



Performance Evaluation of Machine Learning Models for Predictive Maintenance in Gas Treatment Plants under Unscaled and Noisy Data Conditions

Femi Adeoye Alabi¹

Department of Electrical/Electronic Engineering, Bells University of Technology, Ota, Ogun State

falabious@yahoo.com

Bamidele Stephen Omoyajowo² (Ph.D)

Department of Science Education, Faculty of Education Lead City University

omoyajowo.bamidele@lcu.edu.ng

<https://orcid.org/0009-0009-7824-2580>

Adetiba Gbenga³

Department of Electrical/Electronic Engineering, Bells University of Technology, Ota, Ogun State

aadetiba@bellsuniversity.edu.ng

Abstract

The reliability and efficiency of gas treatment plants are critical to sustaining the global energy supply, particularly in the oil and gas sector where operational equipment is exposed to harsh conditions. Traditional maintenance approaches such as reactive and preventive strategies have proven inadequate due to their limitations in predicting failures before they occur, often resulting in costly downtimes and inefficiencies. This study aimed to evaluate the performance of machine learning models for predictive maintenance in gas treatment plants using operational metering data that included outliers and was not scaled. Historical data from January 2019 to June 2024, obtained from Total Energies EP Nigeria Limited, was used to develop and assess four supervised learning models: Random Forest (RF), Decision Tree (DT), K-Nearest Neighbors (KNN), and Support Vector Machine (SVM). The models were simulated in a Python environment and evaluated using accuracy, precision, recall, and F1-score. Results revealed that RF and DT achieved a perfect classification accuracy of 100%, though this may indicate over fitting and require further validation. KNN showed moderate performance with 86% accuracy, while SVM underperformed significantly at 47%, highlighting its sensitivity to unprocessed data. It is concluded that while some models like SVM may struggle under unstructured data conditions, ensemble-based models like RF and DT offer strong potential for real-world deployment. Based on the findings, it is recommended that gas treatment facilities implement ensemble-based machine learning models to enhance predictive maintenance, reduce unplanned downtimes, and improve overall operational efficiency.

Keywords: Predictive maintenance, gas treatment plant, machine learning, outliers, unscaled data, Random Forest.



Introduction

Energy remains the bedrock of modern development, enabling economic advancement, technological progress, and improved living standards. Over the past century, global energy demand has grown significantly—reaching 552 quadrillion British thermal units in 2016 (Pandey et al., 2020)—driven by rising population, industrialization, and digital transformation. The oil and gas (O&G) sector remains central to global energy supply, meeting 55% of demand in 2016 and projected to provide 57% by 2040 (Ediger et al., 2023).

Among fossil fuels, natural gas stands out for its flexibility, availability, and lower environmental footprint compared to coal and oil. It plays a critical role in electricity generation, industrial processes, and transportation. Its cleaner combustion—producing nearly 50% less carbon dioxide than coal—makes it a strategic fuel in the global transition toward cleaner energy (Mohammad et al., 2021). However, the delivery of high-quality natural gas requires a complex infrastructure, especially through gas treatment plants (GTPs), which purify the raw gas by removing contaminants such as carbon dioxide (CO₂), hydrogen sulfide (H₂S), and water vapor. This ensures both efficiency and safety in downstream applications (Gao et al., 2022; Wilson et al., 2023).

The performance and reliability of GTPs depend heavily on critical equipment—compressors, scrubbers, heat exchangers, and separators—operating under harsh conditions. High pressure, corrosive substances, and temperature extremes contribute to equipment wear and potential failure (Poe & Mokhtab, 2017). Equipment breakdowns not only disrupt production but also pose environmental and safety risks and incur high maintenance and shutdown costs (Al-Janabi, 2020).

To maintain operational reliability, traditional maintenance strategies—reactive maintenance (RM) and preventive maintenance (PM)—have long been employed. RM is often costly and disruptive, as it involves addressing breakdowns after they occur (Ucar et al., 2024). PM, although more proactive, is based on time intervals or past records and may lead to unnecessary part replacements and inflated maintenance costs without guaranteeing avoidance of failures (Yang et al., 2021).

The limitations of RM and PM—particularly in high-risk and high-cost environments like GTPs—highlight the need for a smarter, data-driven approach. Predictive maintenance (PdM) emerges as a solution, leveraging real-time data and analytics to anticipate equipment failures before they happen. By relying on historical and sensor data, PdM helps optimize maintenance schedules, reduce unplanned downtime, and extend equipment lifespan (Achouch et al., 2022).

In recent years, the integration of Industry 4.0 technologies—especially the Internet of Things (IoT), Artificial Intelligence (AI), and Machine Learning (ML)—has significantly enhanced the effectiveness of predictive maintenance. ML models can analyze large volumes of equipment data,



recognize patterns, and predict potential failures, thereby supporting timely interventions and improving decision-making (Silvestri et al., 2020; Arena et al., 2024). ML algorithms such as Support Vector Machines (SVM), Artificial Neural Networks (ANN), Fuzzy Logic, and evolutionary optimization techniques like Genetic Algorithms (GA) and Particle Swarm Optimization (PSO) have been successfully applied to fault detection, system optimization, and diagnostics in engineering systems (Sircar et al., 2021).

Despite the remarkable progress in ML applications for predictive maintenance, there remains a gap in understanding how these models perform under real-world conditions, particularly when datasets contain outliers and are not pre-processed through feature scaling. Equipment sensor data are often noisy, inconsistent, and prone to outliers, which can distort model predictions. Additionally, many industrial datasets are not standardized before analysis, which can affect the accuracy and reliability of machine learning algorithms.

