Journal Pre-proof

Machine Fault Diagnosis Using Random Forest with Recursive Feature Elimination and Cross-Validation

Vetrithangam, Shamik Palit, Anshu Mehta, Gaddam Saranya, Donamol Joseph and Abhinav Pathak

DOI: 10.53759/7669/jmc202505134 Reference: JMC202505134 Journal: Journal of Machine and Computing.

Received 07 March 2025 Revised from 29 April 2025 Accepted 16 June 2025



Please cite this article as: Vetrithangam, Shamik Palit, Anshu Mehta, Gaddam Saranya, Donamol Joseph and Abhinav Pathak, "Machine Fault Diagnosis Using Random Forest with Recursive Feature Elimination and Cross-Validation", Journal of Machine and Computing. (2025). Doi: https:// doi.org/10.53759/7669/jmc202505134.

This PDF file contains an article that has undergone certain improvements after acceptance. These enhancements include the addition of a cover page, metadata, and formatting changes aimed at enhancing readability. However, it is important to note that this version is not considered the final authoritative version of the article.

Prior to its official publication, this version will undergo further stages of refinement, such as copyediting, typesetting, and comprehensive review. These processes are implemented to ensure the article's final form is of the highest quality. The purpose of sharing this version is to offer early visibility of the article's content to readers.

Please be aware that throughout the production process, it is possible that errors or discrepancies may be identified, which could impact the content. Additionally, all legal disclaimers applicable to the journal remain in effect.

© 2025 Published by AnaPub Publications.



Machine Fault Diagnosis Using Random Forest with Recursive Feature Elimination and Cross-Validation

¹Dr.Vetrithangam, ²Shamik Palit, ³Anshu Mehta, ⁴Gaddam Saranya, ⁵Donamol Joseph, ⁶Abhinav Pathak

¹Department of Computer Science & Engineering, Chandigarh University, Punjab, Indi

² Department of Computing Science and Software Engineering, University of Stirling RAK Ca Ras Al Khaimah, United Arab Emirates

> ³Department of Computer Science & Engineering, Chandigarh Un ab, India sity, P

⁴Department of Computer science and Engineering, Narasaraopeta Engineering co re, Narasaraopeta, India

⁵ Department of Computer Applications, Marian College Kuttikkanam A nomous, Kerala, India.

⁶Symbiosis Institute of Computer Studies and Research (SICSR), Symb tional (Deemed University), Pune, Maharashtra, 4110

¹vetrigold@gmail.com, ²shamik1980@gmail.co, ³

⁵donamol.joseph@mari

⁶abhina gsits7@gmail.com sollege

@gmail.com, ⁴gaddamsaranya4@gmail.com,

D.Vetrithangam: vetrigold@gmail.com Correspondence should be addressed

Abstract - In modern industrial environmer early and accurate machine fault diagnosis is crucial for minimizing ensuring operational safety. This research presents a robust fault classification downtime, reducing maintenance costs, and on with Cross-Validation (RFECV) and Random Forest classifiers framework that combines Recursive F .mna to address the challenges of high dim asionality, or rfitting, and limited model generalization. The proposed approach begins with comprehensive data precessing, followed by RFECV to identify and retain the most relevant features, thereby enhancing model effiand curacy. Subsequently, a Random Forest classifier is trained on this optimized feature set to classify four fa o Fa ure, Power Failure, Tool Wear Failure, and Overstrain Failure. By integrating types feature selection with ensem e learnin the framework effectively mitigates high variance and improves robustness under distributions. Experimental results demonstrate that the proposed methodology varying operational nd d tio achieves a high curacy of 99.2% along with improved computational efficiency, making it highly suitable for real-time fault di tions in smart manufacturing systems. sis app

Diagnosis, Random Forest, Recursive Feature Elimination (RFECV), Feature Selection, Keywords ne Fa Predictiv

I. INTRODUCTION

In today's rapidly evolving industrial landscape, the integration of intelligent manufacturing systems has become a cornerstone for achieving operational excellence and competitive advantage. As industries increasingly embrace automation, the deployment of embedded sensors and condition-monitoring technologies has revolutionized how machines are monitored and maintained [1]. Predictive maintenance and fault diagnosis have emerged as essential components within this paradigm, enabling organizations to anticipate equipment failures before they occur, thereby minimizing downtime, reducing maintenance costs, and enhancing safety standards [2]. This shift from traditional reactive or scheduled maintenance to proactive and condition-based approaches relies heavily on advanced data-driven methods capable of extracting meaningful insights from vast amounts of sensor data. Using machine learning, it is now possible to analyze historical equipment data and identify complex patterns that change the way fault diagnosis takes place [3]. Machine learning techniques are widely appreciated for their simplicity, reliability, and fast training capabilities, making them suitable for diagnosing relatively simple systems. Deep Learning approaches, on the other hand, offer powerful end-to-end solutions capable of managing complex systems and compound faults, especially whe large training datasets are available. Transfer learning methods address the critical issues of data scarcity and sample imbalance by enabling knowledge transfer across different operating conditions, machines, or even application domains. Despite these advancements, the implementation of machine learning in real-world fault diagnosis continues to face challenges, particularly as engineering systems grow in complexity [4]. When using knowledge from t and typical situations, ML models become very accurate in detecting small errors and predicting potential malf ctions [5]. The Extreme Random Forest (ERF) method was introduced to enhance feature extraction capabilitie hile reducing computational complexity. In this approach, high-dimensional data is projected into a lower-di using a randomly generated mapping matrix, effectively reducing dimensionality. This process no bnlv lov the computational burden but also improves classification performance after dimensionality red 51 Pre ctive capabilities make it easier for manufacturing lines to change their maintenance processes and lo fter eq more carefully. Even so, using fault diagnosis models in real factories is still very diffigure fte cessary for those using older ML methods to have broad experience in the field and face probl with nputing [7][8]. licient Although deep learning is effective at dealing with difficult and complex data, demands a lot of labelled sual samples, has severe computational needs, and remains unclear for users in terms of t standing how AI affects their operations.

