

ESTIMATION OF MISSING STREAMFLOW DATA USING ANFIS MODELS AND DETERMINATION OF THE NUMBER OF DATASETS FOR ANFIS: THE CASE OF YEŞİLIRMAK RIVER

SAPLIOGLU, K.* – KUCUKERDEM, T. S.

*Department of Civil Engineering, Faculty of Engineering, Süleyman Demirel University
Isparta, Turkey
(phone: +90-246-211-12-13)*

**Corresponding author
e-mail: kemalsaplioglu@sdu.edu.tr*

(Received 22nd Mar 2018; accepted 25th May 2018)

Abstract. Good data analysis is required for the optimal design of water resources projects. However, data are not regularly collected due to material or technical reasons, which results in incomplete-data problems. Available data and data length are of great importance to solve those problems. Various studies have been conducted on missing data treatment. This study used data from the flow observation stations on Yeşilirmak River in Turkey. In the first part of the study, models were generated and compared in order to complete missing data using Artificial Neural Network Fuzzy Inference Systems (ANFIS), multiple regression and Normal Ratio Method. Thus, it is tried to define the usability besides the other model to complete the missing data. Likewise, in the study. It is aimed to define the minimum number of data necessary for the use age of ANFIS. For this purpose in the second part of the study, the minimum number of data required for ANFIS models was determined using the optimum ANFIS model. Of all methods compared in this study, ANFIS models yielded the most accurate results. A 10-year training set was also found to be sufficient as a data set.

Keywords: *Anfis, missing data, multiple regression, normal ratio method, Yeşilirmak*

Introduction

Both the growing population and the rapidly developing industrialization lead to an increased demand for water. The limited availability of resources results in a number of problems in meeting the demand. Exploitation of unused water resources or using existing water resources in an optimum way can be a solution to these problems. Optimal utilization of available water resources, in particular, requires a good analysis of data. Due to the small number of stations in project areas or insufficient data length, some studies have been undertaken to generate new data using existing measurement stations (Aissia et al., 2017). These hydrological studies have mainly focused on precipitation (Sun et al., 2017), evaporation (Güçlü et al., 2017) and river flows (Shiau and Hsu, 2016).

Studies on missing data treatment generally address data correlation (Bakis and Göncü, 2015; Britoa et al., 2017), back-propagation (BP) neural network using Artificial Intelligence (Gao and Wang, 2017), ANFIS models (Dastorani et al., 2010; Mpallas et al., 2010) and models using artificial neural networks (ANN) (Kim and Pachepsky, 2010; Kim et al., 2015). In addition, Fuzzy studies (Coulibaly and Evora, 2007), in which modeling is based on pure expert knowledge, are also important. Some studies on missing data treatment using ANFIS are the completion of missing flow data of the Middle Euphrates basin (Yilmaz and Muttil, 2014), completion of missing precipitation data in Serbia (Petkovic et al., 2016) and Malaysia (Nawaz et al., 2016),

and completion of missing flow data and modelling of sediment transport of Terengganu River, Malaysia (Ismail and Zin, 2017) and Gediz River, Turkey (Ulke et al., 2009).

This study investigated the monthly data of the stations of Yeşilirmak River in the North of Turkey. In this study, Normal Ratio Method which is frequently used in literature and of which usability is proven and the Regression methods enable to reach the results by the relationship between the data are chosen to be able to compare the accuracy of offered ANFIS model. In the first part of the study, multiple regression tests based on interstation correlations were performed. In the second part of the study, an optimum data completion model was selected using ANFIS. In the last part of the study, the number of data required for a correct prediction was searched and the minimum number of data required for reliable estimates was discussed.

Material and Method

The first part of this section of the study will present information and statistics on Yeşilirmak River and its stations. The second part will provide information on the classical method, multiple regression method and ANFIS used in the study.

Yeşilirmak River and stations

The Yeşilirmak basin, one of the 25 basins in Turkey, is located between latitudes $39^{\circ} 30'$ and $41^{\circ} 21'$ and longitudes $34^{\circ} 40'$ and $39^{\circ} 48'$ (Figure 1). The basin is named after Yeşilirmak River. The main river channel of the basin is 519 km in length. The main tributaries of Yeşilirmak River are Kelkit, Çekerek, Çorum, Çat and Tersakan streams. Estimated to be about 3,8 million ha, Yeşilirmak basin is the third largest basin in Turkey (Kurunç et al., 2005).

