

Comparison of Decision Tree and Linear Regression Algorithms of COVID-19 Prediction in Indonesia

Darwin^{1)*}, Dwiky Christian²⁾, Wilson Chandra³⁾, Marlince Nababan⁴⁾

¹⁾²⁾³⁾⁴⁾ Universitas Prima Indonesia, Indonesia

¹⁾darwinblacks@email.com, ²⁾dwikyc70@gmail.com, ³⁾wilsonchmc@gmail.com, ⁴⁾marlince@unprimdn.ac.id

ABSTRACT

COVID-19 is a disease that was first discovered in Wuhan, China and caused the 2019-2020 *coronavirus* pandemic. This virus can cause respiratory tract infections such as flu when infecting humans. According to Ministry of Health of the Republic of Indonesia, the number of confirmed cases of COVID-19 in Indonesia at March 2021 is 1,511,712 with 40,858 deaths and 1,348,330 recovered. For that, Indonesia is declared to have the highest confirmed cases in ASEAN. Several studies have been carried out to handle some cases by using the data mining techniques such as Decision Tree or Linear Regression algorithm, as example to classify the respiratory diseases and predict pregnancy hypertension. In this study, we tried to analyze COVID-19 cases in Indonesia and conducted an experiment of predicting COVID-19 new cases with the Decision Tree (CART) and Linear Regression algorithms. Then we will compare the values of these two algorithms by using R^2 Score to evaluate the prediction performance. The results of this analysis state that DKI Jakarta province has the highest number of positive cases, cures and deaths in Indonesia. The value of the comparison results from the R^2 Score obtained in the Decision Tree algorithm reached 95.69% (training) and 92.15% (testing) while the Linear Regression algorithm reached 79.93% (training) and 77.25% (testing).

Keywords: CART, COVID-19, Data Mining, Decision Tree, Linear Regression.

INTRODUCTION

COVID-19 is a disease caused by a descendant of the new *coronavirus*, 'CO' is taken from the *corona*, 'VI' is virus, and 'D' is disease (Bender, 2020). The disease was first discovered in Wuhan in December 2019, the Capital of China's Hubei Province, and has spread globally around the world since, resulting in the 2019-2020 *coronavirus* pandemic. The World Health Organization (WHO) declared the 2019-2020 *coronavirus* outbreak an International Public Health Emergency on 30 January 2020, and a pandemic on 11 March 2020 (Yezli & Khan, 2020).

The World Health Organization (WHO) named the new virus *Severe Acute Respiratory Syndrome Coronavirus-2* (SARS-CoV-2) and officially designated the disease as *Novel Coronavirus* in humans as *Coronavirus Diseases 2019* or better known as COVID-19. When attacking humans, this virus causes respiratory tract infections such as the flu, to more severe diseases such as Middle East Respiratory Syndrome (MERS) and Severe Acute Respiratory Syndrome (SARS). So far, it has been confirmed that 222 countries have been infected, with confirmed cases of 127,877,462 with a death toll of 2,796,561 (Nabila et al., 2021).

Every day the spread of COVID-19 in Indonesia continues to increase, people are asked to do social distancing to break the chain of the spread of COVID-19 which is spread in various regions in Indonesia (Sindi et al., 2020). The first COVID-19 was reported to have entered Indonesia on March 2, 2020, with 2 confirmed cases. The vast territory of Indonesia allows the need for grouping of parts based on regions in Indonesia (Dwitri et al., 2020). According to data from the Ministry of Health of the Republic of Indonesia as of March 31, 2021, the number of confirmed cases of COVID-19 is 1,511,712 with 40,858 deaths and 1,348,330 recovered. Indonesia is the country with the highest confirmed cases in ASEAN (Annisa, 2021).

LITERATURE REVIEW

Currently the coronavirus pneumonia (COVID-19) caused by it has spread in China and even the world. Accurate

* Corresponding author



This is an Creative Commons License This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License.

simulation and effective prediction with mathematical models on this pandemic development trend have become very important for pandemic prevention and control (Alves et al., 2021; Vadyala et al., 2021; Xiao et al., 2021; Yoo et al., 2020). Several studies have been conducted to predict new cases of COVID-19 in Indonesia using existing algorithms such as Decision Tree and Linear Regression. There are also several cases that have been carried out in this algorithm such as Classification of Respiratory Diseases (Pramadhani & Tedy, 2014), Prediction of Pregnancy Hypertension (Wulandari, 2016) and Classification of Student Data (Sutoyo, 2018). With these cases, it can be said that the algorithm used is quite popular in research circles.

