

1 Prediction of COVID-19 infection risk using personal mobile 2 location data only

3

4

5 Ahreum Jang^{1*}, Sungtae Kim^{1¶}, Hyeongwoo Baek^{1&}, Hyejung Kim^{1%}, Hae-Lee Park^{1%,#a}

6

7 ¹ AI/DX Convergence Business Group, KT

8

9

10 ^{#a}Current Address: KT, 209, Jamsil-ro, Songpa-gu, Seoul, Republic of Korea

11

12

13 * Corresponding author

14 E-mail: ar.jang@kt.com (AJ)

15

16 * Conceptualization, Project administration, Resources, Supervision, Data curation, Formal analysis,

17 Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing –

18 review & editing

19 ¶ Conceptualization, Data curation, Investigation, Writing – original draft, Writing – review & editing

20 & Conceptualization, Data curation, Formal analysis, Methodology, Software, Validation, Visualization,

21 Writing – original draft, Writing – review & editing

22 %Conceptualization, Data curation, Writing – review & editing

23 **Abstract**

24 Predicting an individual's risk of infectious disease is a critical technology in infectious
25 disease response. During the COVID-19 pandemic, identifying and isolating individuals at high risk of
26 infection was an essential task for epidemic control. We introduce a new machine learning model that
27 predicts the risk of COVID-19 infection using only individuals' mobile cell tower location information.
28 This model distinguishes the cell tower location information of an individual into residential and non-
29 residential areas and calculates whether the cell tower locations overlapped with other individuals. It
30 then generates various variables from the information of overlapping and predicts the possibility of
31 COVID-19 infection using a machine learning algorithm. The predictive model we developed showed
32 performance comparable to models using individual's clinical information. This predictive model, which
33 can be used to predict infections of diseases with asymptomatic infections such as COVID-19, has the
34 advantage of supplementing the limitations of existing infectious disease prediction models that use
35 symptoms and other information.

36 **Introduction**

37 Beginning from early 2020, the COVID-19 pandemic caused significant human damage
38 worldwide, and in the absence of treatments and vaccines for the newly emerged infectious disease, each
39 country had no choice but to focus on epidemic prevention through non-pharmaceutical interventions[1].
40 Governments of each country implemented measures to reduce individual infection risk (wearing masks,
41 social distancing, quarantine, lockdown, etc.) to prevent the spread of infectious diseases, and it was
42 necessary to quickly lead individuals at high risk of infection to testing and treatment[2]. Implementing
43 preventive measures for individuals at high risk of infection was one of the important response strategies
44 in the early stages of the COVID-19 pandemic[3].

45 There are two typical ways to identify individuals at high risk of COVID-19 infection. One
46 can identify high-risk individuals through the physical symptoms that appear when infected with
47 COVID-19 and through the potential contact with an infected person. Firstly, screening high-risk
48 individuals through several factors such as fever and respiratory symptoms, the main symptoms of
49 COVID-19. Various machine learning and AI research to predict infection risk through these symptoms
50 have been conducted, and there were models that showed a performance of up to 97.79% Accuracy[4–
51 6]. However, this method has the disadvantage of missing asymptomatic infected individuals.

52 Another way to identify high-risk individuals is to trace those who have come into contact
53 with an infected person. There have been attempts to utilize IT technology for this. Contact tracing
54 mobile applications that track infected individuals and confirm whether they have come into contact
55 with them were developed. These applications, which utilize widely used modern mobile devices and
56 technologies such as Bluetooth and GPS (Global Positioning System), provided important information
57 that could confirm individual locations and determine whether they had contact with infected people.
58 Several open-source technologies emerged, and two companies providing mobile OS (Operation
59 System), Apple and Google, even added technology to check and manage whether they came into
60 contact with a COVID-19 infected individual in the mobile OS, and these digital contact tracing
61 technologies are known to have been used in more than 46 countries[7,8]. These technologies are useful
62 for both governments and individuals, but there is a limitation. These mobile applications must be
63 installed and used by individuals themselves. In countries where it is not mandatory by the government,
64 the actual usage rate of these applications was very low[9]. It's a technology that is hard to see effects if
65 there are few users.

66 South Korea took a different approach. Instead of contact tracing using mobile applications,
67 they collected mobile cell tower location data from mobile carriers to find places visited by infected

68 individuals and those who came into contact with them. It was possible due to South Korea's high
69 mobile usage rate and support of laws and systems. When a COVID-19 confirmed case occurred in
70 South Korea, the government collected the cell tower location data for the 14 days prior to the COVID-19
71 confirmed patient's PCR (Polymerase Chain Reaction) test date. This information was used as validation
72 for epidemiological investigations conducted through interviews with the confirmed cases. Also, when a
73 mass infection occurred with many confirmed cases, people who overlapped with the confirmed case at
74 the cell tower location were deemed at risk of infection and were allowed to get tested. Individuals could
75 not use this information directly and could get a COVID-19 PCR test if they were considered at high
76 risk of infection based on information such as the visit places of infected individuals announced by the
77 government[10]. Generally, it is not possible to know the exact location of an individual or whether they
78 had contact with a specific person using cell tower location information[8]. Nevertheless, in South
79 Korea, cell tower location information was used in epidemiological investigations of infected
80 individuals and in analyzing infection hotspots[11].

81 The aim of this study is to develop a machine learning model that predicts an individual's risk
82 of COVID-19 infection using only cell tower location information. We conducted a study using
83 COVID-19 test results and cell tower location information collected from a mobile application during
84 the period of significant COVID-19 outbreak in South Korea. While it's not possible to determine an
85 individual's exact location and whether they had contact with an infected individual with cell tower
86 location information, we obtained results that we can predict an individual's risk of COVID-19 infection
87 using machine learning techniques. If we can address a few constraints discussed in the conclusion of
88 this study, it can be used as a technology to prevent the spread of a pandemic by complementing the
89 disadvantages of existing methods of finding high-risk infections in an infectious disease pandemic
90 situation.

91 **Methods**

92 **Data**

93 All data were obtained from the SHINE mobile application. SHINE is an application that
94 provides location-based COVID-19 outbreak information and COVID-19 coping strategies based on
95 user-recorded information such as location, gender, age, COVID-19 test results, and vaccination data,
96 etc. The app was launched on October 13, 2021, and operated until March 31, 2023, and was available
97 for use through Apple's AppStore and APK file installation(Android only) for individuals aged 14 and
98 above (Fig 1). Upon registration on the mobile app, users consented online to the use of their de-
99 identified device GPS data and personal information entered in the app for research purposes. Moreover,
100 users who were subscribers of KT (Korea Telecom) also provided online consent to extract and utilize
101 cell tower location data from KT's network data infrastructure. Location data was stored only for the 14
102 days preceding the date when users recorded information such as PCR test results in the app. The
103 COVID-19 PCR test results were uploaded directly to the application by the users themselves, and the
104 application service administrators filtered out inaccurate information by comparing all uploaded test
105 results – verified by documents issued by hospitals or screenshots of text messages with test results
106 featuring the individual's name – with the personal information entered during registration.
107 Consequently, from the app's launch to June 30, 2022, there were 43,270 total users, 21,046 users who
108 registered PCR test results, and there were 17,678,028 cases of mobile cell tower location logs.

