

## Prediction of daily failure rate using the serial triple diagram model and artificial neural network

Burak Kizilöz 

Kocaeli Metropolitan Municipality Water and Sewerage Administration, İzmit, Kocaeli, Turkey  
E-mail: bkiziloz@isu.gov.tr

 BK, 0000-0001-5243-8889

### ABSTRACT

In this study, 41 models used for the prediction of daily failure rates in water distribution networks have been designed via the Serial Triple Diagram Model (STDM) and artificial neural network (ANN) methods. For this purpose, daily failure data measured coordinately in the water distribution system network in Gebze and transferred to the geographic information system (GIS) has been used. The data has been normalized through the min-max technique to scale it at regular intervals and develop the model prediction performance. In this study, certain meteorological variables such as temperature and precipitation have been taken into account as model input for the first time. According to the increasing values of these two variables, it is observed in the model results that daily failure rates tend to increase. The expected model accuracies in failure rate prediction could not be obtained through the suggested ANN models. The higher prediction performances have been obtained through the STDM, a structure that enables visualization of the model results by making inferences. The STDM method is a significant alternative approach to determine the relationship of the variables on the failure and to predict the failure rate. It is predicted that the suggested STDM charts will contribute to the decision-makers, experts, and planners to determine effective infrastructure management. Also, investment planning prioritization will be able to reduce failure rates by interpreting prediction charts.

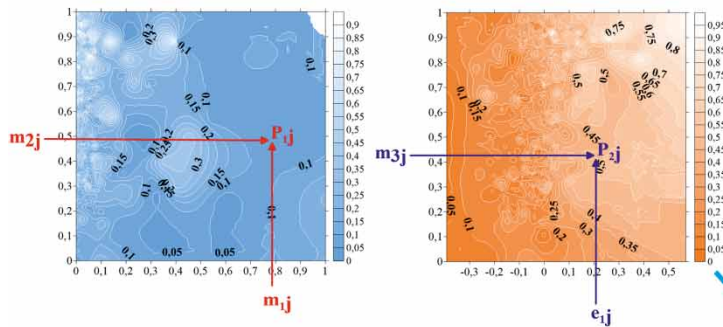
**Key words:** artificial neural network, failure rate, Kriging, serial triple diagram model, water distribution network

### HIGHLIGHTS

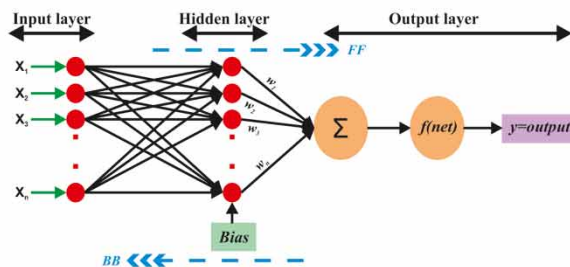
- Counter maps have been used for predicting and interpreting failure rates.
- Temperature and precipitation effect has been directly taken into consideration for the first time.
- The highest accuracy has been achieved with the least input in models.

## GRAPHICAL ABSTRACT

## STDM model structure

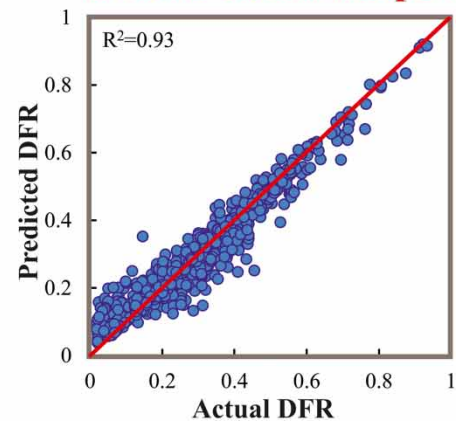


AND



## ANN model structure

## STDM model output



## LIST OF ABBREVIATIONS

The following symbols are used in this paper:

Acronym	Definition
AC	Asbestos-cement
ANN	Artificial neural network
APP	Average pressure of pipe
AWWA	American Water Works Association
BI	Bias
CI	Cast iron
DFR	Daily failure rate
DI	Ductile iron
DMA	District metered area
DMT	Daily mean temperature
DTR	Daily total rainfall
FFBP	Feed forward back propagation
FR	Failure rate
GIS	Geographic information system
GRP	Glass reinforced plastic
HDPE	High-density polyethylene
IWA	International Water Association
LM	Levenberg-Marquardt algorithm
LS-SVM	Least-squares supportive machine
MAE	Mean absolute error
MAP	Mean age of pipe
MARS	Multivariate adaptive regression splines

MDP	Mean diameter of the pipe
MGM	Turkish Meteorological Service
MLP	Mean length of pipe
MSE	Mean square error
NRWR	Non-revenue water ratio
PMS	Pressure management system
PRV	Pressure reducing valve
PVC	Polyvinyl chloride
R <sup>2</sup>	Coefficient of determination
RF	Random forest
SCADA	Supervisory control and data acquisition system
SI	Scatter index
ST	Steel
STD	Serial triple diagram model
SV	Semi-variogram
SVM	Support vector machines
TDM	Triple diagram model
WDN	Water distribution network

## INTRODUCTION

Water utilities are responsible for delivering quality purified drinking water to consumers without any interruption through water distribution networks (WDNs). Failures that may occur in WDNs and their components (such as valves, fire hydrants, pumps, service connection, etc.) bring significant problems for the water utility as they cause a decrease in network pressures, deterioration of water quality, increase in consumer dissatisfaction, and excessive water loss. For this reason, rapid detection and appropriate repair are significant in terms of reducing these problems. Water utilities, domain experts, and decision-makers need to review their current management and operating approaches on WDNs due to environmental, economic, and social impacts of failures. So it is important to determine the rates of failures, their frequencies and densities, and the factors causing the failures (Aydogdu & Fırat 2015). In the literature, failure rate is considered as the ratio of the failure numbers to the network length.

Failures in WDNs can arise out of physical (pipe age, pipe diameter, pipe material, pipe installation, etc.), environmental (climate, pipe location, seismic activity, pipe bedding, groundwater, etc.), and operational (water quality, pressure, flow velocity, leakage, etc.) factors (Trifunović 2012). Researchers have made failure predictions by using different methodologies (Rajani & Kleiner 2001; Pelletier *et al.* 2003; Wilson *et al.* 2017; Robles-Velasco *et al.* 2020). Cooper *et al.* (2000) have made failure predictions in the WDN of London through the model based on probabilistic approach and geographic information system (GIS) by using certain variables such as pipe density function, soil corrosivity class, number of buses per hour, and pipe diameter. The authors believed that model results contribute to prioritization of the operational and investment plans of water utilities to reduce in WDN failures. Ilić (2009) has studied the factors affecting the failure frequency in the WDN of Zagreb, the capital city of Croatia, through the statistical log-linear analysis method. He stated that there are significant parameters affecting failures such as pressure regulation, pipe age, pipe diameter, and pipe material. Carrión *et al.* (2010) have analyzed the recorded failure data in WDNs through the non-parametric and semi-parametric methods to evaluate failure probabilities. The analysis results showed that certain pipe characteristics such as traffic conditions, diameter, material, and length have an impact on failures. Kakoudakis *et al.* (2017) have preferred the evolutionary polynomial regression model to predict the failures in UK WDNs by using pipe age, length, and diameter characteristics. The developed approach contributed to obtaining higher model output accuracies, reducing prediction errors in networks with excessive failures. Winkler *et al.* (2018) have set up the failure model for replacing plans of urban WDNs by using the decision tree learning technique. They indicated that the models can be a good alternative to traditional statistical failure approaches. Wols *et al.* (2019) have researched the effect of meteorological parameters such as air temperature, wind, and drought on failures in drinking WDNs of the Netherlands. At the end of the study, it is stated that temperature is a significant parameter for cast iron (CI) and asbestos-cement (AC) pipes. Failures have increased in high temperature for AC pipes and increased in low temperature for CI pipes. Also, it is stated that winds uprooting trees and drought periods cause failures.

