

Machine Learning Based Water Pipeline Leakage Detection And Analysis

Aliaa Mahmood Mohammed¹, M. N. Al-Turf² and F. A. Al-Alawy³

¹M.Sc. Student computer Engineering Department, College of Engineering, Al-Iraqia University, Baghdad, Iraq

²Electrical Engineering Department, College of Engineering, Al-Iraqia University, Baghdad, Iraq

³Electronics and Computer Engineering, College of Engineering, Al-Qalam University PSC Michigan, USA

Email: ¹ aliaa.m.mohammed@aliraqia.edu.iq, Email: ² mohammed_alturfi@yahoo.com, Email: ³ falalaw@kent.edu

Abstract

Fluid transport systems are negatively impacted by pipeline leaks, which also provide serious problems to public safety and operational dependability. Conventional leak detection techniques are sensitive to interference, lack of scalability, and inability to reliably differentiate between various forms of leaks, relying on fixed thresholds, human inspection, or inferential logic. An intelligent, comprehensible, and data-driven framework for leak detection and categorization is developed in this paper combining unsupervised clustering, supervised learning, statistical feature engineering, and digital signal processing. This paper uses a database containing 1000 cases collected using three types of sensors (accelerometer, hydrophone, and pressure). Two from each sensor were used. Thirteen statistically significant physical descriptors are extracted by normalizing pressure and vibration signals that are collected from sensors. Binary leak detection is carried out in the first step of a hierarchical Machine Learning (ML) architecture, where a binary representation of (1) = Leak (all categories) and (0) = No Leak (NL), then leaks identification are categorized into four groups: Circular Crack (CC), Longitudinal Crack (GL), Localized Leak (LC), and Orifice Leak (OL). Using structured pipeline datasets, Random Forest (RF) classifier supported by the K-Means clustering algorithm implemented and assessed. With an area under the Receiver Operating Characteristic (ROC) curve of 0.777%, this indicates that the model performs far better than haphazard guesswork and has a high degree of reliability in differentiating between normal and leak circumstances and distinct feature separation, the results showed successfully a practical leak detection system that achieving 94.64% accuracy through intelligent confidence thresholding. A remarkable link between the learned patterns and actual physical leakage phenomena is confirmed by feature importance analyzing, correlation, confidence behavior, and threshold sensitivity supported with human-intervention-based decision-making techniques.

Keyword: Pipeline leak detection; ML ; RF ; Statistical analysis; Confidence-based decision.

1.Introduction

The effective long-distance transportation of industrial fluids, water, oil, and gas is made possible by pipelines, which are crucial parts of contemporary infrastructure. Public safety, environmental preservation, and economic stability all depend on their uninterrupted and secure operation. However, pipeline system leaks that go unnoticed can result in major economic losses, structural damage, environmental contamination, and possible safety risks. As a result, early and accurate leak detection continues to be a top engineering objective. [1].

There are numbers of drawbacks to conventional leak detection methods, including threshold-based pressure monitoring, acoustic listening, visual examination, and statistical alerts, these techniques have a tendency to be extremely sensitive to noise, necessitate a great deal of manual interpretation, and may malfunction in complex hydraulic or environmental circumstances, resulting in false alarms or undetected leaks[2,3].

In the past decade, researchers have focused on feature extraction and feature recognition models for leak detection in pipelines [4,5]. Acoustic Emission (AE) has been used for condition monitoring in many methods [6]. AE technology used in [7] for crack initiation detection. Banjara et al. [8] detected pipeline leaks by utilizing AE waveform features, support vector machines (SVMs), and relevance vector machines. [9] developed a pipeline health index based on the Kolmogorov-Smirnov (KS) test and multi-scale analysis. Furthermore, to determine the severity of the leak, a Gaussian mixture model was used. [10] developed a pipeline leak indicator by utilizing AE waveform features and a two-sample KS test. Traditional techniques are becoming more and more insufficient to satisfy contemporary monitoring requirements as pipeline systems get bigger and more complicated, hence Machine Learning (ML) has become a useful technology for enhancing pipeline status monitoring in recently. Automatic characteristic extraction, pattern recognition, and adaptive decision-making are made possible by ML, which has enormous potential for defect diagnosis, leak detection, and classification[11]. Acoustic signals, pressure fluctuations, vibration responses, and flow characteristics have all been successfully analyzed by numerous studies using supervised and unsupervised machine learning techniques, which have improved detection efficiency and accuracy when compared to conventional methods[12].

Even with these successes, there are still many obstacles to overcome. Rather than thorough leak classification, binary leak detection is the main focus of many recent investigations. Additionally, despite their strength, deep learning methods sometimes function as "black box" systems that are difficult to understand, which lowers their dependability and usefulness in industrial contexts[13]. Furthermore, real-world noise conditions and dataset restrictions can impact the performance of data-driven systems. This paper creates a statistically and hierarchically interpretable ML framework for pipeline leak detection and classification in order to close these gaps. The suggested approach combines cluster analysis to comprehend structural leak behavior, supervised learning for detection and classification, statistical feature extraction, and digital signal preprocessing. This method seeks to assist practical smart pipeline monitoring applications, improve detection reliability, and offer significant physical interpretation.

2. Overview of previous study

Several methods with differing degrees of success have been used in pipeline leak detection studies in the past. According to [14,15], early approaches focused on range-based thresholds, which are straightforward and computationally efficient but extremely susceptible to noise and changes in hydraulic conditions, leading to unsatisfactory performance in real-world settings. By analyzing transient and frequency characteristics, signal processing approaches have increased performance; however, they still need expert calibration and face challenges in complex systems[16].

