

# Predicting CoVID-19 community mortality risk using machine learning and development of an online prognostic tool

Ashis Das <sup>Corresp., 1</sup>, Shiba Mishra <sup>2</sup>, Saji Saraswathy Gopalan <sup>1</sup>

<sup>1</sup> The World Bank, Washington, District of Columbia, United States

<sup>2</sup> Credit Suisse Private Limited, Pune, India

Corresponding Author: Ashis Das  
Email address: adas8@worldbank.org

**Background.** The recent pandemic of CoVID-19 has emerged as a threat to global health security. There are very few prognostic models on CoVID-19 using machine learning.

**Objectives.** To predict mortality among confirmed CoVID-19 patients in South Korea using machine learning and deploy the best performing algorithm as an open-source online prediction tool for decision-making. **Materials and methods.** Mortality for confirmed

CoVID-19 patients (n=3,524) between January 20, 2020 and May 30, 2020 was predicted using five machine learning algorithms (logistic regression, support vector machine, K nearest neighbor, random forest and gradient boosting). The performance of the algorithms was compared, and the best performing algorithm was deployed as an online prediction tool.

**Results.** The logistic regression algorithm was the best performer in terms of discrimination (area under ROC curve=0.830), calibration (Matthews Correlation Coefficient=0.433; Brier Score=0.036) and. The best performing algorithm (logistic regression) was deployed as the online CoVID-19 Community Mortality Risk Prediction tool named CoCoMoRP ( <https://ashis-das.shinyapps.io/CoCoMoRP/> ). **Conclusions.** We

describe the development and deployment of an open-source machine learning tool to predict mortality risk among CoVID-19 confirmed patients using publicly available surveillance data. This tool can be utilized by potential stakeholders such as health providers and policymakers to triage patients at the community level in addition to other approaches.

# **Predicting CoVID-19 community mortality risk using machine learning and development of an online prognostic tool**

Ashis Kumar Das, MBBS, MPH, PhD<sup>1#</sup>

Shiba Mishra, BE, PGDBA<sup>2</sup>

Saji Saraswathy Gopalan, PhD, DrPH<sup>1</sup>

1. The World Bank, Washington DC, USA
2. Credit Suisse Private Limited, Pune, India

# Corresponding author:

Ashis Kumar Das

The World Bank, Washington DC, USA.

E-mail: [adas8@worldbank.org](mailto:adas8@worldbank.org)

# Abstract

**Background.** The recent pandemic of CoVID-19 has emerged as a threat to global health security. There are very few prognostic models on CoVID-19 using machine learning.

**Objectives.** To predict mortality among confirmed CoVID-19 patients in South Korea using machine learning and deploy the best performing algorithm as an open-source online prediction tool for decision-making.

**Materials and methods.** Mortality for confirmed CoVID-19 patients (n=3,524) between January 20, 2020 and May 30, 2020 was predicted using five machine learning algorithms (logistic regression, support vector machine, K nearest neighbor, random forest and gradient boosting). The performance of the algorithms was compared, and the best performing algorithm was deployed as an online prediction tool.

**Results.** The logistic regression algorithm was the best performer in terms of discrimination (area under ROC curve=0.830), calibration (Matthews Correlation Coefficient=0.433; Brier Score=0.036) and. The best performing algorithm (logistic regression) was deployed as the online CoVID-19 Community Mortality Risk Prediction tool named CoCoMoRP (<https://ashisdas.shinyapps.io/CoCoMoRP/>).

**Conclusions.** We describe the development and deployment of an open-source machine learning tool to predict mortality risk among CoVID-19 confirmed patients using publicly available surveillance data. This tool can be utilized by potential stakeholders such as health providers and policymakers to triage patients at the community level in addition to other approaches.