Researchers have developed prediction models for failure rate (FR) using different methods. Wang *et al.* (2009) have recommended deterioration models to predict annual FR through the multiple regression method considering certain variables such as diameter, pipe age, pipe type, length, and depth of installation. According to the model results, the network length is

quietly effective on FR prediction. [Tabesh \*et al.\* \(2009\)](#) have predicted the FRs in WDNs in Iran by using the neuro-fuzzy systems, and multivariate regression methods. In this analysis, pipe depth, diameter, age, length, and average pressure variables have been used as input. Researchers indicated that artificial neural network (ANN) models have more accuracy and realistic results compared to the neuro-fuzzy systems and multivariate regression models. [Jafar \*et al.\* \(2010\)](#) have developed some models through the ANNs by using materials of pipes, the diameter of pipes, length of pipes, thickness, age, soil type, location, and pressure variables to predict the failure rate in the WDN of north France. The ANN models can help decision-makers to determine the best strategy for network replacing and maintenance investments. [Asnaashari \*et al.\* \(2013\)](#) have modeled FRs in their studies through the multiple linear regression (MLR) and ANN approaches. The model inputs in both studies are the following variables; pipe age, length, diameter, soil type, break category, pipe material, the year of Cathodic Protection (if applicable), and the year of Cement Mortar Lining (if applicable). The ANN model that was applied by the researchers at Etobicoke WDN in Canada gave better prediction results than MLR with  $R^2 = 0.94$ . On the other hand, [Shirzad \*et al.\* \(2014\)](#) have developed models using the ANN and support vector regression (SVR) methods to predict FRs at Mashhad WDN in Iran. Practitioners have used length, age, the hydraulic pressure of pipes, and installation depth diameter model variables as input, and model results showed that the ANN was a better predictor than the SVR. [Aydogdu & Firat \(2015\)](#) have modeled the FR through the least-squares supportive machine (LS-SVM) and fuzzy clustering approaches by using certain variables such as length, age, and diameter in the WDN of Malatya city. In their study, the FR and the effective factors on this rate have been analyzed. In this context, high FRs have been observed in the pipes with a length in the range of 0–200 m, pipes with a diameter of 110 m, and pipes between the ages of 15 and 20 years. The authors emphasized that using clustering analysis and LS-SVM methods together can be useful to obtain more effective performance in the prediction of FRs. In another study, annual failure rates at Scarborough WDN in Canada have been predicted by [Sattar \*et al.\* \(2019\)](#) through the extreme learning machine (ELM) model. For the model, 50 neurons maximum have been employed in the hidden layer. The prediction models showed that the most effective input variable was pipe diameter. [Kutylowska \(2019\)](#) has predicted the FR through the support vector machines (SVMs) and non-parametric regression methods using operational parameters (distribution pipes, house connections, total network length) in WDNs of Polish cities. The linear kernel function has been recommended to obtain the best prediction model. [Motiee & Ghasemnejad \(2019\)](#) have applied four different regression models (Poisson regression, logistic regression, exponential regression, and linear regression) to predict the FR in the WDN of Tehran by using length, material, pressure, age, and diameter variables. It is clear according to the all-model results that the ductile iron pipes among the others types such as ductile iron, asbestos-cement, cast iron, and high-density polyethylene have a significant role reducing failure number. [Shirzad & Safari \(2020\)](#) have been designed FR prediction models through the random forest (RF) and multivariate adaptive regression splines (MARS) techniques by using certain variables such as hydraulic pressure, pipe diameter, pipe age, pipe installation depth, and average, pipe length in the city network of Mahabad city, located in the north-east of Iran. In the first case study, MARS model outputs had higher prediction performance than RF.

Another approach used for the modern prediction models is the triple diagram method (TDM), a mapping technique based on the Kriging method. Counter charts are used to visualize the collective behavior of two independent variables and a dependent variable, and to provide the most ideal interpolation for prediction ([Özger & Sen 2007](#)). The TDM method enables to design of the most suitable surface model and to interpret the variable behaviors despite the extreme data scattering of the variables ([Kızılöz & Şişman 2021](#)). [Şişman & Kızılöz \(2020\)](#) have used the TDM structure to predict the non-revenue water ratio (NRWR). In another study conducted by the researchers, the NRWR has been predicted with high accuracy through the serial triple diagram model (STDM) ([Kızılöz & Şişman 2021](#)).

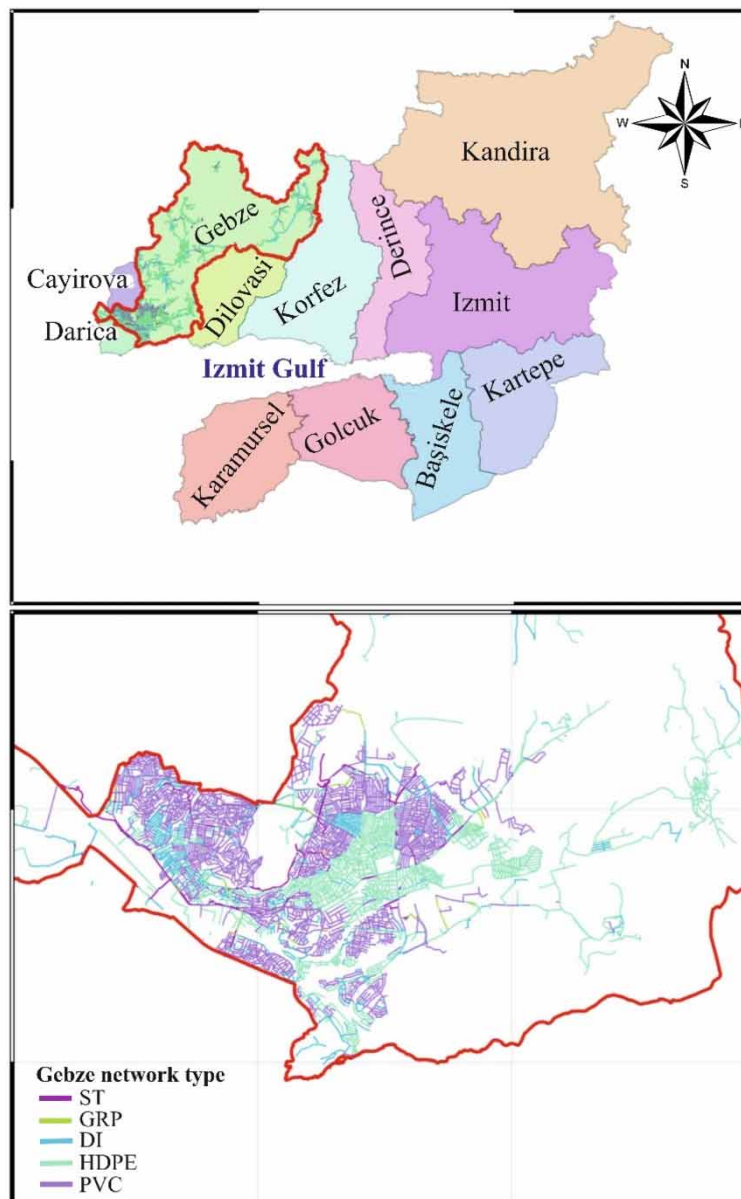
The aim of this study is to predict the FR in WDNs through the mapping technique based on ANNs, an artificially intelligent application, and Kriging methods. For this purpose, (1) the precipitation, temperature, diameter, age, pressure, and length variables have been selected as model inputs to predict the FR; (2) the data set has been normalized to scale model inputs and outputs to a certain scale, and to increase the model performance; (3) ideal ANN and STDMS have been designed by analyzing separately the effect of each variable on FR; (4) the best model accuracy with the least input has been analyzed.

## STUDY AREA AND DATA

The district of Gebze is selected for the modeling, which is the biggest district of Kocaeli city in terms of population density. Gebze, which is a significant trade, transportation, and logistics route between the European and Asian continents, is located

in the east of the Marmara region and the north of İzmit Gulf (Figure 1). The district is located on the border of Istanbul, and has been rapidly growing in population with 392,945 people in 2020. As of 2020, the number of water consumers was 157,412 people (ISU 2020), and its total surface was 584 km<sup>2</sup>. The infrastructure network of Gebze is operated by the ISU (Kocaeli Water and Sewerage Administration). In the region, there are two drinking and domestic water ponds, one drinking water treatment facility, seven natural water springs, 26 water tanks, and two drinking water pumping stations.

The total network length of the modeling area is 1,401 km and the service connection length is 420,78 km. The existing water distribution network consists of pipes with different material properties such as 4.75% steel (ST), 1.28% glass reinforced plastic (GRP), 18.66% ductile iron (DI), 39.78% high-density polyethylene (HDPE), and 35.53% polyvinyl chloride (PVC) (Figure 1). Daily measured failures on WDNs have been recorded by the ISU, Department of Water and Wastewater Technologies, and the Branch Office of Geographical Information Systems. In this modeling study, 5117 data in total, including 731 data for each model variable (measured and recorded daily values between 2019 and 2020) have been used to analyze,



