

1 **Predictive Modeling of COVID-19 Case Growth Highlights Evolving Demographic**
2 **Risk Factors in Tennessee and Georgia**

3

4 Jamieson D. Gray, BS¹, Coleman R. Harris, BS^{1,2}, Lukasz S. Wylezinski, PhD^{1,3}, and
5 Charles F. Spurlock, III, PhD^{1,3,†}

6 ¹Decode Health, Nashville, TN, 37203

7 ²Department of Biostatistics, Vanderbilt University School of Medicine, Nashville, TN,
8 37232

9 ³Department of Medicine, Vanderbilt University School of Medicine, Nashville, TN,
10 37232

11 † Address correspondence to CFS: Charles F. Spurlock, III, PhD, 111 10th Ave South,
12 Suite 102, Nashville, TN, USA, chase@decodehealth.ai

13

14 **Keywords:** SARS-CoV-2, social determinants of health, machine learning,
15 demographics

16

17

18 **Manuscript Word Count:** 791

19 **Pages:** 9

20 **Abstract**

21 The COVID-19 pandemic has exposed the need to understand the unique risk
22 drivers that contribute to uneven morbidity and mortality in US communities. Addressing
23 the community-specific social determinants of health that correlate with spread of
24 SARS-CoV-2 provides an opportunity for targeted public health intervention to promote
25 greater resilience to viral respiratory infections in the future.

26 Our work combined publicly available COVID-19 statistics with county-level
27 social determinants of health information. Machine learning models were trained to
28 predict COVID-19 case growth and understand the unique social, physical and
29 environmental risk factors associated with higher rates of SARS-CoV-2 infection in
30 Tennessee and Georgia counties. Model accuracy was assessed comparing predicted
31 case counts to actual positive case counts in each county. The predictive models
32 achieved a mean r-squared (R^2) of 0.998 in both states with accuracy above 90% for all
33 time points examined. Using these models, we tracked the social determinants of
34 health, with a specific focus on demographics, that were strongly associated with
35 COVID-19 case growth in Tennessee and Georgia counties. The demographic results
36 point to dynamic racial trends in both states over time and varying, localized patterns of
37 risk among counties within the same state.

38 Identifying the specific risk factors tied to COVID-19 case growth can assist
39 public health officials and policymakers target regional interventions to mitigate the
40 burden of future outbreaks and minimize long-term consequences including emergence
41 or exacerbation of chronic diseases that are a direct consequence of infection.

42 **Introduction**

43 In January 2021, Tennessee and Georgia reported over 1,550,000 cases and
44 22,100 deaths due to COVID-19. Hispanic individuals comprise 14% of the states'
45 population but represent 25% of confirmed cases, suggesting race and ethnicity are
46 associated with case growth.¹

47 Combining publicly available COVID-19 data and proprietary social determinants
48 of health (SDOH), which measure certain physical, social, economic, and demographic
49 characteristics, we built and tuned machine learning models to predict COVID-19 case
50 growth in Tennessee and Georgia. We sought to accurately predict COVID-19 case
51 growth and investigate the changing significance of demographic features influencing
52 these predictions. Our approach produced highly accurate forecasts of COVID-19 case
53 growth in both states while uncovering evolving patterns of specific demographic factor
54 importance during a seven-month period. This approach also yielded state- and county-
55 level insights that can inform targeted mitigation efforts to slow respiratory virus spread.

56

57 **Methods**

58 Our approach combined publicly available COVID-19 case, hospitalization and
59 death metrics with county-specific SDOH data.^{2,3} Feature engineering and feature
60 selection were employed to define the data inputs that best represent changes in
61 COVID-19 case growth over time. We lagged (offset case growth over time), windowed
62 (summed or averaged case growth over time), and developed novel time window
63 features (i.e., “days since the 100th COVID-19 case”) using state health department

64 data. SDOH enrichment data, including demographic information, was appended to the
65 engineered features for each county.⁴ The target for predictive modeling was defined as
66 the future relative case growth normalized to the population in Tennessee and Georgia
67 counties from July 2020-January 2021. A grid search of generalized linear and tree-
68 based machine learning models was performed. Briefly, we trained and tested each
69 model using four to six weeks of historical COVID-19 case data and made predictions
70 using the most recent data available. From the ~50 regression models that we built for
71 each timepoint, models were chosen in a survival of the fittest approach comparing
72 statistical and real-world accuracy for predicting COVID-19 case growth.⁵ We identified
73 the top third of each state's counties at highest risk for case growth and assessed our
74 prediction accuracy versus actual case growth over time. Finally, we analyzed each
75 feature's impact at the state- and county-level to understand the demographic features
76 that drove COVID-19 case growth.

77

78 **Results**

79 Candidate models for Tennessee and Georgia achieved excellent metrics across
80 all timepoints including a mean R^2 value of 0.998 (TN and GA), mean Tweedie deviance
81 of 0.003 (TN) and 0.002 (GA), as well as a mean absolute error (MAE) of 0.357 (TN)
82 and 0.337 (GA) (Supplementary Figure 1A). Prediction accuracy was >90% in all
83 models across both states when compared to actual future case growth (Supplementary
84 Figure 1B).