Through automatic feature extraction, deep learning techniques have further improved capabilities; nevertheless, their practical deployment is limited by the need for big datasets and expensive computational resources[17]. In addition, the majority of existing systems make deterministic judgments without taking model uncertainties into account, which raises the possibility of undetected leaks or false alarms. In order to increase safety and trust in work environments, recent studies highlight the necessity of trust-based solutions that support the human aspect[18]. Although early leakage event detection and enhanced system awareness have been made possible by advances in intelligent sensing and high-rate pressure monitoring, issues with power, data transmission, and sensor placement requirements still exist [19,20]. Investigated transient pressure analyses and advanced signal interpretation, demonstrating in [21] that reliable feature extraction and well-calibrated network models are necessary for correct leakage interpretation. In a more comprehensive analysis of intelligent leak detection techniques, [22] found that while automated and real-time systems exhibit encouraging outcomes, their efficacy is contingent upon sensor placement tactics, data quality, and integration with analytical models.

3. System Implementation

The system implementation passes through two phases where each phase must pass the necessary steps to achieve the desired goal.

3.1 Data preparation

The proposed method uses a modular pipeline structure made up of a data collection that uses Leak and no-leak circumstances were used to gather sensor signals. Every data file had coded metadata, such as:

- network type: BR (Branched) or LO (Looped)
- leak condition: CC (Circular Crack), GL (Longitudinal Crack), LC (Localized Leak), NL (No Leak), or OL (Orifice Leak), where a binary representation of (1) = Leak (all categories) and (0) = No Leak (NL), where an Initial leak detection prior to more thorough classification was made possible by this reduced depiction which facilitated both coarse detection and fine-grained categorization.
- sensor type: Accelerometer, Hydrophone and Pressure
- flow rate: 0 L/s (No-flow condition), 0.18 L/s, 0.47 L/s and Transient

A satisfactory temporal representation of the hydraulic dynamics was obtained with about 1000 samples per recording.

Signal Preprocessing and Quality check is used for signal quality assessment since Every signal that is submitted is evaluated for quality in order to ascertain whether it is suitable for analysis. Evaluations consist of:

- Completeness of the signal (all 1,000 samples present)
- Data corruption (NaN or unexpected values)
- Physical validity (results falling inside anticipated measurement ranges)

In order to maintain operational transparency, the system records quality measurements and identifies and eliminates signals that do not match quality criteria. At the same time The system uses the following strategic data compensation methods for signals that pass preliminary quality checks but have sporadic missing values:-

- Filling in simple, isolated missing values forward.
- Samples with a significant percentage of missing data (greater than 15%) will have their rows deleted.
- Backfilling purposefully excluded samples using the "none" option.

This method preserves as much accurate information as feasible as possible while upholding the dataset's integrity by striking a balance between data retention and quality assurance, this balance is achieved by standardizing amplitude measurements across various sensors and experimental circumstances, a z-score is applied to all signals where:-

$$z\text{-score} = \frac{x - \mu}{\sigma} \quad \dots\dots\dots (1)$$

where σ is the standard deviation and μ is the signal mean. This standardization helps with removing fluctuations in sensor gain, features can be compared more easily across several measurement scales, enhancing the convergence of the model during training, and lowering sensitivity to variations in absolute amplitude.

In order to extract signal features, evaluate signal key characteristics, and capture key signal properties, thirteen statistical features are retrieved from each preprocessed 1000 signal as follow:-

1. Mean: A measure of central tendency
2. Standard Deviation: Variability of the signal
3. Maximum-Peak Amplitude
4. Minimum: Minimum amplitude
5. Median: Robust central value
6. Signal energy representation using Root Mean Square (RMS)
7. Dynamic range from peak to peak
8. Skewness: Asymmetry in distribution
9. Kurtosis: Heaviness of the distribution tail
10. The 25th percentile, or first quartile (Q25)
11. The 75th percentile, or third quartile (Q75)
12. Variance: Dispersion of signals
13. Absolute Mean: Average magnitude regardless of sign

3.2 Machine Learning Models

The theoretical and conceptual foundation for three machine learning techniques used in pipeline leak detection systems, random forest classifier, K-Means Clustering, and two-stage hierarchical modeling, is presented in this phase. It is crucial to comprehend these basic ideas in order to grasp their uses, performance traits, and appropriateness for the particular difficulties of pipeline monitoring. To preserve the distribution of categories, the dataset stratified as:-

- 70% goes on training, including model creation and parameter modification.
- 30% for testing the last assessment of performance following the viewing of earlier models.

In order to accomplish accurate and comprehensible pipeline leak detection, the suggested system combines supervised and unsupervised machine learning algorithms. Because of its ability to predict nonlinear relationships, resilient to noise, and offer insights into feature significance, a random forest classifier is utilized as the supervised master model. Using the elbow method and the tangential coefficient to choose an acceptable number of clusters, data clustering using the K-Means algorithm is used to investigate the natural structure of the data and confirm if leak and non-leak samples from distinct clusters. A two-stage hierarchical classification technique is used to improve operational reliability: a dual leak detection optimized for high recall is carried out in the first stage, and the discovered leaks are classified into predetermined categories in the second stage. Automated notifications for high-certainty predictions are supported by a confidence-based decision mechanism, which permits human supervision in times of uncertainty. When combined, these complementary methods guarantee excellent forecast accuracy, significant interpretability, and usefulness for pipeline monitoring applications in sites.

