

Analysis of COVID-19 Virus Spread in Jakarta Using Multiple Linear Regression

Na'il Muta'aly Muhtar¹⁾, Putu Harry Gunawan^{2)*}

^{1,2)}CoE Hemic, School of Computing, Telkom University, Bandung, Indonesia

¹⁾nailmutaaly@student.telkomuniversity.ac.id, ²⁾phgunawan@telkomuniversity.ac.id

Submitted : Jun 26, 2024 | **Accepted** : Jul 1, 2024 | **Published** : Jul 5, 2024

Abstract: COVID-19, first identified in Wuhan, China in December 2019, quickly spread worldwide and was declared a pandemic by WHO in March 2020. Indonesia reported its first case on March 2, 2020, and the pandemic has had a significant impact on the country's economic, social, and health sectors. This study aims to predict the death rate due to COVID-19 in Jakarta using multiple linear regression method. The dataset collected from Andra Farm - Go Green website includes COVID-19 cases recorded in all sub-districts in Jakarta on November 1, 2023. Pre-processing was performed to improve the quality and accuracy of the model. The method used was multiple linear regression. The analysis results show that variables such as total travel and discarded trip have a significant influence in predicting the number of positive cases. The study found that lowering the correlation threshold for selecting independent variables reduced the mean squared error (MSE) and improved model performance, highlighting the importance of variable selection in developing accurate predictive models. These findings provide important insights for the government in making informed decisions regarding post-pandemic healthcare. This research underscores the value of robust data processing and variable selection techniques in enhancing predictive accuracy for public health planning.

Keywords: Covid-19; Correlation matrix; Multiple linear regression; Jakarta; Positive;

INTRODUCTION

The COVID-19 (Coronavirus Disease 2019) was initially discovered in Wuhan, China, and between December 2019 and the beginning of 2020, it expanded to other nations (Qiu et al., 2020). In March 2020, WHO (World Health Organization) designated this outbreak as a pandemic that has an impact on the economy, loss of jobs and income due to quarantine (Padhan et al., 2021). Indonesia is one of the countries affected by the spread of COVID-19 with the first case announced on March 2, 2020 (Sholochah et al., 2024).

COVID-19 has had a significant impact in Indonesia. Across the country, economic, social and health sectors have been impacted by the pandemic. The increase in COVID-19 cases and deaths in the healthcare sector is putting the healthcare system under pressure. Amid rising cases in hospitals across Indonesia, there is a shortage of medical equipment and medical personnel. The Indonesian economy has also been affected by the pandemic, as many businesses have been forced to reduce production or even close, causing job losses and economic uncertainty. As of July 6, 2020, the government of the Republic of Indonesia has reported 64,958 Covid-19 cases, including 3,141 deaths, and 29,919 patients have recovered (Kahar et al., 2020). In the healthcare sector, the number of cases and deaths due to COVID-19 is increasing, putting a huge strain on the healthcare system. One of the regions in Indonesia has confirmed two positive cases of COVID-19 on March 19, 2020 which has been conveyed by the

*Putu Harry Gunawan



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

Gubernur of South Sulawesi that one of them was positively affected from performing Umrah worship and one of them was a student who came from Jakarta (Sumandiyar et al., 2020).

Analyzing COVID-19 is critical as it allows us to understand the widespread impact of the pandemic, even after it has ended. This analysis provides valuable insights into epidemiological trends, the effectiveness of public health interventions, and socio-economic consequences. This understanding helps in preparing for future pandemics and strengthening health systems. The method of using data from Andra Farm - Go Green involves a multiple linear regression algorithm to analyze the dataset obtained. The main objective of this research is to predict the COVID-19 mortality rate in Jakarta using the multivariable regression machine learning method. By using multivariable regression, this research seeks to model and understand the factors that affect the mortality rate due to COVID-19 in Jakarta. This can help the government make more effective decisions regarding health services after the pandemic.