ng settings, and incomplete or noisy In addition, industrial systems have many types of equipment, diff components. As a result of these factors, shifts in the data between h I learns and its use in practice make the the model perform poorly when it meets new or evolving errors. ed issue to maintain fault diagnostic models that are strong, expandable, and responsive on the spot nitations [9]. Including a large number n with of instrument measurements in high-dimensional d redundant or irrelevant information, which may ead reduce the model's accuracy and increase comput 10]. For this reason, Random Forests and other onal rec reme ensemble methods are used widely since they trees together, manage data that contains thousands of multi In Forest models may still be affected by the problem of too features, and give feature importance scores. Yet, I many variables and some of these may not matter for sp ing faults. For this reason, using RFECV enables you to find the best subset of features step by step, removing unimportation ones as it goes, and constantly checks its effects on the model to prevent it from overfitting [1] Random Forest (RF) is a robust ensemble learning method that constructs o improve classification accuracy and model stability compared to multiple decision trees and combines t a single decision tree. Although n techni ues have been explored for fault diagnosis, RF remains a valuable ero peed, ability to handle high-dimensional data, and consistently strong and necessary approach due to its t execu performance in machinery f diag sis tasks[12]. Based on what this research learns, it suggests a strong machine dom Forest classification together with RFECV-based feature selection to fault diagnosis framewoy that learer, and use resources more efficiently. In this system, we aim to separate four enhance accuracy, make rediction Too (No Failure, Pow Failt Vear Failure, and Overstrain Failure) basic failure types that often come up in industrial i es through careful preprocessing of the senses of vibration, torque, the time worked, and ten After the framework uses performance measures from cross-validation to help it remove perat as it progresses. Using the new tools in this feature set, the model is able to ensure accurate and unneces feati prompt f faults. iden cation

beides increasing the correctness of fault classification, the strategy also makes it easier to implement the model in reviewe since it cuts down on model complexity and processing power. That's why intelligent manufacturing setting are excellent places to use it due to the quick decisions and ability to withstand changes in its surroundings. The aim through this research project is to give industry a solid, scalable, and clear method for diagnosing machine faults, so industrial operations become both safer and more efficient than before. The proposed methodology aims to implement Recursive Feature Elimination with Cross-Validation (RFECV) to effectively select the most significant features from the available dataset, which helps reduce dimensionality and enhances both the efficiency and accuracy of the fault diagnosis model. Building on this, a Random Forest classifier is developed and trained using the optimally selected features to accurately classify machine fault types, including No Failure, Power Failure, Tool Wear Failure, and Overstrain Failure, while addressing issues such as overfitting and improving model generalization. Furthermore, this approach tackles the challenges of high variance and limited robustness found in existing machine fault diagnosis

methods by integrating feature selection with ensemble learning techniques, thereby ensuring reliable fault prediction across diverse operating conditions and varying data distributions.

II. LITERATURE REVIEW

Zhao et al.[13] proposed a novel framework named Identification for Fault Diagnosis (I4FD) that integrate regularized data-driven modeling and frequency analysis for machinery fault diagnosis under nonlinear system identification. The framework is designed to mitigate the effects of external environmental changes and in diagnostic accuracy. It introduces a fault diagnosis-oriented regularization (FDoR) technique that incorpora prior physical knowledge through a penalty parameter, making the model specifically tailored for fault d osis applications. Unlike traditional approaches, I4FD supports continuous dynamic modeling using upd model identification, frequency analysis is applied to extract fault-sensitive features. The frame ork ach s an accuracy of 92% on simulation and real-world cases. The advantage of I4FD is its ability t to d amic environments and deliver high accuracy, while a technical gap lies in the computation olexity tential C tection Algorithm tuning challenges of the regularization process. Bode et al. [14] proposed a data ult (FDA) for heat pump systems, addressing the issue of reduced energy efficience al sys failures due to and pote undetected faults in building heating and cooling systems. The model lev data approaches and AI techniques, using features extracted from a comprehensive fault dataset provi by the National Institute of Standards and Technology (NIST). The FDA is trained on this lab-generated data a hen applied to a real-world air-water heat pump system without system modifications. The model achieve in accuracy of 85% on the NIST etailed fault feature analysis from longdataset. The advantage of this approach lies in its cost-efficiency and term monitoring data, which avoids the need for expensive custom vever, the model performs poorly on etun ain shift, data incompleteness, and real-world data, highlighting a technical gap in generalizabi do inadequate fault labeling in practical applications. Brito et al ovel unsupervised framework for fault] proj detection and diagnosis in rotating machinery, addr enge of limited labeled data and the need for model e c interpretability. The approach consists of three extraction (from vibration signals in time and an mod es: fe frequency domains), anomaly-based fault deter ult diagnosis using SHAP for model explainability. n, and To diagnose faults, the model leverages feature in ince scores from SHAP explanations, enabling unsupervised classification and root cause analysis. The proposed thodology demonstrated its effectiveness on three rotating machinery datasets, achieving a maximum unsupervi d classification accuracy of 96.72%, particularly with Ensemble, kNN, and CBLOF algorith The advantages of the proposed model include modularity in algorithm selection, interpretability using SHA accuracy without requiring labeled data. Nevertheless, the weak points in the area are that useful lepends on the quality of the features, and methods such as SHAP me too ally den and Local-DIFFI are computation