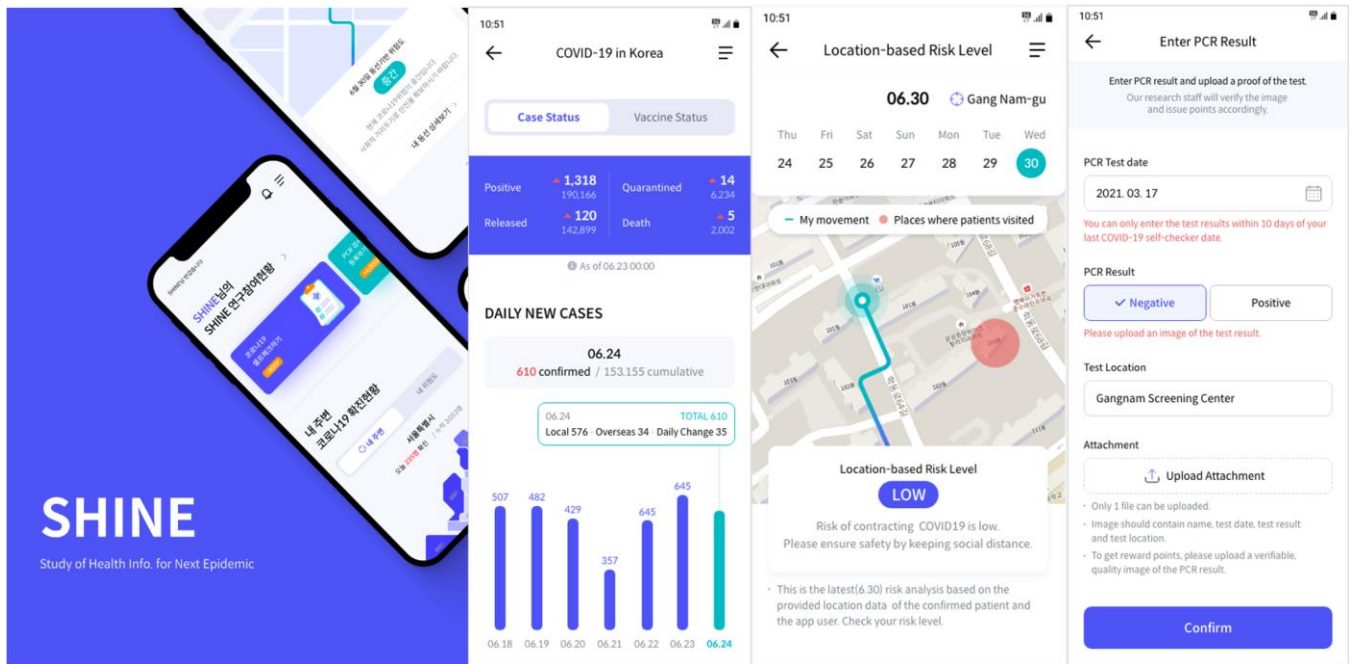


Fig 1. Screenshots of the SHINE mobile application. The language in the image has been translated into English.

This research received exemption from review by the Institutional Review Board of Sungkyunkwan University on November 11, 2022(IRB Number: SKKU 2022-11-014). Subsequently, we accessed location data, PCR test information, gender, and age data, which were collected retrospectively only from users who had consented for research purposes, from October 13, 2021, to June 30, 2022. Our access to this data began on November 28, 2022, and it ensuring all data were non-identifiable at the personal level.

The personal location information used in this study consists of mobile cell tower location data. This is data from regular communication between the subscriber's device and the cell tower, containing information like timestamp, GPS location of the cell tower, and a unique identifier for subscribers[8]. This data has two limitations. The first is that the GPS location of the cell tower is not the precise location of the user but only indicating that the user is within the service radius of the cell

124 tower. The second is that only information from KT subscribers can be used, and KT subscribers
125 represent 31.3% of all mobile users in Korea[12]. These limitations require necessary assumptions and
126 interpretations, which are discussed in the Discussion section.

127 On the other hand, the GPS information from the device collected in the mobile application
128 was not used. Although the SHINE mobile application was set to collect the device's GPS records in the
129 mobile OS's background, there were difficulties in continuously collecting reliable location records.
130 Users often denied providing location information due to battery consumption from running the
131 application in the background. Moreover, the mobile OS regularly sent location access approval or
132 denial notifications to users about applications requesting device location information, leading to limited
133 data tracking users' locations continuously. However, we were able to obtain mobile telecom cell tower
134 location data irrespective of the user's mobile application usage, enabling us to know the reliable
135 location of the user.

136 We extracted only KT users who could use cell tower location data among those who
137 uploaded COVID-19 PCR test results on the mobile app. We further narrowed down the data to include
138 only Seoul residents from January 1, 2022, to June 30, 2022. The reason for these limitations is that this
139 period saw the largest outbreak of COVID-19 in Korea[13], and we had the most COVID-19 PCR test
140 results data, and Seoul is the city with the highest population density in Korea[14]. There was a need to
141 limit the region to areas with high population density to observe overlapping individual locations in the
142 data. After considering the incubation period of COVID-19 and PCR test dates, we used PCR test results
143 and cell tower location data of 837 individuals whose cell tower location data was collected for the
144 seven days before the PCR test for modeling. The demographic information of these 837 individuals is
145 as in Table 1.

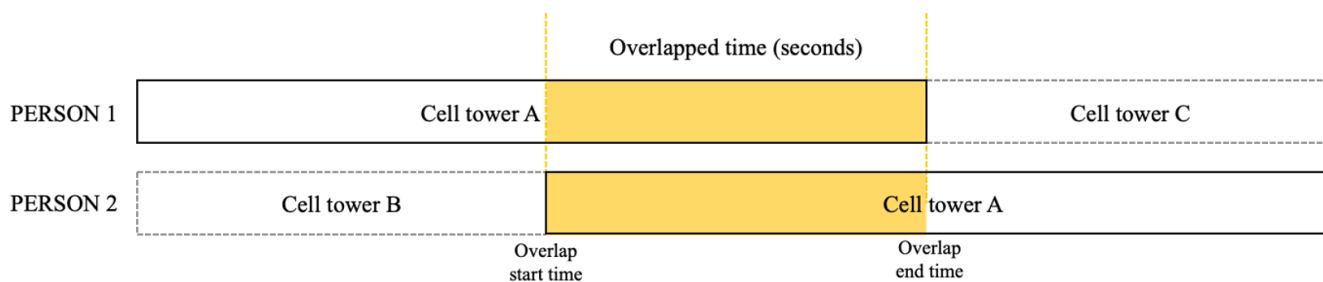
146 **Table 1. Demographic info. of Data**

		COVID-19 PCR test Result	
		Positive	Negative
Total, n(%)		703 (84%)	134 (16%)
Sex, n(%)	Male	225 (26.9%)	43 (5.1%)
	Female	478 (57.1%)	91 (10.9%)
Age, n(%)	10-19	40 (4.8%)	0 (0.0%)
	20-29	232 (27.7%)	39 (4.4%)
	30-39	230 (27.5%)	43 (5.1%)
	40-49	116 (13.9%)	35 (4.2%)
	50-59	57 (6.8%)	14 (1.7%)
	60-69	20 (2.4%)	3 (0.4%)
	70-79	7 (0.8%)	2 (0.2%)
	80-89	1 (0.1%)	0 (0.0%)

