

Title: Machine learning model estimating number of COVID-19 infection cases over coming 24 days in every province of South Korea (XGBoost and MultiOutputRegressor)

Authors: Yoshiro Suzuki\*<sup>1</sup>, Ayaka Suzuki\*<sup>2</sup>, Shun Nakamura, Toshiko Ishikawa, Akira Kinoshita

\*1: Tokyo Institute of Technology

\*2: Ernst & Young ShinNihon LLC (corresponding author)

## Abstract

We built a machine learning model (ML model) which input the number of daily infection cases and the other information related to COVID-19 over the past 24 days in each of 17 provinces in South Korea, and output the total increase in the number of infection cases in each of 17 provinces over the coming 24 days. We employ a combination of XGBoost and MultiOutputRegressor as machine learning model (ML model). For each province, we conduct a binary classification whether our ML model can classify provinces where total infection cases over the coming 24 days is more than 100. The result is Sensitivity =  $3/3 = 100\%$ , Specificity =  $11/14 = 78.6\%$ , False Positive Rate =  $3/11 = 21.4\%$ , Accuracy =  $14/17 = 82.4\%$ . Sensitivity = 100% means that we did not overlook the three provinces where the number of COVID-19 infection cases increased by more than 100. In addition, as for the provinces where the actual number of new COVID-19 infection cases is less than 100, the ratio (Specificity) that our ML model can correctly estimate was 78.6%, which is relatively high. From the above all, it is demonstrated that there is a sufficient possibility that our ML model can support the following four points. (1) Promotion of behavior modification of residents in dangerous areas, (2) Assistance for decision to resume economic activities in each province, (3) Assistance in determining infectious disease control measures in each province, (4) Search for factors that are highly correlated with the future increase in the number of COVID-19 infection cases.

## 1. Introduction

The new-type corona virus (SARS-CoV-2), which suddenly appeared at the end of 2019, caused global epidemic of COVID-19 in 2020. As of May 2020, industry, government, and academia from each country are cooperating to take measures to prevent the spreading of COVID-19 infection and develop therapeutic drugs.

This is the third time in this century that humans have been threatened by corona virus. The first was SARS in 2003, the second was MERS in 2012, and the third is COVID-19. Comparing SARS and MERS with COVID-19, it can be said that computer science including artificial intelligence (AI) and machine learning (ML) has made great progress. In fact, many computer science approaches have been developed to prevent the spreading of COVID-19 infection. For example, AI and ML are used for infection spreading analysis, drug discovery assistance, automatic diagnosis, diagnosis assistance, social trend analysis, and infection route analysis.

This paper proposes an ML technology for infection spreading analysis. To be more specific, we input the number of COVID-19 infection cases per day and other related information in each of 17 provinces in South Korea for the last 24 days, and output the total increase in the number of infection cases in each of 17 provinces over the coming 24 days. The result of conducting binary classification whether the total number of COVID-19 infection cases exceeds 100 over coming 24 days is Sensitivity =  $3/3 = 100\%$ , Specificity =  $11/14 = 78.6\%$ , False Positive Rate =  $3 / 11 = 21.4\%$ , Accuracy =  $14/17 = 82.4\%$ .

The existing ML for analyzing spreading of infectious diseases is introduced in Section 2. However, at present (May 10th, 2020), there is no ML that can estimate the number of COVID-19 infection cases in each province in South Korea with the above-mentioned high accuracy. From the above all, it is demonstrated that there is a sufficient possibility that our ML model

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**  
can support the following four points. (1) Promotion of behavior modification of residents in dangerous areas, (2) Assistance

for decision to resume economic activities in each region, (3) Assistance in determining infectious disease control measures in each region, (4) Search for factors that are highly correlated with the future increase in the number of COVID-19 infection cases.

## 2. Related works

In **Table 1**, we introduce the existing works about ML models which predict the spreading of COVID-19 infection. Also, we show the novelty points of our ML model compared to the existing works.

