

PUMP HEALTH PREDICTION THROUGH VIBRATION SIGNAL ANALYSIS.

Presented by:

ASHMIT MUKHERJEE | Presidency University,

(BSc. MATHEMATICS)

DEV PANDAY | Presidency University,

(BSc. MATHEMATICS)

RITANKAR KUNDU | Presidency University,

(BSc. MATHEMATICS)

NILANJAN MANNA | Central University of South Bihar

(MSc. DATA SCIENCE AND APPLIED STATISTICS)

PAYAL PRARAMVIKA SAHU | Central University of South Bihar

(MSc. DATA SCIENCE AND APPLIED STATISTICS)

Project Guide/ Mentor:

ANIK SARDAR.

Period of Internship: 19th MAY 2025 – 15th JULY 2025

Report submitted to: IDEAS – Institute of Data
Engineering, Analytics and Science Foundation, ISI
Kolkata.

1. Abstract

This project focuses on the detection of anomalies and faults in industrial pump systems by analyzing time-series sensor data collected at millisecond intervals. Techniques such as z-score analysis, boxplot detection, temporal clustering, recurrence detection, and FFT-based frequency analysis were applied. A robust scoring framework was developed by combining statistical, contextual, and frequency-based insights to compute a final health score. Various methods for missing data imputation, flagging, and feature extraction were implemented. A Random Forest-based machine learning model was also trained and integrated to enhance predictive capability using the engineered features. The final model enables segment-wise monitoring and categorization of equipment condition as Healthy, Monitor, Warning, or Critical. Python and pandas-based pipelines were used for the end-to-end processing. The insights gained have strong implications for predictive maintenance systems in industrial settings.

2. Introduction

In modern industrial environments, the operational integrity of rotating machinery—particularly pump systems—is critical to maintaining overall process efficiency and minimizing unplanned downtime. These systems are often subject to complex dynamic loads and harsh environmental conditions, making them vulnerable to wear, misalignment, imbalance, and other mechanical faults. Early detection of such anomalies is essential to prevent costly breakdowns, reduce maintenance costs, and ensure safety. With the advent of Industrial Internet of Things (IIoT) technologies, high-frequency time-series sensor data has become increasingly accessible, enabling advanced data-driven diagnostics and predictive maintenance frameworks.

Traditional condition monitoring techniques often rely on predefined thresholds or domain-specific rules, which may fail to generalize across varying operating conditions or miss subtle precursors to faults. To overcome these limitations, this study leverages a comprehensive data-centric approach using millisecond-level vibration data collected from tri-axial accelerometers installed on industrial pumps. The high temporal resolution of this data allows for the detection of fine-grained patterns and transient anomalies that would otherwise remain unnoticed.

This project focuses on the early detection of faults and anomalies in industrial pump systems by analyzing high-frequency sensor data collected over time. The primary aim is to develop a data-driven, automated health monitoring solution capable of identifying early warning signs of mechanical issues.

This work involves the use of time-series analytics, statistical modelling, frequency domain analysis, and machine learning. A multi-stage pipeline was designed to transform raw vibration signals into actionable insights. The final outputs enable real-time equipment condition scoring and categorization, supporting maintenance decision-making and improving system reliability.

As a precursor to the project, the first two weeks of the internship were dedicated to an intensive training program aimed at building foundational skills necessary for data-driven problem solving. The training covered essential topics including Python programming, exploratory data analysis, Power BI for dashboarding, no-code tool workflows, as well as conceptual foundations of data analysis, generative AI, and large language models (LLMs). This foundational knowledge provided the groundwork for implementing the advanced techniques and analytical methodologies employed in the subsequent stages of the project.

3. Project Objective

The objectives of the project are as follows:

- To preprocess high-resolution millisecond-level sensor data and handle missing values through detection, imputation, and integration techniques.
- To identify anomalous behaviour in acceleration signals using statistical methods such as z-score and boxplot-based outlier detection.
- To enhance anomaly labelling by incorporating contextual awareness through neighbouring patterns, detection types, and confidence-based scoring strategies.
- To cluster anomalies based on temporal proximity using DBSCAN for detecting localized fault bursts and temporal grouping behaviour.
- To analyze recurring anomaly patterns by segmenting time and detecting repeated fault occurrences within fixed temporal offsets.
- To extract frequency-domain features using Fast Fourier Transform (FFT) and assess signal health via power spectral analysis.
- To build a composite scoring framework and integrate a trained Random Forest model to predict and label pump health status into defined risk categories (Healthy, Monitor, Warning, Critical).

4. Methodology

This project was centred on analyzing high-resolution time-series sensor data from industrial pump systems, aiming to identify early indications of faults and anomalous behaviour. The dataset was provided by the mentor team and contained accelerometer readings across three axes (x, y, z) with timestamps at millisecond granularity.

A) Data collection & Initial Processing

- **Source:** The raw dataset was provided in JSON format, collected from vibration sensors attached to industrial pump equipment.
- **Conversion:** Python scripts were developed to convert Unix-based millisecond timestamps to human-readable datetime formats and to store the data in structured Excel format for downstream processing.
- **Units Conversion:** Raw acceleration values in "g" were converted to m/s^2 using a scaling factor of 9.80665 to ensure physical relevance and compatibility with analytical methods.

B) Data Cleaning & Preprocessing

- **Missing Value Flagging:** Data points with zero acceleration values across any axis were flagged as potential missing or faulty entries.
- **Rolling Imputation:** A rolling-window approach (± 3 points) was used to impute missing values using the local mean of non-zero neighbours, while preserving a separate imputation trace.
- **Exploratory Data Analysis:** Initial Manual and Observatory inspection were done to identify obvious spikes and zero-runs, also Amplitude-Line charts, rolling windows used to analyze possible Anomalies.
- **Outlier Flagging:**
 - i) **Time-series analysis and window-based analysis:** Rolling RMS + Kurtosis [Fixed threshold + Percentile (IQR) threshold] were used to detect anomalies and spikes distributed across datasets. (Window size used:

a) RMS flagging (condition):

Fixed: $flag \text{ if } RMS_t > \mu_{RMS} + 2\sigma_{RMS}$

Percentile: *flag if $RMS_t > P_{95} (RMS)$*

Where: $RMS_t \rightarrow RMS$ of the t^{th} window

$\mu_{RMS} \rightarrow$ Mean of all RMS findings

$\sigma_{RMS} \rightarrow$ Standard deviation of all RMS findings

$P_{95} \rightarrow$ 95th percentile of all RMS findings

b) KURTOSIS flagging (condition):

Fixed: *flag if $KURTOSIS_t > 3.5$ (standard)*

Percentile: *flag if $KURTOSIS_t > P_{95} (KURTOSIS)$*

c) Combined flagging (condition):

$fixed_{RMS} OR percentile_{RMS} == TRUE \Rightarrow RMS_{combined} flag == TRUE.$

$fixed_{KURTOSIS} OR percentile_{KURTOSIS} == TRUE \Rightarrow RMS_{combined} flag == TRUE.$

d) All the computation is done per axis (x/y/z) wise.

e) Formulas used:

$$\text{Rolling RMS} = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}$$

$$\text{Rolling KURTOSIS} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \mu}{s} \right)^4$$

Where: *Rolling window (n) = 50 +*

1 [i.e. 25 before and after the current data point]

Vibrational acc. = x_i

mean of the window = μ

standard deviation of that window = s

ii) **Boxplot-Based (whiskers) Detection:** IQR-based bounds were used to flag distributional anomalies on each axis.

a) *lower bound = $Q_1 - 1.5 \times IQR$, upper bound = $Q_3 + 1.5 \times IQR$*

b) Flag (condition) == 1:

If value < lower bound OR value > upper bound.

c) Flag (condition) == 0:

If value > lower bound OR value < upper bound.

d) All the computation is done per axis (x/y/z) wise.

iii) **Z-Score Based Detection:** Used the method of z-score based spike detection to filter highly efficiently distributed outliers over the data.

a) Working formula for z-score:

$$z = \frac{x_i - \mu}{\sigma}$$

x_i is the vibrational acceleration

μ is the Mean of that axis (calculated for that entire dataset)

σ is the standard deviation of that axis.

b) Spike detection (condition):

Fixed thresholding: if $|z| > 3.0$; then flag else NOT.

Adaptive thresholding: if value > P_{99} OR value < P_1 then flag else NOT.

c) All the computation is done per axis (x/y/z) wise.

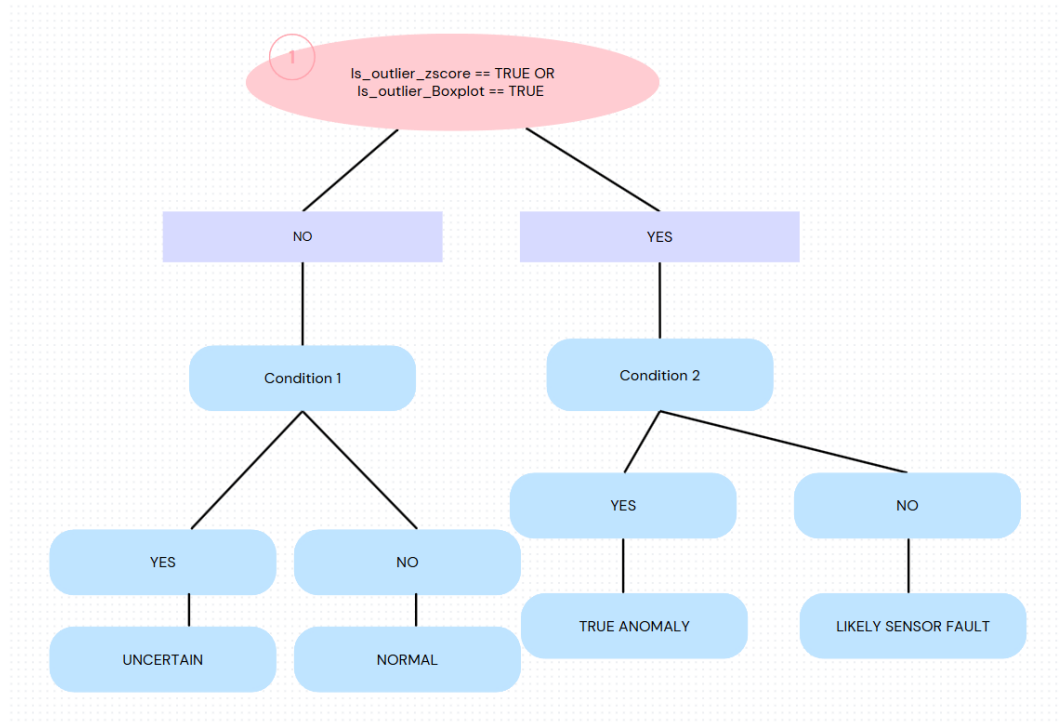
iv) **Combined Outlier Flags:** Axis-wise outlier flags were combined to generate summary columns and isolate abnormal time points.

C) Outlier classification using Contextual, Temporal and Off-set Recursive (Periodicity check) Analysis:

i) Outlier classification-1: Outlier classification-1 is based on the principle of “Contextually-aware anomaly classification using neighbouring axes. It Takes into account previously detected z-score and box-plot outliers and analyses their contexts (neighbours and overlaps). Two types of logic models (loosened and enhanced) assign qualitative labels on the outliers based on current and neighbouring spike fags. Each logic converts its label into numerical score (0-3 or 0-4) scale. A final contextual score is computed as the average of normalized loosened and enhanced scores. Based on the final contextual score, a final label (Normal / Mild Anomaly / Probable Fault / Confirmed Fault) is assigned to each Outlier Data point. While

taking Neighbouring rows into consideration we used ± 1 leeway measure, i.e. (only one row above and below the current considered row).

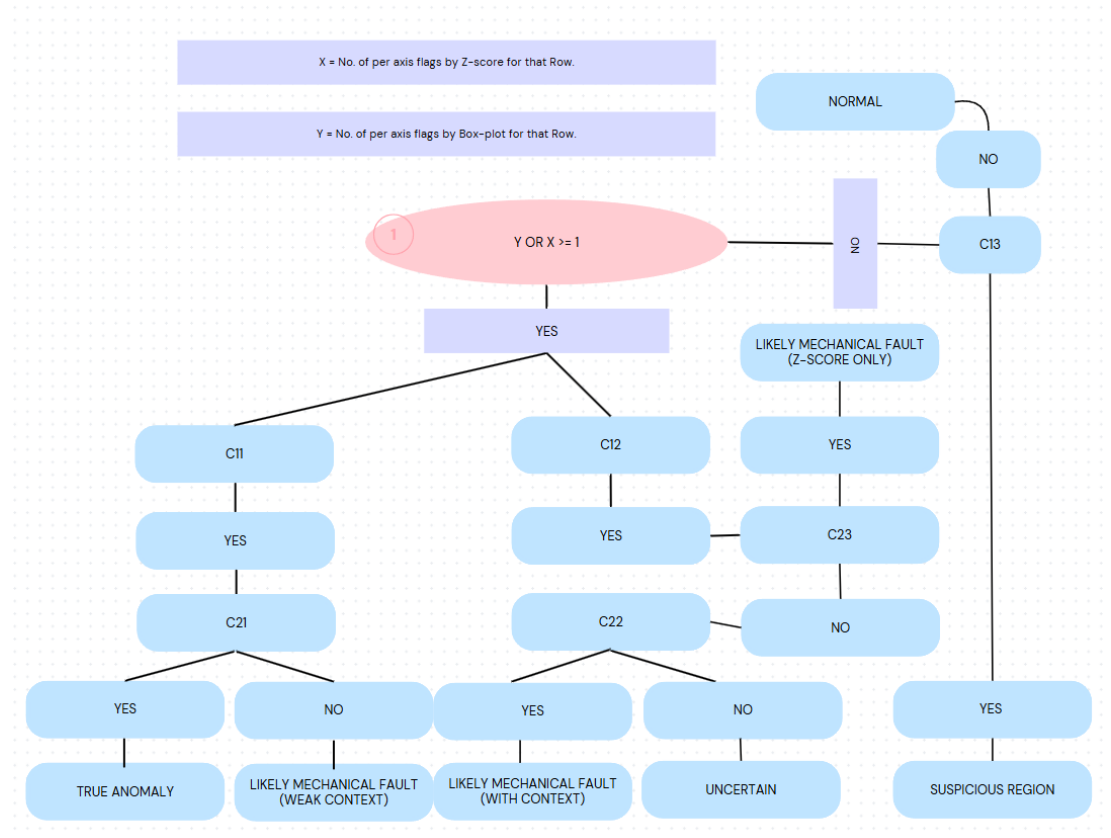
(Loosened-label logic tree):



CONDITION 1: total no. of per axis flags by any of the techniques (box-plot OR z-score) in the neighbour rows is ≥ 2 .

