

A COMPARATIVE STUDY OF MISSING RAINFALL DATA ANALYSIS USING THE METHODS OF INVERSED SQUARE DISTANCE AND ARITHMETIC MEAN

Ekha Yogafanny^{a,b*}, Djoko Legono^a

^aDepartment of Civil and Environmental Engineering, Faculty of Engineering, Universitas Gadjah Mada, Yogyakarta, 55281, Indonesia

^bDepartment of Environmental Engineering, Faculty of Mineral Technology, Universitas Pembangunan Nasional Veteran Yogyakarta, 55283, Indonesia

Article history

Received

03 May 2021

Received in revised form

28 September 2021

Accepted

31 December 2021

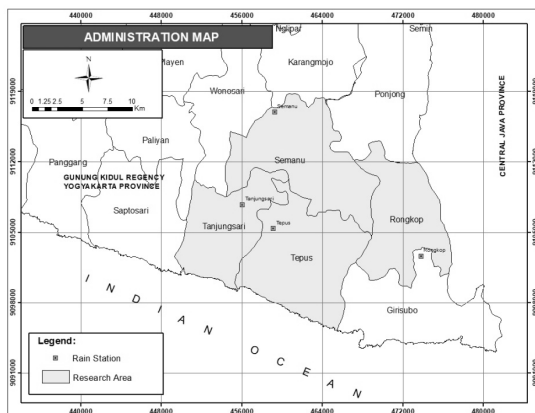
Published online

31 May 2022

*Corresponding author

ekha.yogafanny@mail.ugm.ac.id

Graphical abstract



Abstract

In water resources planning and management, it is essential to have reliable rainfall data. In many cases, rainfall data under the guardian national/ local institution are incomplete. Some data are missing, both monthly and annually. The missing data may persist due to neither damage nor human error. This study aims to estimate the missing rainfall data using two methods, i.e., the inverse square distance and the arithmetic mean methods. The study compared the two mentioned methods using root mean square error (RMSE) and mean absolute error (MAE) and to determine the consistency of rainfall data in all stations using double mass curve analysis. This study utilized the rainfall data from Tepus, Semanu, Rongkop, and Tanjungsari Stations in Gunung Kidul Regency, Yogyakarta Province, Indonesia. The model performance was tested by the root mean square error (RMSE) and mean absolute error (MAE). The rainfall data consistency was determined by double mass curve analysis. The results showed that the arithmetic mean method performed better rather than the inverse square distance method. The smallest RMSE and MAE values in the arithmetic method at the four stations have confirmed the statement. The rainfall data consistency analyzed by the double mass curve is consistent in all stations except Tepus Station.

Keywords: Missing rainfall data, Inversed square distance, Arithmetic mean, RMSE, MAE

© 2022 Penerbit UTM Press. All rights reserved

1.0 INTRODUCTION

Rainfall data are indispensable hydrometeorological data for various purposes, including rainfall characteristics, flood analysis, climate analysis, and water resources management planning [1]. A good quality of rainfall data is highly essential in hydrological and meteorological analysis [2]. Unfortunately, the monthly and annual rainfall data in a region are often incomplete. Insufficient data can cause the deterioration of data's accuracy and lead to ineffective water resources management planning [3]. Therefore, choosing the best method in estimating the missing rainfall data is important. Mostly, the best method that is right in one place is not necessarily correct in another because of seasonal and topographical differences [4].

Damages of the equipment can cause insufficient rainfall data or some missing data. Such conditions need data improvement through systematic methods to estimate the value of the missing rainfall data. Many authors used different ways to estimate the missing data, i.e., Kriging approach [5], regression method, artificial neural network [6], aerial precipitation ratio [7], geometric median [2], etc. Other possible scenarios include the normal ratio [8]–[10], the inversed square distance or reciprocal [8], [11]–[14], arithmetic mean [8]–[10], and the autoregressive model. These methods have advantages and disadvantages, and specific requirements are needed to suit the data's particular conditions and characteristics.