**Table 1** | Introduction of the existing works about ML models which predict the spreading of COVID-19 infection

Reference	Objectives	ML model	Results of only the optimal ML model
G. Pinter, et al., 2020 [4]	To predict the number of COVID-19 infection case and mortality rate.	Hybrid methods of ANN algorithm, i.e., MLP-ICA and ANFIS	<ul style="list-style-type: none"> <li>■ Total infection cases RMSE of MLP-ICA: 167.88 Determination coefficient of MLP-ICA: 0.9971</li> <li>■ Total mortality rate RMSE of MLP-ICA: 8.32 Determination coefficient of MLP-ICA: 0.9986</li> </ul>
A. A. Onovo et al., 2020 [5]	To reveal correlates and patterns of COVID-19 disease outbreak in sub-Saharan Africa (SSA).	Lasso Regression (Adaptive, cvplot and min BIC lasso)	<ul style="list-style-type: none"> <li>■ MSE of Adaptive lasso: <math>2.72 \times 10^{-28}</math></li> <li>■ R-squared of Adaptive lasso: 1.0000</li> </ul>
M. A. M. T. Baldé., 2020 [6]	To predict the future evolution of COVID-19 in Senegal.	SIR model	There was no clear numerical description.
A. Kumar et al., 2020 [7]	To forecast the possible rise in the number of COVID-19 infection cases by considering the daily data of new infection cases.	<ol style="list-style-type: none"> <li>1) Decision Tree algorithm</li> <li>2) Support Vector Machine algorithm</li> <li>3) Gaussian process regression</li> </ol>	There was no clear numerical description.
F. Sattler et al., 2020 [8]	To estimate the risk of COVID-19 infection transmission from the BLE signal strength measurements.	Linear regression	AUC: 0.96
S. Tiwari et al., 2020 [9]	To Predict the number of confirmed cases, recovered cases, and death cases of COVID-19 infection.	In this paper, their method is stated "machine learning approach" and there is no more detailed information.	<p>*We judged the numerical value from figures in this paper as follows.</p> <ul style="list-style-type: none"> <li>■ Confirmation case RMSE: about 2,350 MAE: about 1,850</li> <li>■ Death case RMSE: about 90 MAE: about 58</li> <li>■ Recovered case RMSE: about 36.5 MAE: about 24.5</li> </ul>
L. Magri et al., 2020 [10]	To provide quantitative estimates on the contact, recovery, death rates, etc of COVID-19.	Combination of an ODEsolver and SIRD model	There was no clear numerical description.
R. Sujatha et al., 2020 [11]	To predict the possible trends of COVID-19 impacts in India.	<ol style="list-style-type: none"> <li>1) Linear regression</li> <li>2) MLP</li> </ol>	There was no clear numerical description.
H. Jo et al., 2020 [12]	To analyze the novel coronavirus infection (COVID-19) spreading in South Korea.	SIR model with time-dependent parameters and deep learning	There was no clear numerical description.
M. Paggi, 2020 [13]	To propose a methodological contribution based on machine learning to foster the use of epidemiological models over pure data-driven best-fitting approaches and assess the reliability of their predictions.	A-SIR model	There was no clear numerical description.
M. P.N. et al., 2020 [14]	To present different predictive analytics techniques available for trend analysis, different models and algorithms and their comparison.	Prophet	There was no clear numerical description.
N. S. Punn et al., 2020 [15]	To understand everyday exponential behavior along with the prediction of future reachability of the COVID-19 across the nations.	<ol style="list-style-type: none"> <li>1) Support vector regression (SVR)</li> <li>2) Polynomial regression (PR)</li> <li>3) Standard deep neural network (DNN)</li> <li>4) Recurrent neural networks (RNN) using long short-term memory (LSTM)</li> </ol>	<ul style="list-style-type: none"> <li>■ Confirmed case RMSE of PR: 455.92</li> <li>■ Death case RMSE of PR: 117.94</li> <li>■ Recovered case RMSE of PR: 809.71</li> </ul>

Z. Yang et al., 2020 [16]	To predict the epidemic progression of COVID-19.	1) SEIR model 2) LSTM model, trained on the 2003 SARS data	There was no clear numerical description.
D. Liu et al., 2020 [17]	To forecast COVID-19 activity in Chinese provinces in real-time.	Augmented ARGONet	There was no clear numerical description.
S. Das et al., 2020 [18]	1) To estimate the basic reproduction number R0 at national and state level. 2) To predict COVID-19 infection cases ahead of time.	1) SIR model 2) Statistical Machine Learning (SML) model	■ Number of COVID-19 infection cases from March 25th to April 6th Actual number: 25,486 Prediction by SML model: 21,004
B. M. Ndiaye et al., 2020 [19]	To analyze the coronavirus pandemic in the real world.	1) SIR model 2) Prophet	There was no clear numerical description.
C. Zhou et al., 2020 [20]	1) To propose a high-resolution spatio-temporal (HiRES) model for the risk assessment of epidemic disease with human-to-human transmission. 2) To propose a personal infection risk scoring (HiRES-p) model to obtain objectively quantified risk of infection for every authorized individual.	1) HiRES risk map 2) HiRES-p	1) HiRES risk map: There was no clear numerical description. 2) HiRES-p: A total of 2,186 are detected out of 2,757 confirmed cases, indicating an overall detection rate around 80%.
P. Kumar et al., 2020 [21]	To predict some trajectories of COVID-19 over the coming days.	Auto-Regressive Integrated Moving Average (ARIMA) model	There was no clear numerical description.
R. M. Carrillo-Larco et al., 2020 [22]	To define data-driven clusters of countries.	Unsupervised machine learning algorithms (k-means)	There was no clear numerical description.

