

Statistical Downscaling Model with Jackknife Ridge Regression and Modified Jackknife Ridge Regression to Forecast Rainfall

Sitti Sahrman^{1‡}, Dewi Santika Upa¹

¹Statistics Department, Hasanuddin University, Indonesia

[‡]Corresponding author: sittisahrimansalam@gmail.com

Copyright © 2024 Sitti Sahrman and Dewi Santika Upa. This is an open-access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Statistical downscaling is a transfer function that connects local scale rainfall data with global scale rainfall. Global-scale rainfall can be obtained of the global circulation model output. The global circulation model simulates climate variables in the form of a large-scale grid which causes high correlation between the grids (multicollinearity). The methods used in statistical downscaling modeling to overcome multicollinearity are jackknife ridge regression and modified jackknife ridge regression. The methods are the development of the ridge regression method. This study aimed to forecast local-scale rainfall data in Pangkep Regency (response variable) based on global-scale precipitation from global circulation model output (predictor variables) using the jackknife ridge regression and modified jackknife ridge regression methods. In addition, the K-means cluster technique was used to determine the dummy variable to overcome the heterogeneity of the error variance. The results using training data (1990-2017 period) showed that the modified jackknife ridge regression method was better at explaining the diversity of data, based on a higher coefficient of determination value (68%) and a lower root mean square error value (165.57) than the jackknife ridge regression method (coefficient of determination 67% and root mean square error 167.72). Model validation using testing data (period of 2018) also showed the same results, i.e., the modified jackknife ridge regression was better than the jackknife ridge regression. Furthermore, the addition of dummy variables increased the accuracy of the model in forecasting rainfall data. The addition of dummy variables to the model resulted in a high coefficient of determination (ranging from 94% to 95%) with lower root mean square error values (ranging from 66.60 to 67.69).

Keywords: dummy variables, global circulation model, jackknife ridge regression, modified jackknife ridge regression, statistical downscaling

1. Introduction

Climate change has a very wide impact on people's lives. The increase in the earth's temperature not only has an impact on the rise in the earth's temperature but also changes the climate system which affects various aspects of natural and human life changes, such as the quality and quantity of water, habitat, forests, health, agricultural land, and coastal ecosystems. In addition, the impact of climate change can also affect national salt production which has been very dependent on the

weather. During normal weather, salt production is relatively stable so that it can provide benefits to salt farmers. However, when there is a change in the weather, salt production can decrease or even fail to harvest. This makes the study of climate change very necessary to minimize the losses that might occur.

Pangkep Regency is one of the four districts in South Sulawesi that is the center of salt production. Director of Marine and Fisheries Services Marine and Fisheries Ministry Republic of Indonesia, Rianto Basuki, said that Pangkep Regency was included in 10 salt production centers in Indonesia. Salt production has increased by tens of thousands of tons in Pangkep Regency in 2018. Its production has reached approximately 36 thousand tons when compared to the dry season production in 2017 which only reached 5,700 tons. The increase in production is due to weather factors. Therefore, daily and monthly rainfall prediction is very important information and is needed by the salt farmers (Subhan, 2018).

Rainfall is one of the most important climate elements in Indonesia. This is due to the very high diversity according to time and place. Therefore, studies on climate change are more directed at the rainfall factor. Related to the climate in Indonesia, the process of rain formation in the tropics is the most difficult process to simulate. Until now there has not been a single climate model that is able to simulate rainfall patterns in Indonesia properly (Sahrinan, 2014). Hence, high resolution climate models need to be developed at local scales by utilizing global climate information such as global circulation models.

Global circulation model (GCM) is one method that can simulate climate in the past, present, and forecast climate changes that might occur in the future (Wilby, et al., 2009). GCM simulates global climate variables on each grid (measuring $\pm 2.5^\circ$ or ± 300 km²) for each atmosphere layer which is then used to forecast climate patterns on an annual basis (Wigena, 2006). However, GCM is still on a global scale so it cannot explain the local climate conditions such as local rainfall. Therefore, statistical downscaling (SD) can be used to estimate local scale climate variables by utilizing statistical techniques (Zorita & Storch, 1999). Generally, the approach used in the SD technique is regression analysis in determining the functional relationship of global scale climate variables with local scale climate variables.

