

# Machine Learning Models for Fault Classification in Industrial Motors: A Comparative Analysis

Kommanaboina Mahender

*Department of EEE*

*BITS Pilani*

Hyderabad, India

p20230037@hyderabad.bits-pilani.ac.in

Yash Manohar Agarwal

*Department of EEE*

*BITS Pilani*

Hyderabad, India

f20213120@hyderabad.bits-pilani.ac.in

Siddharth Singh

*Department of EEE*

*BITS Pilani*

Hyderabad, India

f20202029@hyderabad.bits-pilani.ac.in

Nikhil Damwani

*Department of EEE*

*BITS Pilani*

Hyderabad, India

f20201769@hyderabad.bits-pilani.ac.in

Sudha Radhika

*Department of EEE*

*BITS Pilani*

Hyderabad, India

sradhika@hyderabad.bits-pilani.ac.in

**Abstract**—Undetected failures can often lead to complete deterioration and breakdown of a motor in the long term. This study explores the complex process of identifying and predicting faults in three-phase induction machines, which is an essential part of many industrial environments. We induce varying percentages of stator and rotor faults in an industrial motor and record the three-phase current measurements at fixed combinations. Upon pre-processing the data obtained, we apply six machine learning (ML) algorithms to predict the fault type based on the current measurements. We find that the random forests model outperforms the other models in terms of error metrics- accuracy, precision, and F1 score. Moreover, we also determine the predictability of each of the fault types and hypothesize that certain fault types are more predictable than others due to lesser variability in their raw data. We conclude by identifying a need to set up a multi-stage ML-model pipeline that would first filter out no-fault scenarios and effectively predict the fault types more comprehensively.

**Index Terms**—Fault detection, machine learning, industrial motor, rotor fault, stator fault.

## I. INTRODUCTION

Numerous industries, such as the energy generation and manufacturing industries, depend for their operational efficiency upon three-phase induction motors. Renowned for their reliability and efficiency, these motors, however, are not immune to malfunctions that can significantly impact their performance and life [1]. Faults such as rotor bar breakages, stator winding failures, and bearing deterioration are common and can lead to costly downtime and hazardous operational conditions. The early detection and accurate diagnosis of such faults remains a considerable challenge within the field. Traditional diagnostic methods, while useful, often lack the sensitivity and specificity required to identify subtle or emerging faults before they evolve into serious failures [2]. This limitation underscores the pressing need for sophisticated diagnostic techniques capable of capturing the complex dynamics of motor faults [3]. In response to

these challenges, recent advancements in machine learning (ML) have emerged as powerful tools for enhancing fault detection in induction motors [4]. Broadly, there are two ways to handle such faults- prediction (avoidance) based [5], [6] and detection (corrective measures) based [7]–[9]. These have traditionally been analyzed using signal processing techniques such as Fast Fourier Transform (FFT), variations of it, wavelet and spectral analysis. While such techniques are transparent in analysis and reliable, they provide low accuracy and perform very poorly against unprecedented situations. In this paper, the aim is to gather fault analysis data from an induction motor using appropriate sensor mounting techniques and perform supervised machine learning on the data after intensive data pre-processing to predict the type of fault. We then compare and study the results obtained from the models. Moreover, the integration of these advanced analytical techniques into fault diagnosis practices represents a significant step forward in the pursuit of operational excellence and safety in industrial settings. It not only enhances the reliability of fault detection systems but also contributes to the broader field of electrical engineering by advancing our theoretical and practical understanding of induction motor dynamics [10].

In this study, 1 healthy class and 15 fault classes—which include different combinations of rotor and stator faults—are thoroughly analyzed. These classes offer a thorough understanding of the machine’s fault dynamics as they are constructed to reflect a broad range of possible real-world scenarios. The other sections in this paper are organized in the following order: Section II describes the background on industrial motors and the machine learning algorithms used, and Section III discusses the data preparation techniques used. Section IV contains the results obtained. Finally, Section V concludes the study in this paper.

## II. BACKGROUND

Induction motors are the power horses of manufacturing industries. Due to environmental conditions, regular wear and tear, manufacturing and installation defects, and overloading, these motors are subjected to failures. Detection of these faults is an indispensable function for operation safety and fast maintenance. Bearing faults, stator faults, and rotor faults are the frequently occurring types of faults. If the motor continues to run with these faults, a tremendous amount of heat is produced, which results in a decrease in energy efficiency and leads to catastrophic failure. Therefore, detecting these failures earlier reduces maintenance costs and the time needed for maintenance. With the help of mechanical signals like vibration signals, these failures can be detected much before their occurrence [11], [12]. These vibrational signals are usually acquired with the help of accelerometers mounted on the induction motors, and the acquired signals are analyzed in various domains like time response analysis, frequency response analysis, and time-frequency response analysis in which feature extraction and feature classification are the key factors [13].