e learning model to detect and diagnose faults in real time Chen et al. [16] mac in brushless motors, S pr Machines (SVM), Neural Networks (NN), and Random Forests are used (RF). port Vel It collects and c ne orma on from numerous sensors to spot faults and check their degree of severity. offering mects. Experiments prove that NN comes out on top in terms of success rate. SVM and ideas on the RF p ly, each having an accuracy of 95% and 92% respectively, while the best performance was form ery sim 27%. The main benefit of this method is that it improves the reliability, efficiency, and maintenance given f brushless motors in industries. Still, there is a gap in technology when it comes to joining these condition he us implement models in industries, considering they have to work continuously adapting to new the to Its as they happen. Tang et al. [17] proposed an intelligent fault detection system that uses DL for rotating types c hinery that involves bearings, gears and gearboxes, and pumps. The framework tries to find ways to the major problems linked to expert-dependent traditional faults diagnosis methods finding solutions by ov only knowledge and manual work. With the help of deep learning, the framework lets users the automatic discovery of useful features and accurate recognition of types of faults. The model manages to reach an accuracy of 97.75%. It is an effective way to do extract features, since it reduces the amount of manual work. An intervention makes diagnostics more reliable and improves their consistency. However, it faces challenges in generalization, realtime application, and adaptability to unseen fault types, which are highlighted as areas for future research. Gonzalez-Jimenez et al. [18] proposed a machine learning-based fault diagnosis strategy for detecting power connection failures in induction machines, such as high resistance connections (HRC), single phasing faults, and opposite wiring connections. The model is designed to aid maintenance personnel in identifying these faults, particularly those caused by human errors during assembly. Due to the scarcity of real-world failure data, a simulation-driven approach using Software-in-the-Loop (SiL) simulations was adopted to generate synthetic training data. The proposed system achieved an accuracy of 98.5%. Using this approach, it's possible to identify a range of faults even without using real data. Its disadvantage is its dependence on simulations, which may decrease its effective use in real industries.

Tran et al. [19] proposed an IoT-based architecture integrated with machine learning algorithms to enhance cybersecurity in cyber-physical systems (CPS) for industrial electrical machines. The architecture focuses monitoring induction motor status and detecting cyber-attacks in real time. The system uses the Random Fores algorithm for fault detection due to vibration and cyber-attack recognition, achieving an accuracy of 99.03%, which outperforms other ML models in industrial conditions. The infrastructure leverages the CONTACT Element platform to visualize motor faults and fake data signals triggered by detected cyber-attacks on a dashboa The advantage of this model lies in its high detection accuracy, low latency, and clear visualization, making it suit e for cost-effective and secure remote monitoring. However, technical gaps remain in terms of scalabilit industrial networks and robustness under varying attack types. Shubita et al. [20] proposed a machine learnin ased fault diagnosis system that uses acoustic emission (AE) signals for early fault detection in re nachin The system is implemented on an embedded device with IoT connectivity, enabling fault and classification. It achieved an accuracy of 96.1% using a fine decision tree model of this approach is nta its ability to provide accurate and real-time monitoring with minimal later makin for industrial it suit deployment. However, the technical gap lies in the limited exploration of model under varying operational ustr or noisy conditions, which may affect real-world generalization.

Siyuan et al. [21] proposed a duplet classification model combining two Convictional Neural Networks (CNNs) for fault diagnosis in rotating machinery involving both rot ing components. The idea involved through the model was constructed by working on a dataset of 48 m problems created by different faults hine different levels and types of these two parties. CNN architectur d to distinguish between rotor and crea out g having the ability to respond to various external problems w maged. It was possible to achieve the model. A high rate of identifying mixed faults at that the results are highly reliable. Moreover, a pr single-vs-rest approach was built based on CN inform ch known diseases. Four new fault categories, on to this st including those that go unnoticed, were tested y. Its usefulness comes from the fact that it is felt in many parts of society the ability to work in complicated conments and recognize new types of faults. However, there is a technical challenge as using different models for ch type of fault may increase the overall model. Real-time situations can cause major challenges due to lots of calcundons involved. Shao et al. [22] introduced a fault diagnosis method that depends on deep learning (JBN) to detect the main status of induction motors by examining the distribution of their vibration signals (he s made by putting several Restricted Boltzmann Machines (RBMs) on top of each other and training mbines the steps of extracting features and doing the classification ers. It in into one approach, so you do f keer features manually. On data from the machine fault simulator, the ave to accuracy of the classificati Q

Because this way work with tw intermation, the model can learn to structure the data and make the process of finding issues automate and small. Still, getting the right performance from the model requires careful selection of scale and depth, due to which turing hyperparameters and running the model can be difficult.