# Introduction

A novel coronavirus disease 2019 (CoVID-19) originated from Wuhan in China was reported to the World Health Organization in December of 2019.<sup>1</sup> Ever since, this novel coronavirus has spread to almost all major nations in the world resulting in a major pandemic. As of June 08, 2020, it has contributed to more than 7 million confirmed cases and about 404,000 deaths.<sup>2</sup> The first CoVID-19 case was diagnosed in South Korea on January 20, 2020. According to the Korea Centers for Disease Control and Prevention (KCDC), there have been 11,814 confirmed cases and 273 deaths due to CoVID-19 as of June 08, 2020.<sup>3</sup>

In the field of healthcare, accurate prognosis is essential for efficient management of patients while prioritizing care to the more needy. In order to aid in prognosis, several prediction models have been developed using various methods and tools including machine learning.<sup>4-6</sup> Machine learning is a field of artificial intelligence where computers simulate the processes of human intelligence and can synthesize complex information from huge data sources in a short period of time.<sup>7</sup> Though there have been a few prediction tools on CoVID-19, only a handful have utilized machine learning.<sup>8</sup> To the best of our knowledge, by far there is no publicly available online CoVID-19 prognosis prediction tool from the general population of confirmed cases using machine learning. We attempt to apply machine learning on the publicly available CoVID-19 data at the community level from South Korea to predict mortality.

Our study had two objectives, (1) predict mortality among confirmed CoVID-19 patients in South Korea using machine learning algorithms, and (2) deploy the best performing algorithm as an open-source online prediction tool for decision-making.

# **Material and methods**

## **Patients**

Patients for this study were selected from the data shared by Korea Centers for Disease Control and Prevention (KCDC).<sup>3</sup> The timeframe of this study was from the beginning of the detection of the first case (January 20, 2020) through May 30, 2020. In the dataset, there were a total of 4,004 patients. Our inclusion criteria were confirmed CoVID-19 cases with availability of demographic, exposure and diagnosis confirmation features along with the outcome. We excluded patients those had missing features – sex (n=330) and age (n=150), and thus, 3,524 patients were included in the final analysis.

## **Outcome variable**

The outcome variable was mortality and it had a binary distribution – “yes” if the patient died, or “no” otherwise.

## **Predictors**

The predictors were individual patient level demographic and exposure features. They were four predictors - age group, sex, province, and exposure. There were ten age groups as follows below 10 years, 10-19 years, 20-29 years, 30-39 years, 40-49 years, 50-59 years, 60-69 years, 70-79 years, 80-89 years, 90 years and above. Patients represented all 17 provinces of South Korea (Busan, Chungcheongbuk-do, Chungcheongnam-do, Daegu, Daejeon, Gangwon-do, Gwangju, Gyeonggi-do, Gyeongsangbuk-do, Gyeongsangnam-do, Incheon, Jeju-do, Jeollabuk-do, Jeollanam-do, Sejong, Seoul, and Ulsan). Patients were exposed in several settings, such as nursing home, hospital, religious gathering, call center, community center, shelter and apartment, gym facility, overseas inflow, contact with patients and others.

## Statistical Methods

### *Descriptive Analysis*

We performed descriptive analyses of the predictors by respective stratification groups and present the results as numbers and proportions. Potential correlations between predictors were tested with Pearson's correlation coefficient.

### *Predictive Analysis*

We applied machine learning algorithms to predict mortality among CoVID-19 confirmed cases. Machine learning is a branch of artificial intelligence where computer systems can learn from available data and identify patterns with minimal human intervention.<sup>9</sup> Typically, in machine learning several algorithms are tested on data and performance metrics are used to select the best performing algorithm. While selecting the algorithms, we considered commonly used machine learning algorithms in healthcare research that have lower training time as well as lower lag time when built into an online application. Thus, the selected algorithms were – logistic regression, support vector machine, K neighbor classification, random forest and gradient boosting. Using grid search function, we also performed hyperparameter tuning (i.e. selection of the best parameters) for each algorithm (Supplemental Table S1). Logistic regression is best suited for a binary or categorical output. It tries to describe the relationship between the output and predictor variables.<sup>10</sup> In support vector machine (SVM) algorithm, the data is classified into two classes based on the output variable over a hyperplane.<sup>10</sup> The algorithm tries to increase the distance between the hyperplane and the most proximal two data points in each class. SVM uses a set of mathematical functions called kernels, which transform the inputs to required forms. In our SVM algorithm, we used a radial kernel. K Nearest Neighbors (KNN) is a non-parametric approach

that decides the output classification by the majority class among its neighbors.<sup>11</sup> The number of neighbors can be altered to arrive at the best fitting KNN model. Random forest algorithm uses a combination of decision trees.<sup>12</sup> Decision trees are generated by recursively partitioning the predictors. New attributes are sequentially fitted to predict the output. Gradient boosting (GB) algorithm uses a combination of decision trees.<sup>13</sup> Each decision tree dynamically learns from its precursor and passes on the improved function to the following. Finally, the weighted combination of these trees provides the prediction. A decision tree's learning from the precursor and the number of subsequent trees can be respectively adjusted using learning rate and number of trees parameters.

