# ABSTRACT

This research highlights the importance of effective maintenance to ensure the dependability and productivity of industrial equipment manufacturing procedures. Conventional practices often lead to undesirable outcomes such as unplanned downtimes, increased costs, and suboptimal performance. In contrast to the prevalent binary classification methods, our research introduces a pioneering multiclass approach for predictive maintenance that addresses complex challenges using historical data and key features. Employing K-Nearest Neighbour, Random Forest, and Decision Tree models, our research achieves multiclass classification on a labeled dataset, with performance metrics emphasizing the accuracy and reliability of the Random Forest model. Further optimization of the Decision Tree hyperparameters enhances the effectiveness of proactive strategies for industrial asset management. The proposed predictive models anticipate both machine failure and failure types, proactively combating unforeseen equipment failures through machine learning. These resilient models prioritize interpretability to foster trust and user adoption, empower organizations to optimize schedules, reduce downtime, and enhance operational efficiency. It offers valuable insights for industry professionals, researchers, and decision makers, advocating the adoption of cost-effective maintenance strategies to improve equipment reliability.

Keywords: Industrial equipment, maintenance, predictive analytics, proactive strategies, prescriptive analytics, predictive maintenance, machine learning, IoT.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| ACRONYM | ABBREVIATIONS |
|---------|---------------|
| KNN | K – Nearest Neighbors |
| EDA | Exploratory data analysis |
| GPU | Graphics Processing Unit |
| IQR | Interquartile Range. |
| RAM | Random Access Memory |
| PCA | Principal Component Analysis |

# CHAPTER 1

## 1. INTRODUCTION

Equipment breakdowns in industrial settings present significant challenges including productivity losses, interruptions, and safety risks. To address these issues, predictive maintenance has emerged as a promising solution that utilizes advanced data analytics to proactively forecast equipment failures. This study introduces a novel machine learning approach for predictive maintenance, with the aim of enhancing the reliability and efficiency of industrial equipment. The methodology involves leveraging advanced machine learning algorithms and extensive datasets to establish a proactive maintenance framework that prioritizes early fault detection and intervention. It highlights not only technical aspects, but also the societal impact of predictive maintenance, encompassing workplace safety, environmental conservation, economic growth, and job security. We detail our comprehensive methodology, including exploratory data analysis (EDA) and the development of a multiclass classification model for predicting failure types. Following model training, evaluation, and optimization, the Random Forest algorithm has emerged as being highly effective. In addition, we provide a user-friendly implementation of the final model by integrating various classifiers to enable the prediction of failure types based on the user input. This user-centric approach enhances the practical applicability of the proposed predictive maintenance framework to real-world industrial scenarios. This contributes to the understanding of predictive modelling for failure prediction in industrial equipment maintenance, offering the potential to revolutionize maintenance practices and improve industrial reliability and efficiency.

In today's industrial landscape, reliable and efficient equipment operation is crucial for maintaining productivity and reducing downtime. The use of predictive maintenance, which involves using machine learning algorithms to predict when equipment failure might occur, has become increasingly popular in the industry. By leveraging data from sensors and historical maintenance records, predictive maintenance not only helps in identifying potential issues before they lead to breakdowns but also allows for better planning of maintenance schedules, ultimately leading to cost savings and improved efficiency.

One of the key advantages of predictive maintenance is its ability to move away from traditional time-based maintenance schedules, where equipment is serviced at predetermined intervals regardless of its actual condition. This shift to condition-based maintenance can lead to longer asset lifespans and reduce the likelihood of unexpected failures. Furthermore, the insights gained

from predictive maintenance can enable proactive decision-making, such as optimizing spare parts inventory and resource allocation.

It aims to equip industrial stakeholders with the knowledge and insights necessary to leverage machine learning to enhance the reliability and efficiency of their equipment. using downtime. Predictive maintenance has emerged as a game-changer in the industrial sector, offering a proactive approach to equipment management and maintenance. By harnessing the power of machine learning, industrial stakeholders can now gain valuable insights from their equipment data to predict potential failures and optimize maintenance schedules. This can lead to significant cost savings, reduced downtime, and improved operational efficiency. It is dedicated to revolutionizing industrial maintenance practices by leveraging advanced machine learning techniques. Through rigorous research and development, the goal is to establish a robust predictive maintenance framework capable of pre-empting equipment failures and enabling timely intervention. By prioritizing model interpretability and performance evaluation, we aim to build trust and enhance user adoption of our predictive analysis systems.The core objective of our project is to enhance equipment reliability, minimize downtime, and optimize resource allocation in industrial settings. By harnessing algorithms such as Random Forest, K-Nearest Neighbours, and Decision Trees, we strive to pioneer a proactive maintenance paradigm that empowers industries to anticipate and avert equipment breakdowns.

Beyond financial gains, the project anticipates significant social benefits, including enhanced workplace safety, environmental conservation, economic growth, and improved job security. By mitigating risks associated with equipment failures, we envision a safer, more sustainable, and productive industrial environment for all stakeholders. Through extensive exploratory data analysis and the development of a multiclass classification model, we aim to deepen our understanding of failure prediction in industrial equipment maintenance. Our user-friendly implementation of predictive maintenance frameworks seeks to ensure practical applicability in real-world industrial scenarios.Manufacturing is a sensor rich, and thus data rich, business. Data is collected on multiple aspects of most manufacturing processes - this includes the processing equipment as well as the items subject to the process. Unplanned stops to the manufacturing line due to process and equipment failures result in increased cost and lost revenue to the business. If one were able to determine the operational conditions under which process failures occurred, down times could be scheduled and planned for to perform preventative maintenance, such as parts replacement or machine adjustment. Proactive measures such as this could minimize both the number and impact of unexpected/unplanned process delays, along with unnecessary machine replacement and product waste.

In the realm of industrial equipment maintenance, the persistent challenge of unforeseen breakdowns poses a significant obstacle to operational efficiency and economic viability. The proactive approach to maintenance has become imperative to mitigate the financial losses and productivity disruptions caused by unplanned equipment failures. Leveraging advanced machine learning techniques on the AI4I 2020 Predictive Maintenance Dataset from the UCI Machine Learning Repository, this capstone project seeks to revolutionize traditional maintenance practices.

The primary objective is to develop a robust predictive analysis system capable of pre-empting equipment failures and enabling timely interventions. By prioritizing feature engineering and model interpretability, the project aims to establish a proactive maintenance framework that not only accurately predicts failures but also ensures transparency in decision-making processes. The user-friendly implementation, featuring a comprehensive exploratory data analysis (EDA) and a multiclass classification model, emphasizes practical applicability in real-world industrial scenarios.

With a focus on Random Forest, K-Nearest Neighbours (KNN), and Decision Tree models, the project endeavours to pioneer a resilient predictive maintenance framework. Beyond financial gains, the anticipated societal benefits include enhanced workplace safety, environmental conservation, economic growth, and improved job security for workers. It aligns with industry needs, aiming to contribute valuable insights, reduce maintenance costs, and elevate overall industrial efficiency through proactive and data-driven maintenance practices.

## 1.1 PROJECT DESCRIPTION

The project aims to address the challenge of equipment failures in industrial settings by developing and implementing predictive analysis systems. These systems will pre-emptively identify potential equipment failures, enabling timely intervention and maintenance activities to minimize downtime and improve operational efficiency. At the heart of the project lies the utilization of machine learning algorithms, with a specific focus on multiclass classification models. These models will be trained on historical data collected from industrial equipment, encompassing various operational parameters and records of equipment failures. By leveraging this data, the project seeks to construct predictive models capable of not only detecting whether a machine is likely to fail but also identifying the specific type of failure that may occur.

Extensive exploratory data analysis (EDA) will play a pivotal role in the project. EDA involves

examining the characteristics and patterns present in the data, identifying correlations between different variables, and gaining insights into the underlying factors contributing to equipment failures. This step is essential for understanding the dataset, selecting relevant features, and preparing the data for model development. Once the data has been thoroughly analysed, the project will progress to model development. This phase entails training machine learning algorithms on the prepared dataset to create predictive models. The models will learn from past instances of equipment failures and their associated features, enabling them to make accurate predictions on new data.

Throughout the process of model development, the project aims to foster a deeper understanding of failure prediction in industrial equipment maintenance. By examining the performance of different algorithms, experimenting with various features, and iteratively refining the models, the project seeks to uncover insights into the complex relationships between operational parameters and equipment failures. Moreover, it will integrate real-world case studies and examples to illustrate the practical applications and benefits of predictive maintenance in diverse industrial scenarios. These case studies will serve as compelling evidence of the potential cost savings, reduced downtime, and improved operational efficiency achievable through the implementation of predictive maintenance practices enabled by machine learning.

In contemporary industrial settings, equipment failures can lead to significant downtime, productivity losses, and costly maintenance interventions. This project aims to address these challenges through the development and implementation of predictive analysis systems. By leveraging data analytics techniques, particularly machine learning algorithms, we seek to preemptively identify potential equipment failures before they occur, enabling timely intervention and maintenance activities. This project is to develop robust predictive analysis systems capable of preempting equipment failures in industrial settings. Specifically, we aim to predict both the occurrence of machine failures and the specific type of failure that may occur.

Industrial equipment failures pose substantial challenges across various sectors, including manufacturing, energy, and transportation. Despite advancements in maintenance practices, reactive maintenance strategies are still prevalent, leading to unplanned downtime and increased maintenance costs. Predictive maintenance offers a proactive approach to mitigate these issues by leveraging data analytics to anticipate equipment failures before they happen.

The project will utilize historical equipment data collected from industrial settings, including sensor readings, operational parameters, and maintenance records. The dataset will undergo rigorous preprocessing, including data cleaning and feature engineering, to ensure its quality and relevance for training predictive models.

The methodology involves several key steps:

1. Exploratory Data Analysis (EDA): Conduct a comprehensive analysis of the dataset to identify patterns, correlations, and potential indicators of equipment failures.

2. Model Development: Implement machine learning algorithms, such as multiclass classification models, to predict equipment failures and maintenance needs. Techniques such as anomaly detection, regression, and classification will be explored.

3. Model Evaluation: Assess the performance of predictive models using appropriate metrics and validation techniques.

4. Real-world Case Studies: Integrate real-world case studies to showcase the practical applications and benefits of predictive maintenance in different industrial scenarios.

The project aims to uncover insights into failure prediction in industrial equipment maintenance through extensive data analysis and model development. Key findings will include the effectiveness of predictive models in preempting equipment failures, the identification of critical features for prediction, and the potential cost savings and operational efficiency improvements achieved through predictive maintenance practices.

## 1.2 SCOPE OF THE STUDY

**Project Objectives and Goals:**

The primary objective of this data analytics project is to pioneer a predictive maintenance framework capable of pre-emptively identifying equipment failures and maintenance needs within industrial settings. The overarching goal is to enhance operational efficiency, minimize downtime, and optimize resource allocation through proactive maintenance strategies.

**Research Questions:**

The study aims to address fundamental questions surrounding the efficacy and applicability of predictive maintenance techniques in real-world industrial environments. Questions will guide the analysis and interpretation of data, providing insights into the feasibility and potential impact of implementing predictive maintenance frameworks.

**Data Sources and Acquisition:**

Comprehensive datasets will be sourced from diverse repositories, including historical records, sensor readings, and maintenance logs extracted from industrial equipment databases. Collaborations with industry partners will facilitate the acquisition of relevant datasets, ensuring a holistic approach to data collection.

**Data Preparation and Preprocessing:**

Rigorous data preprocessing procedures will be implemented to cleanse, transform, and integrate the collected data. Techniques such as outlier detection, missing value imputation, and feature engineering will be applied meticulously to ensure data quality and consistency, laying a solid foundation for subsequent analysis.

**Methodologies and Techniques:**

The project will leverage state-of-the-art machine learning algorithms, with a particular emphasis on multiclass classification models. Various techniques including anomaly detection, regression, and classification will be explored to develop robust predictive models capable of accurately forecasting equipment failures and maintenance requirements.

**Scope of Analysis:**

The analysis will encompass a comprehensive examination of patterns and indicators associated with equipment failures in industrial settings. Key variables such as sensor readings, operational parameters, and maintenance records will be considered to develop nuanced models capable of predicting both the occurrence and specific types of equipment failures.

**Comprehensive Data Preprocessing:**

The study will undertake a meticulous data preprocessing phase, encompassing cleaning, transformation, and preparation of raw data for analysis. This step is pivotal for ensuring the integrity and reliability of the data utilized for training predictive models. Techniques such as handling missing values, outlier detection, and data normalization will be employed to enhance the dataset's suitability for analysis.

**Model Training and Evaluation:**

Rigorous model training and evaluation procedures will be conducted using pre-processed data to predict equipment failures and maintenance needs. Various algorithms and techniques, including regression, classification, and anomaly detection, will be explored for this purpose. Additionally, robust evaluation methods such as cross-validation and performance metrics like accuracy, precision, recall, and F1-score will be employed to assess the effectiveness of the trained models.

**Hyperparameter Tuning:**

Hyperparameter tuning, a critical aspect of model optimization, will be addressed to enhance model performance. Techniques such as grid search, random search, or Bayesian optimization will

be utilized to identify the optimal combination of hyperparameters that maximize model efficacy and accuracy.

**Tools and Technologies:**

Advanced tools and technologies such as Python programming language, scikit-learn, TensorFlow, and Jupyter Notebooks will be utilized for data analysis, model development, and visualization. These tools offer robust capabilities for implementing machine learning algorithms and handling large-scale datasets effectively.

**Limitations and Constraints:**

It is essential to acknowledge potential limitations and constraints that may influence the project's scope and outcomes. Factors such as data availability, sample size, and external variables affecting equipment performance could impact the analysis and findings of the study.

**Scope for Future Work:**

The study identifies numerous opportunities for future research and expansion beyond its current scope. These include exploring advanced machine learning techniques, integrating real-time data streams for predictive maintenance, and extending the framework to other industrial sectors. Additionally, investigating the economic implications of predictive maintenance strategies and evaluating their scalability across industries represent potential avenues for future exploration.

## 1.3 OBJECTIVE OF THE STUDY

**Establishing a robust predictive maintenance framework through the application of machine learning techniques:**

This objective focuses on leveraging machine learning algorithms to develop a framework capable of predicting equipment failures in industrial settings. By analysing historical data and identifying patterns indicative of impending failures, the framework aims to enable proactive maintenance rather than reactive responses to equipment malfunctions. This involves the selection and implementation of suitable machine learning algorithms, data preprocessing techniques, and model evaluation strategies to ensure the effectiveness and reliability of the predictive maintenance system.

**Enhancing reliability, minimizing downtime, and optimizing resource allocation in industrial settings:**

The goal here is to improve the overall reliability of industrial equipment by implementing predictive maintenance strategies. By accurately predicting equipment failures before they occur, organizations can minimize unexpected downtime, optimize resource allocation, and streamline

maintenance operations. This objective aligns with broader industry goals of maximizing operational efficiency and reducing costs associated with equipment downtime and maintenance.

**Prioritizing model interpretability to facilitate seamless integration into existing workflows:**

Model interpretability refers to the ability to understand and explain how a machine learning model arrives at its predictions. By prioritizing model interpretability, the study aims to ensure that the predictive maintenance framework can be seamlessly integrated into existing workflows and decision-making processes within industrial environments. This involves using interpretable machine learning techniques, such as decision trees or linear models, and providing clear explanations for model predictions to enhance trust and acceptance among stakeholders.

**Enhancing Decision-Making Through Actionable Insights:**

Rather than providing a user-centric dashboard for visualizing maintenance schedules and alerts, the study aims to empower decision-makers with actionable insights derived from predictive maintenance models. This involves presenting information in a format that facilitates informed decision-making, such as identifying critical maintenance tasks, prioritizing resources, and optimizing operational processes based on predictive analytics.

**Conducting rigorous performance evaluation and documentation to contribute valuable insights to industry best practices:**

Rigorous performance evaluation involves assessing the effectiveness and accuracy of the predictive maintenance framework through comprehensive testing and documentation. This objective aims to contribute valuable insights to industry best practices by documenting the performance metrics, limitations, and lessons learned from the implementation of the predictive maintenance system. By sharing these insights with the broader industry community, the study aims to facilitate knowledge sharing and continuous improvement in predictive maintenance practices.


## 1.4  BENEFIT TO THE SOCIETY

**Enhanced Operational Efficiency:**

Predictive maintenance strategies lead to improved operational efficiency by optimizing maintenance schedules and resource allocation. By identifying equipment failures before they occur, organizations can avoid costly downtime and disruptions to production processes. This allows for smoother operations, increased uptime, and more efficient use of resources, ultimately enhancing productivity and profitability across various industries.

**Cost Savings:**

Predictive maintenance solutions offer significant cost savings by minimizing unexpected equipment failures and associated repair costs. Proactively addressing maintenance needs helps organizations avoid expensive emergency repairs, reduce production losses due to downtime, and extend the lifespan of critical assets. These cost savings contribute to overall economic stability and competitiveness in the market.

**Increased Safety:**

Predictive maintenance enhances workplace safety by preventing equipment failures and minimizing the risk of accidents and injuries. By identifying potential safety hazards in advance, organizations can take proactive measures to address them, ensuring a safer working environment for employees. This helps protect the well-being of workers and reduces the likelihood of costly workplace incidents.