LITERATURE REVIEW

Danielle Klinger et al. investigated the association between Bacille Calmette-Guérin (BCG) tuberculosis vaccination and better outcomes for COVID-19 patients. The data they studied came from fifty-five countries that met certain criteria for population size and death rate per million (DPM). Multivariate regression tests involving 23 variables related to economic, demographic, health and pandemic restrictions showed that the year of BCG vaccine administration had a significant negative effect on COVID-19 outcomes. The study found that the young age group (0-24 years) had the strongest correlation with BCG vaccination coverage. In addition, a strong correlation and statistical significance was found with BCG coverage over the past 15 years, while no similar correlation was found for measles and rubella vaccination protocols. This study discovered that BCG vaccination helps reduce the spread and intensity of the COVID-19 pandemic, especially among recently vaccinated populations (Klinger et al., 2020).

This study by Milena Gianfrancesco et al. aims to study the spatial relationship between daily particulate matter (PM) concentrations and COVID-19 mortality rates in China, which was first reported in December 2019 in Wuhan and reached 3,314 deaths on March 31, 2020. Using multiple linear regression methods, a cross-sectional analysis was conducted to evaluate the spatial relationship between COVID-19 mortality rates and daily PM2.5 and PM10 concentrations. The results showed that an increase in PM2.5 and PM10 concentrations led to an increase in COVID-19 mortality rates. This may affect how patients progress from mild symptoms to severe symptoms and may ultimately affect the prediction of COVID-19 patients. (Gianfrancesco et al., 2020).

Research conducted by Mohammed Khaled Al-Hanawi et al public perceptions and awareness have affected Saudi Arabia's adherence to the stringent COVID-19 control procedures. Through the use of an online questionnaire, 3,388 participants provided data for this study that assessed their knowledge, attitudes, and practices on COVID-19. In order to evaluate variations in mean scores and pinpoint the variables influencing their knowledge, attitudes, and practices, univariate and multivariate regression analyses were employed. The most prevalent were favorable attitudes (17.96 average score) and good knowledge (17.96 average score). But older adults outperformed younger ones, and men's knowledge, attitudes, and actions were inferior to women's. These results suggest that certain groups require focused health education programs. (Al-Hanawi et al., 2020).

Research has been conducted by Asmaa S. Qaddoori used multiple linear regression analysis to determine the components contributing to the increase in COVID-19 mortality, with data from hospitals in Salah al-Din Governorate from April 1, 2020 to December 30, 2020. Data were collected from Tikrit Hospital, Baiji Hospital, Shirqat General Hospital, and Samarra Hospital. The results showed that the three main factors leading to an increase in COVID-19-related deaths were age, place of death, and chronic diseases. Chronic diseases worsen the patient's condition, and the older the patient, the higher the risk of death. Health outcomes are also affected by the location of care; care in hospitals with

adequate facilities is more effective than care at home. This suggests that medical care at home cannot match the facilities, medical equipment and expertise of hospital employees. Lastly, about 38% of other factors not examined in this study also influenced mortality, showing how complex the factors contributing to death from COVID-19 are. This study provides important insights into the need for better healthcare facilities and greater attention to chronic diseases to reduce the impact of the COVID-19 pandemic (S. Qaddoori, 2023).

Taylor Jansen et al. in the United States discovered that those over the age of 65 accounted for 81% of all COVID-19 deaths. The purpose of the study was to investigate the relationship between COVID-19 death rates in 208 cities and towns in Connecticut and Rhode Island and asthma rates at the city level. A stepwise analysis was used to look at how these factors interacted. In order to evaluate the relationships between COVID-19 deaths per 100,000 persons and the prevalence of asthma and COPD in individuals 65 years of age and older, bivariate maps were developed. The bivariate analysis's preliminary results indicated a positive relationship between these regions' COVID-19 mortality toll and the prevalence of asthma and COPD. To determine whether COVID-19 deaths are associated with chronic lung disease and other city-level characteristics, multiple linear regression models were used. The multiple linear regression model did not show a significant association between COVID-19 mortality and chronic lung disease at the city level, although the baseline data suggested such an association. However, it was observed that some city-level characteristics, such as age 65 years and older, lower education levels, and African-American and Hispanic ethnicity, were significant predictors of COVID-19 mortality.(Jansen et al., 2022).