[23] posed a fault diagnosis method that combines a two-layer bearing with a hybrid set of data, haib N. The model deals with finding patterns of faults and measurements of crack sizes from vibration along SAL signals th hange ith changing conditions and various fault levels in machines. It is more accurate than SVMs and BPI with curacy of 99.10%. Its real benefit is that it helps find more important features in the vibration data, making asier to classify sounds under changing conditions. Kafeel et al. [24] proposed a fault detection method hachines by studying the vibration signals. This system performs empirical mode decomposition (EMD) rotati ise from the signals and does multi-domain feature extraction to find both the time and frequency features to vibration data collected from healthy and bad induction motors. The extracted features are classified using multiple algorithms including SVM, KNN, Decision Tree, and Linear Discriminant Analysis, with the support vector machine achieving best of 98.2% using а Gaussian kernel the performance accuracy. The advantage of this method lies in the hybrid use of time and frequency features, which enhances the fault discriminative capability of the model. However, a technical gap remains in the generalization of the system across different machine types and operational conditions, which could affect its applicability in broader industrial settings. Hung et al. [25] proposed a system-on-chip (SoC)-based tool wear detection model that leverages deep learning with sensor fusion techniques. The system was trained using vibrational and acoustic signals collected from a three-axis CNC machine operating under various spindle speeds and torque conditions. The inputs to the deep learning model

were frequency spectrum representations of signals from a MEMS microphone and a three-axial accelerometer, with tool flank wear measured via a camera, adhering to ISO 8688-2:1989 standards. The model achieved detection accuracies of 99.7% for the single-sensor model and 87.75% for the fused model when deployed on a Pocket Beagle SoC.

The advantage of this system lies in its real-time detection capability, high accuracy, and cost-efficient embedded implementation. However, it shows reduced performance in the fused model, possibly due to signal integration complexity or variability in machining conditions, indicating a need for more robust fusion strategies. Orrù et al. [26] proposed a simple and easy-to-implement machine learning (ML) model for early fault prediction of centrifugal pumps in the oil and gas industry. The model is based on real-life sensor data including temperature, pressu vibration readings, which are pre-processed and denoised before training. Two algorithms-Support Vector I ichine (SVM) and Multilayer Perceptron (MLP)—were implemented using the KNIME platform. The model achieved d an accuracy of 98.1%, successfully detecting system deviations and issuing fault prediction alerts. The a approach lies in its practical simplicity and effective performance using real industrial data, support ance ng main decision-making. However, the model is still in a preliminary stage, and potential technical ga e the n d for broader validation across different operating conditions and scalability for more complex fat cenario

| Author | Proposed Model | Findings | Challenges |
|-------------------|---|---|--|
| Zhao et al.[13] | Identification for Fault Diagnosis (I4FD) | Achieved 92% accuracy in machinery fault diameters by integrating regular ed MARX modeling and frequency analysis; incorporates of visical powledge via tooR of or commutous dynamic podelar | Computational complexity and tuning difficulties in regularization parameters. |
| Bode et al. [14] | Data-driven Fault Detection Algorithm (FDA) | chieved 85% curacy in a section faults in heat pump systems using AI-based FDAs trained in NIST laboratory data; enabled cansfer to real-world systems without hardware nodifications; leveraged big data and feature extraction for energy- efficient building climate systems. | Poor generalization to real- world data due to domain shift, incomplete data, and fault labeling issues. |
| Brito et al. [15] | Unseperviced Framework Gor Fault Detection an Diagnosis Rotein Cachinery | Achieved 96.72% accuracy in unsupervised classification using Ensemble, kNN, and CBLOF; employs SHAP-based explainability for root cause analysis; effective across three real-world rotating machinery datasets. | Computational cost of interpretability methods (e.g., SHAP, Local-DIFFI); performance sensitivity to the quality of extracted vibration features. |
| Chen en 1 [16] | ML-based fault diagnosis using SVM, NN, and RF | Achieved 97% accuracy with NN, 95% with SVM, and 92% with RF; effectively analyzes fault severity and suggests countermeasures using sensor data. | Real-time integration challenges and limited adaptability to evolving fault patterns. The model faces a high variance issue as it struggles to validate on unseen faults |
| Tang et al. [17] | Deep Learning-Based Intelligent Fault Diagnosis Framework | Achieved 97.75% accuracy in fault classification for rotating machinery components (bearings, gears, pumps) by enabling automatic feature | Generalization issues, real- time implementation constraints, and difficulty adapting to unseen fault types. |

Table 1: A Review of Research on Machine Fault Diagnosis echoques

| | | learning and reducing reliance on manual feature extraction. | The model might lead to overfitting with increased epochs | |
|--------------------------------------|--|---|--|---|
| Gonzalez-Jimenez et al. [18] | ML-Based Fault Diagnosis for Power Connections in IMs | Achieved 98.5% accuracy in diagnosing power connection faults (HRC, single phasing, and opposite wiring) using Software- in-the-Loop (SiL) simulation- generated training data. | Dependency on simulated data may limit real-world generalizability; lacks validation with field datasets. | K |
| Tran et al. [19] | IoT-based architecture with integrated ML (Random Forest) for CPS security and motor fault detection | Achieved 99.03% accuracy in detecting induction motor faults and cyber-attacks using Random Forest; leverages CONTACT Element IoT platform for real- time visualization of motor status and cyber-attack data; offers low latency, high detection accuracy, and clear dashboards. | Scalability access heterogeneous industrial networks and robustrias under diverse acuck scenarios retain open issues. | |
| Shubita et al. [20] | ML-based Fault Diagnosis System using AE on IoT-Enabled Device | Achieved 96.1% accuracy early fault detection of rotating machines using AE signals; implemented on embedded IoT device for real-time monitoring | Limited robustness under roying operational/noisy a ditions; lacks gene lization to real-world environments. | |
| Chen Siyuan et al.[21] | Duplet Classifier using two 1-D CNNs | Achieved 95.93% ccurate in diagnosing mixed facts on rotating machinery; cares to o parallel CNN to o diagnost otor and bearing rults separately; unidated in 48 highine health indition and four new fault type | Increased model complexity due to separate CNNs; computational overhead during real-time deployment. | |
| Shao et al. [22] | Deep Belief Network (DBN)-based Fault Diagnosis | Achieved 99% accuracy in fault diagnosition finduction motors by automatically learning features from vibration signal frequency stributions. Combines feature extraction and classification in a unified deep learning framework using stacked RBMs. | Model performance depends heavily on architecture scale and depth; introduces challenges in hyperparameter tuning and computational complexity. | |
| Kafeel et al. [24] | Fault detection System based on Hybrid machine learn, is more as | The hybrid use of time and frequency features, which enhances the fault discriminative capability of the model | Generalization of the system across different machine types and operational conditions | |
| Hung et al. [23] Orrù et al. [15] | Deep learning with sensor fusion apport Vector Machine | This system provides real-time detection capability. Detecting system deviations and | This model faces integration capability issues This model faces challenges | |
| S | (SVM) | issuing fault prediction alerts" | in broader validation across different operating conditions and in scaling to more complex fault scenarios | |