By comparing the above existing works with our ML model, we found the following novelty points of our ML model.

- The below two additional items were input to our ML model as binary label (0 or 1). The details are described in **Supplementary Table 2**.
  - (1) special measures taken by South Korean government to prevent the spreading of COVID-19 infection
  - (2) date of South Korean legislative election
- By using the XGBoost in combination with the MultiOutputRegressor, multiple objective variables, i.e., the number of COVID-19 infection cases in each of 17 provinces can be output.

### 3. Methods

- Our machine learning (ML) model is a combination of XGBoost [2] and MultiOutputRegressor [3].
- XGBoost

Official web sight of XGBoost: <https://xgboost.readthedocs.io/en/latest/>

XGBoost is a distributed gradient boosting library implements machine learning algorithms under the Gradient Boosting framework. The original XGBoost is a regressor that can only estimate (output) a single objective variable (explained variable). However, by combining XGBoost and MultiOutputRegressor, we can estimate multiple objective variables (multiple explained variables).

- MultiOutoutRegressor

Official web sight of MultiOutoutRegressor: <https://scikit-learn.org/stable/modules/generated/sklearn.multioutput.MultiOutputRegressor.html>

MultiOutputRegressor is one of Scikit-Learn modules that can estimate multiple objective variables

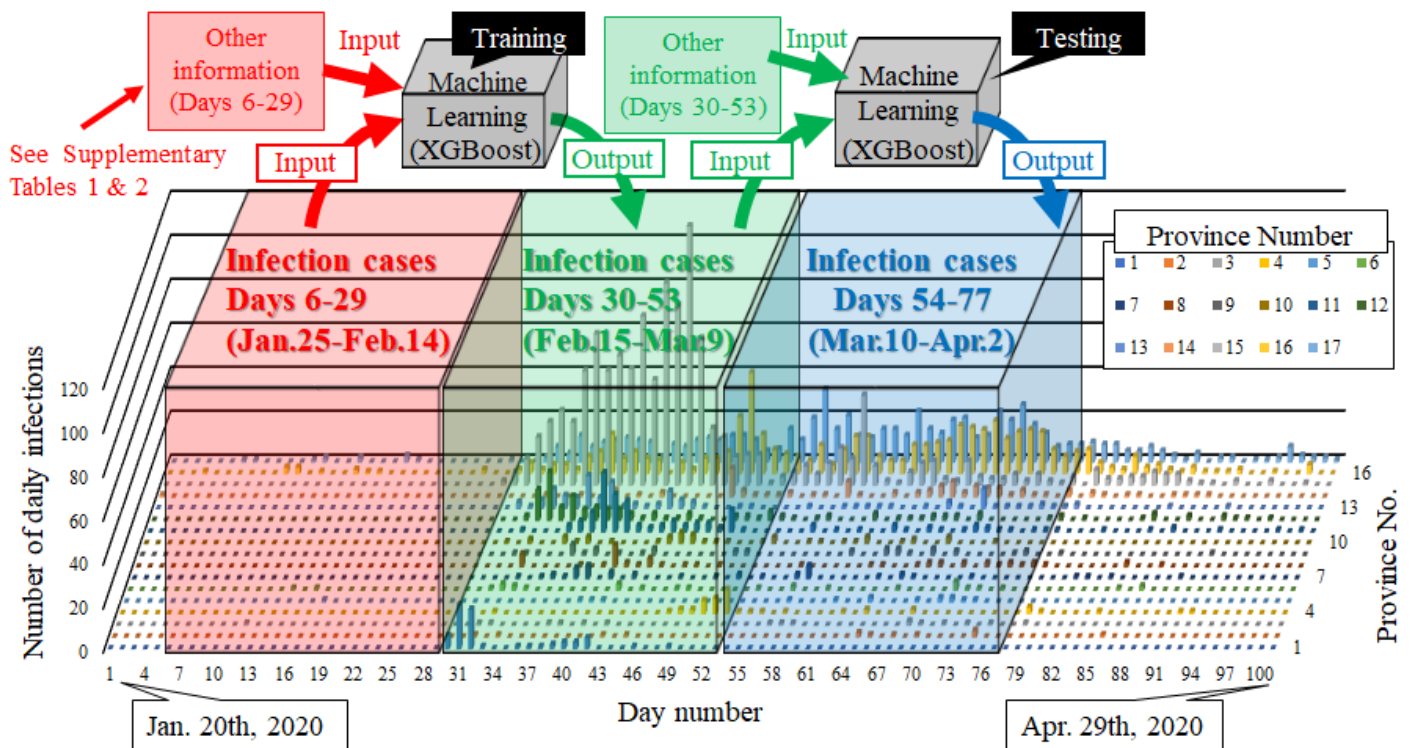
- Task of our ML model: regression of daily infection cases in each of 17 provinces over the coming 24 days
- Loss function for training our ML model: Mean squared error (MSE) of numbers of daily infection cases in 17 provinces of South Korea over the coming 24 days
- Evaluation index for testing: binary classification (whether the ML model can classify provinces where total infection cases over the coming 24 days is more than 100).
- Hyper parameters:
  - I. max\_depth: 6

