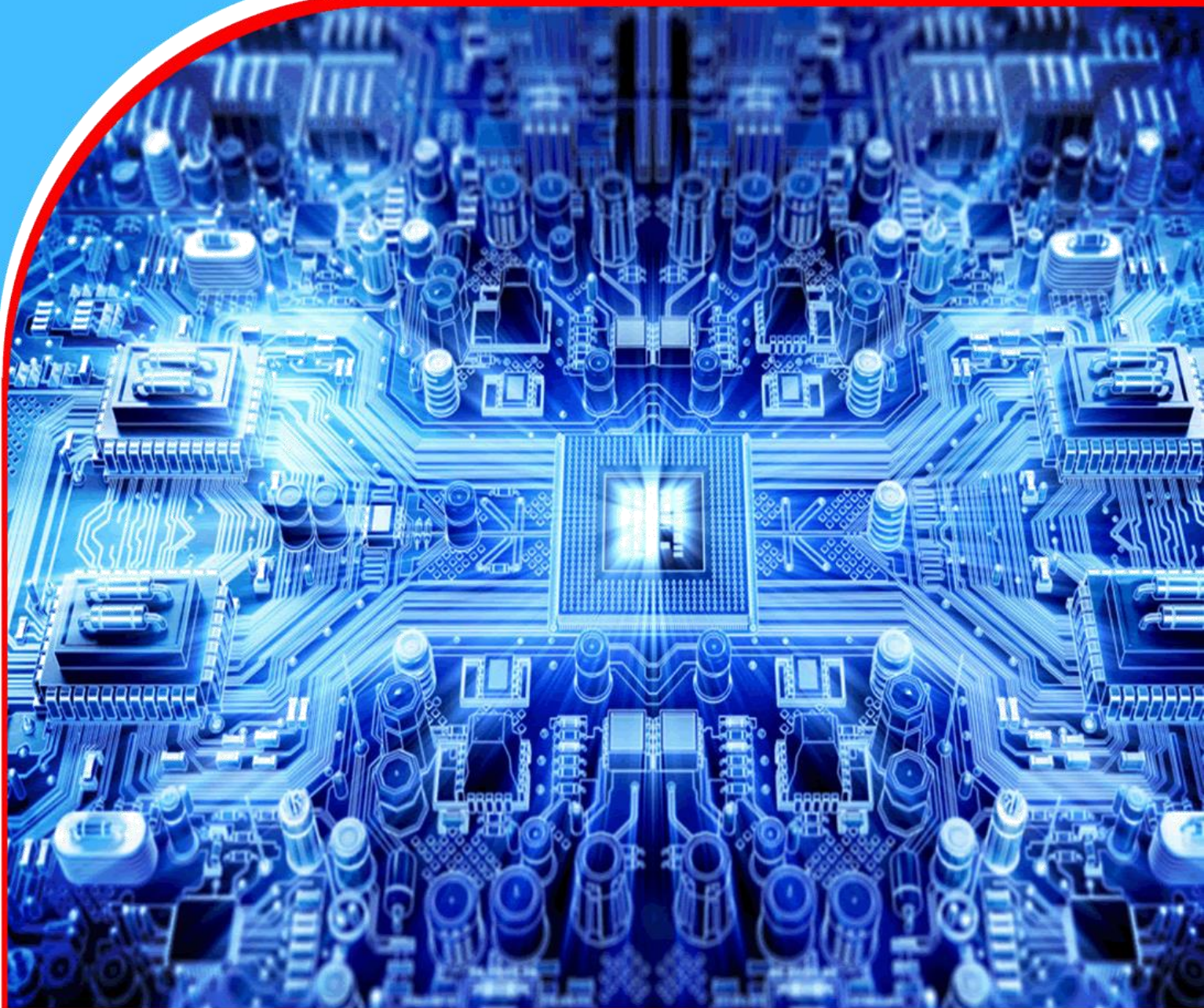


American Journal of Computing and Engineering (AJCE)



**Anomaly Detection in Pipeline Operations Using Unsupervised
and Semi-Supervised Learning**

Pankaj Verma, Krishna Gandhi



Anomaly Detection in Pipeline Operations Using Unsupervised and Semi-Supervised Learning

 Pankaj Verma^{1*},  Krishna Gandhi²

¹Indian Institute of Management, Bangalore (IIM-Bangalore), Bannerghatta Road, Bengaluru, Karnataka, India

²Illinois State University, 100 N University St, Normal, IL 61761, United States



Article history

Submitted 04.05.2023 Revised Version Received 16.06.2023 Accepted 18.07.2023

Abstract

Purpose: Oil, gas, and water transportation is important through pipeline systems which are susceptible to various anomalies such as structural degradation, malfunctions in operations, and leakages. Older physics-based and rule-based methods of monitoring, despite their interpretability, tend to have low sensitivity, flexibility, and scalability. However, the absence of labeled fault data, increasing operational complexity, and non-stationary pipeline conditions create a critical gap in reliable and scalable anomaly detection solutions for real-world deployment. This study addresses this gap by systematically analyzing data-driven unsupervised and semi-supervised learning approaches and their applicability to pipeline monitoring.

Materials and Methods: New developments in unsupervised and semi-supervised learning have made data-driven anomaly detection schemas able to learn typical operational behavior and detect anomalies with little assistance of labeled fault data. This review gives a detailed summary of these methods within pipeline monitoring. Among the methods discussed are distance- and density-based, statistical and subspace methods, and neural network-based methods, including autoencoders and self-organizing maps. Semi-supervised algorithms such as one-class classification and hybrid statistical-learning are also discussed. The review includes the issues of data characteristics, practices of evaluation, interpretability, and real-time implementation.

Findings: The study identifies and discusses a variety of unsupervised and semi-supervised learning techniques that can effectively address the challenges faced by traditional monitoring methods in pipeline systems. It highlights how these data-driven methods are able to detect anomalies by learning typical operational behavior with minimal reliance on labeled fault data. The study also covers important considerations like data characteristics, evaluation practices, and the challenges of implementing these methods in real-time environments.