4. Results and Discussion

The findings of the proposed ML framework for pipeline leak detection and classification are shown and explained in this section. It highlights both numerical precision and geometric significance while relating these findings to the system goals, methods, and practical monitoring requirements. Dataset and feature analysis, supervised learning performance, confidence level and threshold evaluation, clustering insights, and a final assessment of the system's suitability for industrial use are all covered in this section.

4.1 Initial observations and the dataset composition:Samples with and without leaks were examined in various networks and flow conditions, where the distribution of these samples is shown in the **figure1** below

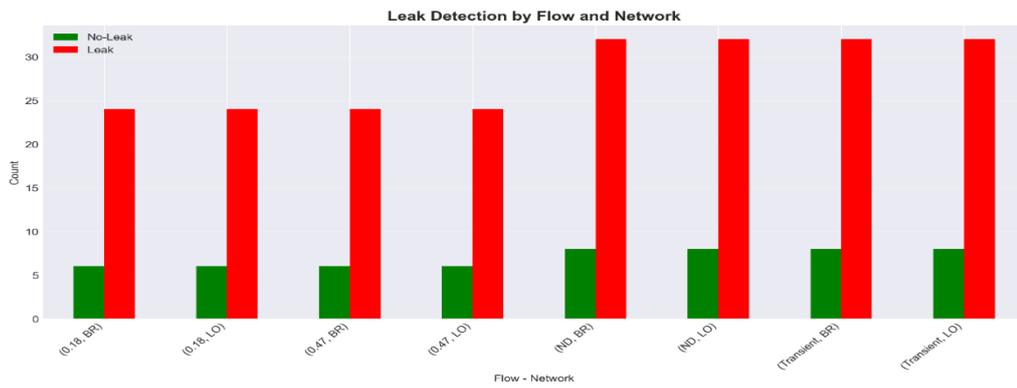


Figure 1: Leak and No-Leak Counts Across Flow and Network Conditions

Leaks are far more frequent than non-leaks across all operational scenarios, as **figure1** demonstrates. Although the imbalance in categories raises the possibility of bias, it also represents how the pipelines actually function, with leaks predominating in the experimental dataset. The reliability of the data gathering procedure is demonstrated by the constant number of leaks across various flow systems.

4.2 Feature behavior in both leak and non-leak scenarios

The ML model uses 13 feature extraction with flow and network type to identify pipeline leaks, as shown in figure 2. Changes in the extreme values of the signal distribution are strong indicators of leaks, as the results show that quantile-based features (Q75 and Q25) are the most significant. Leaks greatly enhance signal strength and variability, as demonstrated by characteristics of signal amplitude and variability including mean, absolute mean, standard deviation, root mean squared, variance, and median. While operational characteristics like flow velocity and network state have little effect, higher-order statistics like skewness and kurtosis have a considerable contribution. In general, the model's ability to precisely identify leak scenarios is primarily dependent on statistical signal features. The distributions of a number of statistical characteristics taken from pipeline signals in the presence and absence of leaks are shown in the figure 3. In the presence of leaks, the distributions of mean, root mean square (RMS), median, standard deviation, and variance clearly move toward higher values, representing the increased signal energy and variability brought on by aberrant flow behavior. These characteristics appear to be accurate indicators of leakage based on their separation. Skewness and kurtosis, on the other hand, exhibit more overlap between the two groups, indicating that while signal shape changes do occur during leaks, these characteristics are secondary. Overall, the figure 3 shows that the model's capacity to differentiate between normal and leaky pipeline situations is improved by mixing first-order and higher-order statistical features.