II.  $n_{estimators}$ : 30

III. The rest of all hyper parameters: default (standard)

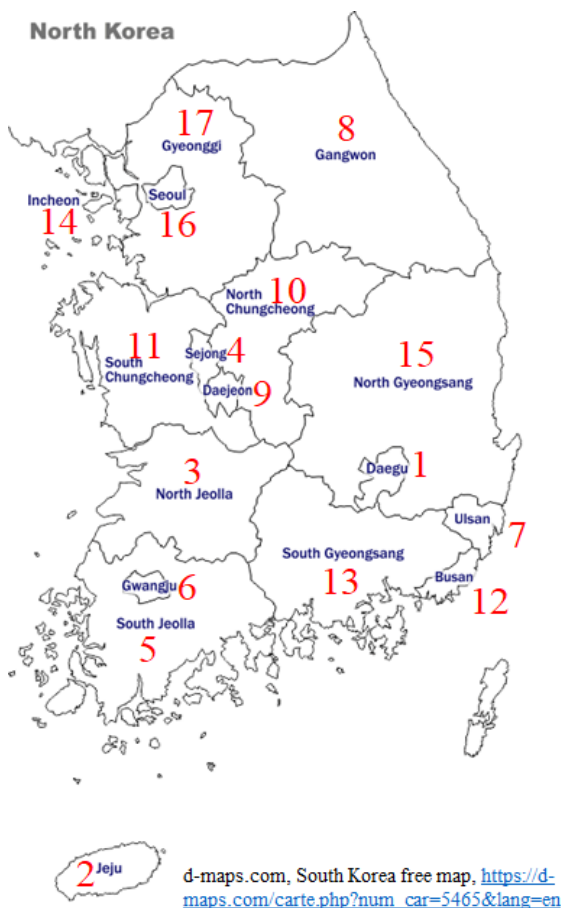
#### 4. Dataset

- Dataset name: Data Science for COVID-19 in South Korea (DS4C)
- Contents: Information of 3,288 COVID-19 infection cases confirmed in South Korea from January 20th to April 29th, 2020. We decided to use this dataset since it contained more detailed information compared to other countries' data. Note that we used only 3,285 cases apart from 3 cases with no mention of infection date.
- Dataset downloaded from: <https://www.kaggle.com/kimjihoo/coronavirusdataset>
- Dataset reported by: Korea Centers for Disease Control and Prevention (KCDC) and 17 provinces in South Korea
- License of the dataset: CC BY-NC-SA 4.0
- Training set: Information related to COVID-19 over the past 24 days (Jan. 25th - Feb. 14th) is input. Total increase in infection cases over the coming 24 days (Feb. 15th - Mar. 9th) in each of 17 provinces is output.
- Testing set: Information related to COVID-19 over the past 24 days (Feb. 15th - Mar. 9th) is input. Total increase in infection cases over the coming 24 days (Mar. 10th - Apr. 2th) in each of 17 provinces is output.
- COVID-19 information input to our ML model is shown in **Supplementary Tables 1 and 2**. In addition to the two tables, the number of daily infection cases in each province of South Korea over the past 24 days (e.g., the red box in **Fig. 1**) is input.



**Fig. 1** | The COVID-19 dataset of South Korea used for training and testing our ML model (a combination of XGBoost and MultiOutputRegressor). This figure shows the number of daily infection cases in each of 17 provinces (South Korea, from January 20th to April 29th in 2020). January 20th is equal to Day 1 and April 29th is equal to Day 101. In training, the information related to COVID-19 over the past 24 days (Jan. 25th - Feb. 14th) is input. Then, the total increase in the number of infection cases over the coming 24 days (Feb. 15th - Mar. 9th) in each of 17 provinces is output. In testing, the information related to COVID-19 over the past 24 days (Feb. 15th - Mar. 9th) is input. Then, the total increase in the number of infection

cases over the coming 24 days (Mar. 10th - Apr. 2th) in each of 17 provinces is estimated. The estimation performances of our ML model for the test set are shown in Fig. 3.