As shown in table 1, the existing fault diagnosis models face several technical challenges, including high computational complexity and difficulties in tuning regularization and hyperparameters. Many models struggle with generalization issues, particularly when validating on unseen fault types or transferring from simulated or laboratory data to real-world scenarios, often due to domain shifts and incomplete or noisy data. Real-time implementation and integration remain problematic, especially for deep learning and ensemble methods with increased model complexity and computational overhead. Industry experts are also very concerned about the ability to scale these networks in many settings and how

they will handle ever-changing threats. Furthermore, knowing how the AI model works is helpful, but it contributes to the model's complexity, and there are usually difficulties for models to maintain their results as faults evolve and work in more types of environments. This research addresses the technical gaps of high variance and overfitting commonly observed in machine fault diagnosis models, focusing on improving robustness and generalization in Random Forest-based predictive maintenance.

III. PROPOSED METHODOLOGY

This section describes the proposed methodology illustrated in the figure 1, which presents a structured method for machine fault classification using a machine learning approach. The process begins with data preprocessing vhich includes steps such as dropping irrelevant columns, label encoding of categorical data, feature and tar and finally, a train-test split to prepare the dataset for modelling. Following preprocessing, a feature se tion is applied using Recursive Feature Elimination with Cross-Validation (RFECV) to identify and retained only th nost significant features, thereby improving both model efficiency and accuracy. The selected feature train used a Random Forest Classifier, a robust ensemble learning algorithm known for its accu erfitting. ience to The classifier is trained to predict different types of equipment failures. For any ne nput i ance e model predicts one of the four possible outcomes: No Failure, Power Failure, Tool Wear Failur ain Failue, thus enabling or Over proactive maintenance and minimizing operational downtime.



1)Input Dataset: Let the raw input dataset be represented by Equation(1)

$$D = \{(x_1, x_1), (x_2, x_2), \dots, (x_N, x_N)\}$$
(1)

where x_i is a vector of features for the i^{th} instance, and y_i is its corresponding label. The dataset has N instances and M initial features.

2)Dropping Irrelevant Columns: As shown in equation (2). This step aims to remove features that do not contribute to the predictive power of the model. denote the set of irrelevant feature indices. After removing these columns, the dataset is transformed into a new feature set, as represented in Equation (2).

$$F_{\text{relevant}} = \{j \mid j \in /F_{\text{irrelevant}}\}$$
(2)

$$D' = \{(x_1', x_1), (x_2', x_2), \dots, (x_N', x_N)\}$$

where x_i' is x_i with columns in $F_{irrelevant}$ removed. As shown in Equation (3), the updated dataset 1 consists of input-output pairs where each x_i' is derived from the original feature vector x_i by excluding the feature adexed in $F_{irrelevant}$. This results in a reduced-dimensional representation that retains on the post 1 want features for model training.

3)Label Encoding: If the target variable is categorical, it needs to be converted up numerical representations. Let $Y = \{y_1, y_2, y_3, ..., y_N (4)$ be the set of original categorical labels. As represented in equation(4), the label encoding maps these to numerical values: L:Y \rightarrow {0,1,2,3} (e.g., "No Failure" \rightarrow , "Power Failure" \rightarrow 1, etc.).

The transformed dataset now has numerical labels:

$$D'' = \{(x_1'', l_1), (x_2'', \dots, u_n'', l_N)\}$$
 (4), where li=L(yi).

(3)

4)Feature and Target Separation: The processes dataset is split into features (X) and the target variable (y). $X = \{x_1'', x_2'', ..., x_N''\}$ (matrix of features) y. $\{x_1, ..., l_N\}$ (vector of target labels)

5) Train-Test Split: The dataset is divided into training and testing sets. Let D_{train}'' and D_{test}'' be the training and testing sets. $D_{train}''=(X_{train}, y_{train})$ D_{test}'' T_{test}, y_{test}

B. Feature Selection

This stage identifies the most relevant set of features.

1)Recursive Feature Elimination with Cross-Validation (RFECV)

RFEC are arsive fits a model and removes the weakest features until the optimal number of features is reached a red on crite-validation performance. Let M_{model} be the base machine learning model. Let K be the number of holes for cross-validation. The process can be described as follows:

🛋 : Initialization

Start with the full set of P features, $F = \{f_1, f_2, ..., f_P\}$.