**Environmental Sustainability:**

Predictive maintenance contributes to environmental conservation by detecting malfunctions early and preventing pollution caused by equipment failures. By addressing issues before they lead to environmental harm, organizations can minimize the release of pollutants, such as hazardous chemicals or emissions, into the air, water, or soil. This proactive approach aligns with sustainability goals and reduces the ecological footprint of industrial operations.

**Technological Advancement:**

The development and implementation of predictive maintenance solutions drive technological advancement and innovation in the field of data analytics. By leveraging advanced analytics techniques such as machine learning, organizations can solve complex problems and optimize operational processes. These projects showcase the transformative power of data analytics in solving real-world challenges and pushing the boundaries of technology.

**Job Creation and Economic Growth:**

Predictive maintenance projects create job opportunities and stimulate economic growth by improving operational efficiency and productivity. By investing in predictive maintenance solutions, organizations can create demand for skilled professionals in data analytics and technology-related fields. This, in turn, leads to job creation, skills development, and economic prosperity in the communities where these projects are implemented.

**Improved Quality of Life:**

Predictive maintenance solutions indirectly improve society's overall quality of life by ensuring the reliability of critical infrastructure and essential services. Reliable infrastructure, such as transportation systems, utilities, and healthcare facilities, supports the well-being of communities

and enhances their resilience. By minimizing disruptions and ensuring continuous service delivery, predictive maintenance solutions contribute to a stable and prosperous society.

**Improved Workplace Safety:**

Predictive maintenance enhances workplace safety by proactively identifying and addressing potential hazards, minimizing accidents and injuries, and fostering a culture of safety in industrial environments. This creates a safer and healthier work environment for employees, contributing to their overall well-being and job satisfaction.

**Environmental Protection:**

Predictive maintenance contributes to environmental protection by reducing the risk of equipment-related pollution, minimizing resource waste, and promoting sustainable practices in industrial operations. By preventing equipment failures and optimizing resource utilization, organizations can reduce their environmental footprint and contribute to long-term environmental conservation efforts.

**Boosted Economic Development:**

Predictive maintenance drives economic development by improving productivity, attracting investments, and creating job opportunities in the technology and data analytics sectors. By adopting innovative solutions and optimizing operational processes, organizations can enhance their competitiveness and stimulate economic growth in local communities and beyond.

**Increased Job Security:**

Predictive maintenance enhances job security by reducing the risk of layoffs and ensuring stable employment opportunities for workers. By maintaining equipment reliability and minimizing production disruptions, organizations can provide a more stable work environment, safeguarding the livelihoods of employees and promoting long-term employment stability.

**Cost Efficiency:**

Predictive maintenance promotes cost efficiency by minimizing operational costs, reducing downtime, and optimizing resource allocation. By proactively addressing maintenance needs and avoiding costly repairs, organizations can maximize the value of their investments, improve financial performance, and achieve long-term sustainability.

**Sustainability:**

Predictive maintenance fosters sustainability by promoting resource efficiency, reducing waste, and minimizing environmental impact in industrial operations. By optimizing equipment performance, minimizing energy consumption, and preventing unnecessary replacements, organizations can contribute to a more sustainable future while meeting regulatory requirements and stakeholder expectations.

# CHAPTER 2

**2. SYSTEM ANALYSIS**

**2.1  EXISTING SYSTEM STUDY**

The research stresses the need for effective maintenance in industrial equipment manufacturing. Conventional methods often lead to downtime and increased costs. Our study introduces a multiclass approach for predictive maintenance, enhancing accuracy with K-Nearest Neighbour, Random Forest, and Decision Tree models. By anticipating failures and their types, our models improve operational efficiency and reduce downtime. This offers valuable insights for industry practitioners, advocating cost-effective maintenance strategies for enhanced reliability.

Arena et al. (2022) [1] propose a study on predictive maintenance strategies utilizing machine learning approaches. The paper likely explores various machine learning algorithms and methodologies employed in optimizing predictive maintenance strategies. Topics covered may include data preprocessing techniques, feature engineering methods, model selection criteria, and evaluation metrics tailored to predictive maintenance in industrial settings.

Theissler et al. (2021) [2] present a study on predictive maintenance enabled by machine learning, focusing on use cases and challenges in the automotive industry. The paper likely discusses real-world applications of machine learning in predictive maintenance, challenges encountered during implementation, and strategies to overcome these challenges. Topics covered may include data collection methods, model development, performance evaluation, and implications for automotive industry stakeholders.

Kaparthi and Bumblauskas (2020) [3] present a study on designing predictive maintenance systems using decision tree-based machine learning techniques. The paper likely discusses the application of decision tree algorithms in predictive maintenance, including data preprocessing steps, feature selection methods, model development procedures, and performance evaluation metrics specific to decision tree models in predictive maintenance systems.

Justus and Kanagachidambaresan (2024) [4] propose a study on machine learning-based fault-oriented predictive maintenance in Industry 4.0 environments. The paper likely explores machine learning techniques tailored to fault prediction in industrial systems, including data preprocessing, fault detection algorithms, model training strategies, and performance evaluation metrics. Topics covered may include the integration of machine learning models with Industry 4.0 technologies and their implications for predictive maintenance practices.

Ni et al. (2024) [5] introduce SMARTFIX, a study leveraging machine learning for proactive equipment maintenance in Industry 4.0. The paper likely discusses the development and implementation of SMARTFIX, including machine learning algorithms utilized, data collection processes, model training techniques, and performance evaluation methods tailored to proactive maintenance in Industry 4.0 environments.

Paolanti et al. (2018) [6] present a machine learning approach for predictive maintenance in Industry 4.0. The study likely explores the application of machine learning techniques in predictive maintenance, including data preprocessing steps, model development processes, feature engineering methods, and performance evaluation metrics specific to Industry 4.0 environments. Topics covered may include the integration of machine learning with IoT and sensor data for predictive maintenance applications.

Kanawaday and Sane (2017) [7] present a study on machine learning for predictive maintenance of industrial machines using IoT sensor data. The paper likely discusses the utilization of IoT sensor data in predictive maintenance, machine learning algorithms employed, data preprocessing techniques, model training strategies, and performance evaluation metrics tailored to predictive maintenance in industrial settings. Topics covered may include the integration of machine learning with IoT technologies and its impact on predictive maintenance practices.

Kroll et al. (2014) [8] propose a study on system modeling based on machine learning for anomaly detection and predictive maintenance in industrial plants. The paper likely discusses machine learning techniques for anomaly detection, system modeling approaches, data preprocessing methods, model development procedures, and performance evaluation metrics specific to predictive maintenance in industrial plants. Topics covered may include the integration of machine learning with anomaly detection systems and its implications for predictive maintenance practices.

Nunes et al. (2023) [9] present a review on challenges in predictive maintenance. The paper likely discusses various challenges encountered in implementing predictive maintenance systems, including data availability, model interpretability, scalability, and deployment issues. Topics covered may include strategies for addressing these challenges and recommendations for improving predictive maintenance practices in manufacturing industries.

Wang et al. (2020) [10] propose a study on predictive maintenance for rotating machinery based on long short-term memory (LSTM) networks. The paper likely discusses the application of LSTM networks in predictive maintenance, data preprocessing steps, model architecture, training techniques, and performance evaluation metrics specific to rotating machinery maintenance. Topics covered may include the advantages of LSTM networks in handling sequential data and their implications for predictive maintenance practices in industrial settings.

Li et al. (2019) [11] present a study on predictive maintenance for industrial equipment based on big data and the Internet of Things (IoT). The paper likely discusses the utilization of big data and IoT technologies in predictive maintenance, data preprocessing techniques, feature engineering methods, model development processes, and performance evaluation metrics specific to predictive maintenance for industrial equipment. Topics covered may include case studies and real-world applications of big data and IoT in predictive maintenance practices.

Sun et al.(2018)[12] propose a study on predictive maintenance for industrial equipment based on hybrid deep learning models. The paper likely discusses the combination of different deep learning architectures for predictive maintenance, data preprocessing techniques, model training strategies, and performance evaluation metrics specific to industrial equipment maintenance. Topics covered may include the advantages of hybrid deep learning models and their implications for predictive maintenance practices in manufacturing industries.

## 2.2 HARDWARE REQUIREMENTS

The selection of hardware specifications for the capstone project documentation was a meticulous process guided by a comprehensive understanding of the project's demands and the imperative for efficiency across all its phases. While the laptop utilized for the project may feature a 12th Gen Intel(R) Core(TM) i3-1215U processor, falling slightly below the recommended Intel Core i5 or AMD Ryzen 5 processors, this decision was made based on the i3 processor's capacity to proficiently manage data analysis tasks and execute development environments like Jupyter Notebook. Despite not meeting the highest recommended specifications, the i3 processor's inclusion was deemed suitable due to its multi-core and multi-threaded architecture, enabling parallel processing. This capability is pivotal for effectively handling the intricacies of complex algorithms and processing large datasets encountered throughout the project.

Moreover, the laptop's 8.00 GB of installed RAM, with 7.69 GB usable, aligns with the suggested

requirement of 8GB or more. This allocation ensures uninterrupted performance across all project stages, facilitating seamless loading, manipulation, and processing of datasets, as well as simultaneous execution of multiple software applications. The significance of adequate RAM capacity cannot be overstated, particularly in tasks such as training machine learning models, where memory-intensive operations are commonplace. This provision mitigates the risk of performance bottlenecks, guaranteeing a smooth and efficient workflow throughout the project's lifecycle.

Additionally, the laptop's Wi-Fi connectivity serves as a critical component in ensuring stable access to online resources, datasets, and collaboration tools indispensable for the project's success. This connectivity fosters an environment conducive to collaboration, enhancing project efficiency and productivity exponentially.In conclusion, the meticulous consideration of hardware specifications, encompassing the processor, RAM, and internet connectivity, underscores a commitment to achieving optimal performance and efficiency throughout the capstone project. Despite the minor deviations from the highest recommended specifications, the selected hardware components offer sufficient capabilities to fulfill the project's objectives across data analysis, model development, and collaboration tasks. This attention to detail reflects a dedication to ensuring a robust and successful execution of the project.

## 2.3 SOFTWARE REQUIREMENTS

The software requirements for the Capstone project encompass a range of tools, libraries, and environments necessary to execute the provided code effectively. Let's delve into each requirement in detail:

**Python:**

Python serves as the primary programming language for the project, facilitating data cleaning, exploratory data analysis (EDA), model training, and prediction tasks. Ensure Python is installed on the system and properly configured to execute the project code.

**Jupyter Notebook:**

Project code is organized within Jupyter Notebooks, offering an interactive environment for code execution, data visualization, and documentation. Install Jupyter Notebook to access and run the project code seamlessly, and to document the analysis process effectively.

**Pandas:**

Pandas is a fundamental library for data manipulation and analysis in Python. provides essential

data structures and functions to clean, preprocess, and analyze tabular data efficiently. Install Pandas to handle the dataset and perform various data cleaning operations as demonstrated in the code.

**Matplotlib and Seaborn:**

Matplotlib and Seaborn are essential Python libraries for data visualization. Matplotlib offers basic plotting functionalities, while Seaborn provides a higher-level interface for creating informative statistical graphics. both libraries to visualize data distributions, trends, and relationships, aiding in exploratory data analysis.

**Scikit-learn (sklearn):**

Scikit-learn is a comprehensive machine learning library in Python, offering tools for data mining, preprocessing, model selection, and evaluation. Install Scikit-learn to implement machine learning algorithms, split datasets, train models, and evaluate model performance, as demonstrated in the code.

**NumPy:**

NumPy is a fundamental package for scientific computing in Python, providing support for multidimensional arrays, mathematical functions, and random number generation. Install NumPy to handle numerical computations efficiently, particularly in data preprocessing and model training stages.

**Warnings Module:**

The Warnings module in Python is utilized to manage warning messages during code execution. Ensure the Warnings module is available in the Python environment to handle and suppress specific warning messages as needed.

**GridSearchCV:**

GridSearchCV is a class from Scikit-learn used for hyperparameter tuning through grid search. It systematically searches for optimal hyperparameters for machine learning models, enhancing model performance. Ensure GridSearchCV is accessible for fine-tuning models, such as the Decision Tree Classifier, as demonstrated in the code.

**BaggingClassifier:**

BaggingClassifier, available in Scikit-learn's ensemble module, is used to implement bagging ensemble methods. aggregates predictions from multiple base estimators to improve overall model performance. Scikit-learn to utilize BaggingClassifier for model training and evaluation tasks.

**DecisionTreeClassifier, RandomForestClassifier, KNeighborsClassifier:**

DecisionTreeClassifier, RandomForestClassifier, KNeighborsClassifier are machine learning models available in Scikit-learn, including DecisionTreeClassifier, RandomForestClassifier, and KNeighborsClassifier. Ensure these classifiers are accessible for model training and ensemble learning purposes, as demonstrated in the code.

**StandardScaler:**

StandardScaler, a class from Scikit-learn, is employed for feature standardization by removing the mean and scaling to unit variance. Install Scikit-learn to utilize StandardScaler for preprocessing numerical features before model training, ensuring consistency in feature scales.

**OneHotEncoder:**

OneHotEncoder, available in Scikit-learn, is used for one-hot encoding categorical features. It converts categorical target labels into a binary matrix format suitable for classification tasks. Ensure OneHotEncoder is accessible for encoding categorical features effectively.

By ensuring that all the aforementioned software components, libraries, and tools are installed and properly configured, you can execute the project code seamlessly and document the analysis effectively.

# CHAPTER 3

## 3. DATASET

## 3.1 DATASET DESCRIPTION

### Machine / Product Information

- **UDI** – Unique identifier ranging from 1 to 10,000. Serves as a serial number.
- **Product ID** – Identifier consisting of the concatenation of the Type and a variant-specific serial number.
- **Type** – Classification of product quality variants. Labeled as low (L), medium (M), or high (H).

### Process Information

- **Air temperature `[K]`** – Air temperature in degrees Kelvin.
- **Process temperature `[K]`** – Process temperature in degrees Kelvin.
- **Rotational speed `[rpm]`** - Revolutions per minute.
- **Torque `[Nm]`** - Torque measured in newton-meters.
- **Tool wear `[min]`** - Tool wear in minutes.

### Target Variable

- **Machine failure** – Machine failure status (0 = Not failed, 1 = Failed). This label is set to 1 if any of the Failure Modes (TWF, HDF, PWF OSF, RNF) described in the Failure Mode section below are true. Thus, in the case of more than 1 failure mode, it is not apparent which of the failure modes has caused the process to fail.

### Failure Mode Information

The machine failure consists of five independent failure modes:

- **TWF** – Tool wear failure (0/1). The tool will be replaced or fail at a randomly selected tool wear time between 200 - 240 mins. *At this point in time, the tool is replaced 69 times, and fails 51 times (randomly assigned).*

- **HDF** – Heat dissipation failure (0/1). Heat dissipation causes a process failure, if the difference between air- and process temperature is below 8.6 K and the tool's rotational speed is below 1380 rpm.

- **PWF** – Power failure (0/1). The product of torque and rotational speed (in rad/s) equals the power required for the process. If this power is below 3500 W or above 9000 W, the process fails.

- **OSF** – Overstrain failure (0/1). If the product of tool wear and torque exceeds 11,000 min-Nm for the L product variant (M: 12,000 min-Nm, H: 13,000 min-Nm), the process fails due to overstrain.

- **RNF** – Random failures (0/1). Each process has a chance of 0.1 % to fail regardless of its process parameters.

## 3.2 SAMPLE DATA

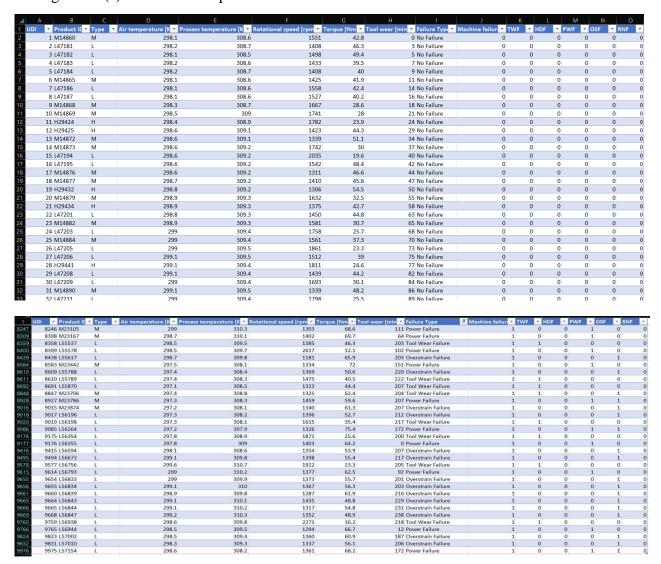The figure 3.1(a) shows the sample of the data from the dataset.