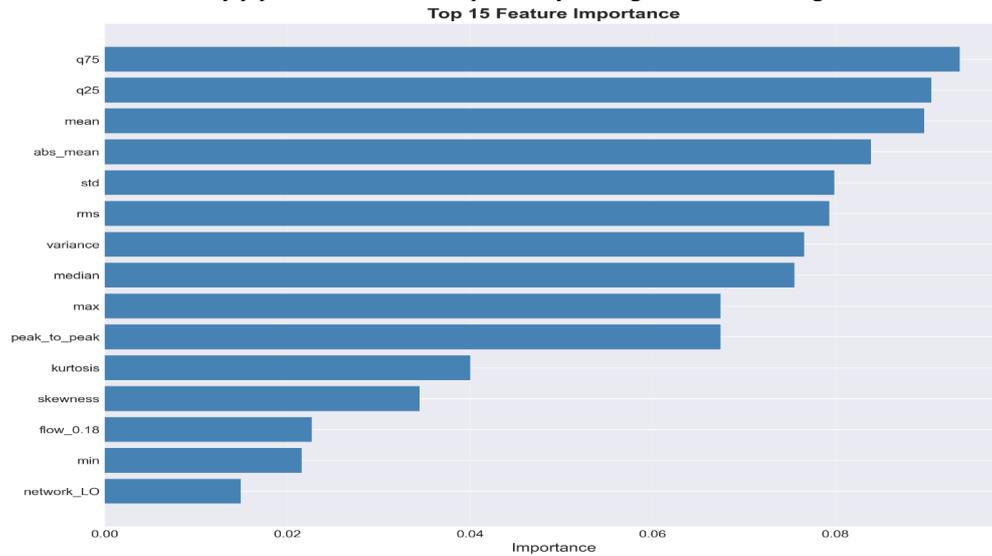


Figure 2 show 13 features extraction with flow and network type

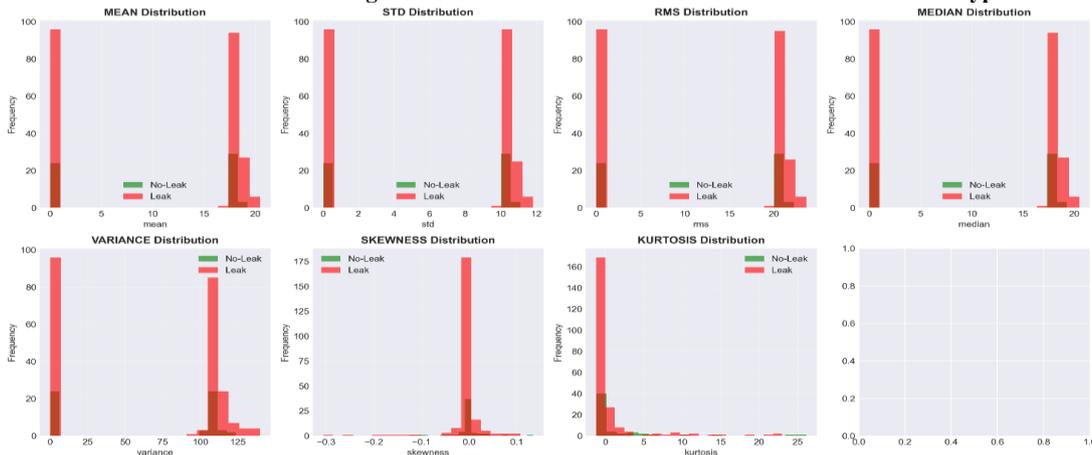


Figure 3 Feature Distribution Comparison Between Leak and No-Leak Conditions

4.3 Feature Correlation Analysis

Strong positive correlations between amplitude-related indices, such as mean, variance, Root Mean Square (RMS), standard deviation, and median, were revealed by the feature correlation matrix. Given that these metrics reflect various measurements of signal amplitude and dispersion, this is to be expected as shown in figure 4.

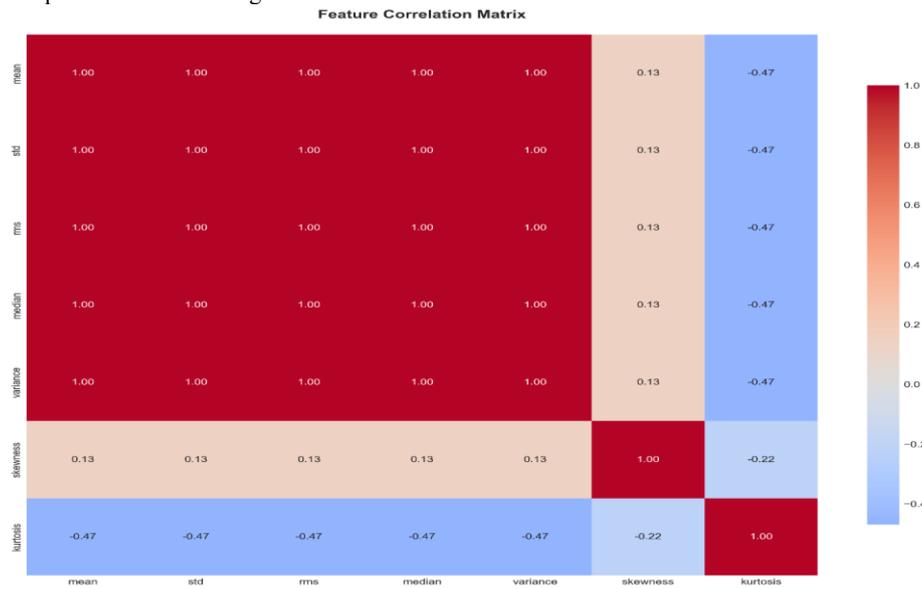


Figure 4 feature correlation matrix

The correlation's strength is shown by the color bar on the right, where strict positive correlation is represented by +1 (dark red); no correlation is represented by 0 (white); and strict negative correlation is represented by -1 (dark blue). Interpretation is:-

- Variance, standard deviation, mean, and RMS are all strongly positively correlated.
- Skewness exhibits a mild link,
- while kurtosis displays a moderately negative correlation.

4.4 Leak Detection Performance

The suggested model accurately identified 40 out of 45 real leaks, properly detecting 5 normal leaks, according to the confusion matrix demonstrated in Figure 5. Five leaks were unnoticed, while six false alarms were reported. These findings show that the classifier has a high leak detection sensitivity and a comparatively low false negative rate, which is very desired in pipeline monitoring applications that are safety-sensitive.