Step 2 : Iteration

The model is trained on the current feature set FFF using K-fold cross-validation applied to the training data X_{train} . During this process, the model's performance measured using metrics such as accuracy or F1-score—is evaluated on each fold. Let S_k represent the score obtained on fold k, and the average score across all folds is calculated as represented in equation (5).

$$\dot{\mathbf{s}} = \frac{1}{\kappa} \sum_{k=1}^{K} \mathbf{S}_{k} \tag{5}$$

After evaluating performance, the feature with the lowest importance, denoted as $f_{weakest}$, is identified and removed from the feature set F. This iterative process continues to refine the model by eliminating the least significant features.

Step 3 :Recursion

Repeat step 2 until an optimal performance is observed or a minimum number of features is reached.

Step 4 :Optimal Feature Set Selection

Select the feature set $F_{selected}$ that yields the highest average cross-validation score. The dataset is the projected onto this selected feature set: $X_{train}'=X_{train}[F_{selected}] X_{test}'=X_{test}[F_{selected}]$

Step 5: Feature Set

The output of the feature selection phase is the reduced set of features, F_s

C. Random Forest Classifier

The selected features are fed into a Random Forest Classifier for predicting the failure ty,

1)Random Forest (RF): An ensemble learning method that constructs a multiplication trees.

Let T be the number of decision trees in the forest. Each tree $t \in [1, ..., T_n]$ is grained as follows:

Step 1 : Bootstrap Aggregating (Bagging)

A random subset of the training data x_{rain}' (with eplacement) is sampled to train each tree. Let this sample be $Dt'=(X_{train},t',y_{train},t)$.

Step 2: Random Feature Subspace

At each node of the decision the, only a random subset of m features is considered for splitting.

Step 3: Tree Construction A decision tr e T_t is grown on Dt'.

Step 4: Training

The Random For t model lenoted as RF, is trained on the selected features of the training data:

Step 5. rediction

RF

For a tw, unseen instance x_{new} from X_{test}' (with features corresponding to $F_{selected}$), each tree t in the forest predicts a y^{hew} =mode(y^{hew} =mode(

ep 6 : Output Classes

The model outputs one of the four predefined failure types: "No Failure", "Power Failure", "Tool Wear Failure", "Overstrain Failure".

IV. RESULTS AND DISCUSSION

A. Dataset Description

The dataset used in this study contains detailed information related to engine performance and failure analysis. It includes variables such as vibration levels, torque, process temperature, air temperature (in Kelvin), engine speed (in RPM), and operational hours. Each entry is uniquely identified by a UDI (Unique Identifier) and is associated with a specific Product ID and engine type, where the type may denote categories such as motor (M) or liquid (L). The dataset also records the type of failure (if any), including specific classifications such as rotational failures, across a total of 500 machines. These attributes enable a comprehensive analysis of engine behavior under varying operational conditions. It can be used in many ways for example, spotting reasons for engine failure, checking for engine temperature, speed, and torque, examining various ngine types, and making forecasts for maintenance. The dataset is available at the following source link: https://www.kaggle.com/datasets/nair26/predictive-maintenance-of-machines. The dataset is split into 75% training and 25% testing.





The figure 2 represents the prediction results of Random Forest classifier with Recursive Feature Elimination and Cross-Validation (RFRFECV) on multi-class machine failure problems. The four labels tested in the model were No Failure, Overstrain Failure, Power Failure, and Tool Wear Failure. It classified 115 instances as No Failure and just one was ruled as Tool Wear Failure. All three cases of Overstrain Failure were grouped under the correct class with 100% correctness. No errors happened in the prediction of Power Failure, as all two instances were accurately classified, and although four instances of Tool Wear Failure were found, the model misclassified one as being from the No Failure class. On the whole, the confusion matrix confirm that the main class is classified very accurately and that all failure categories

are detected well. The findings prove that choosing the right features and training the model correctly worked well. The slight number of cases that were wrongly classified implies that some failure groups may have traits in common with others. Therefore, RFRFECV was a dependable choice for handling data from many types of machinery and for recognizing faults in machines with preventive measures.





The figure 3 illustrates how a Ran classifier worked well when it was trained using RFRFECV to predict om F categories used in this classification problem: No Failure, Overstrain multiple machine failure conditions re fou Failure, Power Failure, and Tool r Failur is model was able to identify 115 of the instances in the "No Failure" category and just one case w marked as involving "Tool Wear Failure." When it comes to the "Overstrain ron Failure" category, the syste any mistakes and identified all the instances correctly. Thus, the model has the did n nal ability to tell between rou e condit ns and certain types of failures. Consequently, the RFRFECV method allowed the his improved the model's precision in spotting and classifying different machine team to pick the failures.



Feature Relationships of Failure Types in the Proposed Model

F

re Relationships that are colored by the type of failure that occurs. Scatter plots in the matrix The Figure 4 ents Fe display the between two variables in the data, and the diagonal plots indicate how each variable is spread. No ection Failure is depi Power Failure in orange, Tool Wear Failure in green, and Overstrain Failure in red on every l as b. graph. The nal and environmental features are related, color-coded according to the 'Failure Type'. Seeing the iy op possible to spot connections between various features and the different forms of failure. KDE estimates are rest make draw the d conal figures, offering views of the distributions of the features individually. Each scatter plot in the offhows the trend between two features. The analysis using pair plots explains the features' distributions alone iagonal relationships with one another, as well as the significant patterns spotted for each failure case. The feature called Type' is a category, with 'Type 1' occurring most often, while both 'Air Temperature' and 'Process Temperature' are narrowly distributed and only take values inside certain ranges, but 'Air Temperature' sometimes drops below these ranges. 'Rotational speed' displays multiple peaks, suggesting varied operating regimes, whereas 'Torque' and 'Vibration Levels' demonstrate unimodal distributions concentrated at lower values with a tail extending to higher levels. 'Operational Hours' presents a broader distribution, with a noticeable peak at lower values potentially indicating newer units or shorter operational cycles. In terms of bivariate relationships, a strong inverse correlation exists between 'Rotational Speed' and 'Torque', where increased rotational speed generally corresponds to decreased torque, a typical characteristic of mechanical systems with