Province No.	Province name
1	Daegu
2	Jeju-do
3	Jeollabuk-do
4	Sejong
5	Jeollanam-do
6	Gwangju
7	Ulsan
8	Gangwon-do
9	Daejeon
10	Chungcheongbuk-do
11	Chungcheongnam-do
12	Busan
13	Gyeongsangnam-do
14	Incheon
15	Gyeongsangbuk-do
16	Seoul
17	Gyeonggi-do

Fig. 2 | Province number allocated for each of 17 Provinces in South Korea.

## 5. Results and discussion

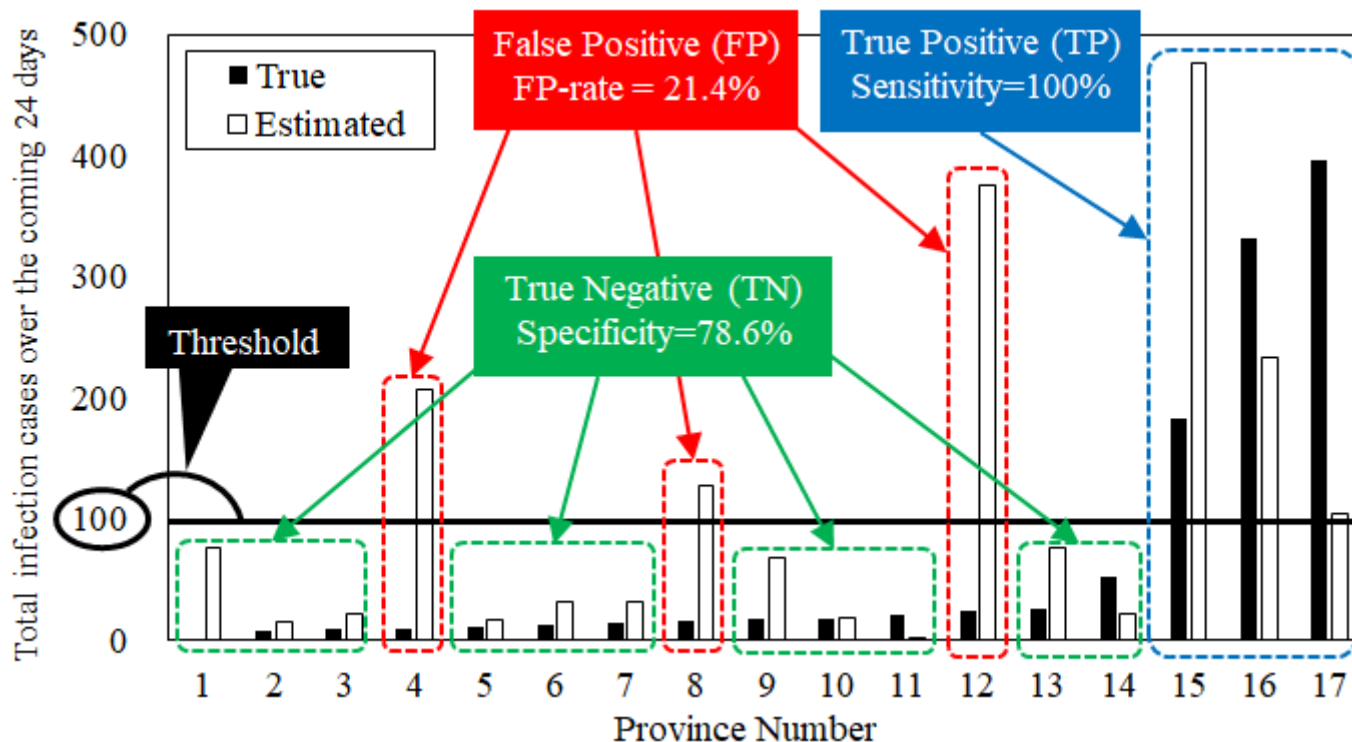


Fig. 3 | Performance validation (test) result of our ML model. The number of daily infection cases and the other information



related to COVID-19 over the past 24 days (Feb. 15th - Mar. 9th) in each of 17 provinces is input. Then, the total increase in the number of infection cases over the coming 24 days (Mar. 10th - Apr. 2th) in each of 17 provinces is estimated. The black bar graph is the true value, and the white bar graph is the estimated value.

The accuracy of the binary classification whether our ML model can classify provinces where total infection cases over the coming 24 days is more than 100 is as follows.

$$\text{Sensitivity} = TP / (TP + FN) = 3/3 = 100\% \text{ (i.e., true positive rate, recall, hit rate)}$$

$$\text{Specificity} = TN / (TN + FP) = 11/14 = 78.6\% \text{ (i.e., true negative rate)}$$

$$\text{False Positive Rate} = FP / (TN + FP) = 21.4\%$$

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN) = 14/17 = 82.4\%$$

where,

TP = True Positive

TN = True Negative

FP = False Positive

FN = False Negative

**Sensitivity=100%** of the binary classification means that we did not overlook the three provinces where the number of COVID-19 infection cases increased by more than 100. In addition, as for the provinces where the actual number of new COVID-19 infection cases is less than 100, the ratio (**Specificity**) that our ML model can correctly estimate is **78.6%**, which is relatively high.