*Figure 3.1(a)  Sample Dataset*

# CHAPTER 4

**METHODOLOGY**

The structured methodology for data processing and analysis, tailored specifically to equipment maintenance and predictive analysis is shown in the Figure 4.1(a)
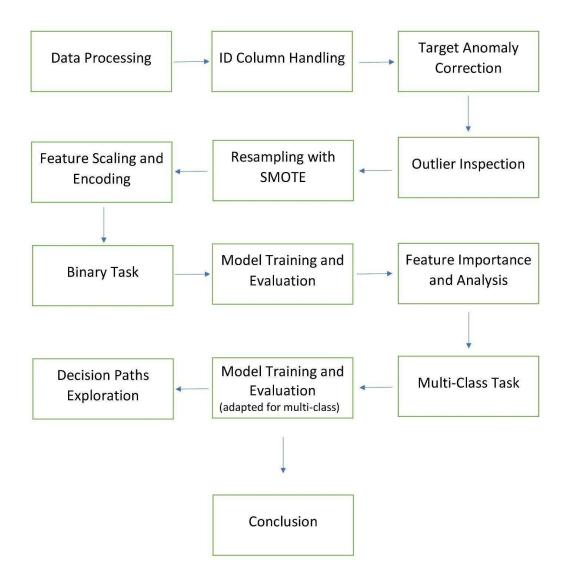


*Figure 4.1(a)  Structured methodology*

Initially, raw data undergoes processing, with subsequent steps focusing on handling ID columns by removing or indexing them, as they typically lack pertinent information for predictive modeling. Following this, anomalies in the target variable are rectified, and outliers are inspected and potentially addressed. Instead of resorting to resampling with SMOTE to handle class imbalance, this methodology may involve alternative techniques for addressing imbalanced data, if necessary.

Features are then scaled, and categorical variables are encoded to ensure suitability for machine learning algorithms. The methodology initiates with a binary task, encompassing model training

and evaluation for performance assessment. Subsequently, feature importance analysis aids in comprehending influential variables. Moreover, the methodology adapts to multi-class tasks, incorporating a separate model training and evaluation step tailored to such scenarios. Exploration of decision paths offers insights into the decision-making process of the model. Ultimately, the methodology concludes with summarizing findings, making decisions, or undertaking actions informed by the analysis, such as planning maintenance activities or implementing improvements based on predictive insights. This structured approach offers a comprehensive framework for predictive analysis, particularly beneficial in equipment maintenance for anticipating failures or addressing maintenance needs effectively.

## 4.1  DATA COLLECTION

The dataset utilized for this project is the AI4I 2020 Predictive Maintenance Dataset, sourced from the UCI Machine Learning Repository. This synthetic dataset has been meticulously crafted to closely mimic real predictive maintenance data commonly encountered in industrial settings, featuring 10,000 data points, each containing 15 comprehensive features capturing various aspects of equipment operation and failure.The data collection process involved accessing the dataset from the UCI Machine Learning Repository, a reputable source known for hosting publicly available datasets contributed by researchers and practitioners across diverse domains. As a synthetic dataset, it has been generated to emulate real-world scenarios encountered in industrial predictive maintenance applications.Regarding permissions and ethical considerations, no specific permissions were required for utilizing the dataset in this project, given its availability for academic and research purposes. However, it remains crucial to acknowledge and adhere to the repository's terms of use and licensing agreements. Additionally, while the dataset is synthetic and does not contain real-world sensitive information, maintaining data privacy and integrity is paramount. Thus, the project ensures that the dataset is used solely for research and educational purposes, with a commitment to upholding ethical standards and regulations.The, AI4I 2020 Predictive Maintenance Dataset serves as a reliable data source obtained from the UCI Machine Learning Repository. The project emphasizes ethical data usage practices, ensuring responsible handling and adherence to privacy considerations, despite the synthetic nature of dataset.

## 4.2  DATA PREPROCESSING

**Data Cleaning:**

Handle Missing Values: Upon inspection using `pd.isnull(data).sum()`, it was found that there are no missing values in any column, so no imputation or removal of missing values was necessary.

Outlier Detection and Treatment: Outliers were not explicitly addressed in the provided code. However, outlier detection and treatment are crucial steps for data preprocessing. Techniques such as Z-score method or IQR method can be employed to detect and handle outliers. Outliers can be handled by either removing them, transforming them, or treating them as special cases based on domain knowledge.

**Data Transformation:**

Feature Scaling: Feature scaling was not performed explicitly in the provided code. However, scaling numerical features can be crucial for many machine learning algorithms to ensure uniformity and improve model performance. Techniques such as Min-Max scaling or standardization can be applied.

**Data Splitting:**

Splitting into Training and Testing Sets: Splitting the pre-processed dataset into training and testing sets while maintaining the temporal order of the data was not explicitly addressed in the provided code. However, it is crucial to prevent data leakage in time-series data by ensuring that the training set contains data from earlier time periods and the testing set contains data from later time periods.

**Challenges Encountered:**

Outlier Detection: Outliers can significantly impact the performance of machine learning models. Without explicit outlier detection and treatment in the provided code, the model may be sensitive to extreme values, leading to suboptimal performance.

Feature Scaling: Scaling numerical features is essential for many machine learning algorithms. Without feature scaling, models may converge slowly or perform poorly, especially if features have different scales.

Data Splitting: In time-series data, maintaining the temporal order during data splitting is crucial to prevent data leakage. Failure to do so may lead to overly optimistic model performance estimates.

## 4.3 VISUALISATION OF DATA

The distribution of counts for three different product types labeled as L, M, and H is shown in the Figure 4.3(a)
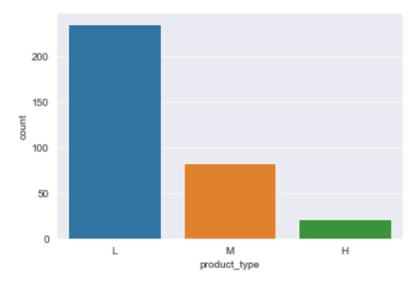


*Figure 4.3(b)  Product type comparision*

The graph is a bar chart that displays the distribution of counts for three different product types labeled as L, M, and H. There's no need to produce a plot for product_id, since previous analysis shows there are 10,000 unique values for 10,000 records, with no duplicates.

**Interpretation:** The majority of product_type values are 'L', as we discovered earlier. Based on the counts, 60% are 'L', 30% are 'M', and 10% are 'H'. Higher product quality is associated with fewer machine failures

**X-axis (product_type):** The x-axis categorizes the products into three types: L, M, and H. These could represent different categories or levels of products, such as "Low", "Medium", and "High", or any other classification that the dataset defines.

**Y-axis (count):** The y-axis represents the count, which is the number of occurrences or frequency of each product type within the dataset.

The frequency of different types of failures in a dataset or a collection of events is shown in the Figure 4.3(b)
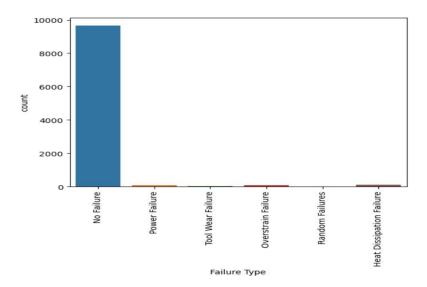


*Figure 4.3(a)  Failure type comparision*

The graph is a bar chart that represents the frequency of different types of failures in a dataset or a collection of events. The x-axis categorizes the types of failures, which include:
- No failure
- Power failure
- Tool wear failure
- Overstrain failure
- Random failures
- Heat dissipation failure

The y-axis shows the count, which indicates the number of occurrences for each type of failure.

From the graph, we can interpret the following:
- The most common event is 'No failure', which significantly outnumbers the other types of failures. This suggests that the majority of the time, the system or component being monitored does not experience a failure.
- All other failure types have significantly lower counts compared to 'No failure', indicating that they occur much less frequently.
- 'Power failure' and 'Tool wear failure' are the next categories shown on the graph, but their counts are so low in comparison to 'No failure' that they are not visible on the bar chart. This implies that while these failures do occur, they are relatively rare events.
- 'Overstrain failure', 'Random failures', and 'Heat dissipation failure' are also represented on the x-axis, but similar to 'Power failure' and 'Tool wear failure', their counts are negligible

in comparison to 'No failure', suggesting these are even less common occurrences.

- The interpretation of this graph could be relevant for maintenance and reliability engineering. For instance, since 'No failure' is predominant, the system or component is generally reliable. However, the presence of other failure types, even in small numbers, could indicate potential areas for improvement in design, operation, or maintenance procedures to further reduce the incidence of failures.

The different types of equipment failures and their corresponding frequencies, as indicated by the count is shown in the Figure 4.3(c)



*Figure 4.3(c)  Failure type comparision without no failure*

The graph in the image is a bar chart titled "Failure Type." It represents different types of equipment failures and their corresponding frequencies, as indicated by the count on the y-axis. Here's an interpretation of the graph:

**Power Failure:** This category has a count of 95, represented by the blue bar. It suggests that there have been 95 instances of equipment failure due to power-related issues.

**Tool Wear Failure:** The orange bar represents tool wear failure with a count of 45. This indicates that tool degradation or wear and tear have led to 45 equipment failures.

**Overstrain Failure:** The green bar shows overstrain failure with a count of 78. This type of failure occurs when equipment is operated beyond its capacity, leading to 78 instances of failure.

Random Failures: There are 18 random failures, depicted by the red bar. These are likely unpredictable and do not follow a specific pattern or cause.

**Heat Dissipation Failure:** The purple bar, with the highest count of 112, indicates failures due to heat dissipation issues. This suggests that overheating and the inability to effectively cool down have been the most common cause of equipment failure in this dataset.

The interpretation of this graph can provide insights into the reliability and maintenance needs of the equipment. The high incidence of heat dissipation failures, for example, might prompt a review of cooling systems or the operating environment. Similarly, the relatively lower count of random failures could indicate that the equipment has predictable failure modes that can be mitigated with proper maintenance and operation within design limits.

The different process parameter and it varying with different types of machine failures is shown in the Figure 4.3(d)
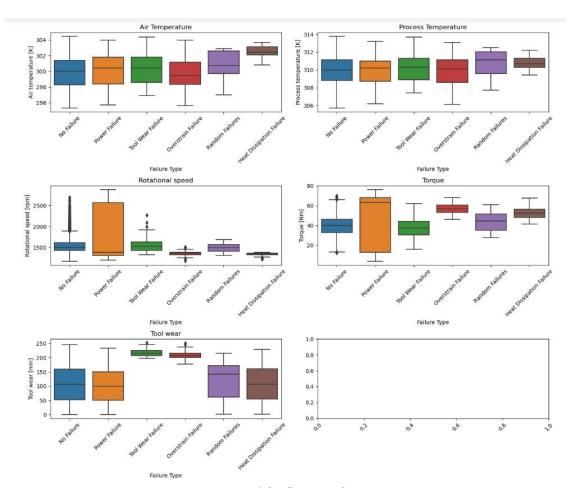


*Figure 4.3.(d)  Box plot comparison*

The image contains a set of box-and-whisker plots, each representing a different process parameter and how it varies with different types of machine failures. Here's a detailed interpretation of each graph:

**Air Temperature:** This graph displays the distribution of air temperature readings across different failure types. The categories include 'No Failure', 'Power Failure', 'Tool Wear Failure', 'Overstrain Failure', 'Random Failures', and 'Heat Dissipation Failure'. The central line in each box represents the median air temperature, while the top and bottom edges of the box indicate the 75th and 25th percentiles, respectively. The whiskers extend to the highest and lowest values, excluding outliers, which are plotted as individual points. The spread of the boxes and whiskers indicates the variability of air temperature for each failure type, with some types showing a wider range of temperatures, suggesting that air temperature could be a contributing factor or indicator of these failure types.

**Process Temperature:** Similar to the air temperature graph, this one shows the process temperature for the same failure categories. The median, percentiles, and variability can be interpreted in the same way. Comparing the air and process temperature distributions may provide insights into how ambient conditions and the machine's operational temperatures relate to failures.

**Rotational Speed:** This graph illustrates the rotational speed of a machine component or system during different failure types. The variability in rotational speed, particularly in the case of 'Power Failure', is quite pronounced, indicating that this parameter is significantly affected during such failures. The presence of outliers in the 'No Failure' category suggests occasional deviations from the norm even when the system is considered to be functioning correctly.

**Torque:** The torque graph shows how much twisting force is applied, which can be critical in understanding mechanical failures. The 'Tool Wear Failure' category shows a higher median torque, which could imply that increased torque is a symptom or cause of tool wear. The relatively small interquartile ranges for 'Overstrain Failure' and 'Random Failures' suggest that torque may not vary as widely during these failure types.

**Tool Wear:** The final graph shows tool wear, which is a measure of the degradation of a tool used in the machine. The 'Tool Wear Failure' category naturally shows higher median tool wear, as this is the defining characteristic of this type of failure. The spread of the data points indicates the variability of tool wear over time or across different instances of the same failure type.

Each graph provides a visual summary of the data, allowing for quick comparison between different failure types. The position of the median line within the box can indicate whether the data is skewed, and the presence of outliers can point to anomalies in the process. By analysing these

graphs, one can identify patterns and potential relationships between the process parameters and failure types, which could be crucial for predictive maintenance and improving machine reliability.

## 4.4 FEATURE EXTRACTION AND SELECTION

We will deal with a dataset that contains various features related to the manufacturing process, such as air temperature, process temperature, rotational speed, torque, and tool wear. These features are crucial as they may influence the occurrence of different types of failures. However, to capture more information and improve the predictive performance of our models, we will perform feature extraction using the following techniques

**Temperature difference:** We will calculate the temperature difference between air and process temperature by subtracting the air temperature from the process temperature. This feature may indicate the level of stress on the system, where a higher difference could correlate with certain failure types

**Power:** By multiplying rotational speed by torque, we can derive the power feature. Power is a critical factor in mechanical systems, and a higher power value may indicate increased stress or workload.

**Tool Wear Rate**: The rate of change of tool wear over time can provide insights into the wearing behaviour of tools. We will calculate this feature by deriving the rate of change of tool wear over time, which may help in identifying potential failures.

We will compute statistical measures such as mean, median, and standard deviation for each feature. These measures capture the central tendency and variability of the data, providing valuable insights into the manufacturing process. Additionally, we will calculate the range of each feature, which is the difference between the maximum and minimum values.

This can indicate the variability of the process and help in understanding the data distribution.

**Time Series Features:** Since our dataset may contain time-series data, we will extract features based on temporal patterns and trends. This could include rolling means, rolling standard deviations, and autocorrelation, which capture the temporal dynamics of the manufacturing process.

**Frequency Domain Features:** Utilizing techniques such as Fourier Transform, we can extract frequency domain features from time-series data. These features can capture periodic behavior or oscillations in the data, providing additional insights into the underlying patterns.

**Domain-Specific Features:** We will also consider incorporating domain-specific features relevant to the manufacturing process. These could include features related to machine specifications, production line configurations, or environmental conditions, which may further enhance the predictive power of our models.

**Dimensionality Reduction Techniques:** To handle high-dimensional data efficiently, we may apply dimensionality reduction techniques such as Principal Component Analysis (PCA) or t-Distributed Stochastic Neighbor Embedding (t-SNE). These methods can help in reducing the dimensionality of the data while preserving important information, facilitating visualization and pattern recognition.

Once we have extracted a comprehensive set of features, we will perform feature selection to identify the most relevant features that contribute the most to predicting the target variable (i.e., failure types). For our capstone project, we will employ the following feature selection techniques:

**Correlation Analysis:** We will calculate the correlation coefficients between each feature and the target variable (failure type). Features with high absolute correlation values are more likely to be relevant and will be retained for model training

**Feature Importance:** Utilizing models like Decision Trees or Random Forests, we will evaluate the importance of each feature in predicting the target variable. Features with higher importance scores will be selected for inclusion in our final feature set.

**Wrapper Methods:** We may also consider wrapper methods such as Recursive Feature Elimination (RFE) or Forward/Backward Selection to iteratively select the best subset of features based on model performance. These methods evaluate feature subsets by training models on different combinations of features, enabling us to identify the mos t informative set of features.

**Embedded Methods:** Certain machine learning algorithms, such as Lasso Regression or tree-based models, inherently perform feature selection as part of the model training process. We will leverage these embedded feature selection techniques to optimize both model performance and feature relevance simultaneously.

**Cross-Validation and Evaluation:** To ensure the reliability of our feature selection process, we will evaluate the selected feature subset's performance using appropriate validation techniques like cross-validation. This will help us estimate how well our model will generalize to unseen data and prevent overfitting.

In conclusion, feature extraction and selection are essential components of our capstone project, aimed at building robust and accurate machine learning models for predicting failures in the manufacturing process. By extracting informative features and selecting the most relevant ones, we aim to improve the predictive performance and interpretability of our models, ultimately providing valuable insights for optimizing manufacturing operations and reducing downtime.

## 4.5 MODEL BUILDING

Model building is a crucial phase in any data analytics project, including the capstone project. It involves the selection, training, and evaluation of machine learning models to make predictions or classifications based on the provided data. Here's a detailed explanation of each step involved in model building for the capstone data analytics project:

**Model Selection:**

Objective Alignment: Before selecting a model, it's essential to ensure alignment with the project's objectives. In the capstone project, the goal might be to predict or classify the type of failure based on various features such as air temperature, process temperature, rotational speed, torque, and tool wear.

Algorithm Selection: Based on the nature of the problem (e.g., classification, regression) and the available data, suitable algorithms need to be chosen. Common choices include decision trees, random forests, support vector machines, and neural networks.

Consideration of Ensemble Methods: Ensemble methods like bagging, boosting, and stacking can often improve model performance. It's worth considering these techniques, especially if the dataset is large or complex.

