

Monetary consequence prediction of hazardous liquid pipelines in US

Khatami, Alireza

Research Assistant, Dept. of Civil, Construction & Environment Engineering, Marquette University, Milwaukee, U.S.

Qindan Huang

Associate Professor, Dept. of Civil, Construction & Environment Engineering, Marquette University, Milwaukee, U.S.

Kiswendsida J. Kere

Research Assistant, Dept. of Civil, Construction & Environment Engineering, Marquette University, Milwaukee, U.S.

ABSTRACT: The objective of this research is to develop monetary consequence prediction model for hazardous liquid pipelines, which can be used in the risk management of aging pipelines. In particular, Lasso regression and artificial neural network (ANN) are adopted with the use of the incident data reported by PHMSA. A feature selection procedure is proposed to pre-filter features that are given for the Lasso model development. The performance of the Lasso and ANN models are compared by using different statistical measures. The results show that the Lasso models overperform over ANN, and the proposed feature selection procedure is effective for Lasso model development to avoid overfitting. The case study is conducted to show the application of the developed failure consequence model.

1. INTRODUCTION

Safety transportation of essential energy commodities through pipeline is vital to ensure an uninterrupted serviceability to modern societies. In the US, around 3 million miles of pipelines have been extended all over the country. A pipeline failure at any point would spark off a big incident, which can jeopardize the safety of human lives, and adversely impact economics and natural environment.

Not many studies have been conducted to quantify the monetary consequence of pipeline failures based on historic data. If one could formulate the correlation between variables (e.g., pipeline properties, location) and failure consequence, it can be used in the decision making of pipeline management for inspection, repair, and replacement. Mostly, the predicted consequence values were obtained by typical assumptions or simply averaging reported historic data for a specific pipeline category. For example,

Nessim and Zhou (2011) estimated the pipeline failure cost by simply taking the average of damage costs over the historic incidents using PHMSA database of large leak and rupture incidents that occurred, and these estimated values were used by Gomez et al. (2014) to determine optimal inspection intervals of pipelines. Abubakirov et al. (2020) also simply assumed one value for cost consequence failure.

Meanwhile, regression has been adopted for quantifying pipeline failure consequence in the past. For example, Belvederesi and Dann (2017) attempted to use the pipe design variables (e.g., class location, diameter, installation year, maximum allowable operating pressure (MAOP)) in a linear regression to predict the damages of injuries, fatalities and property as the consequence of a pipeline incident; and they found less populated areas like lead to more cost of incidents compared to more populated areas, and older installed, larger diameter, or pipelines with higher

MAOP cause more expensive property damages. Simonoff et al. employed a two-step method to predict the cost consequence of a pipeline failure for hazardous liquid (2009) and natural gas transmission and distribution (2010) pipelines, where first the probability of occurrence of a non-zero consequence incidents is assessed using a logistic regression, and then, an ordinary least squares is used to find the property damage. However, their models were not tested against new data; more importantly, no predictor selection was used, resulted in complicated models in their studies.

In the past decade, machine learning has been widely used in engineering for prediction modeling. Neural network as a subfield of machine learning refers to the idea of using successive layers to learn the relationship between the observations by mimicking the way that biological neurons signal to one another in human brain. So far, only few studies (e.g., Parvizesdghy & Zayed (2015)) have adopted this approach for cost consequence prediction of pipeline. In addition, neural network model does not provide an explicit formula.

In summary, to predict cost consequence of pipelines incidents, a brief formula with a good accuracy is still missing, and it is also worthwhile to examine if adopting neural network could offer a better prediction. Therefore, this study adopts both Lasso regression to develop a brief and concise formula as well as a neural network method for a comparison purpose.

Lasso regression is adopted here, since it contains a feature selection process and it evaluates its performance on the prediction during training the model, which results in improved ability over linear regression on unseen data. It is found that the number of features initially provided to Lasso regression plays a role in the feature selection process within the Lasso, which may lead to different Lasso models. Lastly, the prediction performance of the Lasso models and the model based on the neural network model are compared.

