# Prediction of pipe failure rate in Tehran water distribution networks by applying regression models

Homayoun Motiee and Sonya Ghasemnejad

## ABSTRACT

Four statistical models (linear regression, exponential regression, Poisson regression and logistic regression) applied to analyze the variables in pipe vulnerabilities with the objective of finding equations to predict probable future pipe accidents. The most effective variables in pipe failures are material, age, length, diameter and hydraulic pressure. To evaluate these models, the data collected in recent years in the water distribution network of district 1 in Tehran were used, with a total length of 582,702 m of pipes, and 48,500 consumers. The results demonstrate that among the four studied models, the logistic regression model is best able to give a good performance and is capable of predicting future accidents with a higher probability.

**Key words** | distribution water networks, failures, pipes, regression models, root mean square error (RMSE)

**Homayoun Motiee** (corresponding author)
**Sonya Ghasemnejad**
Civil, Water and Environmental Faculty,
Shahid Beheshti University,
Tehran,
Iran
E-mail: h_motiei@sbu.ac.ir

## INTRODUCTION

Accidents and pipe failures in urban water systems, amongst them water distribution networks, create great challenges for water utilities and researchers around the world (Kropp & Herz 2005; Grigg 2012). Winkler et al. (2018) used a boosted technique of decision trees for pipe failure prediction based on existing network data and historical failure records in a medium sized city, and concluded that this method has the best performance for studying real pipe failures in comparison to other models.

Wilson et al. (2017) studied pipe failure prediction for large diameter water mains (greater than 500 mm) within a period of 13 years by a series of statistical models. They concluded that the studied models were capable to predict the failures of an individual pipe or pipe segments with a high probability. Alvisi & Franchini (2010) used two proposed parameterization models by Mailhot et al. (2000) and Le Gat & Eisenbeis (2000) to estimate the pipe failures and concluded that both models have similar performances in an observed period. Savic (2009) and Kakoudakis et al. (2017) developed an evolutionary polynomial regression

model (EPR) with K-means clustering to predict the failure of pipes based on length, diameter, and age.

Inanloo et al. (2016) investigated the risk of vulnerability and pipe failure in distribution networks in Miami, Florida, according to a decision aid GIS-based risk assessment. In their research, Pelletier et al. (2003) used a variety of scenarios with a strategy model to find the number of annual failures in the water networks of three cities. Tabesh et al. (2009) used two models based on the Data-Driven Modeling named Artificial Neural Networks and Neuro-fuzzy systems along with a multivariate regression approach to predict more precisely the failure of pipes and concluded that the Neural network outcomes are more accurate. Kabir et al. (2015) reported that in the year 2000 there were on average 700 water pipe breaks in Canada and USA incurring a cost equivalent to 10 billion Canadian dollars for these two governments. Most current actions to make decisions on a pipe's repair or exchange were based only on the age of the pipe and the number of previous breaks (Shamir & Howard 1979; Kettler & Goulter 1985; Andreou et al. 1987; Le gat &

Eisenbeis 2000; Xu et al. 2018). Debón et al. (2010) compared models with using a receiver operating characteristics (ROC) graph to evaluate them and concluded that the generalized linear regression model is the best model for this subject.

Some researchers in Iran note that over 20–30% of the total revenues of Iranian water and wastewater companies is spent on repair, rehabilitation and corrective actions, and about 30% of accidents occurred in the pipes of the distribution system (Elahipanah 1999; Beigi 2000; Tajrishi & Abrishamchi 2005).

The objective of this research is to compare four statistical regression models that are appropriate for identification of pipe break patterns in a real pilot in a Tehran distribution network. The models used in this research were the linear model, exponential model, the Poisson generalized linear model, and the logistic generalized linear model, while the variables used include pipe diameter, pipe length, pipe age and pipe internal pressure.