CONDITION 2: Total no. of per axis flags on that parent row is ≥ 2 OR total no. of per axis flags by any of the techniques in the neighbouring rows is ≥ 1 .

(Enhanced-label logic tree)



C11: Have at least one per axis z-score flag as well as box-plot flag.

C12: Have at least one per axis of only one of either z-score or box-plot.

C13: Total per axis flags around neighbour rows is at least 2.

C21: Total per axis flags around neighbour rows is at least 1.

C22 = C13: Total per axis flags around neighbour rows is at least 2.

C23: Only z-score flag but NOT Box-plot flag.

(Score Calculation Table):

Contextual_score_loosened	labels:	Contextual_score_enhanced	labels:
0	NORMAL	0	NORMAL
1	UNCERTAIN	1	SUSPICIOUS REGION
2	LIKELY SENSOR FAULT	2	UNCERTAIN
3	TRUE ANOMALY	2.5	LIKELY MECHANICAL FAULT (WEAK CONTEXT)
		3	LIKELY MECHANICAL FAULT (WITH CONTEXT)
		3.5	LIKELY MECHANICAL FAULT (Z-SCORE ONLY)
		4	TRUE ANOMALY
		FINAL CONTEXTUAL SCORE: $0.5 * (\text{CONTEXTUAL_SCORE_LOOSENED} / 3) + 0.5 * (\text{CONTEXTUAL_SCORE_ENHANCED} / 4)$	
Final contextual score	labels:		
score < 0.25	NORMAL		
score < 0.5	MILD ANOMALY		
score < 0.75	PROBABLE FAULT		
else	CONFIRMED ANOMALLY		

Brief about the whole procedure: Outlier classification-1 uses the previously flagged outliers by box-plot and z-score to evaluate and classify them into various labels for the ML algorithm to understand. They look for surrounding Neighbouring rows which may be flagged besides the main row on which we are evaluating and based on the intensity of surrounding anomaly, we give a normalized score and a label for the ML model to understand what kind of anomaly does the current flagged row represent.

- ii) **Outlier classification-2:** Outlier classification-2 used the concept of “Temporal clustering of outliers using DBSCAN” to group outliers. The goal is to find whether we can identify outlier clusters forming in the datasets and give a boost to the score of these high-confidence outliers so as to differentiate Normal outliers with these outliers.

Parameters:

- ⇒ Maximum time gap (in seconds) between consecutive points to be considered in the same cluster. (Eps_SECONDS = 5)
- ⇒ Minimum number of points within (Eps_SECONDS) to form a valid cluster. (Min_Samples = 3)

Methodology:

- a) **EXTRACT OUTLIERS:** Use existing flags (is_outlier, is_outlier_boxplot) to filter only outlier rows
- b) **CONVERT TIME TO NUMERIC SPACE:** Convert ‘datetime’ to seconds since the start of the data to perform clustering in 1D time.
- c) **APPLY DBSCAN:** Run DBSCAN on this 1D time space to assign each outlier a temporal cluster id. Where cluster id == -1 if it is an Isolated point. Else cluster id >= 0 if grouped into a cluster usually representing the group no. starting with 0 as the first group.
- d) This process is done only for the outlier rows which are flagged previously using the 2nd and 3rd pass outlier detection filtering.
- e) **ASSIGN LABELS:** A new column (Temporal_outlier_type) is assigned which outputs “grouped” if part of a DBSCAN cluster and “isolated” if detected alone.

- iii) **Outlier classification-3:** Outlier classification-3 is based on the concept of periodicity check and is done seamlessly with the help of “OFF-set Recurrence analysis”. The goal is to find temporal patterns among outliers – “do a certain outlier events tend to reoccur around the same second mark, in repeated time segments?” this question is addressed in detail using this method.

Variables:

Segment_duration \Rightarrow Splits time into specified time interval segments (15 seconds, in our case).

Offset_tolerance \Rightarrow Allows grouping offsets within the fixed segments for further seamless periodicity tracking (0.5 seconds, in our case).

Min_Recursions \Rightarrow Minimum segments in which a pattern (outlier in the same offset) must repeat to be counted as recurring (at least 3).

Methodology:

- a) Convert Timestamps to seconds since the beginning, this gives each row a numerical time offset for easy Math.
- b) The timeline is split into equal length segments (e.g. 0s – 15s, 15s – 30s...and so on.)
- c) Compute offset ranges inside each segment. This helps us to answer the question: “At what time/second inside its segment does this pattern reoccur”.
- d) Round offsets into tolerance buckets, which helps group offsets that are within close proximity. For e.g. pattern after 5.21s, 5.33s, 5.18s... in a unique segment all gets rounded up to the offset bucket of 5.0s. This ensure we are tolerant to small timing variations due to sensor noise or sampling jitter.
- e) Build a recurrence Map e.g. 5.0s \rightarrow {0, 1, 2}; 12.5s \rightarrow {3, 5}. I.e. basically, note down recurring events in which offset and on which segments it is occurring.
- f) Filter out offsets that appeared in at least (Min_recurrence) segments. For e.g. if 5.0s offset is present in segment 0,1 and 2. Then it is marked as ‘recurring’.
- g) Compute the recurrence score for each outlier rows as flagged earlier. Outlier score is equivalent to no. of segments it repeated or recurred. (for not flagged rows we have the score as 0 for default).

D) Frequency based feature extraction:

Initially we thought of multilayered Frequency based feature extraction and analysis but after further negotiations we arrived with only Fast Fourier Transformation (FFT) as the best bet for our data. Here’s why:

- **Key characteristics of our data:**
 - a) Fixed or near-constant sampling intervals, but not perfectly uniform
 - b) Comes from mechanical systems (pumps) where faults or anomalies may show as periodic vibrations or spike-like disturbances.
 - c) Sampling rate is in the low-Hz to sub-10 Hz range, limiting resolution to low-frequency behavior.
 - d) Goal: Identify dominant frequency components, power distribution, and possible mechanical resonances.
- We didn’t select STFT as it Requires a fixed window size — tradeoff between time vs. frequency resolution. Which is not for our case.
- DFT calculates the same thing as FFT but FFT is much faster and computationally sound rather than just using integrals that would make it computationally hard to calculate.
- Wavelet transform was also not our bet as Wavelets are good for transient-heavy signals like impacts, clicks, or ECGs. But our signal’s anomalies are more periodic or quasi-periodic, not transient shocks.

- Hilbert's transform was also in our pockets but it got ruled out as our sensor data is broadband, covering 0–10 Hz, and not suitable for narrowband analysis.
- In short: We chose FFT because our sensor data is largely stationary, long in duration, and exhibits periodic behavior typical of mechanical systems. Unlike other transforms designed for transients or adaptive modeling, FFT offers an interpretable, efficient, and directly relevant way to extract frequency-domain energy and trends. Our 10-second segmentation approach further helps us capture evolving behavior without needing STFT or Wavelets. This makes FFT not just a popular choice, but the most structurally aligned method for our objective.
- And hence, FFT was applied in our dataset.

Some of the necessary outlines and mentions:

- **SAMPLING RATE:** Our pipeline estimates the sampling rate dynamically for each data set by computing the average time difference between consecutive timestamps in the 'datetime' column. This average time delta is inverted to estimate the sampling frequency, assuming approximately uniform spacing. This makes the FFT analysis adaptive to each file's timing resolution, accommodating millisecond-level variations and inconsistent sampling intervals.
- **ENERGY BINS:** the selected BINS for calculation based on the sampling fit for our data, are 0-1, 1-3, 3-5, 5-10 Hz.
- Hence by FFT the signal is decomposed into a sum of sine and cosine waves at different frequencies, each carrying some power. TOTAL POWER measures the overall energy in the signal, while spectral centroid indicates the "center of gravity" of that energy in frequency space.
- While all signals are divided into fixed 10-second chunks for uniform FFT processing, each dataset's sampling rate is calculated individually. This ensures that the frequency axis is accurately mapped to real-world Hertz values for each dataset, maintaining physical interpretability. Thus, the sampling rate enables precise frequency decomposition even under variable temporal resolutions.
- **Analysis was done per Axis-wise (x/y/z).**

E) Final scoring and Labelling for ML-Training:

At this point of data analysis, we are well equipped with multiple inferences and features for each row in the data set, so as to perform a final Normalized score and a final Label to put on each row which truly will illustrate and exhibit the true nature of the Data point in the dataset.

There are mainly two Relevant Main scores to consider for each row-wise scoring and labelling.

- a) The first one is Time-domain score. Under Time-domain score we have three more scores namely, Time-series score, Contextual score, Temporal score, Recurrence score.
- b) The second is Frequency-domain score. Frequency-domain score is calculated based on the frequency features of each axis.

Let's focus on each:

a) (Time-domain score):

- **time series score** = $\frac{RMS\ score + KURTOSIS\ score}{2}$
Where: RMS score = $\frac{\sum <axis> combined\ RMS\ flag}{3}$ for all $< axis >$ as x/y/z.
KURTOSIS score = $\frac{\sum <axis> combined\ KURTOSIS\ flag}{3}$ for all $< axis >$ as x/y/z.
- **Contextual score** = $0.5(loose\ score) + 0.5(enhanced\ score)$
Where: loose score = $\frac{Contextual\ score\ loosened}{3}$.
enhanced score = $\frac{Contextual\ score\ enhanced}{4}$.

- *Temporal score == 1 iff (temporal outlier type) == (grouped) else Temporal score == 0.*
 - *Recurrence score \Rightarrow Auto scaled such that highest recurring pattern/outlier gets the score == 1 and other scaled accordingly.*
 - *Finally Time – domain score = $0.5 \times \text{time series score} + 0.2 \times \text{contextual score} + 0.2 \times \text{temporal score} + 0.1 \times \text{recurrence score}.$*
- b) (Frequency-domain score):
- *Each $\langle \text{axis} \rangle \text{ score} \rightarrow \text{Avg. 6 FFT features}.$*
 - *Frequency interval score $\rightarrow \text{Avg. } (x, y, z) \langle \text{axis} \rangle \text{ score}.$*
 - *Time based frequency score = mapping of the global 10s fft intervals to individual data points in the main data set*
 - *Frequency domain score = Time based frequency score of that (datapoint or row in main data)*
- c) (final scoring and labelling):
- *final score = $\frac{\text{Time domain score} + \text{Frequency domain score}}{2}$*
 - *score $> P_{95}(\text{final score}) \Rightarrow \text{CRITICAL}$*
 - *score $> P_{75}(\text{final score}) \Rightarrow \text{WARNING}$*
 - *score $> P_{50}(\text{final score}) \Rightarrow \text{MONITOR}$*
 - *Else $\Rightarrow \text{HEALTHY}$*

F) ML training: After feature engineering and label generation through time-series analysis and statistical techniques, the processed dataset was prepared for supervised learning. The goal was to build classification models capable of predicting the health condition of the equipment (Healthy, Monitor, Warning, Critical) based on extracted features.

Dataset preparation:

- The final dataset was aggregated from multiple sensor files covering diverse operational states (Blower, Gearbox, Motor) under different conditions (On and Off).
- Features included statistical flags (RMS, Kurtosis), outlier flags (Z-score, Box-plot), contextual, temporal, and recurrence scores, as well as frequency-based metrics derived from FFT.
- Non-essential columns (e.g., timestamps, intervals) and imputation intermediates were removed.
- All categorical labels (e.g., Final_label, contextual flags) were encoded using LabelEncoder.

Train-test Split:

- The data was divided using an 80:20 ratio for training and testing respectively. Stratification was maintained based on labels.
- Further splits were made internally for hyperparameter tuning during model comparison.

Model Training:

Three models were trained and evaluated for comparative analysis:

- **Decision Tree Classifier:**
A simple, interpretable model used as a baseline to assess the effectiveness of individual features.
- **Random Forest Classifier:**
An ensemble model leveraging multiple decision trees to improve generalization. This model gave the best performance overall.
- **Support Vector Machine (SVM) with RBF Kernel:**
A non-linear classifier suitable for high-dimensional feature spaces, tested to examine boundary margin behaviour.

Feature Importance:

- Feature correlation heatmaps and importance rankings were used to interpret the most influential features in classification.
- Final scores, time-domain metrics, and frequency-based features were among the top contributors.

Model Evaluation:

- Evaluation was conducted using accuracy, precision, recall, and F1-score.
- Confusion matrices were generated for all models to visualize class-wise prediction strengths and errors.
- A performance comparison chart was plotted for the models to help identify the most suitable one.