This study compares and analyses the two commonly used methods in the actual case, i.e., the inverse square distance

(ISD) and the arithmetic mean (AM) method. The inversed square distance is one of the widely used methods because it considers the distances between stations as a control [8], [11]. A previous study found that the data's consistency in estimating missing rainfall data using average ratio and the inversed square methods was reasonably good [11]. The arithmetic mean method is another method commonly used to estimate rainfall data. This method is considered the most straightforward, although it does not consider the distances between observed stations and reference stations. Therefore, the results require confirmation with other methods. The determination of methods in estimating the missing rainfall data follows the physical characteristics of the area and the availability of supporting data and human resources. Therefore, the most appropriate method with the smallest possible error is applied in the study area. The two mentioned methods would consider the availability of supporting data and the ease of processing which, if the results are appropriate, can be easily applied by available human resources.

This study estimates the missing rainfall data from four stations in Gunung Kidul Regency, i.e., Semanu, Rongkop, Tanjungsari, and Tepus Station, with the elevation levels are 191 m, 307 m, 277 m, and 255 m, respectively. The location of the study area is presented in Figure 1. The study duration spans 11 years, from 2008 to 2019. The rainfall data gained from Balai Penyuluhan Pertanian (BPP) did not have data for several months at several stations. The agriculture office in that sub-district mainly uses the rainfall data from each BPP of each sub-district. The complete rainfall data are required for various purposes related to water resource management. Therefore, determining the methods in estimating the missing rainfall data is very important and needs to be done before the data utilization. This study aims to evaluate the missing rainfall data using two methods, i.e., the inverse square distance and the arithmetic mean methods. The study compares the methods mentioned using root mean square error (RMSE) and mean absolute error (MAE). The study also determines the consistency of rainfall data in all stations using double mass curve analysis.

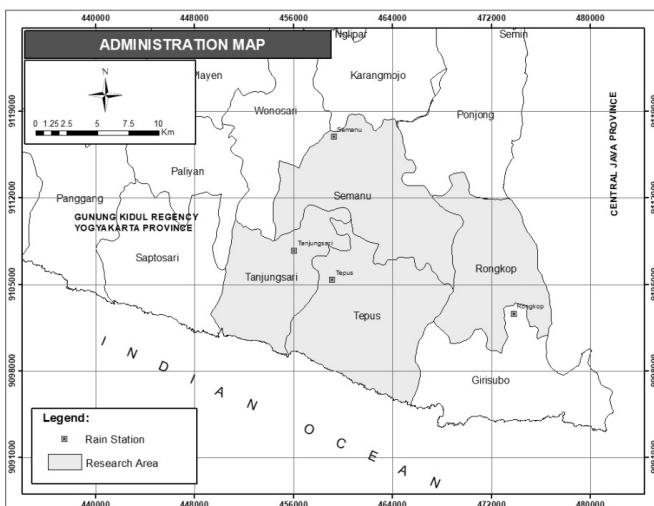


Figure 1 Administration map of the study area

2.0 METHODOLOGY

This study consists of three steps. The first step measures the missing rainfall data using ISD and AM. The second step compares the observed and predicted rainfall data gain from the two mentioned methods using RMSE and MAE to determine the best method applied in the stations. Finally, the third step determines the consistency of complete rainfall data using double mass curve analysis.

In the first step, the missing rainfall data were calculated using the inversed square distances method. The formula used for this method is given in Equation 1:

$$P_x = \frac{\frac{P_A}{(d_{XA})^2} + \frac{P_B}{(d_{XB})^2} + \frac{P_C}{(d_{XC})^2}}{\frac{1}{(d_{XA})^2} + \frac{1}{(d_{XB})^2} + \frac{1}{(d_{XC})^2}} \quad (1)$$

where P_x is the missing rainfall data in station X, P_A is the rainfall data in station A, P_B is the rainfall data in station B, P_C is the rainfall data in station C, d_{XA} is the distance between station X and A, d_{XB} is the distance between station X and B, and d_{XC} is the distance between station X and C.