Supervised learning is a collection of machine learning techniques that involve training models on training sets with definite values of a clearly defined target variable. The goal of supervised learning is to obtain the relationship between different variables that leads to each value of the target variable, essentially creating a function that accepts the values of the variables as input and returns the target variable's value. While training, the algorithm tries to minimize the loss function that quantifies the deviation of predicted values from the actual values. Supervised learning is helpful in situations where the goal is classification, like in this paper, where we have 15 types of fault classes and a healthy class. We have used the following 6 machine learning models to learn from the training data provided and then classify each data point of the test set as a fault type.

### A. Logistic Regression

Logistic regression is a supervised machine learning algorithm that is useful for binary and multiple classification. It uses the sigmoid function that is defined as:

$$f(x) = \frac{1}{1 + e^{-x}}$$

The logistic function's range is from 0 to 1, and the sets of values it returns are then classified according to the type of classification we require. In the case of multiple input variables, a variable  $y$  is defined as the linear combination of these variables and used instead of  $x$ . The coefficients used in this linear combination are known as the regression coefficients.

$$y = \beta_0 X_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

We group and map the outputs of the logistic function to different values of the target variables and define regions for our 16 target classes.

### B. Support Vector Machine

Support vector machine (SVM) is a powerful algorithm for classification. It tries slicing a multidimensional space using a hyperplane to obtain clearly defined regions where the probability of finding a corresponding value of the target variable is high. Hyperplanes are mathematically defined curves that split an  $n$ -dimensional space (' $n$ ' being the number of independent variables) to separate different classes of data points with the minimum number of outliers possible. If there are too many outliers or the data cannot be split, it uses a function to introduce a new variable where spatial groups can be formed. Such a non-linear function that redefines data allowing us to efficiently split it using hyperplanes, is known as a kernel. Different types of kernels can be used, such as 'linear', 'sigmoid', 'poly', etc. In our paper, we use the radial basis function (rbf) kernel, defined as:

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right), \quad \sigma > 0$$

### C. K-Nearest Neighbors

The K-nearest neighbors algorithm classifies data on the intuitive principle that data points of a class are situated closer to each other than they are to other classes. The algorithm accounts for classes of a constant number of points closest to a given point and predicts the target variable's class accordingly. The proximity can be measured by using a distance function, for which multiple functions are available. The most famous are the Euclidean, Manhattan, and Minkowski distances. In this paper, we use the standard Euclidean distance for our predictions:

$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

Once these distances have been used to find the K-nearest neighbors, either the majority is used as the predicted value or the distances of each point can be taken into account, giving points that are closer to the focus more importance than those that are further away.

### D. Gaussian Naive Bayes

The Gaussian Naive Bayes algorithm assumes that each variable contributes independently to the target variable's value. The algorithm also assumes that the input variables are not somehow dependent. It calculates the probability of each predicted value of the target variable using each independent variable and multiplies all obtained probabilities for a given classification to obtain a final class-wise probability distribution of the target variable. Even when the features are not fully independent, as in our case, where statistical features usually have some degree of dependence on each other, Bayesian learning has been shown to work well. Additionally, it works with any number of features. The class with the highest probability is chosen as the predicted output. Each independent probability is calculated by assuming a

Gaussian or Normal probability distribution, which is given by:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

#### E. Decision Tree

A decision tree is a classification algorithm that can be understood as a series of conditions that sequentially eliminate classes for a data point by using pre-learned conditions. They are structured like the tree data structure, with the root acting as incoming data and leaves as classes. Decision trees traverse the data set and figure out conditions that strictly do or do not occur in each of the classes. After defining these conditions, they essentially work like nested if-else statements and classify incoming data based on the combinations of conditions they satisfy. They use the ID3 algorithm to maximize the "information gain" as follows:

$$\text{Gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

We use Gini impurity for the fault classification task as follows:

$$\text{Gini}(T) = 1 - \sum_{i=1}^c p_i^2$$

#### F. Random Forest

A random forest uses numerous decision trees that independently classify incoming data, and the most frequent classification is selected as the prediction. Random forests are a very effective measure against overfitting, which occurs when a decision tree makes too many relations and becomes very accurate for the training set while sacrificing accuracy for evaluation sets and other inputs. The inclusion of multiple random trees overcomes this overfitting.

### III. DATA PREPARATION

Raw three-phase current measurements ( $I_1$ ,  $I_2$ ,  $I_3$ ) were obtained from a 3-phase induction machine. The setup was designed to emulate operational conditions under various fault scenarios [11], [13]. This initial step is crucial for ensuring the relevance of data to real-world applications, where accurate and comprehensive data collection underpins the success of subsequent analysis. The data was systematically categorized into 16 classes, representing different combinations of rotor and stator faults (including no-fault conditions), with each class comprising 500 data files. This classification allows for a nuanced analysis of fault dynamics, facilitating the development of targeted diagnostic and predictive models.