GCM output climate variables have large dimensions and high correlations between the grids. This allows the occurrence of multicollinearity in the data. Multicollinearity is an ill condition which causes the multiplication of the predictor matrix ($X'X$) does not meet the classical assumptions of regression. Multicollinearity causes the standard error of estimating the regression parameters to be large and the value of the confidence interval to be wide so that the use of the least squares method (LSM) becomes invalid (Montgomery & Peck, 1992). Ridge Regression (RR) is one method that can be used to overcome multicollinearity problems. RR is a modified method of LSM by adding the bias constant (k) to the diagonal matrix. Furthermore, some of the development of the RR method namely Generalized Ridge Regression (GRR) by Hoerl & Kennard (1970) and Jackknife Ridge Regression (JRR) by Singh et al. (1986). GRR adds several biased constants to the estimation of the parameters. Meanwhile, according to research by Singh et al. (1986), JRR eliminates one data and repeats as many samples as there are to avoid bias in the RR estimator so that the JRR estimator's Mean Square Error (MSE) is smaller than the GRR estimator. Then, Batah et al. (2008) introduced the parameter estimation method, i.e., Modified Jackknife Ridge Regression (MJR) which is a combination of GRR and JRR. MJR is a parameter estimation method that combines ideas from GRR and JRR regression.

Research on rainfall models in Indonesia has been carried out by several

researchers before. Sahriman et al. (2019) used Liu-type to forecast monthly rainfall in Pangkep Regency as a salt center. In addition, Sahriman et al. (2019) added dummy variables based on hierarchical and non-hierarchical clustering techniques in SD modeling to forecast rainfall. Devita et al. (2014) estimated parameters using JRR in overcoming multicollinearity. Batah et al. (2008) estimated the regression parameters using the MJR method.

This research used the SD model with the JRR and MJR methods to model rainfall data. Furthermore, the addition of dummy variables based on the K-means cluster technique was used to improve the performance of the SD model with the JRR and MJR methods. The dummy variables were used as additional predictor variables. The aim of this research was to obtain forecast rainfall in Pangkep Regency based on the SD model with the JRR and MJR methods. The accuracy of rainfall forecasting results was assessed based on the mean square error and correlation value between these methods.

2. Methods

The data used in this research was a global scale of monthly precipitation period of 1996-2017 (GCM outputs climate models inter-comparison project that measures 8x8 grid) and rainfall Pangkep Regency local scale (period of 1996-2018). The geographical position at 119.57 ° EL until 129.37 ° WL and -14.83 ° SL until 5.17 ° NL above the area of the Pangkep Regency area was used as GCM domain. There were 67 predictor variables used in this study, derived from the domain size 8x8 grid GCM data and three dummy variables. Rainfall in Pangkep Regency was used as predictor variables in this study. The data used for modeling was from the 1966-2017 period, and the data for 2018 was used for model validation. Multicollinearity identification on the data of precipitation using the variance inflation factors (VIF),

$$VIF_j = \frac{1}{1-R_j^2}, \quad j = 1, 2, \dots, 64$$

where R_j^2 is coefficient of determination of the regression result of predictor variable j with other predictor variables. In addition, the dummy variables are determined from the k-means cluster technique.

The SD model in this study uses the JRR and MJR methods. According to Singh et al. (1986), the JRR method is used to overcome bias in the RR method. Estimation of the regression parameters in the jackknifed method is done by eliminating one data and repeating it as much as the size of the data (Iskandar et al., 2013).

$$y_{-i} = X_{-i}^* \alpha + \varepsilon^*$$

where y_{-i} is the vector y with the i -th value removed, X_{-i}^* is the X^* matrix with each i -th row removed, and ε^* is the error vector with the i -th coordinate removed.

JRR parameters are notated, $\hat{\alpha}_{JR}$

$$\hat{\alpha}_{JR} = [I + A^{-1}K] \hat{\alpha}_{GR}$$

where $\hat{\alpha}_{GR} = (I - A^{-1}kI)\alpha_{LS}$ is the parameter of the GRR, $A = (X^*X^* + kI)$ is a fixed value, k is a constant value, X^* is X matrix that has been transformed by centering and scaling bias. Because $\alpha = T'\beta$ and $\beta = T\alpha$ with T is an orthogonal $p \times p$ -sized matrix whose elements are the vector eigenvalue of $X'X$, then the JRR coefficient can be formulated as follows,

$$\hat{\beta}_{JR} = T\hat{\alpha}_{JR}$$

According to Khurana et al. (2014), the variance and MSE of the JRR method can be obtained,

$$\begin{aligned} \text{Var}(\hat{\alpha}_{JR}) &= \sigma^2(\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2})(\mathbf{X}^*\mathbf{X}^*)^{-1}(\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2})' \\ \text{MSE}(\hat{\alpha}_{JR}) &= \sigma^2(\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2})(\mathbf{X}^*\mathbf{X}^*)^{-1}(\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2})' + \mathbf{K}^2\mathbf{A}^{-2}\alpha\alpha'\mathbf{A}^{-1}\mathbf{K} \end{aligned}$$