This study aims to evaluate the performance of machine learning models for predictive maintenance in gas treatment plants under realistic data conditions—specifically in the presence of outliers and without feature scaling. By addressing this gap, the study seeks to contribute to the development of more robust and practical AI-based maintenance strategies that enhance operational efficiency and reliability in the oil and gas sector.

Statement of the Study

Gas treatment plants (GTPs) play a crucial role in ensuring that natural gas meets the stringent quality standards required for safe and efficient use. However, these facilities operate under harsh and demanding conditions, exposing critical components such as compressors, heat exchangers, separators, and scrubbers to high pressures, extreme temperatures, and corrosive gases (Poe & Mokhtab, 2017). This operational stress leads to periodic wear and tear, thereby increasing the risk of equipment failure. Maintaining continuous and reliable operation of this equipment is essential to meet the growing global demand for natural gas, yet traditional maintenance strategies have proven inadequate in addressing this challenge.

Conventional maintenance practices, including reactive and preventive maintenance, fall short in ensuring the sustained reliability of GTP operations. Reactive maintenance addresses equipment issues only after a failure has occurred, often resulting in prolonged downtime, high repair costs, and potential safety hazards (Abidi et al., 2022). Preventive maintenance, although more forward-looking, relies heavily on historical data and scheduled servicing, which may not accurately reflect the real-time condition of the equipment. This can lead to unnecessary maintenance activities, resource wastage, and operational inefficiencies (Yang et al., 2021).



The limitations of these traditional approaches are especially critical in gas treatment plants, where unplanned equipment failure can disrupt energy supply and lead to substantial financial and environmental consequences. The need for a more accurate, responsive, and cost-effective maintenance solution is therefore urgent. Predictive maintenance (PdM), driven by artificial intelligence (AI) and machine learning (ML), provides a promising alternative by utilizing real-time data to anticipate equipment faults before they occur. This proactive approach can significantly reduce maintenance costs, enhance system reliability, and support the uninterrupted operation of GTPs (Arena et al., 2024).

This study is therefore designed to bridge the gap in existing maintenance strategies by developing an AI-driven predictive maintenance alert system tailored for gas treatment plants. The aim is to optimize maintenance practices, minimize unplanned downtimes, and ensure continuous, reliable operation of these critical facilities in the oil and gas industry.

Objective of the Study

The objective of this study is to assess the effectiveness and robustness of machine learning models in predicting equipment failures under data conditions involving outliers and unscaled features.

Methodology

This study focuses on developing a predictive maintenance alert system for gas treatment plants using historical operational metering data obtained from Total Energies EP Nigeria Limited. The data, collected between January 2019 and June 2024, consists of detailed metering records from various components of the gas treatment process. To build the predictive system, supervised machine learning classification algorithms were employed, including Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and K-Nearest Neighbors (KNN). These algorithms were used to predict when maintenance is likely to be required for the plant's equipment. The model development and simulations were conducted using the Jupyter Notebook Python environment. Before model training, the raw dataset underwent thorough preprocessing. As the data was unstructured and contained significant outliers, it required cleaning and transformation. This process began with feature engineering to derive useful attributes, such as a 'maintenance due' indicator, based on specific equipment health parameters. Following this, feature selection was applied to remove irrelevant, redundant, or highly correlated variables to avoid model overfitting and improve prediction accuracy. Correlation analysis was also performed to guide the selection of the most impactful features. To address the presence of outliers and prepare the data for effective analysis, normalization techniques such as z-score scaling and robust scaling were applied. This step ensured that the range and distribution of data did not distort the learning process. Finally, the performance of each predictive model was evaluated using standard

metrics—accuracy, precision, recall, and F1-score—to determine the most reliable algorithm for forecasting maintenance needs in gas treatment operations.

Results

The operational gas metering data collected from Total Energies EP Nigeria Limited, covering the period from January 2019 to June 2024, was subjected to exploratory data analysis to gain a clearer understanding of the dataset prior to model development. Descriptive statistical measures such as mean, standard deviation, minimum, maximum, as well as the 25th, 50th (median), and 75th percentiles were computed and are presented in Table 1. To ensure accuracy and prevent bias in the predictive models, the distribution patterns of all variables were also examined.

Table 1: Statistical Summary of the Collected Gas Data

	mean	std	min	25%	50%	75%	max
KSm ³ (PAY)	8322.32	1345.15	203.92	7759.92	8723.09	9053.49	10296.89
KSm ³ (Check)	8230.11	1503.62	1.00	7637.77	8663.20	9008.75	10293.21
Dev KSm ³	0.01	0.08	-0.18	0.00	0.00	0.00	1.00
Tonne (PAY)	6742.55	1136.03	163.34	6216.40	7183.96	7481.04	8247.80
Tonne (CHECK)	6638.92	1289.25	1.00	6117.84	7132.46	7451.08	8244.86
Dev Ton	0.02	0.09	-0.18	0.00	0.00	0.00	1.00
GJ (PAY)	349854.30	58904.00	7899.21	322874.69	367294.99	380796.171	440840.75
GJ (CHECK)	347207.21	64244.26	1.00	320968.22	364981.08	379440.62	440683.46
Dev GJ	0.01	0.08	-0.19	0.00	0.00	0.00	1.00
TGC	100.00	0.02	99.98	99.99	100.00	100.00	100.38