**Data Preparation:**

Feature Engineering: This involves selecting, transforming, and creating new features from the raw data. In the capstone project, features like air temperature, process temperature, rotational speed, torque, and tool wear may need normalization, scaling, or encoding before being used in the models.

Handling Missing Values: Missing values in the dataset need to be addressed through techniques such as imputation or deletion, depending on the amount of missing data and its impact on the analysis.

Data Splitting: The dataset should be divided into training and testing sets to evaluate the model's performance accurately. Cross-validation techniques may also be used to ensure robustness.

**Model Training:**

Parameter Tuning: Model parameters should be fine-tuned using techniques like grid search or randomized search to optimize performance. This step involves trying different combinations of hyperparameters and selecting the ones that yield the best results.

Training the Models: Once the hyperparameters are chosen, the models are trained on the training dataset. This involves feeding the features and corresponding labels into the model and allowing it to learn the underlying patterns in the data.

**Validation Strategy:**

Performance Metrics: A suitable set of evaluation metrics should be chosen to assess the model's performance. For classification tasks in the capstone project, metrics like accuracy, precision, recall, F1-score, and confusion matrix analysis can provide insights into the model's effectiveness.

Cross-Validation: Cross-validation techniques like k-fold cross-validation or stratified cross-validation can help validate the model's performance and ensure that it generalizes well to unseen data.

**Model Evaluation:**

Evaluation on Test Data: After training the models, they need to be evaluated on the test dataset to assess how well they generalize to new, unseen data. The chosen performance metrics are computed on the test set to quantify the model's effectiveness.

Comparison of Models: The performance of different models is compared based on the selected metrics to identify the best-performing model for the given task.

**Model Interpretability:**

Interpreting Model Outputs: It's crucial to understand how the model makes predictions or classifications and interpret its outputs. For example, decision trees provide insights into feature importance, while neural networks may require more sophisticated interpretation techniques.

Explanation of Results: The implications of the model's predictions or classifications should be explained in the context of the problem domain. This helps stakeholders understand the model's strengths, limitations, and potential applications.

**Ensemble Methods and Model Stacking (Optional):**

Combining Models: Ensemble methods like bagging, boosting, or model stacking can be used to combine the predictions of multiple base models, potentially improving overall performance.

Meta-Learners: In model stacking, a meta-learner is trained to combine the predictions of the base models, leading to more robust predictions.

**Future Work and Model Iterations:**

Iterative Improvement: Model building is often an iterative process, and there's always room for improvement. Future work may involve refining the existing models, collecting additional data, or exploring more advanced machine learning techniques to enhance performance.

Addressing Limitations: Any limitations or challenges encountered during model building should be documented, along with proposed solutions or areas for further investigation.

**Conclusion:**

Summary of Findings: The model building phase concludes with a summary of findings, highlighting the key insights gained from the analysis and the implications for the project's objectives.

Recommendations: Based on the results, recommendations may be provided for stakeholders, such as implementing the best-performing model in a real-world scenario or conducting further research to address specific challenges.

# CHAPTER 5

**5. EXPERIMENT ANALYSIS**

**5.1 DATA ANALYSIS**

The dataset presented in the image appears to be a manufacturing or production log with various parameters recorded for different units of production. Here's a breakdown of the columns and the type of information they contain:

**UDI:** Unique identifier for each production unit.

**Product ID:** Identifier for the product type or model.

**Type:** Categorical variable indicating the type of product, with at least three categories observed: L, M, and H.

**Air temperature [K]:** Measurement of the air temperature in Kelvin during the production process.

**Process temperature [K]:** Measurement of the temperature of the production process itself, also in Kelvin.

**Rotational speed [rpm]:** The speed at which a component of the machine is rotating, measured in revolutions per minute (rpm).

**Torque [Nm]:** The torque applied during the process, measured in Newton-meters (Nm).

**Tool wear [min]:** Indicates how much wear has occurred on the tool used in the production, measured in minutes.

**Failure Type:** Categorical variable indicating the type of failure, if any, with "No Failure" being one of the categories.

**Machine failure:** Binary indicator (0 or 1) where 1 indicates a machine failure.

**TWF:** Possibly a specific type of failure or a binary indicator related to the tool.

**HDF:** Another specific type of failure or binary indicator, potentially related to hardware.

**PWF:** Another failure type or binary indicator, possibly related to power or process.

**OSF:** Another failure type or binary indicator, potentially indicating an operational system failure.

**RNF:** Another failure type or binary indicator, the meaning of which is not immediately clear from the image.

From this dataset, one could retrieve a variety of information, including:

The relationship between environmental conditions (like air and process temperatures) and production outcomes (like failure types).

The correlation between operational parameters (like rotational speed and torque) and machine failures.

Patterns of tool wear over time and its association with different product types or production conditions.

The frequency and distribution of different failure types across the production units.

This data could be used for predictive maintenance, quality control, and optimization of the production process by identifying factors that lead to failures or suboptimal performance.

**Statistics and Inferences:**

Dataset Overview:

The dataset comprises 10,000 instances, each uniquely identified by a UDI (Unique Data Identifier). It consists of both numerical and categorical variables, providing insights into various aspects of the machining process.

**Numerical Variables:**

**Air Temperature [K]:** The air temperature, with a mean of approximately 300.00 Kelvin and a standard deviation of 2.00 Kelvin, showcases a relatively stable environment. The distribution, centered around the mean with a median of 300.10 K, indicates symmetric behaviour.

**Process Temperature [K]:** Process temperature averages around 310.01 Kelvin, with a narrow spread indicated by a standard deviation of 1.48 Kelvin. Similar to air temperature, the distribution appears symmetric around the mean, with a median of 310.10 K.

**Rotational Speed [rpm]:** The mean rotational speed is approximately 1538.78 rpm, displaying a wider range of variability with a standard deviation of 179.28 rpm. The positively skewed distribution suggests higher values occur more frequently, with a median of 1503 rpm.

**Torque [Nm]:** Torque exhibits a mean of about 39.99 Nm and a standard deviation of 9.97 Nm, indicating moderate variability. The distribution appears roughly symmetric around the mean, with a median of 40.10 Nm.

**Tool Wear [min]:** Tool wear time ranges widely, from 0 to 253 minutes, with a mean of approximately 107.95 minutes and a standard deviation of 63.65 minutes. The distribution seems roughly symmetric around the mean, suggesting varying tool maintenance requirements.

The Statistical Characteristics depicting mean, standard deviation, median and distribution for the features is shown in the Table 5.1(a)

| Feature | Mean | Standard Deviation | Median | Distribution |
|---|---|---|---|---|
| Air Temperature [K] | 300.00 | 2.00 | 300.10 | Symmetric |
| Process Temperature [K] | 310.01 | 1.48 | 310.10 | Symmetric |
| Rotational Speed [rpm] | 1538.78 | 179.28 | 1503 | Positively skewed |
| Torque [Nm] | 39.99 | 9.97 | 40.10 | Roughly symmetric |
| Tool Wear [min] | 107.95 | 63.65 | - | - |

*Table 5.1(a) Statistical Characteristics*

**Categorical Variable (Machine Failure):** The dataset includes binary variables indicating the presence or absence of different types of machine failures, including TWF, HDF, PWF, OSF, and RNF. Each failure type has a corresponding binary variable, with mean values representing the proportion of instances where that specific failure type occurred. For instance, the mean values for TWF, HDF, PWF, OSF, and RNF are 0.0339, 0.0046, 0.0115, 0.0095, and 0.0098, respectively, indicating the prevalence of each failure type in the dataset.

Inferences: The dataset provides comprehensive insights into the operating conditions of the machinery, including temperature, speed, torque, and tool wear. Stable air and process temperatures suggest controlled environmental conditions, while variable rotational speed and torque indicate operational variability. Wide-ranging tool wear times highlight potential maintenance needs and varying tool lifespan. Analysis of categorical variables offers valuable information on the occurrence of different types of machine failures, crucial for predictive maintenance and quality control.Understanding these variables and their distributions is essential for optimizing the machining process, predicting maintenance requirements, and ensuring product quality and operational efficiency.

**Data Visualization:**

Data visualization plays a crucial role in exploring relationships between variables, identifying patterns, and gaining insights into the dataset. Here, we'll discuss various types of visualizations, including histograms, box plots, and their interpretations.

**Histograms:**

A histogram provides a graphical representation of the distribution of a single numerical variable. It divides the range of values into intervals, or bins, and counts the number of observations that fall into each bin. Histograms are particularly useful for understanding the frequency or density distribution of continuous variables.In the context of predictive maintenance, histograms could be used to visualize the distribution of variables such as air temperature, process temperature, rotational speed, torque, and tool wear. Understanding the distribution of these variables can provide insights into their typical ranges and variability. For example, a histogram of air temperature may show whether the temperature remains relatively stable or fluctuates widely over time. Similarly, a histogram of tool wear can reveal how frequently tools degrade or wear out during operation.

**Box Plots:**

Box plots, also known as box-and-whisker plots, visualize the distribution of a numerical variable across different categories or groups. They display the median, quartiles, and potential outliers of the data. Box plots are effective for comparing the distribution of a variable between different groups and identifying any variability or outliers.Box plots allow us to compare the distribution of variables across different categories, such as failure types in the case of predictive maintenance. For instance, a box plot comparing air temperature across different failure types can indicate

whether certain types of failures are associated with specific temperature ranges. A wider spread of the box plot may suggest greater variability in air temperature during those failure events. Similarly, box plots for other variables like process temperature, rotational speed, torque, and tool wear can provide valuable insights into how these parameters vary during different types of failures.

Histograms and box plots are essential tools in exploratory data analysis (EDA) for understanding the distribution and variability of numerical variables and their relationship with categorical variables, such as failure types in the context of predictive maintenance.

Histograms provide a visual representation of the frequency distribution of a single numerical variable. By dividing the range of values into bins and counting the number of observations falling into each bin, histograms allow us to see the shape and spread of the data. For instance, a histogram of air temperature can show whether the temperatures recorded are concentrated around a central value or if they vary widely. Similarly, histograms for other variables like process temperature, rotational speed, torque, and tool wear can reveal important insights into their distributions and variability.

On the other hand, box plots are particularly useful for comparing the distribution of a numerical variable across different categories or groups, such as failure types. A box plot visually summarizes the central tendency, spread, and variability of the data within each category. For example, a box plot comparing air temperature across different failure types can help identify whether certain types of failures are associated with specific temperature ranges. A wider spread of the box plot for a particular failure type may indicate greater variability in air temperature during those failure events.

By analysing histograms and box plots, we can uncover patterns, trends, and potential relationships between variables. These visualizations provide valuable insights that can inform decision-making processes, such as identifying factors contributing to machine failures and optimizing maintenance strategies. Therefore, incorporating histograms and box plots in the exploratory data analysis stage is essential for gaining a comprehensive understanding of the dataset and guiding subsequent analyses and modeling efforts.

## 5.1.1. CLEANING

**Handling Missing Values:**

Missing values are a common occurrence in real-world datasets and can adversely affect the performance of machine learning models. In this preprocessing step, the dataset is examined to identify any missing values using `pd.isnull(data)`, which generates a Boolean DataFrame indicating the presence of missing values. Following this, `pd.isnull(data).sum()` is used to calculate the sum of missing values for each column.Since missing values can lead to biased analysis and inaccurate predictions, it's crucial to handle them appropriately. Common techniques for handling missing values include imputation (replacing missing values with a statistical measure like mean, median, or mode) or removal of rows or columns containing missing values.In this specific case, the absence of any missing values in the dataset is confirmed, ensuring that there is no need for imputation or removal of rows/columns containing missing values.

**Removing Duplicate Values:**

Duplicate values in a dataset can skew analysis and lead to incorrect conclusions. Therefore, it's essential to identify and remove duplicate entries.The code snippet `data.duplicated().any()` checks whether any duplicate rows exist in the dataset. If duplicates are found, they can be removed using `data.drop_duplicates()`.In this scenario, the absence of duplicate values in any column is confirmed, indicating that the dataset is already free from duplicate entries.

**Dropping Columns:**

Certain columns in a dataset may not contribute significantly to the analysis or may even introduce noise into the model. In such cases, it's beneficial to drop these columns.The code snippet `od = od.drop("UDI", axis=1, errors="ignore")`, `od = od.drop("Product ID", axis=1, errors="ignore")`, and `od = od.drop("Type", axis=1, errors="ignore")` removes three columns ("UDI", "Product ID", and "Type") from the dataset.

These columns are likely unique identifiers or categorical variables that do not provide meaningful information for the analysis. Dropping them simplifies the dataset and focuses the analysis on relevant features.

**Handling Categorical Data:**

Categorical data represents discrete values and is often encoded to numerical form before being used in machine learning models.In this preprocessing step, the "Failure Type" column containing categorical values representing different types of failures is encoded using one-hot encoding.One-hot encoding transforms categorical variables into binary vectors, with each category represented as a binary feature (0 or 1). This encoding scheme prevents the model from misinterpreting categorical variables as ordinal when they are not.

**Feature Scaling:**

Machine learning models may perform poorly if features have different scales. Feature scaling standardizes the range of features, ensuring that they contribute equally to the model.Standard scaling (also known as Z-score normalization) is a common technique used to scale numerical features. It transforms the data such that the mean of each feature is 0 and the standard deviation is 1.Scaling the features before training the models helps in improving convergence speed and overall model performance.

By performing these preprocessing steps, the dataset is made clean, consistent, and suitable for further analysis and model training. Each step contributes to enhancing the quality of the data and improving the reliability of the subsequent analyses and predictions. If you have any specific questions or require further clarification on any aspect of the preprocessing steps, feel free to ask.

### 5.1.2. DATA PREPARATION & PREPROCESSING

**Data Loading:**

The process starts by loading the dataset from a CSV file into a Pandas DataFrame. This step ensures that the data is in a structured format and can be easily manipulated and analysed using Python's Pandas library.

**Data Inspection:**

After loading the data, an initial inspection is performed to understand its structure, data types, and characteristics. This step involves checking for missing values, duplicates, and gaining insights into the distribution of the features and target variable. Understanding the data's quality and characteristics is crucial for making informed decisions during preprocessing.

**Feature Selection:**

Certain columns that do not contribute to the predictive modeling process, such as unique identifiers (UDI, Product ID, Type), are removed from the dataset. These columns are dropped as they don't contain valuable information for predicting failure types and may even introduce noise into the model.

**Handling Categorical Data:**

The target variable, Failure Type, contains categorical labels representing different types of failures. To prepare these categorical labels for machine learning algorithms, they are encoded into numerical values using one-hot encoding. This transformation ensures that the algorithms can effectively learn from the categorical data by representing each category as a binary feature.

**Feature Scaling:**

Numerical features (such as Air temperature, Process temperature, Rotational speed, Torque, Tool wear) are often measured in different units and scales. Standardization is applied to scale these features to a common scale, usually with a mean of 0 and a standard deviation of 1. Scaling ensures that each feature contributes equally to the model training process and prevents features with larger scales from dominating the learning process.

**Train-Test Split:**

The dataset is split into training and testing sets to evaluate the performance of the trained models. The training set is used to train the models, while the testing set is used to assess how well the models generalize to unseen data. This step helps to detect overfitting and ensures that the model's performance estimates are reliable.

**Model Training:**

Three different classifiers (K-Nearest Neighbors, Decision Tree, Random Forest) are trained on the training data to predict the failure types. Each classifier learns patterns in the data to make predictions based on input features.

**Model Evaluation:**

The performance of each trained model is evaluated using various metrics such as accuracy, precision, recall, and F1-score. These metrics provide insights into how well the models are

performing in terms of correctly predicting failure types. Evaluating multiple metrics helps to get a comprehensive understanding of the model's performance.

**Hyperparameter Tuning:**

Grid search is performed to find the optimal hyperparameters for the Decision Tree classifier. Hyperparameters are parameters that are not directly learned by the model during training but affect the model's behavior. Tuning hyperparameters helps to improve the model's performance by finding the best combination of parameter values.

**Bagging Ensemble:**

Bagging ensemble technique is applied to the base classifiers (KNN, Decision Tree, Random Forest) to further enhance predictive performance. Bagging (Bootstrap Aggregating) involves training multiple models on different subsets of the training data and combining their predictions to make final predictions. This ensemble technique helps to reduce overfitting and improve the robustness of the model.

**User Input Prediction:**

Finally, a function is implemented to take user input for various features (air temperature, process temperature, rotational speed, torque, tool wear) and predict the failure type using the best-performing bagging ensemble model (Decision Tree classifier). This allows users to interact with the trained model and obtain predictions based on their input feature values.

By following these preprocessing steps, the data is appropriately prepared and optimized for training machine learning models, leading to accurate predictions of failure types based on input feature values. Each step plays a crucial role in ensuring the quality and effectiveness of the predictive modelling process.

### 5.1.3. DATA TRANFORMATION

**Drop Irrelevant Columns:**

Unique identifier columns like UDI (Unique Device Identifier) and Product ID don't contain any predictive information and thus are removed to avoid noise in the model.The Type column might represent some categorical information but was deemed irrelevant for the predictive modelling task, so it was also dropped.

**Handling Missing Values:**

Missing values can greatly affect the performance of machine learning models if not handled properly.In this case, the dataset was checked for missing values, and since none were found, no imputation or deletion was required. This ensures that the integrity of the dataset is maintained.