Bata et al. (2008) proposed a new estimator, i.e., the MJR method which combines the GRR and JRR ideas of Singh et al. (1986). The MSE value of the MJR method is smaller than the GRR and JRR methods (Khurana et al., 2014). The MJR parameter, $\hat{\alpha}_{MJR}$, is obtained by

$$\hat{\alpha}_{MJR} = [\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2}][\mathbf{I} - \mathbf{K}\mathbf{A}^{-1}]\hat{\alpha}_{LS}$$

where $\alpha = \mathbf{T}'\beta$ and $\beta = \mathbf{T}\alpha$, so that the MJR coefficient can be formulated

$$\hat{\beta}_{MJR} = \mathbf{T}\hat{\alpha}_{MJR}$$

where Variance and MSE can be obtained by,

$$\begin{aligned} \text{Var}(\hat{\alpha}_{MJR}) &= \sigma^2\mathbf{W}(\mathbf{X}^*\mathbf{X}^*)^{-1}\mathbf{W}' \\ \text{MSE}(\hat{\alpha}_{MJR}) &= \sigma^2\mathbf{W}(\mathbf{X}^*\mathbf{X}^*)^{-1}\mathbf{W}' + \mathbf{K}\Phi\mathbf{A}^{-1}\gamma\gamma'\mathbf{A}^{-1}\Phi'\mathbf{K} \end{aligned}$$

and $\mathbf{W} = (\mathbf{I} - \mathbf{K}^2\mathbf{A}^{-2})(\mathbf{I} - \mathbf{K}\mathbf{A}^{-1})$ and $\Phi = [\mathbf{I} + \mathbf{K}\mathbf{A}^{-1} - \mathbf{K}\mathbf{A}^{-2}]$.

3. Results and Discussion

3.1 Multicollinearity Identification

Multiple regression analysis assumes that there is no multicollinearity in the predictor variables. Multicollinearity in data can be detected through a significant VIF value (VIF > 10). The analysis showed that the precipitation data of GCM had VIF values ranging from 5.01 to 2942.24. This indicated that there was high multicollinearity between adjacent GCM data grids. Therefore, SD modeling was used with the JRR and MJR methods.

3.2 Determination Dummy Variables with K-Means

K-means method is a clustering method to utilize the concept of a centroid or midpoint. There are several k centroids that can be set at the beginning. Sahriman et al. (2019) used this method to determine dummy variables, and the results increased the accuracy of rainfall forecasting. Based on the results of this research, rainfall data with 4 clusters was the optimal grouping. The Partial Least Square method produced a plot between component scores of Y and components scores of X (Figure 1). The rainfall formed 4 groups based on the colors that appeared. Group 1 is rainfall with an intensity of more than 1019 mm/month. Group 2 is rainfall with an intensity of 608-1019 mm/month. Group 3 is rainfall with an intensity of 233-607 mm/month. Group 4 is rainfall with an intensity of 0-232 mm/month. Thus, the number of dummy variables used in SD modeling is 3 variables (D_1 , D_2 , dan D_3).