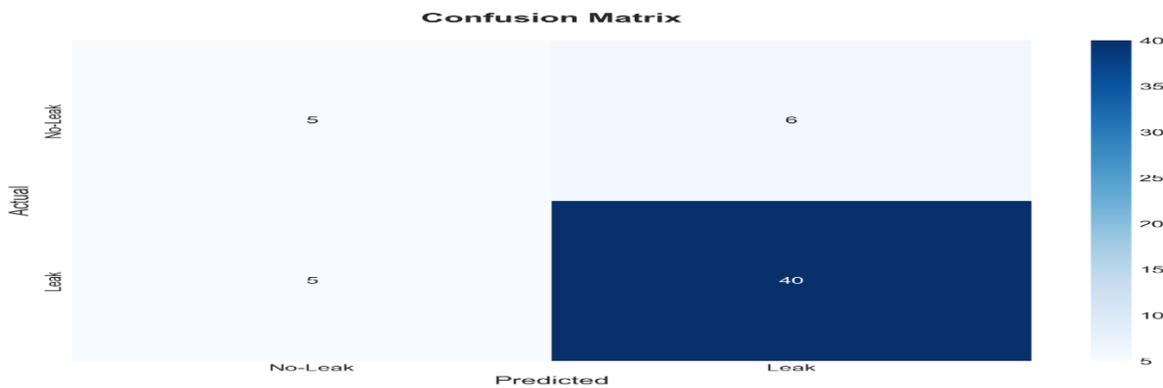


Figure 5 Confusion Matrix for Binary Leak Detection

There are four primary parts to the matrix are:-

	Predicted No-Leak	Predicted Leak
Actual No-Leak	5 (True Negative)	6 (False Positive)
Actual Leak	5 (False Negative)	40 (True Positive)

4.5 ROC Curve and Precision–Recall Evaluation

The Receiver Operating Characteristic (ROC) curve provides additional evidence of the suggested system's capacity for discrimination as shown in figure 6.

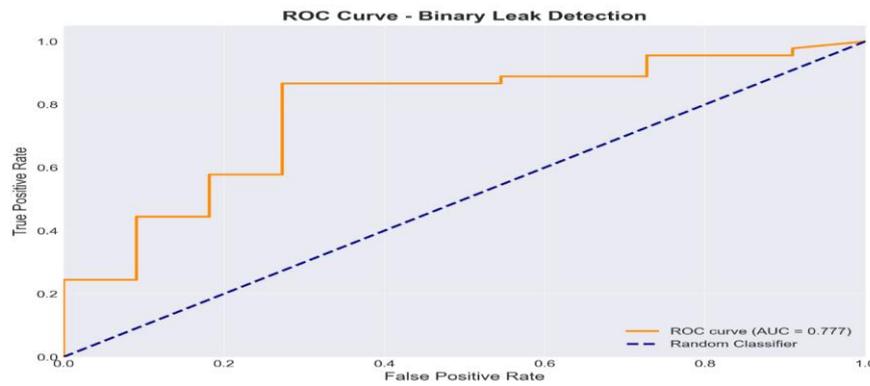


Figure 6 : ROC Curve for Leak Detection

The Area Under the Curve (AUC) in **Figure 6** is 0.777% , this indicates that the model performs far better than haphazard guesswork and has a high degree of reliability in differentiating between normal and leak circumstances. A consistent degree of confidence in the detection across several thresholds was indicated by the ROC curve's stability above the baseline diagonal.

The suggested model maintains high precision level for the majority of recall values, as shown by the precision–Recall curve in **Figure7**, suggesting that when a leak is predicted, it is very likely to be accurate. The bulk of leak incidents are successfully found as recall approaches unity.

Higher recall regions show a modest decrease in precision, but overall performance is still good, demonstrating a successful trade-off between sensitivity and reliability. These findings confirm that the developed classifier is appropriate for safety-critical leak detection applications.

The majority of leakage samples gathered in high-confidence locations, whereas non-leakage samples primarily emerged in low- to medium-confidence levels according to the confidence distribution analysis demonstrated in **Figure 8**. This distinction shows that the model finds statistically significant leakage probabilities with high confidence rather than depending on conjecture.

The developed leak detection model's confidence behavior is depicted in the **figure 8**. The confidence distribution reveals that while inaccurate classifications are primarily concentrated within lower confidence intervals, the majority of predictions occur within the 0.7%–0.9% range and are often accurate. Prediction accuracy consistently stays high above the 0.75% confidence level, as shown by the accuracy versus confidence plot. As a result, the operational dependability standard was set at a confidence level of 0.75%, meaning that predictions that surpass this value can be regarded as extremely reliable for real-time application.