#### ***Evaluation of the performance of the algorithms***

We split the data into training (80 percent) and test cohorts (20 percent). Initially, the algorithms were trained on the training cohort and then were validated on the test cohort (new data) for determining predictions. The data was passed through a 10-fold cross validation where the data was split into training and test cohorts at 80/20 ratio randomly ten times. The final prediction came out of the cross-validated estimate. As our data was imbalanced (only 2.1% output were with the condition against 97.9% without), we applied two oversampling techniques called synthetic minority oversampling technique (SMOTE) and adaptive synthetic (ADASYN) method to enhance the learning on the training data.<sup>14,15</sup> SMOTE creates synthetic samples from the minority class (cases with deaths in our data) according to feature space similarities between nearest neighbors.<sup>14</sup> ADASYN adaptively generates synthetic samples based on their difficulty in learning.<sup>16</sup>

The performance of the algorithms were evaluated for discrimination, calibration and overall performance. Discrimination is the ability of the algorithm to separate out patients with the

mortality risk from those without, where as calibration is the agreement between observed and predicted risk of mortality. An ideal model should have the best of both discrimination and calibration. We tested discrimination with area under the receiver operating characteristics curve (AUC) and calibration with Matthews correlation coefficient. A receiver operator characteristic (ROC) curve plots the true positive rate on y-axis against the false positive rate on x-axis.<sup>17</sup> AUC is score that measures the area under the ROC curve and it ranges from 0.50 to 1.0 with higher values meaning higher discrimination. Matthews correlation coefficient (MCC) is a measure that takes into account all four predictive classes – true positive, true negative, false positive and false negative.<sup>18</sup> Brier score simultaneously account for discrimination and calibration.<sup>17</sup> A smaller Brier score indicates better performance. We also estimated accuracy, sensitivity and specificity. Accuracy is a measure of correct classification of death cases as death and survived cases as survived.<sup>17</sup> Sensitivity is a measure of correctly predicting death among all those who died, whereas specificity is a measure of correctly predicting survival among all those who survived. In addition, relative influence of the predictors with the output was estimated using mean decrease Gini coefficients (MDG) in the random forest algorithm.<sup>19</sup> MDG quantifies which predictor contributed most to the classification accuracy.

The statistical analyses were performed using Stata Version 15 (StataCorp LLC. College Station, TX), Python programming language Version 3.7.1 (Python Software Foundation, Wilmington, DE, USA); e1071 and caret packages of R programming language Version 3.6.3 (R Foundation for Statistical Computing, Vienna, Austria). The web application was built using the Shiny package for R and deployed with Shiny server.

# Results

## Patient profile

The profile of the patients is presented in Table 1. Out of 3,524 confirmed patients, a slightly more than half were females (55.1%). Among the age groups, the maximum patients were from 20-29 years (24.4%), followed by 50-59 years (17.7%), 30-39 years (14%), 40-49 years (13.7%), and 60-69 years (12%). Gyeongsangbuk-do (35.1%), Gyeonggi-do (23.5%) and Seoul (16%) provinces together presented the maximum patients. Considering the source/mode of infection, the largest group had unknown mode (39.3%) followed by direct contact with patients (29.8%) and from overseas (17.4%). According to this available data source, there were 74 deaths accounting for 2.1 percent of the patients.

The correlation coefficients among the predictors ranged from -0.12 to 0.22. Using the random forest algorithm, we estimated the relative influence of the predictors (figure 1). Age was the most important predictor followed by exposure, sex and province.

## Performance of the algorithms

Table 2 presents the performance metrics of all algorithms – logistic regression, support vector machine, K nearest neighbor, random forest and gradient boosting. The area under receiver operating characteristic curve (AUC) ranged from 0.644 to 0.830 with the best score for the logistic regression (SMOTE) algorithm. Similarly, logistic regression (SMOTE) performed the best on Matthews correlation coefficient. It was in the middle for the performance on Brier score. The accuracy of all algorithms was very similar with random forest (SMOTE) performing the

best (0.972) and K nearest neighbor with the least score (0.924). Considering all the performance metrics, logistic regression (SMOTE) was the best performing algorithm.