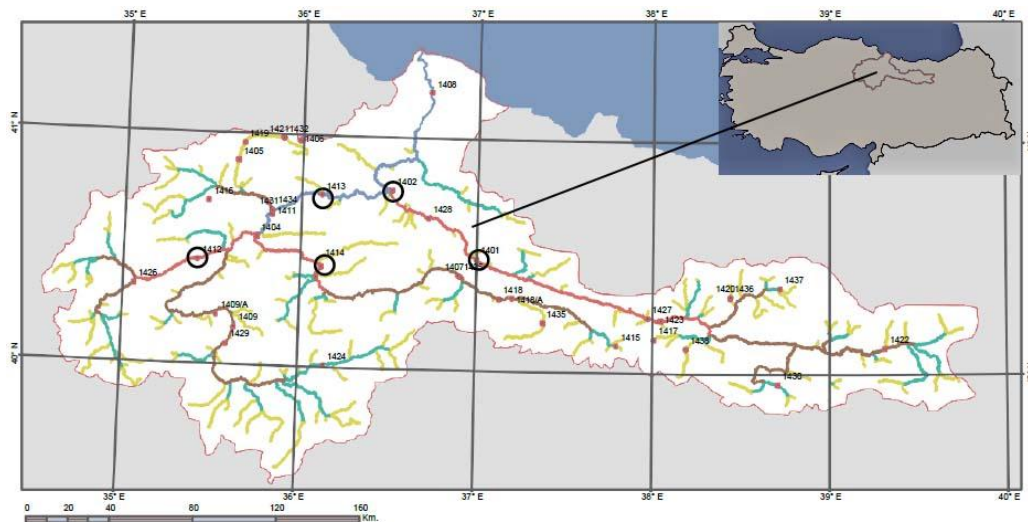


Figure 1. Site location map of Yeşilirmak Basin

Stations No 1401, 1402, 1412, 1413 and 1414 of Yeşilirmak River were used in the study. *Table 1* shows the statistics of the stations. *Table 2* summarizes the correlation between the stations.

Table 1. Statistical analysis of data from stations

	1401	1402	1412	1413	1414
Y- coordinate	40°28'42''	40°46'18''	40°27'06''	40°44'40''	40°26'03''
X- coordinate	36°59'56''	36°30'45''	35°25'03''	36°06'43''	36°07'05''
Precipitation Area(km²)	10048.8	33904.0	3668.8	21667.2	5409.2
Altitude	375	190	530	301	510
Mean	71.96	150.01	7.35	63.88	25.68
Standard Error	3.07	5.17	0.37	2.34	0.84
Median	35.45	102.50	4.21	45.30	19.65
Standard Deviation	83.65	136.11	8.65	54.95	19.71
Kurtosis	4.27	2.81	6.10	3.41	4.79
Skewness	2.06	1.66	2.24	1.70	1.81
Max	5.23	13.50	0.02	2.47	2.46
Min	548.00	791.00	59.10	350.00	139.00
Number of Data	744	692	536	552	546
Confidence Interval (95,0%)	6.02	10.16	0.73	4.59	1.66

Table 2. Correlation between data from stations

Stations	1401	1402	1412	1413	1414
1401	1				
1402	0.909	1			
1412	0.516	0.784	1		
1413	0.702	0.926	0.885	1	
1414	0.539	0.565	0.432	0.518	1

Methods

Missing data treatment using Normal Ratio Method

In this method, each input data is divided by its annual average value, and these values are multiplied by the average of the station (average of data) whose missing data are to be completed. All input values obtained in the last stage of the calculation are summed, and divided by the number of inputs so that the missing data are completed (Ismail and Zin, 2017).

$$Q_e = (Q_1 * \frac{Q_e(ort)}{Q_1(ort)} + Q_2 * \frac{Q_e(ort)}{Q_2(ort)} + \dots + Q_n * \frac{Q_e(ort)}{Q_n(ort)})/n \quad (\text{Eq.1})$$

where Q is the flow rate and n is the number of input stations.

Multiple regression analysis

Multiple regression analysis is a statistical method for determining the mathematical dimension of the relationship between variables affecting each other. The value to be estimated using the equation formulated based on multiple regression analysis is written in the form of a function of values affecting it (Sun and Trover, 2018).

$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + U \quad (\text{Eq.2})$$

where Y is the dependent (estimated) variable, X is the independent (explanatory) variable, β is the regression coefficient, k is the number of input parameters and U is the error term.

Multiple linear regression analysis can be used when data are normally distributed, the relationship between independent variables and dependent variable is linear, and error variance for each independent variable is constant (Hair et al., 2009).

