —
455 UL1TR001857: The Clinical and Translational Science Institute (CTSI) • University of Pennsylvania —
456 UL1TR001878: Institute for Translational Medicine and Therapeutics • University of Rochester —
457 UL1TR002001: UR Clinical & Translational Science Institute • University of Southern California —
458 UL1TR001855: The Southern California Clinical and Translational Science Institute (SC CTSI) •
459 University of Vermont — U54GM115516: Northern New England Clinical & Translational Research
460 (NNE-CTR) Network • University of Virginia — UL1TR003015: iTHRIV Integrated Translational health
461 Research Institute of Virginia • University of Washington — UL1TR002319: Institute of Translational
462 Health Sciences • University of Wisconsin-Madison — UL1TR002373: UW Institute for Clinical and
463 Translational Research • Vanderbilt University Medical Center — UL1TR002243: Vanderbilt Institute

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

464 for Clinical and Translational Research • Virginia Commonwealth University — UL1TR002649: C.
465 Kenneth and Dianne Wright Center for Clinical and Translational Research • Wake Forest University
466 Health Sciences — UL1TR001420: Wake Forest Clinical and Translational Science Institute •
467 Washington University in St. Louis — UL1TR002345: Institute of Clinical and Translational Sciences •
468 Weill Medical College of Cornell University — UL1TR002384: Weill Cornell Medicine Clinical and
469 Translational Science Center • West Virginia University — U54GM104942: West Virginia Clinical and
470 Translational Science Institute (WVCTSI) □ Submitted: Icahn School of Medicine at Mount Sinai —
471 UL1TR001433: ConduITS Institute for Translational Sciences • The University of Texas Health Science
472 Center at Tyler — UL1TR003167: Center for Clinical and Translational Sciences (CCTS) • University of
473 California, Davis — UL1TR001860: UC Davis Health Clinical and Translational Science Center •
474 University of California, Irvine — UL1TR001414: The UC Irvine Institute for Clinical and Translational
475 Science (ICTS) • University of California, Los Angeles — UL1TR001881: UCLA Clinical Translational
476 Science Institute • University of California, San Diego — UL1TR001442: Altman Clinical and
477 Translational Research Institute • University of California, San Francisco — UL1TR001872: UCSF
478 Clinical and Translational Science Institute □ NYU Langone Health Clinical Science Core, Data Resource
479 Core, and PASC Biorepository Core — OTA-21-015A: Post-Acute Sequelae of SARS-CoV-2 Infection
480 Initiative (RECOVER) □ Pending: Arkansas Children’s Hospital — UL1TR003107: UAMS Translational
481 Research Institute • Baylor College of Medicine — None (Voluntary) • Children’s Hospital of
482 Philadelphia — UL1TR001878: Institute for Translational Medicine and Therapeutics • Cincinnati
483 Children’s Hospital Medical Center — UL1TR001425: Center for Clinical and Translational Science and
484 Training • Emory University — UL1TR002378: Georgia Clinical and Translational Science Alliance •
485 HonorHealth — None (Voluntary) • Loyola University Chicago — UL1TR002389: The Institute for
486 Translational Medicine (ITM) • Medical College of Wisconsin — UL1TR001436: Clinical and
487 Translational Science Institute of Southeast Wisconsin • MedStar Health Research Institute — None
488 (Voluntary) • Georgetown University — UL1TR001409: The Georgetown-Howard Universities Center
489 for Clinical and Translational Science (GHUCCTS) • MetroHealth — None (Voluntary) • Montana State
490 University — U54GM115371: American Indian/Alaska Native CTR • NYU Langone Medical Center —
491 UL1TR001445: Langone Health’s Clinical and Translational Science Institute • Ochsner Medical Center
492 — U54GM104940: Louisiana Clinical and Translational Science (LA CaTS) Center • Regenstrief
493 Institute — UL1TR002529: Indiana Clinical and Translational Science Institute • Sanford Research —
494 None (Voluntary) • Stanford University — UL1TR003142: Spectrum: The Stanford Center for Clinical
495 and Translational Research and Education • The Rockefeller University — UL1TR001866: Center for
496 Clinical and Translational Science • The Scripps Research Institute — UL1TR002550: Scripps Research
497 Translational Institute • University of Florida — UL1TR001427: UF Clinical and Translational Science
498 Institute • University of New Mexico Health Sciences Center — UL1TR001449: University of New
499 Mexico Clinical and Translational Science Center • University of Texas Health Science Center at San
500 Antonio — UL1TR002645: Institute for Integration of Medicine and Science • Yale New Haven Hospital
501 — UL1TR001863: Yale Center for Clinical Investigation