Unique Contribution to Theory, Practice, and Policy: This review provides a thorough evaluation of emerging data-driven anomaly detection methods, contributing to the theoretical understanding of how unsupervised and semi-supervised learning can be applied in pipeline monitoring. The study's practical contribution lies in its exploration of real-world applicability, offering insight into methods that can enhance the sensitivity and scalability of anomaly detection in pipeline systems. For policy, the research suggests future directions, including enhanced feature learning, concept drift adaptation, and integration with digital twins, which aim to improve the trustworthiness and efficiency of anomaly detection in pipeline operations.

Keywords: *Monitoring Pipelines; Anomaly Detection; Unsupervised Learning; Semi-Supervised Learning, One-Class SVM, SCADA Systems, Time-Series Data, Fault Detection*

INTRODUCTION

The pipeline systems form a significant part of the current energy and utility infrastructure, as they are used to conduct the large-scale transportation of oil, natural gas, and water on long-distance routes [1]. They run 24 hours in different environments and operational conditions and have many sensors attached to them, which measure the pressure, flow rate, temperature, and vibration among others. The infrastructure in the pipeline should operate reliably to maintain consistent supply, safety of the population, and environmental safety [2]. Even small failures may become serious events, which will lead to massive economic losses, environmental degradation, and endangering human lives.

In ensuring integrity of the pipeline operations, anomaly detection has the key role to play [3]. Structural degradation, malfunctions or external disturbances may result in anomalies and usually these are observed as minor anomalies in sensor measurements before escalating into a serious failure. Modern pipeline monitoring systems, as shown in Fig. 1, are based on the constant sensor measurements that are converted into processed and analyzed data to detect abnormal behavior at its early stages [4]. It is important that such anomalies are detected in time to avoid leakages, ruptures and unexpected service disruption.

Traditional methods of pipeline monitoring have been majorly based on the use of physics-based models and rule-based mechanisms [5]. The mass balance analysis, negative pressure wave detection and threshold-based alarm systems are such techniques that are commonly used in supervisory control and data acquisition (SCADA) systems. Although these techniques are based on established physical principles and provide some interpretability, they tend to be weak in detecting onset faults and they are prone to sensor noise, parameter uncertainty and modeling errors. Moreover, rule-based systems are commonly very expert-intensive to set up as well as may not be able to cope with different operational environments.

Data-driven anomaly detection techniques have received growing popularity in efforts to address these shortcomings as the quantity of past sensor data and the development of machine learning techniques has increased. Specifically, unsupervised and semi-supervised learning methods have become appealing to pipeline monitoring, since they imply that no labels of faults are required, which are limited in reality due to the operation of pipeline systems [6]. The objectives of these approaches are to acquire normal operational behavior directly out of data and determine deviations of the obtained behavioral baseline as possible anomaly to provide better adaptability and robustness.

A typical learning-based pipeline monitoring workflow consists of sensor data acquisition, preprocessing steps such as filtering and normalization, and the application of anomaly detection models that generate anomaly scores or alarms for decision support, as depicted in Fig. 1. This integration of sensing, data analytics, and operational decision-making forms the foundation of contemporary pipeline anomaly detection systems and motivates the focus of this review.

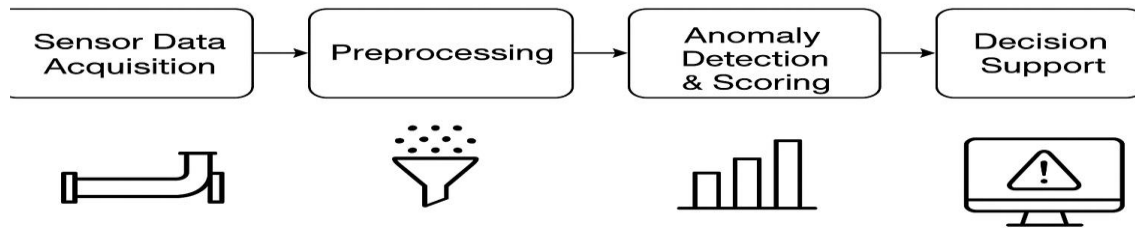


Figure1: Pipeline Anomaly Detection Workflow

The objective of this review is to provide a comprehensive overview of unsupervised and semi-supervised learning techniques for anomaly detection in pipeline operations. The review aims:

1. Describe common pipeline anomalies and associated data characteristic
2. Examine traditional monitoring approaches and their limitations,
3. Analyze learning-based anomaly detection methods applied to pipeline systems,
4. Discuss evaluation practices, practical challenges, and future research directions identified in the existing literature.

LITERATURE REVIEW

Pipeline Operations and Anomaly Types

This part gives preliminary background information in the domain of pipeline systems and the nature of the anomalies that occur during the running of the systems. Interpretation of the results of the anomaly detectors and the choice of the unsupervised or semi-supervised learning methods require a good understanding of pipeline configurations, sensing infrastructure, and fault properties [7].

Overview of Pipeline Systems

The systems of pipelines are usually categorized depending on the kind of the transported medium and their application in the system of supplies [8]. Transportation of crude oil, refined products of petroleum and water occurs by liquid pipelines whereas compressible media, such as natural gas, are transported by gas pipelines. These discrepancies result in different flow patterns and transient dynamics which has a direct impact on the statistical characteristics of sensor measurements applied in anomaly detection.

Pipelines are further classified into transmission and distribution systems in terms of how they are operated [9]. The transmission pipelines are used with high pressure in long distance and relatively steady flow and distribution pipelines distribute products among final consumers in a complex and branched network with frequent flow variations. The greater diversity of the distribution systems usually makes it hard to determine whether the changes in the operation are normal or abnormal.