ANFIS

Developed by Jang (1993), ANFIS is a modeling method that combines Fuzzy Logic and YSA models. Different from Fuzzy Logic, ANFIS is based on the use of data for the automatic acquisition of rules. ANFIS structure uses artificial neural networks' learning ability and fuzzy logic inference, and therefore, it is more successful than when artificial neural networks model or fuzzy logic is used alone. When input and output values are known, ANFIS determines all possible rules or allows them to be generated using input and output values (Figure 2). ANFIS structure consists of five layers: fuzzification layer, rule layer, normalization layer, defuzzification layer and summation layer (Figure 3). The first and fourth layers are adaptable (Nayak et al., 2004; Jang, 1993). ANFIS models can be especially used in the conditions if there are lots data and if they are organized.

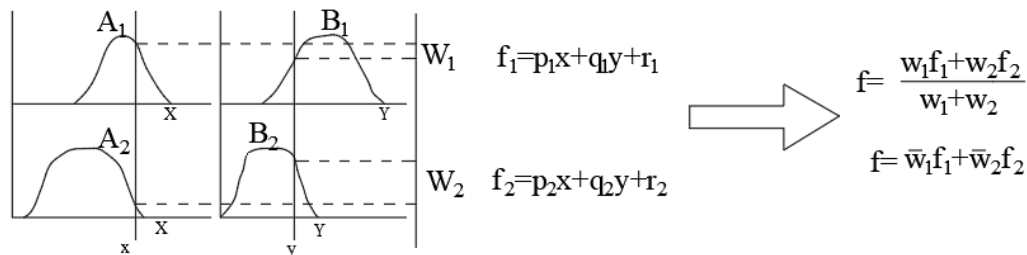


Figure 2. Fuzzy Inference System

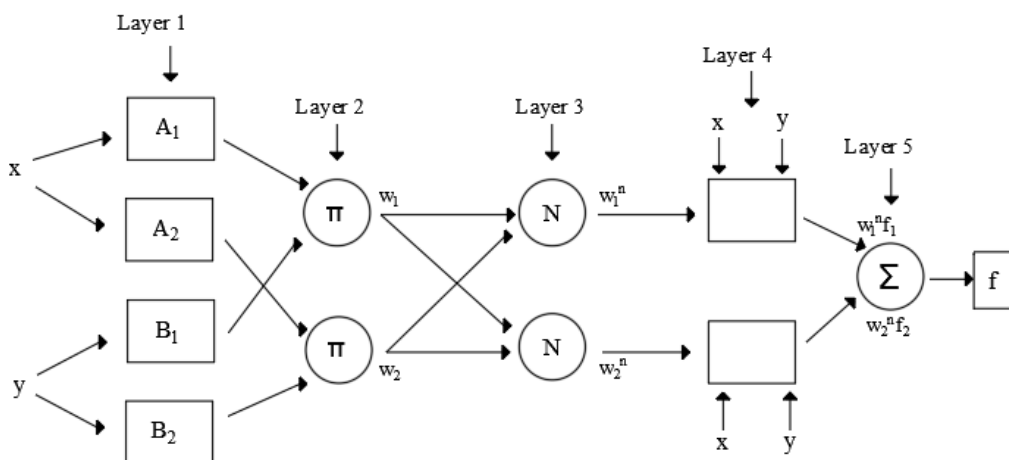


Figure 3. ANFIS architecture

Results

Models were developed using Yeşilirmak River data for the estimation of missing data of stations 1402 and 1413. Two-input-one-output models and four-input-one-output models were developed to complete the missing data of station 1402. The two-input-one-output models were used as the output of station 1402. Stations 1413 and 1401 connected to station 1402 on the left- and right-hand sides, respectively, were used to estimate station 1402. In addition to these stations, stations 1412 and 1414 connected to station 1413 on the left- and right-hand sides, respectively, were used to estimate station 1413 in the four-input-one-output models. A two-input-one-output model was developed, and stations 1414 and 1412 were used for the estimation of missing data of station 1413. In the models developed for stations 1402 and 1413, classical and multiple regression models were constructed and compared as well as the ANFIS method. In the last part of the study, the minimum number of data required to reach the correct result using ANFIS models was obtained.

First data set models

The aim of these models was to complete the missing data of station 1402. For this, the data of stations 1413 and 1401 were used. Of 540 data, the first 400 were used for training and the remaining for testing. In addition to ANFIS models generated by changing the number of sets of input parameters, classical method and the multiple regression model were used to compare the results (*Table 3*).