Next, we evaluate the accuracy of our ML model from the perspective of the regression task, not the binary classification. The ratio that our ML model can estimate the increase of COVID-19 infection cases in each province over the coming 24 days when the maximum permissible error is set to 100 infection cases is as follows.

$$\text{Another accuracy when the maximum permissible error is set to 100 infection cases} = 12/17 = 70.6\%$$

From the above all, it is demonstrated that there is a sufficient possibility that our ML model can support the following four points.

- (1) Promotion of behavior modification of residents in dangerous areas
- (2) Assistance for decision to resume economic activities in each region
- (3) Assistance in determining infectious disease control measures in each region
- (4) Search for factors that are highly correlated with the future increase in the number of COVID-19 infection cases.

#### Notes for the performance of our ML model

- It is pointed out that the actual number of positives may be higher than this dataset because PCR tests may be insufficient. If this point is correct, there is a possibility that the performance of our ML model (**sensitivity = 100%**, **specificity = 78.6%**, **false Positive Rate = 21.4%**) may change.
- There is a possibility that the current input information may not contain important information. For example, in this method, population, population density, temperature, humidity, weather, regulation of economic activity, degree of land (urban area, depopulated area, industrial area, forest area), etc. of each province are not input. By inputting this information, the estimation performance might be improved.
- In this paper, both input and output of our ML model are set to 24 days, but this is not always optimal. Performance may change by changing the number of days.

## 6. Conclusions

We built a machine learning model (ML model) which input the number of daily infection cases and the other information related to COVID-19 over the past 24 days in each of 17 provinces in South Korea, and output the total increase in the number of infection cases in each of 17 provinces over the coming 24 days. We employ a combination of XGBoost and MultiOutputRegressor as machine learning model (ML model). We trained the ML model by setting loss function as the number of daily infections in the 17 provinces (i.e., regression task). We tested the ML model in terms of binary classification (whether the ML model can classify provinces where total infection cases over the coming 24 days is more than 100). As a result, Sensitivity =  $3/3 = 100\%$ , Specificity =  $11/14 = 78.6\%$ , False Positive Rate =  $3/11 = 21.4\%$ , Accuracy =  $14/17 = 82.4\%$ . Sensitivity = 100% means that we did not overlook the three provinces where the number of COVID-19 infection cases increased by more than 100. In addition, as for the provinces where the actual number of new COVID-19 infection cases was less than 100, the ratio (Specificity) that our ML model could correctly estimate is 78.6%, which is relatively high. From the above all, it is demonstrated that there is a sufficient possibility that our ML model can support the following four points.

- (1) Promotion of behavior modification of residents in dangerous areas
- (2) Assistance for decision to resume economic activities in each region
- (3) Assistance in determining infectious disease control measures in each region
- (4) Search for factors that are highly correlated with the future increase in the number of COVID-19 infection cases.

## 7. References

- [1] Dataset: <https://www.kaggle.com/kimjihoo/coronavirusdataset>
- [2] XGBoost: <https://xgboost.readthedocs.io/en/latest/>
- [3] MultiOutputRegressor: <https://scikit-learn.org/stable/modules/generated/sklearn.multioutput.MultiOutputRegressor.html>
- [4] Pinter, G.; Felde, I.; Mosavi, A.; Ghamisi, P.; Gloaguen, R. COVID-19 Pandemic Prediction for Hungary: A Hybrid Machine Learning Approach. Preprints 2020, 2020050031 (doi: 10.20944/preprints202005.0031.v1).
- [5] A. Onovo, A. Atobatele, A. Kalaiwo, C. Obanubi, E. James, P. Gado, G. Odezugo, D. Magaji, D. Ogundehin, M. Russell, Using Supervised Machine Learning and Empirical Bayesian Kriging to reveal Correlates and Patterns of COVID-19 Disease outbreak in sub-Saharan Africa: Exploratory Data Analysis, <https://www.medrxiv.org/node/78792.external-links.html>
- [6] M. A. M. T. Baldé, Fitting SIR model to COVID-19 pandemic data and comparative forecasting with machine learning, <https://www.medrxiv.org/content/10.1101/2020.04.26.20081042v1>
- [7] A. Kumar, Farhan M. Khan, R. Gupta, H. Puppala, Preparedness and Mitigation by projecting the risk against COVID-19 transmission using Machine Learning Techniques, <https://www.medrxiv.org/content/10.1101/2020.04.26.20080655v1>
- [8] F. Sattler, J. Ma, P. Wagner, D. Neumann, M. Wenzel, R. Schäfer, W. Samek, Klaus-Robert Müller, T. Wiegand, Risk Estimation of SARS-CoV-2 Transmission from Bluetooth Low Energy Measurements, <https://arxiv.org/abs/2004.11841>
- [9] S. Tiwari, S. Kumar, K. Guleria, Outbreak Trends of Coronavirus Disease–2019 in India: A Prediction, <https://www.cambridge.org/core/journals/disaster-medicine-and-public-health-preparedness/article/outbreak-trends-of-coronavirus-disease2019-in-india-a-prediction/76090B13B7FDD2C96920A81CAF608264/core-reader>
- [10] L. Magri, N. A. K. Doan, First-principles machine learning modelling of COVID-19, arXiv:2004.09478
- [11] Sujatha, R.; Chatterjee, Jyotir; Hassaniien, Aboul ella (2020): A machine learning methodology for forecasting of the COVID-19 cases in India. TechRxiv. Preprint. <https://doi.org/10.36227/techrxiv.12143685.v1>