The Decision Tree and Linear Regression algorithms have been used for various classifications and predictions, but currently there is no comparison of the effectiveness of these two algorithms in the case of predicting the spread of COVID-19 in Indonesia. Therefore, a comparison is needed between the value of the Decision Tree algorithm and Linear Regression to be able to test its prediction performance in the case of predicting the spread of COVID-19 in Indonesia.

Based on this background, the researchers are interested in conducting research entitled "Comparison of the Values of the Decision Tree Algorithm and Linear Regression to the Predicted Case of the Spread of COVID-19 in Indonesia".

METHOD

In this study, researchers collected the data needed for the data mining process. The data used in this study was sourced from *Kaggle*, namely the Indonesian COVID-19 dataset (<https://www.kaggle.com/hendratno/covid19-indonesia>). *Kaggle* is a website that provides various datasets used for data science purposes (Sodik et al., 2020). The data taken has a total of 16283 rows and 42 attributes/columns. The data from this selection will be used in the data mining process.

Before entering the data processing stage, it is necessary to conduct an analysis of the data that has been taken to be able to check the clarity and completeness of the data in the study. The purpose of data analysis is to determine the quality of the data to be tested and to draw hypotheses related to the data to be tested. The steps in analyzing the data are: 1) Data Preprocessing, at this stage is the stage for cleaning the results of data selection. Before the data mining process can be carried out, it is necessary to carry out a data cleaning stage. The cleaning of the data carried out in this study is to remove data attributes that are not needed in the prediction process later and to delete data that is *null* or N/A. 2) Data Transformation, the next stage is the data transformation stage. At this stage, the data normalization process will be carried out. The purpose of this normalization is so that the data, which are especially numeric, are in the range 0 to 1 so that the data distribution is not too far away. Meanwhile, categorical data will be converted to numeric so that it can be matched in the next data mining process. 3) Data Mining, this stage is the process of looking for interesting patterns or information in the data that has been selected and transformed. The algorithm used in this research is the Decision Tree and Linear Regression algorithm. This algorithm predicts new cases of COVID-19 that will come based on the total number of ongoing cases, deaths and recovered of COVID-19 in Indonesia.

The Python experiment is the implementation stage of the analysis and prediction of new COVID-19 cases in Indonesia where testing will be carried out in the form of a program written in Python version 3.7.12 on the *Google Colab* platform. Python is claimed to be a programming language that combines capabilities, capabilities, with a clear code syntax, also equipped with a large and comprehensive standard library functionality (Baharuddin et al., 2019).

System testing is the stage of executing the software system to determine whether the application of data mining matches the system specifications. Testing this system is in the form of verifying whether the program that has been made can be run according to the specifications and designs that have been made.

The evaluation stage is the final stage in this research. This stage is carried out to obtain the quality and effectiveness of the model whether it has met the objectives and has solved the problem in this research, as well as making decisions regarding the use of the results from data mining.

* Corresponding author



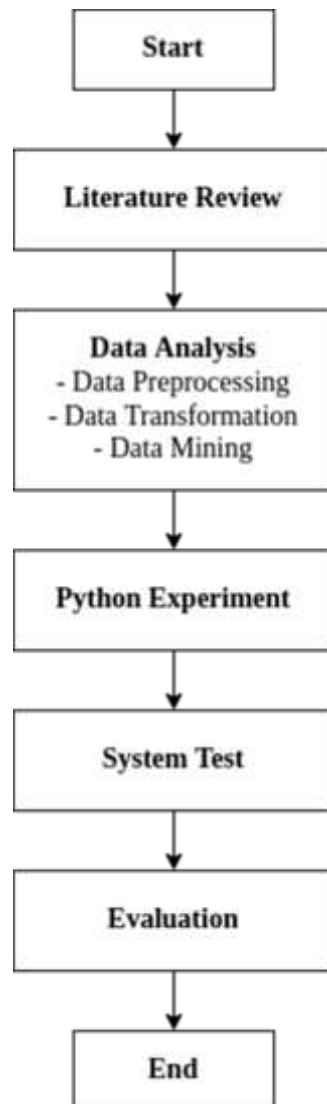


Fig. 1 Research Method
Diagram

RESULT

In this study, the data used in the analysis of COVID-19 is from January 8, 2020 to July 9, 2021. There are several attributes used in the analysis and prediction of COVID-19 in Indonesia, such as Date, Location, New Cases, New Recovered, Total Cases, Total Deaths, Total Recovered, Total Active Cases. The following is a sample of the COVID-19 dataset in Indonesia from January 2020 to July 2021.