The equation of the classical method is:

$$X_{1402} = (2,072 * X_{1401} + 2,215 * X_{1413})/2 \quad (\text{Eq.3})$$

The equation of the multiple regression model is:

$$X_{1402} = 0,896 * X_{1401} + 1,178 * X_{1413} + 5,37 \quad (\text{Eq.4})$$

Table 3. Training and testing data results of models developed for the first data set

Models	Training Data		Testing Data		
	R ²	Mean Squared Error %	R ²	Mean Squared Error %	
ANFIS Models	3-3	0.976	11.73	0.977	17.36
	4-4	0.977	11.66	0.978	17.38
	5-5	0.983	10.99	0.979	17.55
	6-6	0.980	11.13	0.978	16.60
	7-7	0.979	11.12	0.961	17.00
	8-8	0.983	11.09	0.959	17.12
Normal Ratio Method	0,970	11.79	0.955	18.85	
Multiple Regression	0,972	13.54	0.960	18.62	

The results of this part of the study show that ANFIS models are not superior to the classical and multiple regression models but that all ANFIS models yield better results than the other two methods. The models in which each input has 5 subsets are the optimum models. The speed of the training phase of the model is also noteworthy. The

comparison of the values obtained from the optimum ANFIS model with the observed values shows that the errors of both the minimum and maximum flow values are very few (*Figure 4*).

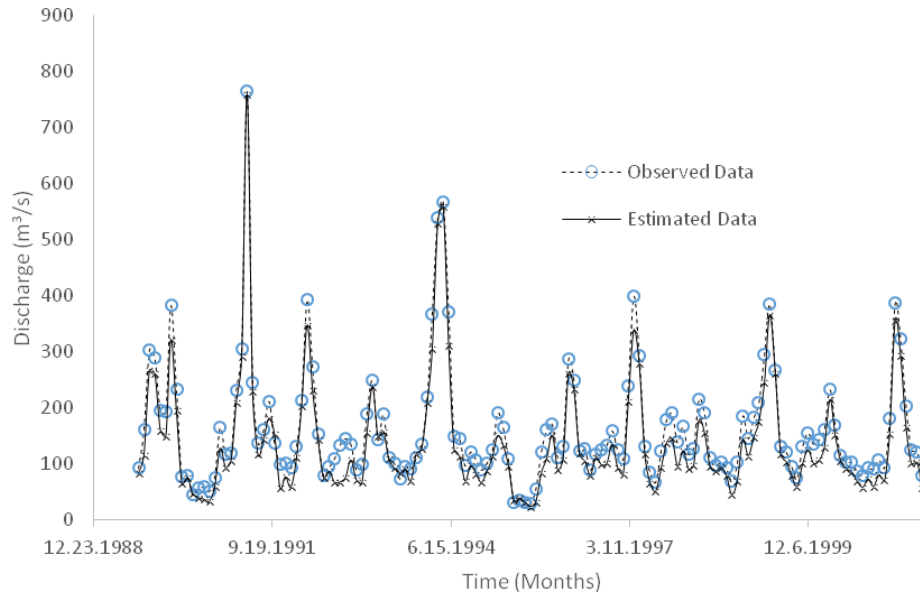


Figure 4. Comparison of observed data and estimated data of Anfis (5-5) model test data for the first data set

Second data set models

These models also aimed to complete the missing data of station 1402. To achieve this, the data of stations 1412 and 1414 as well as those of 1413 and 1401 were used. Of 504 data, the first 405 were used for training and the remaining for testing. In addition to ANFIS models generated by changing the number of sets of input parameters, classical method and multiple regression model were used to compare the results (*Table 4*).

The equation of the classical method is:

$$X_{1402} = (2,107 * X_{1401} + 19,53 * X_{1412} + 2,26 * X_{1413} + 5,44 * X_{1414})/4 \quad (\text{Eq.5})$$

The equation of the multiple regression model is:

$$X_{1402} = 0,921 * X_{1401} + 1,036 * X_{1412} + 1,090 * X_{1413} - 0,251 * X_{1414} + 11,071 (\text{Eq.6})$$

The results show that ANFIS models provide more accurate results than the classical model and worse results than multiple regression models. The models with 4 inputs are quite slow, especially when the number of subsets of inputs is greater than 5. The models with an increasing number of inputs are much slower than multiple regression models. The comparison of the values of the optimum ANFIS model with the observed values shows that although the largest error is at the minimum flow values, this error is quite small at the maximum flow values (*Figure 5*).