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

502 **References**

- 503 Bennett, T. D., Moffitt, R. A., Hajagos, J. G., Amor, B., Anand, A., Bissell, M. M., Bradwell, K. R.,
504 Bremer, C., Byrd, J. B., Denham, A., DeWitt, P. E., Gabriel, D., Garibaldi, B. T., Girvin, A. T.,
505 Guinney, J., Hill, E. L., Hong, S. S., Jimenez, H., Kavuluru, R., ... National COVID Cohort
506 Collaborative (N3C) Consortium. (2021). Clinical Characterization and Prediction of Clinical
507 Severity of SARS-CoV-2 Infection Among US Adults Using Data From the US National COVID
508 Cohort Collaborative. *JAMA Network Open*, 4(7), e2116901.
509 <https://doi.org/10.1001/jamanetworkopen.2021.16901>
- 510 Blacketer, C. (2021, January 11). *Chapter 4 The Common Data Model*. The Book of OHDSI.
511 <https://ohdsi.github.io/TheBookOfOhdsi/>
- 512 CDC. (2022, June 22). *Nearly One in Five American Adults Who Have Had COVID-19 Still Have “Long*
513 *COVID.”* https://www.cdc.gov/nchs/pressroom/nchs_press_releases/2022/20220622.htm
- 514 CDC. (2023, January 4). *Long COVID - Household Pulse Survey—COVID-19*.
515 <https://www.cdc.gov/nchs/covid19/pulse/long-covid.htm>
- 516 Cohen, J., & van der Meulen Rodgers, Y. (2023). An intersectional analysis of long COVID prevalence.
517 *International Journal for Equity in Health*, 22, 261. <https://doi.org/10.1186/s12939-023-02072-5>
- 518 Crook, H., Raza, S., Nowell, J., Young, M., & Edison, P. (2021). Long covid—Mechanisms, risk factors,
519 and management. *BMJ*, 374, n1648. <https://doi.org/10.1136/bmj.n1648>
- 520 Cutler, D. (2022, July 22). *The Economic Cost of Long COVID: An Update*.
521 https://scholar.harvard.edu/files/cutler/files/long_covid_update_7-22.pdf
- 522 Cutler, D. M., & Summers, L. H. (2020). The COVID-19 Pandemic and the \$16 Trillion Virus. *JAMA*,
523 324(15), 1495–1496. <https://doi.org/10.1001/jama.2020.19759>
- 524 Glasheen, W. P., Cordier, T., Gumpina, R., Haugh, G., Davis, J., & Renda, A. (2019). Charlson
525 Comorbidity Index: ICD-9 Update and ICD-10 Translation. *American Health & Drug Benefits*,
526 12(4), Article 4.

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

- 527 Global Burden of Disease Long COVID Collaborators. (2022). Estimated Global Proportions of
528 Individuals With Persistent Fatigue, Cognitive, and Respiratory Symptom Clusters Following
529 Symptomatic COVID-19 in 2020 and 2021. *JAMA*, 328(16), 1604–1615.
530 <https://doi.org/10.1001/jama.2022.18931>
- 531 Haendel, M. A., Chute, C. G., Bennett, T. D., Eichmann, D. A., Guinney, J., Kibbe, W. A., Payne, P. R.
532 O., Pfaff, E. R., Robinson, P. N., Saltz, J. H., Spratt, H., Suver, C., Wilbanks, J., Wilcox, A. B.,
533 Williams, A. E., Wu, C., Blacketer, C., Bradford, R. L., Cimino, J. J., ... N3C Consortium.
534 (2021). The National COVID Cohort Collaborative (N3C): Rationale, design, infrastructure, and
535 deployment. *Journal of the American Medical Informatics Association: JAMIA*, 28(3), 427–443.
536 <https://doi.org/10.1093/jamia/ocaa196>
- 537 Herman, E., Shih, E., & Cheng, A. (2022). Long COVID: Rapid Evidence Review. *American Family*
538 *Physician*, 106(5), 523–532.
- 539 Kharroubi, S. A., & Diab-El-Harake, M. (2022). Sex-differences in COVID-19 diagnosis, risk factors and
540 disease comorbidities: A large US-based cohort study. *Frontiers in Public Health*, 10, 1029190.
541 <https://doi.org/10.3389/fpubh.2022.1029190>
- 542 Khullar, D., Zhang, Y., Zang, C., Xu, Z., Wang, F., Weiner, M. G., Carton, T. W., Rothman, R. L., Block,
543 J. P., & Kaushal, R. (2023). Racial/Ethnic Disparities in Post-acute Sequelae of SARS-CoV-2
544 Infection in New York: An EHR-Based Cohort Study from the RECOVER Program. *Journal of*
545 *General Internal Medicine*, 38(5), 1127–1136. <https://doi.org/10.1007/s11606-022-07997-1>
- 546 Klaassen, F., Chitwood, M. H., Cohen, T., Pitzer, V. E., Russi, M., Swartwood, N. A., Salomon, J. A., &
547 Menzies, N. A. (2022). *Changes in population immunity against infection and severe disease*
548 *from SARS-CoV-2 Omicron variants in the United States between December 2021 and November*
549 *2022* (p. 2022.11.19.22282525). medRxiv. <https://doi.org/10.1101/2022.11.19.22282525>
- 550 Mancini, D. M., Brunjes, D. L., Lala, A., Trivieri, M. G., Contreras, J. P., & Natelson, B. H. (2021). Use
551 of Cardiopulmonary Stress Testing for Patients With Unexplained Dyspnea Post–Coronavirus
552 Disease. *JACC: Heart Failure*, 9(12), 927–937. <https://doi.org/10.1016/j.jchf.2021.10.002>