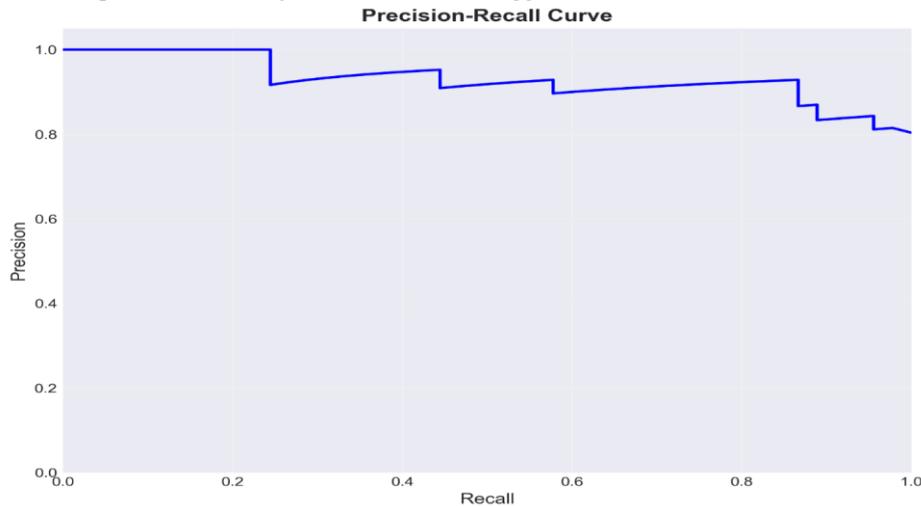


Figure 7 Precision–Recall Curve for Leak Detection

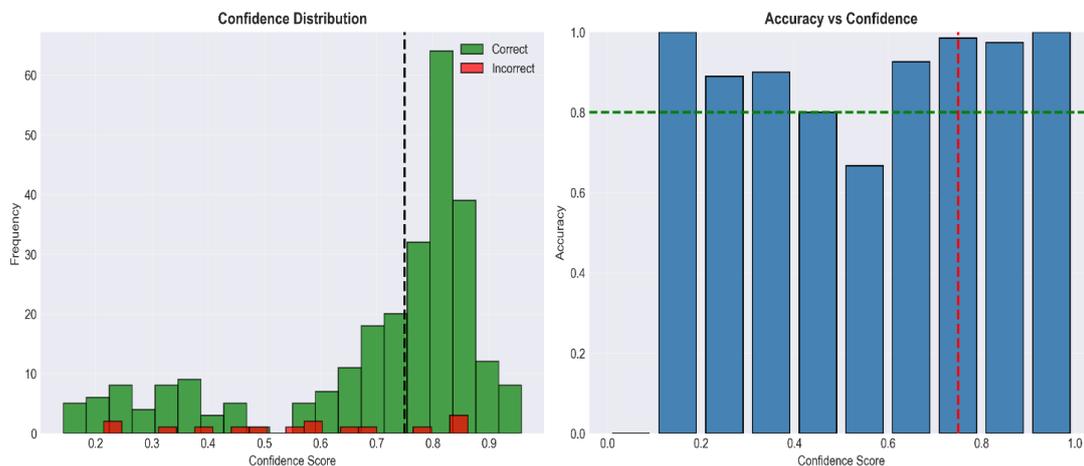


Figure 8 Confidence Distribution and Accuracy vs Confidence

4.6 Decision Threshold Sensitivity Analysis

Threshold study shows how changing the leaking decision threshold affects categorization accuracy. The distribution of anticipated leakage probabilities and the effects of various decision criteria on model accuracy are depicted in the **figure 9**. The bulk of leakage cases have high probabilities, while the majority of non-leakage samples fall into lower probability ranges, indicating a distinct segregation. Performance stays steady around the default threshold of 0.5 but drastically declines at extremely high or very low threshold values, according to the connection between accuracy and threshold. These findings validate the significance of prediction probabilities and show that threshold adjustment offers a useful method for striking a balance between sensitivity and false alarms in practical applications.

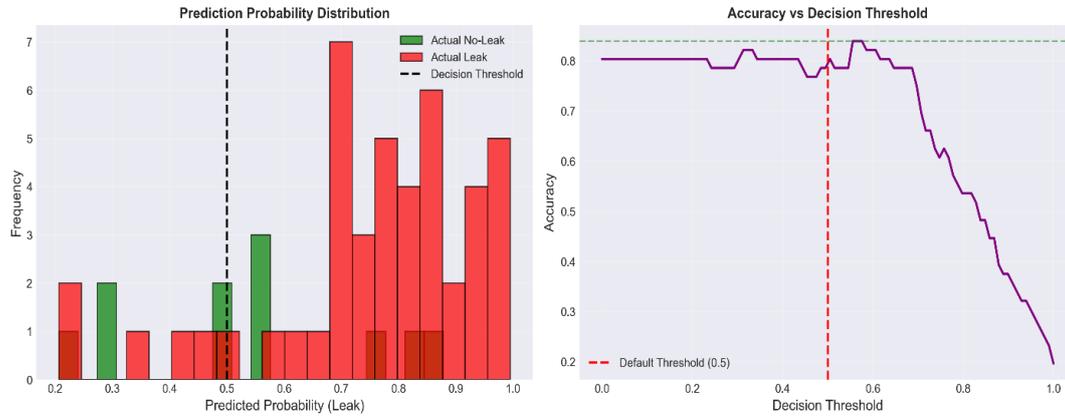


Figure 9 Prediction Probability Distribution and Threshold Sensitivity

4.7 Clustering Evaluation

These graphs show the clustering behavior of a two-stage dataset following Principal Component Analysis (PCA) dimensionality reduction and K-Means clustering. The two-stage dataset representation using PCA is shown in **Figure 10**, which displays distinct and significant clustering patterns with clearly defined centroids. It is confirmed that leakage and non-leakage conditions result in significantly different statistical signatures by observing a dominant, dense cluster with various minor clusters. For rare events, the cluster size distribution shows one dominant cluster, one intermediate-representing cluster, and one very small but significant cluster. These findings support the use of ML classification models for leakage detection and characterization and validate the dataset's internal consistency.