2. RESEACH METHODOLOGY

2.1. Lasso Basics

In regression, a linear model is used to describe the linear relationship between response, y , and multiple features (or predictors, x_j) as shown:

$$y = \beta_0 + \sum_{j=1}^p \beta_j x_j + \sigma \varepsilon \quad (1)$$

where $\beta = \{\beta_0, \beta_1, \dots\}$ = unknown parameters, $\sigma \varepsilon$ = random error term that is independent of \mathbf{x} with mean of zero. Typically, in the least square regression, β are estimated by minimizing the sum of squared residuals (RSS), i.e.,

$$\min (RSS) = \min \left[\sum_{i=1}^n (y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{j,i})^2 \right] \quad (2)$$

where the hat refers to the estimated unknown parameters and i refers to the i th observation data set. In the least square, the unknown parameters are estimated solely based on the training data (that are the ones used in the minimization process to obtain β); thus, it does not guarantee a good prediction for the unseen data. In fact, the accuracy of a predictive model should be based on unseen data.

On the other hand, in Lasso, the unknown parameters are estimated by minimizing a sum of RSS with the consideration of a constrain on the sum of absolute values of coefficients, shown as:

$$\min(RSS); \text{ subjected to } \sum_{j=1}^p |\beta_j| < s \quad (3)$$

where s = tuning parameter. When s is large, then this constraint is not very restrictive and the impact value of s on the coefficient estimates is small. In fact, if s is set to be infinite, the coefficient estimates from Eq. (3) are the same as the least square estimation.

In Lasso, the optimal s is determined by utilization of cross validation on the training data (i.e., the data used for model development) when the mean squared error (MSE) reaches the minimum. During the cross validation, the training data are split into two different size subsets: one set for estimating β , and the other for evaluating the current model with the estimated β . This process is repeated multiple times by splitting the subsets randomly. With the optimal s , the final values of β are estimated using the whole training data. Since Lasso approach incorporates

the evaluation of a model over unseen data during the cross-validation process, it provides better accuracy (i.e., better prediction for the unseen data) compared with least square approach.

In addition, in Lasso, when a variable does not contribute to the prediction in unseen data, β_j is estimated to be zero. This is a nice feature because it automatically drops out unimportant variables. However, it is found that when too many unimportant variables are used initially in the model, Lasso is not able to remove all of them, which eventually results in overfitting.

In many cases, higher orders and interaction among features are usually considered as input predictors in regression. Thus, one needs to deal with a large number of predictors in the model development. To ensure Lasso produces an accurate and concise model, a feature selection is proposed in this study to pre-filter out some of the insignificant variables.

2.2. Proposed feature selection for Lasso

The following describes the proposed feature selection procedure to eliminate the insignificant predictors prior to Lasso analysis.

First, the impact of each predictor on target response is checked individually by null hypothesis in which the associated p -value is used to indicate the correlation between the predictor and response. In another word, if p -value of the variable is smaller than a threshold (e.g., $\alpha = 0.05$), it provides compelling evidence that the variable is related to the response.

In the second step, the selected predictors in the previous step are ordered by the R^2 of the model with each predictor alone. Then a forward selection is used, where a predictor is added one at time starting with the one with a higher R^2 . If the contribution of the added predictor results in an improvement of R^2 is above a certain value (e.g., $\Delta R^2 = 5\%$), then the added predictor will remain on the model; otherwise, it will not be added to the model. All the survived predictors will be the ones for the Lasso model development.

In summary, the number of predictors remained after this proposed procedure depends on the two thresholds (i.e., α and ΔR^2) used in the

two steps, respectively. If the criterion is set stricter (meaning a small value of α and a higher value of ΔR^2), less predictors will survive; otherwise, more predictors will be selected.