The other method used here was the arithmetic mean method. The formula used for this method is given in Equation 2:

$$P_x = \frac{P_A + P_B + P_C}{n} \quad (2)$$

P_x is the missing rainfall data in station X, P_A is the rainfall data in station A, P_B is the rainfall data in station B, P_C is the rainfall data in station C, and n is the total reference stations.

The second step aims to measure the RMSE and MAE on each method from all stations. The results were then compared to determine the best method applied to the station. The RMSE is one of the standard statistical metrics to measure the model's error in meteorology and climate research [15]–[18]. The other statistical approach to evaluate the model performance is MAE. This method was also used in other study to measure the consistency of rainfall data [17], [18]. The lower value of RMSE and MAE, the more accurate is the rainfall data. Therefore, the recommended method for each station is obtained from the lowest value of both RMSE and MAE. The RMSE and MAE are given in Equation 3 and 4:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_o - y_p)^2}{n}} \quad (3)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |(y_o - y_p)| \quad (4)$$

The y_o is the observed rainfall data, y_p is the predicted rainfall data, and n is the rainfall data.

The third step is measuring the consistency of the annual rainfall data on each station using double mass curve analysis. The annual rainfall is obtained from the monthly data filled by the two compared methods. The analysis results could show the validity and consistency of the data to be used for other purposes. The formulation used in double mass curve analysis is presented in Equation 5, and a graphical representation of the method is available in Figure 2:

$$Y_Z = \frac{tg \alpha}{tg \alpha_o} Y \quad (5)$$

Yz is the corrected data, Y is the observed data, tg α is the slope before the breakpoint, and tg αo is the slope after the breakpoint. The scattered points resulted from the double mass curve analysis represents the consistency of the data. The more the points converge on a straight line, the more consistent is the rainfall data.

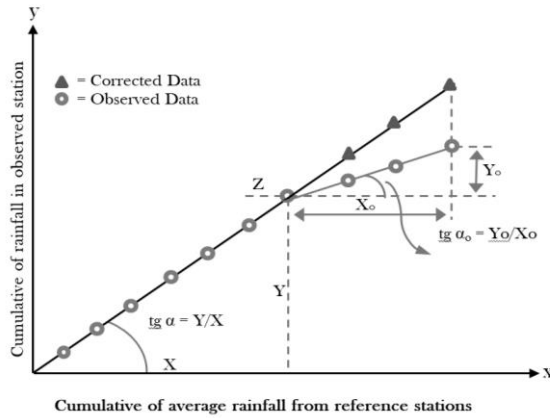


Figure 2. Double mass curve analysis

3.0 RESULTS AND DISCUSSION

3.1 Estimation of Missing Rainfall Data

The four mentioned rainfall stations are adjacent to Gunung Kidul Regency, the southeastern part of Yogyakarta Province. Table 1 shows the distances between these stations as well as their correlations. The rainfall data gained from BPP are presented in Table 2.

Table 1 Distance between stations and its correlation

Station Name	Tanjungsari	Rongkop	Semanu	Tepus
	Distance (km) Correlation			
Tanjungsari	0 0	25 0.89	17 0.85	5.6 0.84
Rongkop	25 0.89	0 0	25 0.82	23 0.83
Semanu	17 0.85	25 0.82	0 0	18 0.78

Table 2. The average monthly rainfall from 2008 – 2019

Station	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Semanu (ISD)	315	261	253	160	106	76	10	6	24	58	219	270
Semanu (AM)	301	245	248	155	90	67	9	5	22	56	219	266
Tepus (ISD)	421	338	272	191	101	138	48	7	65	81	245	369
Tepus (AM)	421	343	269	190	91	111	40	6	63	79	245	361
Tanjungsari (ISD)	367	240	290	166	76	76	21	6	54	68	257	308
Tanjungsari (AM)	367	240	290	166	69	62	19	5	49	66	257	308
Rongkop (ISD)	394	212	271	164	77	60	31	20	55	51	183	273
Rongkop (AM)	394	212	271	164	70	54	28	18	50	51	183	273