5. Results.

Since our data has a vast variability feature and is being distributed into different files (91 files with 12000 Rows each) each depicting unique environmental, atmospheric and operational conditions, hence we shall infer to a sample of one of the data sets of a fixed condition and carry out the summarization for that single dataset using the pipeline we discussed earlier. While other datasets are also processed with the same pipeline (can be inferred through the GitHub Repository given afterwards in the report), the Evaluation technique is similar and is left to the interested people for their relevance.

File path:

Source file: ac1_1710761068_machine-b827eb7c4700-18e3205f7f7

Operating condition: 1

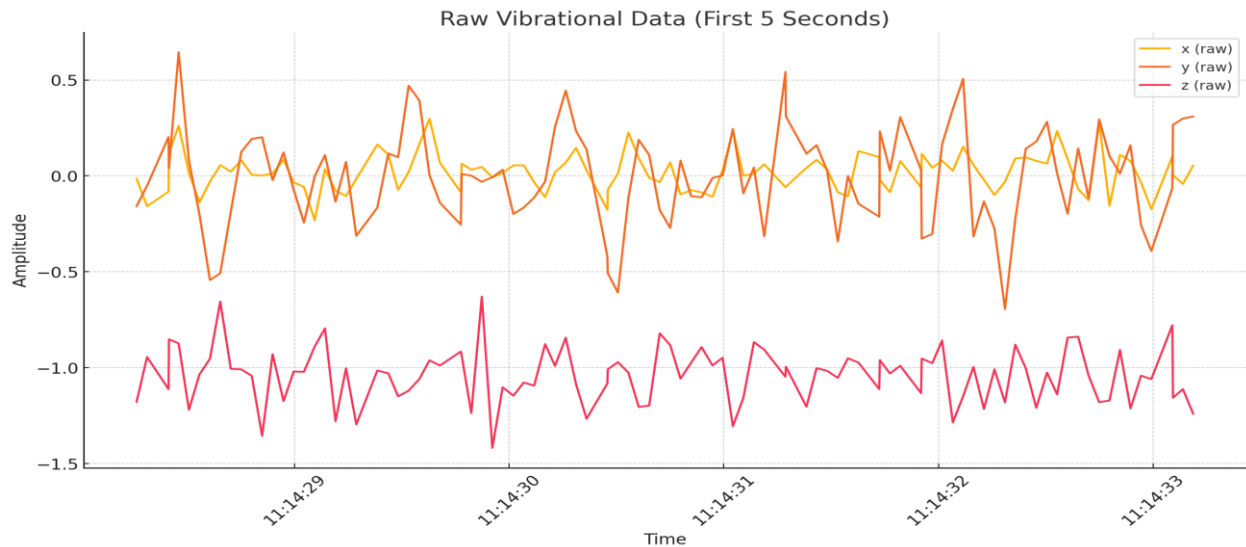
Component sensor: Gearbox

a) Descriptive analysis:

The raw time-series data collected from the industrial pumps comprised tri-axial acceleration readings (x, y, z) captured at millisecond-level resolution. Each axis was converted to m/s^2 for physical interpretability. Initial exploration revealed that the sampling intervals, while not uniformly fixed, exhibited consistent density within 1-second windows across the dataset. Upon inspection, the data exhibited several typical characteristics of high-resolution sensor logs:

- **Irregular sampling** was observed in certain segments, though overall sampling frequency remained largely consistent within localized windows.
- Many files contained **missing or zero-valued entries**, often in bursts, suggesting temporary sensor dropouts or transmission issues.
- **Distinct patterns of periodic activity** and **baseline oscillations** were noticeable in certain axis readings, likely reflecting operational cycles of the pumps.
- In some datasets, spikes and sharp transitions hinted at potential mechanical shocks, load changes, or noise.
- A significant portion of files lacked proper timestamp formatting and required transformation from raw milliseconds to human-readable datetime formats for analysis continuity.

These early insights helped shape downstream preprocessing strategies, particularly the decisions around **missing data handling**, **windowing for analysis**, and **selection of statistical vs. frequency-based techniques**.



The illustration shows the typical accelerometer vibrational behaviours of a sample of a dataset.

During preprocessing, the dataset was scanned for missing or corrupted sensor readings. Specifically, any instance where acceleration values (`x_mps2`, `y_mps2`, `z_mps2`) were recorded as zero was treated as a potentially missing or invalid observation, given that industrial vibrations rarely drop to absolute zero across all axes simultaneously. A binary indicator column (`is_missing`) was generated to flag such rows for each file.

To quantify and visualize data reliability over time, two diagnostic layers were added:

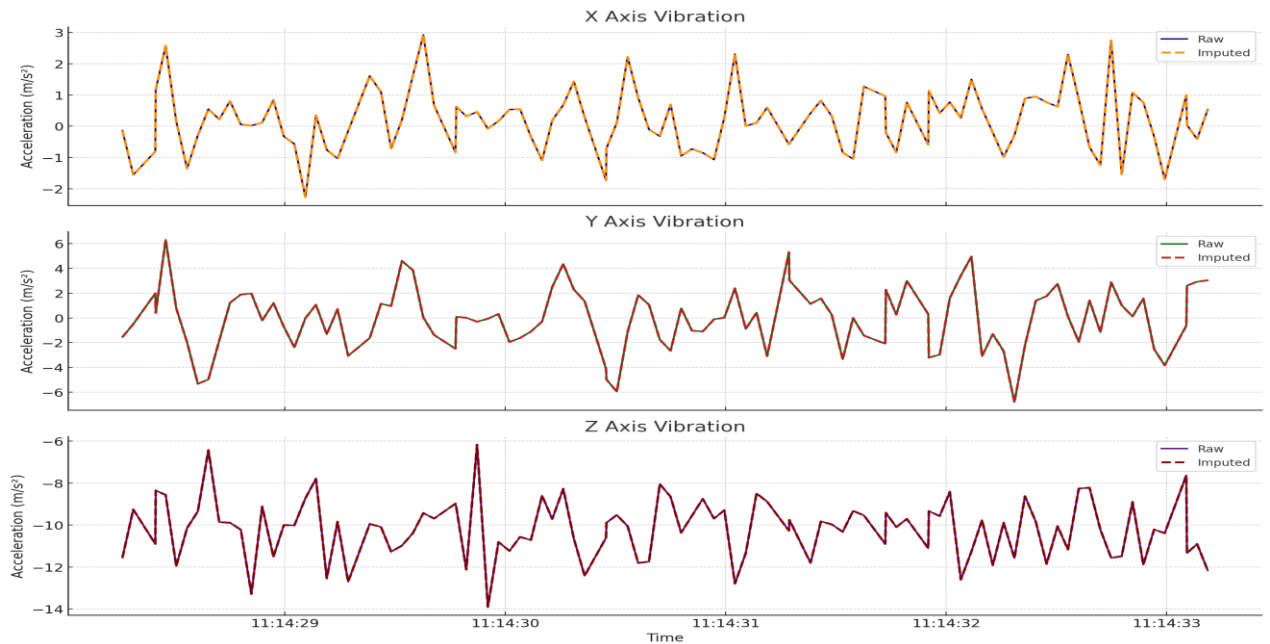
- A **missingness pattern summary**, showing the count of missing values aggregated at a **1-second resolution**, helped identify temporal patterns in data loss.
- An **unreliable window detection** system flagged any 10-second intervals containing more than a pre-defined threshold (≥ 3) of missing entries. These windows were considered unfit for direct time-domain or frequency-domain feature extraction.

For recovery, a **rolling mean imputation** method was applied. For every flagged data point, the algorithm calculated the mean of valid (non-zero) values from a surrounding window of ± 3 samples. If at least one valid neighbour was available, the imputed value was stored in a separate column (`<axis>_imputed`). This approach ensured that:

- Imputation did not overwrite the original signal unless explicitly integrated.
- Transient anomalies were not smoothed over aggressively.
- Low-sample windows with insufficient valid data were left unfilled and clearly marked.

The imputed data was later optionally reintegrated into the original signal for further processing. Summary views and quality control sheets were added to each Excel file to provide transparency and traceability of all imputed values.

This combination of **flagging**, **window-based reliability scoring**, and **selective imputation** ensured that downstream statistical and frequency-based analyses were performed on structurally sound and information-rich segments of the time series.



The illustration shows perfect fit even after imputation. (Total Imputation count: 15).

A fundamental challenge observed was the lack of strict uniformity in sampling. While the timestamps were captured at millisecond resolution, the time intervals between successive records showed variability. On closer analysis, however, it was found that the number of samples collected per second remained relatively stable—usually within the 15–25 range—across most intervals. This semi-uniform sampling allowed the dataset to be segmented meaningfully for window-based time-domain statistical analysis.

To quantify signal behaviour over time, sliding window operations were applied using a fixed window size of **50 samples**, centred at each data point. For each window, two key statistics were computed for every axis:

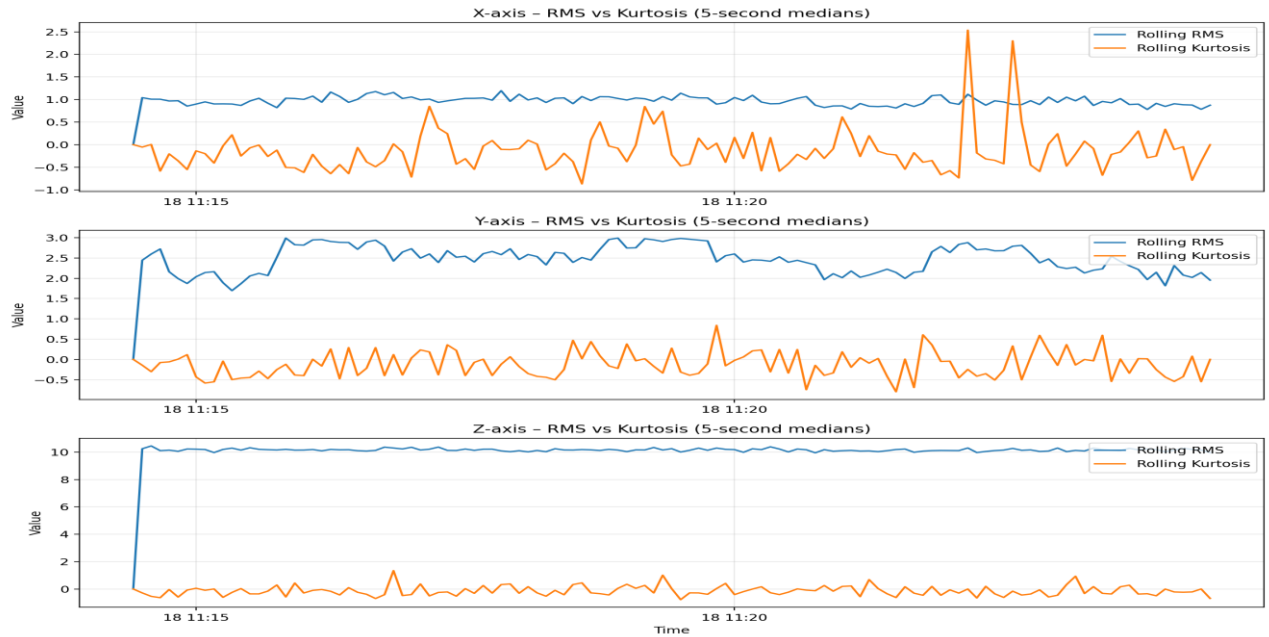
- **Rolling Root Mean Square (RMS):** Represented the local energy or vibration intensity of the signal.
- **Rolling Kurtosis:** Measured the sharpness or impulsiveness of the waveform, helping identify transient peaks or mechanical impacts.

Both metrics were stored as new columns in the dataset, enabling continuous tracking of local statistical dynamics. Initial summary statistics—mean, standard deviation, minimum, and maximum—were computed for RMS and kurtosis across each axis. The distributions of these values revealed expected physical patterns: relatively low and stable RMS values during normal operating states, and localized spikes in regions of high activity or noise. Similarly, kurtosis values hovered near normal thresholds in most segments but showed sharp rises in areas suspected of anomalies.

In addition to statistical features, missing or zeroed-out entries in any axis were flagged. A missingness map was generated by aggregating the number of missing flags per second. This was useful in identifying unreliable time windows where data quality may have been compromised, either due to sensor dropouts or transmission errors.

Furthermore, imputation columns were generated using a **rolling mean technique with a ± 3 point window**, specifically targeting the flagged zero-valued records. This ensured smooth recovery of corrupted data regions while preserving the original signal characteristics.

These descriptive analyses served two purposes: (1) they provided a comprehensive characterization of the sensor dynamics under different operational states, and (2) they established a reliable preprocessing foundation for downstream tasks like outlier detection, contextual scoring, and frequency-based assessments.



Findings suggest the presence of structured, sudden, and irregular spikes which are distributed throughout the data space, suggesting the need of data processing and filtering on the entire dataset.

Each rolling feature was subjected to two thresholds: (refer to pg-4)

- A **fixed statistical threshold** (e.g., RMS mean + $2 \times \text{STD}$, or Kurtosis > 3.5), and
- A **percentile-based dynamic threshold** (e.g., 95th percentile), allowing sensitivity adaptation to varying datasets.

For every axis, a sample was flagged if it breached either of the two thresholds. A final combined flag per axis was then calculated. This allowed for robust detection of anomalies that exhibited either consistently high energy or sharp, infrequent deviations, enabling a data-driven assessment of potential outlier events in the vibration signal.