Table 2 presents the skewness values of the features, which indicate the extent to which each feature deviates from a symmetrical (normal) distribution. Visual representations, including histograms and kernel density estimation (KDE) plots shown in Figures 1 and 2, further illustrate these distributions. Key features such as KSm³ (PAY), KSm³ (CHECK), Tonne (PAY), Tonne (CHECK), GJ (PAY), and GJ (CHECK) demonstrated negative skewness, indicating a shift to the left of a normal distribution, with skew values of -1.814, -2.254, -1.937, -2.152, -1.652, and -2.132 respectively. Conversely, features like Dev KSm³, Dev Ton, Dev GJ, and TGC displayed strong positive skewness, with values of 11.109, 8.312, 11.858, and 16.880 respectively. Only the "Days" feature appeared normally distributed, with a skewness value of 0.015.

Table 2: Distribution of the collected gas data based on skewness parameters

S/N	Features	Skewness
1	Days	0.0153
2	KSm3 (PAY)	-1.814
3	KSm3 (CHECK)	-2.254
4	Dev KSm3	11.109
5	Tonne (PAY)	-1.937
6	Tonne (CHECK)	-2.152
7	Dev Ton	8.312
8	GJ (PAY)	-1.652
9	GJ (CHECK)	-2.132
10	Dev GJ	11.858
11	TGC	16.880

The raw gas plant data obtained did not include a direct indicator for when maintenance was due. To support the development of a predictive model, a new target variable labeled “maintenance_due” was created. Figure 1 indicates that maintenance is conducted frequently to avoid unexpected equipment failures. However, such frequent interventions are costly and not economically sustainable in the long term.

To optimize the prediction of maintenance requirements, a correlation matrix was generated between all available features and the newly created maintenance_due variable, as shown in Figure 2. Based on this analysis, four features—Dev KSm³, Dev Ton, Dev GJ, and TGC—were identified as being most strongly associated with maintenance events and were therefore selected for the model.