* Corresponding author



This is an Creative Commons License This work is licensed under a
Creative Commons Attribution-NoDerivatives 4.0 International License.

Table 1 Sample Dataset of COVID-19 in Indonesia

| Date | Location | New Cases | New Recovered | Total Cases | Total Deaths | Total Recovered | Total Active Cases |
|------------|------------|-----------|---------------|-------------|--------------|-----------------|--------------------|
| 2020-05-12 | Jawa Barat | 54 | 4 | 1658 | 116 | 316 | 1226 |
| 2020-07-07 | Maluku | 26 | 33 | 830 | 17 | 414 | 399 |
| 2020-10-31 | Riau | 127 | 231 | 14798 | 336 | 11208 | 3255 |
| 2021-01-04 | Papua | 36 | 124 | 13306 | 150 | 7623 | 5533 |
| 2021-02-23 | Sum. Utara | 139 | 107 | 23894 | 820 | 20684 | 2390 |

COVID-19 Analysis per Province

In this section, the analysis conducted on COVID-19 cases in Indonesia is by province. The following are the results of the analysis of COVID-19 in Indonesia by province which can be seen in the images below.

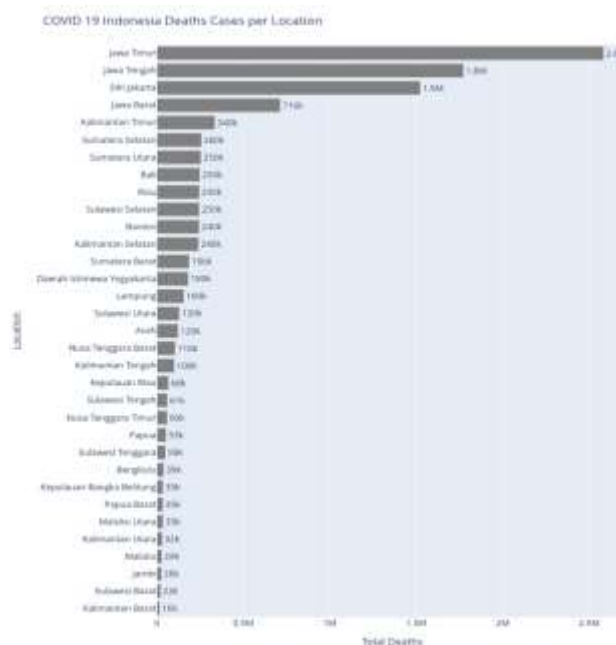


Fig. 2 COVID-19 Death Cases per Province

* Corresponding author



This is an Creative Commons License This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