Notes: ISD (inversed square distance); AM (arithmetic mean)

Tepus	5.6 0.84	23 0.83	18 0.78	0 0
-------	------------	-----------	-----------	-------

According to the correlation value derived from the original monthly rainfall data (n = 132), the two-tailed significance level of 0.05, and the critical value of 0.171, all correlation values presented in Table 1 are larger than the critical value. It means that there is a significant correlation between the stations. A short distance between two stations does not always yield a strong correlation and vice versa. It indicates that the distance does not directly affect the rainfall data. However, other factors may affect the rainfall intensity, such as hilly topography on karst landforms, type of equipment, etc.

The missing rainfall data spread throughout the year in several rainfall stations were estimated using the inversed square distance, and arithmetic mean methods used two to three closest stations as references. An example of calculation step for Semanu Station in March 2009 with inversed square distance method is presented below:

$$P_{Sem} = \frac{\frac{P_{Tep}}{(d_{Sem-Tep})^2} + \frac{P_{Tan}}{(d_{Sem-Tan})^2} + \frac{P_{Rong}}{(d_{Sem-Rong})^2}}{\frac{1}{(d_{Sem-Tep})^2} + \frac{1}{(d_{Sem-Tan})^2} + \frac{1}{(d_{Sem-Rong})^2}}$$

$$P_{Sem} = \frac{\frac{145}{(18)^2} + \frac{188}{(17)^2} + \frac{172}{(25)^2}}{\frac{1}{(18)^2} + \frac{1}{(17)^2} + \frac{1}{(25)^2}} = 168 \text{ mm}$$

An example of estimating the missing rainfall data with the arithmetic mean method is explained below for Semanu Station in March 2009.

$$P_{Sem} = \frac{P_{Tep} + P_{Tan} + P_{Rong}}{n}$$

$$P_{Sem} = \frac{145 + 188 + 172}{3}$$

$$P_{Sem} = 168 \text{ mm}$$

After calculating all the missing rainfall data using the formulas in Equations 1 and 2, the calculated rainfall data were obtained and used to measure the average monthly rainfall from 2008 to 2019 for the four stations. Table 2 compares the average monthly rainfall data between inversed square distance and arithmetic mean methods for the missing data estimation.

Root Mean Square Error (RMSE) and Mean Absolute Error (MAE)

These two statistical tests were used to measure the performance of ISD and AM to be used as a model to predict the missing rainfall data. A small difference between the observed and predicted rainfall data, RMSE and MAE will display the smallest values. The results of these tests can be seen in Table 3.

Table 3. The RMSE and MAE values for two methods on four stations

Station Name	Methods	RMSE	MAE
Tepus	ISD	106.21	61.35
	AM	99.26	57.74
Tanjungsari	ISD	74.05	44.12
	AM	70.02	44.39
Rongkop	ISD	86.88	55.66
	AM	76.75	47.21
Semanu	ISD	79.68	46.04
	AM	71.99	41.11

It can be seen from Table 3 that the RMSE values gain from AM represent the smallest value compare to those from ISD. This trend is also represented by the MAE value that mostly the MAE values obtained from AM display the smallest value compare to those on ISD. From this calculation, it can be concluded that AM is the most recommended method to be applied at four stations, even though it leaves a quite large error value for rainfall data.

Consistency Checking

The annual data is obtained from monthly rainfall data whose data have been completed using the ISD and AM methods. Based on estimating the missing rainfall data from the previous section, annual rainfall data were measured from the four stations from 2008 to 2019, as tabulated in Table 4.

Data consistency testing is necessary to ensure the validity of data that we had previously calculated and whether the data are reliable for other purposes. The annual rainfall from ISD data in Table 4 was tested for consistency, and the results are presented in Figure 3.