**Figure 1** | General view and network type of Gebze.

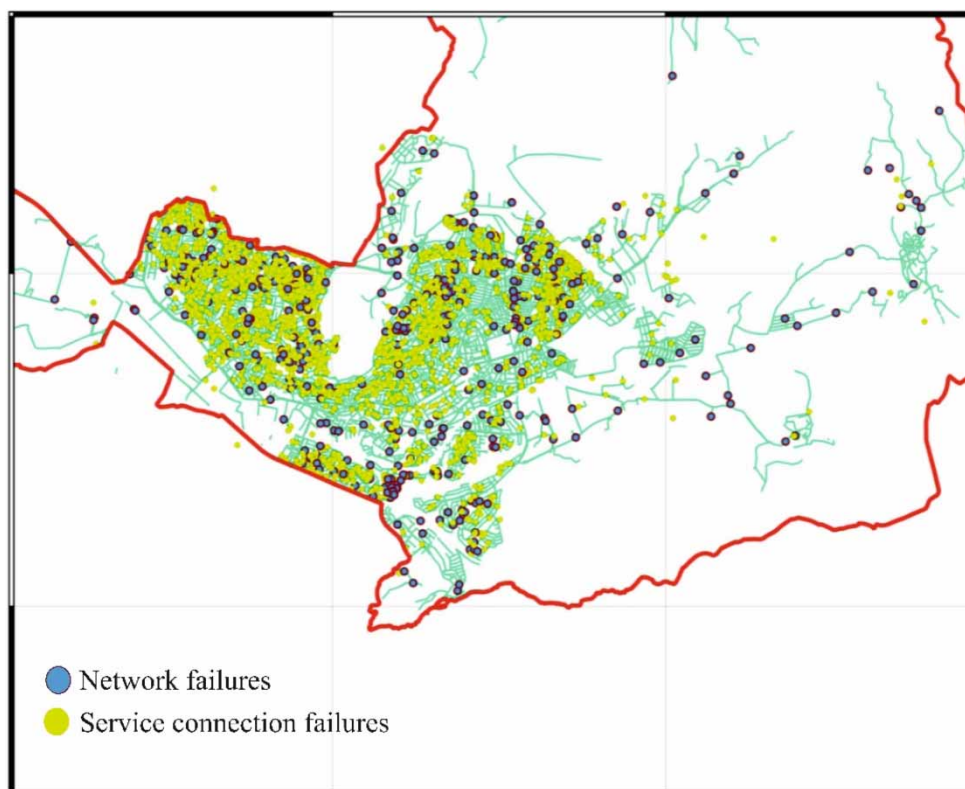


map, and make inferences on FRs. Each independent variable of the WDN measured daily at failure points has been given in Table 1. Different variables are available in this modeling data such as the mean diameter of the pipe (MDP), the mean age of pipe (MAP), the average pressure of pipe (APP), mean length of pipe (MLP), and daily FR (DFR). Also, certain meteorological variables such as daily mean air temperature (DMT) and daily total rainfall (DTR) have been used for the first time as model input and their relationships with FR have been analyzed. Daily mean temperature and precipitation data have been obtained from the Turkish Meteorological Service (MGM). Statistical properties of all variables used as model inputs and output are available in Table 1.

In the study area, 5,986 failures have occurred within the two-year periods (Figure 2). The 866 failures have occurred in the network (14.46%) and 5,120 in service connections (85.54%). On the other hand, the distribution of the failure network is as follows; 344 of PVC, 17 of GRP, 43 of DI, and 460 of HDPE. In this study, as in many literature studies, it is seen that most of

**Table 1** | Variable summary in the study model

Variables	Unit	Min	Mean	Max	Median	Std.D.
MDP	mm	32	50.79	441	39.60	38.15
MAP	year	2	19.02	33	19.50	4.13
APP	m	39.60	50.78	61.88	50.77	3.39
MLP	km	10.13	30.17	55.37	29.76	4.76
DMT	°C	−0.90	15.42	26.40	15.60	6.99
DTR	mm	0	1.44	45	0	4.25
DFR	–	0.02	0.27	0.93	0.25	0.17



**Figure 2** | Spot distribution of failures occurred in the study area.

the failures in the WDN have occurred in service connections. The service connection failure rate in total failure has been calculated as Nicolini *et al.* (2014) 58%, Aydogdu (2014) 60%, and Boztaş *et al.* (2019) 77.4%.

The standard water balance budget suggested by American Water Works Association (AWWA) and International Water Association (IWA) is available in Table 2 for the region of Gebze. According to this table, the total amount of water entering the system is 30,777,924 m<sup>3</sup>/year, the amount of water loss is 6,670,870 m<sup>3</sup>/year, and the amount of unbilled unmetered consumption is 409,884 m<sup>3</sup>/year (Table 2). It is seen in the table that the water loss amount arising out of failures is at the level of 1.33%. The related water utility tries to reduce this level up to 1% and water loss rate up to 20% through DMA, PM, minimum night flow monitoring, and active leakage control methods. If failures repeat intensively at the same point despite the above-mentioned activities, old networks and service connections will be replaced.

## METHODOLOGY

In this study, the DFR has been predicted via ANNs and STDMS. The stages of the model studies are given below respectively:

1. The daily failure number of the WDN of Gebze that was measured and recorded coordinately for model study between the years of 2019 and 2020 has been calculated.
2. The min-max normalization technique has been used to increase model prediction accuracy and measure the data on a certain scale.
3. Different combinations of model have been analyzed to predict the FR with the least input with the help of model input number increase, from one to three. The best prediction model has been obtained by the combination of two inputs.
4. Some statistical indicators such as R, MSE, BIAS, MAE, and SI have been used to evaluate model performances.
5. Model predictions have been made on the basis of the Kriging interpolation technique.
6. An STDMS has been used to increase the input number and accuracy in prediction charts.
7. The prediction accuracy of the STDMS has been evaluated through the same statistical indicators.
8. The FR prediction model with the highest accuracy has been determined by comparing the ANN and STDMS model results.

### Min-max normalization

Certain values on the input variable data set used for model applications can be bigger or lower than zero. Especially, extreme values affect negatively model results due to the differences between the input values. To remove any undesirable condition, the data can be normalized by scaling in the same interval. In the literature, there are different normalization methods preferred frequently such as Statistical or Z-Score, Min-Max, Median, Sigmoid, and Statistical Column.

In this study, the min-max normalization method, a method linearly normalizes the data, has been used to convert data into numerical values between 0 and 1. The most significant advantage of this method is to protect all relationships in data (Jaya-lakshmi & Santhakumaran 2011). The min-max method is available below (Equation (1)).

$$x_n = \frac{x_i - x_{min}}{x_{max} - x_{min}} \quad (1)$$

**Table 2** | Water balance components of Gebze in 2020

System input volume (SIV)	Authorized consumption	Billed authorized consumption	Billed meter consumption	76.58%	Revenue water
		76.75%	Billed unmetered consumption	0.17%	
30,777,924 m <sup>3</sup> /year	Water losses	Unbilled authorized consumption	Unbilled meter consumption	0.25%	Non-revenue water (NRW)
		1.58%	Unbilled unmetered consumption	1.33%	
		Apparent losses	Unauthorized consumption	0.62%	
		4.98%	Authorized consumption errors	4.36%	
%100	21.67%	Real losses	Leakage on transmission and distribution mains and service connections	16.65%	23.25%
		16.69%	Leakage and overflows at storage tanks	0.04%	

where,  $x_n$ , represents the normalized data,  $x_i$ , the input value,  $x_{min}$ , the lowest value in input data set, and  $x_{max}$  represents the highest value in input data set.

### ANNs

An ANN is a machine learning method used frequently for solving nonlinear complex engineering problems. The ANN is a system modeled by inspiring the biological behaviors of neurons in the human brain. A typical ANN structure is composed of three independent layers; input, hidden, and output (Figure 3). Also, there are five basic elements in the ANN structure named as inputs, weights, transfer function, activation function, and outputs besides the existing layers.

The information from the input layer ( $x_1, x_2, \dots, x_n$ ) is transferred to the transfer function being multiplied by weights. The weights ( $w_1, w_2, \dots, w_n$ ) show the importance of the information coming into the artificial cell and their effects on neurons. Also, the correct determination of the weights increases the learning performance of ANN models. The transfer function is used to calculate the net input coming into the cell as a result of multiplying the input and weights. The transfer function frequently used is given in Equation (2).

$$f(net) = \sum_{i=1}^n w_i x_i + bias \quad (2)$$

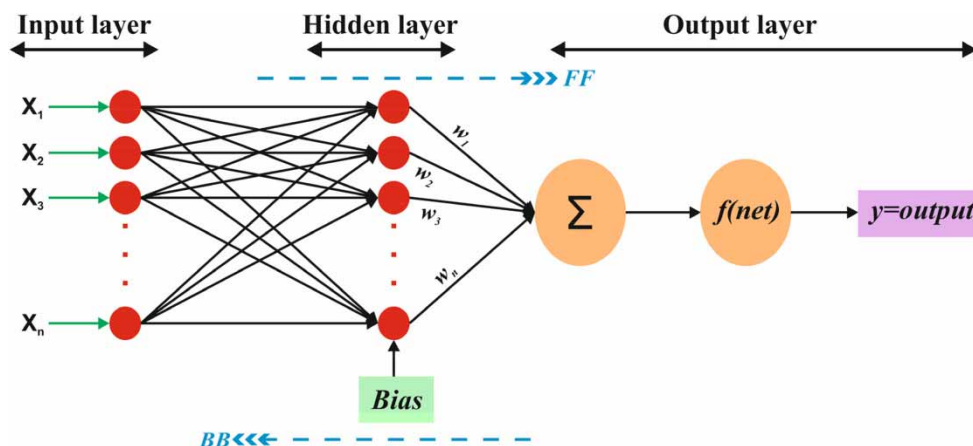
where  $x_i$ , represents the inputs,  $w_i$ , the weight values,  $n$ , the total input number, and  $f(net)$  represents the net input value.

The model output is calculated by passing  $f(net)$  obtained as a result of transfer function through the activation function. In the literature, there are various activation functions such as linear, sigmoid, hyperbolic tangent, step, and ramp. In this study, the sigmoid activation function that is used frequently has been preferred. The model output obtained by using this function is available below (Equation (3)).

$$y = f(net) = \frac{1}{1 + \exp[-(net)]} \quad (3)$$

Another issue after the determination of the ANN model structure is to train outputs to be produced in response to inputs in the network. The forward transfer to obtain target output throughout the model inputs is called Feed-Forward and the back-transfer to the input layer of the network, in case of any difference between target output and actual output, is called back-propagation (Kizilöz 2021). To reduce this difference at the acceptable levels in the stage of back-propagation, the weight and bias values of the network are changed iteratively, and thus, the train process of the network is completed. In addition, the weights that are assigned randomly at first are repeated continuously until achieving the required values. The Equation below is used to arrange the weights.

$$E_w = \frac{1}{2} \sum_{k=1}^n (y_p - z_a)^2 \quad (4)$$



**Figure 3** | The ANN structure used in the study.



where,  $E_w$ , indicates the sum of square errors,  $n$ , the weight number,  $y_p$ , the prediction, and  $z_a$  indicates the actual output values.