147 **Data preprocessing**

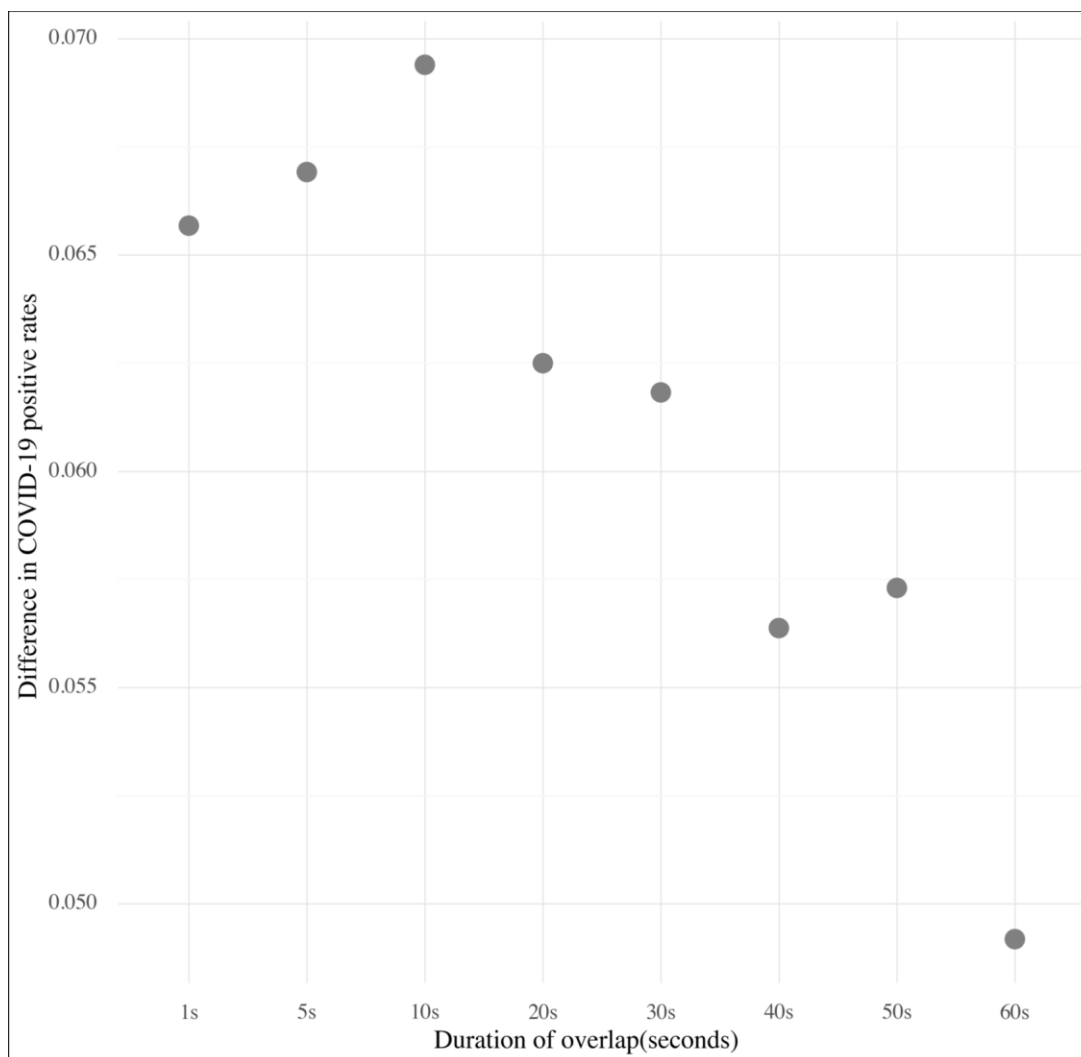
148 First, we distinguished whether the cell tower location from each individual's seven-day
149 location data was a residential area or a non-residential area. We designated the cell tower location that
150 appeared most frequently from 10 pm to 7 am the next day for each individual as the residential cell
151 tower location. We then added a variable to each individual's cell tower location record to distinguish
152 whether it was a residential or non-residential area, anticipating that the infection characteristics would
153 vary depending on this classification. We also extracted cell tower location records only from 9 AM to
154 10 PM. To assess the possibility of contact with others, as shown in Fig 2, for each individual's cell
155 tower location record, we checked the overlapping time with others at the same cell tower. If the overlap
156 was more than 10 seconds, we marked it as an overlapping. The overlap criterion of 10 seconds was
157 selected as it provided the largest difference in COVID-19 infection rates between user groups with at
158 least one overlapping record and those without any overlapping records (Fig 3). Furthermore, in
159 checking the overlap time at the same cell tower with others, we did not refer to the others' PCR results
160 information. This is because the others' COVID-19 infection status cannot be confirmed when
161 predicting an individual's COVID-19 infection. The overlapping time information at cell tower locations

162 for each individual over seven days was summarized into six types of information by distinguishing the
163 type of cell tower location (residential or outing area), and the characteristics of all predictor variables
164 used for modeling are as in Table 2. And all predictor variables were transformed using the natural
165 logarithm to reduce data skewness.



166

167 **Fig 2. Method for calculating overlaps.** In each individual's location records, it was counted as an
168 overlap if they were located at the same cell tower as another person for 1 second or more.



169

170 **Fig 3. Difference in COVID-19 positive rates based on the duration of overlap (in seconds).**
 171 Difference in COVID-19 positive rates by comparing user groups with overlapping records and those
 172 without, based on each duration of overlap.

173 **Table 2. Characteristics of features**

Variables		n	Mean	SD ^a
Number of Overlaps	Outside	837	38.62	±126.01
	Residence	837	15.01	±63.71
Number of Overlapped people	Outside	837	1.4	±1.46
	Residence	837	0.37	±0.64
Total Overlapped Time ^b (seconds)	Outside	837	888.98	±2956.47
	Residence	837	302.01	±1353.17
Max Overlapped Time ^c (seconds)	Outside	837	75.59	±199.36
	Residence	837	25.32	±77.47
Average Overlapped Time ^d (seconds)	Outside	837	0.33	±0.86
	Residence	837	0.27	±0.91

Variables		n	Mean	SD ^a
Minimum Overlapped Time ^e (seconds)	Outside	837	9.16	±12.59
	Residence	837	3.64	±11.13
Number of overlapped locations(outside)		837	4.17	±11.65

174 ^aSD = Standard Deviation

175 ^bTotal Overlapped Time = Sum of all time overlapped with others over 7 days.