Contemporary pipeline systems are equipped by various sensors interrelated by the means of supervisory control and data acquisition (SCADA) [9]. Some of the common sensor modalities are pressure, flow rate, temperature, acoustic, and vibration, which submit continuous multivariate time-series measurements. These measurements represent both stationary behavior and non-stationary events, which are the main input to learning-based anomaly detection models.

Typical Problems in The Operation of the Pipeline

Abnormalities of the pipeline may either be as a result of physical degradation, operational disturbances, accidental and external occurrences [10]. All types of anomalies will leave specific signatures in sensor data that will influence various measurements and time trends. To analyze the issue of anomaly detection, a well-organized list of these anomalies and the indicators they provide is necessary to associate the anomaly types with the sensors that are impacted and the normal behavior of the signal.

Structural and Mechanical Abnormalities

The gradual wear and tear or breakdown of pipeline parts are structural and mechanical anomalies. One of the most widespread degradation mechanisms is corrosion, which causes a decrease in the thickness of walls and declines the structural integrity in the long run. Cyclic loading of pressure, thermal stress or defects in materials may lead to cracks and fatigue failures, whereas long term accumulation of stress or failure of joints and weld lines have some relationship with installation problems or build-to-last [11].

Such anomalies are usually slow and might not present sudden sensor read alterations. Their signatures, as shown in Table 1, tend to be more of higher signal variance, small pressure losses, or long-term trends instead of sudden deviations, and thus they are hard to tell using simple threshold based techniques.

Table 1: Pipeline Anomaly Types, Causes, Affected Sensors and Signal Characteristics

Anomaly Category	Typical Causes	Affected Sensors	Observable Signal Characteristics
Structural and mechanical anomalies	Corrosion, cracks, fatigue, joint failures	Pressure, flow, acoustic	Gradual pressure loss, increased variance, long-term trends
Operational and process anomalies	Blockages, valve malfunctions, pump failures	Pressure, flow, temperature	Abrupt or intermittent deviations, flow imbalance
Small leak events	Minor wall defects, seal degradation	Pressure, flow, acoustic	Localized pressure drop, subtle flow discrepancy
Catastrophic rupture events	Severe structural failure, external damage	Pressure, flow	Sudden and large pressure loss, sharp flow reduction

Operational and Process Anomalies

The issues of operational and process anomalies are connected with a deviation of normal operating conditions or the failures of active elements. This can be blocked due to debris/deposits, valves that are defective, limiting flow and failure of pumps or compressors causing changes in pressure and throughput. These anomalies may happen abruptly and are often related to the modifications of the control actions or the equipment state.

Abrupt change in the pressure or flow measurements is a common occurrence in sensor data and is the indication of operational anomalies. But such similar trends can as well be caused by valid operational changes, like demand changes or even maintenance. This similarity in signal properties is a significant issue in the context of detecting anomalies and inspires the

application of data driven methods that have the potential to characterize normal variability in operation.

Peak and Rupture Events

The most important type of pipeline anomaly is leak and rupture events because of the safety and environmental impacts of these events. Minor leaks usually form slowly and they might not generate strong sensor observables hence small pressure drops or imbalanced flows. Leakages at early stages are thus very difficult to detect and may need sensitive detection methods which could pick up insidious alterations of normal behavior.

By contrast, catastrophic ruptures are distinguished by abrupt and considerable pressure losses and disruptions of flows in more than one sensing position. The magnitude and spatial area of noticeable alterations to the signals is directly connected to the severity of the event and therefore there is a requirement of detection techniques to be functional over the variety of anomaly scales.

Data Characteristics and Challenges in Pipeline Monitoring

The nature of the data produced by the monitoring systems and the real-world difficulties related to the application of the technology have a potent effect on the issue of effective anomaly detection in the pipeline operations [12]. The temporal and statistical characteristics of pipeline sensor data are quite complicated and restrict the usefulness of fully supervised learning methods and induce the application of unsupervised and semi-supervised ones, which currently are more capable of learning normal behaviour with minimal labeling conditions.

MATERIALS AND METHODS

Characteristics of Pipeline Sensor Data

The pipeline monitoring systems are systems that continually gather data and information on various distributed sensors placed on the pipeline. Consequently, the data are intrinsically multivariate time-series and the different sensor channels like pressure, flow rate, and temperature are highly correlated. The interdependencies indicate the physical dynamics underlying the fluid transport, and they have to be reflected in the anomaly detection models in order to prevent spurious alarms [13].

Most operational environments have relatively high sampling rates of pipeline sensors in order to be able to detect abnormal events in a timely manner. High-frequency measurements permit observing the temporary phenomena, like pressure waves or a sudden change of flow, but cause an expansion of the amount of data and computation needed. Both normal and abnormal operating conditions can have temporary variations as demonstrated in Fig. 2 and it is difficult to differentiate between significant deviations and natural changes.

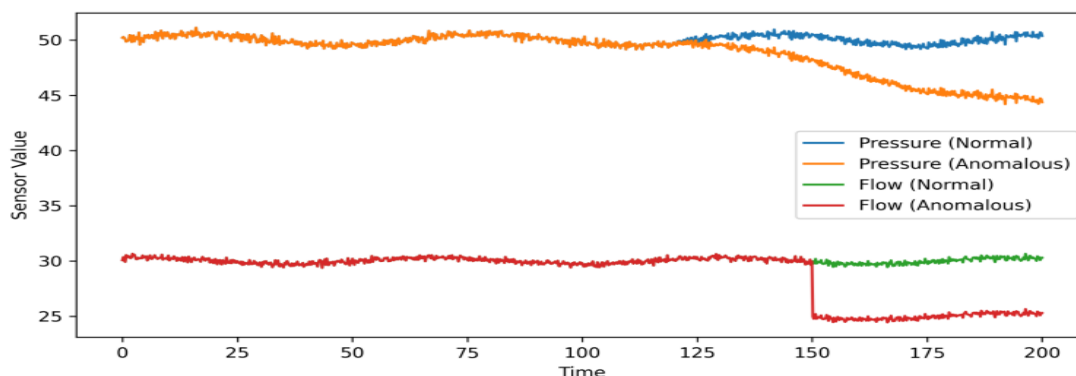


Figure 2: Illustration of Normal Versus Anomalous Pipeline Sensor Time-Series

The existence of long-term temporal dependencies is another characteristic feature of the pipeline data. The pipeline systems are usually on the move all through and the anomalies are not always unexpected but can occur over a long period [14]. Slow changes in signal behavior may be due to structural degradation, small leaks, and sensor drift and may last hours or days and even months. Such long-term patterns are critical in making an accurate anomaly decision and challenging to models whose usage is based on short-term statistics.