Rajani Kumari et al examine the coronavirus outbreak first reported in late December 2019, which has infected more than 7 million people and resulted in more than 0.40 million deaths worldwide. In India, the first case was found on January 30, 2020. By June 6, 2020, the number of cases had exceeded 0.24 million. This study provides a new in-depth forecasting model and predicts the number of confirmed COVID-19 cases, recoveries, and deaths in India. For prediction, the model uses correlation coefficient and multiple linear regression, and to improve accuracy, it also uses autoregression and autocorrelation methods. With an R-squared value of 0.9992, the predicted results show excellent agreement with the actual values. The results suggest that social distancing and lockdown are two important components in reducing the spread of COVID-19 (Kumari et al., 2021).

Sushma Dahal et al. investigates how geographic differences in excess mortality rates in Mexico during the COVID-19 pandemic are affected by demographics, meteorology and health elements. Using the Serfling regression model, excess mortality from all causes in 32 Mexican states was estimated. Using multiple linear regression, the relationship between excess mortality and various sociodemographic, climatic, and health factors was studied. The study identified groups of states with different excess mortality growth trends using functional data analysis. The results project that the overall excess mortality rate in Mexico will reach 39.66 per 100,000 people by April 10, 2021. The states of Chiapas (12.72) and Oaxaca (13.42) in the southeastern part of the country showed the lowest excess mortality rates, while Mexico City (106.17) and Tlaxcala (51.9%) showed the highest excess mortality rates. A positive association ($P < 0.001$) between excess mortality and factors such as aging index, marginalization index, and average household size was found in the final adjusted model ($R^2 = 77\%$). In addition, four separate groups were found to have similar trends in excess mortality. In conclusion, the distribution of aging index, marginalization index, and average household size explained the differences in excess mortality rates in Mexico, with the central states experiencing the highest excess mortality rates (Dahal et al., 2021).

METHOD

In this section, each researcher is expected to provide the latest contribution to solving the existing problems. Researchers can also use flowcharts, pictures, and diagrams to show how to solve the problem.