- [12] H. Jo, H. Son, S. Y. Jung, H. J. Hwang, Analysis of COVID-19 spreading in South Korea using the SIR model with time-dependent parameters and deep learning, <https://www.medrxiv.org/content/10.1101/2020.04.13.20063412v1>
- [13] M. Paggi, Simulation of Covid-19 epidemic evolution: are compartmental models really predictive?, <https://arxiv.org/abs/2004.08207>
- [14] Mahalle, P.N.; Sable, N.P.; Mahalle, N.P.; Shinde, G.R. Predictive Analytics of COVID-19 Using Information, Communication and Technologies. Preprints 2020, 2020040257 (doi: 10.20944/preprints202004.0257.v1).
- [15] N. S. Punna, S. K. Sonbhadra, S. Agarwal, COVID-19 Epidemic Analysis using Machine Learning and Deep Learning Algorithms, <https://www.medrxiv.org/content/10.1101/2020.04.08.20057679v1>
- [16] Z. Yang, Z. Zeng, K. Wang, Sook-San Wong, W. Liang, M. Zanin, P. Liu, X. Cao, Z. Gao, Z. Mai, J. Liang, X. Liu, S. Li, Y. Li, F. Ye, W. Guan, Y. Yang, F. Li, S. Luo, Y. Xie, B. Liu, Z. Wang, S. Zhang, Y. Wang, N. Zhong, J. He, <http://jtd.amegroups.com/article/view/36385/html>
- [17] D. Liu, L. Clemente, C. Poirier, X. Ding, M. Chinazzi, J. T Davis, A. Vespignani, M. Santillana, A machine learning methodology for real-time forecasting of the 2019-2020 COVID-19 outbreak using Internet searches, news alerts, and estimates from mechanistic models, <https://arxiv.org/abs/2004.04019>
- [18] S. Das, Prediction of COVID-19 Disease Progression in India: Under the Effect of National Lockdown, <https://arxiv.org/abs/2004.03147>
- [19] B. Mbaye Ndiaye, L. Tendeng, D. Seck, Analysis of the COVID-19 pandemic by SIR model and machine learning techniques for forecasting, <https://arxiv.org/abs/2004.01574>
- [20] C. Zhou, W. Yuan, J. Wang, H. Xu, Y. Jiang, X. Wang, Q. H. Wen, P. Zhang, Detecting Suspected Epidemic Cases Using Trajectory Big Data, <https://arxiv.org/abs/2004.00908>
- [21] P. Kumar, H. Kalita, S. Patariya, Y. D. Sharma, C. Nanda, M. Rani, J. Rahmani, A. S. Bhagavathula, Forecasting the dynamics of COVID-19 Pandemic in Top 15 countries in April 2020: ARIMA Model with Machine Learning Approach, <https://www.medrxiv.org/content/10.1101/2020.03.30.20046227v2>
- [22] Carrillo-Larco RM and Castillo-Cara M. Using country-level variables to classify countries according to the number of confirmed COVID-19 cases: An unsupervised machine learning approach [version 1; peer review: 1 approved]. Wellcome Open Res 2020, 5:56 (<https://doi.org/10.12688/wellcomeopenres.15819.1>)

## 8. Declaration of competing interest

On behalf of all authors, the corresponding author states that there is no conflict of interest.

## 9. Contributions of authors

Ayaka Suzuki, the corresponding author, contributed to the software engineering and developed all the codes.

Yoshiro Suzuki managed the project, analyzed the results, and wrote the paper.

The other authors helped designing this research project. Their contributions are almost equal to each other.

**Supplementary Table 1** | Natural number (NOT label) information input to our ML model. In addition to the information in this table, both information shown in **Supplementary Table 2** and information depicted in **Fig. 1** are input to our ML model.