2.3. Neural Network

Typically, a neural network consists of an input layer, hidden layers, and one output layer to capture the relationship between the response to the inputs. Each layer takes the layer input matrix, \mathbf{x} , and return the output matrix, \mathbf{y}_{layer} , through the following transformation:

$$\mathbf{y}_{layer} = g(\mathbf{x} \cdot \mathbf{w} + \mathbf{b}) \quad (4)$$

where $g(\cdot)$ = an activation function to transform the dot product between \mathbf{x} , and the tensor weight of the layer (i.e., \mathbf{w}), with an addition a bias vector, \mathbf{b} . During the training process, \mathbf{w} and \mathbf{b} are adjusted to minimize a loss function value.

In this context, the loss function measures the difference between the actual (y) and predicted values (\hat{y}). In general, there are two types of loss functions for regression and classification. In a regression setting, two of the most popular loss functions are MSE and Mean Absolute Error (MAE). MSE calculated the average of squared differences between y and \hat{y} , while the MAE measures the average of the absolute distance between y and \hat{y} .

The activation function, $g(\cdot)$ in Eq. (4) decides a neuron should be activated based on its importance, and helps the network to learn complex patterns in the data through accounting for the interaction effects between variables and nonlinear behavior. It also has an impact on the convergence. Without using this function, the output model is a simple linear model. The typical non-linear activation functions include Sigmoid, Tanh and ReLU, which are different in terms of output range, and their derivatives. In particular, ReLU function can activate a certain number of neurons during learning and have non-zero derivatives (while Sigmoid and Tanh have close-to-zero gradients for values beyond -3 to 3), making it more computationally efficient. Here, the ReLU function shown below is adopted:

$$g(z) = \begin{cases} z & z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

3. APPLICATION

In US, the Pipeline and Hazardous Materials Safety Administration (PHMSA) requires operators to submit a report of the incident within 30 days for any of pipeline systems. To date, PHMSA has collected comprehensive information that includes the pipeline basic information (e.g., material properties, pipe geometries, installation year), year of incident, number of injuries, fatalities, incident cause. (U.S. Department of Transportation, 2022).

The PHMSA incident database indicates that equipment failure and corrosion are two of the most influential sources of failures in hazardous liquid (HL) pipelines. Thus, this research aims to develop a monetary consequence prediction model for corroded steel pipelines in HL pipelines. Due to the inconsistencies in the report format between old and newer datasets, only the data after 2010 are utilized in this study.

In general, the consequence reported includes property damage, fatality, and injury. In the PHMSA report, the property damage includes cost of public and non-operator private property damage, commodity lost, operator's property damage and repairs, emergency response, environmental remediation and other costs. Since the cost of emergency response and environmental remediation are heavily dependent on the specific pipeline environment and the "other costs" item is not specified in the PHMSA database, these three costs are thus excluded. Thus, the response used in the model development is the sum of cost of public and non-operator private property damage, commodity lost, and operator's property damage and repairs.

In the incident report, the following are used as the features for the model development: type of commodity released, pipeline geometry properties (e.g., diameter, wall thickness, depth of soil cover, etc.), pipeline material properties (e.g., age, material type, yielding strength), estimated volume of released commodity, location class,

incident ignition/explosion occurrence, and pipe pressure (including accident pressure, maximum allowable operating pressure). Noted that there are missing values in some features reported; thus, those incidents with missing values are eliminated for the analysis. As a result, the total number of incidents used for the analysis is 420.

4. ANALYSIS AND RESULTS

4.1. Model development

To understand the sensitivity of the final model to the number of the initial features provided to Lasso, two sets of predictors are generated using the proposed feature selection procedure described in Section 2.2. In particular, with setting $\alpha = 5\%$, we adopted two thresholds for ΔR^2 : 2% as a stricter criterion that results in 11 predictors and 5% as a looser criterion that results in 49 predictors.

