

The Analysis Of Bayesian Bootstrap Binary Logistic Regression In Modeling The Recovery Rate Of Covid-19 Patients

Lili Rahmita Sari*, Ferra Yanuar, Dodi Devianto

Mathematics Department,

Faculty of Mathematics and Natural Sciences, Andalas University,

Campus of UNAND Limau Manis Padang-25163, Indonesia



Abstract— This study applied the Bayesian binary logistic regression method and the Bootstrap method to model the recovery rate of COVID-19 patients. The aim is to model the recovery rate of COVID-19 patients in order to identify the symptoms and factors affecting the recovery rate of COVID-19 patients. Data obtained from the M. Djamil Hospital, Padang City and Andalas University Hospital on COVID-19 patients in West Sumatra in 2020 were used as the case data. The case data were randomly divided with proportion of 80% training data and 20% testing data. The training data with Bayesian binary logistic regression and perform parameter estimation were later analyzed for testing data using the Bootstrap method. The parameter significance results show that there is one predictor variable that significantly affects the recovery rate of COVID-19 patients, namely patients aged 0 – 59 years. The Bayesian binary logistic regression method used in the modeling has been accurate based on the performance test of the algorithm that has been used with the Bootstrap method. This study proves that the estimated value with Bayesian binary logistic regression is at the 95% Bootstrap confidence interval. The results of the classification model for the recovery rate of COVID-19 patients show good performance by producing high accuracy, sensitivity, and precision values in identifying patients. Therefore, it can be concluded that Bayesian binary logistic regression and the Bootstrap method can be used to model the recovery rate of COVID-19 patients as they produce high classification accuracy.

Keywords— COVID-19; Bayesian binary logistic regression; Bootstrap method

I. INTRODUCTION

The *Coronavirus Disease* (COVID-19) pandemic is an outbreak of a new type of disease originating in Wuhan, China. It is highly Infections with symptoms range from mild to severe symptoms that cause death to many. The World Health Organization (WHO) [1] states that COVID-19, infections in the human respiratory tract, has become a global health problem. According to Zhai [2] COVID-19 spreads to humans through transmission from wild animals sold illegally at a wholesale seafood specialty market in Wuhan, China. Therefore, this virus spreads quickly through human-animal contact, and spreads widely between humans and humans.

Patients infected with COVID-19 show different recovery ranges for each individual. The recovery is influenced by age, comorbidities (comorbidities) and symptoms experienced by the patient. Those with old age, comorbidities, and severe symptoms experience a longer recovery rate. Das et al [3] found that the older age group with comorbidities tended to have a longer recovery time compared to the younger age group. Tolossa et al [4] found that, on average, elderly patients experience comorbidities, a fever, and a longer recovery time than those who were younger and/or had no comorbidities. Thus, it is important to model the recovery rate of COVID-19 patients in order to identify the symptoms and factors affecting the recovery rate of patients infected with COVID-19.

Patients infected with COVID-19 have two possibilities, namely recovery and death. This indicates that there are two categories (dichotomous) that will be examined in this case, namely recovery or death after being infected with COVID-19. In this case, the relationship between the two variables was not examined so that ordinary linear regression analysis could not be used but the binary logistic regression method was used instead. The method used examined the relationship between two variables, namely predictor variables and dichotomous response variables [5].

Maximum likelihood estimation (MLE) is a parameter estimation method based only on information from sample data. In addition to MLE, there is a parameter estimation method not only based on information from sample data, but also based on initial information about the distribution of the parameters to be estimated, namely the Bayesian method [6]. Furthermore, in order to test the performance of the algorithm constructed in the application of the Bayesian binary logistic regression method, a parameter estimation method by performing a random resampling technique from the original sample, namely the Bootstrap method was also employed [7].

II. CASE DESCRIPTION

2.1 Data

The case data used is data on COVID-19 patients in West Sumatra in 2020 obtained from the M. Djamil General Hospital, in Padang City and The Hospital of Andalas University from March to December 2020. The Y response variable in this study is the recovery rate of COVID-19 patients. Patients infected with COVID-19 that are declared recovery is stated with the terms “Y= 1” (81%), and “Y=0” (19%) if the patient infected with COVID-19 dies “Y=0” (19%). There are three predictor variables used, including gender (female and male categories), age level (age categories 0 - 59 years, and > 59 years, and also based on the number of comorbidities (categories not having comorbidities, having 1-3 comorbid, and >3 comorbid).