In the experimental setup, to avoid information loss and ensure high fidelity in fault diagnosis, we captured the data using a NI myDAQ USB6210 multifunction I/O device, boasting a 16-bit analog input with a sampling rate of 250 kS/s, 4 digital inputs, and 4 digital outputs. The acquisition was conducted at a 10 kHz sampling rate, which is crucial for capturing the higher frequency components present in stator current under the influence of multi-component failure. This

TABLE I: Nomenclature for fault combinations

Rotor Fault Percentage	Stator Fault Percentage			
	25%	50%	75%	0%
25%	Fault-A	Fault-B	Fault-C	Fault-D
50%	Fault-E	Fault-F	Fault-G	Fault-H
75%	Fault-I	Fault-J	Fault-K	Fault-L
0%	Fault-M	Fault-N	Fault-O	Healthy

elevated sampling rate facilitates the collection of nuanced data that is essential for distinguishing between healthy and faulty induction motors during subsequent analyses.

We perform classification on one healthy condition and 15 types of rotor and stator faults. Four types of rotor faults are induced at levels of 0%, 25%, 50%, and 75%. For each of these rotor fault levels, stator faults are also induced at 0%, 25%, 50%, and 75%. For each stator fault condition, we have 500 test files. Each test file contains rows with features such as  $I_1$ ,  $I_2$ ,  $I_3$ ,  $I_1(\text{PSD})$ ,  $I_2(\text{PSD})$ , and  $I_3(\text{PSD})$ . The data is stored in the format shown in Table I.

For the pre-processed dataset of 16 fault types, each with 500 test files (a total of 8000 files), we extract and transform the original features into statistical measures: mean, median, mode, variance, standard deviation, and kurtosis. As a result, the new dataset consists of 36 features and 8000 rows, each corresponding to one of the test files. Therefore, all fault classes are represented equally, with 500 instances each, effectively avoiding class imbalance.

### IV. RESULTS

After obtaining a dataset with 36 features and a target variable representing the fault type (numbered from 1 to 16 for the 15 faults and one healthy state), and after thoroughly analyzing the data obtained by refining the sensor mounting techniques for data extraction, we perform supervised machine learning (ML) to predict the fault type given new rows of data. We deploy the algorithms described in Section II for a comparative analysis. The results reveal that the best-performing model is random forests. All models, in decreasing order of F1 score, are mentioned in Fig. 1.

Additionally, to assess the predictability of each fault type, we analyze the results presented in Table II by calculating the average error metrics for each of the fault types across all machine learning algorithms. The fault types Fault-B, Fault-F, Fault-H, Fault-I, Fault-J, and Fault-N, demonstrate excellent prediction performance (Accuracy, Precision, and F1 Score near 1.00). The second-lowest prediction accuracy is observed for the no-fault condition (Healthy State). Although there is no class imbalance among the 16 target states, an underlying imbalance exists between the no-fault condition (Healthy State) and the faults (all other 15 combinations). This suggests an internal class imbalance issue, which we propose addressing through an improved ML pipeline in

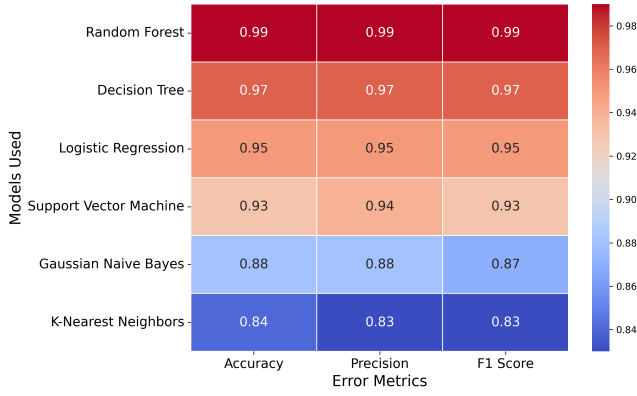


Fig. 1: Comparative analysis of ML algorithms used

TABLE II: A comparison between the fault types - how well each of them can be predicted using the ML models

Model Used	Error Metrics		
	Accuracy	Precision	F1 Score
Fault-A	0.81	0.95	0.87
Fault-B	0.99	0.99	0.99
Fault-C	0.92	0.96	0.94
Fault-D	0.86	0.87	0.86
Fault-E	0.98	0.99	0.99
Fault-F	0.99	0.99	0.99
Fault-G	0.90	0.90	0.90
Fault-H	0.99	0.99	0.99
Fault-I	0.99	0.99	0.99
Fault-J	0.99	0.99	0.99
Fault-K	0.78	0.75	0.76
Fault-L	0.87	0.93	0.90
Fault-M	0.82	0.93	0.87
Fault-N	0.99	0.99	0.99
Fault-O	0.89	0.92	0.90
Healthy	0.84	0.77	0.80

Section V. The poorest prediction accuracy is observed for Fault-K, which exhibits the highest severity due to a 75% rotor fault and a 75% stator fault, significantly increasing the likelihood of motor failure. Despite its critical importance, current ML models struggle to predict this fault accurately. Our analysis also indicates that inducing a 50% fault in either the stator or rotor results in high error metric values for the ML models. The 0% (no fault) condition shows a wide variation in results, suggesting that ML models are unreliable when healthy conditions are used.