85 Demographics produced variable trends at both the state- and county-level. The
86 two most populous counties in Tennessee, Shelby and Davidson, revealed an identical
87 pattern of importance for Native American demographics in determining future case
88 growth while exhibiting differences among the Asian demographic. Shelby County
89 displayed a gradual increase in importance in the Asian demographic while Davidson
90 County saw a more pronounced spike between October and November. Comparing
91 demographic importance at the Tennessee state-level versus individual counties yields
92 similar patterns (Non-Hispanic White) as well as contrasting trends (African American).
93 Further, Tennessee's stable Hispanic demographic trend differed from the individual
94 counties' more acute fluctuation of importance (Figure 1A).

95 Additionally, similarities and differences in demographic trends extend across
96 state borders. While the Hispanic demographic displayed the most meaningful
97 importance in Tennessee during July and August, Georgia saw a similar increase in
98 importance starting in September. Comparison of the two states' top demographic
99 drivers showed a potential macro-pattern in which the most important driver for one
100 state often preceded its rise to top importance in the other (Figure 1A and Figure 1B).

101

102 **Discussion**

103 This analysis of community-specific relationships among SDOH and COVID-19
104 case growth in Tennessee and Georgia discovered localized, evolving patterns of risk,
105 highlighting the quantitative differences in state- and county-level case growth, and the
106 qualitative differences in important demographic factors that influence spread of

107 infection. These patterns can shift dramatically month-to-month, increasing or
108 decreasing over time and vary significantly by geography; even among similarly sized
109 counties within a state or between two neighboring states.

110 Identifying the specific risk drivers across the country during a pandemic can
111 assist decision-makers in protecting especially vulnerable populations through targeted
112 interventions. Closing the loop to address these risk factors can also enhance
113 community resilience to future viral respiratory infections.⁶

114 Applications of this approach extend beyond acute respiratory infection to chronic
115 disease outcomes including those that are a consequence of COVID-19. A growing
116 percentage (>10%) of patients infected with SARS-CoV-2 develop long-COVID.⁷ These
117 patients experience prolonged, debilitating symptoms months after infection and
118 emergence or exacerbation of chronic illness. Thus, targeted approaches to mitigate
119 spread of disease can lessen future acute and chronic disease burden.

120

121

122

123

124

125

126

127

128 **Acknowledgements**

129 This work was supported by Decode Health, Inc. and grants from the National Institutes
130 of Health (AI124766, AI129147 and AI145505). Dr. Spurlock had full access to all of the
131 data in the study and takes responsibility for the integrity of the data and the accuracy of
132 the data analysis. Dr. Spurlock devised the concept and study design. All authors took
133 part in acquisition, analysis and interpretation of the data along with drafting and
134 revising the manuscript.

135

136 **Conflict of interest statement**

137 Authors Gray, Wylezinski and Spurlock are shareholders in Decode Health, Inc.
138 (Nashville, TN). Decode Health develops artificial intelligence approaches to predict
139 chronic and infectious disease risk in patient populations.

140

141

142

143

144

145

146

147

148 References

- 149 1. The COVID Tracking Project. Racial Data Dashboard.
150 <https://covidtracking.com/race/dashboard>. Accessed November 25, 2020,
- 151 2. Johns Hopkins University Coronavirus Resource Center. COVID-19 United States
152 Cases by County. <https://coronavirus.jhu.edu/us-map>. Accessed November 25, 2020,
- 153 3. Vest JR, Ben-Assuli O. Prediction of emergency department revisits using area-level
154 social determinants of health measures and health information exchange information. *Int J Med*
155 *Inform.* 09 2019;129:205-210. doi:10.1016/j.ijmedinf.2019.06.013
- 156 4. Kolak M, Bhatt J, Park YH, Padrón NA, Molefe A. Quantification of Neighborhood-Level
157 Social Determinants of Health in the Continental United States. *JAMA Netw Open.* Jan
158 2020;3(1):e1919928. doi:10.1001/jamanetworkopen.2019.19928
- 159 5. Muhlestein WE, Akagi DS, Chotai S, Chambless LB. The impact of presurgical
160 comorbidities on discharge disposition and length of hospitalization following craniotomy for
161 brain tumor. *Surg Neurol Int.* 2017;8:220. doi:10.4103/sni.sni_54_17
- 162 6. Graves E, Weiss A, Rickles M, Colyer E. Community Resiliency to COVID-19 in a
163 Subset of US Communities. [https://about.sharecare.com/wp-](https://about.sharecare.com/wp-content/uploads/2020/08/National_COVID-whitepaper_PROOF_08.28.20.pdf)
164 [content/uploads/2020/08/National_COVID-whitepaper_PROOF_08.28.20.pdf](https://about.sharecare.com/wp-content/uploads/2020/08/National_COVID-whitepaper_PROOF_08.28.20.pdf). Accessed
165 November 25, 2020,
- 166 7. Greenhalgh T, Knight M, A'Court C, Buxton M, Husain L. Management of post-acute
167 covid-19 in primary care. *BMJ.* 08 2020;370:m3026. doi:10.1136/bmj.m3026