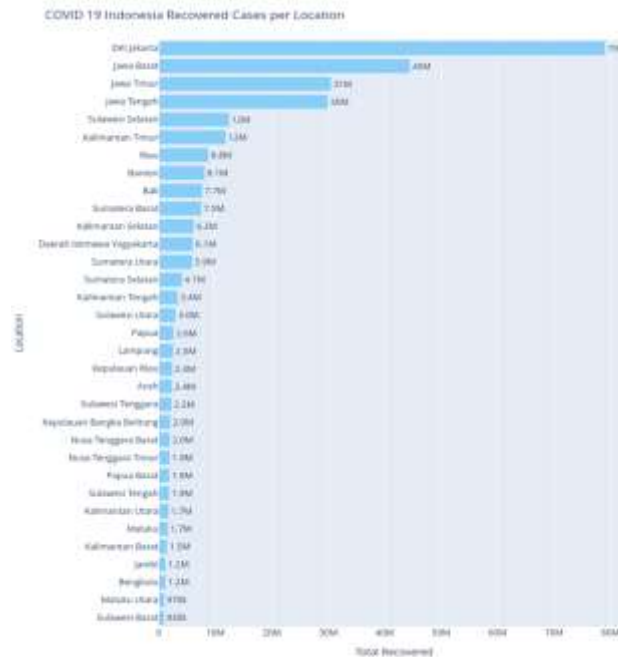


Fig. 4 COVID-19 Recovered Cases per Province

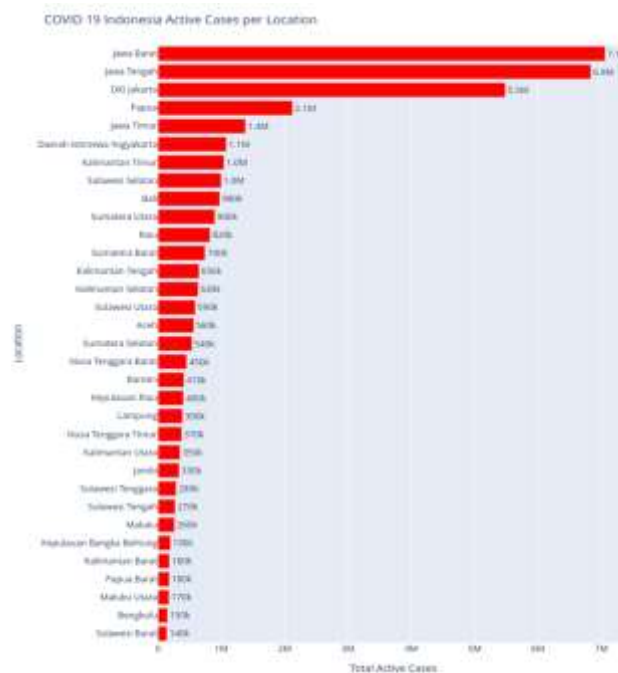


Fig. 3 COVID-19 Positive Cases per Province

* Corresponding author



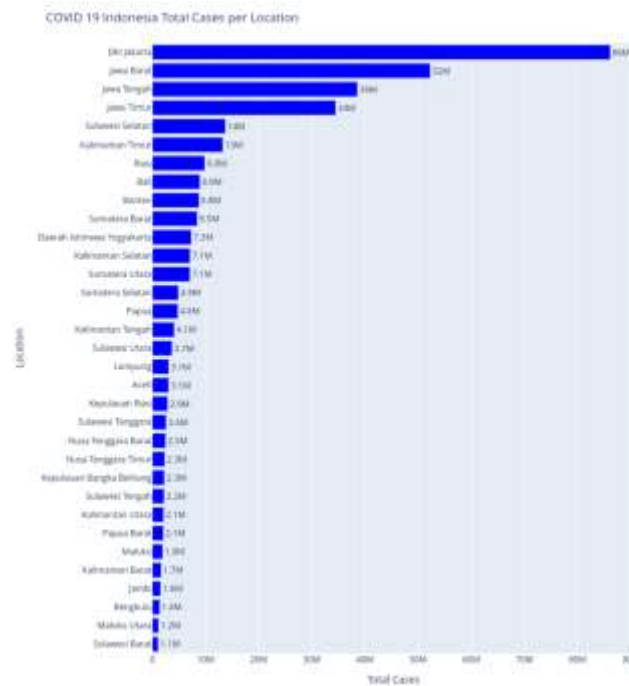
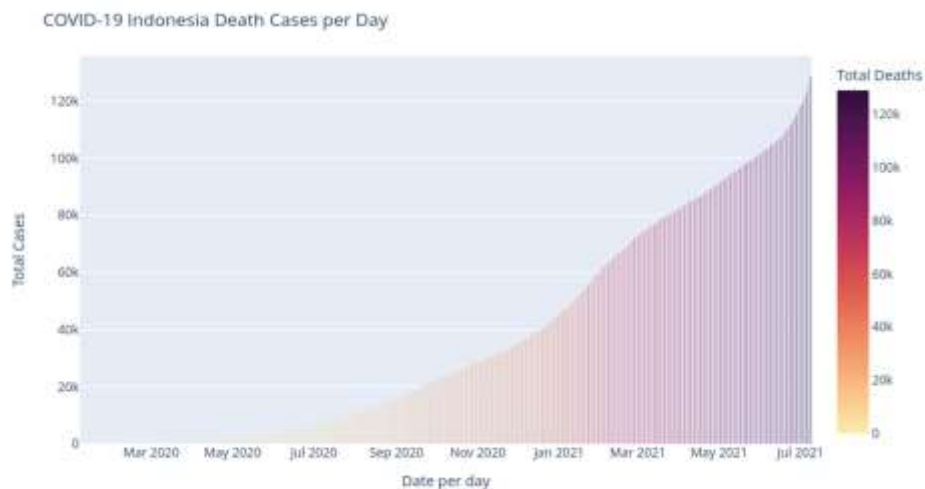


Fig. 5 COVID-19 All Cases per Province

COVID-19 Analysis per Day

In this section, the analysis carried out on COVID-19 cases in Indonesia is by day. The following are the results of the analysis of COVID-19 in Indonesia by day which can be seen in the images below.