Table 4. Training and testing data results of models developed for the second data set

Models	Training Data		Testing Data		
	R ²	Mean Squared Error %	R ²	Mean Squared Error %	
ANFIS Models	3-3-3-3	0.933	28.96	0.916	26.92
	4-4-4-4	0.986	11.27	0.919	16.94
	5-5-5-5	0.991	10.30	0.918	18.14
	6-6-6-6	0.991	9.85	0.919	19.09
Normal Ratio Method	0.873	33.79	0.863	28.61	
Multiple Regression	0.976	18.26	0.975	16.73	

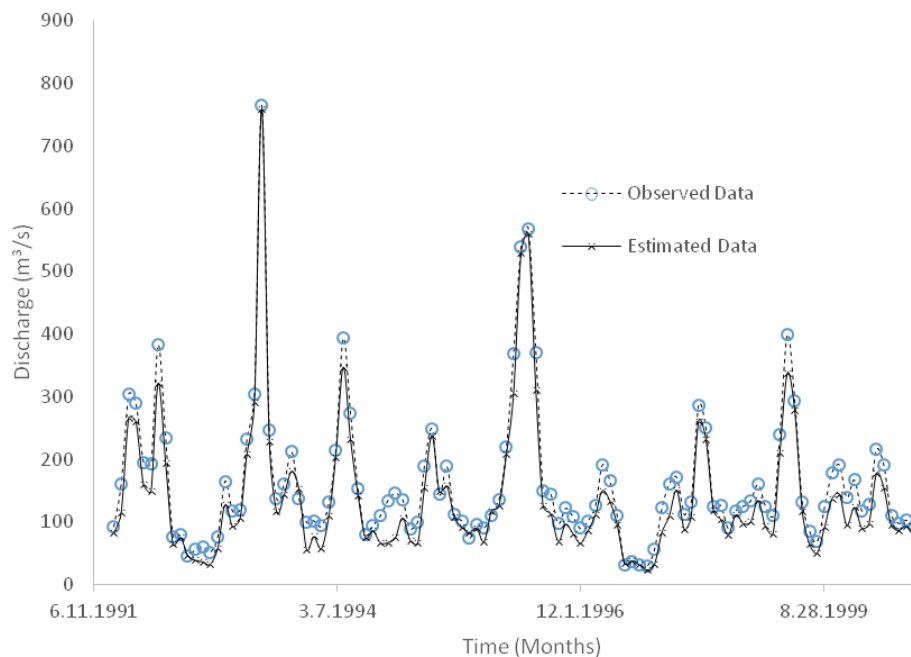


Figure 5. Comparison of observed data and estimated data of Anfis (5-5) model test data for the second data set

Third data set models

The aim of these models was to complete the missing data of station 1413. For this, the data of stations 1412 and 1414 were used. Of 504 data, the first 405 were used for training and the remaining for testing. The results were compared using the classical method and multiple regression model as well as ANFIS models generated by changing the number of sets of input parameters (Table 5).

The equation of the classical method is:

$$X_{1413} = (2,407 * X_{1414} + 8,641 * X_{1412})/2 \quad (\text{Eq.7})$$

The equation of the multiple regression model is:

$$X_{1413} = 0,404 * X_{1414} + 4,904 * X_{1412} + 17,576 \quad (\text{Eq.8})$$

Table 5. Training and testing data results of models developed for the third data set

Models		Training Data		Testing Data	
		R ²	Mean Squared Error %	R ²	Mean Squared Error %
ANFIS Models	3-3	0.890	20.29	0.865	23.24
	4-4	0.892	19.35	0.872	22.93
	5-5	0.881	18.64	0.879	22.56
	6-6	0.865	19.88	0.854	23.07
Normal Ratio Method		0.764	27.41	0.694	36.60
Multiple Regression		0.905	34.26	0.757	46.74

The results show that ANFIS models provide more accurate results than the classical and multiple regression models. Although the results are not as good as those in the first data set, they remain within acceptable error limits. The models in which each input has 5 subsets are the optimum models. The comparison of the values of the optimum ANFIS model with the observed values shows that the errors of both the minimum and maximum flow values are very few (*Figure 6*).