Based on the graph obtained from the double mass curve analysis, one curve shows inconsistency, proven by a breakpoint on the cumulative rainfall graph at Tepus Station (Figure 3d). These graphs mean that the annual rainfall data are consistent in Semanu, Tanjungsari, and Rongkop Station, while inconsistent in Tepus Station. Therefore, this study recommends taking those three consistent rainfall data for any further uses.

Table 5 shows the estimation results of missing rainfall data using the arithmetic mean, presented in annual rainfall information. This information was then used to examine the data consistency of each station to the reference stations.

Figure 4 illustrates the consistency as analyzed by the double mass curve.

All curves portray good data consistency except for the data from Tepus Station, similar to the double mass curve analysis from the previous subsection. It is found that the rainfall data in Tepus Station were inconsistent, which could be seen in Figure 4d so that it is not recommended to be used further.

Based on the estimation of RMSE and MAE on these two methods, the arithmetic mean method shows satisfactory results compared to those using the inversed square distance. The other studies found that the inversed square distance method was good due to the use of distance data between stations in the process [12], [14], [19], [20]. However, in this study, the results show the opposite. It might happen because, in this case, there is no correlation between the two rain stations with the distance between them. This study has a high correlation value of rainfall data at two stations far apart and vice versa. It is possible because other factors can affect rainfall intensity in the study area, including the topography of karst hills, type of equipment, etc. Not all regions are ideal for using the same method in estimating rainfall data. The appropriate method is determined based on each area's characteristics and the availability of supporting reference station data.

Table 4. Annual rainfall data from 2008 to 2019 (inversed square distance)

YEAR	SEMANU	TEPUS	TANJUNGSARI	RONGKOP
2008	1,423	2,319	1,486	1,600
2009	1,848	2,309	1,437	1,063
2010	2,075	2,773	1,804	1,440
2012	1,008	1,353	1,785	1,462
2013	2,173	2,955	1,885	2,028
2014	1,745	2,478	1,574	1,918
2015	1,685	2,081	1,698	1,847
2016	2,684	2,612	2,923	2,901
2017	2,302	2,689	3,188	2,564
2018	1,772	1,858	1,827	1,572
2019	1,557	1,057	1,298	1,061

Table 5. Annual rainfall data from 2008 to 2019 (arithmetic mean method)

Year	SEMANU	TEPUS	TANJUNGSARI	RONGKOP
2008	1,376	2,319	1,486	1,595
2009	1,603	2,309	1,437	1,063
2010	1,758	2,773	1,804	1,440
2012	1,008	1,353	1,785	1,462
2013	2,118	2,926	1,868	2,028
2014	1,745	2,478	1,574	1,918
2015	1,685	2,137	1,698	1,847
2016	2,684	2,522	2,923	2,901
2017	2,302	2,679	3,188	2,564
2018	1,772	1,858	1,827	1,572
2019	1,557	1,057	1,298	1,061

