

A Data Driven Machine Learning Model for Fault Detection in Industrial Pipelines

Namrata Vyas¹, Pallavi Bagde²
Department of CSE, SDBCT, Indore^{1,2}

Abstract: Industry 4.0 has marked a paradigm shift in the way production and manufacturing operates. Process monitoring and fault diagnosis are important for the safety and reliability of industrial processes especially for smart manufacturing. As a data-driven process monitoring methodology, multivariate statistical analysis techniques, such machine learning based approaches have become extremely critical for automation. Pipelines carrying oil and gas are essential for a nation's economic sustainability. In order to maximise their function and prevent product losses during the transportation of petroleum products, they must be carefully inspected. However, they are susceptible to failure, which could have negative effects on the environment, financial loss, and safety. Therefore, evaluating the pipe's state and quality would be crucial. Despite being time-consuming and expensive, a number of inspection procedures are used to assure the safety of pipelines. However, due to the time consumption and error prone nature of manual inspection, data driven models are being explored for forecasting failure in oil and gas pipelines. The proposed work presents a Bayesian Regularized classifier for forecasting failure and the results show that the proposed approach outperforms the existing baseline techniques in terms of forecasting accuracy.

Keywords:- Industry 4.0, oil pipelines; failure forecasting, Bayesian Regularization, Confusion Matrix, Classification Accuracy.

I. Introduction

Industry 4.0 has been defined as “a name for the current trend of automation and data exchange in manufacturing technologies, including cyber-physical systems, the Internet of things, cloud computing and cognitive computing and creating the smart factory. The advent Fourth Industrial

Revolution, Industry 4.0 conceptualizes rapid change to technology, industries, and societal patterns and processes in the 21st century due to increasing interconnectivity and smart automation. Coined

popularly by the World Economic Forum Founder and Executive Chairman, Klaus Schwab, it asserts that the changes seen are more than just improvements to efficiency, but express a significant shift in industrial capitalism.

The Four Industrial Revolutions

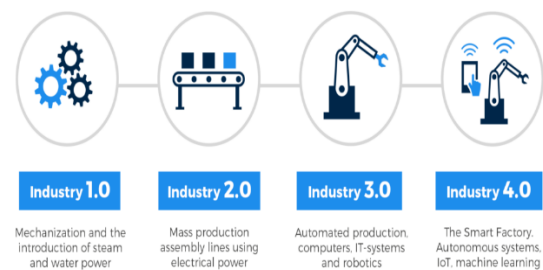


Fig.1. Concept of Industry 4.0

Industrial control systems are the mechanism (or brain) behind automated machine independence and motion. This is the technology that allows for industrial processes to be automated. The primary make-up of a control system is the control loop. Examples of control systems range from the very simplest, a discrete controller, to the complex SCADA system which manages all levels of a business's manufacturing processes and geographical locations.

II. Failures in Oil and Gas Pipelines

Petroleum products are transported by pipelines, the backbone of the oil and gas sector, in a range of locations (such as onshore or offshore) [1]. The first oil pipeline was built in Pennsylvania in 1879, and it had a diameter of 6 inches and a length of 109 miles [2]. In 120 nations around the world, more than 2 million miles of pipeline have been constructed. US pipelines account for 65% of the world's total length, with Canada and Russia following closely behind at 8% and 3%, respectively [3]. About 75% of the total length of the pipeline is shared by these three nations

[4]. There will be 491 operational oil pipelines in use by 2020 [5]. The Asia-Pacific region is home to over 46% (19,122 miles) of the world's oil and gas pipelines, while Canada is only expected to contribute 6%. External corrosion of insulated pipelines transporting hot products has been a major issue in the past, particularly in the 70s and 80s with several failures reported in any one year [6]. The problem was inherent to the design of these lines. Over time most such lines have been taken out of service (only 59 km remains today from a peak of over 1100 in the late 70s) and the issue disappeared with them, with only 2 cases recorded in the last 20 years [7].