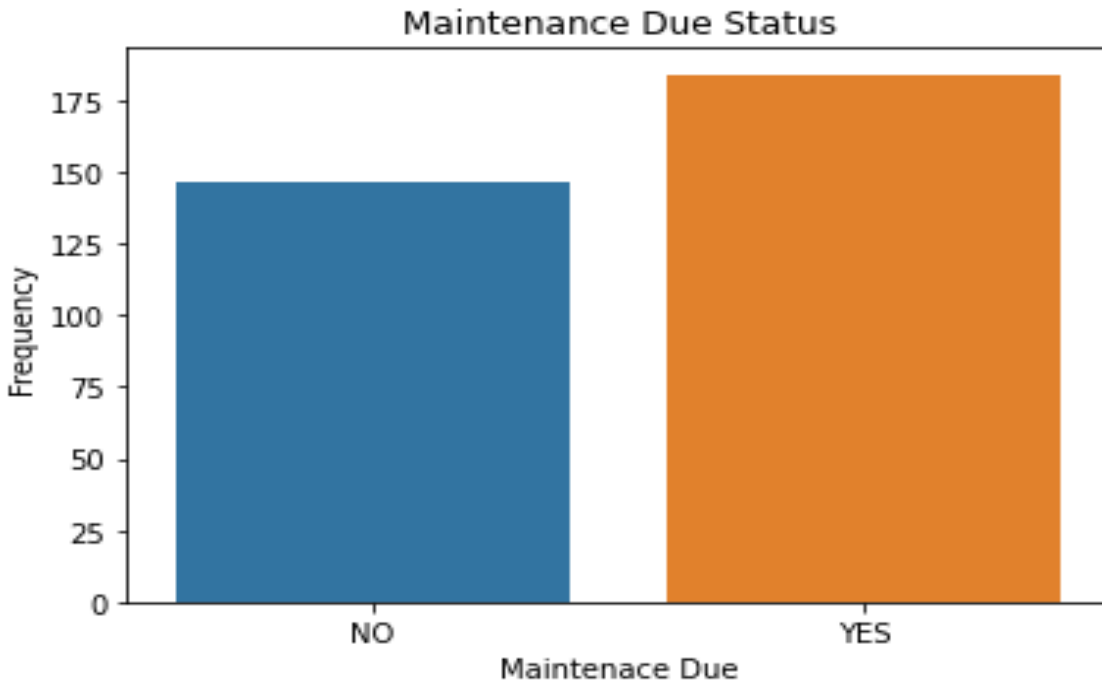


Figure 1: Maintenance Due Status of the Gas Treatment Plants

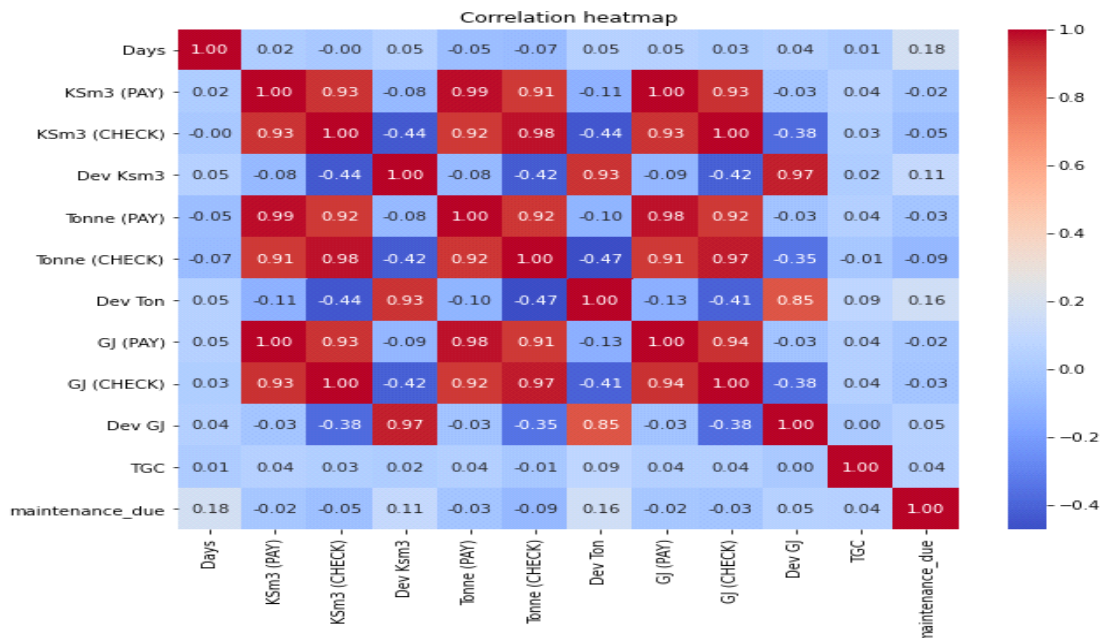


Figure 2: Confusion matrix showing the correlation between the variables



In line with the aim and objective of this study, the performance of four supervised machine learning models—Random Forest (RF), Decision Tree (DT), Support Vector Machine (SVM), and K-Nearest Neighbors (KNN)—was evaluated using operational metering data from a gas treatment plant. The evaluation was conducted under realistic data conditions involving the presence of outliers and without applying feature scaling, to assess the models' robustness and predictive effectiveness in forecasting maintenance needs.

The models were tested using confusion matrices and key performance metrics, including accuracy, precision, recall, and F1-score. Results showed that both the RF and DT models achieved an accuracy of 100%, with zero misclassifications, indicating perfect classification performance under the given conditions. However, such ideal results are considered overly optimistic and not typical of real-world systems, suggesting possible overfitting or sensitivity to the specific dataset.

The KNN model achieved a reasonably good performance with an accuracy of 86%, correctly classifying 57 instances and misclassifying 9. This indicates that KNN was moderately affected by the presence of outliers but did not reach the reliability level required for deployment in critical gas plant operations.

In contrast, the SVM model performed poorly, with an accuracy of only 47%, correctly classifying just 31 instances and misclassifying 35. This significant drop in performance highlights the model's sensitivity to unscaled features and outliers, rendering it unsuitable for predictive maintenance in its current form.

The results reveal that while RF and DT models are less affected by data irregularities, their idealized performance warrants caution. KNN shows promise but requires further optimization, and SVM is not robust enough under these data conditions. These findings underscore the importance of data quality and preprocessing in machine learning applications for predictive maintenance and affirm the need for models that balance accuracy with generalizability in real-world environments.