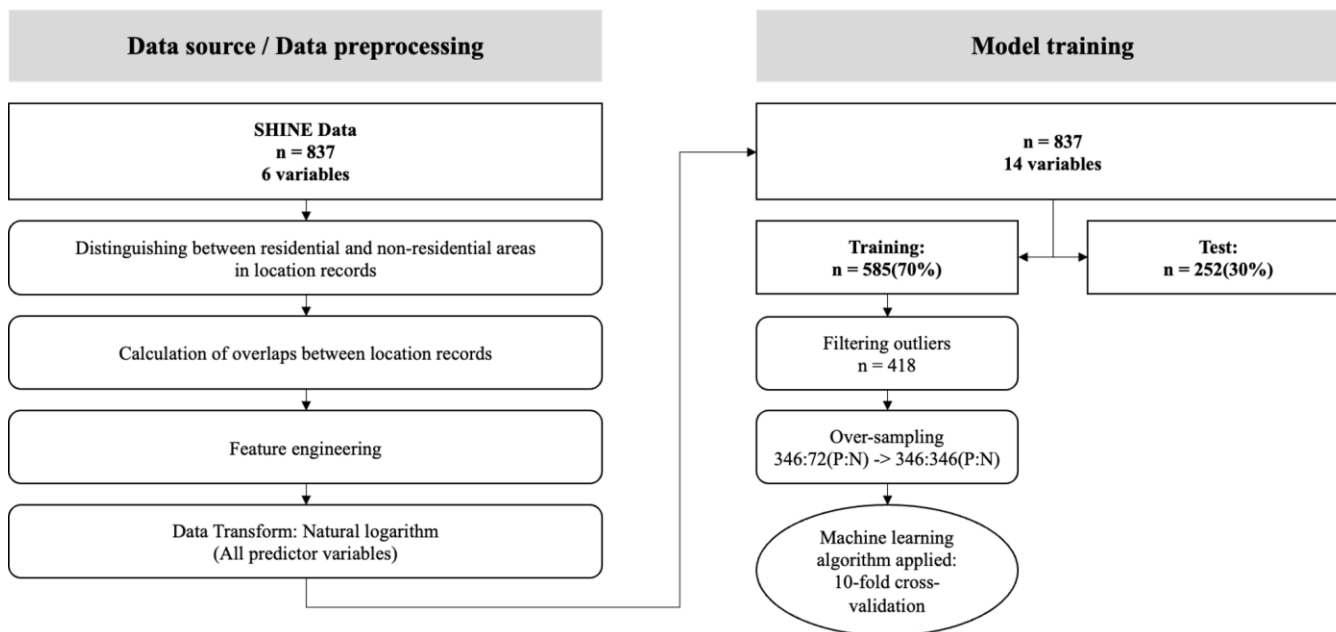
176 ^cMax Overlapped Time = Max duration of time overlapped with others in a single event over 7 days.

177 ^dAverage Overlapped Time = Average duration of time overlapped with others per event over 7 days.

178 ^eMinimum Overlapped Time = Minimum duration of time overlapped with others per event over 7 days.

179 **Development of model for prediction**

180 To train the model, the data was divided into a Training data set (585/837, 70%) and Test data
181 set (252/837, 30%). It was ensured that the ratio of COVID-19 PCR test results was maintained during
182 this division. To enhance the model's generalization performance, records that corresponded to outliers
183 were removed from the Training data set prior to model training. Outliers were determined by
184 calculating the IQR (Interquartile Range) for each variable; values smaller than 1.5 times the first
185 quartile or larger than 1.5 times the third quartile were considered as outliers. Subsequently, data with
186 negative COVID-19 PCR test results was subjected to Random over-sampling, so the ratio of positive to
187 negative results for the target variable, the COVID-19 PCR test result, became 1:1. We experimented
188 with machine learning models like Logistic Regression, a binary classification model, XGBoost
189 Classifier, and Random Forest Classifier. These models are frequently used for classification and were
190 used in this study to develop a COVID-19 infection prediction model. Python (version 3.8.10)'s scikit-
191 learn (version 1.0.2) library's GridSearchCV was used to train the Training data set using 10-Fold Cross
192 Validation. The entire process from data preprocessing to machine learning model training is as
193 described in Fig 4.



194

195 **Fig 4. Workflow of analysis**

196 Results

197 Table 3 shows the performance results measured in the Test data set with five metrics for the
198 three trained models. Logistic regression exhibited the lowest performance across all metrics, while
199 XGBoost showed the highest Accuracy and Sensitivity. In contrast, the Random Forest model
200 demonstrated high Specificity and Precision, outperforming both the Logistic Regression and XGBoost
201 models in terms of Specificity. There was no significant difference in AUC(Area Under the Curve)
202 performance between Random Forest and XGBoost, but among the models tested, Random Forest had
203 the best AUC. The ROCs (Receiver Operating Characteristic curves) for the three models are shown in
204 Fig 5.

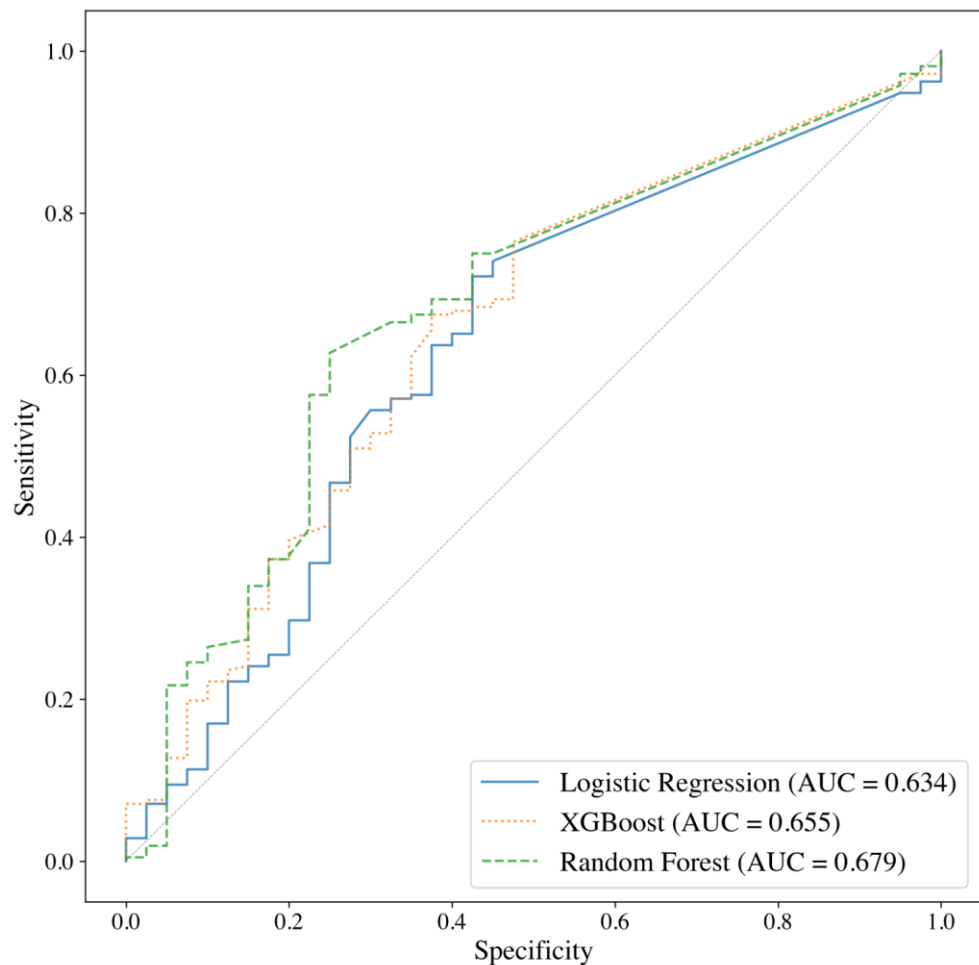
205 The rank of importance of predictor variables in the Random Forest model, which had the
206 highest AUC, is shown in Fig 6. The variable with the highest importance was the Total Overlapped
207 Time in non-residential areas. Also, variables related to non-residential areas had higher importance than

208 those related to residential areas. This finding suggests that activities conducted outside residential areas
209 can provide significant insights for predicting COVID-19 infection.