Practical Challenges

Although there is a rich supply of sensor data in large quantities, there are various practical issues that render anomaly detection in pipeline monitoring complicated, which additionally helps to motivate the use of unsupervised and semi-supervised methods of learning.

Data Imbalance and Rarity of Anomalies

Abnormalities in the functioning of pipelines are rare by nature when compared to the normal functioning of the pipeline [15]. The normal behavior dominates most historical data sets and within the data set there is a low proportion of data representing faults/failures. In addition, there is a general lack of labeled examples of anomalies, since many of the faults are either not seen at all or not marked. This extreme disparity of classes makes the traditional supervised learning algorithms less relevant and unsupervised and semi-supervised, that emphasize modeling of normal behavior, more feasible in the real-world pipeline surveillance.

Noise, Drift, and Missing Data

Measurement of pipeline sensors is often interfered with by noise caused by environmental factors, electromagnetic interference and sensor resolution limitations. Sensors can also tend to drift or degrade over time giving rise to gradual variations in the properties of the signals which is not related to the actual faults in the pipes. Besides that, communication issues or maintenance procedures may lead to the loss of streams of data or incompleteness of the same [16].

The development of small anomalies, especially of the pressure and flow signals, may be obscured by noise and drift. Powerful anomaly detection algorithms should hence be in a position to tell real abnormal behavior and sensor induced artifacts without using a large amount of labeled data to calibrate the algorithm.

Non- Stationarity and Operational Variability

A wide variety of conditions, such as start-up and shutdown periods, fluctuation of load, change of seasonal demand, etc. are associated with the operation of pipeline systems. With these operating regimes, sensor data contain non-stationarity, with the statistical properties of signals potentially drastically changing with time. The patterns resulting due to changes in flow setpoints, valve settings or pump operation may appear as anomalies unless adequately explained [17].

Normal transitions in operation can produce similar signal deviations as those produced by real faults. This overlap makes it hard to detect anomalies, and introduces the necessity of learning-based methods, which can improve their response to changing operating conditions and reflect complicated temporal variations, without being used to fixed thresholds or fixed models.

Traditional Approaches to Pipeline Anomaly Detection

Prior to the introduction of learning methods, pipeline anomaly detection used mostly physics-based models and rule-driven monitoring systems. Such classical approaches are the basis of numerous industrial pipeline monitoring systems and offer valuable insight on the reasons to

embrace data-driven approaches. This section examines the major groups of traditional methods and explains their limitations inherent in them.

Physics-Based and Model-Driven Methods

Physics-based and model-driven approaches make use of mathematical models of the fluid dynamics and conservation laws in identifying anomalies in the pipeline activities. Mass balance methods are among the most commonly used methods; they are grounded on the principle of mass conservation. These processes are used to compare the actual inflow and outflow of a given pipeline with the set point to detect any differences that could be a leak or a loss. Mass balance methods are comparatively easy to apply, and efficient in identifying big leakages; nonetheless, their functionality may worsen when there is measurement error, momentary operating circumstances, and inaccurate flow measurements.

Negative pressure wave (NPW) methods are another widely known type of physics-based methods. These techniques take advantage of the pressure wave of abrupt occurrences like rupture or leakage that travel along the pipeline and are registered by pressure sensors. With the help of the analysis of the time and magnitude of these pressure waves, NPW methods are in a position to determine the position and the strength of the anomalies [18]. Although NPW methods are quite effective in the high-level speed of detecting the catastrophic failure, they are not efficient when it comes to detecting the small or gradually developing leaks that do not produce the high-pressure transients. The relative merits and working nature of these physics-based methods.

Threshold-Based Systems and Rule-Based Systems

Another commonly used type of conventional pipeline monitoring methods is the rule-based and threshold-based systems. These systems are based on preset rules or alarm settings to sensor measurement like pressure, flow rate and temperature [19]. In the case of static threshold systems, an alarm goes off when a measurement has surpassed a set limit, whereas in the case of adaptive threshold systems the limits are adjusted according to past history or even operating circumstances.

The simplicity, transparency and ability to integrate with the current SCADA systems make threshold-based approaches appealing. Nevertheless, the choice of thresholds is frequently the domain-specific expertise and hand-tuning. These systems are vulnerable to the false alarm when the normal operation variability and transient condition make the measurements go beyond the set limits and they are also susceptible to fail to identify subtle anomalies that are not crossing the threshold level.

Limitations of the Traditional Techniques

Although they are commonly used, conventional techniques of pipeline anomaly detection have a number of shortcomings limiting their usefulness in complex and current pipeline networks. Physics-based algorithms are very susceptible to sensor noise, parameter uncertainty, and model errors that may result in erratic detection in non-ideal operation circumstances. They tend to begin to fail in transient conditions like start-up or shutdown when they must be operating under assumptions of steady-state behaviors that are not true.

The rule based as well as threshold based systems are simple and easy to interpret, but they lack scalability and flexibility. With the expansion of the pipeline networks, the variety of rules and thresholds necessary to perform the monitoring processes effectively is also increasing significantly, and system maintenance becomes challenging. Moreover, these approaches do not have the capability to learn based on the past trends and respond to changing trends of operation which is more clearly observed when non-stationarity and the rapidly moving

operating trends are involved. The weaknesses as shown in Table 2 have prompted the development of learning-based solutions to anomaly detection that can reveal the pattern of complex data and minimize the use of rule development that is manually programmed.