## METHOD

This research has applied the four mentioned regression models to analyze the variables in pipe vulnerabilities. In regression, a search is made for an equation which can express the relationship between the variables and on the basis of which it is possible to make the necessary predictions or estimations (Scheidegger et al. 2015; Nishiyama & Filion 2013; Rausand & Arnljot, 2004). Explanations of the four models are given below.

### The linear regression model

In this model it is assumed that variable Y is a function of the descriptive variable $X_i$.

$$Y = \beta_0 + \sum_{i=1}^{n} \beta_i X_i + \varepsilon \tag{1}$$

in which $\beta_0$ and $\beta_i$ are the constants (regression parameters) that are estimated, $\varepsilon$ is the error value, assuming that errors of zero average and unknown variance have normal distribution and are independent (Montgomery et al., 2012).

The linear relation between the numbers of breaks in a given section of pipe and its age was first proposed by (Kettler & Goulter 1985) and is given in Equation (2).

$$NB = k_0 \times Age \tag{2}$$

in which NB represents the number of breaks per year at any pipe section, $k_0$ is the regression parameter and Age is the age of the pipe at the first break.

In the current research changes were made in the initial form of Equation (1), and parameters such as pipe material, pipe length, pipe internal pressure and pipe diameter were added; the initial equation included only time. The improved model is shown in Equation (3).

$$NB = \beta_0 + \beta_1(L) + \beta_2(P) + \beta_3(D) + \beta_4(Age) + \beta_5(AC)$$
$$+ \beta_6(DI) + \beta_7(GCI) + \beta_8(PE) + \beta_9(PVC) \tag{3}$$

The variables and the dimensions in all equation are as follows: NB: number of failures; D: pipe diameter (mm); Age: pipe age (years); P: pressure (atmosphere); L: pipe length (meters). The pipe materials are: DI: Ductile iron; AC: asbestos cement; GCI: Cast iron; PE: poly ethylene. It should be mentioned that the numerical value for pipe materials in all equations is either 0 or 1. For example, if a pipe is in asbestos cement (AC) material , the value is 1, but 0 for other materials.

### Exponential regression model

This is a type of non- linear regression model and the general formula is:

$$y = f(x, \beta) + \varepsilon \tag{4}$$

where y is the dependent variable, f(x, $\beta$) is the non-linear function with the parameters of $\beta 0$, $\beta 1$, … , and $\varepsilon$ is the amount of remaining error. Shamir & Howard (1979) applied non-linear regression analysis to find the exponential relation between the age of pipes and breaks. Their model is shown in Equation (5):

$$N(t) = N(t_0) \times e^{A(t+g)} \tag{5}$$

where; N(t): the number of breaks (NB) in the unit length per year; $N(t_0)$: the number of breaks in the unit length at the year of pipe installations; t: the time between the break in a year and the year of previous break; g: the pipe age at time t; and A: the coefficient rate of break $yr^{-1}$.

Some researchers developed Equation (5) so it was capable of considering all the variables in pipe failures (Andreou *et al.*, 1987; Tabesh *et al.* 2009). In this study the exponential regression model is not just dependent on age but also on all the variables according to Equation (11).

## Poisson generalized linear model

The Poisson model is commonly used for regression analysis of the failures in infrastructural systems (Agresti & Kateri, 2011; Guikema & Davidson (2006). This model is based on a function $\mu\ (\vec{x}_i, \vec{\beta})$; if $\vec{x}_i = [x_{1i}, x_{2i}, \dots, x_{ni}]$ is the vector of covariates from the $i^{th}$ section of the system (i = 1,2, ...,m), $\vec{\beta}$ is the vector of the regression coefficients, and $y_i$ is the independent variable according to Equations (6) and (7).

$$E\ (y_i|x_i) = \mu\ \left(\ \vec{x}_i\ ,\ \vec{\beta}_i\ \right) \qquad (6)$$

The Poisson generalized linear model is presented here in the form of Equation (7):

$$\begin{aligned} Log(\mu) = {}& \beta_0 + \beta_1(D) + \beta_2(AC) + \beta_3(GCI) + \beta_4(DI) \\ & + \beta_5(PE) + \beta_6(PVC) + \beta_7(L) + \beta_8(Age) + \beta_9(P) \end{aligned}$$
$$(7)$$

where Y represents the number of failures (NB), $\beta i$, the regression coefficients and the independent variables are as explained above.