2.2 Method

The stages in this research can be seen in Figure 1 below:

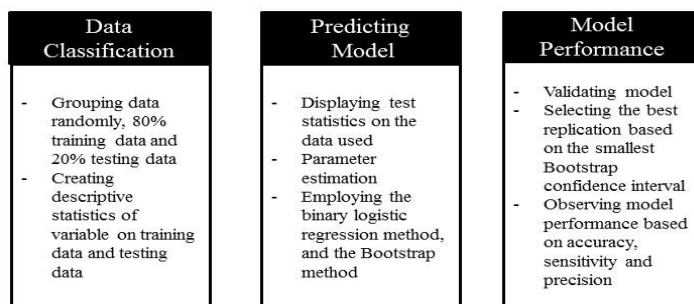


Figure 1: Research Method

In this study, data on cases of COVID-19 patients were used with a total of 458 patients with COVID-19. In this data, the researcher subjectively divided the data into two parts and the data are taken randomly, namely 80% training data and 20% testing data [8]. Training data were later used to build a model, while testing data were used to validate the model obtained. The following is a description of the training data and testing data of patients with COVID-19 in Tables 1 and 2.

Table 1: Description of COVID-19 Patient Training Data

Variable	Category	Frequency	Percentage(%)
Gender	Woman	209	57.10
	Man	157	42.80
Age	Age 0-59 Year	261	71.30
	Age >59 Year	105	28,60
Comorbid	Zero Comorbid	145	39.60
	Comorbid 1-3	220	60,10
	Comorbid >3	1	0.27

Table 2: Description of COVID-19 Patient Data Testing

Variable	Category	Frequency	Percentage(%)
Gender	Woman	45	48,90
	Man	47	51.08
Age	Age 0-59 Year	74	80.43
	Age >59 Year	18	19.56
Comorbid	Zero Comorbid	30	32.60
	Comorbid 1-3	44	47,80
	Comorbid >3	18	19.56

The variables used in this study were categorical data, and were first transformed using dummy variables. The dummy variable for the gender category is female (X_{1D1}), and as a comparison to the gender is male. The age level category with dummy variables is the age level 0 – 59 years (X_{2D1}), and the patient is >59 years old as a comparison. The dummy variables for the number of comorbid categories were having no comorbidities (X_{3D1}), having 1 - 3 comorbidities (X_{3D2}), and the comparison category of having > 3 comorbids.

In modeling research data on the recovery rate of COVID-19 patients, regression analysis techniques were used. Regression analysis is a statistical analysis that models the causal relationship between the response variable and the predictor variable. There are two regression analyzes, namely simple and multiple linear regression analysis, however in this study the response variable has a dichotomous categorical measurement scale or nominal scale , therefore, the analysis cannot be used [5]. One of the regression methods used is binary logistic regression analysis. Binary logistic regression analysis was combined with the Bayesian method for estimating the model parameters [9], [10], and [3]. Bayesian method is parameter estimation based on likelihood function (1) and the distribution of priors (2).

$$L(p) = \prod_{i=1}^n (p^{x_i} (1 - p)^{1-x_i}) \tag{1}$$

where $p(x)$ = the probability that a person is declared recovered from COVID-19. Prior distribution is assumed to be normally distributed with mean (μ) and variance (σ^2), written with $X \sim N(\mu, \sigma^2)$, with the following formula

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \left\{ \exp -\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right\} \tag{2}$$

The combination of the likelihood function and the prior distribution forms the posterior distribution. The estimation of model parameters was obtained by determining the mean and variance of the posterior distribution, and mathematically the formula for the parameters of the posterior distribution model was obtained as follows.

$$L(\beta) = \prod_{i=1}^n \left(\frac{\exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})} \right)^{x_i} \times \left(1 - \left(\frac{\exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})} \right) \right)^{1-x_i} \tag{3}$$