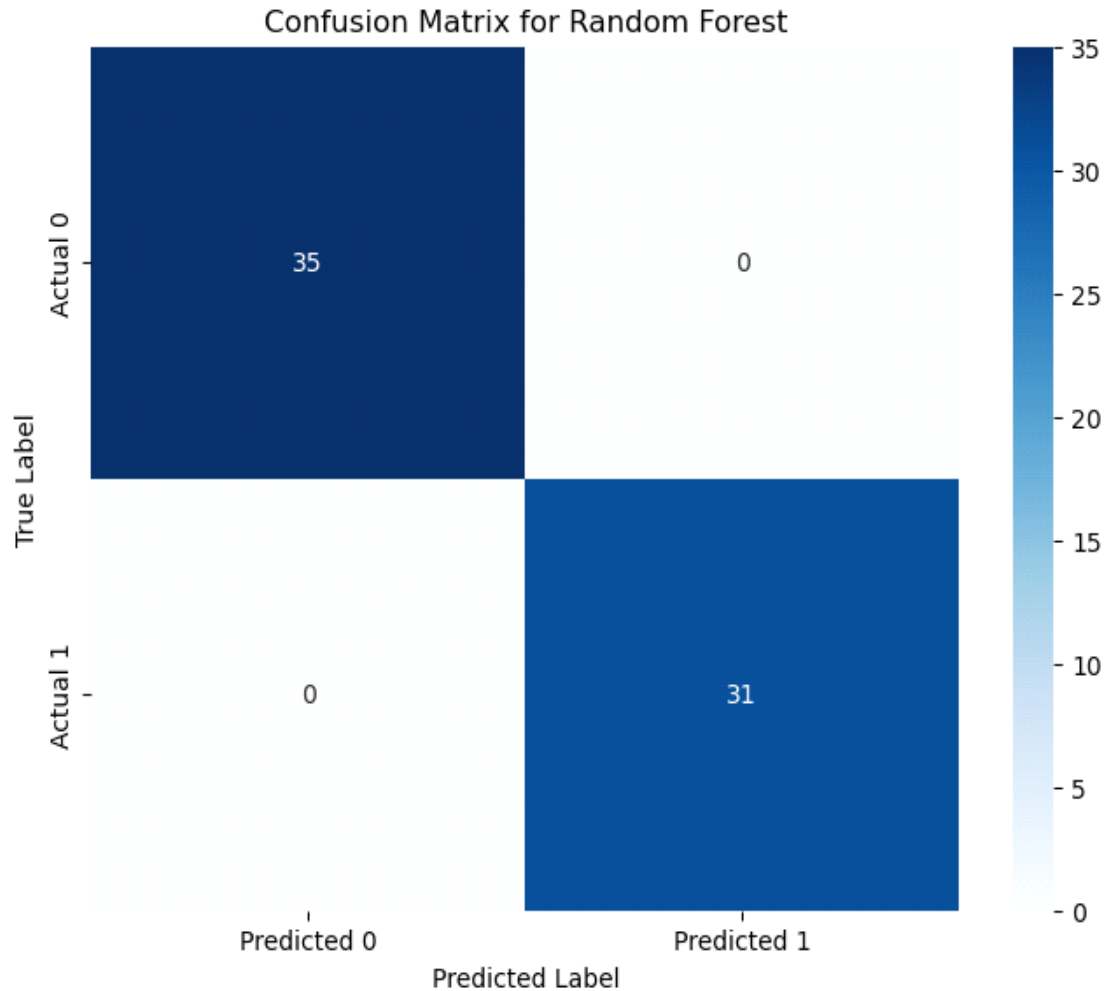


Figure 3: Confusion matrix of RF based maintenance alert system

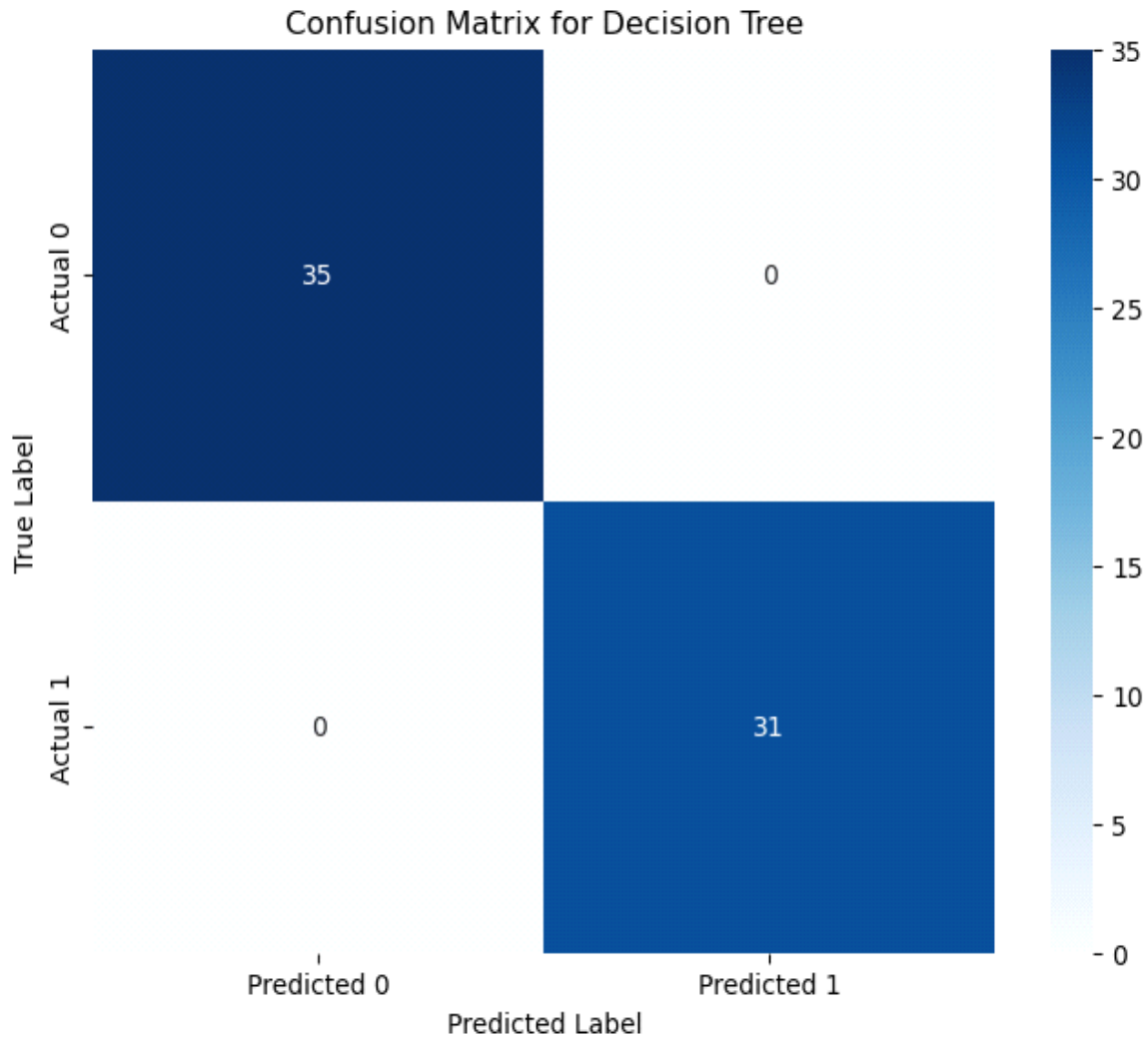


Figure 4: Confusion matrix of DT based maintenance alert system

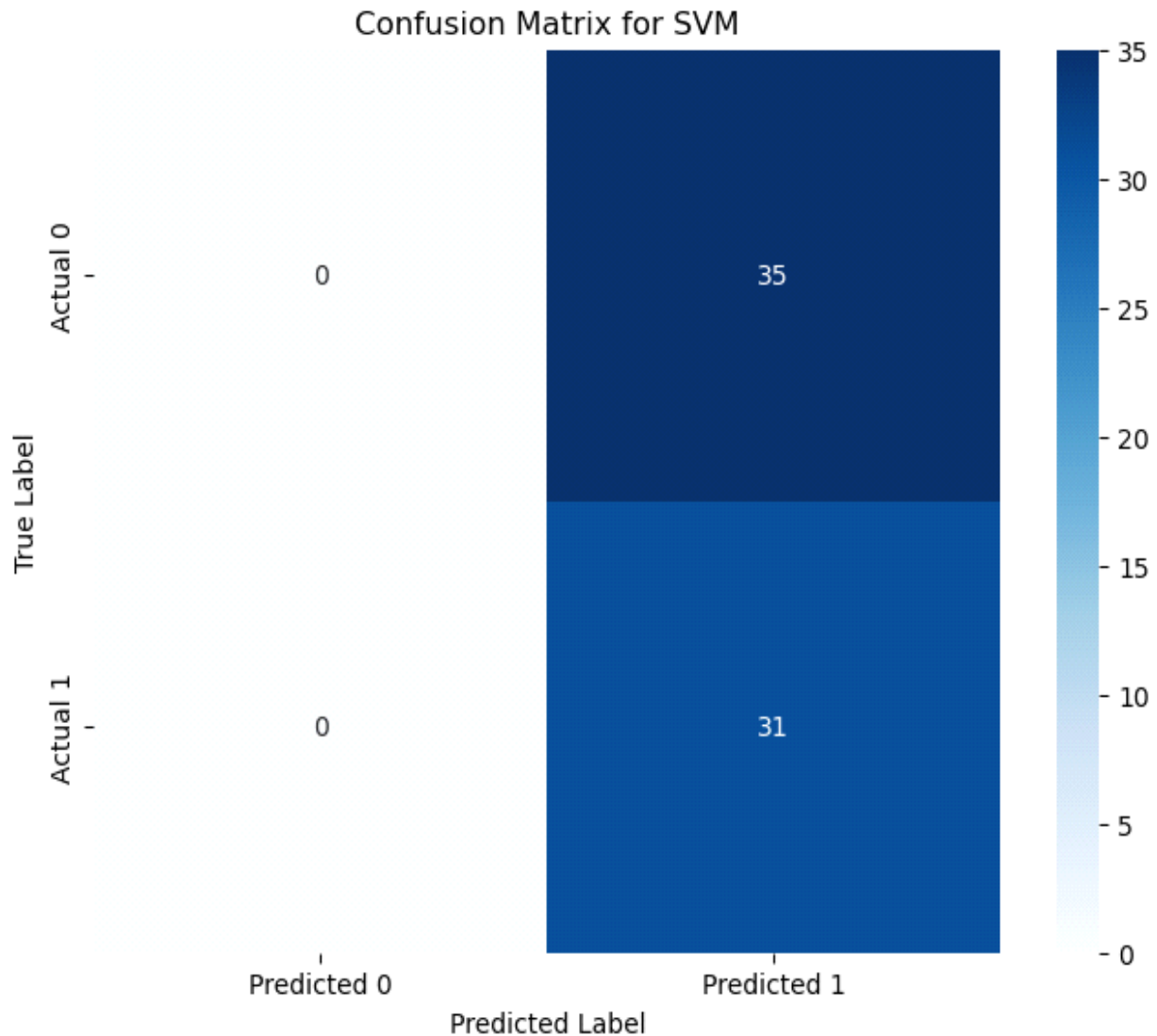


Figure 5: Confusion matrix of SVM based maintenance alert system

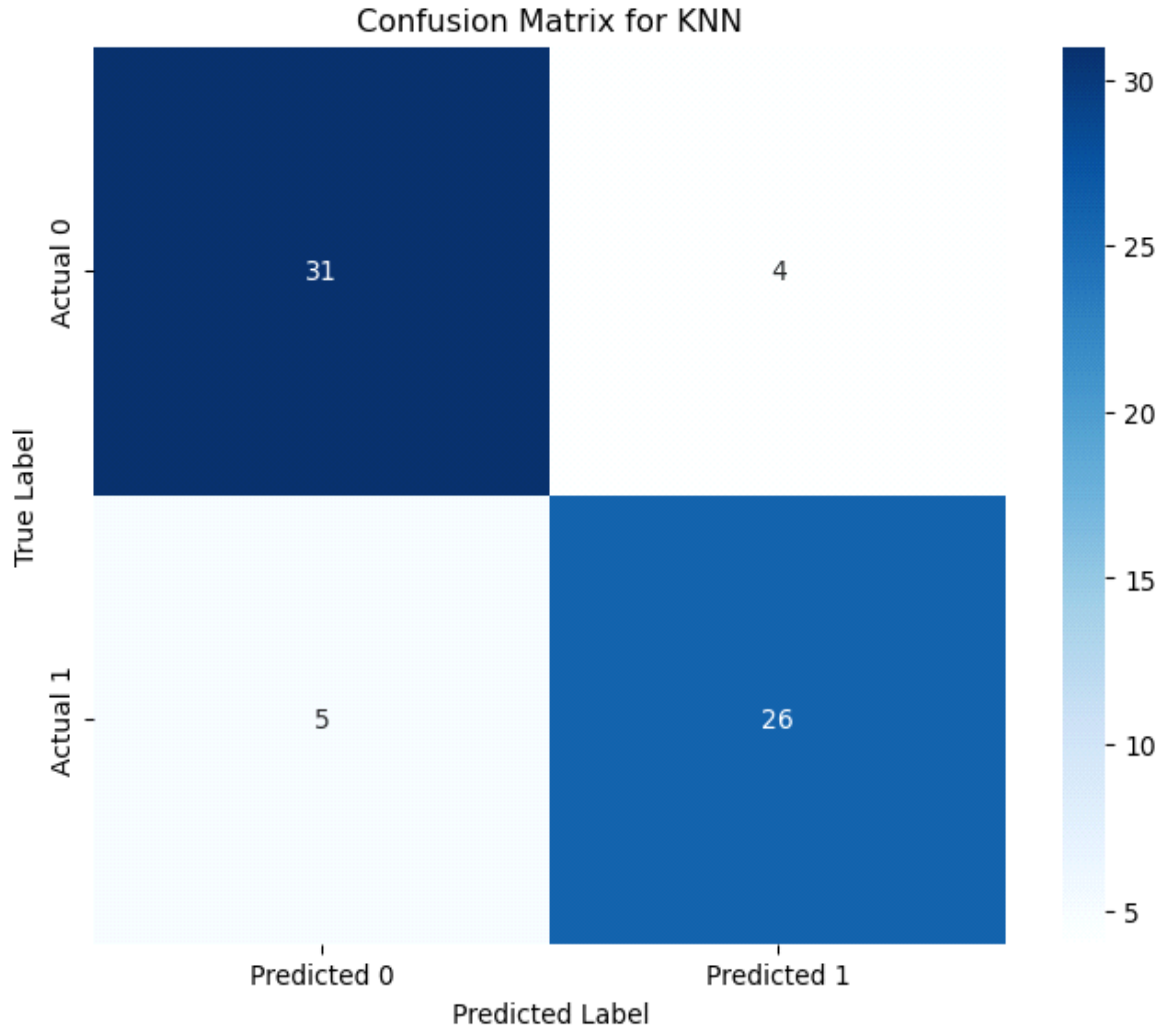


Figure 6: Confusion matrix of KNN based maintenance alert system

Table 3: Performance Comparison of the Models with Outliers and Without Scaling

Metrics	RF	DT	SVM	KNN
Accuracy	100	100	47	86
Precision	100	100	47	87
Recall	100	100	100	84
F1-Score	100	100	64	85

Discussion of Findings

The findings from this study reveal critical insights into the performance and reliability of machine learning models used for predictive maintenance in gas treatment plants, particularly under real-world conditions involving data outliers and without feature scaling. The aim was to test the



robustness and applicability of different algorithms in predicting maintenance needs based on historical metering data.

The results showed that Random Forest (RF) and Decision Tree (DT) models performed exceptionally well, each achieving 100% accuracy in classification. While this may appear ideal, such perfect outcomes are often considered overly optimistic in practical applications, suggesting that the models may be too finely tuned to the specific dataset used and may lack generalizability. K-Nearest Neighbors (KNN) followed with a respectable 86% accuracy, showing moderate resilience to the presence of outliers. On the other hand, Support Vector Machine (SVM) underperformed significantly, with an accuracy of just 47%, highlighting its sensitivity to unscaled and noisy data.