The data is split into two: training dataset (that contains randomly 80% of the total data) and test dataset (the remaining 20%). With the application of Lasso, the number of predictors reduces to 11 and 29, respectively. Note that with the 11 initial predictors, Lasso does not further eliminate predictors, indicating that the proposed feature selection procedure does effectively select significant predictors. Meanwhile, when 49 initial predictors are given for Lasso, Lasso is only capable reduces the number of predictors to 29. The Lasso models achieved based on the 11 and 49 initial predictors are called "moderate" model and "flexible" model, respectively. As shown later in this Section, the moderator model outperforms the flexible model; thus, only the moderate model formulation is provided here:

$$\ln(y) = \beta_0 + \beta_1 \cdot \ln(x_1) + \beta_2 \cdot (\ln(x_1))^2 + \beta_3 \cdot \sqrt{x_2} + \beta_4 \cdot x_2 + \beta_5 \cdot x_2^{3/2} + \beta_6 \cdot \sqrt{x_3} + \beta_7 \cdot \ln(x_4) + \sigma \varepsilon \quad (6)$$

where y = predicted cost (in dollar), x_1 = volume of unintentional release of commodity (in number of barrels), x_2 = specified minimum yielding stress of the pipe (in psi), x_3 = maximum allowable pressure during operation (in psi), x_4 = pipe diameter (in inch), and σ = standard deviation of the model error, $\sigma \varepsilon$, and ε = standard normal

random variable. Note that the transformation used in Eq. (6) for all x_i are obtained through box-cox transformation so that their distributions are closer to a normal-like distributions to achieve a better prediction. Using the training data, the estimated predictor coefficients are summarized in Table 1 and Table 2. Note that area category and commodity type are also selected variables in the Lasso model, their impact are incorporated through β_0 as shown in Table 1.

Table 1. Statics of intercepts (β_0)

Area Category	Commodity type	μ	σ	<i>C.O.V</i>
High Population Area	Crude Oil	8.69	1.06	0.12
	HVL or Other Flammable Fluid	9.74	1.07	0.10
	Refined or Petroleum Products	9.47	1.07	0.11
Low Population Area	Crude Oil	7.83	1.08	0.13
	HVL or Other Flammable Fluid	8.88	1.09	0.12
	Refined or Petroleum Products	8.61	1.09	0.12

Table 2. Statistics of other predictor coefficients

Co-efficient	Mean (μ)	Standard Deviation (σ)	Coefficient of Variance
β_1	13.02	1.55	0.119
β_2	6.28	1.60	0.180
β_3	3.16	1.8	0.569
β_4	-0.38	1.54	-0.24
β_5	4.57	1.52	0.332
β_6	0.045	0.009	0.16
β_7	0.49	5	3.06
σ	1.54	-	-

When applying ANN, the training dataset used for Lasso is further split into 60% and 20% for the model training and fitting evaluation, respectively, and the rest of 20% data is used for testing. Here, MSE is adopted as the loss function.

As results, the best model achieved by the ANN contains two hidden layers with 256 and 128 neurons within first and second layer, respectively.

4.2. Performance comparison

The performances of the three developed models (i.e., Lasso moderate and flexible models, and ANN model) on the same testing data are compared in this Subsection.

Figure 1 indicates the scatter plots of the predicted and actual failure cost in natural logarithmic space. For a perfect prediction, the data should fall on the 45-degree line. As shown in Figure 1, the points predicted by all three models are evenly distributed around the 45-degree line, indicating all three models are unbiased. The spread for the moderate model is slightly smaller, indicating the prediction variance of the moderate model is smaller than the other two model.

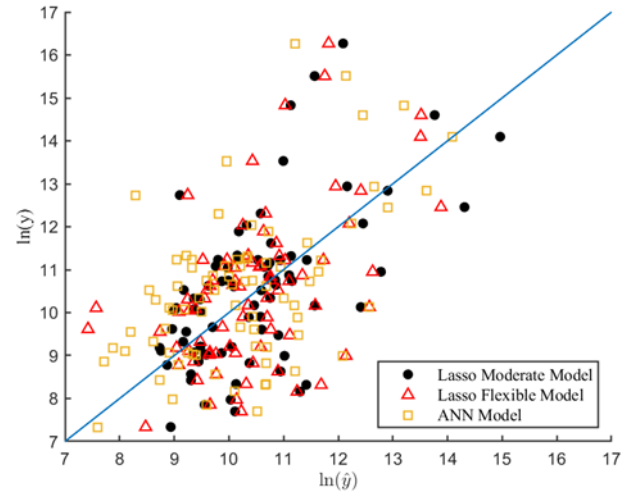


Figure 1. Predictions by three developed models vs. actual values of cost consequence in log space