* Corresponding author

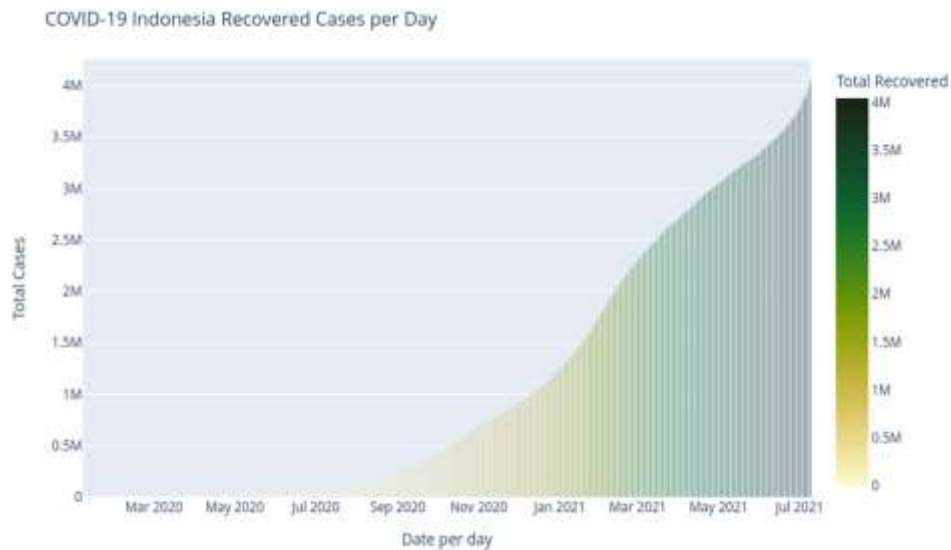


Fig. 7 COVID-19 Recovered Cases per Day



Fig. 8 COVID-19 Positive Cases per Day

* Corresponding author



This is an Creative Commons License This work is licensed under a
Creative Commons Attribution-NoDerivatives 4.0 International License.

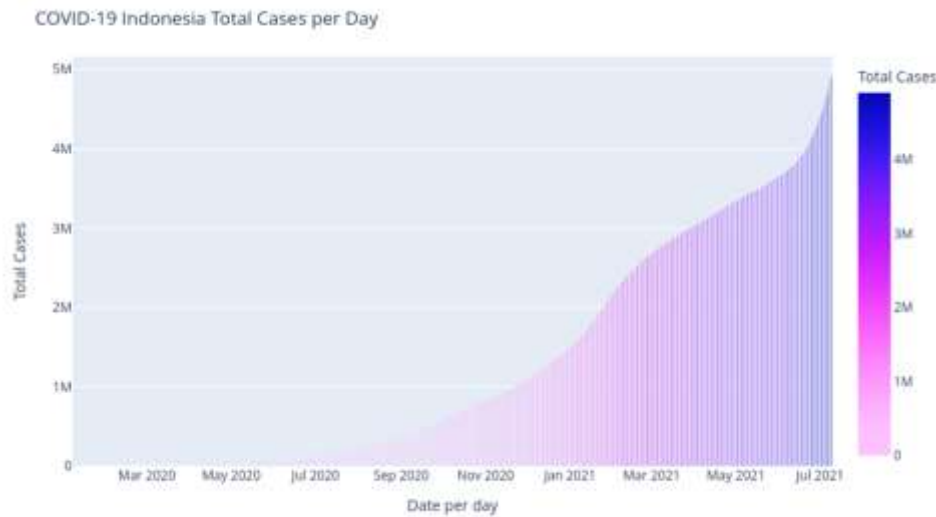


Fig. 9 COVID-19 All Cases per Day

Data Transformation

Before the data to be tested is entered into the algorithm, a data transformation is needed which will perform a normalization technique on numeric or data converting data into numbers for categorical data so that they can be processed into the Decision Tree and Linear Regression algorithm (Shalabi et al., 2006). Several predictor variables that will be used in the Decision Tree and Linear Regression algorithms include Location, New Recovered, Total Cases, Total Deaths, Total Recovered, Total Active Cases and the criteria variable is New Cases.

Data that has numerical properties will be normalized using min-max normalization into the range $0 \leq X' \leq 1$ where X' is the result of the normalization value.

$$X' = \frac{X - X_{max}}{X_{max} - X_{min}}, \quad (1)$$

where X' is normalization value, X is original value, X_{max} is maximum original value and X_{min} is minimum original value. The results of the normalization of numerical data obtained on the attributes of New Recovered, Total Cases, Total Deaths, Total Recovered and Total Active Cases can be seen below.