## The logistic generalized linear mode

In most cases the main concern for water supply utilities is to check for occurrence/non-occurrence of failure in a pipe during a specific period of time and not the actual number of failures. The dependent variable is a binary variable, which at the time of pipe break during a specific period returns an answer of 1. It is not necessary in this model for the independent variable to have a normal distribution. The independent variables could have the value of

1 with the probability of $P$ or the value of 0 with the probability of $(1 - P)$.

The relation between the independent and dependent variables in the logistic model is not linear. A logistic regression function, which is the transformation of $P$ to its LOGIT is shown in Equation (8):

$$P = \frac{e^{\alpha+\beta 1 x 1 + \beta 2 x 2 + \dots + \beta i x i}}{1 + e^{(\alpha+\beta 1 x 1 + \beta 2 x 2 + \dots + \beta i x i)}} \qquad (8)$$

where $\alpha$ is the regression's constant parameter, $\beta_i$ are the regression coefficients for descriptive variables and $x_i$ are the independent variables. LOGIT is a form of substitute for this model, wherein the joint function is a LOGIT one.

The proposed generalized linear model is shown as Equation (9):

$$\begin{aligned} Logit[P(x)] = {}& Log\left[\frac{p(x)}{1 - p(x)}\right] \\ = {}& \alpha + \beta_0(D) + \beta_1(AC) + \beta_2(GCI) \\ & + \beta_3(DI) + \beta_4(PE) + \beta_5(PVC) \\ & + \beta_6(Age) + \beta_7(L) + \beta_8(P) \end{aligned}$$
$$(9)$$

where $P(x)$ is the break probability, $1 - P(x)$ is the no-break probability, $\alpha$ is the width from the beginning and the $\beta_i$ represent the estimated regression parameters.