The Levenberg-Marquardt (LM) algorithm, a back-propagation algorithm, has been used for the training of ANN models in this study. This algorithm is a repetitive method used for non-linear least-square problems. The LM algorithm is a combination of gradient descent and Gauss-Newton methods. This algorithm is preferred frequently by researchers due to its speed and stability in the training of prediction models (Kizilöz *et al.* 2015; Şişman & Kizilöz 2020; Kizilöz 2021).

Determination of ideal hidden layer and neuron numbers in the ANN structure is also significant in achieving model predictions. In the literature, ANN applications with single and multi-hidden layers have been studied commonly by model designers to solve complex engineering problems (Zealand *et al.* 1999; Rajurkar *et al.* 2004; Şen *et al.* 2021; Kizilöz *et al.* 2022). As the hidden layer with a higher number preferred in most studies cannot develop the model prediction, the ANN model structure with a single hidden layer has been preferred in this study due to its easy application and prevalence (Kizilöz *et al.* 2015; Kizilöz 2021). On the other hand, there is no mathematical approach to determining the most effective neuron number in the hidden layer (Kizilöz 2021). This number is determined generally through the trial-and-error method (Kizilöz 2021). Also, in this modeling study, MDP, MAP, APP, MLP, daily mean temperature DMT, and DTR variables have been used as input, and DFR has been used as output.

### STDM

Differently from classical statistical methods, geo-statistics is a statistical method taking into account the relationship between the points and their coordinates based on the stationary randomness in the theory functions. The first stage of geo-statistical analysis is to make the semi-variogram (SV) analysis. This function is explained as the variance of the difference between two variables that are apart from each other as  $h$  distance. Also, changes based on the distance of spatial variables are identified by this function. The SV is available below (Equation (5)) for all measured space.

$$\gamma(h) = \frac{1}{2N(h)} \sum_{i=1}^{N(h)} [Z(x_i) - Z(x_j + h)]^2 \quad (5)$$

where  $\gamma(h)$  means the SV value,  $h$ ,  $i$  and  $j$  the distance between the points,  $N(h)$  the dot pair number in  $h$  length,  $Z(x_i)$  the measured value of the variable in  $i$  point, and  $Z(x_j + h)$  the measured value of the variable in  $j$  point.

The analysis of the SV is completed by adapting a theoretical model to the experimental SV and by making the goodness-of-fit test determining model parameters. The theoretical SV models used frequently in hydrological and hydro-meteorological studies are linear, spherical, Gaussian, exponential, and circular. In this study, the linear SV model structure has been preferred (Equation (6)).

$$\gamma(h) = C_0 + h \left( \frac{C}{A_0} \right) \quad (6)$$

where,  $C_0$  is the nugget variance;  $C$  is the structural variance;  $A_0$  is the range parameter;  $h$  is the vector distance between the observation pairs.

The interpolation technique applied after the theoretical SV structure has been mathematically determined to predict variable values from measured points to unmeasured points is called Kriging method (Krige 1951; Matheron 1963). In general, the prediction is made through the weighted average of the know values. The basic equation used for the Kriging method has been given below.

$$Z^0(x_0) = \sum_{i=1}^n \varepsilon_j z(x_i) \quad (7)$$

where,  $n$  is the number of points to be used in the Kriging prediction;  $z(x_i)$ , the measured values;  $\varepsilon_j$ , the weight coefficients; and  $Z^0(x_0)$ , the Kriging value predicted in the  $x_0$  point.

The Kriging approach is based on the least error squares method and is known as the best linear unbiased estimator. The weights determined through this approach, which depend on the SV and the spatial position of the data, are calculated as the mean of the difference between the prediction and actual values is zero and the prediction error variance is the smallest.

Researchers take advantage of various Kriging methods according to the study area and data structure such as Indicator Kriging (Armstrong 1998), Universal Kriging (Lark *et al.* 2014), Simple Kriging (Elbasiouny *et al.* 2014), Co-Kriging (Chica-Olmo *et al.* 2014), Ordinary Kriging (Sanusi *et al.* 2014), Block Kriging, and Punctual Kriging methods (Şişman & Kızılöz 2020). In this study, the Point Kriging method has been preferred to predict data at unmeasured points.

Model predictions can be made through the contour maps based on the above-mentioned variogram analysis and different Kriging methods. While changes of two variables each other that are composed of one dependent variable and one independent variable can be seen easily in coordinates in the X and Y planes, the behavior between them can be seen through these maps when the number of independent variables increases to three (Sirdas & Şen 2003). Model predictions with two inputs and one output are made by using input variables in the X and Y coordinates and output variables in the Z coordinate of these maps that also called the TDM structure. If the model results are not as accurate as desired, the predictions are repeated through the second TDM structure that is formed by the error terms (difference between the measured value and prediction) obtained after the modeling and new independent variable addition (Kızılöz & Şişman 2021). Thus, the model accuracy can be increased by designing a model structure with four inputs and one output through the error terms obtained from the first TDM and three independent input variables measured (Kızılöz & Şişman 2021). Sequential models can be repeated until achieving the prediction accuracy demanded by the model designer. The calculation steps of these maps are called the STDM by Kızılöz & Şişman (2021), and a sample model structure is available in Figure 4.

$$P_{1j} = f [m_{1j}, m_{2j}] \quad (8)$$

$$P_{2j} = f [m_{3j}, e_{1j}] \quad (9)$$

$$P_{(1j)} = a_{1j} \cdot m_{1j} + b_{1j} \cdot m_{2j} \quad (10)$$

$$e_{1j} = M_j - P_{1j} \quad (\text{measured DFR} - \text{predicted DFR}) \quad (11)$$

$$P_{2j} = c_{1j} \cdot m_{3j} + d_{1j} \cdot e_{1j} \quad (12)$$

$m_{ij}$  =  $ij$ . measured input data of  $i$ . series as MAP, APP, and etc. ( $i = 1, 2, 3$ );  $j = 1, 2, \dots, n$  where  $n$  is the record length,

$M_j$  = measured DFR as output data,  $j = 1, 2, \dots, n$  where  $n$  is the total number of data,

$P_{ij}$  =  $ij$ . model output as prediction DFR, ( $i = 1, 2$ )

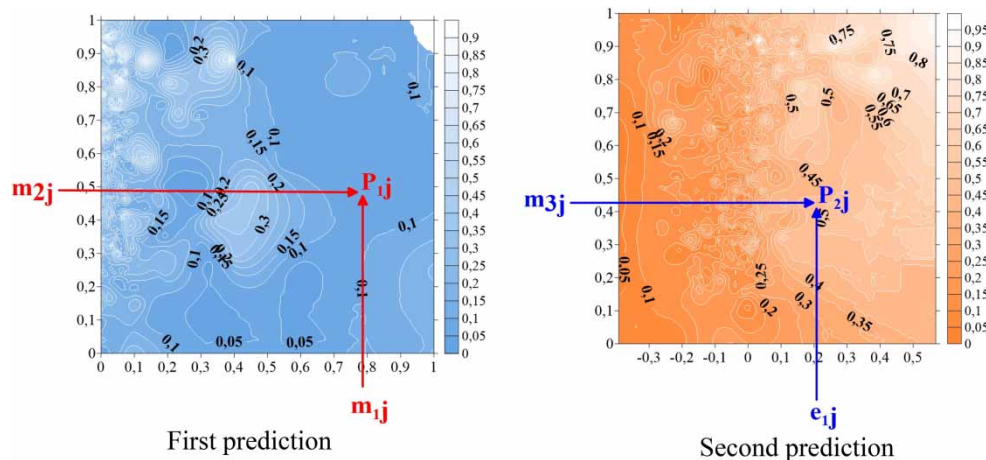
$e_{1j} = 1j$ . 1. model error terms

$e_{2j} = 2j$ . 2. model result errors

$a_{1j}$ ,  $b_{1j}$ ,  $c_{1j}$  and  $d_{1j}$  = model constant

### Model assessment criteria

Some statistical performance indicators have been used to evaluate the prediction accuracy of models, for example, coefficient of determination ( $R^2$ ), mean square error (MSE), root mean square error (RMSE), bias (BI), mean absolute error



**Figure 4** | STDM structure for model predictions.

(MAE), and scatter index (SI).

$$R^2 = \left[ \frac{\sum_{i=1}^n (A_i - \bar{A})(P_i - \bar{P})}{\sqrt{\sum_{i=1}^n (A_i - \bar{A})^2 \sum_{i=1}^n (P_i - \bar{P})^2}} \right]^2 \quad (13)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (A_i - P_i)^2 \quad (14)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (A_i - P_i)^2} \quad (15)$$

$$BI = \frac{1}{n} \sum_{i=1}^n (A_i - P_i) \quad (16)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n (|A_i - P_i|) \quad (17)$$

$$SI = \frac{RMSE}{\bar{A}} \quad (18)$$

where,  $n$  is the total number of data points,  $A_i$  is the actual data, and  $P_i$  is the predicted data. Finally,  $\bar{A}$  and  $\bar{P}$  are the means of the actual and prediction data.