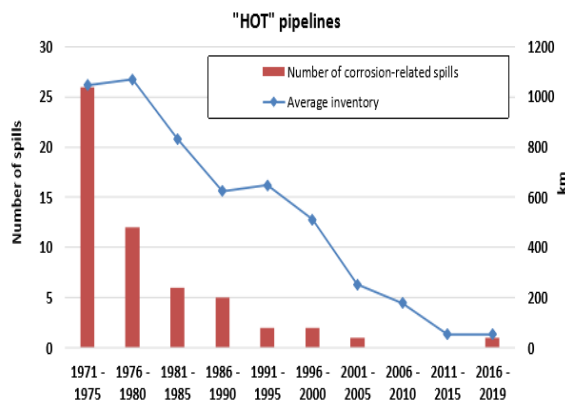


Fig.2 Number of recent hot pipe spills

Inspection techniques have been applied to discover pipeline anomalies and flaws without shutting down production [8]. In order to overcome the significant cost and time required by these inspection techniques, numerous studies have been undertaken to examine the condition, diagnose failure causes, and anticipate the residual lives of pipelines. Some failure prediction models were founded on subjective assessment, making them susceptible to different opinions [9]. Due to the size and complexity of the data being shared, it is almost infeasible for manual detection of faults and anomalies as it would consume large man hours and would also be less accurate [10]. Therefore, it has become mandatory to design automated systems which can detect faults/failures in very less time and with high accuracy. Since the data size to be analysed by time critical applications is enormous indeed, therefore the

conventional statistical techniques prove to be infeasible to detect fault detection with high level of accuracy, which primarily leads the focus to machine learning tools for the same.

III. System Design using Regression Learning based Bayesian Regularized ANN

Neural networks, with their remarkable ability to derive meaning from complicated or imprecise data, can be used to extract patterns and detect trends that are too complex to be noticed by either humans or other computer techniques. Other advantages include [11]:

1. **Adaptive learning:** An ability to learn how to do tasks based on the data given for training or initial experience.
2. **Self-Organization:** An ANN can create its own organization or representation of the information it receives during learning time.
3. **Real Time Operation:** ANN computations may be carried out in parallel, and special hardware devices are being designed and manufactured which take advantage of this capability.

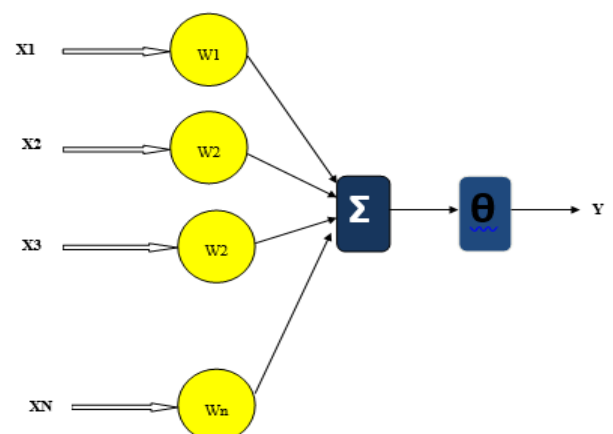


Fig.3 Mathematical Model of Neural Network

The output of the neural network is given by:

$$y = f(\sum_{i=1}^n X_i W_i) \quad (1)$$

Where,

X_i represents the signals arriving through various paths,

W_i represents the weight corresponding to the various paths and

f represents the activation function.

It can be seen that various signals traversing different paths have been assigned names X and each path has been assigned a weight W . The signal traversing a particular path gets multiplied by a corresponding weight W and finally the overall summation of the signals multiplied by the corresponding path weights reaches the neuron which reacts to it according to the bias Θ . Finally it's the bias that decides the activation function that is responsible for the decision taken upon by the neural network. The activation function ϕ is used to decide upon the final output. The learning capability of the ANN structure is based on the temporal learning capability governed by the relation [12]:

$$w(i) = f(i, e) \quad (2)$$

Here,

$w(i)$ represents the instantaneous weights

i is the iteration

e is the prediction error

The weight changes dynamically and is given by:

$$W_k \xrightarrow{e, i} W_{k+1} \quad (3)$$

Here,

W_k is the weight of the current iteration.