## APPLYING THE MODELS IN A REAL PILOT

The selected pilot is district 1 of Tehran's north region with a total length of 582,702 m pipes, and 48,500 consumers (Figure 1). Given the large number of subscribers, this district has a high record of accidents. For example, during the period from July 2004 to December 2007, more than 65,000 cases were registered, of which more than 25,000 were registered alone in 2007 (IWWC 2017). Given the existing constraints in data collection, the parameters of pipe material, diameter, length, pressure and age were selected. The pipes in the study area are mainly of ductile (DI), asbestos cement (AC), cast iron (GCI), and poly ethylene (PE), having an average age of about 30 years.
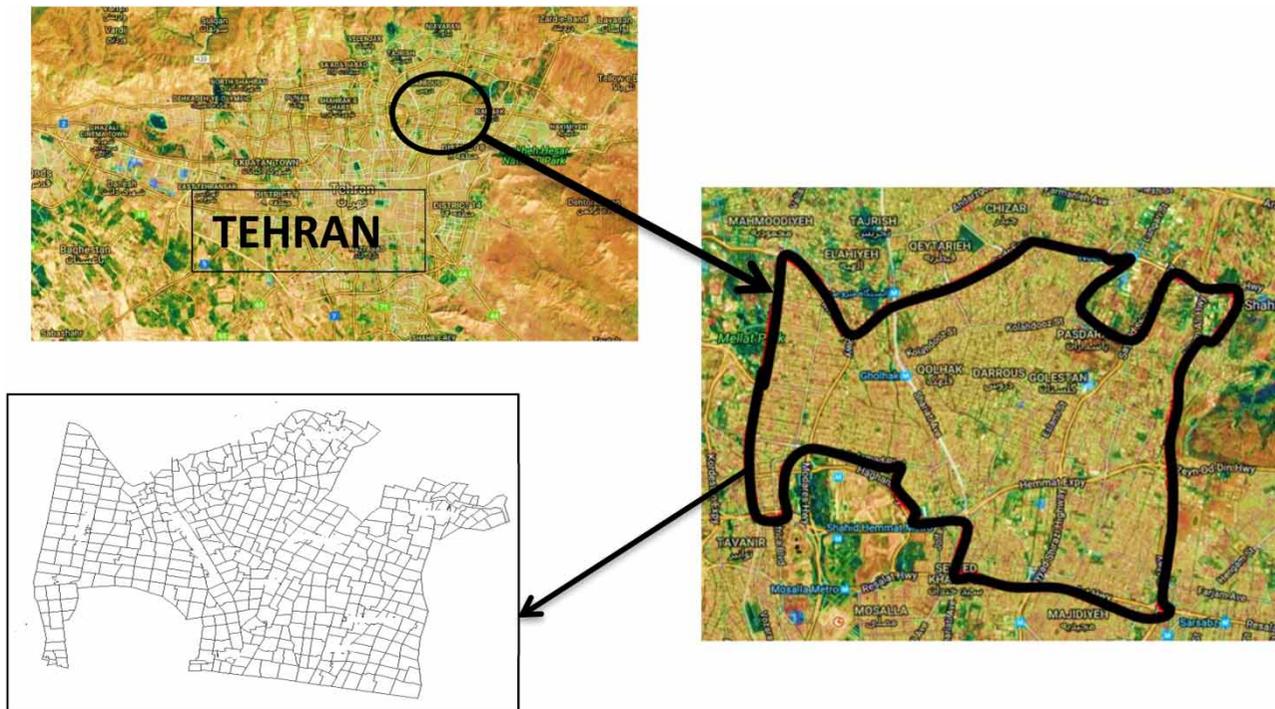
**Figure 1** | The scope of the pilot area, mainly district 1 of the Tehran water distribution network.

## RESULTS AND DISCUSSION

The four described regression models were applied to the gathered data of the study pilot to find the number and probability of pipe failures. The open source Statistical R software (Hothorn & Everitt 2014) was used for this study. It should be mentioned that for every model, different cases were analyzed. For example as shown in Table 1, six cases were studied for then linear regression model. The methodology of using the data for every case was as follows: 70% of the data randomly was used for training and the remaining 30% for testing and validation.

### Linear model results

This method was used according to *p*-value and the level of significance was defined as 0.05. Variables with a lower significance level (*p*-value > 0.1) were eliminated during the training of data in the next case. This process continued until all variables available in the model had a high significance level (*p*-value < 0.001).

A summary of the process stages is shown in Table 1. Finally the most appropriate linear model was the one shown in Equation (10):

$$NB = 1.024 + 0.00073L - 1.293\,DI \qquad (10)$$

The results of *p*-value of the pipe diameter have a very low significant level and therefore in case 2 this variable was eliminated from covariates. The highest *p*-value of parameters existing in Equation (10) was related to variables L and DI but lower for L.

Figure 2(a) shows the number of failures predicted versus their actual number per year.

### Exponential model results

The calibration of a non-linear regression model requires the initial values for the model's parameters (Montgomery *et al.*, 2012). The selection of the initial values was made through trial and error and then they were applied as the initial values of parameters. Two exponential cases were used in this research

**Table 1** │ Results of significance test of parameters in linear model

|  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 | Case 6 |
|---|---|---|---|---|---|---|
| Width from beginning | **** | **** | **** | **** | **** | **** |
| L | **** | **** | **** | **** | **** | **** |
| P | * | * | * | . | . | – |
| D | . | – | – | – | – | – |
| Age | . | . | – | – | – | – |
| AC | . | . | . | – | – | – |
| DI | **** | **** | **** | **** | **** | **** |
| GCI | . | . | . | . | – | – |
| PE | NA | – | – | – | – | – |
| R² | 0.1913 | 0.1906 | 0.1902 | 0.1891 | 0.1881 | 0.1856 |
| Degree of freedom | 771 | 772 | 773 | 774 | 775 | 776 |