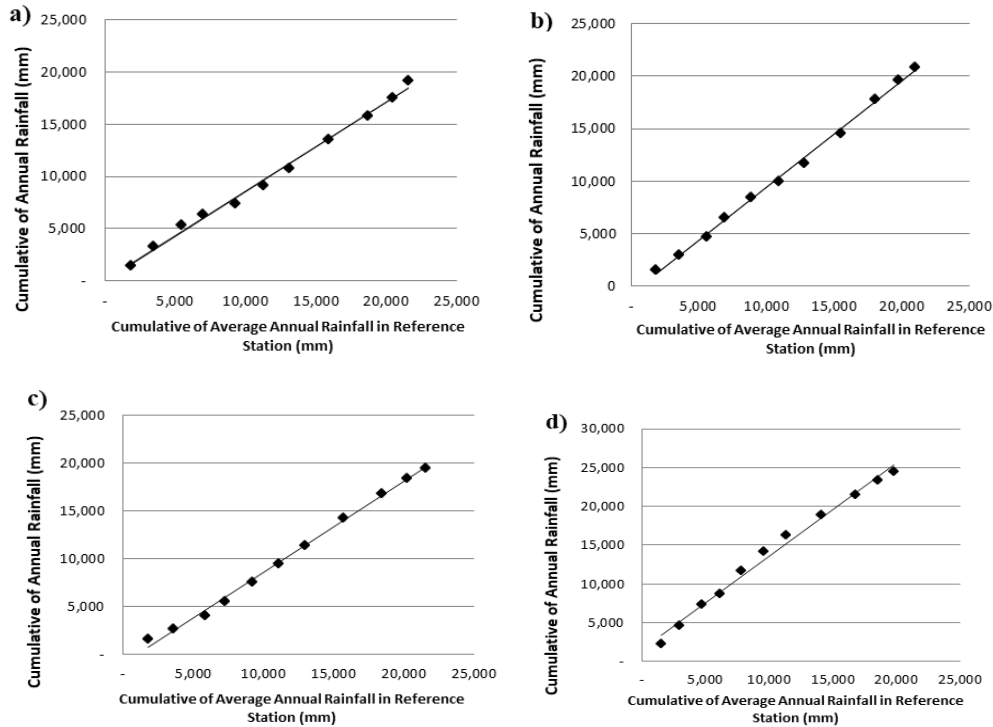


Figure 3. Double mass curve for a) Semanu, b) Tanjungsari, c) Rongkop, and d) Tepus stations based on inversed square distance

The data consistency test from the two methods represented the same pattern that the rainfall data in all stations are good except Tepus Station. The distribution of rainfall data at Tepus Station shows a deflection that comes out of a straight line,

indicating an inconsistency in the rainfall data. Therefore, this study recommends using rainfall data from Semanu, Tanjungsari, and Rongkop Station for further uses.

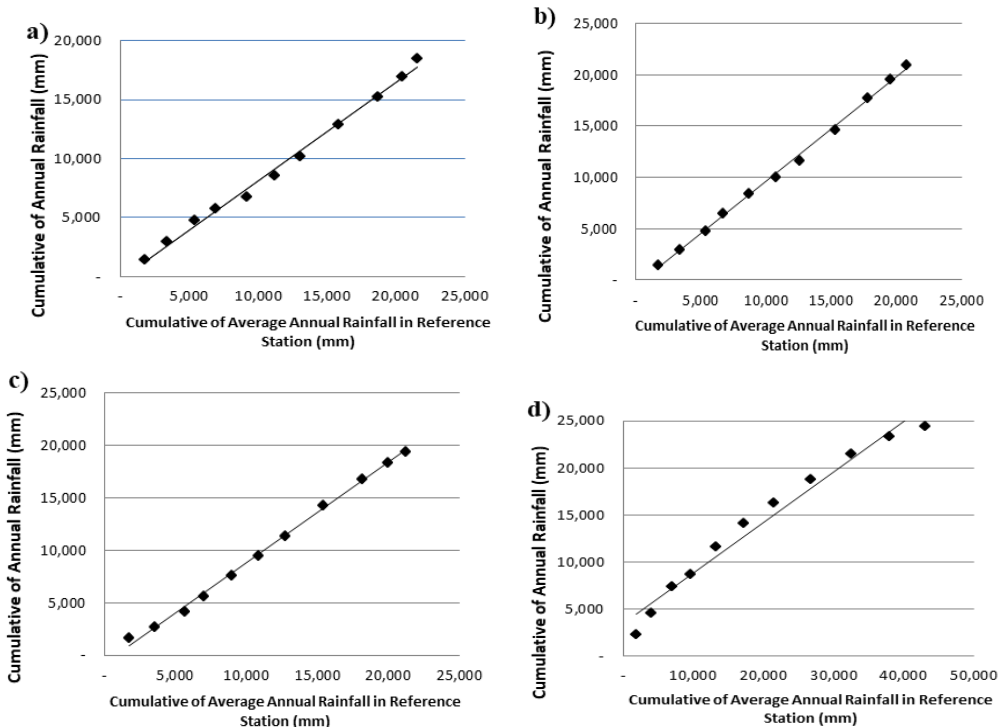


Figure 4. Double mass curve for a) Semanu, b) Tanjungsari, c) Rongkop, and d) Tepus Stations based on Arithmetic Mean Method