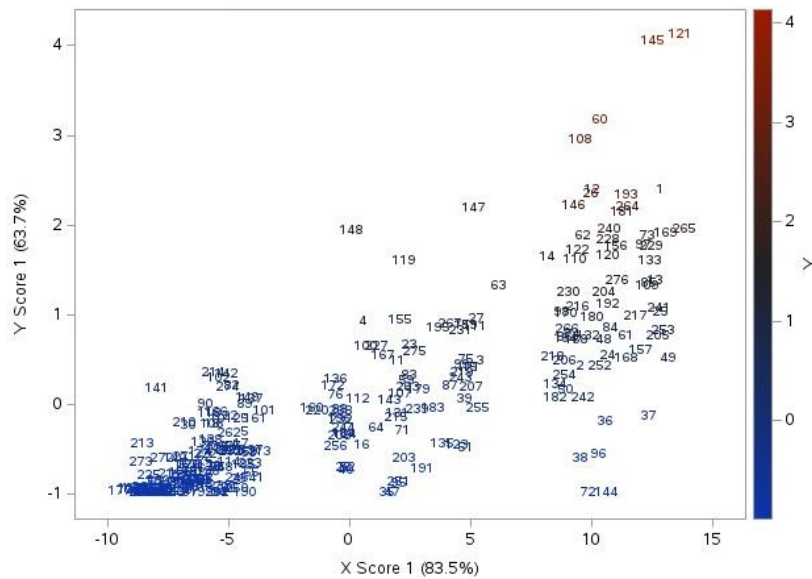


Figure 1. Groups of Rainfall Data.

Table 1 represented the dummy variables of D_1 , D_2 , and D_3 . The dummy variables were determined based on the rainfall data group from the K-means method. Group 1 had dummy values of $D_1 = D_2 = D_3 = 0$. Furthermore, group 2 had $D_1 = 1$ dan $D_2 = D_3 = 0$. Group 3 had $D_2 = 1$ dan $D_1 = D_3 = 0$. Meanwhile, group 4 had $D_1 = D_2 = 0$ dan $D_3 = 1$. Rainfall data for the period of 1997-2018 was grouped into 4. Groups 1 and 2 respectively consisted of 4 and 32 observations which generally occur in December and January. Group 3 consisted of 92 observations which generally occur in November, February, March, and April. Meanwhile, group 4 consisted of 136 observations and generally occur in May to October. Thus, the intensity of rainfall in Pangkep Regency generally occurs at low intensity.

Table 1. Dummy Variables

Time	Y	D_1	D_2	D_3
Jan-97	593	0	1	0
Feb-97	401	0	1	0
Mar-97	28	0	0	1
⋮	⋮	⋮	⋮	⋮
Jan-11	1519	0	0	0
Feb-11	967	1	0	0
⋮	⋮	⋮	⋮	⋮
Nov-18	245	0	1	0
Dec-18	918	1	0	0

3.3 Statistical Downscaling with Jackknife Ridge Regression and Modified Jackknife Ridge Regression

Jackknife Ridge Regression (JRR) and Modified Jackknife Ridge Regression (MJR) are methods used to overcome multicollinearity in precipitation of GCM data. The initial stage in SD modeling using the JRR and MJR methods is the determination of the bias constant (k). By using $\hat{\sigma}^2 = 0.001566$ of MKT, the value of the bias constant k was equal to 0.00071. Table 2 showed the results of SD modeling using the RR, JRR, and MJR methods on training data. Table 2 presented the results of the SD model using the RR, JRR, and MJR methods with training data. Table 2 explained that the MJR method was better at explaining data diversity than the RR and JRR methods. The R^2 value of the MJR method was higher (68%) than that of the RR and JRR methods. Additionally, the RMSE value (165.57) of the MJR method was lower. Furthermore, the addition of dummy variables to each model (RR, JRR, MJR) could increase the accuracy of the model, with R^2 values ranging from 68% to 95%. The addition of a dummy variable could reduce the RMSE values of the model around 59.52-100.04. In general, the MJR model with dummy variables was the best model based on R^2 and RMSE values.

Table 2. The R^2 and RMSE of RR, JRR, and MJR Models

Method	R^2	RMSE
RR	56%	190.93
JRR	67%	167.72
MJR	68%	165.57
RR Dummy	68%	134.93
JRR Dummy	94%	67.69
MJR Dummy	95%	66.60

Residual diagnosis was performed on the RR, JRR, and MJR models. The plot between the residuals and the estimated rainfall values from the RR model showed a heterogeneous variance of the residuals (Figure 2). The residual scatter of the RR model produced a diverging pattern. The residual pattern at high rainfall was more diverse than that at low rainfall. This was also shown by the residual pattern of the JRR and MJR models (Figure 2). Therefore, dummy variables were added to the SD model to overcome the heterogeneous variance of the residuals. The residual diagnostics of the RR, JRR, MJR models with dummy variables showed 4 groups of residuals. The formation of these group was caused by the addition of dummy variables in the model. However, the residual patterns of the JRR dummy and MJR dummy models in Figure 3 were more homogeneous than the residual patterns of the RR, JRR, and MJR models without dummy variables. Thus, the addition of dummy variables to the JRR and MJR models produced a model with a homogeneous residual variance.