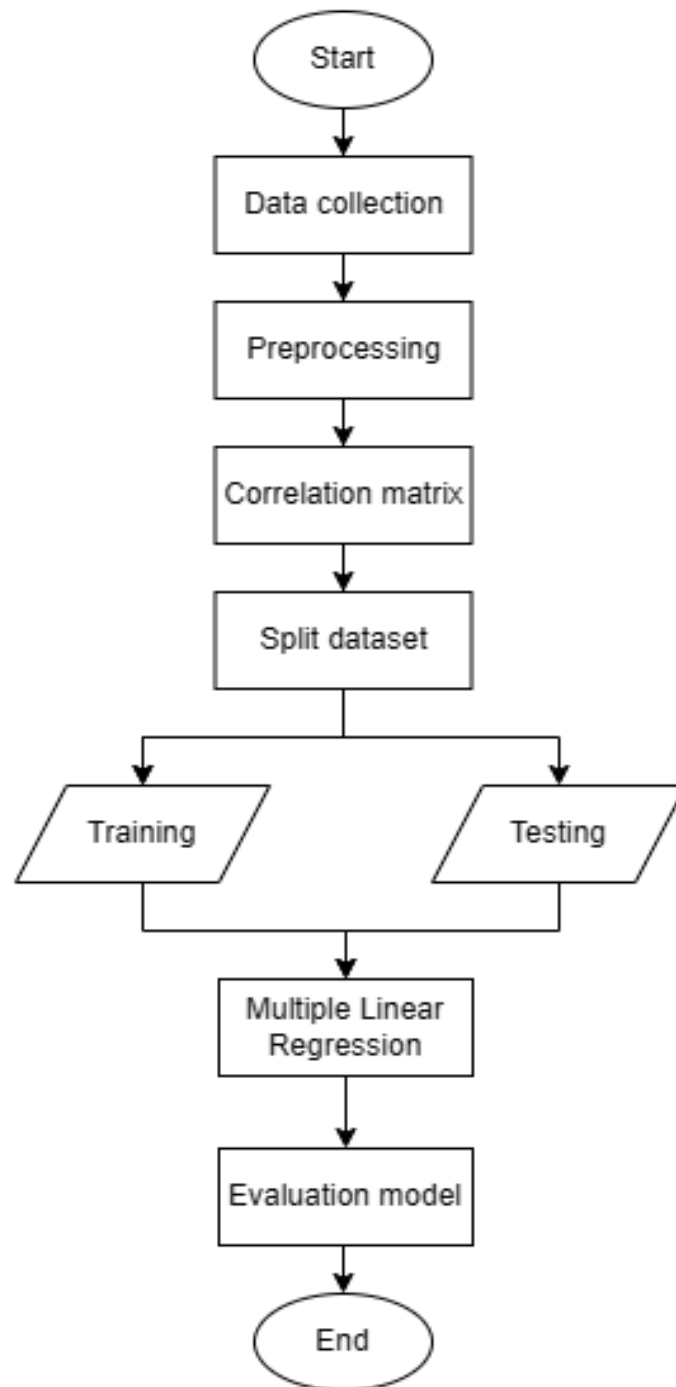


Fig. 1 Research Desain

Dataset Collection

The dataset for this research was obtained through the Andra Farm - Go Green website, which provides data for each sub-district in DKI Jakarta. This dataset focuses on Covid-19 cases recorded on November 1, 2023. The use of this dataset is very important because it allows a more in-depth, effective, and representative analysis related to the level of Covid-19 spread throughout Jakarta.

Preprocessing

To improve accuracy in the analysis, a series of important preprocessing steps were performed. One of the main steps is to remove some columns that are not used in further analysis. These columns may

not be relevant to the purpose of the analysis or may not contribute significantly to the final model results. In addition, columns containing a value of 0 were also removed. These columns with 0 values may indicate missing or irrelevant data, which could compromise the quality of the analysis and model predictions. By removing these columns, the dataset becomes cleaner and more focused, allowing the model to be trained with more representative and high-quality data. These preprocessing steps are essential to ensure that the data used in the model is relevant, free from distractions, and able to provide more accurate and reliable results.

Correlation matrix

Correlation matrices are a very useful tool in describing the correlation properties of Boolean mappings directly. This matrix presents the relationship between variables in a form that is easy to analyze and understand, especially in the context of linear cryptanalysis (Daemen et al., 2020). Linear cryptanalysis is a technique used to solve cryptographic systems by analyzing the linear relationship between plaintext, ciphertext, and key.

In this context, the correlation matrix R can be used to identify and exploit weaknesses in encryption algorithms. Each element in the R correlation matrix represents the correlation between two Boolean variables, which can be calculated using the Pearson correlation coefficient. Equation (1) for calculating the Pearson correlation between two Boolean variables X and Y is (Zhang et al., 2023):

$$\rho_{X,Y} = \frac{\text{Cov}(X,Y)}{\sigma_X\sigma_Y} = \frac{E[(X-\mu_x)(Y-\mu_y)]}{\sigma_X\sigma_Y} \quad (1)$$

Where $\rho_{X,Y}$ is the correlation coefficient between X and Y , $\text{Cov}(X,Y)$ is the covariance between X and Y , and $\sigma_X\sigma_Y$ is the standard deviation of X and Y .