4.0 CONCLUSION

The smallest RMSE and MAE values in the arithmetic method at the four stations have confirmed the statement. The rainfall data consistency analyzed by the double mass curve is consistent in all stations except Tepus Station. The rainfall data consistency analyzed by the double mass curve is reliable in all stations except Tepus Station for both methods. This study recommends using rainfall data from Semanu, Tanjungsari, and Rongkop Station for further uses.

Acknowledgement

This research is fully supported by The Indonesia Endowment Funds for Education (LPDP). The authors fully acknowledged Department of Civil and Environmental Engineering, Universitas Gadjah Mada Indonesia for the supports which makes this research well done.

References

- [1] R. D. Kurniawan, R. Hadiani, and Setiono. 2017. Mengisi Data Hujan Yang Hilang dengan Metode Autoregressive dan Metode Reciprocal dengan Pengujian Debit Kala Ulang (Studi Kasus di Das Bakalan). *e-Jurnal MATRIKS Tek. SIPIL*. 5(4): 1315 – 1323. DOI: <https://doi.org/10.20961/mateksi.v5i4.36912>.
- [2] S. N. Z. A. Burhanuddin, S. M. Deni, and N. M. Ramli. 2015, February. Geometric Median for Missing Rainfall Data Imputation. *AIP Conference Proceedings*. 1643: 113–119. DOI: <https://doi.org/10.1063/1.4907433>.
- [3] [A. Bagiawan, S. M. Yuningsih, and D. Windatiningsih. 2011. Pengujian Data Hidrologi dalam Rangka Peningkatan Efektivitas dan Efisiensi Pengelolaan Sumber Daya Air. *Jurnal Sumber Daya Air*. 7(1): 1–17. DOI: <https://doi.org/10.32679/jsda.v7i1.379>.
- [4] M. Hasanpour Kashani and Y. Dinpashoh. 2012. Evaluation of Efficiency of Different Estimation Methods for Missing Climatological Data. *Stochastic Environmental Research and Risk Assessment*. 26(1): 59–71. DOI: <https://doi.org/10.1007/s00477-011-0536-y>.
- [5] N. I. Jaman and S. K. Adhikary. 2020. A Positive Kriging Approach for Missing Rainfall Estimation. *Proceedings of The 5th International Conference on Civil Engineering for Sustainable Development*. 1–10. ISBN: 9789843487643
- [6] R. J. Kuligowski and a. P. Barros. 1998. Using Artificial Neural Networks to Estimate Missing Rainfall Data. *Journal of the American Water Resources Association*. 34(6): 1437–1447. DOI: <https://doi.org/10.1111/j.1752-1688.1998.tb05443.x>.
- [7] R. P. De Silva, N. D. K. Dayawansa, and M. D. Ratnasiri. 2007, May. A Comparison of Methods Used in Estimating Missing Rainfall Data. *Journal of Agricultural Sciences*. 3(2): 101-108. DOI: <https://doi.org/10.4038/jas.v3i2.8107>.
- [8] F. Prawaka, A. Zakaria, and S. Tugiono. 2016, September. Analisis Data Curah Hujan yang Hilang dengan Menggunakan Metode Normal Ratio, Inversed Square Distance, dan Cara Rata-Rata Aljabar (Studi Kasus Curah Hujan Beberapa Stasiun Hujan Daerah Bandar Lampung). *Journal Rekayasa Sipil dan Desain*. 4(3): 397–406. <http://journal.eng.unila.ac.id/index.php/jrsdd/article/view/418/pdf>.
- [9] A. M. Armanuos, N. Al-Ansari, and Z. M. Yaseen. 2020. Cross Assessment of Twenty-One Different Methods for Missing Precipitation Data Estimation. *Atmosphere*. 11(389): 1-34. DOI: <https://doi.org/10.3390/ATMOS11040389>.