210 Table 4 summarizes the algorithms, data, sample sizes, and performance of existing studies
211 related to COVID-19 infection prediction. The six studies we reviewed all primarily utilized individual
212 demographics and clinical information, and performance ranged from 0.689 to 0.98 based on AUC
213 (limited to studies that disclosed AUC). Compared to these studies, our developed model's performance
214 was relatively low. However, despite our model utilizing only individual location records for
215 predictions, its AUC did not significantly differ from that of models incorporating individual symptoms
216 [15]. When comparing the results of our Random Forest model with the results of this study, our model
217 showed lower AUC and Sensitivity, but higher Specificity and Precision.

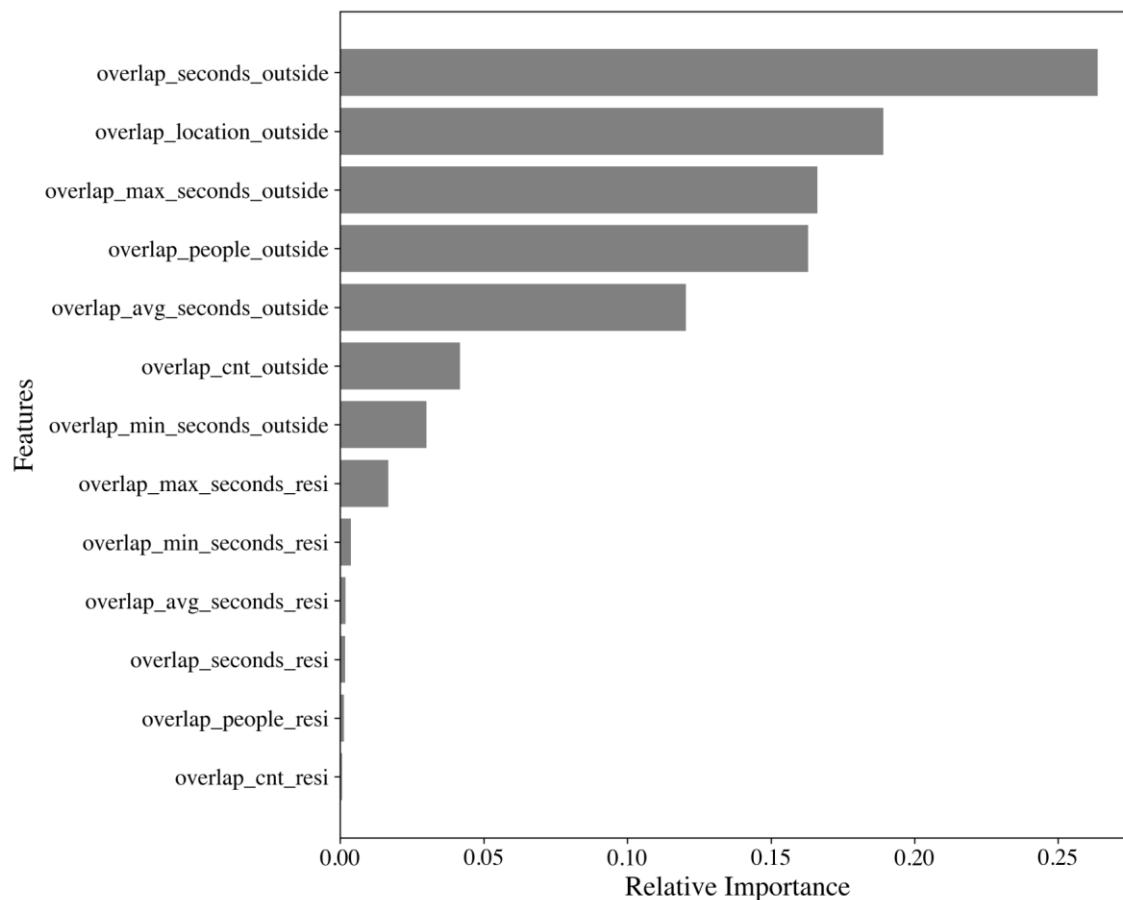
218 **Table 3. Performance of machine-learning algorithms and logistic regression.**

Metrics	Logistic Regression	XGBoost	Random Forest
Accuracy	0.635	0.690	0.647
Sensitivity	0.637	0.722	0.627
Specificity	0.625	0.525	0.750
Precision	0.879	0.876	0.897
AUC	0.634	0.655	0.679



219

220 **Fig 5. Receiver operating characteristic curves with corresponding AUC values.** AUC values for
221 each model are also presented in Table 2.



222

223 **Fig 6. Variable importance plots of COVID-19 PCR result predictors for Random forest.**

224 **Table 4. Comparison with the Prediction model using symptoms**

Reference	ML/AI methods	Types of Data	Sample	Performance
[16]	Support Vector Machine	Clinical, laboratory features, Demographics	556	Accuracy: 77.5% Specificity:78.4% AUC: 0.98(testing dataset)
[17]	Random forest	Clinical, Demographics	253	Accuracy: 95.95% Specificity: 96.95%
[18]	Logistic regression, Random Forest, Decision tree, Linear SVM, Naive Bayes, Gradient boosting classifier	Clinical, Demographics	5,434	Accuracy: 97.79%, Sensitivity: 0.99, Precision 0.97
[15]	Logistic Regression	Clinical	378	AUC: 0.6891, Specificity: 58.6%, Sensitivity: 64.7%, Precision: 43.1%
[19]	Random Forest	Clinical, Demographics	1,653	AUC: 0.788, recall: 0.799, FPR: 0.38
[20]	Logistic Regression (LASSO)	Clinical, Demographics	143,531	AUC: 0.78

225 **Discussion**

226 We have developed a prediction model that addresses the limitations of existing methods in
227 predicting high-risk individuals for COVID-19 infection. We did this by creating a variable that
228 represents the possibility of contact with others using only an individual's cell tower location
229 information, and then developing an infection risk prediction model using machine learning algorithms.
230 Though our model showed lower performance compared to the infection prediction model using
231 symptom information, our results indicate that an individual's COVID-19 infection status can be
232 predicted to a certain degree without relying on explicit symptoms or contact tracing information. This
233 could potentially address the limitations of existing studies, which struggle to predict the infection risks
234 of asymptomatic carriers, and of mobile contact tracing applications with a low user base. Moreover, as
235 shown in Fig 6, our study supports the general characteristic of infectious diseases that a higher
236 possibility of contact with others leads to a higher risk of infection.

237 Another point worth discussing is the overlap information we used for model development.
238 We created a variable for overlap information if the cell tower locations of each individual overlapped
239 with others for more than 10 seconds. However, an overlap of cell tower locations does not necessarily
240 imply direct contact between two individuals. Furthermore, the overlap of cell tower locations for 10
241 seconds does not confirm the transmission of COVID-19. The data we used represents only a very small
242 part of all Korean citizens' cell tower location records, and it is not possible to conclude that an
243 individual was infected with COVID-19 from the people they overlapped locations with based on this
244 information alone. We believe this data does not suggest that an individual contracted COVID-19 from
245 someone with whom they overlapped locations, but rather indirectly indicates they were in a location
246 with a high risk of COVID-19 infection. In our sample of 837 people, the places where locations
247 overlapped are likely to have seen many more overlapping individuals, thus increasing potential contact

248 points. As shown in Fig 6, we believe the risk of COVID-19 infection increases when an individual
249 frequently overlaps cell tower locations outside their residence, thereby increasing potential contact with
250 others.