168

169

170

171

172

173

174

175

176

177

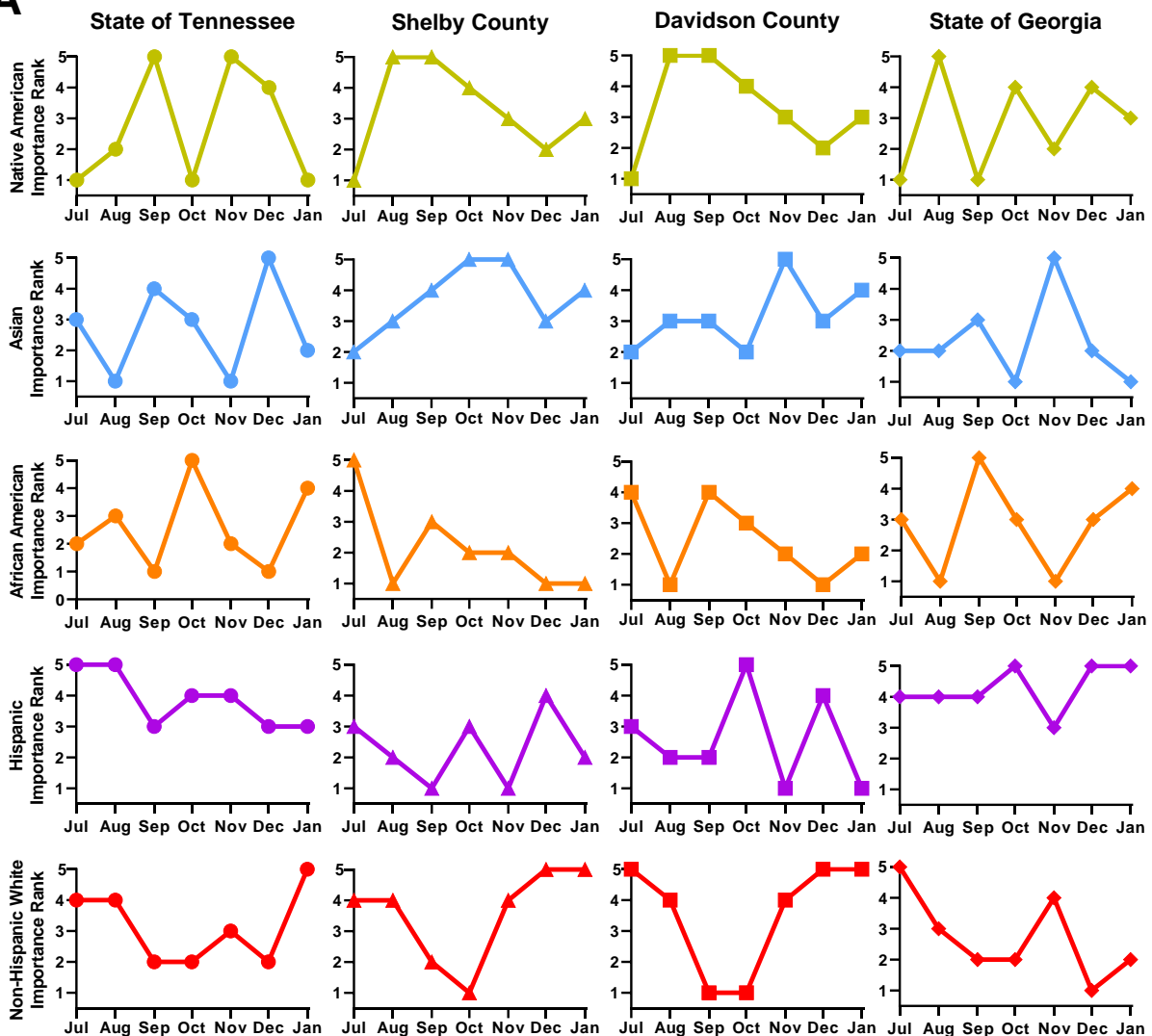
178

179 **Figure Title/Legend**

180 **Figure 1: Influence of demographic features linked to COVID-19 case growth**
181 **exhibit dynamic shifts over time in Tennessee and Georgia. (A)** Relative rank of
182 demographic feature importance across top predictive models are reported for the entire
183 state of Tennessee (●) and the two most populous counties in Tennessee, Shelby
184 County (▲) and Davidson County (■) as well as the state of Georgia (◆). A score of 5
185 on the importance rank indicates the most important demographic feature relative to the
186 other four demographic features. Groups include Native American (●), Asian (●),
187 African American (●), Hispanic (●), and Non-Hispanic White (●). **(B)** Differences in the
188 rank of demographic feature importance in Tennessee and Georgia over time. The color
189 of the bubble (TN ●; GA ●) indicates the state that exhibited a higher importance rank
190 of the specific demographic feature for predicting COVID-19 case growth. Black dots
191 (●) designate months where the two states displayed the same importance rank for an
192 individual demographic feature. The size of the bubbles shows the difference in
193 importance of each demographic feature between the two states. Larger bubbles
194 connote greater difference in importance.

195

A



B