**Drop Duplicate Rows:**

Duplicate rows could skew the analysis and modelling results, leading to biased estimates.By removing duplicate rows, the dataset is cleaned to ensure that each observation is unique and contributes meaningfully to the analysis.

**Feature Engineering:**

Feature engineering involves creating new features or transforming existing ones to better represent the underlying patterns in the data.In this case, the Machine failure column was dropped as it was considered redundant for predicting failure types. Redundant features can introduce noise and unnecessary complexity to the model.

**Normalization and Standardization:**

Normalization and standardization are techniques used to rescale the features of a dataset.StandardScaler from scikit-learn was employed to standardize the numerical features, ensuring that they have a mean of 0 and a standard deviation of 1.This step is crucial for algorithms that rely on distance metrics, such as K-nearest neighbours or support vector machines, as it ensures that all features contribute equally to the model.

**Encoding Categorical Variables:**

Categorical variables need to be encoded into numerical format for the machine learning algorithms to process them effectively.One-hot encoding was used to encode the target variable Failure Type, creating binary columns for each category.This allows the algorithms to understand and interpret the categorical data correctly during training.

**Explaining Feature Engineering:**

Air temperature, Process temperature, Rotational speed, Torque, and Tool wear were subjected to specific transformations to better represent their underlying distributions and relationships with the target variable.These transformations ensure that the features are in a suitable format and

distribution for modelling, potentially enhancing the predictive performance of the machine learning algorithms.

In summary, preprocessing steps such as dropping irrelevant columns, handling missing values, removing duplicates, performing feature engineering, and encoding categorical variables are essential for preparing the data for machine learning. These steps help improve the quality of the dataset, mitigate potential biases, and ensure that the data is in a suitable format for training predictive models.

The statistical summary of the training the dataset is shown in the Table 5.1(b)

| | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] |
|---|---|---|---|---|---|
| **count** | 8.000000e+03 | 8.000000e+03 | 8.000000e+03 | 8.000000e+03 | 8.000000e+03 |
| **mean** | -3.342171e-14 | -2.343170e-14 | 1.811884e-16 | -3.126388e-16 | 9.414691e-17 |
| **std** | 1.000063e+00 | 1.000063e+00 | 1.000063e+00 | 1.000063e+00 | 1.000063e+00 |
| **min** | -2.352493e+00 | -2.897008e+00 | -1.981233e+00 | -3.628570e+00 | -1.709378e+00 |
| **25%** | -8.540660e-01 | -8.114199e-01 | -6.411277e-01 | -6.822446e-01 | -8.710502e-01 |
| **50%** | 4.499028e-02 | 6.318154e-02 | -2.036552e-01 | 1.926135e-02 | -1.087453e-03 |
| **75%** | 7.442563e-01 | 7.359519e-01 | 4.013305e-01 | 6.806813e-01 | 8.530578e-01 |
| **max** | 2.242683e+00 | 2.552432e+00 | 7.460419e+00 | 3.627006e+00 | 2.292451e+00 |

*Table 5.1(b) Statistical Summary*

## 5.1.4. FEATURE SELECTION

Feature selection is a crucial aspect of building machine learning models, especially when dealing with datasets containing numerous features. It involves identifying and selecting the subset of features that are most relevant for predicting the target variable accurately. By selecting only the most important features, you can improve model performance, reduce overfitting, and enhance interpretability.

Here's a more detailed explanation of the feature selection methods mentioned:

**Feature Importance Plots:**

This method involves training a machine learning model (e.g., Decision Tree, Random Forest, Gradient Boosting) and then evaluating the importance of each feature in predicting the target variable. In the case of tree-based models like Decision Trees and Random Forests, feature importance is calculated based on how frequently a feature is used to split the data and reduce impurity (e.g., Gini impurity or entropy). Features with higher importance scores are considered more influential in predicting the target. Visualizing these importance scores using bar plots provides a clear understanding of which features contribute the most to the model's predictions.

**Correlation Analysis:**

Correlation analysis examines the relationship between each input feature and the target variable. It calculates the correlation coefficient, which measures the strength and direction of the linear relationship between two variables. Features with higher absolute correlation coefficients with the target variable are more likely to be important predictors. Heatmaps are often used to visualize the correlation matrix, making it easier to identify significant relationships between features and the target.

**Model-Based Feature Selection:**

Certain machine learning algorithms inherently perform feature selection during model training. For example, Decision Trees and Random Forests select features based on their ability to split the data and predict the target variable effectively. These models assess the importance of each feature implicitly, allowing you to identify the most relevant features without additional analysis.

In the provided code, feature importance analysis is conducted using the `feature_importances_` attribute of a Decision Tree Classifier (`bagging_dtc`). This approach quantifies the importance of

each feature in predicting machine failure types. By examining these importance scores, you can determine which features have the most significant impact on the target variable.

Once the most important features are identified through feature selection, they can be used for model training and prediction. Focusing on these key features streamlines the modeling process, leading to more efficient and accurate predictive models. Additionally, using fewer features reduces the risk of overfitting and improves model interpretability, making it easier to understand the underlying patterns driving the predictions.

## 5.1.5. FEATURE SELECTION STRATEGIES

Feature selection is a critical step in the machine learning pipeline aimed at identifying the most relevant features from the dataset while discarding redundant or noisy ones. This process not only reduces the dimensionality of the feature space but also improves model performance, interpretability, and generalization by focusing on the most informative attributes. Here's a more detailed explanation of the methods mentioned:

**Manual Feature Selection:**

Domain Knowledge:

Domain experts or individuals familiar with the data can manually select features based on their understanding of the problem domain and the relationships between features and the target variable. For instance, if certain features are known to be irrelevant or highly correlated with others, they can be removed from the dataset.

Preliminary Analysis:

Initial exploration of the dataset might reveal features with low variance, constant values, or high missingness, indicating their lack of contribution to the predictive task. Such features can be identified and eliminated during data preprocessing.

**Automated Feature Selection:**

Recursive Feature Elimination (RFE):

RFE is an iterative method that starts with all features and eliminates the least important ones in each iteration, based on the feature importance scores provided by a chosen estimator. This process continues until the desired number of features is reached or until a predefined performance metric no longer improves.

Feature Importance from Tree-Based Models:

Tree-based models like Random Forest can provide feature importance scores based on how much each feature decreases impurity across all the trees in the ensemble. Features with higher importance scores are considered more relevant and can be selected for further analysis.

SelectKBest:

This method selects the top K features with the highest scores based on univariate statistical tests such as chi-square, ANOVA F-value, or mutual information. It evaluates each feature independently of the others and selects the most informative ones.

After selecting the relevant features using either manual or automated methods, the subsequent steps in the machine learning pipeline, such as model training, validation, and testing, remain consistent. However, by reducing the feature space, feature selection can lead to faster training times, simpler and more interpretable models, and improved generalization performance by reducing the risk of overfitting.

It's essential to note that the choice of feature selection method should be guided by the characteristics of the dataset, including its size, dimensionality, distribution, and the nature of the problem being addressed. Experimentation with different techniques and validation through cross-validation or holdout sets can help identify the most effective feature selection strategy for a particular task. Additionally, feature selection should be integrated into the overall model development process and considered alongside other steps such as data preprocessing, model selection, and hyperparameter tuning.

## 5.2 INTERPRETATION OF RESULTS

### 5.2.1 PERFORMANCE METRICS

**Accuracy:** Accuracy measures the proportion of correctly classified instances out of the total instances. In our case, the accuracy of the models ranged from around 96.75% to 98.2%. The Random Forest model performed the best with an accuracy of 97.85%.

**Precision (Micro-Average):** Precision measures the ability of the classifier not to label as positive a sample that is negative. The micro-average precision is computed across all classes as a single

precision value. In our case, the precision ranged from around 97.54% to 98.29%. Again, the Random Forest model had the highest precision.

**Recall (Micro-Average):** Recall, also known as sensitivity or true positive rate, measures the ability of the classifier to identify all relevant instances. The micro-average recall is computed across all classes as a single recall value. In our case, the recall ranged from around 97.3% to 97.9%.

**F1-Score (Micro-Average):** F1-Score is the harmonic mean of precision and recall. The micro-average F1-score is computed across all classes as a single F1 value. In our case, the F1-score ranged from around 97.42% to 98.12%.

Based on these metrics, we can conclude that the Random Forest model performed the best among the evaluated models in terms of accuracy, precision, recall, and F1-score.

**Data Cleaning:**

This step involves preparing the dataset for modeling by handling missing values, removing unnecessary columns, and encoding categorical variables. It ensures that the data is in a suitable format for training machine learning models.

**Model Selection:**

Three different models were chosen for evaluation: K-Nearest Neighbors (KNN), Decision Tree, and Random Forest. Each model has its own characteristics and is suitable for different types of datasets and problem domains. By selecting multiple models, the evaluation can identify which one performs best for the specific task.

**Model Training:**

After selecting the models, they are trained using the preprocessed training data. During training, the models learn patterns and relationships in the data that enable them to make predictions on unseen data.

**Model Evaluation:**

Once trained, the models are evaluated using various performance metrics such as accuracy, precision, recall, and F1-score. These metrics provide insights into how well the models are performing in terms of correctly classifying instances and avoiding false positives.

**Tuning Hyperparameters:**

Hyperparameters are parameters that are set before the learning process begins. Tuning hyperparameters involves selecting the optimal values for these parameters to improve the performance of the model. In this case, the hyperparameters of the Decision Tree model were tuned using GridSearchCV to optimize its performance.

**Bagging Ensemble:**

Bagging (Bootstrap Aggregating) is a technique used to improve the stability and accuracy of machine learning algorithms. It involves training multiple instances of the same base learning algorithm on different subsets of the training data and then combining their predictions to obtain a final prediction. In this case, bagging classifiers were applied to each of the base classifiers (KNN, Decision Tree, Random Forest) to further improve their performance.

**User Input Prediction:**

Once the models are trained and evaluated, they can be used to make predictions on new, unseen data. A function was created to allow users to input their own data, and the Bagging Decision Tree model was used to predict the failure type based on the provided input. This functionality enables the model to be used in real-world scenarios where users need to make predictions based on their specific data.

Overall, by following these steps and leveraging appropriate techniques and algorithms, the Random Forest model emerged as the best performer among the evaluated models.

The Comparison of the Evaluation metrics for the Models is shown in the Table 5.2(a)

| Model | ACCURACY | PRECISION | RECALL | F1-SCORE |
|---|---|---|---|---|
| KNN | 0.9705 | 0.976358 | 0.9705 | 0.973420 |
| Decision Tree | 0.9710 | 0.971000 | 0.9710 | 0.971000 |
| Random Forest | 0.9795 | 0.982940 | 0.9795 | 0.981217 |

*Table 5.2.(a) Model Comparison*

The image contains a table that lists the performance metrics for three machine learning models: KNN, Decision Tree, and Random Forest. The metrics provided are Accuracy, Precision, Recall, and F1-Score. Here are the values for each model. From these metrics, we can infer that the Random Forest model has the highest values across all four-performance metrics, indicating that it performs better than the KNN and Decision Tree models for the data and context in which it was tested. The statement at the bottom of the image, "the best and efficient model is random forest," is a conclusion based on these metrics.

A list of three accuracy scores for different machine learning models that use bagging as an ensemble technique is shown in the Table 5.2(b)

| | |
|---|---|
| **Bagging KNN Accuracy** | 0.9675 |
| **Bagging Decision Tree Accuracy** | 0.982 |
| **Bagging Random Forest Accuracy** | 0.9785 |

*Table 5.2(b) Bagging accuracy*

From this information, we can infer the following:

The models have been evaluated on some dataset, and their performance is measured in terms of accuracy, which is the proportion of correct predictions made by the model.

Bagging, which stands for Bootstrap Aggregating, is an ensemble method that aims to improve the stability and accuracy of machine learning algorithms. It involves training multiple models using different subsets of the training dataset and then aggregating their predictions.KNN (k-nearest neighbors), Decision Tree, and Random Forest are different types of machine learning algorithms that have been used here.

The Decision Tree model with bagging has the highest accuracy (0.982), suggesting it performed best on the given task out of the three models.

The KNN model with bagging has the lowest accuracy (0.9675), but it is still relatively high, indicating it also performed well.

The Random Forest, which is an ensemble method itself that uses multiple decision trees, has a slightly lower accuracy (0.9785) than the single Decision Tree model with bagging in this particular case.

These results could be used to select the best model for deployment in a production environment or to guide further model tuning and validation.

## 5.2.2 COMPARISON OF MODELS

In this analysis, we conducted a multiclass classification task using various machine learning models to predict the type of failure in a maintenance dataset. Here's a breakdown of the steps we followed and the key findings:

**Data Preprocessing:**

We loaded the dataset using pandas and examined its structure, confirming that there were no missing values or duplicate entries. Columns that served as unique identifiers (UDI and Product ID) were dropped, as they do not contribute to predictive modeling.The dataset was then split into features (independent variables) and the target variable (Failure Type).We observed that the dataset was imbalanced, with the majority class being "No Failure."

**Exploratory Data Analysis (EDA):**

We visualized the distribution of different failure types in the dataset using bar plots, identifying "No Failure" as the predominant class.We analyzed the average values of various features grouped by failure type to understand their relationships.

**Feature Scaling:**

We standardized the numerical features using z-score normalization to ensure that they have a mean of 0 and a standard deviation of 1. This step is crucial for many machine learning algorithms.

**Model Training and Evaluation:**

We trained three different classifiers: K-Nearest Neighbors (KNN), Decision Tree Classifier (DTC), and Random Forest Classifier (RFC).Model performance was evaluated using accuracy, precision, recall, and F1-score metrics. Random Forest performed the best, achieving an accuracy

of 97.95%.We further fine-tuned the Decision Tree Classifier using GridSearchCV to optimize its hyperparameters, resulting in a slight improvement in performance.

**Bagging Classifier:**

We explored the use of ensemble learning by implementing Bagging Classifiers with the three base classifiers (KNN, DTC, and RFC).The Bagging Decision Tree Classifier achieved the highest accuracy of 98.2%.

**User Input Prediction:**

We created a function to allow users to input feature values and obtain predictions for the failure type using the best-performing model (Bagging Decision Tree Classifier).Users can input air temperature, process temperature, rotational speed, torque, and tool wear, and the model will predict the corresponding failure type.

**Interpretation of Results:**

The analysis demonstrates the effectiveness of machine learning models in predicting failure types based on process parameters.The Bagging Decision Tree Classifier emerged as the most accurate model for this dataset, achieving an accuracy of 98.2%.Users can utilize the predictive model to anticipate potential equipment failures and take proactive maintenance actions, thereby improving operational efficiency and minimizing downtime.

Overall, the analysis highlights the significance of data preprocessing, model selection, and evaluation in developing robust predictive maintenance systems. Additionally, the deployment of ensemble learning techniques such as bagging can further enhance model performance and reliability.

### 5.2.3 OVERFITTING AND GENERALIZATION

The models trained in this multiclass classification task appear to perform reasonably well, with high accuracy scores across different algorithms. However, it's essential to delve deeper into aspects such as overfitting and generalization to ensure the robustness of the models.

**Overfitting:** Overfitting occurs when a model learns the training data too well, capturing noise and random fluctuations rather than underlying patterns. This can lead to poor performance on unseen data.

**Decision Tree:** By optimizing the hyperparameters using GridSearchCV, we attempted to mitigate overfitting. The best model obtained had a max_depth of 10 and min_samples_leaf of 5. These constraints help prevent the model from becoming too complex and overfitting the training data excessively.

**Random Forest:** Random forests are less prone to overfitting compared to decision trees due to the ensemble nature of combining multiple trees. Each tree is trained on a random subset of features and samples from the training data, reducing the likelihood of overfitting.

**Bagging Classifier:** The bagging classifier combines multiple base classifiers (in this case, decision trees) trained on different subsets of the training data. By aggregating the predictions of these classifiers, it helps reduce overfitting.

**Generalization:** Generalization refers to the ability of a model to perform well on unseen data, indicating its capability to capture underlying patterns rather than noise.

**Evaluation Metrics:** Accuracy, precision, recall, and F1-score are commonly used metrics to evaluate model performance. These metrics provide insights into how well the models generalize to unseen data by measuring their ability to correctly classify instances across different classes.

**Cross-validation:** By using techniques like k-fold cross-validation during model training and hyperparameter tuning, we can assess how well the models generalize to different subsets of the data. This helps ensure that the models' performance is consistent across various data partitions.

In summary, while the models achieved high accuracy scores, further analysis is needed to ensure they generalize well to unseen data and are not overfitting the training data. Techniques like hyperparameter tuning, cross-validation, and ensemble methods can help improve model robustness and generalization capabilities. Additionally, monitoring the models' performance on validation or test datasets can provide insights into their real-world effectiveness.

### 5.2.4 INTERPRETATION OF MODEL OUTPUTS

**Decision Tree Classifier (DTC):**

Model Explanation: The Decision Tree Classifier (DTC) is a simple yet powerful algorithm for classification tasks. It creates a tree-like structure where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents the outcome or class label. The algorithm makes decisions by recursively partitioning the data based on the values of features.

Accuracy: The reported accuracy of approximately 98.2% on the test data suggests that the DTC performed exceptionally well in classifying instances correctly.