Multiple linear regression

Multiple regression analysis is a linear regression model that includes one continuous variable and k (two or more) independent variables (Muthahharah et al., 2021). This method predicts how the value of certain variables will change in response to changes in other variables (Alita et al., 2021). It makes it possible to model and measure how different independent variables affect the dependent variable. Equation (2) for multiple linear regression is as follows (Alita et al., 2021):

$$Y = \beta_0 + \beta_1X_1 + \beta_2X_2 + \dots + \beta_nX_n + \epsilon \quad (2)$$

Where Y is the predicted dependent variable, β_0, \dots, β_n are regression coefficients, dan X_0, \dots, X_n are independent variables

RESULT

In Figure 2, the correlation heatmap shows the relationship between various variables associated with the number of positive COVID-19 cases. From this analysis, it can be seen that some variables have a high correlation with the number of positive cases, such as treated, recovered, died, suspected, suspected isolated, probable died, close contact, close contact isolated, and close contact discarded. This suggests that an increase in the number of patients treated, recovered, died, or close contacts isolated correlates with an increase in the number of positive cases. In contrast, some variables such as total trips, quarantined trips, discarded trips, suspected died, suspected discarded, and probable total showed low or negative correlations with the number of positive cases, suggesting that they have less influence in determining the number of positive cases and may not be relevant to include in the prediction model.