Table 3 compares the three models in terms of four statistical measures: residual standard error (*RSE*), R^2 , *MAE*, and standard deviation of residual (σ), based on training and test data separately. In particular, *RSE* is used to show the average amount the predicted value deviates from the actual values, and is calculated by

$$RSE = \sqrt{RSS/(n - p - 1)} \quad (7)$$

where $RSS = \sum (y_i - \hat{y}_i)^2$, n = number of data, and p = number of variables. While *RSE* provides lack of fit of the model to the data, R^2 offers an

alternative measure to show the proportion of variance explained in the model, defined as:

$$R^2 = 1 - RSS/TSS \quad (8)$$

where $TSS = \sum(y_i - \bar{y}_i)^2$. The advantage of R^2 over RSE is that it always lies between the 0 and 1. However, it is challenging to determine what is good value for R^2 . In addition, MAE and σ are used to measure the average of absolute error and variability of the error, respectively.

As shown in Table 3, when using training data, flexible model and moderate model have approximately same values of the four statistical measures. However, these two Lasso models have much higher R^2 , and lower values of MAE and σ than the ANN model. These observations show that that the performance of both Lasso models is similar but they are much better than ANN model for the data used in the model development.

When using the testing data that is not used in the model development, the moderate model has smaller values of RSE , MAE and σ and a greater R^2 compared to the flexible model. This shows that the 29 variables used in the flexible model contain insignificant variables that only add error to the predictions, meaning this model overfits the training data. The better performance of the moderate model shows that the proposed feature selection process with a stricter criterion is effective for a Lasso model development.

In addition, compared with the moderate model, the ANN model still performs worse. The reason that ANN fails to produce a good model could be due to the limited training data (i.e., 420) available in this study. One could conclude that for limited observations, regression approach may have more advantages.

Table 3. Prediction accuracy comparison of three developed models

Model	Number of Variables	Training data				Test data			
		RSE	R^2	MAE	STD of Residuals, σ	RSE	R^2	MAE	STD of Residuals, σ
Flexible Lasso	29	0.77	0.42	1.11	1.43	0.88	0.27	1.31	1.68
Moderate Lasso	11	0.74	0.43	1.1	1.43	0.80	0.36	1.2	1.57
ANN	N. A	N. A	0.1	1.23	1.60	N. A	0.25	1.34	1.70

4.3. Performance of moderate Lasso model

Based on the results shown above, the moderate model has the best performance. Figure 2 shows the scatter plot of prediction by the moderate model vs. the actual values in original space using the 75-percentile of training and test datasets, where the dashed lines refer to 1 standard deviation of residuals from the perfect prediction (i.e., 45-degree line). While the variance is constant in the log space as shown in Figure 1, the prediction variation increases with the increase of the property damages as shown in Figure 2.

It is also worth comparing the developed moderate Lasso model with a linear regression model developed by Restrepo et al. (2009) for the failure consequence cost prediction of hazard liquid pipeline, where 37 variables were used and the model has R^2 and RSE of 0.358 and 2.031 in logarithmic space, respectively. Our

moderate Lasso model is able to achieve 0.43 of R^2 and 0.74 of RSE with only 11 variables using the training data. In addition, while the performance of the model by Restrepo et al. (2009) was not investigated on unseen data, the predictability of proposed model in this study is confirmed on the test dataset.

5. CASE STUDY

The developed moderate model is applied to a case study to investigate the property damage cost impact on life-cycle cost (LCC) estimation of steel pipelines exposed to corrosion. Herein, the expected total life cycle cost, $E[C_T]$, that excludes the initial construction cost, consists of inspection cost (C_{in}), repair cost (C_r) and cost of failure based on the estimated property damage by the moderate Lasso model. Note that the probability of failure is impacted by the repair action that is impacted by the time of inspection

and repair criteria. Kere and Huang (2023) describe the details of how the three components of costs are calculated based on a decision tree model to evaluate all possible events.