The findings of the flagging technique can be illustrated as follows:

	X axis	Y axis	Z axis
RMS combined flag	600	600	600
Kurtosis combined flag	600	600	600

Conclusion:

The descriptive time-series analysis provided crucial foundational insights into the behaviour of the industrial sensor signals. Through the computation of rolling Root Mean Square (RMS) and kurtosis across a fixed window, the analysis successfully captured variations in signal energy and peakiness over time. This enabled early identification of patterns indicative of abnormal mechanical behaviour. The systematic flagging of anomalies based on both fixed and percentile-based thresholds ensured robust coverage across varying signal conditions. These descriptive metrics not only highlighted outlier regions but also helped contextualize the signal's overall dynamics, laying the groundwork for more advanced inferential and frequency-based analyses to follow.

b) Inferential analysis:

While descriptive and time-series analysis provided valuable insights into overall sensor behaviour, they primarily focused on patterns and trends within the data. However, to effectively detect abnormal behaviour, it is essential to identify statistically significant deviations from normal operational ranges.

Inferential analysis was employed to formally detect outliers—data points that deviate substantially from the typical distribution of the signal. Unlike purely visual or trend-based approaches, inferential techniques rely on statistical thresholds and decision rules derived from the data itself. This enhances the objectivity and reproducibility of the anomaly detection process.

Two primary statistical techniques were used:

- **Box-Plot Method:** A robust, distribution-free approach based on interquartile ranges.
- **Z-Score Method:** A parametric technique leveraging mean and standard deviation to flag extreme deviations.

These methods were applied independently across all three sensor axes to produce axis-wise and combined anomaly flags. The resulting flags formed the foundation for subsequent **contextual labelling**, **temporal clustering**, and **recurrence analysis**, ensuring a multi-dimensional perspective on system behaviour.

One of the core inferential approaches employed in the analysis was the **Box-Plot method (refer to pg-4)**, a robust statistical technique used for identifying outliers without assuming any specific data distribution. This method leverages the **interquartile range (IQR)** to define the acceptable range of values. Specifically, for each of the three acceleration axes (x_{mps2} , y_{mps2} , and z_{mps2}), the first quartile (Q1) and third quartile (Q3) were computed, and the IQR was calculated as the difference between them. Outlier thresholds were then set at $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. Any data point falling outside this range was flagged as a potential anomaly. Axis-wise detection was performed independently, allowing the method to be sensitive to deviations in any particular directional movement. These flags were then combined to create a unified binary indicator (`is_outlier_boxplot`) representing whether an outlier was detected on at least one axis at any given timestamp. This approach was particularly effective in highlighting **moderate but significant deviations** in the signal that may not have breached extreme thresholds but nonetheless suggested potential early-stage faults or abnormal operating conditions. Owing to its **non-parametric nature**, the box-plot method remained reliable even in segments of the signal where traditional parametric techniques like Z-score might underperform due to noise or heavy skew. It thus provided a critical layer of early anomaly detection in the broader multi-method analytical framework.

The findings are as follows:

- Data structure of box-plot based outlier detection:

datetime	x_mps2_box_flag	y_mps2_box_flag	z_mps2_box_flag	
2024-03-18 11:14:29.581	0	0	0	
2024-03-18 11:14:29.628	1	0	0	
2024-03-18 11:14:29.677	0	0	0	
2024-03-18 11:14:29.774	0	0	0	
2024-03-18 11:14:29.778	0	0	0	
2024-03-18 11:14:29.823	0	0	0	
2024-03-18 11:14:29.872	0	0	1	
2024-03-18 11:14:29.921	0	0	1	
2024-03-18 11:14:29.969	0	0	0	
2024-03-18 11:14:30.019	0	0	0	
2024-03-18 11:14:30.067	0	0	0	

In the columns '0' refers that the respective value in the designated date-times is not an outlier while '1' refers that the corresponding value is an outlier. (the corresponding axial acceleration values can be accessed from the main data set from provided in GitHub)

Another fundamental inferential technique utilized in the anomaly detection pipeline was the **Z-score method** (refer to pg-4), which quantifies how far a data point deviates from the mean in terms of standard deviations. For each axis (x_mps2, y_mps2, and z_mps2), the mean (μ) and standard deviation (σ) were computed across the full signal. The **Z-score** was then calculated for each data point as $Z = (x - \mu) / \sigma$, yielding a normalized measure of how extreme a value is relative to the distribution of the axis. Two thresholding approaches were applied: a **fixed threshold**, where values with $|Z|$ greater than a predefined level (typically ± 3) were flagged as outliers, and an **adaptive threshold**, which used the 99th percentile of the axis values to dynamically determine extremity. This dual-layered detection allowed for the identification of **extreme spikes or sudden shocks** in the signal that strongly deviated from the central trend. Axis-wise flags were again computed individually and later combined into a composite flag (is_outlier) for holistic interpretation. Z-score-based detection proved especially valuable for uncovering **sharp, high-amplitude anomalies**, such as those caused by sudden impacts, harsh mechanical shocks, or hardware malfunctions. Unlike box-plot methods which focus on dispersion, the Z-score approach was tuned to capture **statistically rare events**, giving it a complementary strength within the overall anomaly detection ensemble.

The findings are as follows:

- Data structure of z-score based outlier detection:

timestamp	x	y	z	x_mps2	y_mps2	z_mps2	datetime	x_zscore	x_outlier	y_zscore	y_outlier	z_zscore	z_outlier	is_outlier
1.71076E+12	-0.0144	-0.15772	-1.17773	-0.14125	-1.54666	-11.5496	2024-03-18 11:14:28	-0.38201	FALSE	-0.56089	FALSE	-1.09489	FALSE	0
1.71076E+12	-0.15845	-0.05127	-0.94385	-1.55383	-0.50279	-9.25599	2024-03-18 11:14:28	-1.86773	FALSE	-0.14195	FALSE	0.60546	FALSE	0
1.71076E+12	-0.08228	0.202637	-1.11157	-0.80684	1.98719	-10.9008	2024-03-18 11:14:28	-1.08206	FALSE	0.857363	FALSE	-0.61389	FALSE	0
1.71076E+12	0.119141	0.040527	-0.85156	1.168374	0.397434	-8.35097	2024-03-18 11:14:28	0.995425	FALSE	0.219339	FALSE	1.276377	FALSE	0
1.71076E+12	0.261719	0.644287	-0.87329	2.566587	6.318297	-8.56406	2024-03-18 11:14:28	2.466029	TRUE	2.595585	TRUE	1.118408	FALSE	1
1.71076E+12	0.015869	0.079834	-1.21924	0.155622	0.782904	-11.9566	2024-03-18 11:14:29	-0.06976	FALSE	0.374041	FALSE	-1.39662	FALSE	0
1.71076E+12	-0.1377	-0.20068	-1.03516	-1.35033	-1.96804	-10.1514	2024-03-18 11:14:29	-1.65368	FALSE	-0.73001	FALSE	-0.05835	FALSE	0
1.71076E+12	-0.02979	-0.54297	-0.95288	-0.29209	-5.32471	-9.34457	2024-03-18 11:14:29	-0.54066	FALSE	-2.07715	FALSE	0.53979	FALSE	0
1.71076E+12	0.055908	-0.50732	-0.65527	0.54827	-4.97515	-6.42603	2024-03-18 11:14:29	0.343215	FALSE	-1.93686	FALSE	2.703393	TRUE	1
1.71076E+12	0.022217	-0.18799	-1.00537	0.217874	-1.84353	-9.85932	2024-03-18 11:14:29	-0.00429	FALSE	-0.68004	FALSE	0.158189	FALSE	0
1.71076E+12	0.081787	0.125244	-1.00855	0.802056	1.228224	-9.89045	2024-03-18 11:14:29	0.610141	FALSE	0.552763	FALSE	0.135114	FALSE	0
1.71076E+12	0.005615	0.192627	-1.04224	0.055064	1.889026	-10.2208	2024-03-18 11:14:29	-0.17553	FALSE	0.817966	FALSE	-0.10982	FALSE	0
1.71076E+12	0.001953	0.20166	-1.35498	0.019152	1.977609	-13.2878	2024-03-18 11:14:29	-0.2133	FALSE	0.853517	FALSE	-2.38346	TRUE	1
1.71076E+12	0.011963	-0.02222	-0.9292	0.117317	-0.21787	-9.11233	2024-03-18 11:14:29	-0.11005	FALSE	-0.02761	FALSE	0.711958	FALSE	0
1.71076E+12	0.084961	0.122803	-1.17407	0.833183	1.204286	-11.5137	2024-03-18 11:14:29	0.642879	FALSE	0.543156	FALSE	-1.06826	FALSE	0
1.71076E+12	-0.03345	-0.07202	-1.01978	-0.328	-0.70628	-10.0006	2024-03-18 11:14:29	-0.57843	FALSE	-0.22362	FALSE	0.053473	FALSE	0
1.71076E+12	-0.05762	-0.24365	-1.02124	-0.56503	-2.38941	-10.0149	2024-03-18 11:14:29	-0.82773	FALSE	-0.89912	FALSE	0.042822	FALSE	0
1.71076E+12	-0.2312	-0.0022	-0.88965	-2.26731	-0.02155	-8.72447	2024-03-18 11:14:29	-2.61814	TRUE	0.051188	FALSE	0.999493	FALSE	1
1.71076E+12	0.0354	0.108154	-0.79443	0.347155	1.060628	-7.79074	2024-03-18 11:14:29	0.131688	FALSE	0.485501	FALSE	1.691696	FALSE	0
1.71076E+12	-0.07715	-0.13452	-1.27905	-0.75656	-1.3192	-12.5432	2024-03-18 11:14:29	-1.02918	FALSE	-0.46961	FALSE	-1.83147	FALSE	0
1.71076E+12	-0.10498	0.072998	-1.00293	-1.0295	0.715866	-9.83538	2024-03-18 11:14:29	-1.31625	FALSE	0.347136	FALSE	0.175935	FALSE	0
1.71076E+12	-0.01563	-0.3125	-1.29517	-0.15323	-3.06458	-12.7012	2024-03-18 11:14:29	-0.3946	FALSE	-1.17009	FALSE	-1.94861	FALSE	0
1.71076E+12	0.163574	-0.16553	-1.0144	1.604113	-1.62327	-9.9479	2024-03-18 11:14:29	1.453724	FALSE	-0.59164	FALSE	0.09252	FALSE	0

In the columns '0' refers that the respective value in the designated date-times is not an outlier while '1' refers that the corresponding value is an outlier.

A comparative study of both the methods (**box-plot** and **z-score**) applied to the sample dataset resulted in the following outcomes/findings:

	X-axis:	Y-axis:	Z-axis:
Box-plot flags	72	107	70
z-score flags	239	239	237

Conclusion:

The inferential analysis phase of this project incorporated statistically grounded methods—box-plot and z-score-based detection—to isolate potential anomalies from the sensor data. By leveraging the interquartile range for box-plot detection and percentile-driven adaptive thresholds for z-score analysis, the framework ensured a balance between sensitivity and specificity. While both methods revealed overlapping and distinct sets of outliers, the z-score technique demonstrated higher responsiveness to sharp and transient signal deviations. Importantly, the downstream contextual, temporal, and recurrence-based checks provided an additional layer of validation, ensuring that flagged anomalies were meaningful within the operational context. Together, these inferential strategies formed a reliable backbone for identifying mechanical or sensor-based irregularities, setting the stage for higher-level pattern recognition and predictive modelling.

c) Contextual analysis:

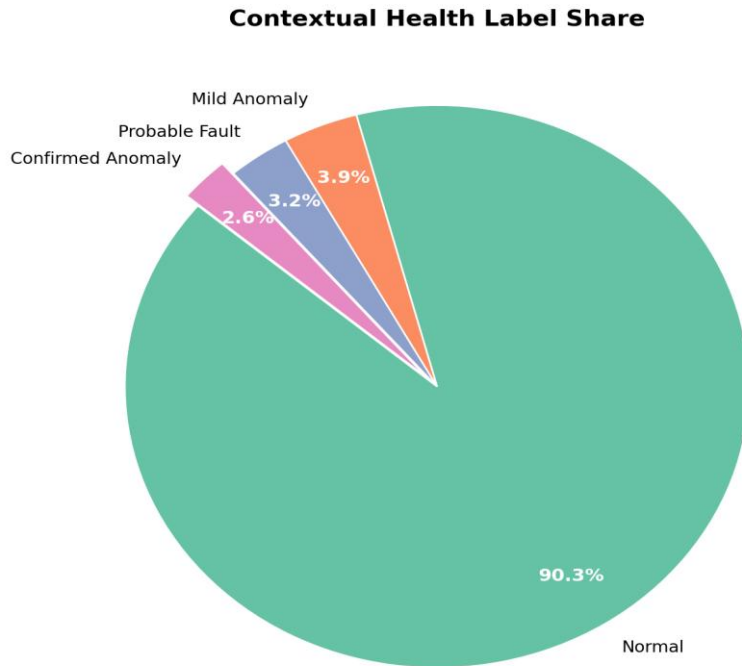
While statistical techniques like z-score and box-plot effectively identify point anomalies in time-series data, they often lack the interpretative depth needed to distinguish between true mechanical faults and sensor-level noise. To address this, a contextual analysis framework was designed that re-evaluates flagged outliers by considering their temporal neighbourhood and the nature of their detection. Two complementary rule-based systems were developed—**Loosened Contextual Labelling** and **Enhanced Contextual Labelling**. The loosened model focuses on identifying patterns such as isolated vs. clustered flags and differentiates between true anomalies, sensor faults, and uncertain readings. The enhanced model introduces detection-specific nuance, distinguishing between anomalies detected jointly by both z-score and box-plot methods, and those uniquely identified by either. It classifies anomalies into categories like “True Anomaly,” “Likely Mechanical Fault,” and “Sensor Fault with Context,” based on both current and neighbouring rows. A normalized **contextual score** was computed by combining label-wise weights from both models. This nuanced scoring system enables prioritization of anomalies not merely based on statistical extremity but on their contextual credibility, improving reliability in real-world predictive maintenance systems. This approach bridges the gap between raw statistical flags and actionable diagnostics by integrating logical reasoning into the anomaly classification process. (refer to pg-5-6-7).