These results are in line with prior studies in the literature. For instance, Cardoso and Luis (2021) emphasized the value of predictive maintenance and the role of machine learning in anticipating equipment failure. In their study, Random Forest and Artificial Neural Networks emerged as top performers when applied to a large industrial dataset, demonstrating their strength in handling complex predictive tasks using real-time telemetry, maintenance logs, and error reports. Similarly, in the present study, RF also showed high predictive capacity, reinforcing its suitability for maintenance tasks where decision accuracy is critical.

Arena et al. (2022) further validated the importance of maintenance as a key factor in ensuring operational continuity and efficiency. Their work on revamping topping plant components showed the effectiveness of using historical maintenance alerts and structured data preprocessing to support failure prediction. Though they adopted a different modeling approach, using regression techniques such as ARIMA, their emphasis on clean, structured data for model accuracy aligns with the present study's recognition that outliers and unscaled features can dramatically impact prediction performance—particularly for models like SVM.

In line with Aliyu et al. (2022), who emphasized the use of diagnostic systems for real-time health monitoring and failure forecasting, this study also supports the notion that a well-structured machine learning model, fed with reliable operational data, can greatly improve predictive maintenance capabilities. The integration of health monitoring data and predictive analytics—as demonstrated in this study—lays the foundation for timely intervention, reducing downtime and extending the lifespan of critical equipment.

Furthermore, Amangeldy et al. (2024) highlighted the advantages of integrating machine learning into oil well monitoring systems via SCADA technologies to enhance predictive maintenance. This supports the argument that intelligent systems embedded within industrial operations can significantly improve efficiency and reduce the risk of unexpected failures. The predictive alert



system developed in this study contributes to this ongoing technological evolution by providing a decision-support tool that can help maintenance teams act proactively.

In summary, the findings confirm that while some machine learning models such as RF and DT show high effectiveness in predicting maintenance needs—even under data challenges—others like SVM are more vulnerable and less reliable without appropriate preprocessing. This underscores the importance of selecting robust algorithms and ensuring that operational data is adequately prepared, especially when deploying predictive maintenance systems in critical infrastructures like gas treatment plants.