W_{k+1} is the weight of the subsequent iteration.

(i) Regression Learning Model

Regression learning has found several applications in supervised learning algorithms where the regression analysis among dependent and independent variables is needed [13]. Different regression models differ based on the kind of relationship between dependent and independent variables, they are considering and the number of independent variables being used. Regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a relationship between x (input) and y (output). Mathematically,

$$y = \theta_1 + \theta_2 x \quad (4)$$

Here,

x represents the state vector of input variables

y represents the state vector of output variable or variables.

Θ_1 and Θ_2 are the coefficients which try to fit the regression learning models output vector to the input vector.

By achieving the best-fit regression line, the model aims to predict y value such that the error difference between predicted value and true value is minimum. So, it is very important to update the θ_1 and θ_2 values, to reach the best value that minimize the error between predicted y value (pred) and true y value (y). The cost function J is mathematically defined as [14]:

$$J = \frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2 \quad (5)$$

Here,

n is the number of samples

y is the target

pred is the actual output.

(ii) Gradient Descent in Regression Learning

To update θ_1 and θ_2 values in order to reduce Cost function (minimizing MSE value) and achieving the best fit line the model uses Gradient Descent. The idea is to start with random θ_1 and θ_2 values and then iteratively updating the values, reaching minimum cost. The main aim is to minimize the cost function J .

(iii) Bayesian Regularization

The Bayesian Regularization (BR) algorithm is a modified version of the LM weight updating rule with an additional advantage of using the Bayes's theorem of conditional probability for a final classification.

The weight updating rule for the Bayesian Regularization is given by:

$$w_{k+1} = w_k - (J_k J_k^T + \mu I)^{-1} J_k^T e_k \quad (6)$$

Here,

w_{k+1} is weight of next iteration,

w_k is weight of present iteration
 J_k is the Jacobian Matrix
 J_k^T is Transpose of Jacobian Matrix
 e_k is error of Present Iteration
 μ is step size
 I is an identity matrix.

The decision making approach of the Bayesian Classifier can be understood graphically using the graph theory approach. The approach for computing the probability among different disjoint sets can be understood using the set theory approach shown in the subsequent steps. The figures clearly depict the decision to be taken in cases of different overlapping data value categories.

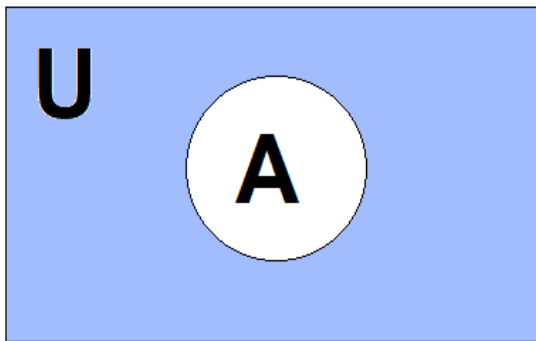


Fig.4 Universal Set Containing a Subset 'A'

Let us assume that the Bayesian Regularization algorithm needs to categorize the set A among multiple subsets in the superset U, for the time being in which only A exists exclusively.