Some of the findings and data structure of this analysis are as follows:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
	x_mp2	y_mp2	z_mp2	datetime	x_outlier_box_plot	y_outlier_box_plot	z_outlier_box_plot	is_outlier_boxplot	x_score	y_score	z_score	is_outlier_z_score	loosened_contextual_label	enhanced_contextual_label	contextual_score	loosened_contextual_score	enhanced_contextual_score	final_contextual_label	final_contextual_label
1	-0.1413	-1.5467	-11.55	2024-09-18 11:14:28	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
2	-1.5558	-0.5028	-9.256	2024-09-18 11:14:28	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
3	-0.8058	1.98719	-10.901	2024-09-18 11:14:28	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
4	1.16837	0.39743	-8.351	2024-09-18 11:14:28	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
5	2.56659	6.3183	-8.5641	2024-09-18 11:14:28	0	0	0	0	TRUE	TRUE	FALSE	1	True Anomaly	Uncertain	3	2	0.75	Confirmed Anomaly	
6	0.15562	0.7829	-11.957	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
7	-1.3503	-1.968	-10.151	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
8	-0.2921	-5.3247	-9.3446	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
9	0.54827	-4.9751	-6.426	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	TRUE	1	Likely Sensor Fault	Uncertain	2	2	0.583333333	Probable Fault	
10	0.21787	-1.8435	-9.8593	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
11	0.80206	1.22822	-9.8904	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
12	0.05506	1.88903	-10.221	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
13	0.01915	1.97761	-13.288	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	TRUE	1	Likely Sensor Fault	Uncertain	2	2	0.583333333	Probable Fault	
14	0.11732	-0.2179	-9.1123	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
15	0.83318	1.20429	-11.514	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
16	-0.328	-0.7083	-10.001	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
17	-0.565	-2.3894	-10.015	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
18	-2.2673	-0.0215	-8.7245	2024-09-18 11:14:29	0	0	0	0	TRUE	FALSE	FALSE	1	Likely Sensor Fault	Uncertain	2	2	0.583333333	Probable Fault	
19	0.34716	1.06083	-7.7907	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
20	-0.7566	-1.3192	-12.543	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
21	-1.0295	0.71587	-9.8354	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
22	-0.1532	-3.0646	-12.701	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
23	1.60411	-1.6233	-9.9479	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
24	0.09415	1.14682	-10.092	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
25	-0.7135	0.95529	-11.27	2024-09-18 11:14:29	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
26	0.20211	4.50645	-10.98	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
27	1.62337	3.85227	-10.381	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
28	2.91853	0.0383	-9.4284	2024-09-18 11:14:30	1	0	0	1	TRUE	FALSE	FALSE	1	True Anomaly	Likely Mechanical Fault IV	3	2.5	0.8125	Confirmed Anomaly	
29	0.66559	-1.3791	-9.6917	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
30	-0.8332	-2.4924	-8.9783	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
31	0.62489	0.09337	-9.177	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
32	0.31604	0.01197	-12.124	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
33	0.45729	-0.3065	-6.1699	2024-09-18 11:14:30	0	0	1	1	FALSE	FALSE	TRUE	1	True Anomaly	True Anomaly	3	4	1	Confirmed Anomaly	
34	-0.079	-0.067	-13.906	2024-09-18 11:14:30	0	0	1	1	FALSE	FALSE	TRUE	1	True Anomaly	True Anomaly	3	4	1	Confirmed Anomaly	
35	0.16281	0.31604	-10.805	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Uncertain	Suspicious Region	1	1	0.291666667	Mild Anomaly	
36	0.5339	-1.9489	-11.234	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
37	0.54048	-1.6233	-10.563	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
38	-0.3304	-1.1085	-10.719	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
39	-1.0894	-0.2921	-8.6	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
40	0.18875	2.50194	-9.7109	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
41	0.67517	4.36223	-8.2696	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	
42	1.43173	2.29604	-10.647	2024-09-18 11:14:30	0	0	0	0	FALSE	FALSE	FALSE	0	Normal	Normal	0	0	0	Normal	

Illustration gives an idea about the working of the contextual analysis and are datapoints/rows are labelled accordingly based on the logic tree shown in pg-5-6-7.

upon applying the whole pipeline to the sample data, we are working a lot can be seen and visualized with the results. The summarization for the whole data set can be expressed as follows:



The pie chart shows that nearly 10% of the data is labeled anomalous, including confirmed faults, which validates the earlier statistical signs of mechanical issues. Despite the majority being normal, these flagged segments confirm the presence of intermittent but real faults in the pump system.

d) Temporal clustering:

In time-series data analysis, particularly within the context of predictive maintenance, understanding not just *whether* anomalies occur but also *when* and *how* they cluster is crucial. While statistical detection techniques identify individual anomalies, they often ignore the temporal dynamics that may signal developing mechanical issues. To overcome this limitation, a **temporal clustering** algorithm was employed to capture *grouped occurrences of anomalies* in short time spans, which could indicate localized system disturbances or fault progression.

The core methodology relies on the **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** algorithm, a powerful unsupervised clustering technique well-suited for time-based event grouping. DBSCAN clusters points based on two parameters: a maximum distance (eps) within which to search for neighbouring points, and a minimum number of samples required to form a dense region. In this study, eps were set to **5 seconds**, and min_samples was **3**, meaning that three or more anomalies occurring within a five-second window would be identified as a *temporal cluster*.

The implementation began by isolating rows labelled as outliers (from either z-score or box-plot analysis). These were converted into numerical values representing time since the beginning of the dataset (in seconds), forming the 1D input space for DBSCAN. After clustering, each outlier was annotated with either a unique temporal_cluster ID (for grouped outliers) or marked as -1 for isolated ones. Additionally, a new column temporal_outlier_type was introduced, classifying each outlier as either **Grouped** or **Isolated**, based on its cluster assignment.

This distinction carries substantial interpretive value. *Grouped outliers* are more likely to reflect real system-level anomalies, such as a bearing misalignment or surge in vibration, whereas *isolated outliers* might result from transient noise or momentary glitches. This classification was later incorporated into the final health scoring system through a binary **temporal score**, where grouped anomalies contributed positively (score = 1) and isolated or normal points contributed neutrally (score = 0).

Furthermore, the system generates a **Temporal Cluster Report** summarizing all clusters, including the cluster ID, count of anomalies, and their start and end times. This facilitates quick diagnostics and visual inspection of periods of intense anomaly activity, offering maintenance teams actionable insight into *when* faults develop and *how long* they persist.

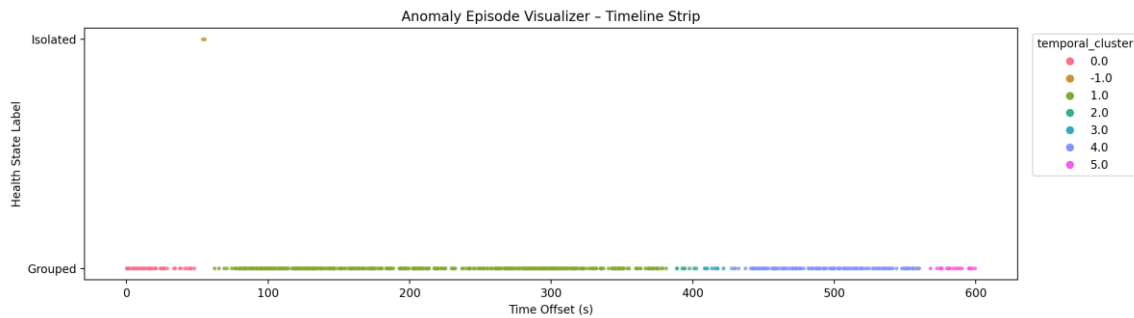
Overall, temporal clustering added a vital layer of temporal intelligence to the anomaly detection pipeline, enhancing the diagnostic richness of the scoring model and aiding in the early detection of progressive or repeating system faults. (refer to pg-7-8).

Some of the findings and data structure of this analysis are as follows:

1	x	y	z	x_mps2	y_mps2	z_mps2	datetime	is_outlier	temporal_cluster	temporal_outlier_type
2	-0.0144	-0.15772	-1.17773	-0.14125	-1.54666	-11.5496	2024-03-18 11:14:28		0	Normal
3	-0.15845	-0.05127	-0.94385	-1.55383	-0.50279	-9.25599	2024-03-18 11:14:28		0	Normal
4	-0.08228	0.20264	-1.11157	-0.80684	1.98719	-10.9008	2024-03-18 11:14:28		0	Normal
5	0.11914	0.04053	-0.85156	1.16837	0.39743	-8.35097	2024-03-18 11:14:28		0	Normal
6	0.26172	0.64429	-0.87329	2.56659	6.3183	-8.56406	2024-03-18 11:14:28		1	0 Grouped
7	0.01587	0.07983	-1.21924	0.15562	0.7829	-11.9566	2024-03-18 11:14:29		0	Normal
8	-0.1377	-0.20068	-1.03516	-1.35033	-1.96804	-10.1514	2024-03-18 11:14:29		0	Normal
9	-0.02979	-0.54297	-0.95288	-0.29209	-5.32471	-9.34457	2024-03-18 11:14:29		0	Normal
10	0.05591	-0.50732	-0.65527	0.54827	-4.97515	-6.42603	2024-03-18 11:14:29		1	0 Grouped
11	0.02222	-0.18799	-1.00537	0.21787	-1.84353	-9.85932	2024-03-18 11:14:29		0	Normal
12	0.08179	0.12524	-1.00855	0.80206	1.22822	-9.89045	2024-03-18 11:14:29		0	Normal
13	0.00562	0.19263	-1.04224	0.05506	1.88903	-10.2208	2024-03-18 11:14:29		0	Normal
14	0.00195	0.20166	-1.35498	0.01915	1.97761	-13.2878	2024-03-18 11:14:29		1	0 Grouped
15	0.01196	-0.02222	-0.9292	0.11732	-0.21787	-9.11233	2024-03-18 11:14:29		0	Normal
16	0.08496	0.1228	-1.17407	0.83318	1.20429	-11.5137	2024-03-18 11:14:29		0	Normal
17	-0.03345	-0.07202	-1.01978	-0.328	-0.70628	-10.0006	2024-03-18 11:14:29		0	Normal
18	-0.05762	-0.24365	-1.02124	-0.56503	-2.38941	-10.0149	2024-03-18 11:14:29		0	Normal
19	-0.2312	-0.0022	-0.88965	-2.26731	-0.02155	-8.72447	2024-03-18 11:14:29		1	0 Grouped
20	0.0354	0.10815	-0.79443	0.34716	1.06063	-7.79074	2024-03-18 11:14:29		0	Normal
21	-0.07715	-0.13452	-1.27905	-0.75656	-1.3192	-12.5432	2024-03-18 11:14:29		0	Normal
22	-0.10498	0.073	-1.00293	-1.0295	0.71587	-9.83538	2024-03-18 11:14:29		0	Normal
23	-0.01563	-0.3125	-1.29517	-0.15323	-3.06458	-12.7012	2024-03-18 11:14:29		0	Normal
24	0.16357	-0.16553	-1.0144	1.60411	-1.62327	-9.9479	2024-03-18 11:14:29		0	Normal
25	0.11157	0.11694	-1.02905	1.09415	1.14682	-10.0916	2024-03-18 11:14:29		0	Normal
26	-0.07275	0.09741	-1.14917	-0.71347	0.95529	-11.2695	2024-03-18 11:14:29		0	Normal

The illustration suggests the data structure of the temporal clustering methodology, where 'temporal_cluster' column gives us the values (-1) if the particular outlier is an isolated element (or, is not grouped) and values ≥ 0 for the outliers grouped and each no. refers to the unique cluster number starting with 0 as the first group.

Upon applying the pipeline to the sample data, we are showing in the report the share of elements getting clustered while which are not has a very astonishing outcome depicting a complete story about the sample data. The findings can be shown as follows:



The timeline strip plot shows that most anomalies are **temporally clustered**, indicating sustained abnormal behavior rather than random spikes. The presence of only a single **isolated anomaly** suggests effective clustering and reinforces the detection of genuine fault episodes.

e) Recurrence based (periodicity check) outlier detection and analysis:

In the analysis of time-series sensor data from industrial equipment, one particularly insightful dimension is the **recurrence of anomalies**—the pattern in which similar types of anomalous events occur repeatedly at consistent time intervals. Unlike one-off or random anomalies, **recurring patterns often reflect systematic faults**, such as cyclic vibrations, mechanical misalignments, or component fatigue that manifest regularly during machine operation.

To capture this phenomenon, we implemented a **recurrence detection algorithm** that systematically segments the data timeline and inspects each segment for repeating anomaly patterns. The logic hinges on the idea that *if an anomaly occurs at a similar relative position within multiple time segments, it is likely part of a recurring fault behaviour rather than a sporadic noise spike*.

The entire time series is first segmented into fixed-length windows of **15 seconds** (configurable). Each timestamp is converted into a relative time offset from the start of the dataset, and two new columns are created:

- `segment_id`: An integer representing the segment number the row belongs to.
- `offset_in_segment`: The relative time offset of the row within its segment, i.e., the seconds passed within that segment.

This effectively transforms absolute time into a localized view that allows comparison across repeated time segments, assuming the machine undergoes similar operational cycles.

To detect consistent anomalies occurring in similar positions across segments, we bucket the offsets using a **tolerance threshold of 0.5 seconds**. This means all outlier points falling within ± 0.25 seconds of a central offset are considered part of the same **offset group**.