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

- 553 Marshall, J. C., Murthy, S., Diaz, J., Adhikari, N. K., Angus, D. C., Arabi, Y. M., Baillie, K., Bauer, M.,
554 Berry, S., Blackwood, B., Bonten, M., Bozza, F., Brunkhorst, F., Cheng, A., Clarke, M., Dat, V.
555 Q., de Jong, M., Denholm, J., Derde, L., ... Zhang, J. (2020). A minimal common outcome
556 measure set for COVID-19 clinical research. *The Lancet Infectious Diseases*, 20(8), e192–e197.
557 [https://doi.org/10.1016/S1473-3099\(20\)30483-7](https://doi.org/10.1016/S1473-3099(20)30483-7)
- 558 Nasserie, T., Hittle, M., & Goodman, S. N. (2021). Assessment of the Frequency and Variety of
559 Persistent Symptoms Among Patients With COVID-19: A Systematic Review. *JAMA Network*
560 *Open*, 4(5), e2111417. <https://doi.org/10.1001/jamanetworkopen.2021.11417>
- 561 Nittas, V., Gao, M., West, E. A., Ballouz, T., Menges, D., Wulf Hanson, S., & Puhan, M. A. (2022). Long
562 COVID Through a Public Health Lens: An Umbrella Review. *Public Health Reviews*, 43,
563 1604501. <https://doi.org/10.3389/phrs.2022.1604501>
- 564 Observational Health Data Sciences and Informatics. (2018). *COVID-19 Clinical Data Warehouse Data*
565 *Dictionary Based on OMOP Common Data Model Specifications Version 5.3*. National Center
566 for Advancing Translational Sciences (NCATS).
567 https://ncats.nih.gov/files/OMOP_CDM_COVID.pdf
- 568 Patient Led Research Collaborative. (2023). *About the Patient-Led Research Collaborative*. Patient Led
569 Research Collaborative – for Long COVID. <https://patientresearchcovid19.com/>
- 570 Pfaff, E. R. (2022, March 11). *Latest Phenotype* [GitHub]. National-COVID-Cohort-Collaborative.
571 [https://github.com/National-COVID-Cohort-](https://github.com/National-COVID-Cohort-Collaborative/Phenotype_Data_Acquisition/wiki/Latest-Phenotype)
572 [Collaborative/Phenotype_Data_Acquisition/wiki/Latest-Phenotype](https://github.com/National-COVID-Cohort-Collaborative/Phenotype_Data_Acquisition/wiki/Latest-Phenotype)
- 573 Sudre, C. H., Murray, B., Varsavsky, T., Graham, M. S., Penfold, R. S., Bowyer, R. C., Pujol, J. C.,
574 Klaser, K., Antonelli, M., Canas, L. S., Molteni, E., Modat, M., Jorge Cardoso, M., May, A.,
575 Ganesh, S., Davies, R., Nguyen, L. H., Drew, D. A., Astley, C. M., ... Steves, C. J. (2021).
576 Attributes and predictors of long COVID. *Nature Medicine*, 27(4), Article 4.
577 <https://doi.org/10.1038/s41591-021-01292-y>