Conclusion

This study set out to evaluate the performance of selected machine learning models—Random Forest, Decision Tree, K-Nearest Neighbors, and Support Vector Machine—in predicting maintenance needs in gas treatment plants using operational data containing outliers and without feature scaling. The findings revealed that Random Forest and Decision Tree models performed exceptionally well, achieving perfect classification accuracy, though such results may not reflect real-world conditions due to the risk of overfitting. K-Nearest Neighbors showed moderate accuracy and resilience, while Support Vector Machine was highly affected by unprocessed data, resulting in poor performance. The study underscores the importance of choosing models that are both accurate and robust when applied to real operational data, which is often unstructured and contains irregularities. It also highlights the need for proper data preparation and model evaluation to ensure reliability and practical application. Overall, the results demonstrate the potential of machine learning models, particularly ensemble-based approaches, in supporting predictive maintenance strategies that can reduce downtime, improve efficiency, and enhance the operational reliability of gas treatment plants.

Recommendation

It is recommended that gas treatment plants adopt ensemble-based machine learning models, such as Random Forest and Decision Tree, for predictive maintenance, as they demonstrated strong performance and reliability even under challenging data conditions.



References

- Abidi, M. H., Mohammed, M. K., & Alkhalefah, H. (2022). Predictive Maintenance Planning for Industry 4 . 0 Using Machine Learning for Sustainable Manufacturing.
- Achouch, M., Dimitrova, M., Ziane, K., Karganroudi, S. S., Dhouib, R., Ibrahim, H., & Adda, M. (2022). applied sciences On Predictive Maintenance in Industry 4 . 0 : Overview , Models , and Challenges.
- Al-janabi, Y. T. (2020). An Overview of Corrosion in Oil and Gas Industry : Upstream , Midstream , and Downstream Sectors.
- Aliyu, R., Mokhtar, A. A., & Hussin, H. (2022). Prognostic Health Management of Pumps Using Artificial Intelligence in the Oil and Gas Sector : A Review.
- Amangeldy, B., Tasmurzayev, N., Shinassylov, S., & Mukhanbet, A. (2024). Integrating Machine Learning with Intelligent Control Systems for Flow Rate Forecasting in Oil Well Operations. 343–359.
- Arena, S., Florian, E., Sgarbossa, F., Endre, S., & Zennaro, I. (2024). Engineering Applications of Artificial Intelligence A conceptual framework for machine learning algorithm selection for predictive maintenance. 133(October 2023). <https://doi.org/10.1016/j.engappai.2024.108340>
- Arena, S., Manca, G., Murru, S., Orrù, P. F., Perna, R., & Recupero, D. R. (2022). Data Science Application for Failure Data Management and Failure Prediction in the Oil and Gas Industry : A Case Study.
- Cardoso, D., & Luis, F. (2021). *Application of Predictive Maintenance Concepts Using Artificial Intelligence Tools*.
- Ediger, V. Ş., & Berk, I. (2023). Future availability of natural gas: Can it support sustainable energy transition?. Resources Policy, 85, 103824. <https://doi.org/10.1016/j.resourpol.2023.103824>
- Gao, L., Wang, J., Binama, M., & Li, Q. (2022). The Design and Optimization of Natural Gas Liquefaction Processes : A Review.
- Mohammad, N., Widad, W., Ishak, M., & Mustapa, S. I. (2021). Natural Gas as a Key Alternative Energy Source in Sustainable Renewable Energy Transition : A Mini Review. 9(May), 1–6. <https://doi.org/10.3389/fenrg.2021.625023>
-



- Pandey, Y. N., Rastogi, A., Kainkaryam, S., Bhattacharya, S., & Saputelli, L. (2020). *Machine Learning in the Oil and Gas Industry* *Machine Learning in the Oil and Gas Industry*.
- Poe, W. A., & Mokhatab, S. (2017). Introduction to Natural Gas Processing Plants. In *Modeling, Control, and Optimization of Natural Gas Processing Plants*. <https://doi.org/10.1016/b978-0-12-802961-9.00001-2>
- Silvestri, L., Forcina, A., Introna, V., Santolamazza, A., & Cesarotti, V. (2020). Computers in Industry Maintenance transformation through Industry 4 . 0 technologies : A systematic literature review. *Computers in Industry*, 123, 103335. <https://doi.org/10.1016/j.compind.2020.103335>
- Sircar, A., Yadav, K., Rayavarapu, K., Bist, N., & Oza, H. (2021). Application of machine learning and artificial intelligence in oil and gas industry. *Petroleum Research*, 6(4), 379–391. <https://doi.org/10.1016/j.ptlrs.2021.05.009>
- Ucar, A., Karakose, M., & Kirimca, N. (2024). *Artificial Intelligence for Predictive Maintenance Applications :*
- Wilson, E. F., Taiwo, A. J., Akintola, J. T., & Chineme, O. M. (2023). A Review on the Use of Natural Gas Purification Processes to Enhance Natural A Review on the Use of Natural Gas Purification Processes to Enhance Natural Gas Utilization. *April*. <https://doi.org/10.11648/j.ogce.20231101.13>
- Yang, H., Li, W., & Wang, B. (2021). Joint optimization of preventive maintenance and production scheduling for multi-state production systems based on reinforcement learning. *Reliability Engineering and System Safety*, 214(May), 107713. <https://doi.org/10.1016/j.res.2021.107713>