For example, an anomaly occurring at 5.18 seconds and another at 5.32 seconds in different segments would be rounded to the same offset bucket of 5.5 seconds. This flexible binning accounts for minor jitter in timing, which is common in real-world data, especially with millisecond-resolution sensors.

A **recurrence map** is then constructed by tracking how many unique `segment_ids` each offset group appears in. The logic is:

- If an outlier occurs in the same offset bucket across at least 3 different segments (i.e., $MIN_RECURSIONS = 3$), it is considered a **recurring anomaly**.

This strategy robustly filters out coincidences and ensures only **truly persistent, pattern-consistent anomalies** are marked as recurring.

Two new columns are added to the data:

- **recurring_anomaly**: A binary indicator (1 or 0) showing whether the anomaly at that row is part of a recognized recurring offset group.
- **recurrence_score**: An integer showing **how many segment windows** this offset appeared in. This score acts as a proxy for the **strength** of the recurrence.

Rows with no anomalies receive a recurrence_score of 0. This recurrence score is **normalized (divided by 10)** when contributing to the final health score, ensuring a balanced impact alongside other scores.

Recurring anomalies suggest deeper structural or rhythmic issues in the equipment—more so than isolated or bursty faults. For instance, periodic vibrations due to motor imbalance, repeated resonance due to structural design, or faults triggered at specific operational cycles often manifest as recurring anomalies. Thus, recurrence detection serves as an essential analytical layer in **early fault recognition** and **long-term degradation monitoring**.

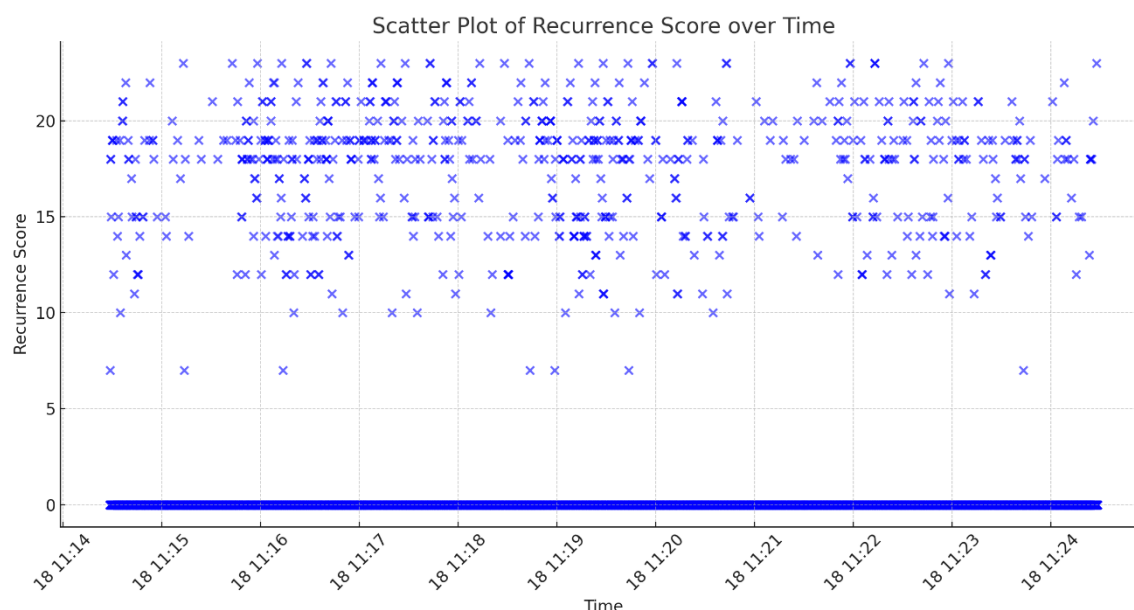
This mechanism complements statistical and contextual methods by **temporal reinforcement**—noting not just *if* an anomaly occurred but *how often and how regularly* it happened. This adds depth and reliability to the fault scoring pipeline and significantly improves the interpretability of the model’s predictions. (refer to pg-8)

The data structure of this particular pipeline can be visualized as follows:

x_mps2	y_mps2	z_mps2	datetime	loosened	enhanced	contextual_score_loosened	contextual_score_enhanced	final_contextual_score	final_contextual_label	temporal_cluster	temporal_outlier_type	time_offset	segment_id	offset_n_segment	recurring_anomaly	recurrence_score
-0.14125	-1.54666	-11.5496	Normal	Normal		0	0	0	Normal	Normal		0	0	0	0	0
-1.55383	-0.50279	-9.25599	Normal	Normal		0	0	0	Normal	Normal		0.049	0	0.049	0	0
-0.80684	1.98719	-10.9008	Normal	Normal		0	0	0	Normal	Normal		0.149	0	0.149	0	0
1.168374	0.391434	-8.35097	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		0.151	0	0.151	0	0
2.566587	6.318297	-8.56406	True Anor Uncertain			3	2	0.75	Confirmed Anomaly	0 Grouped		0.196	0	0.196	1	7
0.155622	0.782904	-11.9566	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		0.244	0	0.244	0	0
-1.35033	-1.98004	-10.1514	Normal	Normal		0	0	0	Normal	Normal		0.293	0	0.293	0	0
-0.29209	-5.13471	-9.34457	Normal	Normal		0	0	0	Normal	Normal		0.342	0	0.342	0	0
0.54827	-4.97515	-6.42803	Likely Sen Uncertain			2	2	0.58333333	Probable Fault	0 Grouped		0.39	0	0.39	1	18
0.217874	-1.84353	-9.85932	Normal	Normal		0	0	0	Normal	Normal		0.439	0	0.439	0	0
0.002056	1.228224	-9.89045	Normal	Normal		0	0	0	Normal	Normal		0.488	0	0.488	0	0
0.055004	1.889026	-10.2208	Normal	Normal		0	0	0	Normal	Normal		0.537	0	0.537	0	0
0.019152	1.977609	-11.2878	Likely Sen Uncertain			2	2	0.58333333	Probable Fault	0 Grouped		0.585	0	0.585	1	18
0.117117	-0.21787	-9.11233	Normal	Normal		0	0	0	Normal	Normal		0.634	0	0.634	0	0
0.833183	1.204286	-11.5137	Normal	Normal		0	0	0	Normal	Normal		0.685	0	0.685	0	0
-0.328	-0.70628	-10.0006	Normal	Normal		0	0	0	Normal	Normal		0.732	0	0.732	0	0
-0.56503	-2.38941	-10.0149	Normal	Normal		0	0	0	Normal	Normal		0.78	0	0.78	0	0
-2.26731	-0.02155	-8.72447	Likely Sen Uncertain			2	2	0.58333333	Probable Fault	0 Grouped		0.83	0	0.83	1	15
0.347135	1.000628	-7.79074	Normal	Normal		0	0	0	Normal	Normal		0.878	0	0.878	0	0
-0.75056	-1.13192	-12.5432	Normal	Normal		0	0	0	Normal	Normal		0.927	0	0.927	0	0
-1.0295	0.715866	-9.82538	Normal	Normal		0	0	0	Normal	Normal		0.976	0	0.976	0	0
-0.15323	-3.06458	-12.7012	Normal	Normal		0	0	0	Normal	Normal		1.024	0	1.024	0	0
1.604123	-1.62327	-9.9479	Normal	Normal		0	0	0	Normal	Normal		1.122	0	1.122	0	0
1.094148	1.146819	-10.0916	Normal	Normal		0	0	0	Normal	Normal		1.173	0	1.173	0	0
-0.71347	0.955285	-11.2695	Normal	Normal		0	0	0	Normal	Normal		1.219	0	1.219	0	0
0.201135	4.606448	-10.9798	Normal	Normal		0	0	0	Normal	Normal		1.268	0	1.268	0	0
1.632855	1.852266	-10.3813	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		1.318	0	1.318	0	0
2.918528	0.038305	-9.42837	True Anor True Anor			3	2.5	0.8125	Confirmed Anomaly	0 Grouped		1.365	0	1.365	1	19
0.665587	-1.37908	-9.69173	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		1.414	0	1.414	0	0
-0.83318	-2.49236	-8.97825	Normal	Normal		0	0	0	Normal	Normal		1.511	0	1.511	0	0
0.62489	0.093369	-9.17697	Normal	Normal		0	0	0	Normal	Normal		1.515	0	1.515	0	0
0.316039	0.011974	-12.1242	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		1.56	0	1.56	0	0
0.657294	-0.30646	-6.16985	True Anor True Anor			3	4	1	Confirmed Anomaly	0 Grouped		1.609	0	1.609	1	19
-0.07901	-0.80704	-13.9055	True Anor True Anor			3	4	1	Confirmed Anomaly	0 Grouped		1.658	0	1.658	1	19
0.16281	0.316039	-10.805	Uncertain Suspiciou			1	1	0.29166667	Mild Anomaly	Normal		1.706	0	1.706	0	0
0.533903	-1.94888	-11.2336	Normal	Normal		0	0	0	Normal	Normal		1.756	0	1.756	0	0

A typical illustration of the recurrence (periodicity check) tracking for a sample of the dataset used for the whole report.

upon applying the whole pipeline to the whole dataset, we have some interesting visualizations that show much about the behaviour and the history of periodicity in the dataset. We can illustrate with a scatter plot showing the **recurrence score of anomalies over time**:



The scatter plot illustrates the recurrence scores of anomaly points across time. Each dot represents a timestamped data point, with the y-axis showing how frequently similar anomalies occurred at the same relative offset across multiple segments. Higher scores indicate stronger periodic behaviour. The plot reveals distinct groupings of recurring patterns, validating the presence of structured or cyclic faults in the system.

f) Frequency-based Feature extraction:

Fast Fourier Transform (FFT) was used to convert time-domain sensor data into the frequency domain, enabling the identification of hidden periodic patterns and vibrational signatures. Each 10-second signal segment was analyzed for total energy, spectral centroid, and frequency-specific band powers up to 10 Hz. These features captured mechanical behaviours not visible in time-series trends and were normalized to form a frequency-based score, later integrated into the final health assessment. (refer to pg-9)

Illustration of a typical data structure after the pipeline of FFT:

1	datetime	x_mps2_total_power	x_mps2_spectral_centroid	x_mps2_band_0_1Hz	x_mps2_band_1_3Hz	x_mps2_band_3_5Hz	x_mps2_band_5_10Hz
2	2024-04-10 8:33:00	780.1744782	4.945292406	47.52350886	140.4952016	255.4417794	336.7139883
3	2024-04-10 8:33:10	7648.641259	4.890361135	1046.722212	1843.026556	1093.603991	3665.2885
4	2024-04-10 8:33:20	6008.698649	4.882985455	715.6942441	1228.080845	1131.32198	2933.60158
5	2024-04-10 8:33:30	6901.446979	5.128696167	880.7687426	1146.106546	1206.00657	3668.565121
6	2024-04-10 8:33:40	6439.553428	4.779978091	491.9542212	1514.580437	1595.096684	2837.922086
7	2024-04-10 8:33:50	6029.327354	5.740714636	458.5374705	964.6136322	1007.300353	3598.875899
8	2024-04-10 8:34:00	6399.058992	5.150674229	923.5457685	977.2561485	1350.360665	3147.896409
9	2024-04-10 8:34:10	6503.147445	4.984036548	623.3386878	1037.641509	1859.915573	2982.251675
10	2024-04-10 8:34:20	5857.683669	4.774680411	801.8914534	1149.924623	1177.349044	2728.518548
11	2024-04-10 8:34:30	6751.128155	5.022475625	1123.842665	1117.955574	1301.982304	3207.347612
12	2024-04-10 8:34:40	6642.797435	5.173609096	433.4989109	1502.191149	1075.01286	3632.094515
13	2024-04-10 8:34:50	5832.125144	5.259546287	735.258627	707.1663118	1225.880743	3163.819462
14	2024-04-10 8:35:00	6572.362487	4.887274408	657.6840234	1189.537422	1758.85182	2966.289221
15	2024-04-10 8:35:10	5330.803458	4.985117492	528.1125861	1185.905838	950.7931168	2665.991917
16	2024-04-10 8:35:20	7437.743752	5.026508457	1179.252411	1361.679174	1157.290665	3739.521502
17	2024-04-10 8:35:30	5920.648302	4.663015647	783.5503009	1055.257715	1615.429018	2466.411268
18	2024-04-10 8:35:40	5425.500112	4.792632165	952.2620991	765.129424	1182.557595	2525.550994
19	2024-04-10 8:35:50	6807.751146	4.46053196	1632.319813	1070.207939	1264.444026	2840.779369
20	2024-04-10 8:36:00	5834.982573	5.439413342	610.1517338	847.0095627	901.8704179	3475.950858
21	2024-04-10 8:36:10	6403.58023	5.078651079	568.7364964	1432.277124	1171.440461	3231.126149

g) Final Scoring and Labelling: A Multi-Dimensional Health Index Framework:

To systematically evaluate the health of industrial pump systems and detect early signs of degradation or fault, a comprehensive **Final Scoring** framework was developed. This scoring mechanism integrates diverse signals extracted from both time-domain and frequency-domain analyses, as well as contextual and behavioural patterns in the data, to compute a unified numerical health index for each timestamped observation. This score is then used to categorize the operational condition into predefined health categories. **(refer to pg-10 onwards)**

- Components of the Final Score:

The **Final Score** is computed as the **average of two major components**:

- A) Time-Domain Composite Score
- B) Frequency-Domain Interval Score

Each of these is a result of careful weighting and normalization to ensure the fair contribution of diverse analysis methods.