'****' denotes that the *p*-value of the considered parameter is between 0 and 0.001.
'***' denotes that the *p*-value of the considered parameter is between 0.001 and 0.01.
'**' denotes that the *p*-value of the considered parameter is between 0.01 and 0.05.
'*' denotes that the *p*-value of the considered parameter is between 0.05 and 0.01.
'.' denotes that the *p*-value of the considered parameter is between 0.1 and 1.
'_' denotes that the considered parameter was not used in the model.

and the investigated exponential regression is shown in Equation (11).

$$NB = exp\,(0.00092L - 0.13P - 0.0049D$$
$$- 0.79\,AC - 0.895DI + 0.905GCI$$
$$+ 0.86(PE)) \tag{11}$$

According to the parameters estimated in Equation (11) for the existing variables, as the pipe length increases the number of failures also increases, while there is an inverse relation between the pipe diameter and number of failures. Furthermore, the use of ductile iron pipe significantly reduced the number of accidents.

Figure 2(b) is a chart of the number of predicted breaks versus the number of observed breaks. The MSE for zero and non-zero breaks in the analysis were 0.54 and 2.28, respectively, as shown in Table 2.

## Poisson generalized linear model results

As in the linear and exponential models, in this model the process began with the level of significance set on 0.05. The summary of the significance results of parameters of the Poisson generalized linear are shown in Table 2.
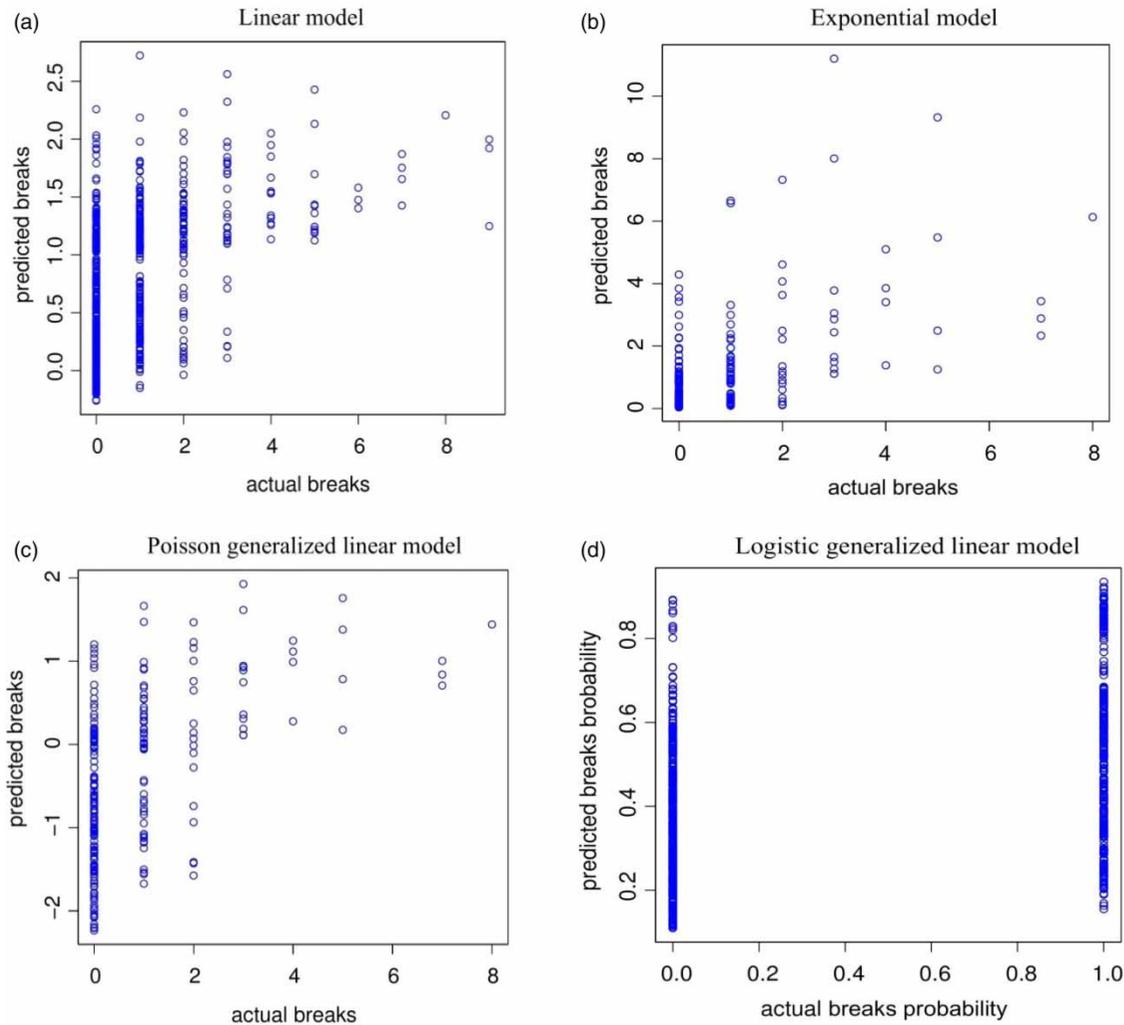
The resulting model is shown in Equation (12).