251 In order to apply the findings of this research, the following prerequisites need to be fulfilled.
252 The government or the institution intending to utilize these research results should be able to collect cell
253 tower location information from mobile carriers. The Korean government was able to collect cell tower
254 location information of confirmed cases without the individual's consent through legal procedures.
255 While it is not necessary for many users to install and use the contact tracing app, the system and
256 technology should allow for the collection of individual cell tower location information.

257 Additionally, we foresee the following further research to enhance our findings. First,
258 important conclusions were derived from a relatively small sample of 837 people within a limited period
259 and geographic scope. However, these results emphasize the need for further research based on a larger
260 dataset. If the study is expanded to include a significantly larger number of people, more substantial
261 conclusions could be drawn. Second, it would be beneficial to investigate whether the risk of infection
262 from diseases other than COVID-19 can also be predicted using cell tower location information. For
263 infectious diseases with varying characteristics, the accuracy of infection prediction using cell tower
264 location information may vary, and this is important consideration. Lastly, we believe that by
265 incorporating individual symptom information, a more accurate infectious disease prediction model
266 could be developed.

267 **Acknowledgements**

268 This research is based on the researchs "A Next Generation Surveillance Study for Epidemic
269 Preparedness(INV-006404)" which was funded by the Bill & Melinda Gates Foundation, and "COVID-

270 19 Self-Risk Evaluation with Digital Contact Tracing(RF-TAA-2020-D07)” which was funded by the
271 RIGHT Foundation. The findings and conclusions contained within are those of the authors and do not
272 necessarily reflect positions or policies of the Bill & Melinda Gates Foundation.

273 **References**

- 274 1. Perra N. Non-pharmaceutical interventions during the COVID-19 pandemic: A review. *Physics*
275 *Reports*. 2021;913: 1–52. doi:<https://doi.org/10.1016/j.physrep.2021.02.001>
- 276 2. Wake RM, Morgan M, Choi J, Winn S. Reducing nosocomial transmission of COVID-19:
277 Implementation of a COVID-19 triage system. *Clin Med*. 2020;20: e141–e145.
- 278 3. Ayouni I, Maatoug J, Dhouib W, Zammit N, Fredj SB, Ghammam R, et al. Effective public
279 health measures to mitigate the spread of COVID-19: A systematic review. *BMC Public Health*.
280 2021;21: 1015. doi:10.1186/s12889-021-11111-1
- 281 4. Alballa N, Al-Turaiki I. Machine learning approaches in COVID-19 diagnosis, mortality, and
282 severity risk prediction: A review. *Informatics in Medicine Unlocked*. 2021;24: 100564.
283 doi:<https://doi.org/10.1016/j.imu.2021.100564>
- 284 5. Jamshidi M, Lalbakhsh A, Talla J, Peroutka Z, Hadjilooei F, Lalbakhsh P, et al. Artificial
285 intelligence and COVID-19: Deep learning approaches for diagnosis and treatment. *IEEE Access*.
286 2020;8: 109581–109595. doi:10.1109/ACCESS.2020.3001973
- 287 6. Lalmuanawma S, Hussain J, Chhakchhuak L. Applications of machine learning and artificial
288 intelligence for covid-19 (SARS-CoV-2) pandemic: A review. *Chaos, Solitons & Fractals*. 2020;139:
289 110059. doi:<https://doi.org/10.1016/j.chaos.2020.110059>
- 290 7. Bradford L, Aboy M, Liddell K. COVID-19 contact tracing apps: a stress test for privacy, the
291 GDPR, and data protection regimes. *Journal of Law and the Biosciences*. 2020;7.
292 doi:10.1093/jlb/ljaa034
- 293 8. Grantz KH, Meredith HR, Cummings DAT, Metcalf CJE, Grenfell BT, Giles JR, et al. The use
294 of mobile phone data to inform analysis of COVID-19 pandemic epidemiology. *Nature*
295 *Communications*. 2020;11: 4961. doi:10.1038/s41467-020-18190-5
- 296 9. Shahroz M, Ahmad F, Younis MS, Ahmad N, Kamel Boulos MN, Vinuesa R, et al. COVID-19
297 digital contact tracing applications and techniques: A review post initial deployments. *Transportation*
298 *Engineering*. 2021;5: 100072.
- 299 10. Kang J, Jang YY, Kim J, Han S-H, Lee KR, Kim M, et al. South korea’s responses to stop the
300 COVID-19 pandemic. *American Journal of Infection Control*. 2020;48: 1080–1086.
301 doi:<https://doi.org/10.1016/j.ajic.2020.06.003>

- 302 11. Lee M, MANSU K, Jaejin Yi SH kyuhwan moon. Big data based epidemic investigation support
 303 system using mobile network data. The Korea Journal of BigData. 2020;5: 187–199.
 304 doi:10.36498/kbigdt.2020.5.2.187
- 305 12. KT. 2022 form 20-f annual report. 2022. Available:
 306 <https://www.sec.gov/ix?doc=/Archives/edgar/data/892450/000119312523123967/d436251d20f.htm>
- 307 13. Park MJ, Choi JH, Cho JH. Estimation of the effectiveness of a tighter, reinforced quarantine for
 308 the coronavirus disease 2019 (COVID-19) outbreak: Analysis of the third wave in south korea. Journal
 309 of Personalized Medicine. 2023;13. doi:10.3390/jpm13030402
- 310 14. Jun M-J, Kim JI, Kwon JH, Jeong J-E. The effects of high-density suburban development on
 311 commuter mode choices in seoul, korea. Cities. 2013;31: 230–238.
 312 doi:<https://doi.org/10.1016/j.cities.2012.06.016>
- 313 15. Bhattacharya A, Ranjan P, Kumar A, Brijwal M, Pandey RM, Mahishi N, et al. Development
 314 and validation of a clinical symptom-based scoring system for diagnostic evaluation of COVID-19
 315 patients presenting to outpatient department in a pandemic situation. Cureus. 2021;13: e13681.
- 316 16. Sun L, Song F, Shi N, Liu F, Li S, Li P, et al. Combination of four clinical indicators predicts the
 317 severe/critical symptom of patients infected COVID-19. J Clin Virol. 2020;128: 104431.
- 318 17. Wu J, Zhang P, Zhang L, Meng W, Li J, Tong C, et al. Rapid and accurate identification of
 319 COVID-19 infection through machine learning based on clinical available blood test results. medRxiv.
 320 2020. doi:10.1101/2020.04.02.20051136
- 321 18. Pal M, Parija S, Mohapatra RK, Mishra S, Rabaan AA, Al Mutair A, et al. Symptom-based
 322 COVID-19 prognosis through AI-based IoT: A bioinformatics approach. Biomed Res Int. 2022;2022:
 323 3113119.
- 324 19. Sudre CH, Lee KA, Lochlainn MN, Varsavsky T, Murray B, Graham MS, et al. Symptom
 325 clusters in COVID-19: A potential clinical prediction tool from the COVID symptom study app. Sci
 326 Adv. 2021;7: eabd4177.
- 327 20. Kennedy B, Fitipaldi H, Hammar U, Maziarz M, Tsereteli N, Oskolkov N, et al. App-based
 328 COVID-19 syndromic surveillance and prediction of hospital admissions in COVID symptom study
 329 sweden. Nature Communications. 2022;13: 2110.