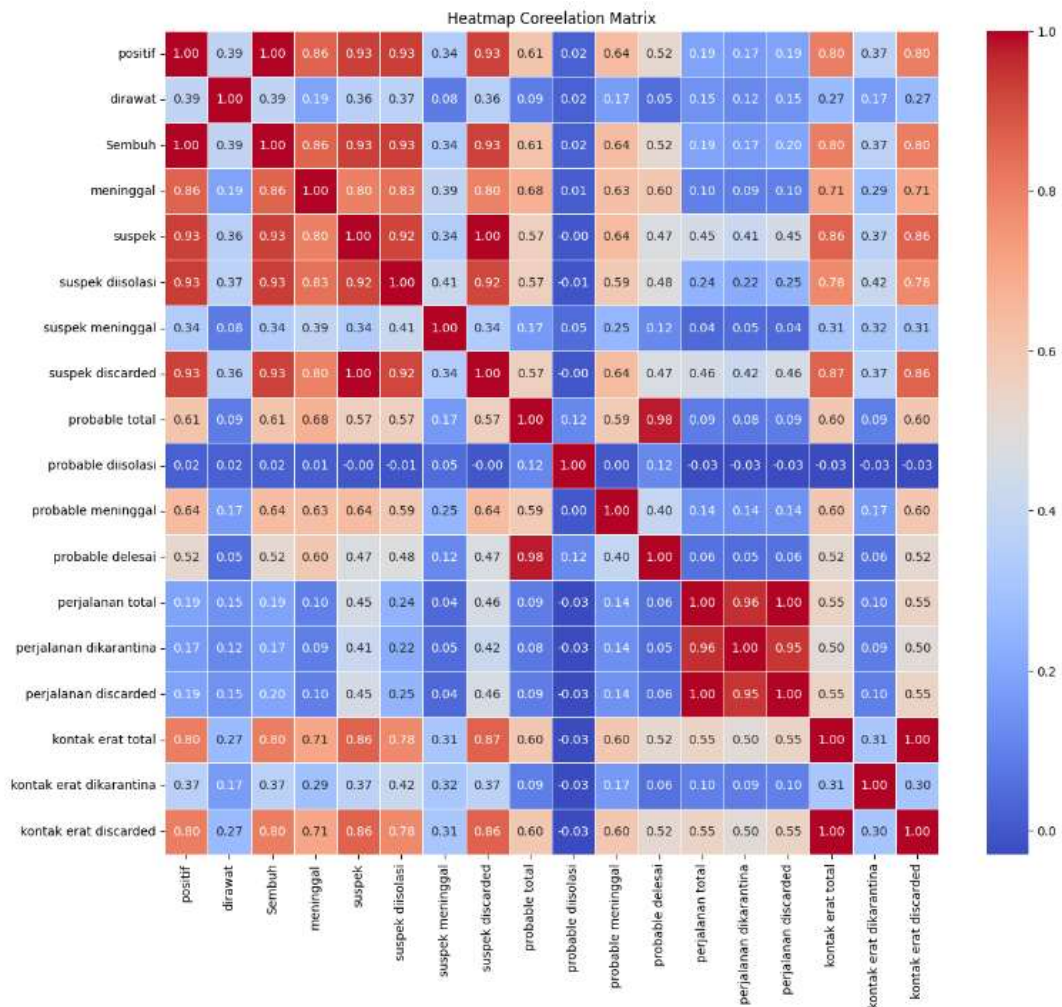


Fig. 2 Result Correlation Matrix

Overall, this correlation analysis provides valuable insights into the relationship between various variables and the number of positive cases. Some variables such as recovered, died, suspect, isolated suspect, and close contact had high correlations with the number of positive cases, indicating the importance of these variables in predicting the number of cases. In contrast, variables such as total trips, quarantined trips, and isolated probables had low correlations, suggesting that they have less influence in the predictive model. These results can be used to develop a more accurate predictive model by including variables that have a high correlation and excluding variables that are less relevant.

This study evaluated the effect of the number of independent variables in a linear regression model on the mean squared error (MSE) and variable coefficients. Four linear regression models were built based on "positive" correlation values between the independent variables and the target variable as the selection criteria. Each model had a different correlation threshold ranging from 0.40 to 0.33 to determine which variables to include in the model.

The correlation threshold in the first test was below 0.40. The regression model produced the highest value of all tests, with a Mean Squared Error (MSE) of 1.41. The variables suspected dead, probable isolated, quarantined travel and quarantined close contact had negative coefficients, indicating that an increase in the number of treated patients significantly increased the number of positive cases. The variables of total trips and discarded trips also have large coefficients, indicating that an increase in the number of treated patients significantly increases the number of positive cases.

Table 1
 Testing using 7 variable

| Variabel | Coefficient |
|---------------------------|---------------|
| Treated | 2.421451e+00 |
| Suspected dead | -1.334418e-01 |
| Probable isolated | -5.711537e-01 |
| Total travel | 1.383382e+06 |
| Quarantined travel | -1.064489e+05 |
| Discarded trip | -1.282096e+06 |
| Quarantined close contact | -4.813104e-01 |

In the Table 2, the correlation threshold dropped to less than 0.38, which caused the MSE to drop significantly to 0.08. The variables total trips and discarded trips still have large coefficients, indicating a significant influence on the number of positive cases. The variables probable isolated had larger negative coefficients, suggesting a more substantial reduction in the number of positive cases as the number of probable isolated increased; and the variables suspected death, quarantined trips, and close contact during quarantine still had negative coefficients indicating that the reduction in the number of positive cases had an impact.

Table 2. Testing using 6 variable

| Variabel | Coefficient |
|---------------------------|----------------|
| Treated | 0 |
| Suspected dead | -0.096899 |
| Probable isolated | -0.891669 |
| Total travel | 978438.097068 |
| Quarantined travel | -75293.733126 |
| Discarded trip | -906796.132003 |
| Quarantined close contact | -0.251066 |

In the Table 3, the correlation threshold dropped further to below 0.36, which caused the MSE to drop to 0.07. The variables total trips and discarded trips still had large coefficients, indicating a significant influence on the number of positive cases. The variable's likelihood of isolation suspected death, and quarantined trips still had negative coefficients, indicating a smaller influence on the number of positive cases.

Table 3. Testing using 5 variable

| Variables | Coefficient MLR |
|---------------------------|-----------------|
| Treated | 0 |
| Suspected dead | -0.105254 |
| Probable isolated | -0.765610 |
| Total travel | 630061.240298 |
| Quarantined travel | -48486.652284 |
| Discarded trip | -583926.277004 |
| Quarantined close contact | 0 |

In the Table 4, the correlation threshold was lowered to less than 0.33, resulting in the lowest MSE of 0.06. It can be seen that the variables of total trips and discarded trips still have large coefficients and have a strong influence on the number of positive cases. The variables of trips likely to be quarantined and probable isolated still have negative coefficients, indicating a decreasing influence on the number of positive cases.