$$Log\mu = 0.7158 + 0.00089\,L - 0.144\,P - 0.0034\,D$$
$$- 1.585\,DI \tag{12}$$

By setting the level of significance on 0.05 the variables of pressure and ductile pipe material have the highest level of significance from a statistical aspect. In this model the number of failures increases with the increase in the length of the pipe and reduces with the increase in diameter. Moreover the use of ductile iron pipes lead to a reduction in the number of breaks. Table 2 shows the MSE for this model. Figure 2(c) shows the number of observed failures versus their predicted number.

The results showing in Table 2 demonstrates that all three models perform better in predicting the zero breaks than the non-zero ones.

By considering the log likelihood (LL), Deviance, AKAIKE Information Criterion (AIC) Akaike (1974), and Generalized Cross Validation (GCV), (Golub *et al.* 1979), the model with the highest value in (LL) and the lower value in (AIC), Deviance and (GCV) would be the better one. Based on these results the Poisson model is better

**Figure 2** │ Results analyses for linear (a), exponential (b), Poisson (c) and probability of failures in logistic generalized linear (d) models.

than two other models, while the quality of the linear and exponential models are very similar.

In general, the models used up to this point in the research are not suitable predictors for the number of non-zero breaks. To address the weakness of the previous models in predicting non-zero breaks, the logistic generalized linear model was used.

## Logistic generalized linear model results

The logistic model was tailored for working with binary data of 0 and 1. If, during a period of recording collected data, a pipe has experienced at least one break, its return variable would be 1 otherwise 0. The available data were calibrated with four logistic cases, amongst which, the

**Table 2** │ MSE, LL, AIC, deviance, GCV values for linear, exponential and poisson regression models

| Regression model | MSE | MSE (zero failure) | MSE (non-zero failure) | LL (Log likelihood) | AIC | Deviance | GCV |
|---|---|---|---|---|---|---|---|
| Linear | 1.41 | 0.65 | 2.37 | −1240.4 | 2488.8 | 1101.8 | 1.42 |
| Exponential | 1.3 | 0.54 | 2.28 | −1209.08 | 2434.16 | 1016.7 | 1.32 |
| Poisson | 1.3 | 0.55 | 2.28 | −900.75 | 1811.5 | 991.3 | 1.31 |

fourth proved to be the most appropriate. The results of calibration show that pipe length, pipe pressure, pipe age, ductile and cast iron variables are significant at the level of 0.05. Equation (13) presents the final model using the estimated parameters.

$$Logit\ [P(x)] = 0.0009L\ -\ 0.2P\ +\ 0.025Age - .76DI$$
$$+\ 1.38GCI \tag{13}$$

As can be observed, the number of failures in the model increases with an increase in the length and age of the pipe. Moreover, the use of ductile iron reduces the number of failures and by contrast, the use of cast iron pipes increases the number of breaks. The results of significance for the logistic general model are shown in Table 3. Based on Table 3, in case 4 the model's deviance and AIC are 936.6 and 946 respectively, which shows a better performance in comparison with other cases.