## RESULTS AND DISCUSSION

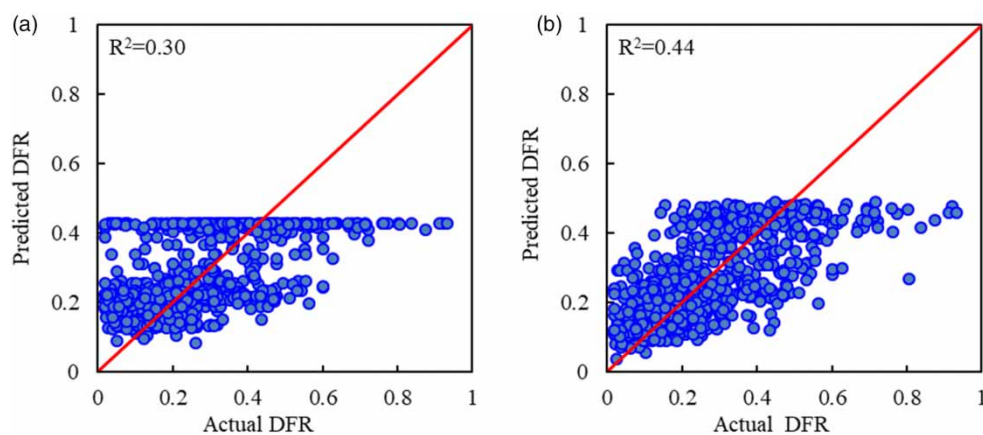
Before starting model applications, 731 raw data belonging to the MDP, MAP, APP, MLP, DMT, DTR variables that were measured and recorded daily between the 2019 and 2020 time period have been reduced in the range of [0 1] standardizing them through the min-max normalization method to scale the data in a certain range and increase model performances.

After the normalization process, various combinations have been formed through the ANNs to predict daily failure rates (DFRs). The model application has been made by using the Neural Network Fitting Toolbox in the MATLAB software. In this study, the three-layer feed-forward back propagation network structure, FFBP, with a single input, hidden, and output layer has been preferred (Figure 3). Also, the linear activation functions in the output layer and sigmoid function in the hidden layer have been used. The model data has been classified randomly; 60% training (439 data), 30% validation (239 data), and 10% testing (73 data). Before the training process, each model which has irregular weights and biases at the beginning has been initialized again (Kizilöz *et al.* 2015; Şişman & Kizilöz 2020; Kizilöz 2021). The Levenberg-Marquardt back propagation (trainlm) that is preferred frequently due to its speed and sensitivity has been used for the training algorithm of the ANN models in the study (Kizilöz *et al.* 2015; Şişman & Kizilöz 2020; Kizilöz 2021). The neuron number in the hidden layer that is effective on the model accuracy has been defined via the trial-and-error methods. Firstly, a single neuron was used in the hidden layer, and then the model performances were evaluated by increasing the neuron numbers one by one. According to the results, the best model performance was obtained when five neurons were used in the hidden layer. Therefore, it was decided to have five neurons in the hidden layer.

To analyze the effect of some independent variables such as MDP, MAP, APP, MLP, DMT, DTR, ANN models with one input and one output have been designed and their model performances have been given in detail in Table 3. According to the model results, the most effective variables are DMT and MAP, respectively. The DMT prediction model has the highest coefficient of determination ( $R^2 = 0.30$ ) and the lowest mean square error ( $MSE = 0.021$ ). It is seen in Figure 5(a) that the scatterings are too excessive with 52.64%, according to the 1:1 straight line (45°) used for the comparison of actual and predicted data of the DMT model. A similar effect of temperature, which is a meteorological variable, on the failure has been also revealed by certain researchers such as Wols & Van Thienen (2014), Pietrucha-Urbanik (2015), and Wols *et al.* (2019). The determination coefficient of the MAP model, another independent variable effective on DFR, is 0.23 and its mean square

**Table 3** | Accuracy of ANN models with one input variable

Model no	Model inputs	R <sup>2</sup>	MSE	BIAS	MAE	SI%
ANN-1.1	DMT	0.30	0.021	−0.008	0.113	52.64
ANN-1.2	MAP	0.23	0.023	−0.004	0.119	54.98
ANN-1.3	MDP	0.18	0.024	−0.002	0.123	56.49
ANN-1.4	MLP	0.15	0.025	−0.005	0.125	57.73
ANN-1.5	APP	0.12	0.026	−0.002	0.128	58.70
ANN-1.6	DTR	0.10	0.028	0.004	0.135	61.13

**Figure 5** | Scatter graphs for (a) DMT and (b) MAP-DMT models.

error is only 0.023. In the studies conducted by Wang *et al.* (2009), Shi *et al.* (2013), and Aydogdu & Firat (2015) on the MAP variable, it is also emphasized that the network age is an effective factor on the FR. The ANN models with a one-input that have been applied in this study showed again that the MAP, MDP, APP, and NLP variables are effective on the DFR. Additionally, meteorological variables such as DMT and DTR have been used as model input in this study for the first time to make DFR predictions.

The predictions belonging to the twelve different combinations that formed model structures with single input and single output have been sorted in Table 4 according to some statistical performance indicators such as R<sup>2</sup>, MSE, BI, MAE, and SI to increase model accuracies of DFRs. According to this sorting, the best prediction model has been obtained through the MAP-DMT combination (Figure 5(b)). When the prediction model results with two inputs are compared with the model results with single input given in Table 3, a performance development is observed over 46%.

Various model combinations with three inputs have been designed increasing input number to develop prediction accuracy of ANN models. Twelve prediction models evaluated according to the five different performance criteria are available below in Table 5. When these model results are compared with model results with two inputs, there is no remarkable improvement. The determination coefficients of the models given in Table 4 change between 0.22 and 0.44. These coefficients change between 0.21 and 0.40 in Table 5. In brief, the desired predictive accuracies have not been obtained through the model results with three inputs.

When thirty different models with one, two, and three inputs that applied to predict the DFR through the ANN approach are analyzed together with R<sup>2</sup>, MSE, BI, MAE, and SI statistical performance indicators, it is seen that the best model prediction has been achieved by the MAP-DMT model combination. The best actual-prediction relationship with one and two inputs is given together in Figure 5. Model structure with single input is insufficient for predictions. On the other hand, the expected development of prediction is not also obtained through the model structure with three inputs. When

**Table 4** | Accuracy of ANN models with two input variables

Model no	Model inputs	R <sup>2</sup>	MSE	BIAS	MAE	SI%
ANN-2.1	MAP-DMT	0.44	0.016	0.005	0.101	46.68
ANN-2.2	DMT-MLP	0.41	0.017	0.002	0.102	48.15
ANN-2.3	MDP-DMT	0.38	0.018	0.005	0.105	49.11
ANN-2.4	MAP-NL	0.35	0.019	-0.004	0.108	50.58
ANN-2.5	MDP-MAP	0.31	0.020	0.001	0.111	51.77
ANN-2.6	AAP-DMT	0.30	0.021	0.003	0.111	52.21
ANN-2.7	MDP-NL	0.28	0.021	-0.003	0.113	52.95
ANN-2.8	MAP-APP	0.28	0.021	0.009	0.113	53.33
ANN-2.9	MDP-APP	0.24	0.022	-0.002	0.119	54.37
ANN-2.10	MAP-DTR	0.24	0.022	0.005	0.116	54.60
ANN-2.11	APP-NL	0.23	0.023	0.007	0.116	54.87
ANN-2.12	MDP-DTR	0.22	0.023	0.007	0.119	55.28

**Table 5** | Accuracy of ANN models with three input variables

Model no	Model inputs	R <sup>2</sup>	MSE	BIAS	MAE	SI%
ANN-3.1	MDP-MAP-DMT	0.40	0.018	0.006	0.104	48.56
ANN-3.2	APP- DMT -MLP	0.39	0.018	-0.001	0.104	48.85
ANN-3.3	DMT -DTR- MLP	0.38	0.018	0.007	0.103	49.43
ANN-3.4	MAP – APP – DMT	0.38	0.018	-0.012	0.106	49.36
ANN-3.5	APP – DMT – DTR	0.33	0.020	-0.001	0.110	51.21
ANN-3.6	MDP – APP – DMT	0.33	0.020	-0.008	0.111	51.33
ANN-3.7	MAP – APP – MLP	0.31	0.020	-0.003	0.111	51.78
ANN-3.8	MDP – MAP – APP	0.27	0.021	0.003	0.115	53.26
ANN-3.9	MAP- APP – DTR	0.25	0.022	0.008	0.117	54.02
ANN-3.10	MDP – MAP – DTR	0.25	0.022	0.002	0.116	54.30
ANN-3.11	MDP – APP – MLP	0.24	0.022	0.007	0.114	54.67
ANN-3.12	MDP – MAP – MLP	0.21	0.023	-0.001	0.119	55.55