constant power output. No direct linear relationship is evident between 'Operational Hours' and either 'Torque' or 'Rotational Speed across the entire dataset, although specific failure types might exhibit localized clustering. 'Air' and 'Process temperatures' show an expected correlation with each other, but their relationships with other operational parameters like 'Torque' or 'Rotational Speed' are less pronounced linearly. Similarly, 'Vibration Levels' show scatter with other features, but no strong linear correlations are immediately apparent across the dataset. Crucially, the coloring by 'Failure Type' illuminate key patterns: 'No Failure' instances, representing the majority, are broadly distributed across all features, forming the primate clusters. 'Power Failure' instances are fewer and tend to cluster in specific regions, such as higher torque values at varying operational hours, or lower rotational speeds combined with higher torque, potentially indicating overload conditions. 'Tota' Wear Failure' events are sparse but more prominent at higher operational hours, consistent with accumulated wear, and also appear at higher 'Vibration Levels', a common symptom of tool degradation. Finally, 'Overstrain Failure' events are vary rare and consistently occur at extremely high 'Torque' values, aligning with the definition of overstrain.

| Author | Proposed Model | Accuracy |
|----------------------------|---|----------|
| Zhao et al.[] | Identification for Fault Diagnosis (I4FD) | 92% |
| Bode et al. [] | Data-driven Fault Detection Algorithm (FDA) | 85% |
| Brito et al. [] | Unsupervised Framework for Fault Detection and Diagnosis Rotating Machinery | 96.72% |
| | | 95% |
| Chen et al. [] | ML-based fault diagnosis using SVM-NI and RF | 97% |
| | | 92% |
| Tang et al. [] | Deep Learning-Based Vielligel Fault Fagnosis Framework | 97.50% |
| Gonzalez-Jimenez et al. [] | ML-Based Fault Diagnois for Power Connections in IMs | 98.50% |
| Tran et al. [] | IoT-based architecture with stegrated ML (Random Forest) for CPS security and motor fault detection | 99.03% |
| Shubita et al. [] | ML-based Lult Diagnosis System using AE on IoT-Enabled Device | 96.10% |
| Chen Siyuan et al.[] | Duple: Classifier using two 1-D CNNs | 95.93% |
| Shao et al. [] | Deep Brief Network (DBN)-based Fault Diagnosis | 99% |
| Kafeel et al. [] | Fault steetion system based on Hybrid machine learning models | 98.20% |
| Orrù et al. [] | Suppor Vector Machine (SVM) | 98.10% |
| Proposed mot | RFRFECV Classifier | 99.20% |

| | Table 2 : Accuracy-based | performance of the | proposed model |
|--|--------------------------|--------------------|----------------|
|--|--------------------------|--------------------|----------------|

Table 2 des formance (in terms of accuracy) of existing fault diagnosis methods used across various domains achiney, induction motors, and cyber-physical systems. Models range from conventional machine such as learning te A as Support Vector Machines (SVM) and Random Forest (RF) to deep learning-based frameworks iques Networks (DBN) and 1-D CNN classifiers. The accuracy ranges from 85% to 99.2%, showing significant like ep Bel ne. The figure 5 visually illustrates the accuracy performance of each fault diagnosis model, with each bar progre ed to distinguish between different techniques. The RFRFECV Classifier, proposed in this study, achieves the uuely curacy of 99.20%, outperforming all other existing approaches. Notably, models such as the IoT-based architecture with integrated machine learning (99.03%), Deep Belief Network (99%), and Hybrid ML models (98.20%) also demonstrate strong performance, reflecting a clear trend toward the adoption of hybrid and deep learning-based solutions for fault liagnosis.



This study presents a comprehensive machine fault diagnosis framework that effectively combines Recursive Feature) and Pendom Forest classification to enhance predictive accuracy and model Elimination with Cross-Validation (RFEC robustness. By systematically select nificant features, the proposed approach reduces dimensionality, əst s mitigates overfitting, and improve efficiency. The Random Forest classifier trained on the optimized omputational feature set demonstrated exceptional rformance, achieving an accuracy of 99.2% in classifying multiple fault types, Wear Failure, and Overstrain Failure. This validates the effectiveness of including No Failure, Powy nble learning in addressing common challenges such as high variance and poor integrating feature selection with en generalization. The framew tness and reliability make it well-suited for real-time fault diagnosis applications s's robi in smart man ultimately contributing to improved operational safety and reduced maintenance on extending this approach to other industrial domains and exploring adaptive methods to costs. Future nav fo handle ev pattern