- [10] C. Yozgatligil, S. Aslan, C. Iyigun, and I. Batmaz. 2013. Comparison of Missing Value Imputation Methods in Time Series: The Case of Turkish Meteorological Data. *Theoretical and Applied Climatology*. 112: 143–167. DOI: <https://doi.org/10.1007/s00704-012-0723-x>.
- [11] A. S. Yusman. 2018. Aplikasi Metode Normal Ratio dan Inversed Square Distance untuk Melengkapi Data Curah Hujan Kota Padang yang Hilang. *Menara Ilmu*. 12(9): 1–9. DOI: <https://doi.org/10.33559/mi.v12i9.947>.
- [12] I. W. Yasa, M. B. Budianto, and N. M. K. Santi. 2015. Analisis Beberapa Metode Pengisian Data Hujan yang Hilang di Wilayah Sungai Pulau Lombok. *Spektrum Sipil*. 2(1): 49–60. <https://spektrum.unram.ac.id/index.php/Spektrum/article/view/42>.
- [13] N. F. A. Radi, R. Zakaria, S. Z. Satari, and M. A. Z. Azman. 2016, February. Spatial Dependence of Extreme Rainfall. *Proceedings of 3rd ISM International Statistical Conference 2016 (ISM-III)*. 1842: 1–11. DOI: <https://doi.org/10.1063/1.4982833>.
- [14] H. P. G. M. Caldera, V. R. P. C. Piyathisse, and K. D. W. Nandalal. 2016. A Comparison of Methods of Estimating Daily Rainfall Data. *Engineer: Journal of the Institution of Engineers, Sri Lanka*. 49(4): 1-8. DOI : <https://doi.org/10.4038/engineer.v49i4.7232>.
- [15] S. M. C. M. Nor, S. M. Shaharudin, S. Ismail, N. H. Zainuddin, and M. L. Tan. 2020. A Comparative Study of Different Imputation Methods for Daily Rainfall Data in East-Coast Peninsular Malaysia. *Bulletin of Electrical Engineering and Informatics*. 9(2): 635–643. DOI: <https://doi.org/10.11591/eei.v9i2.2090>.
- [16] W. Sanusi, W. Z. W. Zin, U. Mulbar, M. Danial, and S. Side. 2017. Comparison of The Methods to Estimate Missing Values in Monthly Precipitation Data. *International Journal on Advanced Science, Engineering and Information Technology*. 7(6): 2168–2174. DOI: <https://doi.org/10.18517/ijaseit.7.6.2637>.
- [17] F. Jahan, N. C. Sinha, M. M. Rahman, M. M. Rahman, M. S. H. Mondal, and M. A. Islam. 2019. Comparison of Missing Value Estimation Techniques in Rainfall Data of Bangladesh. *Theoretical and Applied Climatology*. 136: 1115–1131. DOI: <https://doi.org/10.1007/s00704-018-2537-y>.
- [18] N. Kanda, H. S. Negi, M. S. Rishi, and M. . Shekhar. 2018. Performance of Various Techniques in Estimating Missing Climatological Data Over Snowbound Mountainous Areas of Karakoram Himalaya. *Meteorological Applications*. 25: 337 – 349. DOI: <https://doi.org/10.1002/met.1699>.
- [19] J. Suhaila, M. D. Sayang, and A. A. Jemain. 2008. Revised Spatial Weighting Methods for Estimation of Missing Rainfall Data. *Asia-Pacific Journal of Atmospheric Sciences*. 44(2): 93–104. ISSN: 19767633
- [20] H. Rizky, Y. N. Nasution, and R. Goejantoro. 2019. Analisis Data Curah Hujan yang Hilang Menggunakan Metode Inversed Square Distance. *Proceedings of Seminar Nasional Matematika, Statistika, dan Aplikasinya*. 1: 138–149. e-ISSN : 2657-232X