## V. CONCLUSION

In this study, we explored the application of machine learning models—Logistic Regression, Support Vector Machines, K-Nearest Neighbors, Gaussian Naïve Bayes, Decision Trees, and Random Forests—for fault classification in

industrial motors. We compared the results and evaluated the predictability of each fault type. Our findings suggest that fault types with higher severity levels introduce greater variability into the data, which, in turn, affects model performance and results in lower prediction accuracy. Despite these challenges, the research demonstrates potential for refinement. We identified the need for a robust machine learning pipeline, incorporating cascaded models, to improve prediction accuracy. Specifically, we propose adding a preliminary stage to classify whether a fault is present, which would help distinguish between healthy and faulty states. This step is critical, as the difficulty in predicting healthy states likely affects the performance of models for other fault types. Future work will focus on integrating this initial fault identification stage, with the goal of enhancing the overall accuracy and reliability of the fault classification system.

## REFERENCES

- [1] M. Pineda-Sanchez *et al.*, “Diagnosis of induction motor faults via sound and vibration signals,” *IEEE Transactions on Industrial Electronics*, vol. 57, no. 3, pp. 871–882, 2010.
- [2] J. Faiz and B. Ebrahimi, “Mixed-fault diagnosis in induction motors through vibration and current signature analysis,” *IEEE Transactions on Industrial Electronics*, vol. 56, no. 11, pp. 4715–4722, 2009.
- [3] W. Thomson and M. Fenger, “Current signature analysis to detect induction motor faults,” *IEEE Industry Applications Magazine*, vol. 7, no. 4, pp. 26–34, 2001.
- [4] A. Bellini *et al.*, “Advances in diagnostic techniques for induction machines,” *IEEE Transactions on Industrial Electronics*, vol. 55, no. 12, pp. 4109–4126, 2008.
- [5] E. Leandro, L. E. de Lacerda de Oliveira, J. G. B. da Silva, G. Lambert-Torres, and L. E. B. da Silva, *Predictive Maintenance by Electrical Signature Analysis to Induction Motors*. InTech, Nov 2012.
- [6] V. Kavana and M. Neethi, “Fault analysis and predictive maintenance of induction motor using machine learning,” in *2018 International Conference on Electrical, Electronics, Communication, Computer, and Optimization Techniques (ICECCOT)*, Mysuru, India, 2018, pp. 963–966.
- [7] P. Kumar and A. Hati, “Review on machine learning algorithm based fault detection in induction motors,” *Archives of Computational Methods in Engineering*, vol. 28, pp. 1929–1940, 2021.
- [8] S. Mari, G. Bucci, F. Ciancetta, E. Fiorucci, and A. Fioravanti, “Machine learning for anomaly detection in induction motors,” in *2023 IEEE International Conference on Metrology for eXtended Reality, Artificial Intelligence and Neural Engineering (MetroXRINE)*, Milano, Italy, 2023, pp. 1057–1062.
- [9] T. B. Sawai, S. Mudke, P. V. Kapoor, and U. B. Mujumdar, “Bearing fault detection in induction motor using ensemble learning,” in *2023 5th International Conference on Energy, Power and Environment: Towards Flexible Green Energy Technologies (ICEPE)*, Shillong, India, 2023, pp. 1–6.
- [10] J. Wadibhasme, S. Zaday, and R. Somalwar, “Review of various methods in improvement in speed, power & efficiency of induction motor,” in *2017 International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS)*, Chennai, India, 2017, pp. 3293–3296.
- [11] Y. Liu and A. M. Bazzi, “A review and comparison of fault detection and diagnosis methods for squirrel-cage induction motors: State of the art,” *ISA Transactions*, vol. 70, pp. 400–409, 2017.
- [12] K. Yatsugi, S. E. Pandarakone, Y. Mizuno, and H. Nakamura, “Common diagnosis approach to three-class induction motor faults using stator current feature and support vector machine,” *IEEE Access*, vol. 11, pp. 24 945–24 952, 2023.
- [13] M.-K. Liu and P.-Y. Weng, “Fault diagnosis of ball bearing elements: A generic procedure based on time-frequency analysis,” *Measurement Science Review*, vol. 19, no. 4, pp. 185–194, 2019.