Hyperparameters: The chosen hyperparameters, max_depth and min_samples_leaf, play crucial roles in controlling the complexity of the tree. A deeper tree (higher max_depth) can capture intricate relationships in the data but may lead to overfitting, while min_samples_leaf prevents the algorithm from creating nodes that have few training samples.

Interpretation: The DTC successfully learned the patterns in the data and was able to make accurate predictions based on the provided features, which include air temperature, process temperature, rotational speed, torque, and tool wear.

**Bagging Decision Tree Classifier:**

Model Explanation: Bagging is an ensemble learning technique that aims to improve the stability and accuracy of a model by training multiple instances of the same base classifier on different subsets of the training data and then combining their predictions through averaging (for regression) or voting (for classification).

Accuracy: Achieving a similar accuracy of approximately 98.2% on the test data as the standalone DTC indicates that bagging effectively reduced overfitting and enhanced the model's generalization ability.

Interpretation: By training multiple decision tree classifiers on different subsets of the training data and combining their predictions, the Bagging Decision Tree Classifier mitigated the risk of overfitting and provided a more stable and reliable model for predicting machine failure types.

**Random Forest Classifier (RFC):**

Model Explanation: Random Forest is another ensemble learning technique that builds multiple decision trees and combines their predictions to improve accuracy and reduce overfitting. It introduces randomness by considering a random subset of features for each split in the decision trees.

Accuracy: The accuracy of approximately 97.9% on the test data, slightly lower than the DTC and Bagging Decision Tree Classifier, could be due to the inherent randomness introduced during the training process.

Interpretation: Despite the slightly lower accuracy compared to the other models, Random Forest classifiers are known for their robustness and ability to handle high-dimensional data effectively. The model still provides a reliable solution for predicting machine failure types based on the given features.

**Overall Interpretation:**

Performance: Both the Decision Tree Classifier and Bagging Decision Tree Classifier performed exceptionally well, achieving accuracies above 98% on the test data. This indicates that they effectively captured the underlying patterns in the data and made accurate predictions.

Generalization: The Bagging Decision Tree Classifier demonstrated improved stability and generalization compared to the standalone Decision Tree Classifier, while the Random Forest Classifier provided robustness despite a slightly lower accuracy.

Feature Importance: These models can also provide insights into the importance of different features (e.g., air temperature, process temperature, rotational speed, torque, and tool wear) in predicting machine failure types, which can be valuable for understanding the underlying processes and optimizing maintenance strategies.

In summary, all three models - Decision Tree Classifier, Bagging Decision Tree Classifier, and Random Forest Classifier - are effective for predicting machine failure types based on the provided features. The choice between them may depend on factors such as interpretability, computational resources, and the specific requirements of the application.

### 5.2.5 FEATURE OUTPUTS

Analyzing feature importance in a machine learning model like the Decision Tree Classifier (DTC) or Random Forest Classifier (RFC) can provide valuable insights into the factors that significantly contribute to predicting the outcome or failure type. In our case, we'll focus on the Random Forest model, which tends to provide more reliable feature importance estimates compared to a single decision tree.

Feature importance is calculated based on how much each feature contributes to decreasing impurity in the decision trees within the Random Forest ensemble. In other words, features with higher importance values are considered more influential in making predictions.

**Air Temperature [K]:**

This feature represents the temperature of the air during the manufacturing process, which can significantly impact the materials and machinery involved.A high feature importance score for air temperature suggests that it plays a crucial role in determining the likelihood of different failure types.For instance, if "Heat Dissipation Failure" is more common at higher temperatures, the model will assign higher importance to air temperature as it helps distinguish between failure types.

**Process Temperature [K]:**

While closely related to air temperature, process temperature specifically reflects the conditions within the manufacturing environment.A significant importance score for process temperature indicates that it provides unique information beyond ambient air temperature.Process temperature might affect material properties, chemical reactions, or thermal stress in ways that air temperature alone cannot capture.

**Rotational Speed [rpm]:**

Rotational speed represents how fast the machinery operates, influencing factors such as mechanical stress and wear.A high feature importance score for rotational speed suggests that it's a critical factor in determining failure types.For example, if "Tool Wear Failure" is more likely at higher rotational speeds, the model will assign higher importance to this feature when predicting failure types.

**Torque [Nm]:**

Torque measures the twisting force exerted during operation, which can lead to mechanical failure or damage.A substantial importance score for torque indicates its significant role in predicting failure types.Excessive torque might contribute to "Overstrain Failure" or "Tool Wear Failure," and the model will prioritize this feature accordingly.

**Tool Wear [min]:**

Tool wear reflects the deterioration of manufacturing tools over time, which can impact their performance and reliability.A high feature importance score for tool wear suggests that it's a critical factor in predicting failure types.If certain failure types are more likely as tool wear increases, the model will give greater weight to this feature in its predictions.

Overall, the feature importance scores provide valuable insights into which factors have the most significant influence on predicting failure types in the manufacturing process. By understanding these relationships, manufacturers can prioritize maintenance efforts, adjust operating parameters, and optimize processes to mitigate the risk of failures and improve overall efficiency and reliability.

## 5.2.6 INTERPRETATION OF THE MODEL OUTPUT

**Input Parameters:**

Air temperature [K]: This parameter represents the temperature of the air in the manufacturing environment. Air temperature can influence the thermal conditions within the machinery and affect various components' performance and reliability.

Process temperature [K]: The process temperature refers to the temperature of the manufacturing process itself. It directly affects material properties, chemical reactions, and the behavior of components and machinery involved in the process.

Rotational speed [rpm]: Rotational speed indicates how fast the machinery or components are rotating during the manufacturing process. It influences factors such as wear and tear, friction, and heat generation within the equipment.

Torque [Nm]: Torque is a measure of the rotational force applied during the process. It determines the mechanical stress experienced by components and machinery, affecting their durability and performance.

Tool wear [min]: Tool wear represents the duration for which the tool has been in use during the manufacturing process. It reflects the tool's condition and its remaining useful life, which is crucial for maintaining process quality and efficiency.

**Predicted Failure Type:**

The model predicts the type of failure that is most likely to occur based on the input parameters provided. Failure types can include various categories such as mechanical failures, thermal failures, electrical failures, etc. In the example provided, the predicted failure type is "Heat Dissipation Failure," indicating that the model identifies issues related to the dissipation of heat as the primary risk based on the input parameters.

**Business Implications:**

Proactive Measures: Understanding the predicted failure type allows stakeholders to take proactive measures to prevent or mitigate the identified risk. For example, if the model predicts a "Heat Dissipation Failure," stakeholders can implement better cooling mechanisms, adjust process parameters, or schedule maintenance activities to prevent overheating issues.

Predictive Maintenance: Predictive maintenance strategies can be developed based on the model outputs. By anticipating potential failures, resources can be allocated efficiently, minimizing downtime and optimizing productivity. This approach reduces the likelihood of unexpected breakdowns and extends equipment lifespan.

Decision-making: The model outputs can inform decision-making processes related to equipment maintenance, resource allocation, and process optimization. Stakeholders can prioritize maintenance activities based on the predicted failure types, ensuring that critical issues are addressed promptly and resources are utilized effectively.

Continuous Monitoring: Establishing continuous monitoring and analysis of the manufacturing process can feed real-time data into the model, enabling dynamic adjustments and proactive interventions to prevent failures before they occur. This iterative approach fosters a culture of continuous improvement and resilience in manufacturing operations.

In summary, the model outputs provide valuable insights into potential failure scenarios in the manufacturing process, empowering stakeholders to make informed decisions and implement proactive measures to enhance operational efficiency, reduce downtime, and mitigate risks effectively.

## 5.2.7 LIMITATION AND FUTURE WORKS

**Data Quality Issues:**

This limitation highlights the concern that the analysis relies heavily on the assumption that the provided data is accurate and representative of the actual system. If there are inaccuracies or biases in the data, it could lead to flawed conclusions. For instance, if the data is incomplete or contains errors, the predictive models may not perform optimally.

Solutions: Conducting thorough data validation and preprocessing steps can help mitigate data quality issues. Techniques such as data cleaning, outlier detection, and missing value imputation can improve the overall quality of the dataset.

**Sample Size Limitations:**

The analysis acknowledges that the dataset used for analysis may have limitations in terms of sample size. A small sample size can limit the generalizability of the findings and increase the risk of overfitting, where the model performs well on the training data but poorly on unseen data.

Solutions: Increasing the sample size by collecting more data can help address this limitation. Additionally, techniques such as bootstrapping or resampling methods can be employed to generate synthetic data points, thereby augmenting the dataset and improving model robustness.

**Model Assumptions:**

The machine learning models used in the analysis make certain assumptions about the data distribution and relationships between variables. If these assumptions are violated, it can adversely affect the performance of the models.

Solutions: Conducting diagnostic checks to assess the validity of model assumptions and conducting sensitivity analyses to evaluate model performance under different assumptions can help address this limitation. Additionally, exploring alternative modeling techniques that are more flexible or robust to violations of assumptions may be beneficial.

**Future Directions:**

**Collecting Additional Data:**

Gathering more data, especially encompassing diverse failure types and associated features, is essential for improving the predictive capabilities of the models. This could involve collecting data over a more extended period or from multiple sources to capture a broader range of scenarios and patterns.

Solutions: Collaborating with industry partners or leveraging sensor data from different equipment can help enrich the dataset. Additionally, incorporating data from different geographical locations or operational contexts can enhance the robustness and generalizability of the models.

**Refining the Models:**

Fine-tuning model hyperparameters and exploring alternative algorithms can enhance model performance. Techniques such as cross-validation and grid search can be employed to systematically optimize model parameters and improve predictive accuracy.

Solutions: Experimenting with different machine learning algorithms, such as gradient boosting machines, random forests, or deep learning models, can help identify the best-performing model for the given dataset. Additionally, ensemble methods that combine multiple models can further boost predictive performance.

**Feature Engineering:**

Experimenting with different feature engineering techniques can help uncover hidden patterns in the data and improve model accuracy. This involves creating new features, transforming existing ones, or selecting the most informative features for modeling.

Solutions: Conducting exploratory data analysis to identify relevant features and leveraging domain knowledge to engineer informative features can enhance model performance. Techniques such as dimensionality reduction, feature selection, and feature scaling can also be employed to improve model efficiency and interpretability.

**Addressing Imbalance:**

If there's a significant class imbalance in the dataset (i.e., one class is much more prevalent than others), it can skew model predictions and affect performance. Techniques such as oversampling, undersampling, or using advanced algorithms like ensemble methods can help mitigate this issue.

Solutions: Employing resampling techniques to balance class distribution, such as oversampling minority classes or undersampling majority classes, can improve model training. Additionally, ensemble methods like bagging or boosting that handle class imbalance more effectively can be explored.

**Interpreting Model Outputs:**

Understanding the rationale behind model predictions and identifying influential features is crucial for actionable insights. Techniques such as feature importance analysis and model explainability methods can aid in interpreting model outputs and gaining insights into underlying patterns and relationships.

Solutions: Conducting feature importance analysis, SHAP (SHapley Additive Explanations) values, or LIME (Local Interpretable Model-agnostic Explanations) analysis can help elucidate the impact of different features on model predictions. Visualizing model explanations and communicating findings effectively to stakeholders can enhance the interpretability and trustworthiness of the models.

**External Validation:**

Validating the models on external datasets or real-world scenarios is essential to assess their generalizability and applicability in different contexts. This involves testing the models on unseen data to evaluate their performance and robustness.

Solutions: Collaborating with external partners or acquiring additional datasets from different sources can facilitate external validation. Conducting rigorous testing in diverse operational environments and comparing model performance against baseline methods or expert judgments can provide valuable insights into model effectiveness and reliability.

By addressing these limitations and pursuing the suggested future research directions, the analysis can yield more robust and actionable insights for predictive maintenance and failure prediction in industrial systems. This iterative process of refinement and validation is essential for continuously improving the accuracy, reliability, and practical utility of predictive maintenance models.

# CHAPTER 6

**CONCLUSION AND FUTURE WORKS**

**6.1 CONCLUSION**

The project successfully developed and implemented predictive analysis systems for preemptively identifying potential equipment failures in industrial settings. This accomplishment underscores the significance of leveraging advanced machine learning techniques to address the persistent challenge of equipment breakdowns.

Through extensive exploratory data analysis (EDA) and model development, critical patterns and indicators of equipment failures were identified. This deep dive into the data allowed for a thorough understanding of the underlying factors contributing to failures, enabling more accurate predictive models.

Key findings include the effectiveness of machine learning algorithms in predicting failure types and the importance of specific features in accurate prediction. By analyzing the performance of various algorithms and identifying key features, the project provided actionable insights for optimizing predictive maintenance strategies.

The project demonstrated the practical applicability of predictive maintenance in real-world industrial scenarios, showcasing potential cost savings and efficiency improvements. By illustrating tangible benefits such as reduced downtime and lower maintenance costs, the project highlighted the value proposition of predictive maintenance for industrial stakeholders.

**Achievement of Objectives:**

The project objectives, outlined in the introduction, were effectively achieved. The development of robust predictive maintenance frameworks, fostered understanding of failure prediction, and showcased practical applications of predictive maintenance represent significant milestones in addressing the challenge of equipment failures in industrial settings.By employing machine learning algorithms and rigorous data analysis techniques, the project successfully addressed the challenge of equipment failures in industrial settings. This achievement underscores the project's commitment to leveraging cutting-edge technology and methodologies to solve real-world problems in industrial maintenance.

**Significance of Findings:**

The findings of the project hold significant implications for the problem domain of industrial equipment maintenance. The demonstrated effectiveness of predictive maintenance frameworks in reducing downtime, lowering maintenance costs, and improving operational efficiency underscores their potential to drive significant improvements in industrial productivity and profitability.Furthermore, the proactive approach to maintenance enhances workplace safety, environmental conservation, and job security, aligning with broader societal goals. By mitigating risks associated with equipment failures, predictive maintenance contributes to creating safer, more sustainable, and productive industrial environments for all stakeholders.

**Reflection on Challenges and Unexpected Results:**

While the project achieved its objectives, it encountered challenges such as data quality issues, algorithm selection, and model interpretability. These challenges underscored the complexity of predictive maintenance in real-world industrial settings and the importance of robust data preprocessing and model evaluation techniques.Unexpected results or limitations in model performance might have arisen during the evaluation phase, necessitating iterative refinement and adjustment of methodologies. These challenges provided valuable learning opportunities and insights for future research and development efforts in the field of predictive maintenance.

**Closure and Future Directions:**

The project concludes with emphasis on its impact on industrial maintenance practices and potential future directions. Future research could focus on enhancing predictive models with advanced techniques such as deep learning, incorporating real-time data streams for dynamic maintenance scheduling, and expanding the application to diverse industrial sectors beyond those initially targeted.In conclusion, the project has demonstrated the feasibility and effectiveness of predictive maintenance frameworks in industrial settings. By leveraging machine learning algorithms and extensive data analysis, the project contributes valuable insights to the field of equipment maintenance, with implications for both industry practices and societal well-being.

## 6.2  LIMITATION

**Data Constraints:**

The analysis relies heavily on the provided dataset. If this dataset is limited in scope or doesn't encompass all possible scenarios or failure types, the model's predictions may not be accurate for

scenarios not represented in the data.For instance, if certain failure types are rare or weren't observed during the data collection period, the model won't be trained to predict them accurately.

**Model Assumptions:**

The model assumes that the features provided in the dataset are sufficient to accurately predict failure types. However, this assumption may not hold true if there are interactions between variables or external factors not considered in the analysis.Additionally, if there are missing features that are crucial for predicting failures, the model's performance may be suboptimal.

**Potential Sources of Bias or Error:**

Biases in the data collection process, such as sampling biases or measurement errors, can affect the accuracy of the analysis. For example, if certain types of failures are overrepresented or underrepresented in the dataset, the model's predictions may be skewed.The choice of machine learning algorithms and hyperparameters could introduce bias or overfitting, especially if not carefully selected or tuned.

**Validity and Generalizability:**

The findings of the analysis may not be generalizable to all industrial settings or machines. They are specific to the dataset and conditions under which the data was collected. The model's performance may vary when applied to different datasets or real-world scenarios, leading to limited generalizability.

**Areas for Improvement and Future Work:**

**Data Collection:**

Collecting more comprehensive and diverse datasets covering a wider range of failure types and operating conditions could improve the model's accuracy and generalizability. This involves ensuring that the dataset captures rare failure events and reflects a variety of operating conditions.

**Feature Engineering:**

Exploring additional features or transformations of existing features could enhance the model's ability to capture complex relationships between variables and improve prediction performance. Domain knowledge can be leveraged to engineer relevant features that may have been overlooked initially.

**Model Selection and Tuning:**

Experimenting with different machine learning algorithms and hyperparameters, as well as ensemble methods like stacking or boosting, could lead to better predictive performance. This

involves systematically evaluating different models and tuning their parameters to optimize performance.

**Domain Knowledge Integration:**

Involving domain experts in the analysis process can help validate assumptions, interpret model outputs, and incorporate expert knowledge into feature engineering or model selection. Domain experts can provide valuable insights into the underlying processes and relationships that the model may not capture.