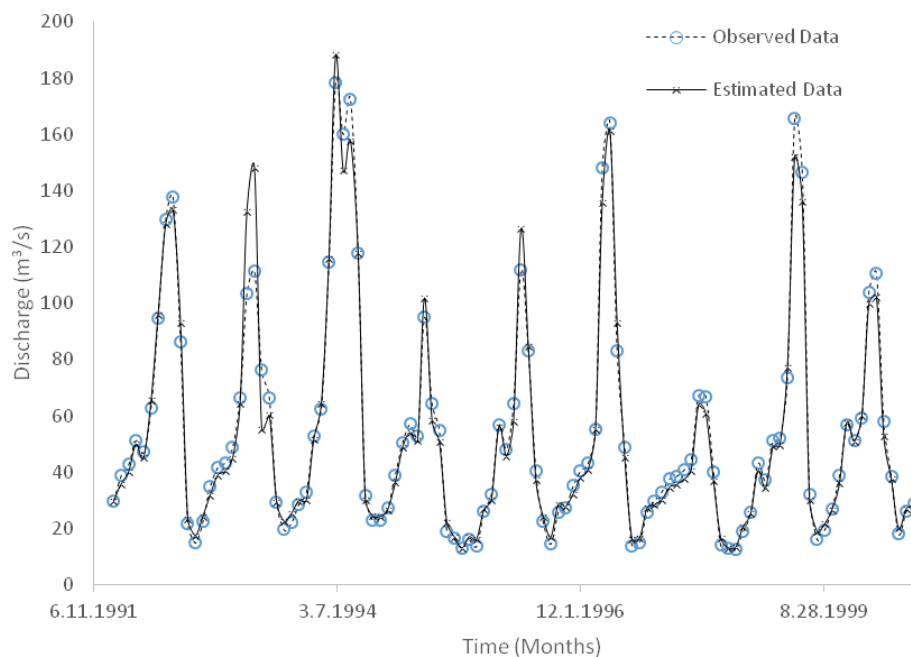


Figure 6. Comparison of observed data and estimated data of Anfis (5-5) model test data for the third data set

Determining the Minimum Number of Data for Anfis Model Training

This part of the study attempted to obtain the minimum number of data required for a reliable ANFIS model training. The model with a 5-5 set of the first data set (the optimum modelling) were the model of choice for this purpose. The models were trained using a 10-year data set and the procedure was repeated year by year. The evaluation of the results is summarized in *Table 6*.

Table 6. Regression and error values for the number of data used for ANFIS model training

Number of Data	Training Data		Testing Data	
	R ²	Mean Squared Error %	R ²	Mean Squared Error %
1- year	0,999	0,74	0,036	169,27
2- year	0,993	13,42	0,071	118,02
3-year	0,992	16,79	0,139	111,06
4- year	0,985	21,38	0,198	36,91
5- year	0,980	22,39	0,628	30,44
6- year	0,981	21,23	0,738	21,92
7- year	0,982	19,77	0,803	20,12
8- year	0,983	18,76	0,844	17,77
9- year	0,984	16,88	0,898	17,74
10- year	0,985	12,27	0,967	17,65
33- year	0,983	10,99	0,979	17,55

The number of data used for ANFIS model training does not affect the regression coefficient of the training data very much (*Table 6*). However, the regression results obtained during the testing of the models show that the results of the 10-year data set and model are very similar to those of the 33-year data and training (*Figure 7*). The error values show that the 8-year data set is sufficient for training (*Figure 8*). In conclusion, a 10-year data set may be sufficient for Anfis model training.