Table 2 Numeric Attribute Transformation

| Location | New Cases | New Recovered | Total Cases | Total Deaths | Total Recovered | Total Active Cases |
|------------|-----------|---------------|-------------|--------------|-----------------|--------------------|
| Bengkulu | 0 | 0.000000 | 0.000437 | 0.001760 | 0.000300 | 0.023848 |
| Jakarta | 552 | 0.013913 | 0.047061 | 0.073561 | 0.038119 | 0.107669 |
| Yogyakarta | 16 | 0.000466 | 0.001633 | 0.002054 | 0.001544 | 0.024814 |
| Jambi | 20 | 0.000067 | 0.000383 | 0.000367 | 0.000231 | 0.024049 |
| Jawa Barat | 44 | 0.002996 | 0.013453 | 0.017235 | 0.008955 | 0.057438 |

* Corresponding author



Data attributes that have categorical properties will be transformed into numbers so that they can be processed by the algorithm used. The results of the transformation of the categorical data attribute, such as Location, can be seen in this table below.

Table 3 Categorical Attribute Transformation

| Location | New Cases | New Recovered | Total Cases | Total Deaths | Total Recovered | Total Active Cases |
|----------|-----------|---------------|-------------|--------------|-----------------|--------------------|
| 3 | 0 | 0.000000 | 0.000437 | 0.001760 | 0.000300 | 0.023848 |
| 4 | 552 | 0.013913 | 0.047061 | 0.073561 | 0.038119 | 0.107669 |
| 5 | 16 | 0.000466 | 0.001633 | 0.002054 | 0.001544 | 0.024814 |
| 6 | 20 | 0.000067 | 0.000383 | 0.000367 | 0.000231 | 0.024049 |
| 7 | 44 | 0.002996 | 0.013453 | 0.017235 | 0.008955 | 0.057438 |

Decision Tree Experiment

A decision tree is a tree structure like a flow chart, where each *internal node* (*non-leaf node*) represents a test on a data attribute, each branch node represents the test result, and each *leaf node* (or *terminal node*) holds a class label. The top node in the tree is the *root node*. There are several models contained in the Decision Tree such as *IDS*, *C4.5*, and *CART*. *CART* stands for *classification and regression trees* (Indah Prabawati et al., 2019). In this research, the algorithm chosen is *regression tree*. The splitting formula used in the *regression tree* is the *Mean Squared Error* (*MSE*).

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (2)$$

where y is actual value, \hat{y} is predicted value and n is total data. Below are the results of training on the Decision Tree algorithm.

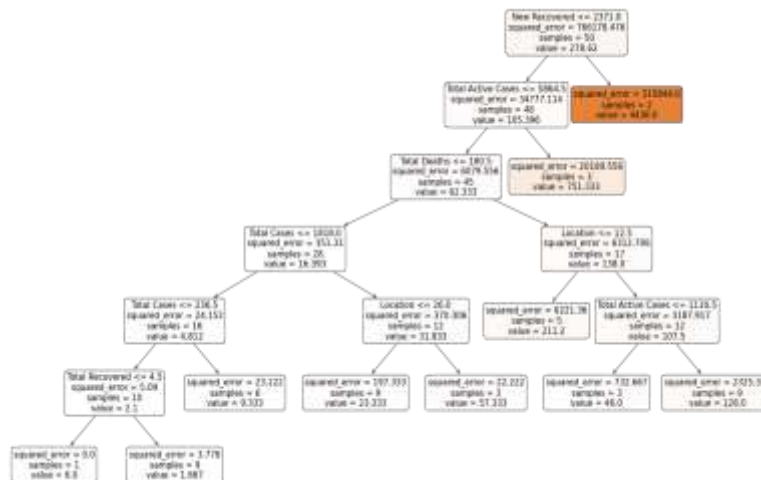


Fig. 10 Decision Tree Visualization

* Corresponding author



The results of the prediction of new cases of COVID-19 in Indonesia with the Decision Tree can then be made into a comparison between the actual new cases and the predicted positive cases as shown in the table below.

Table 4 Comparison of Actual and Predicted Positive Cases of Decision Tree

| Location | Actual New Cases | Predicted New Cases |
|---------------------|------------------|---------------------|
| Bali | 44 | 2 |
| Jambi | 34 | 35 |
| Sumatera Selatan | 53 | 52 |
| Nusa Tenggara Timur | 10 | 47 |
| Kalimantan Barat | 39 | 84 |
| Jawa Timur | 291 | 281 |
| Sulawesi Utara | 2 | 13 |
| Kalimantan Utara | 36 | 24 |

Linear Regression Experiment

Linear Regression is an algorithm used to form a model using *predictor variable* X and *criterion variable* Y . In this study, the algorithm used is Multiple Linear Regression. The way Multiple Linear Regression works is to predict the effect of two or more *predictor variables* X_1, \dots, X_n against one *criterion variable* Y to be able to prove the existence of a functional relationship between these variables (Panggabean et al., 2020).

$$y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n, \quad (3)$$

where y is *criterion variable*, w_0 is *bias*, w_i is *coefficient*, x_i is *predictor variable* and n is total data. Predicted results with Linear Regression which can then be made into a comparison between the actual new cases and predicted new cases as shown in table below.