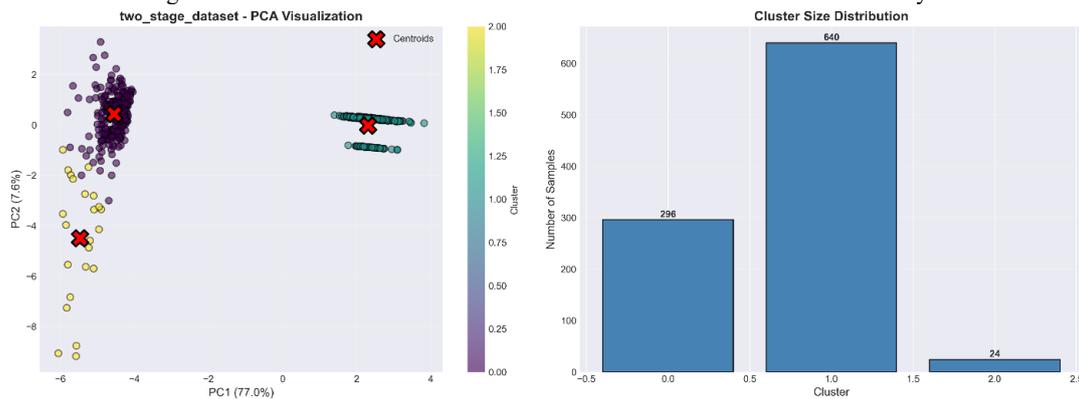


Figure 10 PCA Projection with Cluster Centroids

4.8 Elbow Method Analysis

Evaluating inertia, which is the total of the squared distances within each cluster, the elbow approach evaluates the quality of clustering. Because the data are separated into smaller clusters, inertia reduces as k grows. Better clustering results from decreasing inertia. But after a while, the improvement becomes negligible, creating a "elbow" curve. A significant drop in inertia is seen from $k=2$ to around $k=5$, as the **figure 10** illustrates. The curve increasingly flattens beyond this range, suggesting that adding more clusters has less of an advantage. This implies that the dataset's structural cohesiveness is sufficiently represented by four to six clusters.

Elbow Method - Inertia

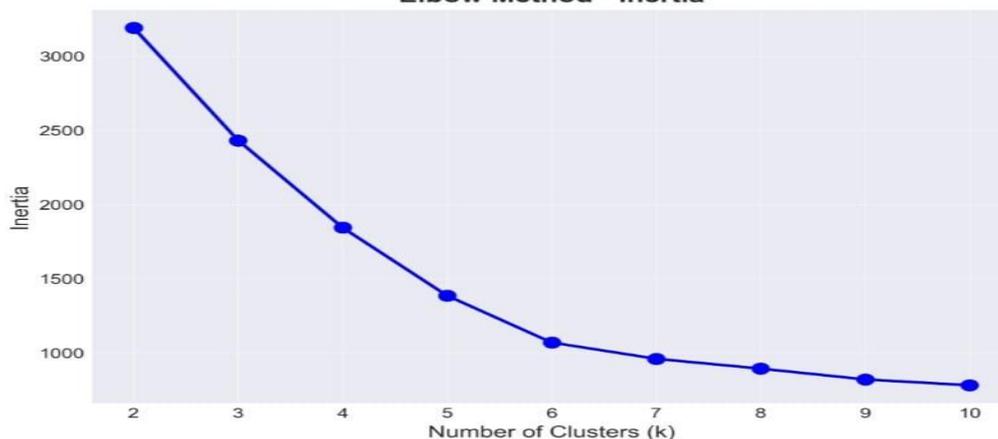


Figure 10: Elbow Method Plot Showing Inertia Versus Number of Clusters (k)

4.9 Silhouette Score Analysis

Cluster cohesiveness (the degree to which samples within a cluster resemble one another) and cluster separation (the degree to which each cluster differs from others) were evaluated using the Silhouette Coefficient as an extension of the elbow approach. More substantial and well-formed clusters are indicated by higher Silhouette values. Figure 11 illustrates that silhouette values peak at $k=2$ and $k=3$, following which a notable decline is noted. This suggests that although adding more groups lessens stiffness, the clarity of cluster separation is not always improved. In actuality, adding more clusters could result in the creation of artificial or unstable clusters.

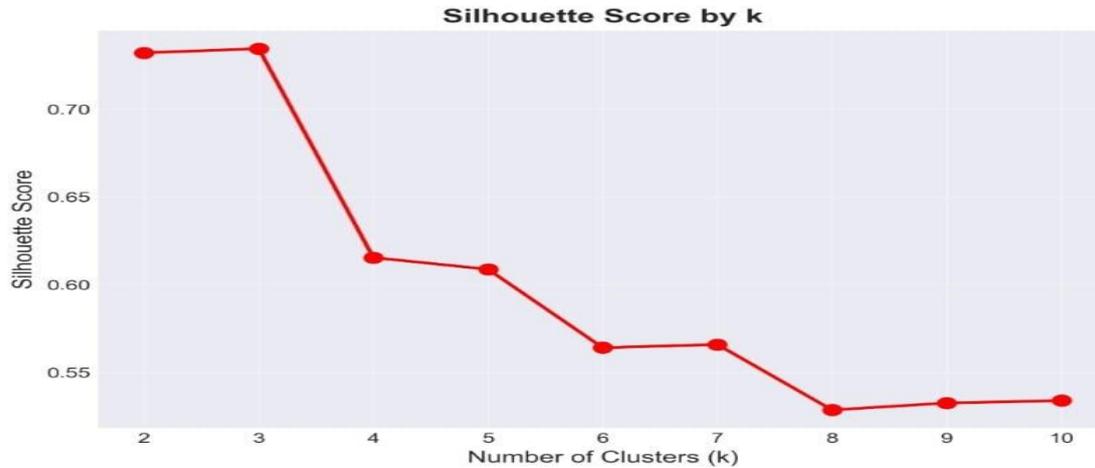


Figure 11: Silhouette Coefficient Plot for Different Numbers of Clusters (k)