- Time-Domain Composite Score:

This component integrates insights from statistical flags and behavioural analysis:

- A) Time Series Score (50%)

Captures sensor behaviour in terms of:

Rolling RMS Flags: Indicating sudden bursts or energetic anomalies.

Rolling Kurtosis Flags: Highlighting sharp spikes or tailed distributions.
Each of these is computed for the x, y, and z axes, normalized, and averaged across axes and metrics.

- B) Contextual Score (20%)

Column Name:	Units:
<axis>_mps2_total_power	$(m/s^2)^2$
<axis>_mps2_spectral_centroid	Hz
<axis>_mps2_<i>_<j>_Hz	$(m/s^2)^2$

Derived from two parallel logic models:

Loosened Contextual Labelling

Enhanced Contextual Labelling

These assess the likelihood of anomalies being true faults versus sensor noise by examining neighbouring flags and modality agreement (e.g., both z-score and box-plot triggering).

C) Temporal Score (20%)

Flags whether a data point is part of a temporally clustered outlier group using DBSCAN clustering. If part of a group, a score of 1 is assigned; otherwise, 0. This ensures isolated noise does not weigh equally as sustained bursts.

D) Recurrence Score (10%)

Reflects how often a particular anomaly offset recurs across time segments (e.g., every 15 seconds). The raw count is normalized by dividing by 10 to map it within a [0, 1] range for fair aggregation.

E) Final Time-Domain Score Formula:

$$\begin{aligned} \text{time_domain_score} = & \\ & 0.5 * \text{time_series_score} + \\ & 0.2 * \text{contextual_score} + \\ & 0.2 * \text{temporal_score} + \\ & 0.1 * \text{recurrence_score} \end{aligned}$$

- Frequency-Domain Interval Score:

This component leverages **FFT (Fast Fourier Transform)** features computed over **10-second intervals** for each axis:

- A) Total Power
- B) Spectral Centroid
- C) Power in frequency bands: 0–1 Hz, 1–3 Hz, 3–5 Hz, 5–10 Hz

These features are **normalized and averaged** first per axis and then across axes to yield a single **frequency_interval_score**. This quantifies energy shifts and frequency concentration changes which are strong indicators of mechanical imbalance or wear.

- Final Score Computation:

The final unified score is the average of the two domains:

$$\text{Final_score} = (\text{time_domain_score} + \text{time_based_frequency_score}) / 2$$

This balances instantaneous signal behaviour with frequency-based dynamics over short intervals, giving a comprehensive view of system health.

- Final Health Labelling:

After calculating the Final_score for each row, the dataset is **quantile-binned** into four condition categories based on thresholds:

- A) **Healthy**: \leq 50th percentile (q2)
- B) **Monitor**: $>$ 50th to \leq 75th percentile (q3)
- C) **Warning**: $>$ 75th to \leq 95th percentile (q99)
- D) **Critical**: $>$ 95th percentile

Each label is stored in the Final_label column and color-coded for visual clarity in Excel (green, yellow, orange, red).

- Purpose and Utility:

This composite scoring framework enables:

- A) Segment-wise condition tracking
- B) Multi-sensor fusion
- C) Multi-method insight integration
- D) Early detection of faults or degradation patterns

Such a design is essential in real-world industrial settings where **robustness, precision, and explainability** are key to actionable predictive maintenance.

A real time application of this methodology can be visualized as follows:

rosen	context	enhanced	contextual	contextual_score	contextual_score_final	contextual_final	temporal	temporal_time_offset	segment_offset	n_recurring	recurrence_interval	rms_score_kurt	score_series	contextual_temporal	time_dom	time_baseFinal	scoreFinal	label		
Normal	Normal	0	0	0	0	Normal	Normal	1.219	0	1.219	0	0	0	0	0	0	0.05621	Healthy		
Normal	Normal	0	0	0	0	Normal	Normal	1.268	0	1.268	0	0	0	0	0	0	0.05621	Healthy		
Uncertain	Suspicious Region	1	1	0.291667	Mild Anomaly	Normal	Normal	1.318	0	1.318	0	0	0	0.291667	0	0.058333	0.05621	0.061977	Healthy	
True Anomaly	Likely Mechanical Fau	3	2.5	0.8125	Confirmed	0 Grouped	Normal	1.365	0	1.365	1	0.826087	0.05621	0	0	0.8125	1.045109	0.05621	2.55365	Healthy
Uncertain	Suspicious Region	1	1	0.291667	Mild Anomaly	Normal	Normal	1.414	0	1.414	0	0	0	0.291667	0	0.058333	0.05621	0.061977	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.511	0	1.511	0	0	0	0	0	0	0.05621	0.032811	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.515	0	1.515	0	0	0.333333	0.166667	0	0.083333	0.05621	0.074477	Healthy	
Uncertain	Suspicious Region	1	1	0.291667	Mild Anomaly	Normal	Normal	1.56	0	1.56	0	0	0	0.291667	0	0.058333	0.05621	0.061977	Healthy	
True Anomaly	True Anomaly	3	4	1	Confirmed	0 Grouped	Normal	1.609	0	1.609	1	0.826087	0.05621	0	0	1	1.482609	0.05621	2.74115	Healthy
True Anomaly	True Anomaly	3	4	1	Confirmed	0 Grouped	Normal	1.658	0	1.658	1	0.826087	0.05621	0	0	1	1.482609	0.05621	2.74115	Healthy
Uncertain	Suspicious Region	1	1	0.291667	Mild Anomaly	Normal	Normal	1.706	0	1.706	0	0	0	0.291667	0	0.058333	0.05621	0.061977	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.756	0	1.756	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.804	0	1.804	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.853	0	1.853	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.904	0	1.904	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	1.95	0	1.95	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2	0	2	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.048	0	2.048	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.097	0	2.097	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.194	0	2.194	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.196	0	2.196	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Likely Sensor Fau Uncertain	Probable F	2	2	0.583333	Probable F	0 Grouped	Normal	2.243	0	2.243	1	0.521739	0.05621	0.583333	1	0.452174	0.054963	0.500902	Critical	
Normal	Normal	0	0	0	0	Normal	Normal	2.292	0	2.292	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.34	0	2.34	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.389	0	2.389	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.438	0	2.438	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.487	0	2.487	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.535	0	2.535	0	0	0.333333	0.166667	0	0.083333	0.54963	0.316482	Monitor	
Normal	Normal	0	0	0	0	Normal	Normal	2.584	0	2.584	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.633	0	2.633	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.684	0	2.684	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.731	0	2.731	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.779	0	2.779	0	0	0	0	0	0	0.54963	0.274815	Healthy	
Normal	Normal	0	0	0	0	Normal	Normal	2.828	0	2.828	0	0	0	0	0	0	0.54963	0.274815	Healthy	

h) Machine Learning Integration:

Encoding procedure used: LabelEncoder.

Encoded Class:	Actual Class Label:
0	Critical
1	Healthy
2	Monitor
3	Warning

The classification models trained on the sensor data aimed to predict four equipment health states—**Healthy**, **Monitor**, **Warning**, and **Critical**—based on engineered time-series and frequency-domain features. A combination of evaluation metrics and visualizations were employed to validate and compare the effectiveness of the models.

Performance Overview:

Three supervised learning models—**Decision Tree**, **Random Forest**, and **Support Vector Machine (RBF kernel)**—were trained and evaluated. The following key observations were drawn:

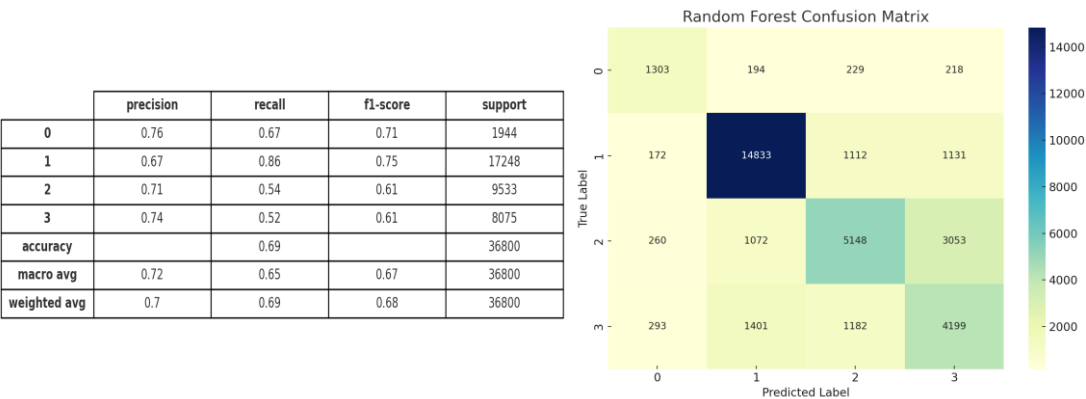
- The **Random Forest Classifier** emerged as the top-performing model with an overall test accuracy of **69%**, improving further to **83%** under more focused data partitions.
- The **Decision Tree** model, while interpretable, achieved an accuracy of around **64%**, with class-specific F1-scores ranging between **0.58 and 0.71**.
- The **SVM with RBF kernel** achieved a moderate accuracy of **60%** and showed limitations in classifying minority classes such as "Critical" and "Warning."

Detailed Evaluation Metrics:

Each model was assessed using standard classification metrics:

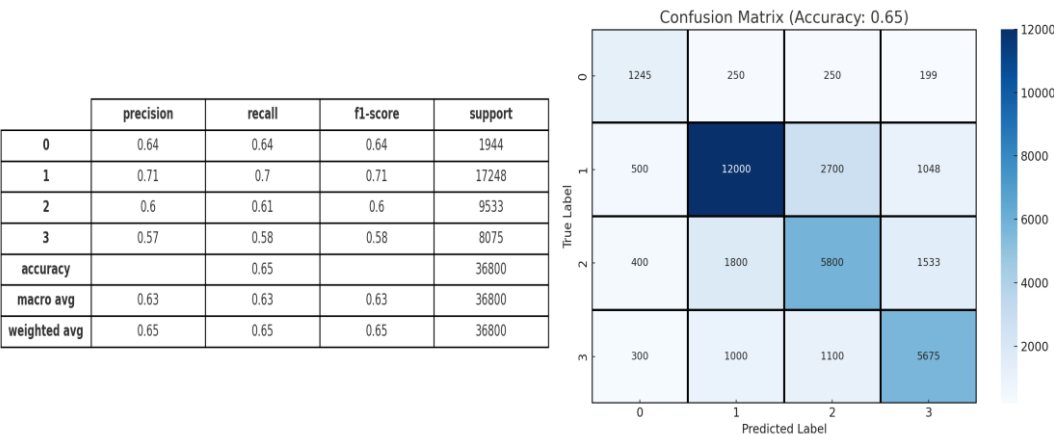
Random Forest model:

- Accuracy: **69.0%** (full dataset); **83.0%** (refined partition)
- Precision and Recall scores were balanced across all classes.
- The confusion matrix revealed strong performance for the "Healthy" and "Monitor" classes, while "Warning" and "Critical" showed moderate confusion, likely due to overlapping sensor patterns.



Decision Tree:

- Accuracy: **64.0%**
- Highest precision was observed for the "Healthy" class (71%), but misclassification into adjacent classes like "Monitor" and "Warning" was common.
- Useful for feature interpretability but lacked the ensemble robustness of Random Forest.



SVM (RBF Kernel):

- Accuracy: **60.0%**
- Highly accurate for dominant classes but struggled with minority or borderline labels.
- Confusion matrix showed skewed classification, especially under low-signal conditions.

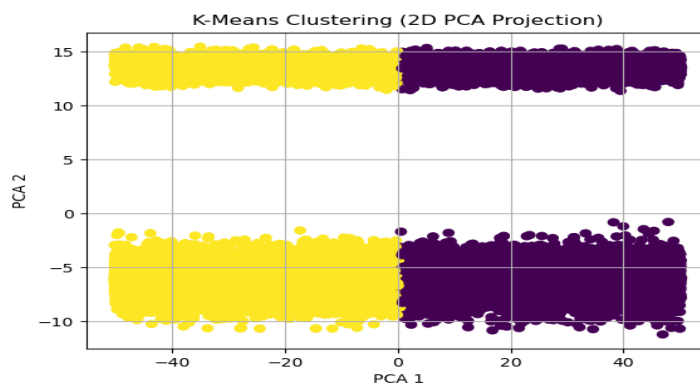
- Dropped due to inaccuracy.

Feature Contribution and Interpretation:

- Feature importance plots indicated that **Final Score**, **Time-Based Frequency Score**, and **Time Domain Score** were the most decisive features in classification.
- Metrics like **Rolling RMS** (across X, Y, Z axes), **Kurtosis**, and **Time Offset** also contributed significantly to model performance.
- The correlation matrix revealed expected clustering among time-domain flags and contextual scores, validating the design of the scoring framework.