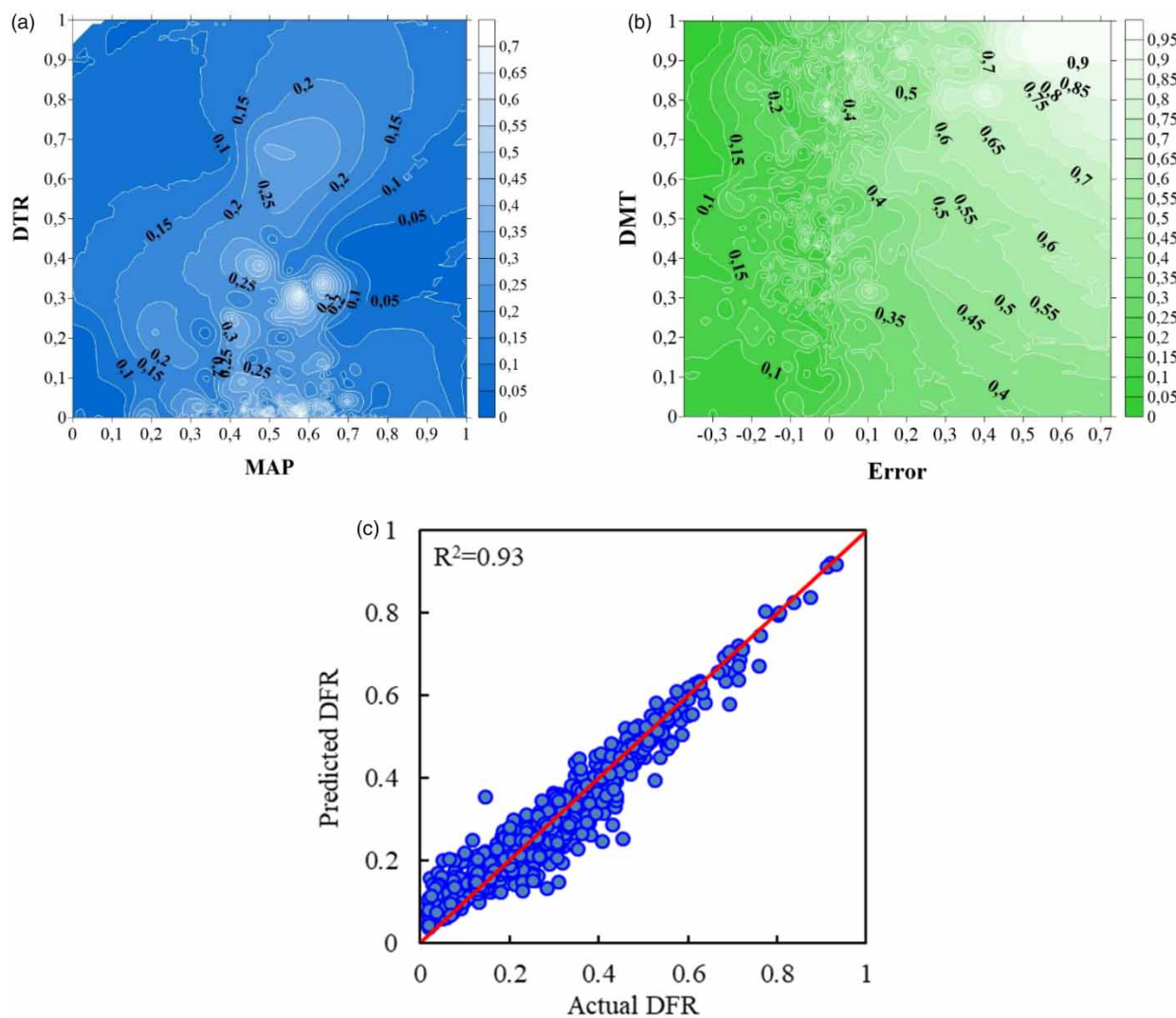
ANN model results are examined all together, it is clear that some variables such as MDP, MAP, APP, MLP, DMT, and DTR are effective factors in the DFR prediction. In a conclusion, after being observed that the desired model performances cannot be obtained through the ANN model approach for the DFR predictions, it has been tried to obtain predictions models with high accuracy through the STDM.

In this study, two-dimensional contour charts with three variables have been used through the Surfer program to make predictions with the STDM structure. It is necessary to have one dependent output and two independent variables (input) to obtain a contour chart through the Surfer program for problem-solving. While the independent variables are shown on the X and Y-axis in the prediction charts, the dependent variable values are shown on the Z-axis with curves depending on independent variables. The STDM structure has been applied for the first time for model charts to predict the DFR. This structure is composed of different model prediction charts that are linked to each other and have two inputs. In the beginning, the first model prediction has been made with two independent variables that are effective on the DFR. In case the model outputs of these charts are not in the desired accuracy, a new model prediction is made through the prediction errors of the first model and with the help of a third independent variable. Thus, an STDM structure with one output and four inputs is formed through the prediction error of the first model and three independent variables. STDM structures can be developed providing forward model flows by each model prediction error being the input of the next model.



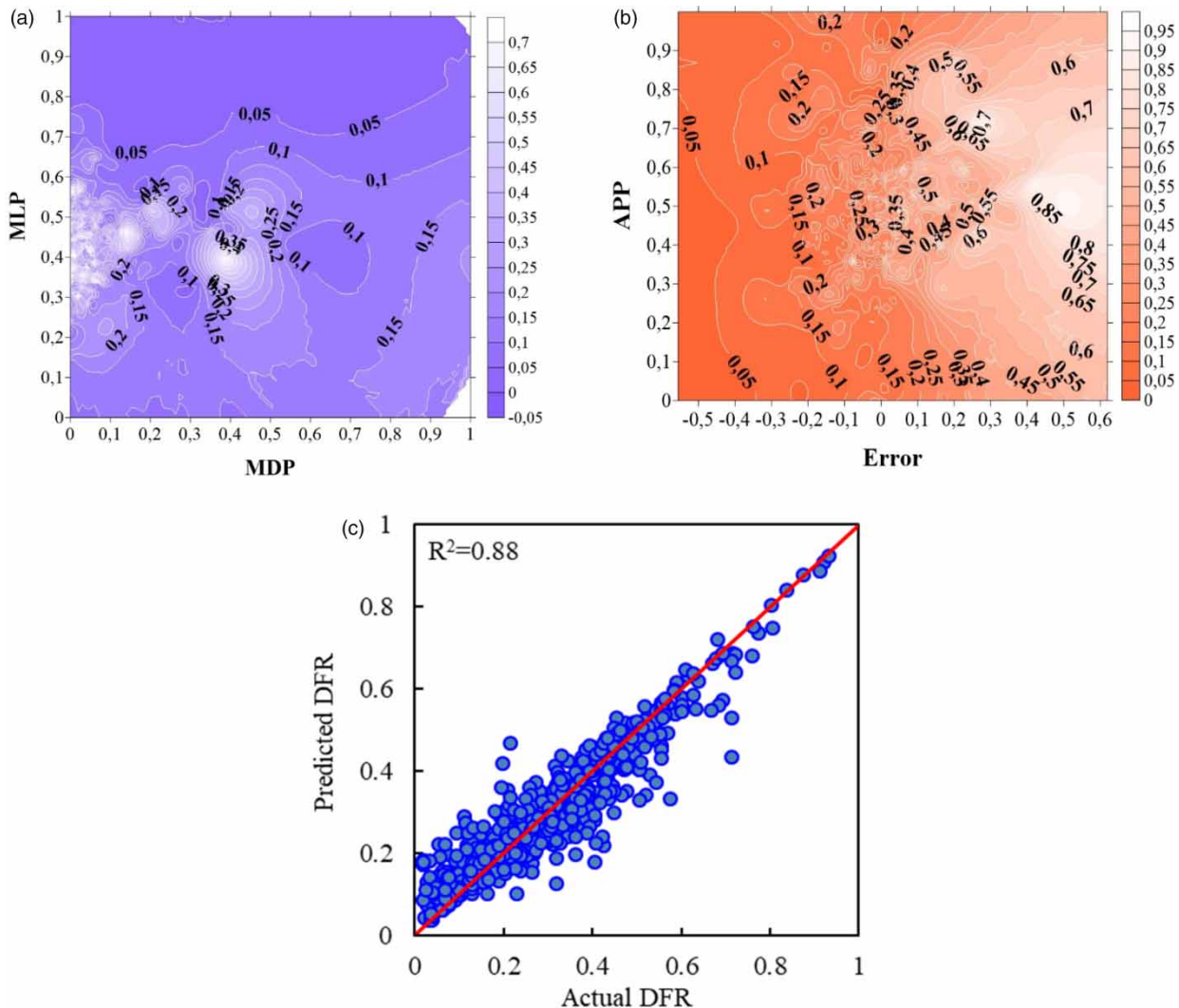
**Table 6** | Accuracy of STDMs with four input variables

Model no	First model inputs	Second model inputs	R <sup>2</sup>	MSE	BIAS	MAE	SI%
STDM.1	MAP-DTR	Error-DMT	0.93	0.002	0.000	0.033	16.62
STDM.2	MDP-APP	Error- DMT	0.92	0.002	0.000	0.034	18.01
STDM.3	APP – DTR	Error- DMT	0.92	0.002	–0.001	0.035	18.39
STDM.4	MAP – APP	Error- DMT	0.90	0.003	0.000	0.039	20.18
STDM.5	MDP – MAP	Error- APP	0.89	0.003	–0.003	0.041	21.50
STDM.6	MAP – MLP	Error- APP	0.88	0.003	–0.003	0.042	21.77
STDM.7	MDP – MLP	Error- APP	0.88	0.003	–0.001	0.044	22.23
STDM.8	DTR – MLP	Error- APP	0.87	0.004	–0.004	0.043	22.71
STDM.9	MDP – DMT	Error- MAP	0.84	0.004	0.001	0.050	25.19
STDM.10	MDP – DTR	Error- APP	0.84	0.005	–0.002	0.049	26.04
STDM.11	DMT -MLP	Error- APP	0.83	0.005	–0.004	0.052	26.24