Figure 2(d) shows probability of prediction failures against probability of actual breaks. Consequently the value of 0 represents pipes without failures probability and the value of 1 with failures shows probability.

## CONCLUSION

Four different statistical models, namely the linear regression, exponential regression, Poisson generalized linear regression and logistic generalized regression were used to study probability of failures in water distribution pipes. The results obtained by the linear model shows that it is not suitable for modeling the reliability of a water distribution network system. The exponential model results showed that with an increase in the length of the pipe the number of failures increases as well, while it decreases with the increase of diameter.

The comparison of the goodness of fit criterion in the linear, exponential and Poisson models shows that the Poisson generalized linear regression model is better at predicting failures.

Overall, according to the results, the logistic model is more appropriate for estimating and predicting the probability of failures in the considered pilot distribution system over the

**Table 3** | Results of significance test of parameters in the logistic generalized linear model

|  | Case 1 | Case 2 | Case 3 | Case 4 |
|---|---|---|---|---|
| Width from beginning | . | – | – | – |
| L | **** | **** | **** | **** |
| P | ** | ** | ** | ** |
| D | . | . | . | – |
| Age | ** | *** | *** | *** |
| AC | . | . | – | – |
| DI | **** | **** | **** | **** |
| GCI | ** | ** | **** | **** |
| PE | NA | – | – | – |
| Deviance | 932.71 | 933.14 | 934.01 | 936.6 |
| AIC | 948.71 | 947.14 | 946 | 946 |

'****' denotes that the *p*-value of the considered parameter is between 0 and 0.001.
'***' denotes that the *p*-value of the considered parameter is between 0.001 and 0.01.
'**' denotes that the *p*-value of the considered parameter is between 0.01 and 0.05.
'*' denotes that the *p*-value of the considered parameter is between 0.05 and 0.01.
'.' denotes that the *p*-value of the considered parameter is between 0.1 and 1.
"_" denotes that the considered parameter was not used in the model.

three and half year data. In all the presented models, the ductile iron pipe variable was an important factor in reducing the number of failures in comparison with other independent variables. However, it should be mentioned that the results demonstrated in this paper have not included pressure data, since the available data were not sufficient.

## REFERENCES

Agresti, A. & Kateri, M. 2011 Categorical data analysis. In: M. Lovric (ed.), *International Encyclopedia of Statistical Science* (2011 edition). Springer, Berlin, Heidelberg, Germany. http://dx.doi.org/10.1007/978-3-642-04898-2_161.

Akaike, H. 1974 A new look at the statistical model identification. *IEEE Transactions on Automatic Control* **19** (6), 716–723.

Alvisi, S. & Franchini, M. 2010 Comparative analysis of two probabilistic pipe breakage models applied to a real water distribution system. *Civil Engineering and Environmental Systems* **27** (1), 1–22.

Andreou, S. A., Marks, D. H. & Clark, R. M. 1987 A new methodology for modelling break failure patterns in deteriorating water distribution systems: theory. *Advances in Water Resources* **10** (1), 2–10.

Beigi, F. 2000 Vulnerability in water distribution in Iran (in Persian). *Water and Environment Journal* **37**, 10–16.

Debón, A., Carrión, A., Cabrera, E. & Solano, H. 2010 Comparing risk of failure models in water supply networks using ROC curves. *Reliability Engineering & System Safety* **95** (1), 43–48.

Elahipanah, N. 1999 A review of water distribution networks in Iran till 2020 (in Persian). *Water and Environment Journal* **28**, 4–15.