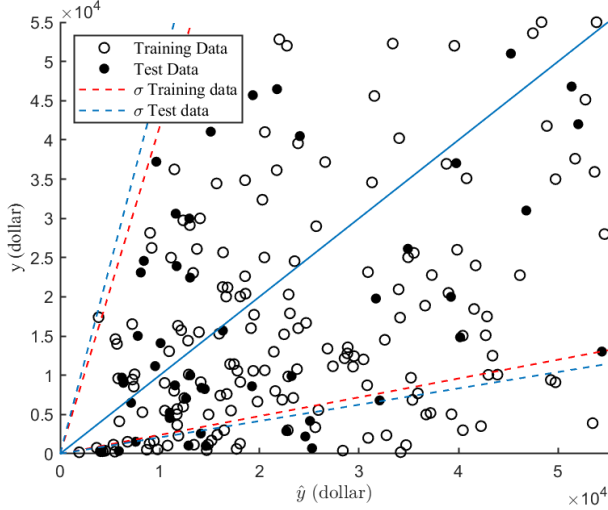


Figure 2. Prediction by moderate Lasso model vs. actual values in original space

In this case study, the failure refers to burst failure (that is when the pressure capacity, C , is less than the operating pressure, D). Thus, the probability of failure, P_f is the conditional probability of being in the failure domain given a set of boundary variables, written as:

$$P_f = \int_{C \leq D} f(\mathbf{X}) d\mathbf{X} \quad (9)$$

where $f(\mathbf{X})$ is the joint probability density function of a vector of random variables, \mathbf{X} . When considering corrosion in pipeline, the capacity weakens as corrosion worsens with time. In particular, the corrosion depth ($d(t)$) is calculated using a power law function, $d(t) = k \cdot t^n$, where k and n refer to defect growth parameters and values are assumed based on Li et al. (2019).

In addition, the unit cost of inspection and repair are assumed to be \$1000 and \$12,500, respectively (Zakikhani et al., 2020; Zhang & Zhou, 2014). Repair action is taken if the defect depth larger than a threshold value. When using Eq. (6) to evaluate the property damage, the variables values listed in Table 4 are used.

Figure 3 shows the calculated LCC of the steel pipeline subjected to corrosion with respect to various inspection interval, Δt , considering a

service life of 20 years and assuming a repair is conducted if the corrosion defect is found to be larger than 10% of wall thickness. In particular, the three curves of the expected LCC are calculated using three values of the property damage cost calculated based on Eq. (6): point estimate and point estimate $\pm \sigma$. The top curve (resulted from point estimate + σ) indicates that when $\Delta t = 7$ years, expected LCC reaches lowest, where the other two curves become flat when $\Delta t \geq 8$ years. The variation in the failure cost can significantly affect the decision-making regarding inspection intervals. Therefore, the accuracy of a cost consequence prediction model is critical in the risk management of deteriorating pipelines.

Table 4. Variable values used in Eq. (6)

Variable	Value
x_1	4914 barrels
x_2	42,000 psi
x_3	275 psi
x_4	35 inches
Commodity type	Crude oil
Population	High

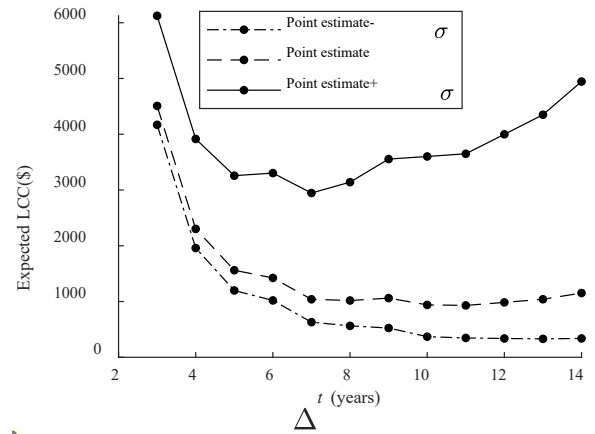


Figure 3. Expected LCC based on three values of property damage cost

6. CONCLUSIONS

This study aims to develop a prediction model of cost consequence of hazardous liquid steel pipelines failures due to corrosion based on incident data collected by PHMSA using two