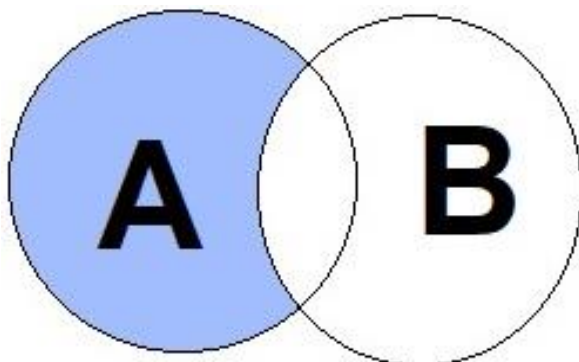


Fig.5 Probability of Exclusive Occurrence of 'A'

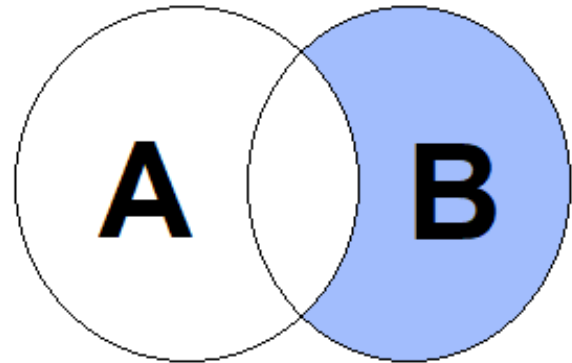


Fig.6 Probability of Exclusive Occurrence of 'B'

Figures above depict the probability of exclusive occurrence of events A and B respectively.

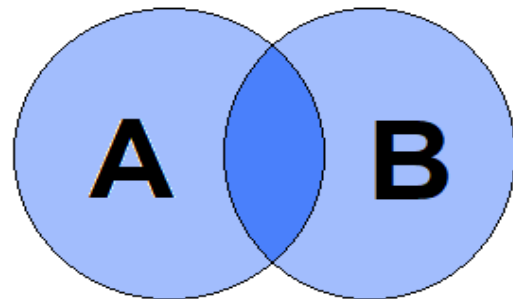


Fig.7 Probability of Union of A and B

Moreover for the predictive classification of ant data set, the Baye's Rule is followed, which is given by:

$$P \frac{A}{B} = \frac{P(A).P \frac{B}{A}}{P(B)} \quad (7)$$

Here,

$P \frac{A}{B}$ is the probability of occurrence of A given B is true.

$P \frac{B}{A}$ is the probability of occurrence of B given A is true.

$P(B)$ is the probability of occurrence of B

$P(A)$ is the probability of occurrence of A

In the present case the, 70% of the data has been taken for training and 30% of the data has been taken for testing.

The conditional probability of the sentiment can be also seen as an overlapping event with the classification occurring with the class with maximum conditional probability. The mathematical formulation for the above mentioned probabilistic approach can be understood as follows:

Let there be ‘N’ classes of data sets available in the sample space ‘U’.

Let the conditional probability of each of such sets be given by:

$$P(\frac{A}{H}), \quad P(\frac{B}{H}), \quad \dots\dots\dots P(\frac{N}{H}). \quad (8)$$

The BR algorithm tries to find out the maximum among the probabilities:

$$P(\max) = \begin{pmatrix} P(\frac{A}{U}) \\ P(\frac{B}{U}) \\ \vdots \\ P(\frac{N}{U}) \end{pmatrix} \quad (9)$$

The maximum value of the probability decides the classification of a dataset into a particular category. Assuming that X attains the maximum in such a sample space:

$$P_{max} = X \quad (10)$$

Here.

$$P\left(\frac{X}{U}\right) = P \frac{X}{\prod_{i=1}^n U_i}$$

$\prod_{i=1}^{i=n} U_i$ represents the conditional probability cumulative for all possible data set classes in the sample space U

X is the maximum probability corresponding to a particular data set and n is the total number of classes of categorization.

IV. Evaluation Parameters

Since errors can be both negative and positive in polarity, therefore its immaterial to consider errors with signs which may lead to cancellation and hence inaccurate evaluation of errors. Therefore we consider mean square error and mean absolute percentage errors for evaluation. The system accuracy can be evaluated in terms of the mean square error which is mathematically defined as:

$$mse = \frac{1}{n} \sum_{l=1}^N (X - X')^2 \quad (11)$$

Here,

X is the predicted value and

X' is the actual value and n is the number of samples.