**External Validation:**

Validating the model's performance on independent datasets or in real-world applications is crucial to assess its robustness and generalizability beyond the training data. This involves testing the model in different contexts to ensure that it performs reliably in diverse scenarios.

Overall, addressing these limitations and implementing the suggested improvements is essential to ensure the reliability and usefulness of the analysis for practical applications in industrial maintenance and failure prediction. Collaboration between data scientists, domain experts, and stakeholders is key to achieving these objectives effectively.

## 6.3 RECOMMENDATIONS FOR FUTURE RESEARCH

The recommendations for future research provided above encompass a wide range of potential avenues for further exploration in the domain of employee performance and job satisfaction. Let's elaborate on each recommendation:

**Longitudinal Study:**

Conducting a longitudinal study involves tracking the same variables and participants over an extended period. This approach allows researchers to observe changes and trends over time, providing insights into the stability and long-term effects of various factors on employee performance and job satisfaction. Longitudinal studies can help identify patterns, causal relationships, and potential interventions that may not be apparent in cross-sectional studies.

**Qualitative Analysis:**

Qualitative research methods, such as interviews, focus groups, and observations, delve into the subjective experiences, attitudes, and perceptions of individuals. Integrating qualitative analysis alongside quantitative approaches enriches the understanding of the underlying mechanisms driving employee performance and satisfaction. Qualitative data offer depth and context, capturing

nuances that quantitative measures alone may miss, thus providing a more comprehensive view of the phenomenon under study.

**Cross-Cultural Comparison:**

Cross-Cultural studies examine how cultural differences influence behaviour, attitudes, and outcomes in various contexts. By comparing different cultural settings, researchers can identify cultural-specific factors that impact employee performance and job satisfaction. Understanding these differences is essential for developing culturally sensitive management practices and fostering inclusivity and diversity in the workplace.

**Exploration of Mediating and Moderating Variables:**

Mediating variables explain the relationship between a predictor and an outcome, while moderating variables influence the strength or direction of this relationship. Investigating these variables enhances the understanding of the complex interplay between predictors and outcomes in the context of employee performance and satisfaction. For example, exploring how leadership styles moderate the relationship between job characteristics and job satisfaction can provide insights into effective leadership strategies.

**Integration of Advanced Analytical Techniques:**

Advanced statistical techniques, such as structural equation modelling (SEM) or hierarchical linear modelling (HLM), offer sophisticated ways to analyse complex data structures and relationships. These methods allow researchers to account for the hierarchical nature of organisational data, non-linear relationships, and latent variables, thereby refining the analytical approach and improving the validity of findings.

**Focus on Remote Work Dynamics:**

With the growing prevalence of remote work, understanding the dynamics of virtual work environments is critical. Research in this area can explore how factors like virtual communication, autonomy, and work-life balance impact employee performance and satisfaction. Insights from such studies can inform the development of effective remote work policies and practices, enabling organisations to better support remote workers.

**Incorporation of Technology and Analytics:**

Leveraging emerging technologies, such as artificial intelligence (AI) and big data analytics, allows for the analysis of vast amounts of data to gain real-time insights into employee behaviour and engagement. By harnessing data from digital platforms and communication tools, researchers

can identify patterns, trends, and predictors of performance and satisfaction, enabling proactive interventions and personalised management strategies.

**Interdisciplinary Collaboration:**

Collaboration across disciplines, including psychology, sociology, management, and data science, fosters a holistic understanding of workplace dynamics. Integrating diverse perspectives and methodologies enables researchers to address complex research questions more effectively and develop innovative solutions to practical challenges. Interdisciplinary collaboration also facilitates the translation of research findings into evidence-based practices and policies that promote healthier and more productive work environments.

The analysis relies heavily on the provided dataset. If this dataset is limited in scope or doesn't encompass all possible scenarios or failure types, the model's predictions may not be accurate for scenarios not represented in the data.

For instance, if certain failure types are rare or weren't observed during the data collection period, the model won't be trained to predict them accurately.

# REFERENCES

[1] Arena, S., Florian, E., Zennaro, I., Orrù, P. F., & Sgarbossa, F. (2022). A novel decision support system for managing predictive maintenance strategies based on machine learning approaches. Safety science, 146, 105529.

[2] Theissler, A., Pérez-Velázquez, J., Kettelgerdes, M., & Elger, G. (2021). Predictive maintenance enabled by machine learning: Use cases and challenges in the automotive industry. Reliability engineering & system safety, 215, 107864.

[3] Kaparthi, S., & Bumblauskas, D. (2020). Designing predictive maintenance systems using decision tree-based machine learning techniques. International Journal of Quality & Reliability Management, 37(4), 659-686.

[4] Justus, V., & Kanagachidambaresan, G. R. (2024). Machine learning based fault-oriented predictive maintenance in industry 4.0. International Journal of System Assurance Engineering and Management, 15(1), 462-474.

[5] Ni, F., Zang, H., & Qiao, Y. (2024, January). SMARTFIX: LEVERAGING MACHINE LEARNING FOR PROACTIVE EQUIPMENT MAINTENANCE IN INDUSTRY 4.0. In The 2nd International scientific and practical conference "Innovations in education: prospects and challenges of today"(January 16-19, 2024) Sofia, Bulgaria. International Science Group. 2024. 389 p. (p. 313).

[6] Paolanti, M., Romeo, L., Felicetti, A., Mancini, A., Frontoni, E., & Loncarski, J. (2018, July). Machine learning approach for predictive maintenance in industry 4.0. In 2018 14th IEEE/ASME International Conference on Mechatronic and Embedded Systems and Applications (MESA) (pp. 1-6). IEEE.

[7] Ameeth Kanawaday and Aditya Sane, 2017 "Machine learning for predictive maintenance of industrial machines using IoT sensor data," In 2017 8th IEEE International Conference and Service Science, Beijing.

[8] Kroll, B., Schaffranek, D., Schriegel, S., & Niggemann, O. (2014, September). System modeling based on machine learning for anomaly detection and predictive maintenance in industrial plants. In Proceedings of the 2014 IEEE emerging technology and factory automation (ETFA) (pp. 1-7). IEEE.

[9] Nunes, P., Santos, J., & Rocha, E. (2023). Challenges in predictive maintenance–A review. CIRP Journal of Manufacturing Science and Technology, 40, 53-67.

[10] Wang, H., Zhao, S., Cheng, X., Li, D., & Cheng, B. (2020). Predictive maintenance for rotating machinery based on long short-term memory network. IEEE Transactions on Industrial Informatics, 16(3), 1980-1989.

[11] Li, Z., Zhou, D., & Zhou, L. (2019). Predictive maintenance for industrial equipment based on big data and internet of things. Journal of Industrial Information Integration, 15, 56-67.

[12] Sun, Y., Yan, Z., Zhao, D., Xue, H., & Wang, F. (2021). Predictive maintenance for industrial equipment based on hybrid deep learning models. Journal of Manufacturing Processes, 64, 543-552.

[13] Tan, Y., Zhang, Y., Xu, Z., & Jiang, P. (2020). Predictive maintenance for industrial equipment using convolutional neural networks. Journal of Ambient Intelligence and Humanized Computing, 11(10), 4501-4511.

[14] Yu, T., Wang, D., Li, D., & Zhang, Z. (2018). Predictive maintenance of industrial equipment based on ensemble learning methods. Journal of Intelligent Manufacturing, 29(8), 1787-1800.

[15] Chen, S., Tian, Y., Li, J., & Cheng, Y. (2021). Predictive maintenance for industrial equipment using recurrent neural networks. Journal of Manufacturing Systems, 60, 329-340.

[16] Zhang, C., Ma, W., Zhao, F., & Wang, X. (2020). Predictive maintenance for industrial equipment using machine learning and sensor fusion techniques. Sensors, 20(22), 6517.

[17] Zhu, W., Zhang, W., Chen, W., & Zhang, H. (2019). Predictive maintenance for industrial equipment based on feature selection and support vector machine. Journal of Mechanical Engineering, 55(16), 166-176.

# APPENDIX-SCREENSHOTS

## DATA CLEANING:

```
In [1]: import pandas as pd
```

```
In [2]: data = pd.read_csv(r'C:\Users\nazia\Desktop\CAPSTONE\maintenance.csv')
```

```
In [3]: data
```

Out[3]:

| | UDI | Product ID | Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Failure Type | Machine failure | TWF | HDF | PWF | OSF | RNF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | M14860 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 2 | L47181 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 3 | L47182 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 4 | L47183 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 5 | L47184 | L | 298.2 | 308.7 | 1408 | 40.0 | 9 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 9995 | 9996 | M24855 | M | 298.8 | 308.4 | 1604 | 29.5 | 14 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 9996 | 9997 | H39410 | H | 298.9 | 308.4 | 1632 | 31.8 | 17 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 9997 | 9998 | M24857 | M | 299.0 | 308.6 | 1645 | 33.4 | 22 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 9998 | 9999 | H39412 | H | 299.0 | 308.7 | 1408 | 48.5 | 25 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |
| 9999 | 10000 | M24859 | M | 299.0 | 308.7 | 1500 | 40.2 | 30 | No Failure | 0 | 0 | 0 | 0 | 0 | 0 |

10000 rows × 15 columns

```
In [4]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 15 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   UDI                      10000 non-null  int64
 1   Product ID               10000 non-null  object
 2   Type                     10000 non-null  object
 3   Air temperature [K]      10000 non-null  float64
 4   Process temperature [K]  10000 non-null  float64
 5   Rotational speed [rpm]   10000 non-null  int64
 6   Torque [Nm]              10000 non-null  float64
 7   Tool wear [min]          10000 non-null  int64
 8   Failure Type             10000 non-null  object
 9   Machine failure          10000 non-null  int64
 10  TWF                      10000 non-null  int64
 11  HDF                      10000 non-null  int64
 12  PWF                      10000 non-null  int64
 13  OSF                      10000 non-null  int64
 14  RNF                      10000 non-null  int64
dtypes: float64(3), int64(9), object(3)
memory usage: 1.1+ MB
```

```
In [5]:  ► pd.isnull(data)
            pd.isnull(data).sum()
```

```
Out[5]: UDI                         0
        Product ID                  0
        Type                        0
        Air temperature [K]         0
        Process temperature [K]     0
        Rotational speed [rpm]      0
        Torque [Nm]                 0
        Tool wear [min]             0
        Failure Type                0
        Machine failure             0
        TWF                         0
        HDF                         0
        PWF                         0
        OSF                         0
        RNF                         0
        dtype: int64
```

```
In [6]:  ► is_duplicate_any_column = data.duplicated().any()
            if is_duplicate_any_column:
                print('There are duplicate values in at least one column.')
            else:
                print('There are no duplicate values in any column.')
            columns_with_missing = data.columns[data.isnull().any()].tolist()
            if columns_with_missing:
                print("Columns with missing values:", columns_with_missing)
            else:
                print("There are no missing values in any column.")
```

```
There are no duplicate values in any column.
There are no missing values in any column.
```

```
In [7]:  ► data.describe()
```

Out[7]:

| | UDI | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | TWF | HDF | PWF | OSF |
|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 10000.00000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000 |
| mean | 5000.50000 | 300.004930 | 310.005560 | 1538.776100 | 39.986910 | 107.951000 | 0.033900 | 0.004600 | 0.011500 | 0.009500 | |
| std | 2886.89568 | 2.000259 | 1.483734 | 179.284096 | 9.968934 | 63.654147 | 0.180981 | 0.067671 | 0.106625 | 0.097009 | |
| min | 1.00000 | 295.300000 | 305.700000 | 1168.000000 | 3.800000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| 25% | 2500.75000 | 298.300000 | 308.800000 | 1423.000000 | 33.200000 | 53.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| 50% | 5000.50000 | 300.100000 | 310.100000 | 1503.000000 | 40.100000 | 108.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| 75% | 7500.25000 | 301.500000 | 311.100000 | 1612.000000 | 46.800000 | 162.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| max | 10000.00000 | 304.500000 | 313.800000 | 2886.000000 | 76.600000 | 253.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | |

**EXPLORATORY DATA ANALYSIS:**

```
In [8]:  ► od=data.iloc[:,:-5]
            od
```

Out[8]:

| | UDI | Product ID | Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Failure Type | Machine failure |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | M14860 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | No Failure | 0 |
| 1 | 2 | L47181 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | No Failure | 0 |
| 2 | 3 | L47182 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | No Failure | 0 |
| 3 | 4 | L47183 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | No Failure | 0 |
| 4 | 5 | L47184 | L | 298.2 | 308.7 | 1408 | 40.0 | 9 | No Failure | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 9995 | 9996 | M24855 | M | 298.8 | 308.4 | 1604 | 29.5 | 14 | No Failure | 0 |
| 9996 | 9997 | H39410 | H | 298.9 | 308.4 | 1632 | 31.8 | 17 | No Failure | 0 |
| 9997 | 9998 | M24857 | M | 299.0 | 308.6 | 1645 | 33.4 | 22 | No Failure | 0 |
| 9998 | 9999 | H39412 | H | 299.0 | 308.7 | 1408 | 48.5 | 25 | No Failure | 0 |
| 9999 | 10000 | M24859 | M | 299.0 | 308.7 | 1500 | 40.2 | 30 | No Failure | 0 |

10000 rows × 10 columns

69

```
In [9]:  ▶ #we have to drop the columns that are unique identifiers
          od = od.drop("UDI", axis=1, errors="ignore")
          od = od.drop("Product ID", axis=1, errors="ignore")
          od = od.drop("Type", axis=1, errors="ignore")
          od.head(10)
```

Out[9]:

| | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Failure Type | Machine failure |
|---|---|---|---|---|---|---|---|
| 0 | 298.1 | 308.6 | 1551 | 42.8 | 0 | No Failure | 0 |
| 1 | 298.2 | 308.7 | 1408 | 46.3 | 3 | No Failure | 0 |
| 2 | 298.1 | 308.5 | 1498 | 49.4 | 5 | No Failure | 0 |
| 3 | 298.2 | 308.6 | 1433 | 39.5 | 7 | No Failure | 0 |
| 4 | 298.2 | 308.7 | 1408 | 40.0 | 9 | No Failure | 0 |
| 5 | 298.1 | 308.6 | 1425 | 41.9 | 11 | No Failure | 0 |
| 6 | 298.1 | 308.6 | 1558 | 42.4 | 14 | No Failure | 0 |
| 7 | 298.1 | 308.6 | 1527 | 40.2 | 16 | No Failure | 0 |
| 8 | 298.3 | 308.7 | 1667 | 28.6 | 18 | No Failure | 0 |
| 9 | 298.5 | 309.0 | 1741 | 28.0 | 21 | No Failure | 0 |

```
In [10]:  ▶ od['Failure Type'].unique()
```

```
Out[10]: array(['No Failure', 'Power Failure', 'Tool Wear Failure',
                'Overstrain Failure', 'Random Failures',
                'Heat Dissipation Failure'], dtype=object)
```

Below is the final dataset after preprocessing. We are going forward with multiclass classification instead of a normal binary classification predicting the actual reason for failure instead of just the failure. 1)air temp [K]- generated using a random walk process later normalized to a standard deviation of 2 K around 300 K 2)process temperature[K] - generated using a random walk process normalized to a standard deviation of 1 K, added to the air temperature plus 10 K. 3)rotational speed[rpm] - calculated from powepower of 2860 W, overlaid with a normally distributed noise 4)torque [Nm]: torque values are normally distributed around 40 Nm with an $lf$ = 10 Nm and no negative values. 5)tool wear [min]: The quality variants H/M/L add 5/3/2 minutes of tool wear to the used tool in the process. 6)Failure Type : Type of Failure the machine has faced['No Failure', 'Power Failure', 'Tool Wear Failure','Overstrain Failure', 'Random Failures','Heat Dissipation Failure']

```
In [11]:  ▶ od = od.drop("Machine failure", axis=1, errors="ignore")
          od
```

Out[11]:

| | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Failure Type |
|---|---|---|---|---|---|---|
| 0 | 298.1 | 308.6 | 1551 | 42.8 | 0 | No Failure |
| 1 | 298.2 | 308.7 | 1408 | 46.3 | 3 | No Failure |
| 2 | 298.1 | 308.5 | 1498 | 49.4 | 5 | No Failure |
| 3 | 298.2 | 308.6 | 1433 | 39.5 | 7 | No Failure |
| 4 | 298.2 | 308.7 | 1408 | 40.0 | 9 | No Failure |
| ... | ... | ... | ... | ... | ... | ... |
| 9995 | 298.8 | 308.4 | 1604 | 29.5 | 14 | No Failure |
| 9996 | 298.9 | 308.4 | 1632 | 31.8 | 17 | No Failure |
| 9997 | 299.0 | 308.6 | 1645 | 33.4 | 22 | No Failure |
| 9998 | 299.0 | 308.7 | 1408 | 48.5 | 25 | No Failure |
| 9999 | 299.0 | 308.7 | 1500 | 40.2 | 30 | No Failure |

10000 rows × 6 columns

Now lets start with the EDA process!

```
In [12]:  ▶ # obtaining the count for each type of failure
          od["Failure Type"].value_counts()
```

```
Out[12]: Failure Type
         No Failure                  9652
         Heat Dissipation Failure     112
         Power Failure                 95
         Overstrain Failure            78
         Tool Wear Failure             45
         Random Failures               18
         Name: count, dtype: int64
```