approaches, Lasso regression and ANN. Compared with linear regression, Lasso provides better accurate forecasts over unseen data, and it has an ability in selection of important variables. However, too many initial input features given to Lasso still cause overfitting. Thus, a feature selection process is proposed to filter some of the insignificant variables for Lasso model development. An ANN model is also used for a comparison purpose. The findings are summarized as follows:

- The area population, type of commodity transported, volume of released commodity during the incident, yielding stress of the pipe, pipe diameter and maximum allowable operating pressure (MOP) are the selected features for property damage cost prediction.
- When using the incident data of hazard liquid pipeline, Lasso develops a more accurate model than ANN, which may be due to the limited data available for the model training.
- The proposed feature procedure as a pre-processing stage is critical for Lasso model development, which can effectively prevent overfitting.
- The results of a case study on the LCC of a corroding pipeline shows that the variation in the property damage cost prediction can lead to different optimal inspection intervals, confirming that the accuracy of a cost consequence prediction model is critical in the risk management of deteriorating pipelines.

7. REFERENCES

- Abubakirov, R., Yang, M., & Khakzad, N. (2020). A risk-based approach to determination of optimal inspection intervals for buried oil pipelines. *Process Safety and Environmental Protection*, 134, 95-107.
- Belvederesi, C., Dann, M. R., Copelli, S., Raboni, M., Ragazzi, M., Rada, E., Torretta, V., Zhijin, Y., Hu, W., & Pundt, H. (2017). Statistical analysis of failure consequences for oil and gas pipelines.
- Gomes, W. J., & Beck, A. T. (2014). Optimal inspection planning and repair under random crack propagation. *Engineering structures*, 69, 285-296.
- Kere, K.J., & Huang, Q. (2023). "Risk management strategies and probabilistic failure pressure model development for pipelines with crack-like defect," *14th International Conference on Applications of Statistics and Probability in Civil Engineering, ICASP14* (submitted).
- Li, C.-Q., Baji, H., & Yang, W. (2019). Optimal inspection plan for deteriorating structural members using stochastic models with application to buried pipelines. *Journal of Structural Engineering*, 145(11), 04019119.
- Parvzsedghy, L., & Zayed, T. (2015). Consequence of failure: Neurofuzzy-based prediction model for gas pipelines. *J. Perform. Constr. Facil.*, 30(4), 04015073.
- Restrepo, C. E., Simonoff, J. S., & Zimmerman, R. (2009). Causes, cost consequences, and risk implications of accidents in US hazardous liquid pipeline infrastructure. *International Journal of Critical Infrastructure Protection*, 2(1-2), 38-50.
- Simonoff, J. S., Restrepo, C. E., & Zimmerman, R. (2010). Risk management of cost consequences in natural gas transmission and distribution infrastructures. *Journal of Loss Prevention in the Process Industries*, 23(2), 269-279.
- U.S. Department of Transportation. (2022). *Pipeline and Hazardous Materials Safety Administration (PHMSA)*. <https://www.phmsa.dot.gov/>
- Zakikhani, K., Nasiri, F., & Zayed, T. (2020). Availability-based reliability-centered maintenance planning for gas transmission pipelines. *International Journal of Pressure Vessels and Piping*, 183, 104105.
- Zhang, S., & Zhou, W. (2014). Cost-based optimal maintenance decisions for corroding natural gas pipelines based on stochastic degradation models. *Engineering structures*, 74, 74-85.
- Zhou, W., & Nessim, M. A. (2011). Optimal Design of Onshore Natural Gas Pipelines. *Journal of Pressure Vessel Technology*, 133(3). <https://doi.org/10.1115/1.4002496>