The classification accuracy can be computed as:

$$Ac = \frac{TP+TN}{TP+TN+FP+FN} \quad (12)$$

Here,

- **TP:** True Positive
- **(TN):** True Negative
- **(FP):** False Positive
- **(FN):** False Negative

A high value of the accuracy indicates that the proposed algorithm is effective in the performance of the classification problem at hand.

V. Results:

The proposed data driven model needs the analysis of the data to be annotated initially in the categories of failure and non-failure of the pipeline. The data is randomly divided into 70% and 30% for training and testing, respectively.

[illegible]

Fig.8 Raw dataset

The data set parameters which are used for the training process are:

1. Pipeline location
2. Pipeline type
3. Oil Sub-type

4. Liquid name
5. Location in latitude and longitude
6. Cause category
7. Cause Sub-category
8. Intentional Release (in barrels)
9. Unintentional release (in barrels)
10. Liquid recovery (in barrels)
11. Liquid loss (in barrels)
12. Liquid Ignition
13. Liquid Explosion
14. Pipeline shutdown/failure (target)

- 2) The performance evaluation parameter is the mean square error
- 3) The training algorithm is the Bayesian Regularization algorithm

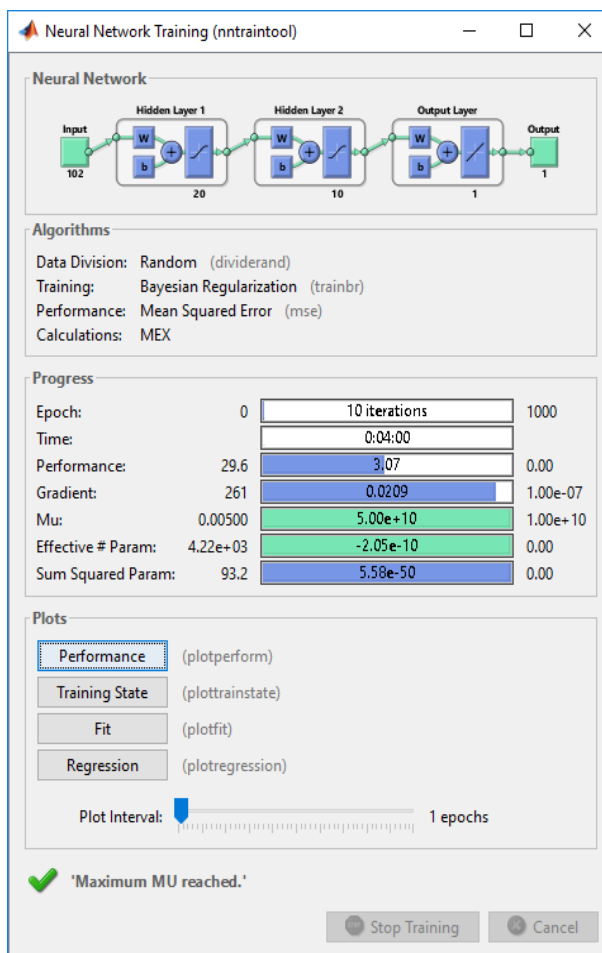


Fig.9 Designed Neural Network and its training parameters

The following attributes of the designed neural network are:

- 1) The number of hidden layers has been limited to 2 in order to limit the complexity of the algorithm.