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

578 Tsampasian, V., Elghazaly, H., Chattopadhyay, R., Debski, M., Naing, T. K. P., Garg, P., Clark, A.,
579 Ntatsaki, E., & Vassiliou, V. S. (2023). Risk Factors Associated With Post–COVID-19
580 Condition: A Systematic Review and Meta-analysis. *JAMA Internal Medicine*.
581 <https://doi.org/10.1001/jamainternmed.2023.0750>

582 Webb Hooper, M., Nápoles, A. M., & Pérez-Stable, E. J. (2020). COVID-19 and Racial/Ethnic
583 Disparities. *JAMA*, 323(24), 2466–2467. <https://doi.org/10.1001/jama.2020.8598>

584 World Health Organization. (2021, October 6). *A clinical case definition of post COVID-19 condition by*
585 *a Delphi consensus, 6 October 2021*. [https://www.who.int/publications-detail-redirect/WHO-](https://www.who.int/publications-detail-redirect/WHO-2019-nCoV-Post_COVID-19_condition-Clinical_case_definition-2021.1)
586 [2019-nCoV-Post_COVID-19_condition-Clinical_case_definition-2021.1](https://www.who.int/publications-detail-redirect/WHO-2019-nCoV-Post_COVID-19_condition-Clinical_case_definition-2021.1)

587

588 **Tables**
 589 **Table 1: Analytic Sample**

	Cognitive Symptom Cluster		Fatigue Symptom Cluster		Respiratory Symptom Cluster	
	Positive N (%)	Negative N (%)	Positive N (%)	Negative N (%)	Positive N (%)	Negative N (%)
Total	2,937 (0.2%)	1,231,182 (99.8%)	57,483 (4.7%)	1,176,636 (95.3%)	30,668 (2.5%)	1,203,451 (97.5%)
Age	60.1 (SD: 17.22)	41.4 (SD: 20.72)	44.1 (SD: 19.2)	41.3 (SD: 20.8)	47.0 (SD: 22)	41.3 (SD: 20.68)
Age group:						
0-9	< 20 ‡	79,370 ¶	1,217 (1.5%)	78,161 (98.5%)	2,383 (3%)	76,995 (97%)
10-19	< 45 ‡	125,311 ¶	4,841 (3.9%)	120,516 (96.1%)	1,754 (1.4%)	123,603 (98.6%)
20-29	115 (0.1%)	184,951 (99.9%)	8,904 (4.8%)	176,162 (95.2%)	2,873 (1.6%)	182,193 (98.4%)
30-39	179 (0.1%)	188,841 (99.9%)	9,583 (5.1%)	179,437 (94.9%)	3,409 (1.8%)	185,611 (98.2%)
40-49	451 (0.2%)	194,235 (99.8%)	10,029 (5.2%)	184,657 (94.8%)	4,561 (2.3%)	190,125 (97.7%)
50-59	474 (0.3%)	186,883 (99.7%)	9,173 (4.9%)	178,184 (95.1%)	5,407 (2.9%)	181,950 (97.1%)
60-69	623 (0.4%)	150,740 (99.6%)	7,393 (4.9%)	143,970 (95.1%)	5,297 (3.5%)	146,066 (96.5%)
70-79	681 (0.8%)	87,738 (99.2%)	4,597 (5.2%)	83,822 (94.8%)	3,665 (4.1%)	84,754 (95.9%)
80+	363 (1.1%)	33,110 (98.9%)	1,746 (5.2%)	31,727 (94.8%)	1,319 (3.9%)	32,154 (96.1%)
Sex:						
Female	1,751 (0.3%)	690,796 (99.7%)	39,362 (5.7%)	653,185 (94.3%)	18,164 (2.6%)	674,383 (97.4%)
Male	1,186 (0.2%)	540,035 (99.8%)	18,111‡ (3.3%)	523,112 (96.7%)	12,501‡ (2.3%)	528,718 (97.7%)
Other	0 (0%)	351 (100%)	< 20 ‡	337 ¶	< 20 ‡	352 ¶
Race:						
White	2,163 (0.3%)	845,108 (99.7%)	41,449 (4.9%)	805,822 (95.1%)	21,293 (2.5%)	825,978 (97.5%)
Black or African American	412 (0.2%)	172,825 (99.8%)	8,068 (4.7%)	165,169 (95.3%)	5,186 (3%)	168,051 (97%)
Asian	56 (0.2%)	25,719 (99.8%)	929 (3.6%)	24,846 (96.4%)	619 (2.4%)	25,156 (97.6%)
Native Hawaiian or Pacific Islander	< 20 ‡	1,866 ¶	61 (3.3%)	1,805 (96.7%)	38 (2%)	1,828 (98%)
Other	< 20 ‡	10,393 ¶	470 (4.5%)	9,937 (95.5%)	286 (2.7%)	10,121 (97.3%)
Missing/Unknown	289 (0.2%)	175,274 (99.8%)	6,506 (3.7%)	169,057 (96.3%)	3,246 (1.8%)	172,317 (98.2%)

medRxiv preprint doi: <https://doi.org/10.1101/2024.02.08.24302528>; this version posted May 31, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