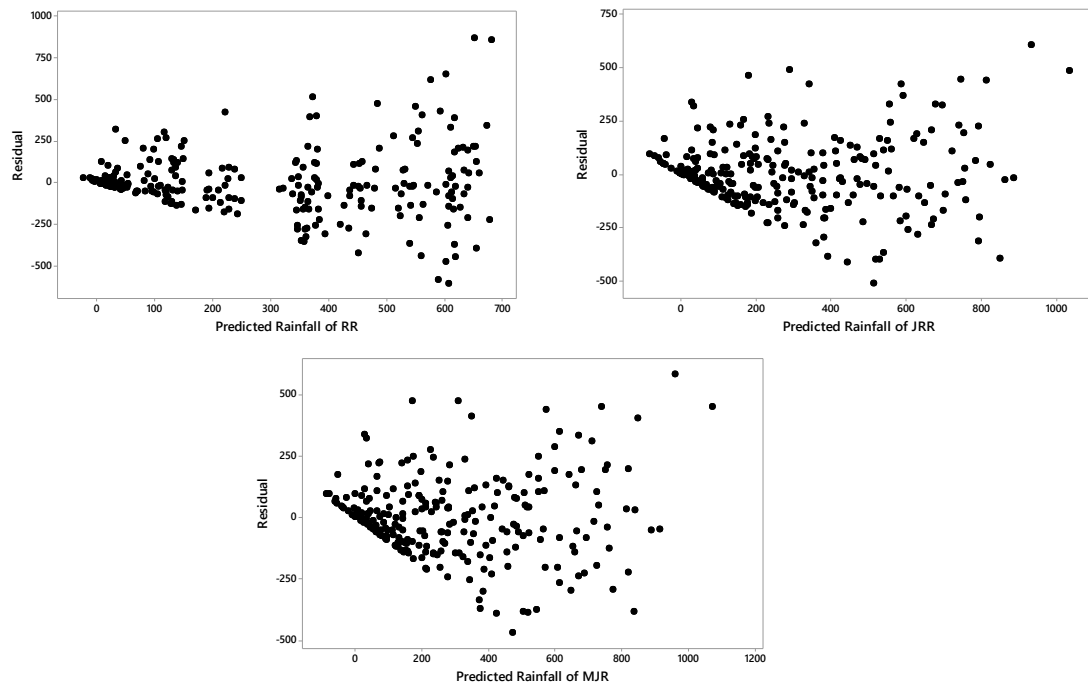


Figure 2. Residual Plots of RR, JRR, and MJR Models without Dummy Variables

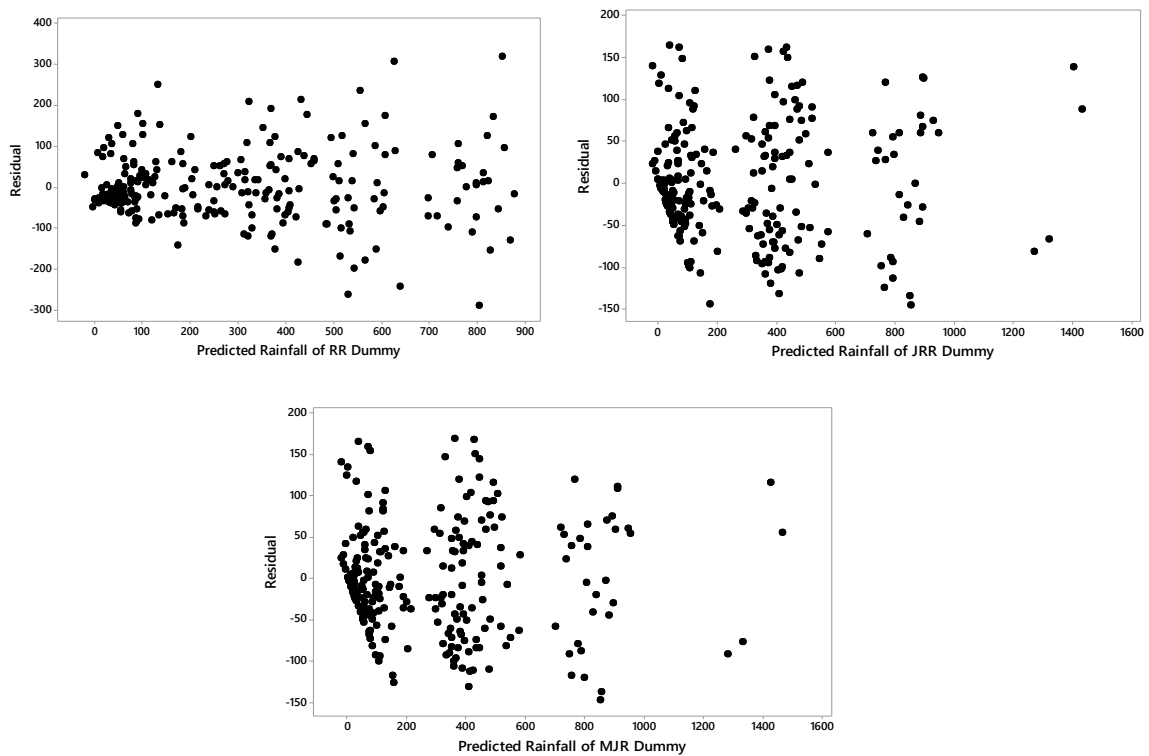


Figure 3. Residual Plots of RR, JRR, and MJR Models with Dummy Variables

Estimator of the best model of each method could be written as follows:

- RR Dummy Method

$$\hat{y} = -51.17 + 2.39 X_1 + \dots - 1.70 X_{64} + 205.38 D_1 - 35.85 D_2 - 205.10 D_3$$
- JRR Dummy Method

$$\hat{y} = 1339.10 - 29.93 X_1 + \dots - 2.25 X_{64} - 522.19 D_1 - 896.37 D_2 - 1164.22 D_3$$

- MJR Dummy Method

$$\hat{y} = 1374.47 - 36.97 X_1 + \dots + 0.56 X_{64} - 552.44 D_1 - 926.15 D_2 - 1195.97 D_3$$

3.4 Model Validation and Selection

Modeling of rainfall data showed that the SD model with the RR, JRR, and MJR methods produced variance from the heterogeneous residual model. The residual plot of the model formed a diverging pattern. The addition of dummy variables to the model showed that the RR Dummy, JRR Dummy, and MJR Dummy methods produced a relatively homogeneous variance of residuals. However, the MJR Dummy was better than the RR Dummy and JRR Dummy methods. Next, model validation used rainfall data from the 2018 period. Model validation was used to test the model in forecasting data. The statistic used was the correlation between predicted rainfall and actual rainfall. The correlation value describes the suitability between the model estimator and the new data. In addition, the root mean square error of prediction (RMSEP) value is used to assess the accuracy of the model in predicting data.