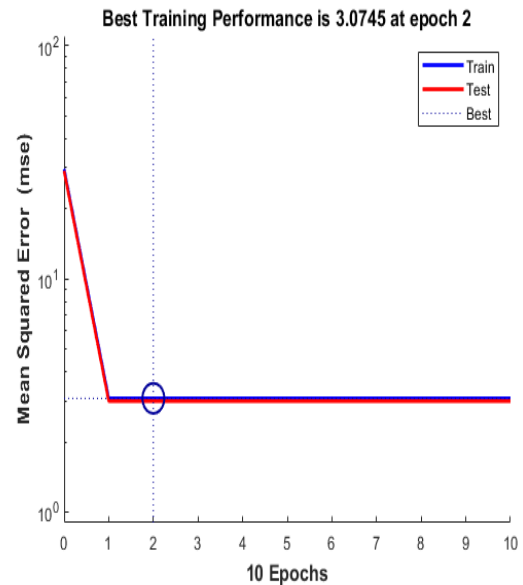


Fig.10 Training and epoch performance of the proposed system

The variation of the mean squared error as a function of the number of epochs is shown in the above figure. It can be seen that the MSE stabilizes at a value of 3.0745.

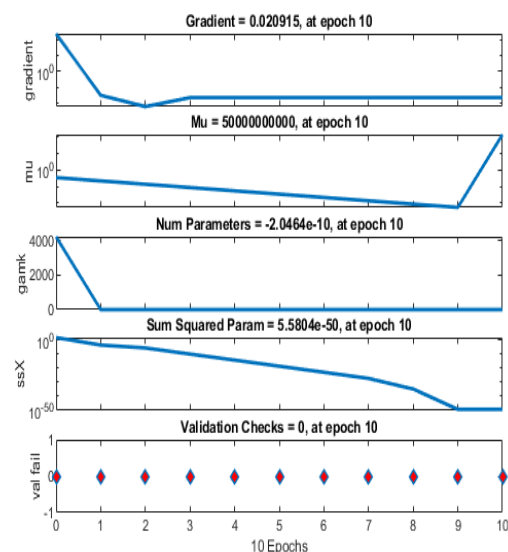


Fig.11 Training States as a function of number of epochs.

The variation in the training states such as the step size (μ), gradient, number of effective parameters and sum squared parameters has been shown in figure 9. The validation has also been shown. The gradient (g) and step size (μ) mathematically defined as:

$$g = \frac{\partial e}{\partial w} \quad (13)$$

Here,

e represents the error

w represents the weight

$$\mu = w_{k+1} - w_k \quad (14)$$

Here,

k represents the present iteration

$k+1$ represents the subsequent or next iteration

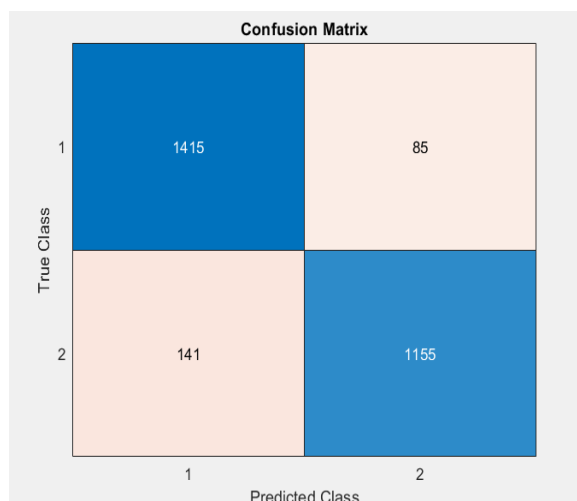


Fig.12 Confusion Matrix

Figure 12 depicts the confusion matrix for the proposed system depicting the TP, TN, FP and FN values respectively. The accuracy is computed as:

$$Ac = \frac{1415 + 1155}{1415 + 1155 + 85 + 141} = \frac{1545}{1796} = 91.91\%$$

It can be clearly seen that the proposed work attains much higher accuracy of 91.91% compared to 85% of previous work [1]. This can be attributed to the regression learning based BR trained ANN design which has a steep descent of error compared to the naïve Baye's classifier or the conventional Bayesian Regularization algorithm.