Ethnicity:						
Hispanic or Latino	295 (0.2%)	167,080 (99.8%)	7,691 (4.6%)	159,684 (95.4%)	3,664 (2.2%)	163,711 (97.8%)
Not Hispanic or Latino	2,408 (0.3%)	946,774 (99.7%)	45,876 (4.8%)	903,306 (95.2%)	24,966 (2.6%)	924,216 (97.4%)
Ethnicity Missing or Unknown	234 (0.2%)	117,328 (99.8%)	3,916 (3.3%)	113,646 (96.7%)	2,038 (1.7%)	115,524 (98.3%)
COVID-19 vaccination:						
Yes	84 (0.3%)	24,999 (99.7%)	1,034 (4.1%)	24,049 (95.9%)	569 (2.3%)	24,514 (97.7%)
No	2,853 (0.2%)	1,206,183 (99.8%)	56,449 (4.7%)	1,152,587 (95.3%)	30,099 (2.5%)	1,178,937 (97.5%)
Smoking status:						
Current or former smoker	529 (0.4%)	146,835 (99.6%)	9,472 (6.4%)	137,892 (93.6%)	5,486 (3.7%)	141,878 (96.3%)
Non-smoker	2,408 (0.2%)	1,084,347 (99.8%)	48,011 (4.4%)	1,038,744 (95.6%)	25,182 (2.3%)	1,061,573 (97.7%)
Hypertension	29 (0.7%)	4,101 (99.3%)	329 (8.0%)	3,801 (92.0%)	172 (4.2%)	3,958 (95.8%)
Obesity	590 (0.5%)	115,173 (99.5%)	9,581 (8.3%)	106,182 (91.7%)	5,281 (4.6%)	110,482 (95.4%)
Immunocompromised	< 20 ‡	138 ¶	< 20 ‡	129 ¶	< 20 ‡	128 ¶
Pregnant	< 20 ‡	2,488 ¶	216 (8.7%)	2,276 (91.3%)	56 (2.2%)	2,436 (97.8%)
CCI Score	2.7 (SD: 3.11)	0.8 (SD: 1.73)	1.4 (SD: 2.35)	0.7 (SD: 1.7)	1.6 (SD: 2.51)	0.7 (SD: 1.71)

590 ‡ To comply with N3C policy, counts below 20 are displayed as < 20, and in this case, additional values must be skewed by up to 5 to render it
 591 impossible to back-calculate precise counts fewer than 20 for the following categories: Age Group 0-9, Sex Other, Native Hawaiian or Pacific
 592 Islander, Race Other, and Pregnant.

593 ¶ This proportion is one of the two columns that sum up to one. Reporting it would enable the calculation of a cell size < 20. Therefore we mark it
 594 as too small to quantitatively report.

595 Abbreviations: SD: standard deviation; CCI: Charlson Comorbidity Index

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

596 **Table 2: Model Performance Results**

Long COVID Outcome	Model Type	Performance Metric	Train Sample	Test Sample
Respiratory	Logistic	AUC	0.616	0.617
		F1	0.596	0.077
		Accuracy	0.616	0.660
		Precision	0.628	0.041
		Sensitivity	0.568	0.572
		Specificity	0.664	0.662
	Random Forest	AUC	0.625	0.624
		F1	0.601	0.081
		Accuracy	0.625	0.679
		Precision	0.642	0.043
		Sensitivity	0.565	0.566
		Specificity	0.685	0.682
Fatigue	Logistic	AUC	0.604	0.602
		F1	0.607	0.123
		Accuracy	0.604	0.594
		Precision	0.602	0.068
		Sensitivity	0.613	0.611
		Specificity	0.594	0.593
	Random Forest	AUC	0.603	0.599
		F1	0.578	0.127
		Accuracy	0.603	0.656
		Precision	0.617	0.072
		Sensitivity	0.545	0.537
		Specificity	0.662	0.661
Cognitive	Logistic	AUC	0.723	0.734
		F1	0.717	0.013
		Accuracy	0.723	0.743
		Precision	0.734	0.007
		Sensitivity	0.702	0.725
		Specificity	0.745	0.743
	Random Forest	AUC	0.732	0.730
		F1	0.732	0.013
		Accuracy	0.732	0.723
		Precision	0.732	0.006
		Sensitivity	0.732	0.737
		Specificity	0.731	0.723