In this study, the identification of the posterior distribution in estimating the model parameters was completed using a numerical approach, namely the MCMC method with the Gibbs sampling algorithm with the help of WinBugs software. The next stage is to perform a convergence test on each parameter. The convergence test of the model parameters was carried out using density plots , trace plots and MC error . If each parameter has converged, it can be said that the alleged model has met the established convergence criteria, so that the model can be accepted. A significance test was later performed on each model

parameter. The predictor variables that significantly affect the response variable means that it can also be said that the corresponding predictor variable has an influence on the recovery rate of COVID-19 patients. Thus the final model of the recovery rate of COVID-19 patients for case data in the study is

$$\hat{p}(x) = \frac{\exp(0,8739+1,144X_{2D1})}{1+\exp(0,8739+1,144X_{2D1})} \tag{4}$$

Furthermore, in order to determine the performance of the algorithm used in the Bayesian binary logistic regression method, the Bootstrap method was used. The application of the Bootstrap method aims to ensure the estimated value obtained, and to test the algorithm used whether it has produced the correct estimated value, where the test is carried out by comparing the 95% Bootstrap confidence interval. In this study, the smallest Bayesian logistic Bootstrap confidence interval was found at 25 replications, namely 1,8481.

Table 3: Confusion Matrix of Two Categories

Actual	Prediction	
	Positive	negative
True	True Positive (TP)	False Negative (FN)
False	False Positive (FP)	True Negative (TN)

This study also validates the model by calculating the classification accuracy based on the Bayesian binary logistic regression model equation (4). If the probability value $\hat{p}(x)$ is greater than or equal to 0.5 then the patient's probability status is grouped into the recovery category, whereas if the probability $\hat{p}(x)$ is less than 0.5 then the patient's discharge status is grouped into the death category [9]. Furthermore, the classification grouping can be estimated as a measure of model performance obtained by calculating accuracy, sensitivity and PPV (Positive Predicted Value) [8] based on the confusion matrix Table 3.

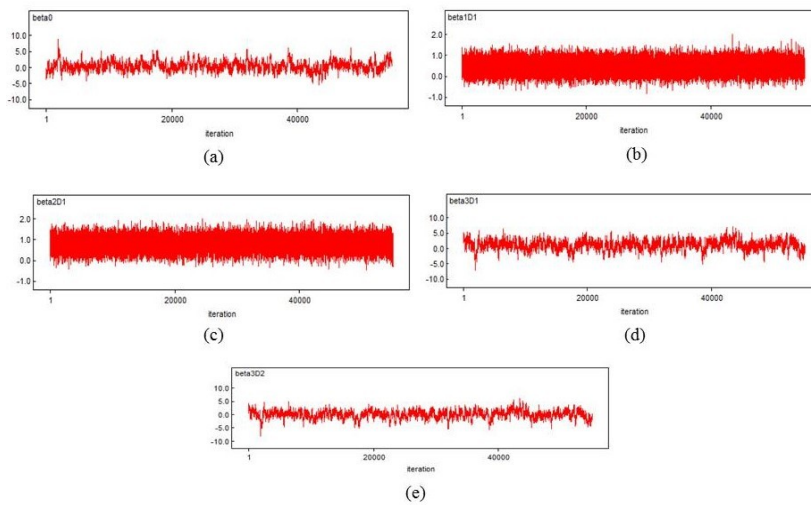
$$\begin{aligned}
 Accuracy &= \frac{TP+TN}{TP+FP+FN+TN} \\
 Sensitivity &= \frac{TP}{TP+FN} \\
 Precision &= \frac{TP}{TP+FP}
 \end{aligned} \tag{6}$$

III. RESULT AND DISCUSSION

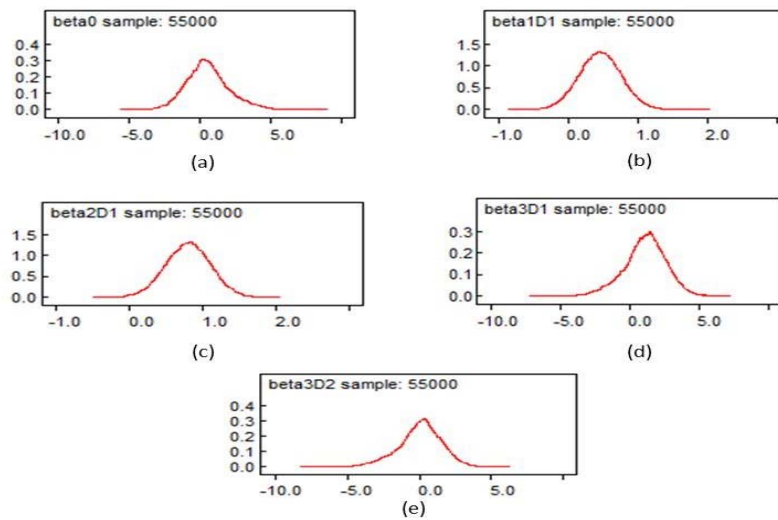
3.1 Classification with Bayesian Binary Logistics Regression (Stage I)