Table 5 Comparison of Actual and Predicted Positive Cases of Linear Regression

| Location | Actual New Cases | Predicted New Cases |
|---------------------|------------------|---------------------|
| Bali | 44 | 75 |
| Jambi | 34 | 46 |
| Sumatera Selatan | 53 | 19 |
| Nusa Tenggara Timur | 10 | 66 |
| Kalimantan Barat | 39 | 74 |
| Jawa Timur | 291 | 348 |
| Sulawesi Utara | 2 | 3 |
| Kalimantan Utara | 36 | 11 |

Decision Tree and Linear Regression Comparison

Prediction results of COVID-19 cases in Indonesia can be evaluated using the R^2 Score. The value of R^2 is generally used as a quantity to estimate the percentage variance of the criterion variable in a linear relationship with the predictor variable (Renaud & Victoria-Feser, 2010).

* Corresponding author



$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}, \quad (4)$$

where RSS is *residual sum of squares*, TSS is *total sum of squares*, y is actual value, \hat{y} is predicted value and \bar{y} is mean value. The results of the comparison of the R^2 Score between the Decision Tree algorithm and Linear Regression are as shown in the table below.

Table 6 R^2 Score Comparison of *Decision Tree* and Linear Regression

| Model | Training R^2 Score (%) | Testing R^2 Score (%) |
|----------------------|--------------------------|-------------------------|
| <i>Decision Tree</i> | 95.69% | 92.15% |
| Regression Linear | 79.93% | 77.25% |

DISCUSSIONS

For further research, there are several suggestions that can be used as material for consideration in the study, such as further research on COVID-19 in Indonesia can be carried out with more recent data, the application of data mining can use *cross validation* techniques to find the best *hyperparameter* values from these algorithms and try eliminating *outlier* data before training the data into the algorithms in order to improve the performance of the algorithms.

CONCLUSION

Based on the results of research conducted, it can be concluded as follows: DKI Jakarta province has the most total cases (positives, deaths, recovered) of COVID-19 from January 2020 to July 2021, positive cases of COVID-19 in Indonesia on July 9, 2021 reached 737 thousand, while the death cases reached 129 thousand and the recovery cases reached 4 million, the R^2 Score value from the predictions of the Decision Tree algorithm for training data and testing data reached 95.69% (training) and 92.15% (testing) while the R^2 Score value from the Linear Regression algorithm reached 79.93% (training) and 77.25% (testing). From the results of this R^2 Score, it can be concluded that the Decision Tree algorithm has better performance than the Linear Regression algorithm.

REFERENCES

- Alves, M. A., Castro, G. Z., Oliveira, B. A. S., Ferreira, L. A., Ramírez, J. A., Silva, R., & Guimarães, F. G. (2021). Explaining machine learning based diagnosis of COVID-19 from routine blood tests with decision trees and criteria graphs. *Computers in Biology and Medicine*, 132. <https://doi.org/10.1016/j.compbiomed.2021.104335>
- Annis, D. (2021). Situasi Terkini Perkembangan Novel Coronavirus (COVID-19) 19 Agustus 2021. *Infeksi Emerging Kementerian Kesehatan RI*, 1–4. <https://covid19.kemkes.go.id/situasi-infeksi-emerging/situasi-terkini-perkembangan-coronavirus-disease-covid-19-20-agustus-2021>
- Baharuddin, M. M., Azis, H., & Hasanuddin, T. (2019). Analisis Performa Metode K-Nearest Neighbor Untuk Identifikasi Jenis Kaca. *ILKOM Jurnal Ilmiah*, 11(3), 269–274. <https://doi.org/10.33096/ilkom.v11i3.489.269-274>
- Bender, L. (2020). Pesan dan Kegiatan Utama Pencegahan dan Pengendalian COVID-19 di Sekolah. *Unicef*, 1–14.
- Dwitri, N., Tampubolon, J. A., Prayoga, S., Ilmi Zer, F., & Hartama, D. (2020). Penerapan Algoritma K-Means Dalam Menentukan Tingkat Penyebaran Pandemi Covid-19 Di Indonesia. *Jti (Jurnal Teknologi Informasi)*, 4(1), 101–105.
- Indah Prabawati, N., Widodo, & Ajie, H. (2019). Kinerja Algoritma Classification And Regression Tree (Cart) dalam

* Corresponding author



This is an Creative Commons License This work is licensed under a
Creative Commons Attribution-NoDerivatives 4.0 International License.