*Putu Harry Gunawan



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

Table 4. Testing using 4 variable

| Variabel | Coefficient |
|---------------------------|----------------|
| Treated | 0 |
| Suspected dead | 0 |
| Probable isolated | -0.716765 |
| Total travel | 585876.717041 |
| Quarantined travel | -45086.779245 |
| Discarded trip | -542976.702308 |
| Quarantined close contact | 0 |

The MSE analysis of Table 5 shows that lowering the correlation threshold for independent variables can significantly lower the MSE, indicating an improvement in model performance. For example, the MSE dropped from 1.41 to 0.06 when the correlation threshold dropped from <0.40 to <0.33. This shows that the improvement in prediction accuracy and the risk of overfitting are reduced by reducing the number of variables present in the model.

Table 5. MSE analysis

| Testing | Threshold Correlations | Mean Squared Error |
|---------|------------------------|--------------------|
| 1 | corr < 0.40 & 0 | 1.41 |
| 2 | corr < 0.38 & 0 | 0.08 |
| 3 | corr < 0.36 & 0 | 0.07 |
| 4 | corr < 0.33 & 0 | 0.06 |

DISCUSSIONS

The results of this test show that lowering the correlation threshold for selecting independent variables significantly reduces MSE and improves model performance. For example, if the correlation threshold decreases from less than 0.40 to less than 0.33, the MSE decreases from 1.41 to 0.06. Reducing the number of variables in the model also helps reduce the risk of overfitting, but we can see that some variables still have large coefficients and have a strong influence on the prediction. The variable's number of trips and abandoned trips consistently have large coefficients across tests, suggesting that they play an important role in influencing the number of positive cases.

Decreasing the correlation threshold reduces the number of independent variables included in the model, which helps improve model efficiency without losing important information. The impact analysis shows that fewer variables with lower correlation can lead to more accurate predictions, indicated by a significant decrease in MSE. In addition, the reduction in variables included helps avoid multicollinearity, which can improve the stability and interpretability of the model.

CONCLUSION

This study shows that selecting independent variables based on correlation values can affect the performance of linear regression models in predicting the number of positive cases. Lowering the correlation threshold from 0.40 to 0.33 resulted in improved model accuracy, as indicated by a decrease in MSE. Some variables such as total trips and discarded trips consistently showed a significant effect on predicting the number of positive cases, indicating the importance of these variables in the model. These results provide valuable insights for the development of more accurate and stable predictive models in the future.