Guikema, S. D. & Davidson, R. A. 2006 Modelling critical infrastructure reliability with generalized linear (mixed) models (PSAM-0009) In: *Proceedings of the Eighth International Conference on Probabilistic Safety Assessment & Management (PSAM)*. ASME Press, USA.

Golub, G. H., Heath, M. & Wahba, G. 1979 Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics* **21** (2), 215–223.

Grigg, N. S. 2012 *Water, Wastewater, and Stormwater Infrastructure Management*. CRC Press, California, USA.

Hothorn, T. & Everitt, B. S. 2014 *A Handbook of Statistical Analyses Using R*. CRC Press, California, USA.

Inanloo, B., Tansel, B., Shams, K., Jin, X. & Gan, A. 2016 A decision aid GIS-based risk assessment and vulnerability analysis approach for transportation and pipeline networks. *Safety Science* **84**, 57–66.

Iran Water and Wastewater Company (IWWC) 2017 Tehran Water and Wastewater Website. http://t1ww.tpww.ir/fa/p8/boarder (in Persian).

Kabir, G., Tesfamariam, S., Francisque, A. & Sadiq, R. 2015 Evaluating risk of water mains failure using a Bayesian belief network model. *European Journal of Operational Research* **240** (1), 220–234.

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.

Kettler, A. J. & Goulter, I. C. 1985 An analysis of pipe breakage in urban water distribution networks. *Canadian Journal of Civil Engineering* **12** (2), 286–293.

Kropp, I. & Herz, R. 2005 Schadensprosnosemodelle fur die zustandsbewertung von Leitungsnetzen der wasservesorgung. Wasserwirtschaft, Wassertechnik No. 5s. 10 ff, Berlin.

Le Gat, Y. & Eisenbeis, P. 2000 Using maintenance records to forecast failures in water networks. *Urban Water* **2** (3), 173–181.

Mailhot, A., Pelletier, G., Noël, J. F. & Villeneuve, J. P. 2000 Modeling the evolution of the structural state of water pipe networks with brief recorded pipe break histories: Methodology and application. *Water Resources Research* **36** (10), 3053–3062.

Montgomery, D. C., Peck, E. A. & Vining, G. G. 2012 *Introduction to Linear Regression Analysis (Vol. 821)*. John Wiley & Sons, New York, USA.

Nishiyama, M. & Filion, Y. 2013 Review of statistical water main break prediction models. *Canadian Journal of Civil Engineering* **40** (10), 972–979.

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.

Rausand, M. & Arnljot, H. Ã. 2004 *System Reliability Theory: Models, Statistical Methods, and Applications (Vol. 396)*. John Wiley & Sons, New York, USA.

Savic, D. A. 2009 The use of data-driven methodologies for prediction of water and wastewater asset failures. In: P. Hlavinek, C. Popovska, J. Marsalek, I. Mahrikova & T. Kukharchyk (eds), *Risk Management of Water Supply and Sanitation Systems*. Springer, Dordrecht, The Netherlands, pp. 181–190.

Scheidegger, A., Leitão, J. P. & Scholten, L. 2015 Statistical failure models for water distribution pipes–A review from a unified perspective. *Water Research* **83**, 237–247.

Shamir, U. & Howard, C. D. 1979 An analytic approach to scheduling pipe replacement. *Journal-American Water Works Association* **71** (5), 248–258.

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.

Tajrishi, M. & Abrishamchi, A. 2005 *Water Conservation, Reuse, and Recycling: Proceeding of Iranian-American Workshop-228-280*. National Academy of Science, Washington, DC, USA. https://www.nap.edu/read/11241/chapter/16#228.

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**, 1–10.

Xu, Q., Qiang, Z., Chen, Q., Liu, K. & Cao, N. 2018 A superposed model for the pipe failure assessment of water distribution networks and uncertainty analysis: a case study. *Water Resources Management* **32** (5), 1713–1723.