# **Online CoVID-19 mortality risk prediction tool – CoCoMoRP**

The best performing model – logistic regression (SMOTE) was deployed as the online mortality risk prediction tool named as “CoVID-19 Community Mortality Risk Prediction” – CoCoMoRP” (<https://ashis-das.shinyapps.io/CoCoMoRP/>). Figure 2 presents the user interface of the prediction tool. The web application is optimized to be conveniently used on multiple devices such as desktops, tablets, and smartphones.

The user interface has four boxes to select input features as drop-down menus. The features are sex (two options – male and female), age (ten options – below 10 years, 10-19 years, 20-29 years, 30-39 years, 40-49 years, 50-59 years, 60-69 years, 70-79 years, 80-89 years, 90 years and above), province (all 17 provinces – Busan, Chungcheongbuk-do, Chungcheongnam-do, Daegu, Daejeon, Gangwon-do, Gwangju, Gyeonggi-do, Gyeongsangbuk-do, Gyeongsangnam-do, Incheon, Jeju-do, Jeollabuk-do, Jeollanam-do, Sejong, Seoul, Ulsan), and exposure (nine options – nursing home; hospital; religious gathering; call center; community center, shelter and apartment; gym facility; overseas inflow; contact with patients; and others).

The user has to select one option each from the input feature boxes and click the submit button to estimate the CoVID-19 mortality risk probability in percentages. For instance, the tool gives a CoVID-19 mortality risk prediction of 94.1% for a male patient aged between 80 and 89 years from Seoul province coming in contact with patient as the exposure.

# Discussion

The CoVID-19 pandemic is a threat to global health and economic security. Recent evidence for this new disease is still evolving on various clinical and socio-demographic dimensions.<sup>20–22</sup> Simultaneously, health systems across the world are constrained with resources to efficiently deal with this pandemic. We describe the development and deployment of an open-source artificial intelligence informed prognostic tool to predict mortality risk among CoVID-19 confirmed patients using publicly available surveillance data. This tool can be utilized by potential stakeholders such as health providers and policy makers to triage patients at the community level in addition to other approaches.

One major limitation of this tool is unavailability of crucial clinical information on symptoms, risk factors and clinical parameters. Recent research has identified certain symptoms, preexisting illnesses and clinical parameters as strong predictors of prognosis and severity of progression for CoVID-19.<sup>22–24</sup> These crucial pieces of information are not publicly available so far in the surveillance data, so the tool could not be tested to include these features. Inclusion of these additional features may improve the reliability and relevance of the tool. Therefore, we urge the users to balance the predictions from this tool against their own and/or health provider’s clinical expertise and other relevant clinical information. The second limitation pertains to lack of availability of the complete data. According to the reports, there were 11,814 confirmed cases and 273 deaths (case fatality rate 2.3%) due to CoVID-19 in South Korea as of June 08, 2020. However, our analysis using the publicly released database found 3,529 cases and 74 deaths (case fatality rate 2.1%) until May 30, 2020. Though the case fatality rates are similar, our

analysis uses respectively about a third and a fourth of totally reported cases and deaths. As more data are released publicly, we would continue to update our analyses and the web-application.

## Conclusions

We tested multiple machine learning models to accurately predict deaths due to CoVID-19 among confirmed community cases in the Republic of Korea. Using the best performing algorithm, we developed and deployed an online mortality risk prediction tool. To the best of our knowledge, our CoVID-19 community mortality risk prediction tool is the first of its kind. Our tool offers an additional approach to informing decision making for CoVID-19 patients.

## Acknowledgements

We are grateful to Korea Center for Disease Control and Prevention for making this data publicly available. The views expressed in the paper are that of the authors and do not reflect that of their affiliations. This particular work was conducted outside of the authors' organizational affiliations.