Table 2: Comparison of Traditional Pipeline Monitoring Techniques

Method Category	Representative Techniques	Strengths	Limitations
Physics-based methods	Mass balance, negative pressure wave	Physically interpretable, effective for large leaks and ruptures	Sensitive to noise, limited performance during transients, poor detection of small leaks
Rule-based methods	Static thresholds, adaptive thresholds	Simple implementation, easy integration with SCADA	High false alarm rates, manual tuning required, limited adaptability

Unsupervised Learning for Pipeline Anomaly Detection

The unsupervised forms of learning have attracted much interest in pipeline anomaly detection since they are capable of modeling normal working conditions without the need of labeled faults information. Such techniques are applied to learn the latent patterns and correlations in normal sensor data and detect deviation as possible anomalies. The section examines the major unsupervised methods, such as distance- and density-based ones, statistical and subspace methods, and neural network-based models [20].

FINDINGS

Distance and Density-Based Methods

Distance- and density-based methods recognize abnormalities with the help of the relative location of individual data points with the rest of the data. Such strategies are especially good at identifying points that are very distinct to patterns of normal operations [21].

k-Means and Clustering-Based Detection

K-Means clustering algorithm divides the data into a specific number of clusters one of which depicts common operational regimes. The identification of anomalies is done in terms of distance to cluster centroids or in terms of calculating cluster deviations. The observations that are very far off all the clusters are regarded as outliers as they represent an aberrant behaviour in the measurements of pressure, flow or temperature.

k-Nearest Neighbors (k-NN) and Local Outlier Factor (LOF)

The local density estimation is utilized in the k-NN and LOF methods. The anomaly score of a data point is calculated as comparing the density of the data point with its neighbors. The points that are found in areas with lower density than the area surrounding them are marked as anomalous. These methods work effectively in the detection of global and local anomalies and have been used successfully with multivariate pipeline sensor data.

Subspace and Statistical Methods

Statistical/subspace techniques are to model the normal behaviour of the system by trying to capture correlations and variances of multiple sensor channels. The fact that the learned subspace is not adhered to means that it might be an anomaly.

Principal Component Analysis (PCA)

PCA provides a lower-dimensional subspace of sensor data that can be used to reflect normal operating variance in high-dimensional sensor data. Deviations of this subspace are then measured by statistical measures like the Squared Prediction Error (SPE) and the T² statistic of Hotelling. The points that are characterized by high SPE or T² are considered as anomalies and it is possible to detect the small fault that can not cause big deviation in the sensor readings individually.

Independent Component Analysis (ICA)

ICA decomposes multivariate sensor signals in statistically independent components. ICA can also be used to identify abnormal behaviors due to a common or coherent source, including simultaneous leaks and transient operations and thereby not detected using PCA-based methods; by examining the residues or reconstruction errors of these components.

Neural Network-Based Unsupervised Models

Neural network approaches for unsupervised anomaly detection focus on learning complex non-linear representations of normal behavior from historical data.

Auto Encoders and Variants

Auto encoders are neural network architecture encoder-decoder neural networks that are exclusively trained on normal operational data. The network is trained to replicate the input signals and the difference between the reconstruction and the input is used as an anomaly score. Large reconstruction errors denote the observations which are not in accordance with normal patterns learned. Fig. 3 shows that auto encoders have developed as the broader choice to detect pipeline anomalies due to the capability to encompass intricate correlations and time sequences among multiple sensor streams of multivariate data.

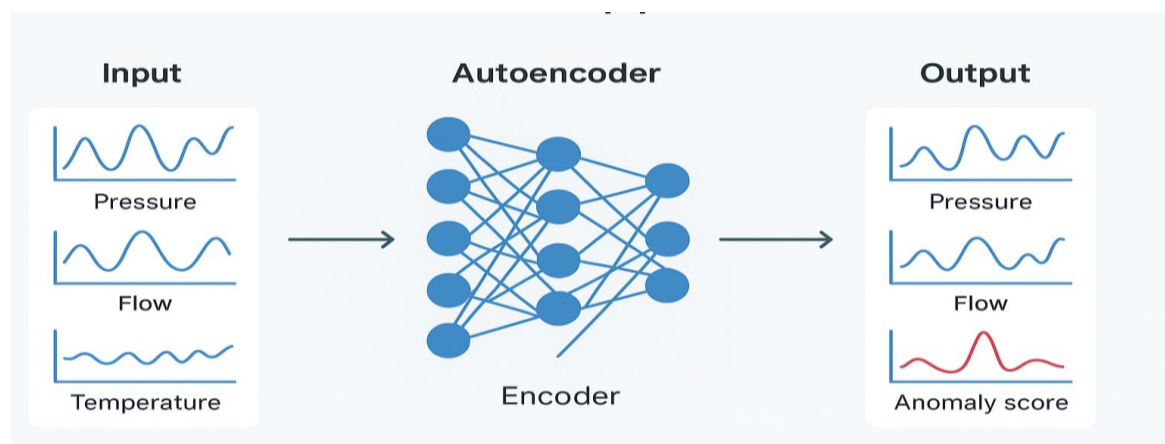


Figure 3: Auto Encoder-Based Anomaly Detection Framework for Pipelines

Self-Organizing Maps (SOM)

SOMs are neural networks which are topology-preserving and therefore map high-dimensional input data onto a low-dimensional grid, but retain the neighborhood relationships. In the training process, the network structures common ways of operation into clusters, and any variation on the learned map is construed as an anomaly. The SOMs are especially helpful

when visualizing multivariate relationships in the pipeline data and finding outliers which do not follow the established patterns.

Semi-Supervised Learning Approaches

Semi-supervised learning techniques are especially useful in pipeline fault detection, where fault data are scarce and labeled, but normal data of operation is even more abundant. These techniques use a small amount of information that is labeled, or purely use normal data to establish limits of intended behavior, and identify deviations that indicate anomalies. This part discusses the one-class classification techniques, the hybrid classification that involves a mix of statistical and learning based models and the early deep semi-supervised classification [22].