REFERENCES

- Al-Hanawi, M. K., Angawi, K., Alshareef, N., Qattan, A. M. N., Helmy, H. Z., Abudawood, Y., Alqurashi, M., Kattan, W. M., Kadasah, N. A., Chirwa, G. C., & Alsharqi, O. (2020). Knowledge, Attitude and Practice Toward COVID-19 Among the Public in the Kingdom of Saudi Arabia: A Cross-Sectional Study. *Frontiers in Public Health*, 8. doi: 10.3389/fpubh.2020.00217
- Alita, D., Putra, A. D., & Darwis, D. (2021). Analysis of classic assumption test and multiple linear regression coefficient test for employee structural office recommendation. *IJCCS (Indonesian Journal of Computing and Cybernetics Systems)*, 15(3), 295. doi: 10.22146/ijccs.65586
- Daemen, J., Govaerts, R., & Vandewalle, J. (2020). Correlation Matrices. *Springer*. Retrieved from <https://lirias.kuleuven.be/retrieve/333387>
- Dahal, S., Luo, R., Swahn, M. H., & Chowell, G. (2021). Geospatial Variability in Excess Death Rates during the COVID-19 Pandemic in Mexico: Examining Socio Demographic, Climate and Population Health Characteristics. *International Journal of Infectious Diseases*, 113, 347–354. doi: 10.1016/j.ijid.2021.10.024
- Gianfrancesco, M., Hyrich, K. L., Al-Adely, S., Al-Adely, S., Carmona, L., Danila, M. I., Gossec, L., Gossec, L., Izadi, Z., Jacobsohn, L., Katz, P., Lawson-Tovey, S., Lawson-Tovey, S., Mateus, E. F., Rush, S., Schmajuk, G., Simard, J., Strangfeld, A., Trupin, L., ... Robinson, P. C. (2020). Characteristics associated with hospitalisation for COVID-19 in people with rheumatic disease: Data from the COVID-19 Global Rheumatology Alliance physician-reported registry. *Annals of the Rheumatic Diseases*, 79(7), 859–866. doi: 10.1136/annrheumdis-2020-217871
- Jansen, T., Lee, C. M., Xu, S., Silverstein, N. M., & Dugan, E. (2022). The Town-Level Prevalence of Chronic Lung Conditions and Death From COVID-19 Among Older Adults in Connecticut and Rhode Island. *Preventing Chronic Disease*, 19. doi: 10.5888/pcd19.210421
- Kahar, F., Dirawan, G. D., Samad, S., Qomariyah, N., & Purlinda, D. E. (2020). The Epidemiology of COVID-19, Attitudes and Behaviors of the Community During the Covid Pandemic in Indonesia. *International Journal of Innovative Science and Research Technology*, 5(8), 1681–1687. doi: 10.38124/ijisrt20aug670
- Klinger, D., Blass, I., Rappoport, N., & Linial, M. (2020). Significantly improved COVID-19 outcomes in countries with higher bcg vaccination coverage: A multivariable analysis. *Vaccines*, 8(3), 1–14. doi: 10.3390/vaccines8030378
- Kumari, R., Kumar, S., Poonia, R. C., Singh, V., Raja, L., Bhatnagar, V., & Agarwal, P. (2021). Analysis and predictions of spread, recovery, and death caused by COVID-19 in India. *Big Data Mining and Analytics*, 4(2), 65–75. doi: 10.26599/BDMA.2020.9020013
- Muthahharah, I., & Fatwa, I. (2021). Modeling The Types of Online Learning Media Using Multiple Linear Regression Analysis. *Jurnal Varian*, 5(1), 39–46. doi: 10.30812/varian.v5i1.1459
- Padhan, R., & Prabheesh, K. P. (2021). The economics of COVID-19 pandemic: A survey. *Economic Analysis and Policy*, 70, 220–237. doi: 10.1016/j.eap.2021.02.012
- Qiu, J., Shen, B., Zhao, M., Wang, Z., Xie, B., & Xu, Y. (2020). A nationwide survey of psychological distress among Chinese people in the COVID-19 epidemic: Implications and policy recommendations. In *General Psychiatry* (Vol. 33, Issue 2). BMJ Publishing Group. doi: 10.1136/gpsych-2020-100213

- S. Qaddoori, A. (2023). Detection of the most important factors affecting the increase in the number of deaths as a result of infection with Covid-19 using the multiple linear regression equation. *Science Archives*, 04(02), 147–153. doi: 10.47587/sa.2023.4212
- Sholochah, H., & Johan, S. (2024). The Effect First Case Covid-19 Announcement on Average Trading Volume Activity of Pharmaceutical Sector Companies. *Jurnal Manajemen Bisnis Dan Kewirausahaan*, 6, 218–224.
- Sumandiyar, A., & Nur, H. (2020). Membangun Hubungan Sosial Masyarakat di Tengah Pandemi Covid-19 di Kota Makassar. *PROSIDING NASIONAL COVID-19*. Retrieved from <https://ojs.literacyinstitute.org/index.php/prosiding-covid19>
- Zhang, H., Liu, H., Ma, G., Zhang, Y., Yao, J., & Gu, C. (2023). A wildfire occurrence risk model based on a back-propagation neural network-optimized genetic algorithm. *Frontiers in Energy Research*, 10. doi: 10.3389/fenrg.2022.1031762