## REFERENCES

1. WHO. WHO Coronavirus disease (COVID-2019) situation reports 2020.
2. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). Available at: <https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6>.
3. KCDC. Korea Centers for Disease Control and Prevention; Seoul, Korea: 2020. 2020
4. Chen, J. H. & Asch, S. M. Machine learning and prediction in medicine-beyond the peak of inflated expectations. *New England Journal of Medicine* (2017). doi:10.1056/NEJMp1702071
5. Qu Y, Yue G, Shang C, Yang L, Zwiggelaar R, Shen Q.. Multi-criterion mammographic risk

- analysis supported with multi-label fuzzy-rough feature selection. *Artif. Intell. Med.* (2019).  
doi:10.1016/j.artmed.2019.101722
6. Lei L, Wang Y, Xue Q, Tong J, Zhou CM, Yang JJ. A comparative study of machine learning algorithms for predicting acute kidney injury after liver cancer resection. *PeerJ* (2020).  
doi:10.7717/peerj.8583
7. Benke, K. & Benke, G. Artificial intelligence and big data in public health. *International Journal of Environmental Research and Public Health* (2018). doi:10.3390/ijerph15122796
8. Wynants L, Van Calster B, Bonten MMJ, Collins GS, Debray TPA, De Vos M, Haller MC, Heinze G, Moons KGM, Riley RD, Schuit E, Smits LJM, Snell KIE, Steyerberg EW, Wallisch C, van Smeden M. Prediction models for diagnosis and prognosis of covid-19 infection: systematic review and critical appraisal. *BMJ* **369**, m1328 (2020).
9. Deo, R. C. Machine learning in medicine. *Circulation* (2015).  
doi:10.1161/CIRCULATIONAHA.115.001593
10. Jiang F, Jiang Y, Zhi H, Dong Y, Li H, Ma S, Wang Y, Dong Q, Shen H, Wang Y. Artificial intelligence in healthcare: Past, present and future. *Stroke and Vascular Neurology* (2017).  
doi:10.1136/svn-2017-000101
11. Raeisi Shahraki, H., Pourahmad, S. & Zare, N. K Important Neighbors: A Novel Approach to Binary Classification in High Dimensional Data. *Biomed Res. Int.* (2017).  
doi:10.1155/2017/7560807
12. Rigatti, S. J. Random Forest. *J. Insur. Med.* (2017). doi:10.17849/in-sm-47-01-31-39.1
13. Natekin, A. & Knoll, A. Gradient boosting machines, a tutorial. *Front. Neurorobot.* (2013).  
doi:10.3389/fnbot.2013.00021
14. Chawla, N. V., Bowyer, K. W., Hall, L. O. & Kegelmeyer, W. P. SMOTE: Synthetic minority over-sampling technique. *J. Artif. Intell. Res.* (2002). doi:10.1613/jair.953
15. Nnamoko, N. & Korkontzelos, I. Efficient treatment of outliers and class imbalance for diabetes prediction. *Artif. Intell. Med.* (2020). doi:10.1016/j.artmed.2020.101815
16. He, H., Bai, Y., Garcia, E. A. & Li, S. ADASYN: Adaptive synthetic sampling approach for imbalanced learning. in *Proceedings of the International Joint Conference on Neural Networks* (2008). doi:10.1109/IJCNN.2008.4633969
17. Huang, Y., Li, W., Macheret, F., Gabriel, R. A. & Ohno-Machado, L. A tutorial on calibration measurements and calibration models for clinical prediction models. *J. Am. Med. Inform. Assoc.* (2020). doi:10.1093/jamia/ocz228
18. Chicco, D. & Jurman, G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics* (2020).  
doi:10.1186/s12864-019-6413-7
19. Xie, J. & Coggeshall, S. Prediction of transfers to tertiary care and hospital mortality: A gradient boosting decision tree approach. *Stat. Anal. Data Min.* (2010). doi:10.1002/sam.10079
20. Sun, P., Lu, X., Xu, C., Sun, W. & Pan, B. Understanding of COVID-19 based on current evidence. *Journal of Medical Virology* (2020). doi:10.1002/jmv.25722
21. Chen H, Guo J, Wang C, Luo F, Yu X, Zhang W, Li J, Zhao D, Xu D, Gong Q, Liao J, Yang H, Hou W, Zhang Y. Clinical characteristics and intrauterine vertical transmission potential of