One-Class Classification Techniques

One-class classification techniques are designed to model the distribution of ordinary operational data, and any form of deviation is considered to be an anomaly. The popular use of these techniques in industrial monitoring is related to the fact that they can be used with a limited or no labeled fault data.

One-Class Support Vector Machine (OCSVM)

OCSVM creates a decision boundary that surrounds the normal data in high dimensional feature space. The algorithm finds the smallest region that includes most of the normal observations so that all new data points that fall beyond this region will be considered as anomalous. OCSVM is especially useful when the data of high-dimensional multivariate sensing of pipelines is available, since it can learn the complicated dependencies of variables using only normal operation samples to train.

Support Vector Data Description (SVDD) Describes Data That is Non-Real and Discrete

SVDD works in a similar manner as OCSVM but forms a hypersphere around normal data in the feature space. The radius and center of the hypersphere are optimized to cover the majority of the normal observations leaving out possible outliers. SVDD is particularly suitable for the representation of small manifolds of normal operation and it works effectively in recognizing infrequent or subtle anomalies that do not conform to the existing patterns.

Hybrid Statistical–Learning Models

Hybrid models are statistical modeling and machine learning approaches used to enhance detection capabilities and strength. Dimensionality reduction and extraction of the important features of multivariate pipeline sensor data can be carried out using PCA then fed into an OCSVM classifier. This pooling of the merits of the two approaches: PCA represents the major variance in typical operation, whereas OCSVM finds variations in the transformed space of characteristics. These hybrid models enhance sensitivity to anomalies and have computational efficiency and interpretability [23].

Early Deep Semi-Supervised Models

The initial methods of semi-supervised anomaly detection by deep learning resort to shallow neural networks that are trained on normal operation data only. These models develop a representation of normal behavior and apply reconstruction or prediction error to identify the deviations. Whereas these models are shallow and lack complexity, in comparison to subsequent deep architectures, they showed that neural networks could be used to complement semi-supervised approaches to pipeline monitoring.

Table 3 presents a comparative summary of the unsupervised and semi-supervised approaches to detect pipeline anomalies and shows that they have variations in data requirements, interpretability, and computational cost. As indicated, semi-supervised methods tend to be

more sensitive to rare anomalies and yet at the same time do not require a high reliance on labeled fault data.

Table 3: Comparison of Unsupervised Vs Semi-Supervised Learning Methods for Pipeline Anomaly Detection

Method Category	Representative Techniques	Data Requirements	Interpretability	Computational Cost
Unsupervised	k-Means, LOF, PCA, ICA, Autoencoders, SOM	Normal data only	Medium	Low to medium
Semi-supervised	OCSVM, SVDD, Hybrid PCA+OCSVM, shallow neural networks	Normal data, limited labeled anomalies	Medium to high	Medium to high

Future research should focus on developing adaptive, theory-informed unsupervised and semi-supervised anomaly detection models that explicitly handle non-stationarity, sensor noise, and rare-event detection, while validating performance on long-term, real-world pipeline datasets.

Evaluation Strategies and Performance Metrics

Measuring the performance of anomaly detection systems is an important attribute towards establishing whether they can be practically applicable in pipeline operations. Best evaluation models can evaluate the accuracy and reliability of models in identifying abnormalities and consider the constraints of scarce labeled fault information and different operating environments.

Data on Benchmarking and Industry Case Studies

The measurement of anomaly detection algorithms is usually based on simulated data and actual industrial data. Simulated pipeline datasets permit controlled experiments where the fault scenarios, the magnitude of the leaks and the variation in the operations can be varied systematically. The datasets can be useful in training algorithms, as well as in comparing the performance of a method across methods in a reproducible manner.

Besides that, real-world complexity of proprietary SCADA data of running pipelines is used, such as sensor noise, non-stationary behaviour and infrequent faults. Such industrial datasets are vital in evaluations to determine the behavior of models under real world scenarios such as presence of subtle or slow growing anomalies. To evaluate the performance of anomaly detection models both types of datasets are utilized to train, validate and test the models in controlled and realistic environments as shown in Fig 4.

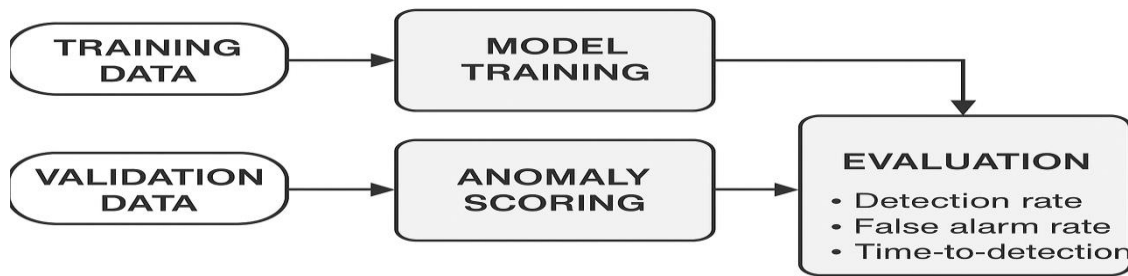


Figure 4: Evaluation Workflow for Pipeline Anomaly Detection Systems

Evaluation Metrics

Several quantitative metrics are commonly used to assess anomaly detection performance:

Detection Rate (True Positive Rate)

Calculates the rate of true anomalies that have been properly identified by the model. Detection rates should be high to make sure that faults are identified in time.

False Alarm Rate (False Positive Rate)

This is a ratio that scales how many normal events are false alarms. False alarm is low, which is valuable to minimize undue interventions and operational interference.

Time-to-Detection

The time-to-Detection measures the time interval between the detection of the anomaly and its beginning. It results in less time to identify a leak, rupture, or damaged equipment.