Table 3 presented the correlation and RMSEP values of rainfall predictions. Predicted rainfall using dummy variables was more accurate than without dummy variables. The RR Dummy method produced predicted rainfall with a high correlation (0.893) and a lower RMSEP (148.500) than the RR method without dummy variables (correlation of 0.785 and RMSEP of 206.515). Furthermore, the prediction of rainfall using the SD model with the JRR Dummy method (RMSEP of 117.631 and a correlation of 0.938) was better than the JRR method without dummy variables (RMSEP of 163.339 and a correlation of 0.876). The same results were shown in the SD model with the MJR Dummy and MJR without the dummy. Additionally, predicted rainfall using the JRR Dummy and MJR Dummy methods produced predictors that were relatively the same based on correlation values. In general, the MJR Dummy method predicted rainfall more accurately based on a lower RMSEP than others.

Table 3. Correlation and RMSEP Values of SD Models

Method	Correlation	RMSEP
RR	0.785	206.515
JRR	0.876	166.188
MJR	0.885	163.339
RR Dummy	0.893	148.500
JRR Dummy	0.938	117.631
MJR Dummy	0.938	117.497

Figure 4 showed a plot between the actual rainfall for the 2018 period and the predicted results using the RR, JRR, and MJR methods. The RR method produced lower predicted rainfall than actual rainfall in the periods of February, March, and December. However, it was higher in January, April, May, and November. The RR model could predict rainfall with high accuracy in June, August, and September. However, this method could not predict the actual rainfall well in the period from January to May and December. The JRR and MJR methods could predict rainfall better than the RR method. The distance between the actual and the predicted results of the JRR and MJR methods was relatively close compared to the RR method. In general, the RR, JRR, and MJR methods failed to capture rainfall

patterns well, except for the period from August to October, which had low rainfall.