In this study, the training data using Bayesian binary logistic regression was analyzed. Classification was carried out by estimating the parameters using the Bayes method. Parameter estimation using Bayes method is parameter estimation by utilizing the likelihood function with prior distribution in order to obtain a posterior distribution. Parameter estimation in this study uses a numerical approach, namely MCMC with the Gibbs sampling algorithm . The sampling process was carried out using 55,000 iterations. The results of the classification of the estimation of model parameters are said to have met the convergence requirements and the model is acceptable, this can be seen based on the trace plot (Figure 1), density plot (Figure 2) and the comparison of MC Error with 5% standard deviation for each model parameter (Table 3). Based on Figure 1, it can be said that each predicted model parameter has met the assumption of convergence because the data distribution has been stable between two parallel horizontal lines. Figure 2 shows that the density plot curve for each parameter shows a normally distributed curve.

Table 3 shows the comparison of the MC Error value with 5% standard deviation. The results of the convergence test show that the MC Error parameter value of each model parameter shows that the MC Error value is less than 5% standard deviation. Thus, it can be said that, each parameter has converged. It can be concluded that the conjecture model has met the established convergence criteria, and that the model can be said acceptable.



Picture 1: Trace Plot Classification with Bayesian Staged Binary Logistic Regression I for Each Model Parameter



Picture 2: Density Classification Plot with Staged Bayesian Binary Logistic Regression I for Each Model Parameter

Table 3: Test Convergence Parameter Model Based on Comparison of MC Error with 5% Standard Deviation

Predicator Variable	Parameter	Standard of Deviation	5% Standar of Deviation	MC Error	Description
Constant	β_0	1,52	0,076	0,075	Convergent
Gender (X_1)					
Female (X_{1D1})	β_{1D1}	0,3016	0,015	0,002	Convergent
Age (X_2)					
0–59 year (X_{2D1})	β_{2D1}	0,3076	0,015	0,002	Convergent
Comorbid (X_3)					
Zero Comorbid (X_{3D1})	β_{3D1}	1,551	0,0077	0,0075	Convergent
1 – 3 Comorbid (X_{3D2})	β_{3D2}	1,521	0,0760	0,075	Convergent

The significance test of each parameter was conducted to examine whether the predictor variables had a significant effect on the response variables. The significance test for each model parameter using the Bayesian method can be seen based on the Wald statistical value (Table 4). The parameter is said to be significant, if the Wald statistic value is greater than the value $\chi^2_{0.05;1} = 3.841$ and the credible intervals for each model parameter (Table 5).

Table 4: Test Significance Parameter Model Stage I ($\chi^2_{0.05;1} = 3.841$)

Variable predictor	Parameter	Mean	Standard Deviation	Statistics Wald	Description
Constant	β_0	0.5145	1.52	0.1145	–
Gender (X_1)					
Gender Female (X_{1D1})	β_{1D1}	0.4601	0.3016	2,327	Not Significant
Age (X_2)					
Age 0- 59 Year (X_{2D1})	β_{2D1}	0.7938	0.3076	6.6596	Significant
Comorbid (X_3)					
Zero Comorbid (X_{3D1})	β_{3D1}	1.106	1.551	0.5084	Not Significant
1-3 Comorbid (X_{3D2})	β_{3D2}	0.0595	1.521	0.0015	Not Significant

Table 4 shows that there is only 1 (one) predictor variable that has a significant effect on the response variable, namely β_{2D1} (patient has an age of 0-59 years), while the other predictor variables have no significant effect on the response variable. The significance test of the model parameters on the Bayesian method can be seen based on the credible intervals which shows that there is only one predictor variable that significantly affects the response variable, namely β_{2D1} (patient has an age of 0-59 years).