When these metrics are put into consideration, researchers will be able to assess the sensitivity and reliability of anomaly detection systems. It also shows a common evaluation process which consists of data partitioning, normal training a model, anomaly scoring and performance measures.

Practical Deployment Considerations

In industrial settings, several practical considerations affect the deployment and usefulness of anomaly detection models:

Alarm Management

Continuous anomaly scoring can generate frequent alarms. Effective thresholding and prioritization strategies are necessary to prevent alarm fatigue among operators.

Human-in-the-Loop Validation

Operator review and expert validation remain essential, especially for subtle or rare anomalies, to ensure reliable decision-making. Integrating machine learning outputs with human expertise enhances overall system trustworthiness and operational safety.

Comparative Discussion, Open Challenges and Future Research Directions

This part summarizes the results of the literature review of the studies of pipeline anomaly detection and identifies the unresolved problems and future research opportunities. It also offers a holistic view to the researchers and practitioners by combining the method comparison, practical constraints, and future prospects.

Comparison of Methods among Studies

Unsupervised and semi-supervised approaches to pipeline anomaly detection have been implemented in a number of different ways with varying strengths and weaknesses. Simpler and denser methods like k-Means, k-NN, and LOF have simplicity and interpretability but might not perform as well with high dimensional data and complicated temporal relationships. Statistical and subspace methods, such as PCA and ICA, can be good at identifying correlated behaviour of sensor behaviours, but can be noisy and can fail to find nonlinear behaviour patterns. Models based on neural networks, especially autoencoders and SOMs, have better ability to capture complex and nonlinear relationships in multivariate sensor measurements though at the expense of increased computational complexity and decreased interpretability.

The semi-supervised techniques, including OCSVM, SVDD and hybrid PCA OCSVM models, are better in detecting infrequent anomalies and take less labelled data. It identifies the trade-offs between unsupervised/ semi-supervised methods in terms of data need, interpretability and cost of computation. In general, the choice of methods is based on a trade-off between the accuracy of detection and the strength, interpretability, and computability.

Interpretability and Trustworthiness

One important issue with the applications of machine learning to pipeline monitoring is that detection outcomes are not easily interpretable. Anomaly alerts should be believed in by operators and maintenance people to make good decisions in time. Although physics-based and threshold-based methods have high interpretability, a significant number of the sophisticated learning-based models are black-box. Increasing explainability, feature attribution, visualization of learned latent representations, or hybrid statistical-learning models must be improved to give operators confidence in the operators, and enable human-in-the-loop decision making.

Scalability and Real Time Constraints

The pipeline networks may extend over several hundreds of kilometers with thousands of sensors creating large amounts of data in real time. Most of the learning based approaches, especially deep neural networks demand huge computational and memory demands which can constrain real time usage. Approaches should then be able to balance the complexity and accuracy of models and efficiency to make sure that timely detection of anomalies occurs within large scale industrial systems.

Future Research Directions

Several avenues were anticipated in the literature for advancing pipeline anomaly detection:

Improved Feature Learning For Time-Series Data

Establishing methods that can extract informative features of high dimensional, multivariate sensor signals automatically to improve performance of detection.

Greater Stability to Concept Drift

Accommodations made to change operation conditions and non-stationary behavior in order to remain reliable with time.

Integration with Digital Twins

Integrating physics-based simulation models with data-driven anomaly detection to enhance accuracy, interpretability and scenario-based predictive control.

Collectively, these issues and research priorities indicate that the methods must be precise, interpretable, computationally efficient, and flexible enough to adapt to the changes in the oil and gas pipelines, as the basis of further advancement in the domain.

CONCLUSION AND RECOMMENDATIONS

Conclusion

The oil, gas and water transportation in the world relies on pipeline systems as one of the most important infrastructures. Their safe and reliable running needs the identification of an anomaly in time, which can be caused by structural degradation, malfunctions and leaks or external disturbances. The conventional pipeline monitoring tools, such as physics-based modeling, mass balance, negative pressure wave analysis, and rule-based thresholding, have interpretability and some reliability in the controlled conditions. Nevertheless, these methods tend to have problems with identifying elusive or changing anomalies, operating under dynamically changing conditions, and supporting large and complex networks of pipelines.

Unsupervised and semi-supervised learning techniques have become a force to reckon with, with the ability to be able to take advantage of the large amounts of normal operation data and few labeled fault data. This review puts into perspective the advantages and disadvantages of major distance- and density-based approaches (k-Means, k-NN, LOF), statistical and subspace approaches (PCA, ICA) and neural network-based models (autoencoders, self-organizing maps). Such methods are specifically useful in representation of complex multivariate relationships, identification of temporal dependencies, and identification of both sudden and gradual anomalies. Further sensitivity to rare and subtle anomalies is achieved by semi-supervised algorithms, such as one-class classification algorithms (OCSVM, SVDD), and by hybrid statistical learning algorithms, which do not trade off flexibility or scalability.

Recommendations

In spite of these developments, there are still a number of issues. Pipeline sensor data is by nature noisy, high dimensional and tends to have missing values, sensor drift as well as non-stationarity occurring as a result of changing operational regimes. Learning-based models should thus be resistant to these data properties whilst being computationally efficient to facilitate real-time detection of large pipeline network systems. Also, the interpretability of complex models is a major issue to the operator trust and human-in-the-loop decision-making. Explainable anomaly scores or model visualization models are more likely to be used in industrial application.

Future prospects The research directions are to focus on more advanced feature extraction and representation learning algorithms on multivariate time-series data, which can detect small or hitherto unknown patterns of anomalies. Algorithms that can deal with concept drift will be critical to ensuring the consistency of performance in the long-term operation. Integration of learning-based models and digital twin models is also a potential solution to integrate physics-informed simulations with data-driven anomaly detection, which provides better predictive performance, scenario analysis, and interpretability. Moreover, scalable architecture and deployment plans with tradeoffs between accuracy, computation efficiency and real time responsiveness will be crucial in the wide scale industrial adoption.