- COVID-19 infection in nine pregnant women: a retrospective review of medical records. *Lancet* (2020). doi:10.1016/S0140-6736(20)30360-3
22. Li B, Yang J, Zhao F, Zhi L, Wang X, Liu L, Bi Z, Zhao Y. Prevalence and impact of cardiovascular metabolic diseases on COVID-19 in China. *Clinical Research in Cardiology* (2020). doi:10.1007/s00392-020-01626-9
23. Li L quan, Huang T, Wang Y qing, Wang Z ping, Liang Y, Huang T bi, Zhang H yun, Sun W, Wang Y. 2019 novel coronavirus patients' clinical characteristics, discharge rate, and fatality rate of meta-analysis. *Journal of Medical Virology* (2020). doi:10.1002/jmv.25757
24. Guan W, Ni Z, Hu Y, Liang W, Ou C, He J, Liu L, Shan H, Lei C, Hui DSC, Du B, Li L, Zeng G, Yuen K-Y, Chen R, Tang C, Wang T, Chen P, Xiang J, Li S, Wang J, Liang Z, Peng Y, Wei L, Liu Y, Hu Y, Peng P, Wang J, Liu J, Chen Z, Li G, Zheng Z, Qiu S, Luo J, Ye C, Zhu S, Zhong N. Clinical Characteristics of Coronavirus Disease 2019 in China. *N. Engl. J. Med.* (2020). doi:10.1056/nejmoa2002032

**Table 1** (on next page)

Sample characteristics

1 **Table 1. Sample characteristics**

Variable	Number	Proportion (%)
<b>Sex</b>		
Female	1,940	55.1
Male	1,584	45.0
<b>Age group (years)</b>		
Below 10	60	1.7
10-19	160	4.5
20-29	859	24.4
30-39	494	14.0
40-49	483	13.7
50-59	625	17.7
60-69	423	12.0
70-79	210	6.0
80-89	162	4.6
90 and above	48	1.4
<b>Province</b>		
Busan	144	4.1
Chungcheongbuk-do	52	1.5
Chungcheongnam-do	146	4.1
Daegu	63	1.8
Daejeon	46	1.3
Gangwon-do	52	1.5
Gwangju	30	0.9
Gyeonggi-do	829	23.5
Gyeongsangbuk-do	1,236	35.1
Gyeongsangnam-do	119	3.4
Incheon	92	2.6
Jeju-do	14	0.4
Jeollabuk-do	20	0.6
Jeollanam-do	19	0.5
Sejong	47	1.3
Seoul	563	16.0
Ulsan	52	1.5
<b>Exposure</b>		
Nursing home	46	1.3
Hospital	37	1.1
Religious gathering	160	4.5
Call center	135	3.8
Community center, shelter and apartment	68	1.9
Gym facility	34	1.0
Overseas inflow	612	17.4
Contact with patients	1,049	29.8
Others	1,383	39.3
<b>Outcome</b>		
Survived	3,450	97.9
Died	74	2.1
Total	3,524	100

2

**Table 2**(on next page)

Performance of the machine learning algorithms

Table 2. Performance of the machine learning algorithms

Algorithm	Oversampling method	Area under ROC curve	Matthews correlation coefficient	Brier score	Sensitivity	Specificity	Accuracy
Logistic regression	SMOTE <sup>#</sup>	0.830	0.433	0.036	0.692	0.968	0.965
	ADASYN <sup>*</sup>	0.823	0.376	0.049	0.692	0.955	0.968
Support vector machine	SMOTE <sup>#</sup>	0.825	0.393	0.045	0.692	0.959	0.970
	ADASYN <sup>*</sup>	0.786	0.345	0.048	0.615	0.958	0.971
K nearest neighbor	SMOTE <sup>#</sup>	0.644	0.253	0.031	0.307	0.981	0.942
	ADASYN <sup>*</sup>	0.759	0.410	0.028	0.538	0.979	0.924
Random forest	SMOTE <sup>#</sup>	0.787	0.351	0.046	0.615	0.959	0.972
	ADASYN <sup>*</sup>	0.787	0.351	0.046	0.615	0.959	0.971
Gradient boosting	SMOTE <sup>#</sup>	0.787	0.351	0.046	0.615	0.959	0.971
	ADASYN <sup>*</sup>	0.787	0.351	0.046	0.615	0.959	0.971