330 Supporting information

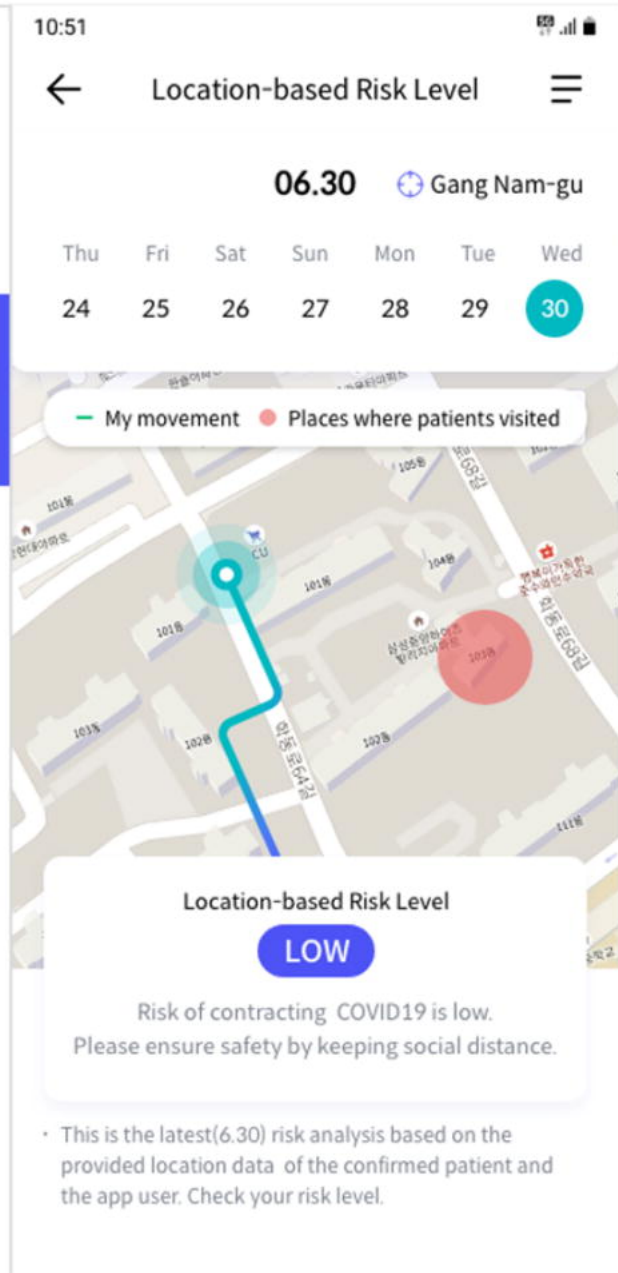
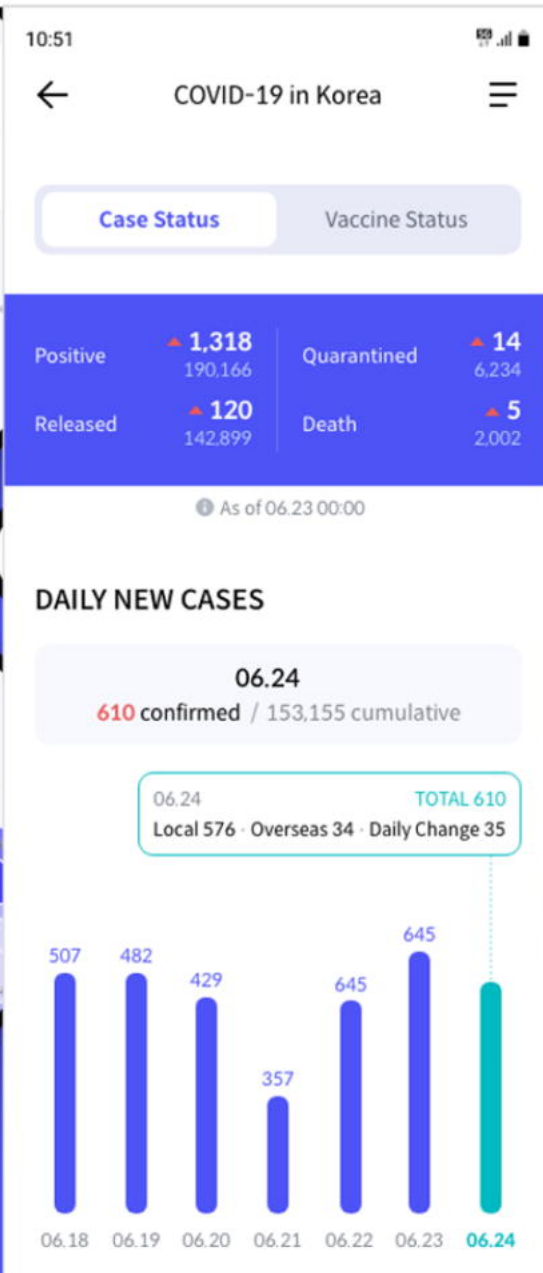
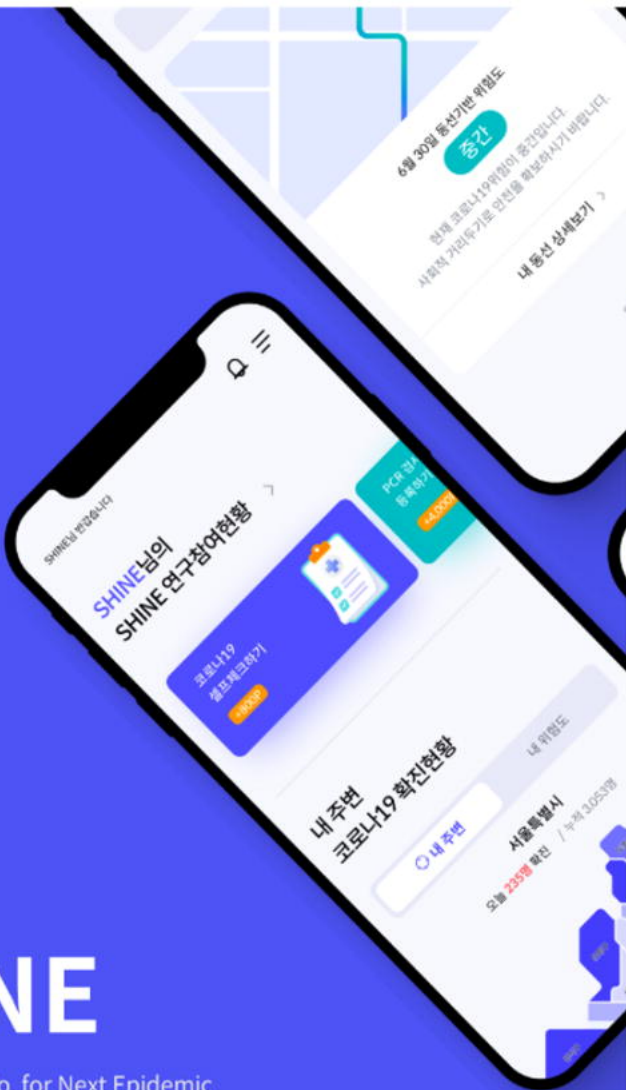
331 S1 Table. Description of software packages, methods and tuning parameters for model 332 development.