Cumulative number of patients all over South Korea
Number of patients all over South Korea by age par day
Number of patients all over South Korea by gender per day (male, female, the other)



Number of overseas inflow patients	
Number of those who have contact(s) with other patient(s)	
Number of isolated patients par day	
Number of released patients par day	
Number of deceased patients par day	
Daily increase in patients who had visited:	Eunpyeong St. Mary's Hospital
	Cheongdo Daenam Hospital
	Bonghwa Pureun Nursing Home
	Gyeongsan Seorin Nursing Home
	Geochang Church
	Shincheonji Church
	Dongan Church
	Onchun Church
	River of Grace Community Church
	Pilgrimage to Israel Milal Shelter
	Suyeong-gu Kindergarten
	Seongdong-gu APT
	Ministry of Oceans and Fisheries
	Gym facility in Cheonan
	Gym facility in Sejong
	Guro-gu Call Center
	Gyeongsan Cham Joeun Community Center
Gyeongsan Jeil Silver Town	
Changnyeong Coin Karaoke	
etc	

**Supplementary Table 2** | Label (NOT variable) information input to our ML model. In addition to the information in this table, both information shown in **Supplementary Table 1** and information depicted in **Fig. 1** are input to our ML model.

Items	Label (0 or 1)	Information sources
Korean government raised the national alert level to Yellow (level 2).	Jan. 20 - Jan. 27: 1 After Jan. 28: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Korean government raised its infectious disease alert level to Orange (level 3).	Jan. 28 – Feb. 22: 1 After Feb. 23: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Korea began banning entry of all foreign nationals who had been to China's Hubei province in the past two weeks.	After Feb. 4: 1 Jan. 20 - Feb. 3: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Korean government requested citizens to refrain from traveling to six countries and regions such as Japan.	After Feb. 11: 1 Jan. 20 - Feb. 10: 0	<a href="https://en.yna.co.kr/view/AEN20200211006000320">https://en.yna.co.kr/view/AEN20200211006000320</a>

Korean government declared 'Special Management Region' in Daegu and Cheongdo.	After Feb. 21: 1 Jan. 20 - Feb. 20: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government raised its infectious disease alert level to Red(level 4) and ordered schools to start the new semester one week later on Mar 9, from Mar 2.	After Feb. 23: 1 Jan. 20 - Feb. 22: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Social distancing approaches have been launched.	After Mar. 1: 1 Jan. 20 - Feb. 29: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government divided confirmed patients into four groups and only the sickest and elderly were sent to hospitals. The young and asymptomatic went to dormitories.	After Mar. 1: 1 Jan. 20 - Feb. 29: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government declared 'Special Management Region' in Gyeongsan.	After Mar. 5: 1 Jan. 20 - Mar. 4: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government permitted the purchase of masks at the mask sales office once a week based on the number at the end of the citizen's birth year (AD).	After Mar. 9: 1 Jan. 20 - Mar. 8: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government applied special entry procedures for those from Japan.	After Mar. 9: 1 Jan. 20 - Mar. 8: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Pilot operation of COVID-19 Epidemiological Investigation Support System was implemented.	Mar. 16 - Mar. 25: 1 Except for the above: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Special entry procedure expanded to all income travelers.	After Mar. 16: 1 Jan. 20 - Mar. 15: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Korean government applied special entry procedures for all foreigners.	After Mar. 19: 1 Jan. 20 - Mar. 18: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
Testing all income traveler from Europe.	After Mar. 22: 1 Jan. 20 - Mar. 21: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
COVID-19 Epidemiological Investigation Support System was officially launched.	After Mar. 26: 1 Jan. 20 - Mar. 25: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
All inbound travelers since April 1 have been required to install the application at entry to monitor symptoms of inbound travelers while also providing them prompt medical advice.	After Apr. 1: 1 Jan. 20 - Mar. 31: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>	
South Korean legislative election was held.	Apr. 15: 1 Except for the above: 0	<a href="https://en.wikipedia.org/wiki/2020_South_Korean_legislative_election">https://en.wikipedia.org/wiki/2020_South_Korean_legislative_election</a>	
Jeju-do has launched a campaign on February 24 to have its citizens voluntarily record their whereabouts on their smartphones using Google Timeline.	Province No. 1	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 2	After Feb. 24: 1 Jan. 20 - Feb. 23: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 3	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 4	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 5	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 6	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 7	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 8	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 9	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
	Province No. 10	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>

Province No. 11	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 12	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 13	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 14	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 15	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 16	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>
Province No. 17	All: 0	<a href="https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf">https://www.ictworks.org/wp-content/uploads/2020/04/Korea-flattening-covid-19-curve.pdf</a>