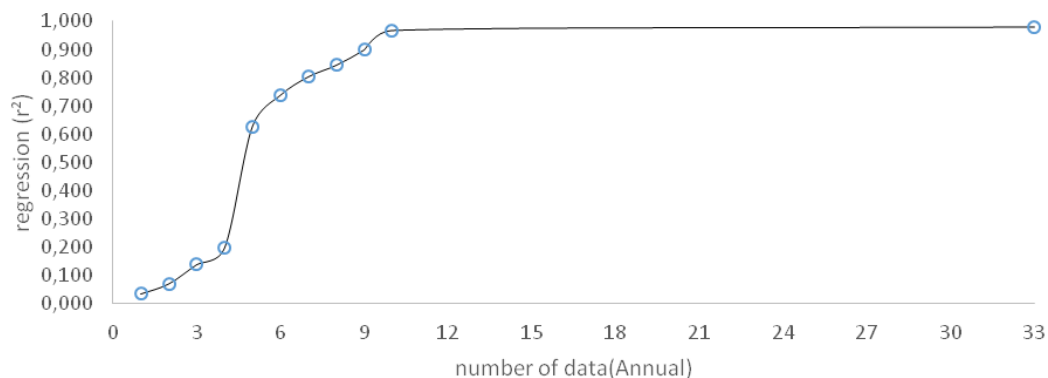


Figure 7. Regression values for the number of data used for ANFIS model training

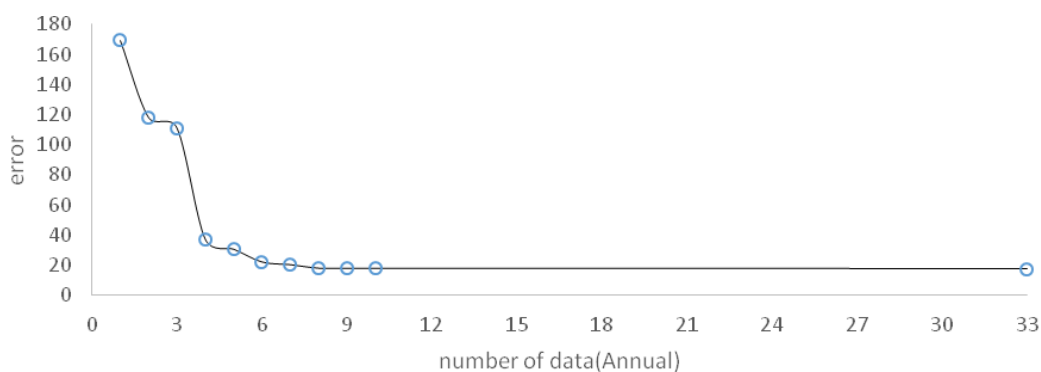


Figure 8. Error values for the number of data used for ANFIS model training

Discussion

The study was repeated with 3 different data sets and ANFIS model behaviors is searched in different models. When the obtained results are analyzed it is seen that the increase in the number of the stations affects ANFIS models adversely; however, this effect isn't so big as Normal Ratio Method has. In contrast to this, it can be said that the increase in the number of the stations affect Multiple Regression Models in a positive way. While the proportion of error is more in the minimum discharge value, it is much smaller in the maximum discharge value. The changing of error in the maximum values isn't effected so much by the changings in the number of stations. One of the biggest lacknesses of ANFIS models is memorizing during the process of the training. To be able to overcome this lackness, some of the data should be allocated for testing. The accuracy rate which has %99 accuracy proportion during the process of training reaches to much lower degrees for testing date. This is a proof for the model that this model is composed via memorizing instead of learning and it can't be used. The closeness between the results obtained out of training and the test results shows that ANFIS is appropriate to complete missed data.

In the second part of the study, it is focused on the number of data which is necessary for being able to use ANFIS models. Especially in the situation of having fewer data, the needed number of data is tried to determine to be able to get reliable results. In the study done for this, it is noticed that there is a difference less than % 1 between the models composed with a decade data and 8 composed with 33 years data. When error test data is inspected, it can be said that 8 years data can be used safely. It is defined that the data is less and 3 years data can't be used. Annual data in the model is the biggest proof for showing to trust the training results. Although the training results are extremely trustable, the test results are weak.

In the situation of increasing the number of data, the results obtained via ANFIS move away the accuracy and so the time or elevation the results becomes longer. Additionally, every input parameter must be balanced in the number of membership functions. Number of membership functions is important for the accuracy of the model. Having many membership functions may affect the accuracy of the model in the negative way. The operation time of the models increases exponentially at the same time. For example while in a model having 2 inputs and having 5 sets of membership functions for each input, 25 rules composes; in a model having 4 inputs and 10 sets of membership functions for each 10^4 rules are needed to be composed. However, this affects the model's operation performance adversely.

Conclusion

Sometimes it is not possible to collect long and coordinated data in order to optimally use water resources projects. The aim of this study was to develop ANFIS models for stations on Yeşilirmak River in order to solve this problem and improve the existing methods. Another aim of the study was to investigate how the number of input parameters and amount of data affect ANFIS models. For this purpose, 3 different data sets were analyzed.

Results show that besides classical and regression models, ANFIS models can be used to complete missing flow data. ANFIS models yield very accurate results especially when the number of input parameters is small. However, multiple regression models yield better results than ANFIS models when the number of input parameters is

large. In addition, it takes ANFIS models longer to achieve results when the number of input parameters increases. Lastly, at least a 10-year data is required for a reliable ANFIS model training phase.

In conclusion, the ANFIS modelling yield accurate results and therefore can be used to complete missing data when the number of input parameters is small and data set is older than 10 years.

ANFIS models must be advocated by using the data of past and they can be used to complete the missing data which is going to happen in the future after proving their usage by looking at the test results. In addition, via this, newly opened stations past data can be supplied. In addition to completing the data of streamflow, it is thought that it can be used to complete the data having great importance in the projects of water sources. Such as rainfall, evaporation, water quality and sediment.