Mengklasifikasikan Lama Masa Studi Mahasiswa yang Mengikuti Organisasi di Universitas Negeri Jakarta. *PINTER: Jurnal Pendidikan Teknik Informatika Dan Komputer*, 3(2), 139–145. <https://doi.org/10.21009/pinter.3.2.9>

- Nabila, Z., Rahman Isnain, A., & Abidin, Z. (2021). Analisis Data Mining Untuk Clustering Kasus Covid-19 Di Provinsi Lampung Dengan Algoritma K-Means. *Jurnal Teknologi Dan Sistem Informasi (JTSI)*, 2(2), 100. <http://jim.teknokrat.ac.id/index.php/JTSI>
- Panggabean, D. S. O., Buulolo, E., & Silalahi, N. (2020). Penerapan Data Mining Untuk Memprediksi Pemesanan Bibit Pohon Dengan Regresi Linear Berganda. *JURIKOM (Jurnal Riset Komputer)*, 7(1), 56. <https://doi.org/10.30865/jurikom.v7i1.1947>
- Pramadhani, E. E., & Tedy, S. (2014). Penerapan Data Mining untuk Klasifikasi Prediksi Penyakit ISPA (Infeksi Saluran Pernapasan Akut) dengan Algoritma Decision Tree (ID3). *Jurnal Sarjana Teknik Informatika*, 2(1), 831–839.
- Renaud, O., & Victoria-Feser, M. P. (2010). A robust coefficient of determination for regression. *Journal of Statistical Planning and Inference*, 140(7), 1852–1862. <https://doi.org/10.1016/j.jspi.2010.01.008>
- Shalabi, L. Al, Shaaban, Z., & Kasasbeh, B. (2006). Data Mining: A Preprocessing Engine. *Journal of Computer Science*, 2(9), 735–739. <https://doi.org/10.3844/jcssp.2006.735.739>
- Sindi, S., Ningse, W. R. O., Sihombing, I. A., Ilmi R.H.Zer, F., & Hartama, D. (2020). Analisis algoritma K-Medoids clustering dalam pengelompokan penyebaran Covid-19 di Indonesia. *Jti (Jurnal Teknologi Informasi)*, 4(1), 166–173.
- Sodik, F., Dwi, B., & Kharisudin, I. (2020). Perbandingan Metode Klasifikasi Supervised Learning pada Data Bank Customers Menggunakan Python. *Jurnal Matematika*, 3, 689–694.
- Supriatna, E. (2020). Wabah Corona Virus Disease (Covid 19) Dalam Pandangan Islam. *SALAM: Jurnal Sosial Dan Budaya Syar-I*, 7(6), 555–564. <https://doi.org/10.15408/sjsbs.v7i6.15247>
- Sutoyo, I. (2018). Implementasi Algoritma Decision Tree Untuk Klasifikasi Data Peserta Didik. *Jurnal Pilar Nusa Mandiri*, 14(2), 217. <https://doi.org/10.33480/pilar.v14i2.926>
- Vadyala, S. R., Betgeri, S. N., Sherer, E. A., & Amritphale, A. (2021). Prediction of the number of COVID-19 confirmed cases based on K-means-LSTM. *Array*, 11, 100085. <https://doi.org/10.1016/j.array.2021.100085>
- Wulandari, R. T. (2016). Model Data Mining sebagai Prediksi Penyakit Hipertensi Kehamilan dengan Teknik Decision Tree. *Scientific Journal of Informatics*, 3(1), 19–26.
- Xiao, S., Cheng, G., Yang, R., Zhang, Y., Lin, Y., & Ding, Y. (2021). Prediction on the number of confirmed Covid-19 with the FUDAN-CCDC mathematical model and its epidemiology, clinical manifestations, and prevention and treatment effects. *Results in Physics*, 20, 103618. <https://doi.org/10.1016/j.rinp.2020.103618>
- Yoo, S. H., Geng, H., Chiu, T. L., Yu, S. K., Cho, D. C., Heo, J., Choi, M. S., Choi, I. H., Cung Van, C., Nhung, N. V., Min, B. J., & Lee, H. (2020). Deep Learning-Based Decision-Tree Classifier for COVID-19 Diagnosis From Chest X-ray Imaging. *Frontiers in Medicine*, 7, 1–8. <https://doi.org/10.3389/fmed.2020.00427>

* Corresponding author



This is an Creative Commons License This work is licensed under a
Creative Commons Attribution-NoDerivatives 4.0 International License.