Table 5: Credible intervals Parameter Model

Variable predictor	Parameter	Mean	Standard Deviation	MC Error	Credible intervals		Description
					Limit On	Limit Lower	
Constant	β_0	0.5145	1.52	0.07589	2,228	3,873	-
Gender (X_1)							
Female (X_{1D1})	β_{1D1}	0.4601	0.3016	0,022	0,1291	1,053	Not Significant

				7 2			
Age (X_2)							
Age 0- 59 Year(X_{2D1})	β_{2D1}	0.7938	0.3076	0.00222 7	0,185	1.396	Significant
Comorbid(X_3)							
Zero Comorbid(X_{3D1})	β_{3D1}	1.106	1.551	0.07586	2,264	3,927	Not Significant
1-3 Comorbid(X_{3D2})	β_{3D2}	0.0595	1.521	0.07591	3,288	2,821	Not Significant

3.2 Classification with Bayesian binary logistic regression (Stage II)

This classification stage is a repeated model analysis involving only the predictor variables that have a significant effect on the response variables from the classification results with Bayesian binary logistic regression stage 1. Model analysis was performed by estimating model parameters. Parameter estimation is carried out with 10,000 iterations until the parameter estimation process has reached convergence. This can be seen based on the trace plot and density plot (Figure 3) and the comparison of the MC error with 5% standard deviation (Table 6).

Figure 3 shows both trace plots showing the distribution of data has stabilized which ranges between two parallel horizontal lines and the resulting density plot is rather good because it has a pattern that tends to be smooth in the form of a normal curve. The results of the MC Error value for each of the resulting model parameters are less than 5% standard deviation, which means that the estimated values for the model parameters have converged, which means the model is acceptable.

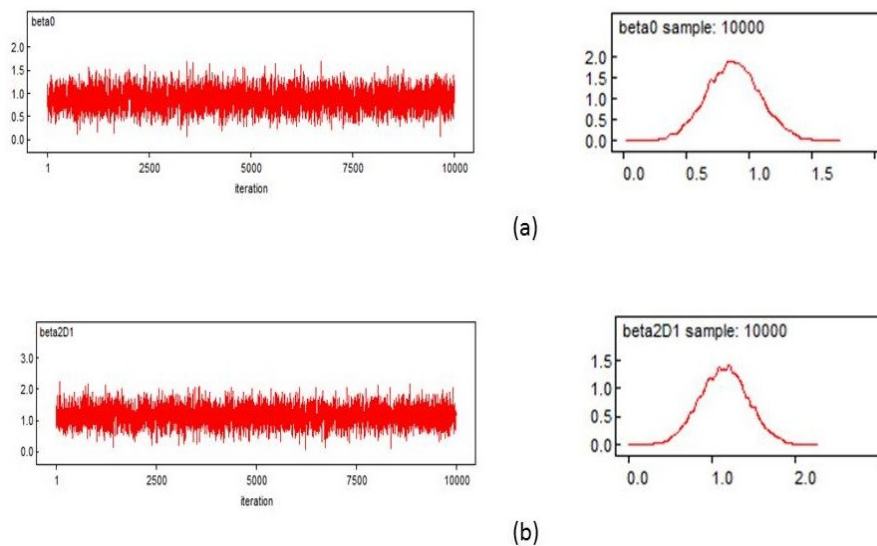


Figure 3: Trace plot and Density Plot on Classification of Parameter Estimation Phase II

Table 6: Test Convergence Parameter Model (Stage II)

Variable predictor	Parameter	Standard Deviation	5% Standard Deviation	MC Error	Description
Constant	β_0	0.2139	0.0106	0.003603	Convergent
Age (X_2)					
0-59 Year (X_{2D1})	β_{2D1}	0.2900	0.0145	0.005042	Convergent

Table 7. Model parameter significance test table ($\chi^2_{0.005,2} = 3,841$)

Variable predictor	Parameter	mean	Standard Deviation	Statistics <i>Wald</i>	Description
Constant	β_0	0.8739	0.2139	16,691	Significant
Age (X_2)					
0-59 Year (X_{2D1})	β_{2D1}	1.144	0.2900	15,561	Significant

The significance test with the Bayesian method can be seen based on the value of the Wald statistical test (Table 7). The significance test of the phase II model parameters was carried out on the predictor variables that were proven to significantly affect the response from the results of the phase I model parameter estimation. The parameter significance testing aimed to see whether the predictor variables had a significant effect on the response variables.