In short, the unsupervised and semi-supervised methods of learning offer a flexible and efficient solution to pipeline anomaly detection to the weaknesses of the classical approaches to monitoring. These approaches, combined with proper evaluation practices, human-in-the-loop validation, and innovative technologies, like digital twins, can potentially make the pipeline systems much safer, reliable, and resilient. The ongoing development of these

strategies offers smarter, more adaptable and reliable anomaly detection systems of highly sophisticated industrial systems.

This study highlights the growing importance of unsupervised and semi-supervised learning approaches for reliable pipeline monitoring in data-scarce environments. The findings provide practical guidance for pipeline operators in selecting scalable and adaptive anomaly detection techniques, while also offering researchers a structured foundation for developing more robust, interpretable, and deployment-ready monitoring solutions.

REFERENCES

- [1] A. B. Klass and D. Meinhardt, "Transporting oil and gas: US infrastructure challenges," *Iowa L. Rev.*, vol. 100, pp. 947, 2014.
- [2] Strogen et al., "Environmental, public health, and safety assessment of fuel pipelines and other freight transportation modes," *Applied Energy*, vol. 171, pp. 266–276, 2016.
- [3] V. N. Onyechi, "Pipeline integrity and risk prevention: Real-time monitoring, structural health analytics, and failure mitigation in harsh operating environments," *Magna Scientia Advanced Research and Reviews*, vol. 3, no. 2, pp. 139–151, 2021.
- [4] Ho, Michael, et al. "Inspection and monitoring systems subsea pipelines: A review paper." *Structural Health Monitoring*, 19.2 (2020): 606-645.
- [5] O. Erge and E. van Oort, "Combining physics-based and data-driven modeling in well construction: Hybrid fluid dynamics modeling," *Journal of Natural Gas Science and Engineering*, vol. 97, p. 104348, 2022.
- [6] M. H. Alobaidi, M. A. Meguid, and T. Zayed, "Semi-supervised learning framework for oil and gas pipeline failure detection," *Scientific Reports*, vol. 12, no. 1, p. 13758, 2022.
- [7] D. Sen et al., "Data-driven semi-supervised and supervised learning algorithms for health monitoring of pipes," *Mechanical Systems and Signal Processing*, vol. 131, pp. 524–537, 2019.
- [8] S. Razvarz, R. Jafari, and A. Gegov, "The importance of pipeline transportation," in *Flow Modelling and Control in Pipeline Systems: A Formal Systematic Approach*, Cham: Springer International Publishing, pp. 1–24, 2020.
- [9] W. Yu et al., "Gas supply reliability assessment of natural gas transmission pipeline systems," *Energy*, vol. 162, pp. 853–870, 2018.
- [10] S. Razvarz, R. Jafari, and A. Gegov, "Flow modelling and control in pipeline systems," *Stud. Syst. Decis. Control*, vol. 321, no. 1, pp. 25–57, 2021.
- [11] D. Zhukov, S. Konovalov, and A. Afanasyev, "Morphology and development dynamics of rolled steel products manufacturing defects during long-term operation in main gas pipelines," *Engineering Failure Analysis*, vol. 109, p. 104359, 2020.
- [12] A. R. Munappy, J. Bosch, and H. H. Olsson, "Data pipeline management in practice: Challenges and opportunities," in *International Conference on Product-Focused Software Process Improvement*, Cham: Springer International Publishing, 2020.
- [13] M. Golshan et al., "Pipeline monitoring system by using wireless sensor network," *IOSR J. Mech. Civ. Eng.*, vol. 13, no. 3, pp. 43–53, 2016.
- [14] S. Anwar et al., "A framework for single and multiple anomalies localization in pipelines," *Journal of Ambient Intelligence and Humanized Computing*, vol. 10, no. 7, pp. 2563–2575, 2019.
- [15] S. Timashev and A. Bushinskaya, "Methods of assessing integrity of pipeline systems with different types of defects," in *Diagnostics and reliability of pipeline systems*, Cham: Springer International Publishing, pp. 9–43, 2016.
- [16] B. Wong and J. A. McCann, "Failure detection methods for pipeline networks: From acoustic sensing to cyber-physical systems," *Sensors*, vol. 21, no. 15, p. 4959, 2021.

- [17] B. Kraszewski, "A study of thermal effort during half-hour start-up and shutdown of a 400 MW steam power plant spherical Y-pipe," *Case Studies in Thermal Engineering*, vol. 21, p. 100728, 2020.
- [18] R. Loubere et al., "ReALE: a reconnection-based arbitrary-Lagrangian–Eulerian method," *Journal of Computational Physics*, vol. 229, no. 12, pp. 4724–4761, 2010.
- [19] J. Tejedor et al., "Machine learning methods for pipeline surveillance systems based on distributed acoustic sensing: A review," *Applied Sciences*, vol. 7, no. 8, p. 841, 2017.
- [20] S. S. Aljameel et al., "An anomaly detection model for oil and gas pipelines using machine learning," *Computation*, vol. 10, no. 8, p. 138, 2022.
- [21] E. Güngör and A. Özmen, "Distance and density based clustering algorithm using Gaussian kernel," *Expert Systems with Applications*, vol. 69, pp. 10–20, 2017.
- [22] C. Spandonidis, P. Theodoropoulos, and F. Giannopoulos, "A combined semi-supervised deep learning method for oil leak detection in pipelines using IIoT at the edge," *Sensors*, vol. 22, no. 11, p. 4105, 2022.
- [23] A. Ayadi et al., "Kernelized technique for outliers detection to monitoring water pipelines based on WSNs," *Computer Networks*, vol. 150, pp. 179–189, 2019.

License

Copyright (c) 2023 Pankaj Verma, Krishna Gandhi



This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

Authors retain copyright and grant the journal right of first publication with the work simultaneously licensed under a [Creative Commons Attribution \(CC-BY\) 4.0 License](https://creativecommons.org/licenses/by/4.0/) that allows others to share the work with an acknowledgment of the work's authorship and initial publication in this journal.