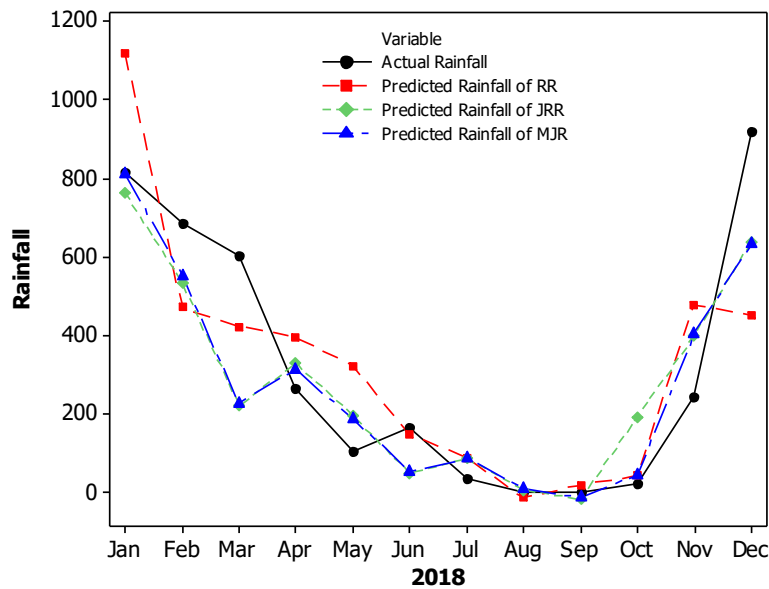


Figure 4. Actual and Predicted Rainfall Plot of SD Models

The JRR Dummy method produced relatively the same predicted rainfall patterns (Figure 5). This method produced predicted rainfall following the actual rainfall pattern compared to the RR Dummy model. Additionally, the distance between the actual rainfall and the predicted rainfall produced was closer than the RR Dummy. Similar to the JRR Dummy method, the predicted rainfall pattern by the MJR Dummy method also showed the same pattern. The distance between the actual rainfall and the predicted rainfall using the JRR Dummy and MJR Dummy methods was relatively close. In general, the SD model with the addition of dummy variables in each method produced predicted rainfall that followed the actual rainfall pattern compared to the SD model without dummy variables. The JRR Dummy and MJR Dummy methods were able to capture rainfall patterns well from January to December. This method could forecast the actual rainfall data with accuracy in the periods of January, May, and July to October. In general, the SD model with the MJR Dummy method was the best model because it could produce a more accurate predicted rainfall.