**Figure 6** | The STD structure with MAP, DTR, Error, and DMT inputs.

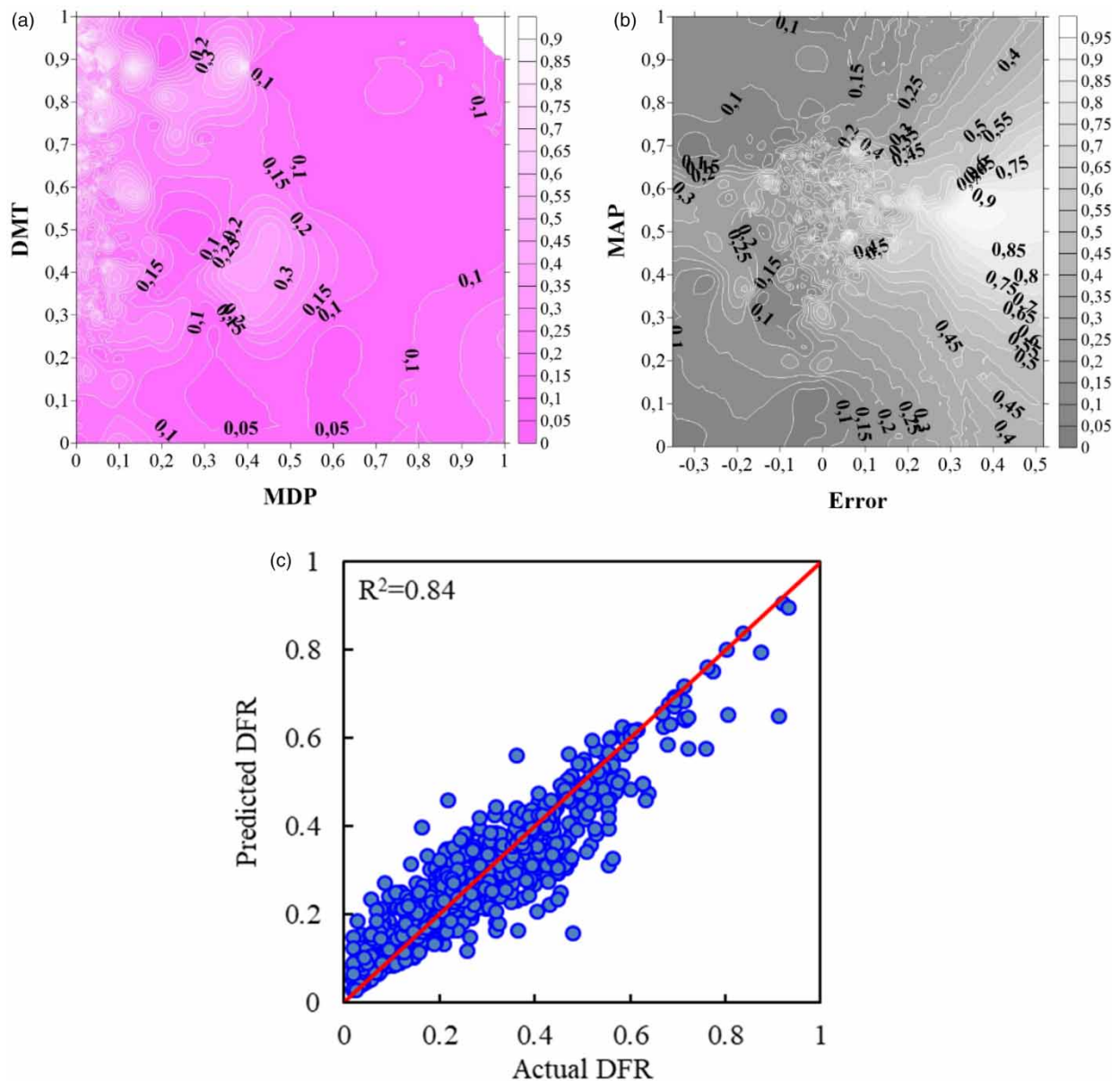
In this study, the statistical performance indicators of eleven models developed by the STD structure are available below in Table 6. According to this table, the  $R^2$  values are between 0.93 and 0.83, MSE values are 0.002–0.005, BI values are –0.004 to 0.001, MAE values are 0.033–0.052, and SI values are in the range of 16.62–26.24%. When the forty-one model accuracies in this study are analyzed together, it is observed that the STD predictions are higher in comparison with the ANN predictions.

The best STD prediction model charts exist in Figure 6. The model charts given in Figure 6(a) and 6(b) have been prepared using MAP, DTR, and DMT variables. Thus, prediction models with four inputs (MAP-DTR-Error-DMT) have been developed for the DFR. According to the model charts, there is a directly proportional relationship between the DMT and DFR. Similarly, the failure increase depending on the increasing temperature has been also revealed in the studies conducted by Wols & Van Thienen (2014) and Wols *et al.* (2019). The failure reasons depending on the temperature are briefly as follows; excessive soil shrinkage, absence of appropriate depth of excavation according to the meteorological conditions, and especially deformation of service connections due to not being covered. The performance of the STD output given in Figure 6(c) has been determined by comparing actual model predictions by referencing the 45° (1:1) straight line shown in red. According to the model charts, the highest prediction model accuracy has been obtained with an  $R^2$  value of 0.93 and an SI value of 16.62%.



**Figure 7** | The STD structure with MDP, MLP, Error, and APP inputs.

Another STDM structure recommended for predicting the DFR is available in Figure 7. The model predictions charts designed using the MDP, MLP, and APP independent variables have been given in Figure 7(a) and 7(b). It is seen in the STDM output prediction chart given in Figure 7(b) that the DFRs tend to increase distinctly depending on the increasing pressure. This problem is at the forefront to be solved for water utilities due to increasing failure rates cause a large amount of water loss. The relationship between the pressure and failure/leaks rates has been also revealed in detail through the studies conducted by Kanakoudis & Muhammetoglu (2014) and Kizilöz (2021). According to these researchers, failure and leakage increases are inevitable under the influence of high pressure. Water utilities and practitioners think that the most effective pressure management method is the PM for reducing failures. Failures can be reduced by providing ideal network operating pressure by separating the WDN into controllable and measurable isolated areas (district metered areas) and placing pressure-reducing valves (PRVs) into the points that are determined as a consequence of hydraulic model studies. According to the STDM chart given in Figure 7(c), the  $R^2$  value is 0.88, and the SI value is 22.23% of the MDP-MLP-



**Figure 8** | The STDM structure with MDP, DMT, Error, and MAP inputs.

Error-APP model. Also, according to the actual values of the STDM chart, values above the 45° straight line show that model predictions are higher than values below this line.

Another prediction model has been made through the MDP, DMT, and MAP variable combination (Figure 8). The increase of DFR can be seen clearly in this model depending on the MAP increase. In the pipes exposed to high pressure and lost their material properties due to age, failures increase. A similar failure increase is seen in excessive temperatures in old service connections that are not buried deep enough or completely uncovered. In Figure 8(c), the actual DFR and model prediction values have been given comparatively referencing the 45° straight line. Scatterings are below the straight line according to this figure. The  $R^2$  value of the DFR prediction model composed of the MDP- DMT-Error-MAP inputs has been calculated as 0.84.

## CONCLUSION

In this study, the ANN and STDM approaches have been used to predict the DFR in the water distribution system in Gebze. The higher prediction performances have been obtained through the STDM, a structure that enables to visualization of the model results by making inferences. The best STDM model has been obtained through the MAP-DTR-Error-DMT combination with four inputs. It is seen that the MDP, MAP, APP, and MLP variables used as model inputs are affected on the DFR in keeping with the previous studies.

In this study, prediction models have been developed by using some meteorological variables for the first time, for instance, DMT and DTR. It is understood that each of these variables is effective on the DFR. A significant DFR increase has been observed through the STDM method in increasing temperatures. Similar increases are also available in MAP and APP variables as such in DFRs. Increases in temperature, age, and pressure values cause serious failures, especially in service connections, so undesired water losses occur. Water utilities can reduce water losses by applying various activities such as district metered area design, pressure management, and active leakage control method.

In the FR predictions, preferring models allowing the interpretation of outputs such as STDM, contributes to the network management strategies to be developed by water utilities, experts, and researchers. Thanks to the prediction model results, water utilities can prioritize their future investments and plans, and thus water losses can be reduced systematically. FR predictions can be made in further studies by taking some variables into account, for instance, depth of excavation, hydraulic radius, speed, pipe thickness, and area of failure point.

## DATA AVAILABILITY STATEMENT

All relevant data are included in the paper or its Supplementary Information.

## CONFLICT OF INTEREST

The authors declare there is no conflict.