Conclusion:

It can be concluded from the previous discussions that pipelines carry petroleum products when compared to rail and roadways. However, pipelines are prone to various failures under diverse circumstances, leading to catastrophic environmental consequences owing to oil spilling as well as substantial economic losses due to production stoppage. The proposed approach uses a Regression Learning based Bayesian Regularized ANN for the forecasting of pipeline failure/shutdown. A set of governing variables have been used to train the BR model. The accuracy of the proposed system has been computed in terms of the true positive and true negative rates. It has been shown that the proposed work attains much higher accuracy of 91.91% compared to 85% of previous work, and thus outperforms the baseline approach in terms of classification accuracy.

References

- [1] Elshaboury N, Al-Sakkaf A, Alfalah D, Abdelkader E (2022). Data-Driven Models for Forecasting Failure Modes in Oil and Gas Pipes, Journal of Processes, MDPI, vol.10, no.2, pp.1-17.
- [2] Said, M., Abdellafou, K.b. & Taouali, O. (2020) Machine learning technique for data-driven fault detection of nonlinear processes. Journal of Intelligent Manufacturing, Springer, vol.31, pp.865–884.
- [3] Li, Y., Carabelli, S., Fadda, E. (2020) Machine learning and optimization for production rescheduling in Industry 4.0. International Journal of Advances in Manufacturing Technology, Springer, vol. 110, pp.2445–2463.
- [4] Carvajal J Soto F, Tavakolizadeh, T, (2019) An online machine learning framework for early detection of product failures in an Industry 4.0 context, International Journal of Production Research, Taylor and Francis, vol. 32, issue-4, pp. 452-465.
- [5] Iqbal R, Maniak T., Doctor F. Karyotis C., (2019) Fault Detection and Isolation in Industrial Processes Using Deep Learning Approaches, in IEEE

Transactions on Industrial Informatics, vol. 15, no. 5, pp. 3077-3084.

[6] Yu W, Dillon T, Mostafa F, Rahayu W, Liu Y, (2019) A Global Manufacturing Big Data Ecosystem for Fault Detection in Predictive Maintenance, IEEE Transactions on Industrial Informatics, vol. 16, no. 1, pp. 183-192.

[7] Paolanti M, Romeo L, Felicetti E, Mancini A, Frontoni E, Loncarski J (2018) Machine Learning approach for Predictive Maintenance in Industry 4.0, 2018 14th IEEE/ASME International Conference on Mechatronic and Embedded Systems and Applications (MESA), 2018, pp. 1-6.

[8] Wang J, Ma Y, Zhang L, Gao R, Wu D (2018) Deep learning for smart manufacturing: Methods and applications, Journal of Manufacturing systems, Elsevier, vol.48, pp.144-156.

[9] Sharp M, Hedberg R (2018) A survey of the advancing use and development of machine learning in smart manufacturing, Journal of Manufacturing systems, Elsevier, vol.48, pp.170-179.

[10] Lopez F, Saez M, Shao Y, (2017) Categorization of Anomalies in Smart Manufacturing Systems to Support the Selection of Detection Mechanisms, IEEE Robotics and Automation Letters, vol. 2, no. 4, pp. 1885-1892.

[11] Lin Y, Hung M, Huang H, Chen C, Yang H, (2017) Development of Advanced Manufacturing Cloud of Things (AMCoT)—A Smart Manufacturing Platform," in IEEE Robotics and Automation Letters, vol. 2, no. 3, pp. 1809-1816.

[12] Kang, H.S., Lee, J.Y., Choi, S. (2016) Smart manufacturing: Past research, present findings, and future directions. International Journal of Precision Engineering and Manufacturing-Green Technology, Springer, vol.3, pp.111–128.

[13] Zhou, Q., Wu, W., Liu, D.; Li, K., Qiao, Q (2016). Estimation of corrosion failure likelihood of oil and gas pipeline based on fuzzy logic approach. Engineering Failure Analysis, Elsevier, vol., pp.48–55.

[14] Zakikhani, K.; Nasiri, F.; Zayed, T. (2021) A failure prediction model for corrosion in gas transmission pipelines. Journal of Risk and Reliability, Institution of Mechanical Engineers (IMEchE), SAGE Publications. vol.235, pp.374–390