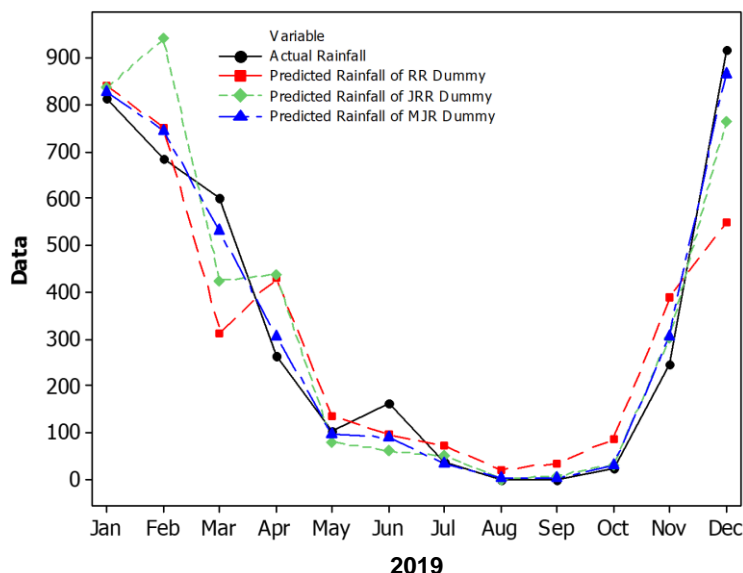


Figure 5. Actual and Predicted Rainfall Plot of SD Models with Dummy

4. Conclusion

The statistical downscaling models with the jackknife ridge regression and the modified jackknife ridge regression methods were utilized to model and forecast rainfall data that contains multicollinearity. The method began with determining the bias constant. The two methods yielded relatively similar results. However, the modeling results using the data period of 1997-2017 showed that the jackknife ridge regression and modified jackknife ridge regression methods with the addition of dummy variables (the coefficient of determination values, R^2 , ranged from 67% to 68% and root mean square error around 167.72) were better than model without dummy variables (R^2 values ranged from 94% to 95% and value root mean square error around 167.72). Model validation used the data period of 2018. The jackknife ridge regression and modified jackknife ridge regression methods with dummy variables produced more accurate predictions of rainfall data (the correlation between actual rainfall data and predicted rainfall was 0.938). However, in general, the modified jackknife ridge regression method with dummy variables was the best model based on the root mean square error of prediction (117.497).

Acknowledgement.

We are grateful to the leaders of Hasanuddin University through the Hasanuddin University LP2M (Institute for Research and Community Service) for providing funding for Unhas internal research. This research is funded through the "Penelitian Dosen Penasehat Akademik/Academic Advisory Lecturer Research (PDPA)" for the 2020 fiscal year. We also thank to all staff who have provided facilities for the completion of this research.

References

- Batah, F. S., Ramanathan, T. V., & Gore, S. D. (2008). The Efficiency of Modified Jackknife and Ridge Type Regression Estimators: A Comparison. *Survey In Mathematics and its Applications*, 3, 111-122.
- Devita, H., Sukarsa, K. G., & Kencana, P. E. (2014). Jackknife ridge regression performance in overcoming multicollinearity. *Mathematics E-Journal*, 3, 146-153.
- Hoerl, A. E., & Kennard, R. W. (1970). Ridge Regression: biased estimation for Nonorthogonal problems. *Technometrics*. 12: 55-67. *Technometrics*, 12, 55-67.
- Iskandar, R., Mara, M. N., & Satyahadewi, N. (2013). Comparison of Bootstrap and Jackknifed Methods in Estimating Regression Parameters to Overcome Multicollinearity. *Mat.Stat and Its Application Scientific Bulletin*.
- Khurana, M., Chaubey, Y., & Chandra, S. (2014). Jackknifing the Ridge Regression Estimator: A Revisit. *Communications in Statistics - Theory and Methods*, 43(24), 5249-5262.
- Montgomery, D. C., & Peck, E. A. (1992). *Introduction To Linear Regression Analysis* Second Edition. New York: John Willey and Sons Inc.

- Sahriman, S. (2014). Statistical Downscaling Model with Time Lag Data of Global Circulation Model to Forecast Rainfall [Thesis]. Bogor (ID): Bogor Agricultural University.
- Sahriman, S., Anisa, & Koerniawan, V. (2019). Liu-type Regression in Statistical Downscaling Model for Forecasting Monthly Rainfall Salt as Producer Region in Pangkep Regency. *Journal of Physics: Conference Series*, 1341 (092021).
- Sahriman, S., Anisa, & Koerniawan, V. (2019). Statistical Downscaling Modeling with Dummy Variables Based on Hierarchical and Nonhierarchical Cluster Techniques to Forecast Rainfall. *Indonesian Journal of Statistics and Its Applications*, 3(3), 295-309.
- Singh, B., Chaubey, Y. P., & Dwivedi, T. D. (1986). An Almost Unbiased Ridge Estimators. *Sankhya*, 48, 342-346.
- Subhan, M. (2018, December 21). SindoNews.com Makassar. Retrieved from SindoNews.com Makassar: <https://makassar.sindonews.com/berita/18506/4/produksi-garam-meningkat-puluhan-ribu-ton-di-pangkep>
- Wigena, A. H. (2006). Statistical modeling of downscaling with persuit projection regression to forecast monthly rainfall [dissertation]. Bogor (ID): Bogor Agricultural University.
- Wilby, R. L., Charles, S. P., Zorita, E., Timbal, B., Whetton, P., & Mearns, L. O. (2009). A review of climate risk information for adaptation and development planning. *Journal of Climatology*, 29, 1193-1215.
- Zorita, E., & Storch, H. V. (1999). The analog method as a simple statistical downscaling technique: comparison with more complicated methods. *J Clim*, 12, 2474-2489.