In Table 7, it is found that the Wald statistical value in parameter β_0 is 16.691 which is greater than the value of 2 table, namely 3.841. In parameter β_{2D1} the Wald statistical value is 15.561 which is greater than the value of 2 table which is 3.841. Thus, it can be concluded that there is only one predictor variable, namely β_{2D1} (patients aged 0 – 59 years) which significantly affected the response variable. This means that the corresponding predictor variable has an influence on the recovery rate of COVID-19 patients.

3.3 Parameter Estimation with Bayesian Binary Logistic Regression with Bootstrap Metode Method

After estimating the model parameters using the Bayesian binary logistic regression method, then the performance of the algorithm used was tested. The algorithm in question is the algorithm used to estimate the parameters of the Bayesian binary logistic regression model whether it has produced the correct estimated value. The test is carried out by comparing the predicted value of Bayesian binary logistic regression with a 95% Bootstrap confidence interval. The estimated model can be accepted because the results of the parameter estimates using the Bayesian binary logistic regression method are at the 95% Bootstrap confidence interval, and it can also be said that the algorithm that has been compiled is able to produce the correct estimated values. This can be seen based on the estimation results of the Bayesian binary logistic regression parameter using the Bootstrap method (Table 8) and the width of the Bayesian binary logistic Bootstrap confidence interval (Table 9).

Table 8 shows that the parameter estimation results with Bayesian binary logistic regression for each model parameter are already in the Bootstrap confidence interval. Thus, it can be said that the estimated value of the model parameters obtained has produced the appropriate value. This means that the algorithm used in the Bayesian binary logistic regression method to estimate the parameters of the COVID-19 patient recovery rate model is correct and is able to produce the appropriate value. In Table 9 it can be seen that the smallest Bayesian logistic Bootstrap confidence interval is at replication 25, which is 1.8481. The width of the Bootstrap confidence interval generated for each replication fluctuated, this was due to iterations carried out on the Bootstrap method itself ([7], [11], [12]).

Table 8: Estimation Results of Bayesian Binary Logistics Regression Parameters Using Bootstrap Method

Replication	Parameter	Score Estimate Logistics Bayesian	Score Estimate Bootstrap	Confidence Interval Bootstrap95%	
				Lower Limit	Upper Limit
25	β_{2D1}	1.144	1.0555	0.1314	1.9796
50	β_{2D1}	1.144	1.0677	0.1048	2,030
100	β_{2D1}	1.144	1.2484	0.1418	2.3549

Table 9. Width of Bayesian Logistic Bootstrap Trust Interval

Replication	Parameter	Confidence Interval Bootstrap 95%		Confidence Interval Width 95% Boost
		Limit Lower	Limit On	
25	β_{2D1}	0.1314	1.9796	1.8481
50	β_{2D1}	0.1048	2,030	1.9258
100	β_{2D1}	0.1418	2.3549	2.2130

3.4 Validating the Bayesian Binary Logistics Regression Model

In this study, model validation was also carried out by calculating the classification accuracy based on the Bayesian binary logistic regression model equation (4). The results of the testing data grouping can be obtained based on Table 3 and the results of the estimated performance measures of the model obtained by calculating the values of accuracy, sensitivity and precision (Positive Presented Value). In this case, the accuracy value of the Bayesian binary logistic regression model is 75%, sensitivity is 100%, and precision (PPV) is 75% for case data of COVID-19 patients. The value of accuracy and sensitivity in general for cases about health or in particular predicting disease must have a high value of accuracy and sensitivity in order to get a more precise or appropriate treatment [8].

3.5 Bayesian Logistics Regression Coefficient Interpretation

Interpretation was used to see how much influence the predictor variable has on the recovery rate of COVID-19 patients, then the odds ratio (OR) is used, the value of which can be seen in Table 10. Based on Table 11, the odds ratio value for the predictor variable for patients aged 0-59 years (X_{2D1}) is 3,1393. This means that COVID-19 patients aged 0 – 59 years have a recovery rate of 3.1393 times compared to COVID-19 patients aged > 59 years.