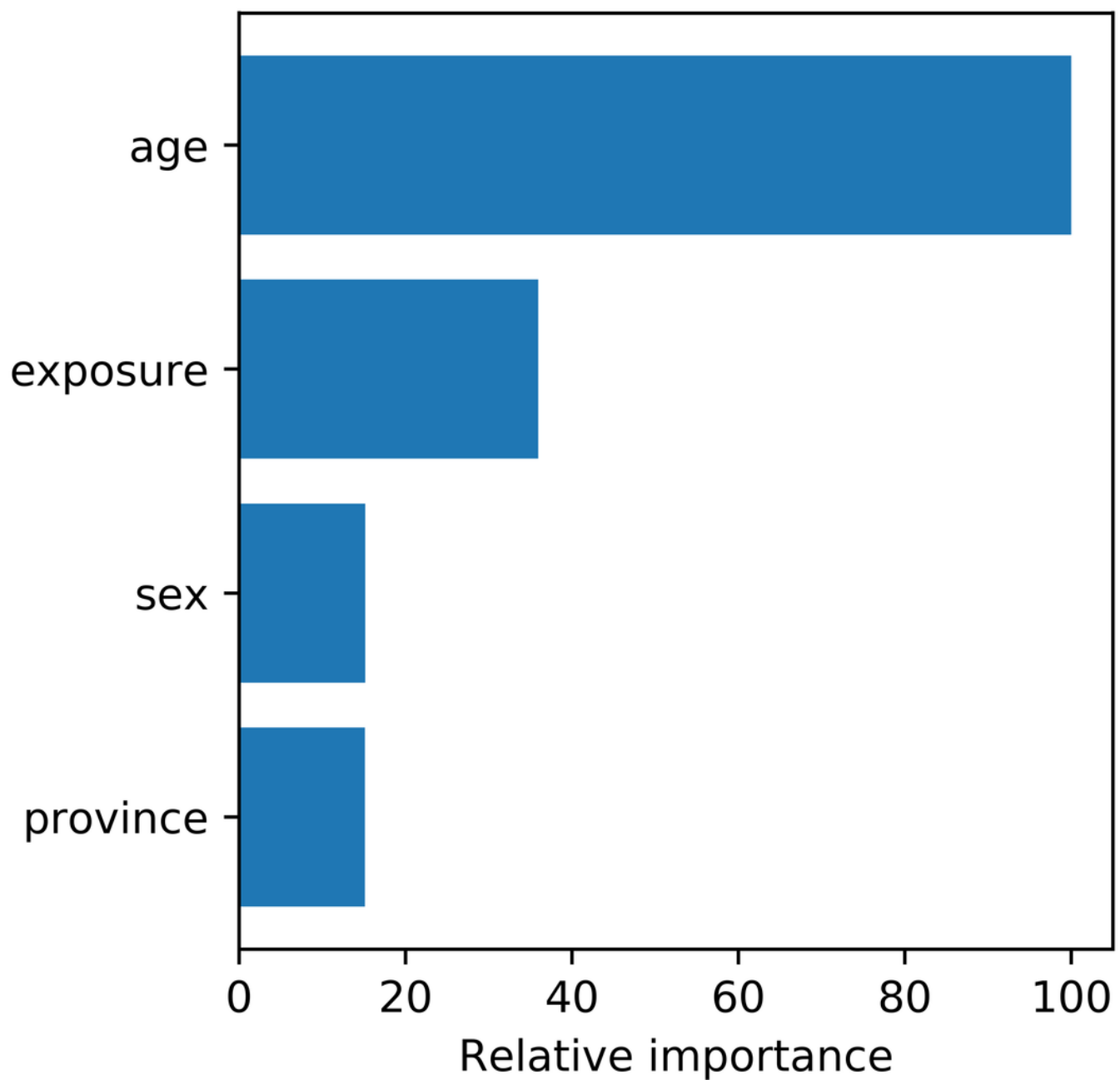
<sup>#</sup> SMOTE – Synthetic minor oversampling technique; <sup>\*</sup> ADASYN – Adaptive synthetic sampling

18

19

# Figure 1

Relative importance of predictors



# Figure 2

CoCoMORP online CoVID-19 Community Mortality Risk Prediction tool

## CoVID-19 Community Mortality Risk Prediction (CoCoMoRP) Tool

(Using Data from Korea Centers for Disease Control and Prevention)

Instructions: Select input values from drop-down menu in the boxes. Then, click the Submit button for predictions.

Sex

Male ▼

Age (Years)

80-89 ▼

Province

Seoul ▼

Exposure

Contact with patients ▼

Submit

## Prediction

Mortality risk:94.1%