597 Abbreviations: AUC: Area under the curve

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

598 **Table 3: Coefficient Values for the Binary Logistic Regression Models**

	Coefficient	Cognitive	Fatigue	Respiratory
Age	0-9	-1.367	-1.070	0.651
	10-19	-0.988	-0.178	-0.167
	20-29	-0.331	0.002	-0.085
	30-39	Reference	Reference	Reference
	40-49	0.663	-0.057	0.166
	50-59	0.781	-0.110	0.351
	60-69	0.945	-0.241	0.472
	70-79	1.598	-0.288	0.522
	80+	1.892	-0.287	0.504
Sex	Female	Reference	Reference	Reference
	Male	-0.211	-0.539	-0.176
	Other	-0.117	-0.147	-0.196
Race	White	Reference	Reference	Reference
	Black or African American	-0.054	-0.178	0.084
	Asian	-0.064	-0.274	0.019
	Native Hawaiian or Pacific Islander	-0.430	-0.462	-0.248
	Other	-0.454	0.111	0.052
	Missing/Unknown	-0.254	-0.263	-0.277
Ethnicity	Not Hispanic or Latino	Reference	Reference	Reference
	Hispanic or Latino	0.084	0.063	0.016
	Ethnicity Missing or Unknown	-0.006	-0.215	-0.234
COVID-19 vaccination	No	Reference	Reference	Reference
	Yes	-0.397	-0.406	-0.338
Smoking status:	Non-smoker	Reference	Reference	Reference
	Current or Former smoker	0.345	0.316	0.386
Hypertension	No	Reference	Reference	Reference
	Yes	0.298	0.130	-0.012
Obesity	No	Reference	Reference	Reference
	Yes	0.431	0.429	0.413
Immunocompromised	No	Reference	Reference	Reference
	Yes	0.214	-0.056	0.092
Pregnant	No	Reference	Reference	Reference
	Yes	0.293	0.286	-0.076
CCI Score		0.237	0.160	0.166

599 Abbreviations: CCI: Charlson Comorbidity Index

White Paper: Can longitudinal electronic health record data identify patients at higher risk of developing long COVID?

600 **Table 4: Impurity-Based Variable Importance Scores for the Random Forest Models**

Variable Importance Score		Cognitive	Fatigue	Respiratory
Age	0-9	0.057	0.124	0.014
	10-19	0.092	0.008	0.061
	20-29	0.064	0.006	0.060
	30-39	Reference	Reference	Reference
	40-49	0.003	0.003	0.003
	50-59	0.007	0.002	0.005
	60-69	0.019	0.002	0.033
	70-79	0.124	0.003	0.044
	80+	0.078	0.002	0.012
Sex	Female	Reference	Reference	Reference
	Male	0.007	0.214	0.008
	Other	0.000	0.000	0.000
Race	White	Reference	Reference	Reference
	Black or African American	0.004	0.003	0.007
	Asian	0.003	0.003	0.001
	Native Hawaiian or Pacific Islander	0.000	0.001	0.001
	Other	0.002	0.001	0.001
	Missing/Unknown	0.010	0.023	0.031
Ethnicity	Not Hispanic or Latino	Reference	Reference	Reference
	Hispanic or Latino	0.006	0.004	0.004
	Ethnicity Missing or Unknown	0.005	0.027	0.019
COVID-19 vaccination	No	Reference	Reference	Reference
	Yes	0.003	0.003	0.003
Smoking status:	Non-smoker	Reference	Reference	Reference
	Current or Former smoker	0.018	0.048	0.061
Hypertension	No	Reference	Reference	Reference
	Yes	0.003	0.002	0.002
Obesity	No	Reference	Reference	Reference
	Yes	0.045	0.128	0.122
Immunocompromised	No	Reference	Reference	Reference
	Yes	0.000	0.000	0.000
Pregnant	No	Reference	Reference	Reference
	Yes	0.001	0.001	0.001
CCI Score		0.445	0.394	0.507

601 Abbreviations: CCI: Charlson Comorbidity Index