Table 10 . Bayesian Binary Logistics Odd Ratio

Predictor Variable	Parameter	mean	Odds Ratio
Age 0-59 Year (X_{2D1})	β_{2D1}	1.144	3.1393

IV. CONCLUSIONS

Based on the results of the analysis carried out, it can be concluded that there is one predictor variable that significantly affects the recovery rate of COVID-19 patients, namely patients aged 0 - 59 years however other predictor variables (gender, number of comorbidities) had no significant effect on the recovery rate of COVID-19 patients.

The Bayesian binary logistic regression method used in modeling the recovery rate of COVID-19 patients at M. Djamil Hospital and Andalas University Hospital, was accurate based on the algorithm test applied with the Bootstrap method. This study proves that the estimated value with Bayesian binary logistic regression is in the 95% Bootstrap confidence interval. The estimated value with 25 replications produces a better estimate value than that in some other replications, such as 50 and 100 replications.

The classification model for the recovery rate of COVID-19 patients with the Bayesian binary logistic regression method has good performance by producing high accuracy values, and has high sensitivity and PPV in identifying patients who are in fact recovered. This means that the model has been effective to accurately classify the patient's recovery rate from the case data of COVID-19 patients at the M. Djamil hospital and the Hospital of Andalas University, in Padang, West Sumatra.

REFERENCES

- [1] World Health Organization, "Coronavirus", 2020. www.who.int. Accessed June 10, 2021.
- [2] Zhai. P., Ding. Y., Wu. X., Long. J., Zhong. Y., and Li. Y. "The Epidemiology, Diagnosis, and Treatment of COVID-19. Elsevier BV and International Society of Chemotherapy." International Journal of Antimicrobial Agents. Vol. 5, pp. 1-12, 2020.
- [3] Das. A. K., and Gopalan. S. S. "Epidemiology of COVID-19 and Predictors of Recovery in the Republic of Korea." Journal of Hindawi: Pulmonary Medicine. (2020), 7291698, 2020.
- [4] Tolossa, T., Wakuma. B., Gebre. D.S., Atomssa. E.M., Getachew. M., Fetensa. G., Ayala. D., and Turi. E. "Time to Recovery from COVID-19 and Its Predictors among Patients Admitted to Treatment Center of Wollega University Referral Hospital (WURH) Western Ethiopia: Survival Analysis of Retrospective Cohort Study." Journal PLOS ONE. Vol. 16, 0252389, 2021.
- [5] Harlan, Johan. "Linear Regression Analysis." Jakarta: Gunadarma, 2018.
- [6] Yanuar, F., Yozza, H., and Rescha, R. V. "Comparison of Two Priors in Bayesian Estimation for Parameter of Weibull Distribution. Indonesian Science and Technology." Vol. 2, pp. 108-113, 2019.
- [7] Efron. B., and Tibshirani, J.R. "An Introduction to the Bootstrap." New York: Champman and Hall, Inc. 1993.
- [8] Yulia. R., Kresnawati. S. E., Dewi, N. R., Zayanti. D. A., and Eliyanti. N. "Diagnosis of Diabetes Mellitus in Women of Reproductive Age using The Prediction Methods of Naive Bayes, Discriminant Analysis, and Logistic Regression." Indonesian Science and Technology Journal. Vol. 6, pp. 96-104, 2021.
- [9] Teti. M S., Yanuar. F., and Yozza. H. "Logistics Regression Analysis with Bayes Estimator to Determine Factors Affecting the Incidence of Low Birth Weight Babies". UNAND Mathematics Journal. Vol. 4, pp. 53-60, 2018.
- [10] Anjullo. B. B., and Haile. T. T. "A Bayesian Binary Logistic Regression Approach in Identifying Factors Associated with

Exclusive Breastfeeding Practices at Arba Minch Town, South Ethiopia.” *Advances in Research*. Vol. 17, pp. 1-14, 2018.

[11] Kariyam, Qairlina. “Prediction of Refractive Disorders Based on the Length of the Eyeball Axis in Axial Myopia Patients Through Bootstrap Regression.” *National Seminar on Mathematics and Mathematics Education 2006 with the theme of Trends in Mathematics Research and Learning in the ICT Era*. 269-279, 2006.

[12] Putra, A., Tiro M. A., and Aidid. M. K. “Bootstrap and Jackknife Methods in Estimating Multiple Linear Regression Parameters.” *Journal of Statistics and Its Application on Teaching and Research*. Vol. 1, pp. 32-39, 2019.