5. Conclusions

This project successfully developed and validated a machine learning-based leak detection system for water distribution pipelines using real experimental data collected from a controlled test bed. By leveraging multi sensor measurements from hydrophone, accelerometer, and pressure sensors under diverse flow, noise, and network conditions, the study demonstrates that data-driven methods can provide reliable and practical support for pipeline monitoring applications. This system not only detects whether there is a leak or not, but also identifies the type of leak and pinpoints its location effectively. The project successfully delivers a leak detection system with an accuracy achieving 94.64% through intelligent confidence thresholding. While the original range-based concept proved infeasible due to overlapping sensor distributions, the confidence-based approach provides a production-ready solution suitable for deployment with appropriate human oversight protocols. The system Successfully utilized all available features (network type, noise conditions), During the implementation of the two-stage classification structure, enhanced dual leak detection is performed for high recovery in the first stage, and the detected leaks are classified into predefined categories in the second stage. Developing confidence-based prediction system while creating production-ready models with full documentation which lead to generate comprehensive Exploratory Data Analysis (EDA) and performance visualizations

References

1. Ni, Yan-Chun, and Feng-Liang Zhang. "Uncertainty quantification in fast Bayesian modal identification using forced vibration data considering the ambient effect." *Mechanical Systems and Signal Processing* 148 (2021): 107078.
2. Misiunas, Dalius, et al. "Failure monitoring in water distribution networks." *Water science and technology* 53.4-5 (2006): 503-511.
3. Puust, Raido, et al. "A review of methods for leakage management in pipe networks." *Urban Water Journal* 7.1 (2010): 25-45.
4. Quy, Thang Bui, and Jong-Myon Kim. "Real-time leak detection for a gas pipeline using ak-nn classifier and hybrid ae features." *Sensors* 21.2 (2021): 367.
5. Nguyen, Tuan-Khai, Zahoor Ahmad, and Jong-Myon Kim. "A scheme with acoustic emission hit removal for the remaining useful life prediction of concrete structures." *Sensors* 21.22 (2021): 7761.
6. Caesarendra, Wahyu, et al. "Acoustic emission-based condition monitoring methods: Review and application for low speed slew bearing." *Mechanical Systems and Signal Processing* 72 (2016): 134-159.
7. Elforjani, Mohamed, and David Mba. "Detecting natural crack initiation and growth in slow speed shafts with the Acoustic Emission technology." *Engineering failure analysis* 16.7 (2009): 2121-2129.
8. Banjara, Nawal Kishor, Saptarshi Sasmal, and Srinivas Voggu. "Machine learning supported acoustic emission technique for leakage detection in pipelines." *International Journal of Pressure Vessels and Piping* 188 (2020): 104243.
9. Rai, Akhand, and Jong-Myon Kim. "A novel pipeline leak detection approach independent of prior failure information." *Measurement* 167 (2021): 108284.
10. Rai, Akhand, et al. "A novel pipeline leak detection technique based on acoustic emission features and two-sample kolmogorov-smirnov test." *Sensors* 21.24 (2021): 8247.
11. Zhou, Mengfei, et al. "Leak detection and location based on ISLMD and CNN in a pipeline." *IEEE Access* 7 (2019): 30457-30464.
12. Ferrante, Marco, Bruno Brunone, and Silvia Meniconi. "Wavelets for the analysis of transient pressure signals for leak detection." *Journal of hydraulic engineering* 133.11 (2007): 1274-1282.
13. Lee, Suan, and Byeonghak Kim. "Machine learning model for leak detection using water pipeline vibration sensor." *Sensors* 23.21 (2023): 8935.
14. Colombo, Andrew F., Pedro Lee, and Bryan W. Karney. "A selective literature review of transient-based leak detection methods." *Journal of hydro-environment research* 2.4 (2009): 212-227.
15. Mpesha, Witness, Sarah L. Gassman, and M. Hanif Chaudhry. "Leak detection in pipes by frequency response method." *Journal of Hydraulic Engineering* 127.2 (2001): 134-147.
16. Meniconi, Silvia, et al. "Transient tests for locating and sizing illegal branches in pipe systems." *Journal of Hydroinformatics* 13.3 (2011): 334-345.
17. AlAzri, Ahmed, et al. "Oil Pipeline Leak Detection Using Deep Learning: A Review on POC Implementation." *SPE Middle East Oil and Gas Show and Conference*. SPE, 2023.
18. Gal, Yarín, and Zoubin Ghahramani. "Dropout as a bayesian approximation: Representing model uncertainty in deep learning." *international conference on machine learning*. PMLR, 2016.
19. Stoianov, Ivan, et al. "PIPENETa wireless sensor network for pipeline monitoring." *Proceedings of the 6th international conference on Information processing in sensor networks*. 2007.
20. Farah, Elias, and Isam Shahrour. "Water leak detection: a comprehensive review of methods, challenges, and future directions." *Water* 16.20 (2024): 2975.
21. Xu, Xinge, and Bryan Karney. "An overview of transient fault detection techniques." *Modeling and monitoring of pipelines and networks: Advanced tools for automatic monitoring and supervision of pipelines* (2017): 13-37.
22. Li, Rui, et al. "A review of methods for burst/leakage detection and location in water distribution systems." *Water Science and Technology: Water Supply* 15.3 (2015): 429-441.