## REFERENCES

- Armstrong, M. 1998 *Basic Linear Geostatistics*. Springer-Verlag, Berlin, Germany.
- Asnaashari, A., McBean, E. A., Gharabaghi, B. & Tutt, D. 2013 Forecasting watermain failure using artificial neural network modelling. *Canadian Water Resources Journal* **38** (1), 24–33.
- Aydogdu, M. 2014 *Analysis of Pipe Failure Occurred in Water Distribution System Using Cluster Method*. MSc Thesis. Inonu University, p. 161. (in Turkish).
- Aydogdu, M. & Firat, M. 2015 Estimation of failure rate in water distribution network using fuzzy clustering and LS-SVM methods. *Water Resources Management* **29** (5), 1575–1590.
- Boztaş, F., Özdemir, Ö., Durmuşçelebi, F. M. & Firat, M. 2019 Analyzing the effect of the unreported leakages in service connections of water distribution networks on non-revenue water. *International Journal of Environmental Science and Technology* **16** (8), 4393–4406.
- Carrión, A., Solano, H., Gamiz, M. L. & Debón, A. 2010 Evaluation of the reliability of a water supply network from right-censored and left-truncated break data. *Water Resources Management* **24** (12), 2917–2935.
- Chica-Olmo, M., Luque-Espinar, J. A., Rodriguez-Galiano, V., Pardo-Igúzquiza, E. & Chica-Rivas, L. 2014 Categorical indicator Kriging for assessing the risk of groundwater nitrate pollution: the case of Vega de Granada aquifer (SE Spain). *Science of the Total Environment* **470–471**, 229–239.
- Cooper, N. R., Blakey, G., Sherwin, C., Ta, T., Whiter, J. T. & Woodward, C. A. 2000 The use of GIS to develop a probability-based trunk mains burst risk model. *Urban Water* **2** (2), 97–103.
- Elbasiouny, H., Abowaly, M., Abu Alkheir, A. & Gad, A.-A. 2014 Spatial variation of soil carbon and nitrogen pools by using ordinary Kriging method in an area of north Nile Delta, Egypt. *CATENA* **113**, 70–78.



- Ilić, K. 2009 The analysis of influential factors on the frequency of pipeline failures. *Water Science and Technology: Water Supply* **9** (6), 689–698.
- ISU (Kocaeli Water and Sewage Administrative General Directorate) 2020 *Kocaeli Water and Sewage Water Administrative: Annual Activity Report for Water Management*. ISU, Kocaeli, Turkey.
- Jafar, R., Shahrour, I. & Juran, I. 2010 Application of Artificial Neural Networks (ANN) to model the failure of urban water mains. *Mathematical and Computer Modelling* **51** (9–10), 1170–1180.
- Jayalakshmi, T. & Santhakumaran, A. 2011 Statistical normalization and back propagation for classification. *International Journal of Computer Theory and Engineering* **3** (1), 1793–8201.
- Kakoudakis, K., Behzadian, K., Farmani, R. & Butler, D. 2017 Pipeline failure prediction in water distribution networks using evolutionary polynomial regression combined with K-means clustering. *Urban Water Journal* **14** (7), 737–742.
- Kanakoudis, V. & Muhammetoglu, H. 2014 Urban water pipe networks management towards non-revenue water reduction: two case studies from Greece and Turkey. *Clean – Soil, Air, Water* **42** (7), 880–892.
- Kizilöz, B. 2021 Prediction model for the leakage rate in a water distribution system. *Water Supply* **21** (8), 4481–4492.
- Kizilöz, B. & Şişman, E. 2021 Non-revenue water ratio prediction with serial triple diagram model. *Water Supply* **21** (8), 4263–4275.
- Kizilöz, B., Çevik, E. & Aydoğan, B. 2015 Estimation of scour around submarine pipelines with Artificial Neural Network. *Applied Ocean Research* **51**, 241–251.
- Kizilöz, B., Şişman, E. & Oruç, H. N. 2022 Predicting a water infrastructure leakage index via machine learning. *Utilities Policy* **75**, 101357.
- Krige, D. G. 1951 A statistical approach to some basic mine valuation problems on the Witwatersrand. *Journal of Southern African Institute of Mining and Metallurgy* **52** (6), 119–139.
- Kutyłowska, M. 2019 Forecasting failure rate of water pipes. *Water Science and Technology: Water Supply* **19** (1), 264–273.
- Lark, R. M., Ander, E. L., Cave, M. R., Knights, K. V., Glennon, M. M. & Scanlon, R. P. 2014 Mapping trace element deficiency by cokriging from regional geochemical soil data: a case study on cobalt for grazing sheep in Ireland. *Geoderma* **226–227**, 64–78.
- Matheron, G. 1963 Principles of geostatistics. *Economical Geology* **58**, 1246–1266.
- Motiee, H. & Ghasemnejad, S. 2019 Prediction of pipe failure rate in Tehran water distribution networks by applying regression models. *Water Science and Technology: Water Supply* **19** (3), 695–702.
- Nicolini, M., Giacomello, C., Scarsini, M. & Mion, M. 2014 Numerical modeling and leakage reduction in the water distribution system of Udine. *Procedia Engineering* **70**, 1241–1250.
- Özger, M. & Sen, Z. 2007 Triple diagram method for the prediction of wave height and period. *Ocean Engineering* **34**, 1060–1068.
- Pelletier, G., Mailhot, A. & Villeneuve, J.-P. 2003 Modeling water pipe breaks – three case studies. *Journal of Water Resources Planning and Management* **129** (2), 115–123.
- Pietrucha-Urbanik, K. 2015 Failure analysis and assessment on the exemplary water supply network. *Engineering Failure Analysis* **57**, 137–142.
- Rajani, B. & Kleiner, Y. 2001 Comprehensive review of structural deterioration of water mains: physically based models. *Urban Water* **3**, 151–164.
- Rajurkar, M. P., Kothiyari, U. C. & Chaube, U. C. 2004 Modeling of the daily rainfall-runoff relationship with artificial neural network. *Journal of Hydrology* **285**, 96–113.
- Robles-Velasco, A., Cortés, P., Muñizuri, J. & Onieva, L. 2020 Prediction of pipe failures in water supply networks using logistic regression and support vector classification. *Reliability Engineering and System Safety* **196**, 106754.
- Sanusi, M. S. M., Ramli, A. T., Gabdo, H. T., Garba, N. N., Heryanshah, A., Wagiran, H. & Said, M. N. 2014 Isodose mapping of terrestrial gamma radiation dose rate of Selangor state, Kuala Lumpur and Putrajaya, Malaysia. *Journal of Environmental Radioactivity* **135**, 67–74.
- Sattar, A. M. A., Ertugrul, Ö. F., Gharabaghi, B., McBean, E. A. & Cao, J. 2019 Extreme learning machine model for water network management. *Neural Computing and Applications* **31** (1), 157–169.
- Şen, Z., Şişman, E. & Kizilöz, B. 2021 A new innovative method for model efficiency performance. *Water Supply* **22** (1), 589–601.
- Shi, W. Z., Zhang, A. S. & Ho, O. K. 2013 Spatial analysis of water mains failure clusters and factors: a Hong Kong case study. *Annals of GIS* **19** (2), 89–97.
- Shirzad, A. & Safari, M. J. S. 2020 Pipe failure rate prediction in water distribution networks using multivariate adaptive regression splines and random forest techniques. *Urban Water Journal* **16** (9), 653–661.
- Shirzad, A., Tabesh, M. & Farmani, R. 2014 A comparison between performance of support vector regression and artificial neural network in prediction of pipe burst rate in water distribution networks. *KSCE Journal of Civil Engineering* **18** (4), 941–948.
- Sirdas, S. & Şen, Z. 2003 Spatio-temporal drought analysis in the Trakya region, Turkey. *Hydrological Sciences Journal* **48** (5), 809–820.
- Şişman, E. & Kizilöz, B. 2020 Artificial neural network system analysis and Kriging methodology for estimation of non-revenue water ratio. *Water Science and Technology: Water Supply* **20** (5), 1871–1883.
- Tabesh, M., Soltani, J., Farmani, R. & Savic, D. 2009 Assessing pipe failure rate and mechanical reliability of water distribution networks using data-driven modeling. *Journal of Hydroinformatics* **11** (1), 1–17.
- Trifunović, N. 2012 *Pattern Recognition for Reliability Assessment of Water Distribution Networks*. PhD Dissertation, UNESCO-IHE Institute for Water Education, Delft Univ. of Technology, Delft, The Netherlands.
- Wang, Y., Zayed, T. & Moselhi, O. 2009 Prediction models for annual break rates of water mains. *Journal of Performance of Constructed Facilities* **23** (1), 47–54.



- Wilson, D., Filion, Y. & Moore, I. 2017 [State-of-the-art review of water pipe failure prediction models and applicability to large-diameter mains](#). *Urban Water Journal* **14** (2), 173–184.
- Winkler, D., Haltmeier, M., Kleidorfer, M., Rauch, W. & Tscheikner-Gratl, F. 2018 [Pipe failure modelling for water distribution networks using boosted decision trees](#). *Structure and Infrastructure Engineering* **14** (10), 1402–1411.
- Wols, B. A. & Van Thienen, P. 2014 [Impact of weather conditions on pipe failure: a statistical analysis](#). *Journal of Water Supply: Research and Technology – AQUA* **63** (3), 212–223.
- Wols, B. A., Vogelaar, A., Moerman, A. & Raterman, B. 2019 Effects of weather conditions on drinking water distribution pipe failures in the Netherlands. *Water Science and Technology: Water Supply* **19** (2), 404–416.
- Zealand, C. M., Burn, D. H. & Simonovic, S. P. 1999 [Short term stream flow forecasting using artificial neural networks](#). *Journal of Hydrology* **214**, 32–48.

First received 2 January 2022; accepted in revised form 23 August 2022. Available online 2 September 2022