REFERENCES

- [1] Aissia, M., Chebana, F., Quarda, T. (2017): Multivariate missing data in hydrology – Review and applications. – *Advances in Water Resources* 110: 299-309.
- [2] Bakıs, R., Goncu, S. (2015): Completion Of Missing Data In Rivers Flow Measurement: Case Study Of Zab River Basin. – *Anadolu University Journal of Science and Technology / A Applied Sciences and Engineering* 16(1): 63-79.
- [3] Britoa, R., Almeida, M., Matos, J. (2017): Estimating flow data in urban drainage using partial least squares regression. – *Urban Water Journal* 14(5): 467-474.
- [4] Coulibaly, P., Evora, N. (2007): Comparison of neural network methods for infilling missing daily weather records. – *Journal of Hydrology* 341: 27-41.
- [5] Dastorani, M., Moghadamnia, A., Piri, J., Ramirez, M. (2010): Application of ANN and ANFIS models for reconstructing missing flow data. – *Environ Monit Assess* 166: 421-434.
- [6] Gao, T., Wang, H. (2017): Testing Backpropagation Neural Network Approach in Interpolating Missing Daily Precipitation. – *Water Air Soil Pollut* 228(404): 2-17.
- [7] Güçlü, Y., Subyani, A., Şen, Z. (2017): Regional fuzzy chain model for evapotranspiration estimation. – *Journal of Hydrology* 544: 233-241.
- [8] Hair, J. F., Black, W. C., Babin, B. J., Anderson, R. E. (2009): *Multivariate Data Analysis*. – Pearson.
- [9] Ismail, W., Zin, W. (2017): Estimation of rainfall and stream flow missing data for Terengganu, Malaysia by using interpolation technique methods. – *Malaysian Journal of Fundamental and applied Sciences* 13(3): 214-218.
- [10] Jang, J. S. R. (1993): ANFIS: Adaptive-Network-Based Fuzzy Inference System. – *IEEE Transactions on Systems, Man, and Cybernetics*.
- [11] Kim, J., Pachepsky, Y. (2010): Reconstructing missing daily precipitation data using regression trees and artificial neural networks for SWAT streamflow simulation. – *Journal of Hydrology* 394: 305-314.
- [12] Kim, M., Beak, S., Ligaray, M., Pyo, J., Park, M., Cho, K. (2015): Comparative Studies of Different Imputation Methods for Recovering Streamflow Observation. – *Water* 7: 6847-6860.
- [13] Kurunç, A., Yürekli, K., Öztürk, F. (2005): Effect of Discharge Fluctuation on Water Quality Variables from the Yeşilirmak River. – *Tarım Bilimleri Dergisi*, 189-195.
- [14] Mpallas, L., Tzimopoulos, C., Evangelidis, C. (2010): Rainfall data calculation using Artificial Neural Networks and adaptive neuro-fuzzy inference systems. – *Sustainable Irrigation Management, Technologies and Policies* 134: 133-144.

- [15] Nawaz, N., Harun, S., Othman, R., Heryansyah, A. (2016): Neuro-Fuzzy Systems Approach To Infill Missing Rainfall Data For Klang River Catchment, Malaysia. – *Jurnal Teknologi* 78(6): 15-21.
- [16] Nayak, P. C., Sudheer, K. P., Rangan, D. M., Ramasastri, K. S. (2004): A neuro-fuzzy computing technique for modeling hydrological time series. – *Journal of Hydrology*: 52-66.
- [17] Petkovic, D., Gocic, M., Shamshirband, S. (2016): Adaptive Neuro-Fuzzy Computing Technique For Precipitation Estimation. – *Facta Universitatis-Series Mechanical Engineering* 14(2): 209-218.
- [18] Shiau, J., Hsu, H. (2016): Suitability of ANN-Based Daily Streamflow Extension Models: a Case Study of Gaoping River Basin, Taiwan. – *Water Resour. Manage.* 30: 1499-1513.
- [19] Sun, S., Leonhardt, G., Sandoval, S., Krajewski, J., Rauch, W. (2017): A Bayesian method for missing rainfall estimation using a conceptual rainfall–runoff model. – *Hydrological Sciences Journal* 62(15): 2456–2468.
- [20] Sun, W., Trover, B. (2018): Multiple model combination methods for annual maximum water level prediction during river ice breakup. – *Hydrological Processes* 32(3): 421-435.
- [21] Ulke, A., Tayfur, G., Ozkul, S. (2009): Predicting Suspended Sediment Loads and Missing Data for Gediz River, Turkey. – *Journal of Hydrologic Engineering* 14(9): 954-965.
- [22] Yilmaz, A., Muttil, N. (2014): Runoff Estimation by Machine Learning Methods and Application to the Euphrates Basin in Turkey. – *Journal of Hydrologic Engineering* 19(5): 1015-1025.