Algorithm	Package / Method	Parameters	Final chosen model
Logistic Regression	scikit-learn / LogisticRegressionCV	solver='liblinear' max_iter=1000	
XGBoost	xgboost / XGBClassifier	n_estimators=range(200, 310, 10) learning_rate=[0.1] max_depth=range(8, 11) alpha=[0, 0.1] gamma=[0, 0.1] lambda=[0, 0.1] subsample=[0.4, 0.5, 0.6, 0.7, 0.8]	alpha=0 colsample_bytree=0.7 gamma=0.1 lambda=0 max_depth=8

Algorithm	Package / Method	Parameters	Final chosen model
		colsample_bytree=[0.4, 0.5, 0.6, 0.7, 0.8] min_child_weight=range(1, 9) objective=['binary:logistic']	n_estimators=260 subsample=0.8
Random forest	scikit-learn / RandomForestClassifier	n_estimators=range(200, 300, 10) max_depth=range(3, 11) max_features=range(4, 13) max_leaf_nodes=range(4, 13) min_samples_leaf=range(5, 11) min_samples_split:range(5, 11)	max_depth=10 max_features=12 max_leaf_nodes=13 min_samples_leaf=7 min_samples_split=5 n_estimators=220 random_state=42

SHINE

Study of Health Info. for Next Epidemic



10:51

← Enter PCR Result ≡

Enter PCR result and upload a proof of the test.
Our research staff will verify the image and issue points accordingly.

PCR Test date

2021. 03. 17

You can only enter the test results within 10 days of your last COVID-19 self-checker date.

PCR Result

Negative Positive

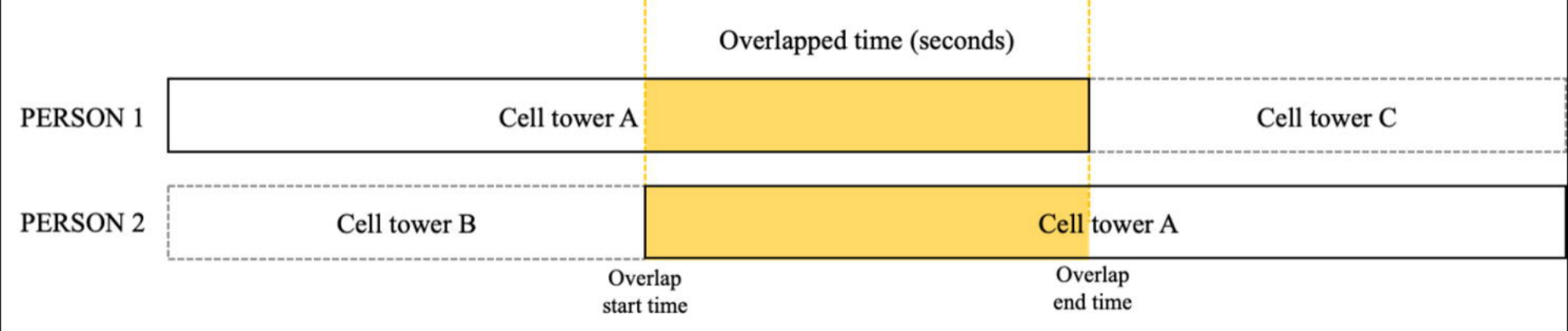
Please upload an image of the test result.

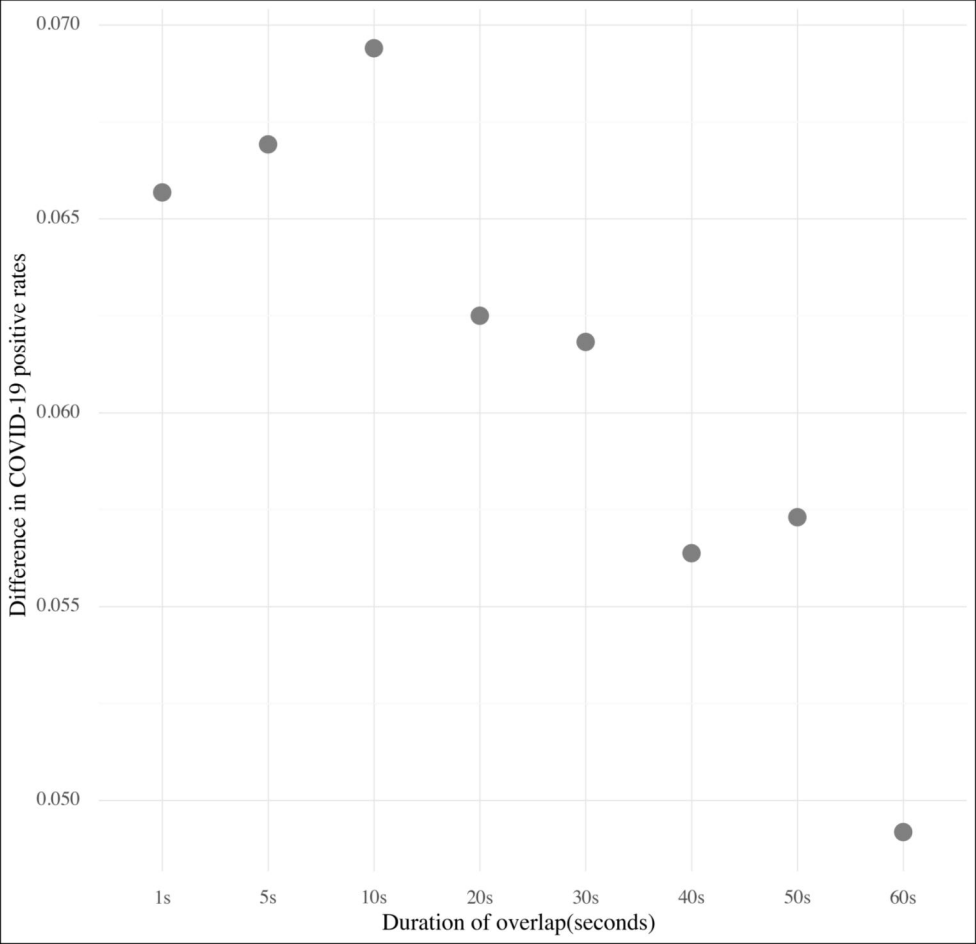
Test Location

Gangnam Screening Center

Attachment

- Only 1 file can be uploaded.
- Image should contain name, test date, test result and test location.
- To get reward points, please upload a verifiable, quality image of the PCR result.





Data source / Data preprocessing

SHINE Data
n = 837
6 variables

Distinguishing between residential and non-residential areas
in location records

Calculation of overlaps between location records

Feature engineering

Data Transform: Natural logarithm
(All predictor variables)

Model training

n = 837
14 variables

Training:
n = 585(70%)

Test:
n = 252(30%)

Filtering outliers
n = 418

Over-sampling
346:72(P:N) -> 346:346(P:N)

Machine learning
algorithm applied:
10-fold cross-
validation