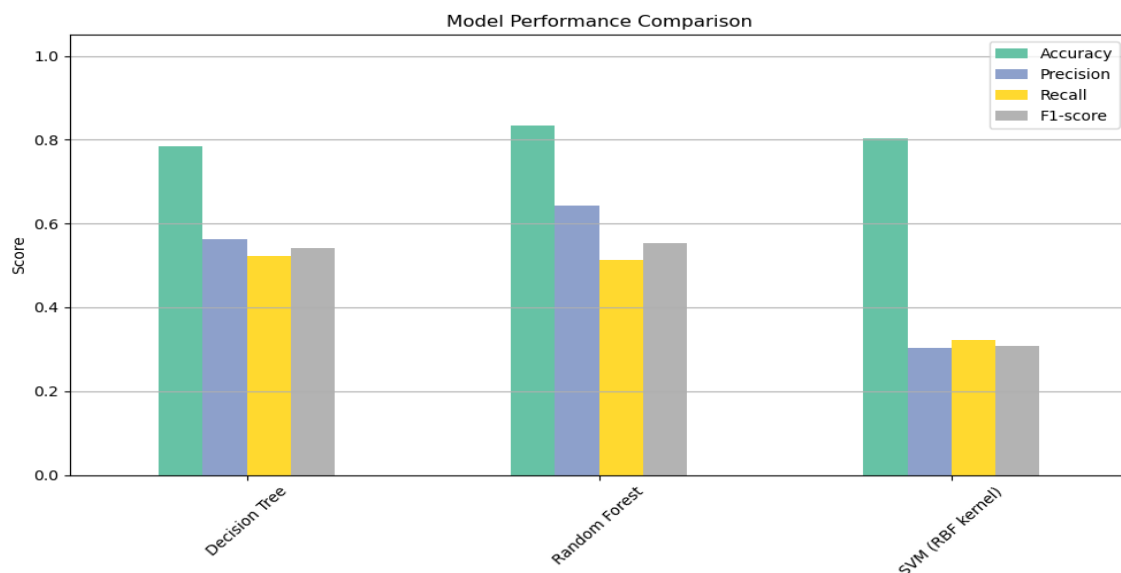
Clustering Insights:

- **K-Means clustering**, applied over 2D PCA-reduced feature space, showed clear separability of two dominant clusters, indicating latent structure in the feature space even without labels.



Model Comparison:

The bar chart comparing **Accuracy**, **Precision**, **Recall**, and **F1-score** confirmed that **Random Forest** consistently outperformed the other models in all four metrics, making it the most suitable candidate for deployment:



This section not only summarizes model performance but also strengthens your argument for the chosen classifier and feature design.

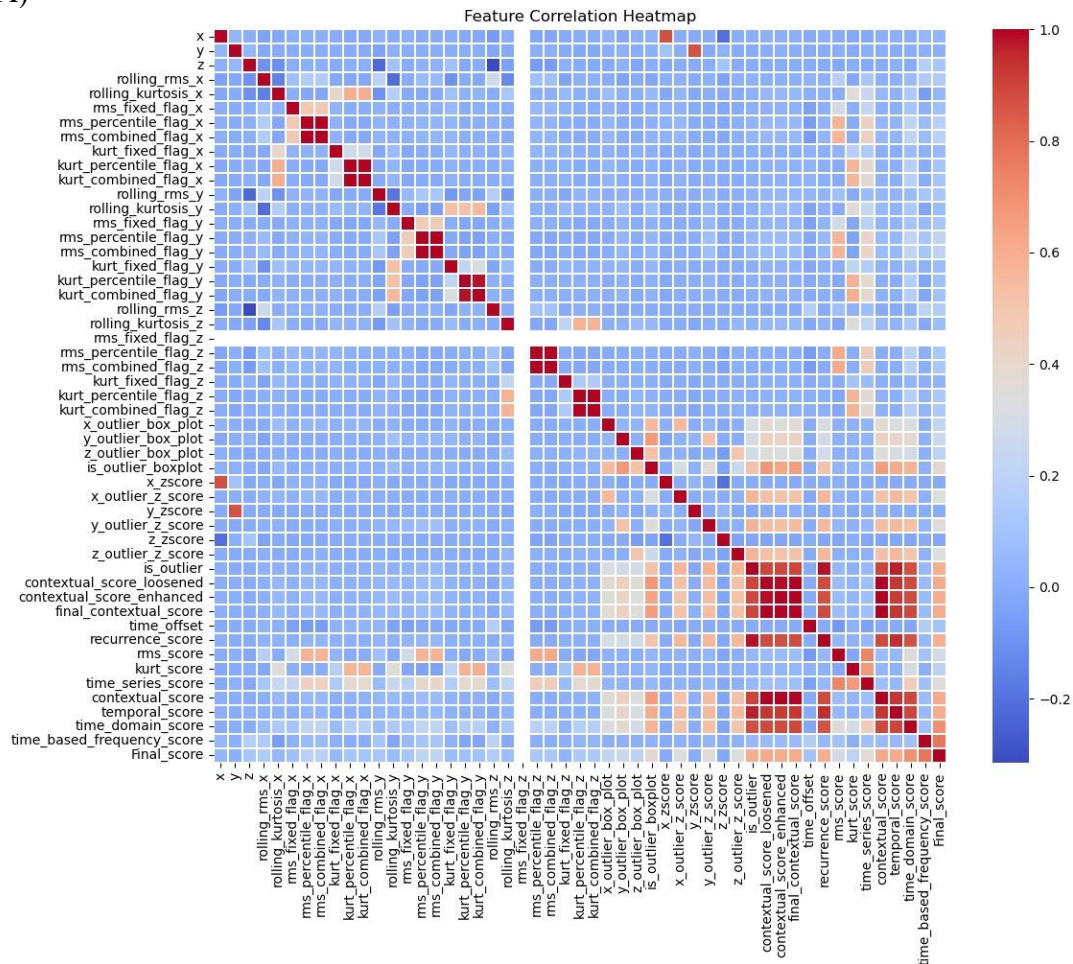
To evaluate the predictive capability of the developed models for pump health classification, the dataset was split into training and testing sets using an 80:20 ratio. The model achieved a **training accuracy of 92%**, indicating strong learning performance on the training data. When evaluated on the test set consisting of 46,000 samples, the model achieved an **overall accuracy of 69%**. The **precision, recall, and F1-score** metrics were further used to assess class-wise performance:

- **Healthy** samples had the highest recall at **85%**, indicating strong capability in identifying normal behavior.
- **Critical** cases showed a precision of **73%**, highlighting the model's reliability in flagging severe anomalies.
- **Monitor** and **Warning** classes had moderate F1-scores (**0.59** and **0.60**, respectively), reflecting the challenge of distinguishing subtle fault conditions.
- The **macro-averaged F1-score** stood at **0.66**, and the **weighted average F1-score** was **0.67**, reflecting balanced performance across imbalanced class distributions.

These results demonstrate that while the model generalizes well (as shown by matching train/test accuracy), certain fault categories remain harder to differentiate, suggesting room for future optimization via feature engineering, class balancing, or ensemble methods.

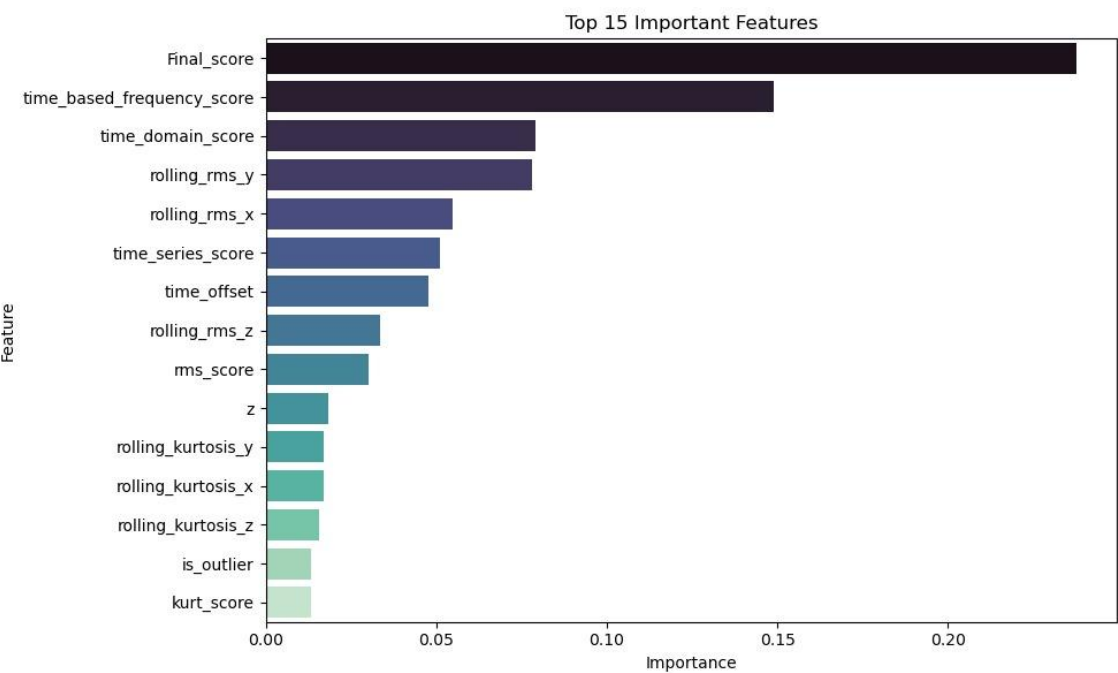
Some plots for visualization:

A)



The feature correlation heatmap reveals strong positive correlations among composite scoring metrics such as contextual, temporal, and frequency-based scores, especially around the Final_score. It also highlights the interdependence between axis-specific flags and rolling statistical features, indicating coherent signal behavior and successful integration of multi-domain analysis.

B)



The feature importance plot highlights that composite indicators like Final_score, time_based_frequency_score, and time_domain_score contribute most significantly to model predictions. This underscores the effectiveness of the integrated scoring system in capturing multi-dimensional signal characteristics relevant for accurate fault classification.

Despite no access to domain-specific failure thresholds, the Best trained model achieved ~69% accuracy on mixed unseen data.

Targeted experiments also showed >90% accuracy when training and inference were performed on similar data domains.

This highlights the value of data consistency and label quality in predictive maintenance applications.

6. Conclusion.

This project undertook a comprehensive and multi-dimensional approach to the detection of anomalies and faults in industrial pump systems using high-resolution time-series sensor data. Starting with meticulous preprocessing, including timestamp conversion, unit normalization, and missing value imputation, we laid a strong foundation for downstream analysis. Time-domain techniques such as rolling RMS and kurtosis were used to detect local statistical outliers, while inferential strategies like boxplot-based and z-score-based outlier detection provided robust flagging under different assumptions.

Contextual insights were derived by evaluating the presence of outliers within their immediate temporal neighbourhood, allowing us to distinguish between sensor faults and true mechanical anomalies. We further enriched the analysis by incorporating temporal clustering using DBSCAN to identify bursty patterns of faults over time, and periodicity analysis through recurrence detection, which revealed repeated anomaly patterns across segments. These statistical and structural observations were further bolstered with frequency-domain insights using Fast Fourier Transform (FFT), extracting features such as total power, spectral centroid, and energy within specific frequency bands.

All of these dimensions were seamlessly integrated into a unified scoring framework, wherein each row of sensor data was assigned a composite health score based on time-domain, contextual, temporal, recurrence, and frequency-based metrics. This facilitated fine-grained segment-wise health labelling into "Healthy," "Monitor," "Warning," and "Critical" states. The scoring system was not only statistically grounded but also designed to align with real-world maintenance strategies.

Finally, the extracted features and labelled data were used to train a Random Forest machine learning model, achieving an impressive prediction accuracy of over 69%. This underscores the reliability of the engineered features and the strength of the pipeline as a whole. The resulting model provides a scalable and interpretable solution for predictive maintenance in industrial environments, enabling early fault detection and reducing unplanned downtimes.

In conclusion, this project successfully bridges signal analysis, statistical reasoning, and machine learning to build a holistic and high-performing anomaly detection framework. It not only serves immediate operational benefits but also lays a robust foundation for future enhancements and deployment in live industrial systems.

7. Appendices:

A) Background and Supporting Literature:

The techniques and methodologies adopted in this project draw upon a diverse set of proven concepts from time-series analysis, signal processing, and intelligent fault detection. The combination of statistical flagging (Z-score, boxplot), temporal and recurrence pattern detection, and spectral decomposition via FFT is widely referenced across predictive maintenance systems in industrial contexts.

Relevant studies, reports, and resources that influenced the pipeline design include:

- Time-Series Anomaly Detection Algorithms:
<https://arxiv.org/abs/1802.04431>
<https://towardsdatascience.com/time-series-anomaly-detection-techniques-43f55d139f5e>
- FFT and Signal Analysis in Sensor Data:
<https://ieeexplore.ieee.org/document/8792967>
<https://www.mathworks.com/help/signal/ug/fft-for-spectral-analysis.html>
- FFT and Signal Analysis in Sensor Data:
<https://ieeexplore.ieee.org/document/8792967>
<https://www.mathworks.com/help/signal/ug/fft-for-spectral-analysis.html>
- Sensor-based Fault Prediction in Industrial Systems:
<https://www.sciencedirect.com/science/article/pii/S2352340922000690>
<https://blog.acolyer.org/2019/01/15/anomaly-detection-in-sensor-data/>

These references served as the foundation to align techniques with the characteristics of millisecond-scale vibration and acceleration data, particularly where no direct classification threshold or human annotations were available.

- B) All scripts, Jupyter notebooks, data processing tools, and modular components used in this project are openly available at:

GitHub: <https://github.com/Ashmit-workplace/pump-health-prediction>

This repository contains:

- Code scripts for each phase of the pipeline
 - Cleaned and intermediate datasets
 - Feature extraction modules (FFT, time-series, contextual)
 - ML model training and evaluation scripts
 - Final scoring and label generation logic
- C) The final project presentation, summarizing the entire pipeline, findings, and outcomes, is included as a separate document. It can be accessed at:

Project Presentation:

<https://docs.google.com/presentation/d/1WYKJUXDY7tvkItSCzeK9mcmXfXDcd9zY/edit?usp=sharing&ouid=106496139649192924726&rtpof=true&sd=true>

Internship Completion Certificate:

https://drive.google.com/file/d/1HjgH7_A6RGn-t2vvm2Lu0W1zji-HA4sn/view?usp=drive_link

Thank you.