References

- [1] Sepulve N. E., & Sinha, J. (2020). Parameter optimisation in the vibration-based machine learning model for accurate and reliable faults diagnosis motating machines. Machines, 8(4), 66.
- [2] Z., Wr g, Y., & Wang, K. S. (2017). Intelligent predictive maintenance for fault diagnosis and prognosis in machine centers: Industry 4.0 scenario Advances in Manufacturing, 5(4), 377-387.
- Xiao, Z., Cheng, Z., & Li, Y. (2021). A review of fault diagnosis methods based on machine learning patterns. 2021 Global Reliability and Reposition and Health Management (PHM-Nanjing), 1-4.
- [4] Cen, J., Yang, Z., Liu, X., Xiong, J., & Chen, H. (2022). A review of data-driven machinery fault diagnosis using machine learning algorithms. Journal of Vibration Engineering & Technologies, 10(7), 2481-2507.
- [5] Fernandes, M., Corchado, J. M., & Marreiros, G. (2022). Machine learning techniques applied to mechanical fault diagnosis and fault prognosis in the context of real industrial manufacturing use-cases: a systematic literature review. Applied Intelligence, 52(12), 14246-14280.
- [6] Luo, J., Liu, Y., Zhang, S., & Liang, J. (2021). Extreme random forest method for machine fault classification. Measurement Science and Technology, 32(11), 114006.
- [7] Patel, R. K., & Giri, V. K. (2016). Feature selection and classification of mechanical fault of an induction motor using random forest classifier. Perspectives in Science, 8, 334-337.

- [8] Susto, G. A., Schirru, A., Pampuri, S., McLoone, S., & Beghi, A. (2014). Machine learning for predictive maintenance: A multiple classifier approach. IEEE transactions on industrial informatics, 11(3), 812-820.
- [9] Xu, G., Liu, M., Wang, J., Ma, Y., Wang, J., Li, F., & Shen, W. (2019, August). Data-driven fault diagnostics and prognostics for predictive maintenance: A brief overview. In 2019 IEEE 15th international conference on automation science and engineering (CASE) (pp. 103-108). IEEE.
- [10] Saucedo-Dorantes, J. J., Delgado-Prieto, M., Osornio-Rios, R. A., & de Jesus Romero-Troncoso, R. (2016). Multifault diagnosis method applied to an electric machine based on high-dimensional feature reduction. IEEE Transactions on industry applications, 53(3), 3086-3097.
- [11] Yang, B. S., Di, X., & Han, T. (2008). Random forests classifier for machine fault diagnosis. Journal of mechanical science and technology, 2 1716-1725.
- [12] Misra, P., & Yadav, A. S. (2020). Improving the classification accuracy using recursive feature elimination with cross-validation. Int. J. Emerg Technol, 11(3), 659-665.
- [13] Zhao, Y., Liu, Z., Yang, Z., Han, Q., & Ma, H. (2025). Machinery fault diagnosis-oriented regularization for nonlinear system identif Framework and applications. Applied Acoustics, 231, 110537.
- [14] Bode, G., Thul, S., Baranski, M., & Müller, D. (2020). Real-world application of machine-learning-based fault detection tract with experimental data. Energy, 198, 117323.
- [15] Brito, L. C., Susto, G. A., Brito, J. N., & Duarte, M. A. (2022). An explainable artificial intelligence approach for unsuper ed rand diagnosis in rotating machinery. Mechanical Systems and Signal Processing, 163, 108105.
- [16] Chen, X., Wang, M., & Zhang, H. (2024). Machine Learning-based Fault Prediction and Diagnosis of Brushless Motor Engineering Acances, 4(3).
- [17] Tang, S., Yuan, S., & Zhu, Y. (2019). Deep learning-based intelligent fault diagnosis methods toward routine machinery. Ieee Access, 8, 9335-9346.
- [18] Gonzalez-Jimenez, D., del-Olmo, J., Poza, J., Garramiola, F., & Sarasola, I. (2021). Machine la ming-based hult determ and diagnosis of faulty power connections of induction machines. Energies, 14(16), 4886.
- [19] Tran, M. Q., Elsisi, M., Mahmoud, K., Liu, M. K., Lehtonen, M., & Darwish, M. M. (2021). Explored setup for online fault diagnosis of induction machines via promising IoT and machine learning: Towards industry 4.0 empowerment. IEEE ccess, 9, 115429-115441.
- [20] Shubita, R. R., Alsadeh, A. S., & Khater, I. M. (2023). Fault detection in rotating machinery based on source in a using edge machine learning. IEEE Access, 11, 6665-6672.
- [21] Chen, S., Meng, Y., Tang, H., Tian, Y., He, N., & Shao, C. (2020). Robust deep learning based segnosis of mixed faults in rotating machinery. IEEE/ASME Transactions on Mechatronics, 25(5), 2167-2176.
- [22] Shao, S. Y., Sun, W. J., Yan, R. Q., Wang, P., & Gao, R. X. (2017). A deep la nine pproch for fault diagnosis of induction motors in manufacturing. Chinese Journal of Mechanical Engineering, 30, 1347-1356.
- [23] Sohaib, M., Kim, C. H., & Kim, J. M. (2017). A hybrid feature model and the p-learning base opearing fault diagnosis. Sensors, 17(12), 2876.
 [24] Kafeel, A., Aziz, S., Awais, M., Khan, M. A., Afaq, K., Idris, S. M. (2021). An expert system for rotating machine fault detection using vibration signal analysis. Sensors, 21(22), 75 .
- [25] Hung, C. W., Lee, C. H., Kuo, C. C., & Zeng, S. X. (200). SoC-base early fail early fail detection system using deep learning for tool wear. Ieee Access, 10, 70491-70501.
- [26] Orrù, P. F., Zoccheddu, A., Sassu, L., Mattia, C., Cozza, R., Ana, S. (2020). Machine learning approach using MLP and SVM algorithms for the fault prediction of a centrifugal pump in the oil and gas to stry. Sustainability, 12(11), 4776.