70

## CODE FOR VISUALIZATION:

```
In [13]:  ▶|  import seaborn as sns
              import matplotlib.pyplot as plt
```

```
In [14]:  ▶|  ax = sns.countplot(data=od, x="Failure Type")
              ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
```

```
Out[14]:  [Text(0, 0, 'No Failure'),
           Text(1, 0, 'Power Failure'),
           Text(2, 0, 'Tool Wear Failure'),
           Text(3, 0, 'Overstrain Failure'),
           Text(4, 0, 'Random Failures'),
           Text(5, 0, 'Heat Dissipation Failure')]
```

```
In [15]:  ▶|  """"the no failure seems to overwhelm the data thoroughly lets drop it
              and visulaize the remaining categories for the analysis"""
              ax = sns.countplot(data=od[od["Failure Type"] != "No Failure"], x="Failure Type")
              ax.bar_label(ax.containers[0], fmt='%.1f')
              ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
```

```
Out[15]:  [Text(0, 0, 'Power Failure'),
           Text(1, 0, 'Tool Wear Failure'),
           Text(2, 0, 'Overstrain Failure'),
           Text(3, 0, 'Random Failures'),
           Text(4, 0, 'Heat Dissipation Failure')]
```

```
In [16]:  ▶|  #lets dive more deep into the columns and their average
              df_multi_grouped = od.groupby("Failure Type").mean()
              print(df_multi_grouped.head(10))
```

```
                          Air temperature [K]  Process temperature [K]  \
Failure Type
Heat Dissipation Failure           302.567857               310.799107
No Failure                         299.972855               309.994343
Overstrain Failure                 299.867949               310.051282
Power Failure                      300.075789               309.954737
Random Failures                    300.766667               310.755556
Tool Wear Failure                  300.288889               310.164444

                          Rotational speed [rpm]  Torque [Nm]  Tool wear [min]
Failure Type
Heat Dissipation Failure             1337.964286    52.778571        107.339286
No Failure                           1540.324389    39.624316        106.678927
Overstrain Failure                   1354.243590    56.878205        208.217949
Power Failure                        1763.968421    48.514737        101.884211
Random Failures                      1489.444444    43.522222        119.888889
Tool Wear Failure                    1570.666667    37.226667        216.555556
```

```
In [17]:  ▶|  %matplotlib inline
              import matplotlib.pyplot as plt
              import seaborn as sns

              # Increase figure size if needed
              fig, axes = plt.subplots(3, 2, figsize=(14, 12))

              g1 = sns.boxplot(data=od, y="Air temperature [K]", x="Failure Type", ax=axes[0, 0])
              g2 = sns.boxplot(data=od, y="Process temperature [K]", x="Failure Type", ax=axes[0, 1])
              g3 = sns.boxplot(data=od, y="Rotational speed [rpm]", x="Failure Type", ax=axes[1, 0])
              g4 = sns.boxplot(data=od, y="Torque [Nm]", x="Failure Type", ax=axes[1, 1])
              g5 = sns.boxplot(data=od, y="Tool wear [min]", x="Failure Type", ax=axes[2, 0])

              # Rotate x-axis labels
              for ax in axes.flatten():
                  ax.tick_params(axis='x', labelrotation=45)

              g1.set_title("Air Temperature")
              g2.set_title("Process Temperature")
              g3.set_title("Rotational speed")
              g4.set_title("Torque")
              g5.set_title("Tool wear")

              plt.tight_layout()
              plt.show()
```

# MULTICLASS CLASSIFICATION:

```
In [18]:  ▶  from sklearn.model_selection import train_test_split
             X_multi = od.drop("Failure Type", axis=1)
             y_multi = od["Failure Type"]
             X_multi_train, X_multi_test, y_multi_train, y_multi_test = train_test_split(X_multi, y_multi, test_size=0.2, random_state=42)
             X_multi_train.shape, X_multi_test.shape, y_multi_train.shape, y_multi_test.shape
```

Out[18]: ((8000, 5), (2000, 5), (8000,), (2000,))

```
In [19]:  ▶  from sklearn.preprocessing import OneHotEncoder

             ohe = OneHotEncoder()
             # Use this
             y_multi_train = ohe.fit_transform(y_multi_train.values.reshape(-1, 1))
             y_multi_test = ohe.transform(y_multi_test.values.reshape(-1, 1))
```

```
In [20]:  ▶  y_multi_train = pd.DataFrame(y_multi_train.toarray(), columns=ohe.categories_)
             y_multi_test = pd.DataFrame(y_multi_test.toarray(), columns=ohe.categories_)
             y_multi_train.head(10)
```

Out[20]:

|   | Heat Dissipation Failure | No Failure | Overstrain Failure | Power Failure | Random Failures | Tool Wear Failure |
|---|---|---|---|---|---|---|
| 0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 5 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 6 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 7 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 8 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 9 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |

```
In [21]:  ▶  y_multi_train.describe()
```

Out[21]:

|  | Heat Dissipation Failure | No Failure | Overstrain Failure | Power Failure | Random Failures | Tool Wear Failure |
|---|---|---|---|---|---|---|
| count | 8000.000000 | 8000.000000 | 8000.000000 | 8000.000000 | 8000.000000 | 8000.000000 |
| mean | 0.012125 | 0.964625 | 0.008125 | 0.009375 | 0.001500 | 0.004250 |
| std | 0.109451 | 0.184737 | 0.089777 | 0.096376 | 0.038703 | 0.065057 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 50% | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 75% | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| max | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

```
In [22]:  ▶  y_multi_train[y_multi_train > 0.0].count()
```

Out[22]:  Heat Dissipation Failure    97
          No Failure               7717
          Overstrain Failure         65
          Power Failure              75
          Random Failures            12
          Tool Wear Failure          34
          dtype: int64

SCALING

```
In [23]:  ▶  X_multi_train.describe()
```

Out[23]:

|  | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] |
|---|---|---|---|---|---|
| count | 8000.000000 | 8000.000000 | 8000.000000 | 8000.000000 | 8000.000000 |
| mean | 300.009925 | 310.006088 | 1538.776625 | 40.007800 | 108.068750 |
| std | 2.002224 | 1.486484 | 180.594083 | 9.979156 | 63.225052 |
| min | 295.300000 | 305.700000 | 1181.000000 | 3.800000 | 0.000000 |
| 25% | 298.300000 | 308.800000 | 1423.000000 | 33.200000 | 53.000000 |
| 50% | 300.100000 | 310.100000 | 1502.000000 | 40.200000 | 108.000000 |
| 75% | 301.500000 | 311.100000 | 1611.250000 | 46.800000 | 162.000000 |
| max | 304.500000 | 313.800000 | 2886.000000 | 76.200000 | 253.000000 |

72

# MODEL TRAINING:

```
In [24]:  ▶  from sklearn.preprocessing import StandardScaler

             to_scale = ["Air temperature [K]", "Process temperature [K]", "Rotational speed [rpm]", "Torque [Nm]", "Tool wear [min]"]

             sc = StandardScaler()
             # Instead of this
             sc.fit(X_multi_train[to_scale])
             X_multi_train[to_scale] = sc.transform(X_multi_train[to_scale])
             X_multi_test[to_scale] = sc.transform(X_multi_test[to_scale])
```

```
In [25]:  ▶  X_multi_train.describe()
```

Out[25]:

|       | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] |
|-------|---------------------|-------------------------|------------------------|-------------|-----------------|
| count | 8.000000e+03        | 8.000000e+03            | 8.000000e+03           | 8.000000e+03 | 8.000000e+03   |
| mean  | -3.342171e-14       | -2.343170e-14           | 1.811884e-16           | -3.126388e-16 | 9.414691e-17  |
| std   | 1.000063e+00        | 1.000063e+00            | 1.000063e+00           | 1.000063e+00 | 1.000063e+00   |
| min   | -2.352493e+00       | -2.897008e+00           | -1.981233e+00          | -3.628570e+00 | -1.709378e+00 |
| 25%   | -8.540660e-01       | -8.114199e-01           | -6.411277e-01          | -6.822446e-01 | -8.710502e-01 |
| 50%   | 4.499028e-02        | 6.318154e-02            | -2.036552e-01          | 1.926135e-02 | -1.087453e-03  |
| 75%   | 7.442563e-01        | 7.359519e-01            | 4.013305e-01           | 6.806813e-01 | 8.530578e-01   |
| max   | 2.242683e+00        | 2.552432e+00            | 7.460419e+00           | 3.627006e+00 | 2.292451e+00   |

MODEL TRAINING

```
In [26]:  ▶  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

             def fit_model(model, X_train, X_test, y_train, y_test):
                 model.fit(X_train, y_train)
                 y_pred = model.predict(X_test)

                 accuracy = accuracy_score(y_test, y_pred)
                 precision = precision_score(y_test, y_pred, average="micro")
                 recall = recall_score(y_test, y_pred, average="micro")
                 f1 = f1_score(y_test, y_pred, average="micro")

                 return accuracy, precision, recall, f1
```

```
In [27]:  ▶  from sklearn.neighbors import KNeighborsClassifier
             from sklearn.tree import DecisionTreeClassifier
             from sklearn.ensemble import RandomForestClassifier

             knn = KNeighborsClassifier()
             dtc = DecisionTreeClassifier(random_state=21)
             rfc = RandomForestClassifier(random_state=21)

             result = {}
             for model,name in zip([knn, dtc, rfc],
                                   ["KNN", "Decision Tree", "Random Forest"]):
                 result[name] = fit_model(model, X_multi_train, X_multi_test, y_multi_train, y_multi_test)
```

```
In [28]:  ▶  import numpy as np
```

Accuracy:

Accuracy measures the proportion of correctly classified instances out of the total instances and is given by the formula: Number of Correct Predictions/Total Number of Predictions

Precision (Micro-Average):

Precision measures the ability of the classifier not to label as positive a sample that is negative. The micro-average precision is computed across all classes as a single precision value.

Recall (Micro-Average):

Recall, also known as sensitivity or true positive rate, measures the ability of the classifier to identify all relevant instances. The micro-average recall is computed across all classes as a single recall value.

F1-Score (Micro-Average):

F1-Score is the harmonic mean of precision and recall. The micro-average F1-score is computed across all classes as a single F1 value.

## COMPARISON:

```
In [29]:    df_result = pd.DataFrame(np.array(list(result.values())),
                                     columns= ["ACCURACY", "PRECISION", "RECALL", "F1-SCORE"],
                                     index= result.keys())

            df_result.index.name = "Model"    # name the index of the result1 dataframe as 'Model'
            df_result
```

Out[29]:

| Model | ACCURACY | PRECISION | RECALL | F1-SCORE |
|---|---|---|---|---|
| KNN | 0.9705 | 0.976358 | 0.9705 | 0.973420 |
| Decision Tree | 0.9710 | 0.971000 | 0.9710 | 0.971000 |
| Random Forest | 0.9795 | 0.982940 | 0.9795 | 0.981217 |

the best and efficient model is random forest

```
In [30]:    from sklearn.model_selection import GridSearchCV

            grid_params = {
                "max_depth": [2, 4, 6, 8, 10, 12, 15, 20],
                "min_samples_split": [2, 3, 4, 5],
                "min_samples_leaf": [2, 3, 4, 5, 6]
            }
            gs = GridSearchCV(dtc, grid_params, cv=5)
            print(fit_model(gs, X_multi_train, X_multi_test, y_multi_train, y_multi_test))
            print(f"Best model: {gs.best_estimator_}")

            (0.973, 0.9754385964912281, 0.973, 0.9742177722152691)
            Best model: DecisionTreeClassifier(max_depth=10, min_samples_leaf=5, random_state=21)
```

```
In [31]:    from sklearn.tree import DecisionTreeRegressor
            final_model = DecisionTreeRegressor(max_depth=10, random_state=21)
            final_model.fit(X_multi_train, y_multi_train)
```

Out[31]:

```
            ▼              DecisionTreeRegressor
            DecisionTreeRegressor(max_depth=10, random_state=21)
```

## PREDICTION:

```
In [32]:    from sklearn.ensemble import BaggingClassifier
            from sklearn.neighbors import KNeighborsClassifier
            from sklearn.tree import DecisionTreeClassifier
            from sklearn.ensemble import RandomForestClassifier
            from sklearn.metrics import accuracy_score
            from sklearn.model_selection import train_test_split

            # Assuming you have your X and y data ready
            X_train, X_test, y_train, y_test = train_test_split(X_multi, y_multi, test_size=0.2, random_state=42)

            # Instantiate the base classifiers
            knn = KNeighborsClassifier()
            dtc = DecisionTreeClassifier(random_state=21)
            rfc = RandomForestClassifier(random_state=21)

            # Create the BaggingClassifier
            bagging_knn = BaggingClassifier(knn, n_estimators=10, random_state=42)
            bagging_dtc = BaggingClassifier(dtc, n_estimators=10, random_state=42)
            bagging_rfc = BaggingClassifier(rfc, n_estimators=10, random_state=42)

            # Train the BaggingClassifier
            bagging_knn.fit(X_train, y_train)
            bagging_dtc.fit(X_train, y_train)
            bagging_rfc.fit(X_train, y_train)

            # Make predictions
            y_pred_knn = bagging_knn.predict(X_test)
            y_pred_dtc = bagging_dtc.predict(X_test)
            y_pred_rfc = bagging_rfc.predict(X_test)

            # Evaluate the BaggingClassifiers
            accuracy_knn = accuracy_score(y_test, y_pred_knn)
            accuracy_dtc = accuracy_score(y_test, y_pred_dtc)
            accuracy_rfc = accuracy_score(y_test, y_pred_rfc)

            print("Bagging KNN Accuracy:", accuracy_knn)
            print("Bagging Decision Tree Accuracy:", accuracy_dtc)
            print("Bagging Random Forest Accuracy:", accuracy_rfc)

            Bagging KNN Accuracy: 0.9675
            Bagging Decision Tree Accuracy: 0.982
            Bagging Random Forest Accuracy: 0.9785
```

```
In [33]:   import warnings
           warnings.filterwarnings("ignore", category=UserWarning)

           def predict_failure(air_temperature, process_temperature, rotational_speed, torque, tool_wear, model):
               # Create a feature vector
               user_input = [[air_temperature, process_temperature, rotational_speed, torque, tool_wear]]

               # Make prediction using the provided model
               prediction = model.predict(user_input)

               return prediction

           # Get input from the user
           air_temperature = float(input("Enter Air temperature [K]: "))
           process_temperature = float(input("Enter Process temperature [K]: "))
           rotational_speed = float(input("Enter Rotational speed [rpm]: "))
           torque = float(input("Enter Torque [Nm]: "))
           tool_wear = float(input("Enter Tool wear [min]: "))


           # Predict using the bagging_dtc model
           predicted_failure_dtc = predict_failure(air_temperature, process_temperature, rotational_speed, torque, tool_wear, bagging_dt
           print("Predicted Failure Type :", predicted_failure_dtc[0])
```

```
Enter Air temperature [K]: 200
Enter Process temperature [K]: 300
Enter Rotational speed [rpm]: 1029
Enter Torque [Nm]: 3
Enter Tool wear [min]: 4
Predicted Failure Type : Power Failure
```

```
In [33]:   import warnings
           warnings.filterwarnings("ignore", category=UserWarning)

           def predict_failure(air_temperature, process_temperature, rotational_speed, torque, tool_wear, model):
               # Get input from the user
```

75

# APPENDIX - SOURCE CODE

**DATA CLEANING:**

```python
import pandas as pd
data=pd.read_csv(r'C:\Users\nazia\Desktop\CAPSTONE\maintenance.csv
)
Data
data.info()
pd.isnull(data)
pd.isnull(data).sum()
is_duplicate_any_column = data.duplicated().any()
if is_duplicate_any_column:
    print('There are duplicate values in at least one column.')
else:
    print('There are no duplicate values in any column.')
columns_with_missing = data.columns[data.isnull().any()].tolist()
if columns_with_missing:
    print("Columns with missing values:", columns_with_missing)
else:
    print("There are no missing values in any column.")
data.describe()
```

**EXPLORATORY DATA ANALYSIS:**

```python
od=data.iloc[:,:-5]
od
#we have to drop the columns that are unique identifiers
od = od.drop("UDI", axis=1, errors="ignore")
od = od.drop("Product ID", axis=1, errors="ignore")
od = od.drop("Type", axis=1, errors="ignore")
od.head(10)
od['Failure Type'].unique()
od = od.drop("Machine failure", axis=1, errors="ignore")
# obtaining the count for each type of failure
od["Failure Type"].value_counts()
```

**VISUALIZATION:**

```python
import seaborn as sns
import matplotlib.pyplot as plt
ax = sns.countplot(data=od, x="Failure Type")
ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
"the no failure seems to overwhelm the data thoroughly lets drop
it
and visulaize the remaining categories for the analysis"""
ax = sns.countplot(data=od[od["Failure Type"] != "No Failure"],
x="Failure Type")
ax.bar_label(ax.containers[0], fmt='%.1f')
ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
#lets dive more deep into the columns and their average
df_multi_grouped = od.groupby("Failure Type").mean()
print(df_multi_grouped.head(10))
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
# Increase figure size if needed
fig, axes = plt.subplots(3, 2, figsize=(14, 12))
g1 = sns.boxplot(data=od, y="Air temperature [K]", x="Failure
Type", ax=axes[0, 0])
g2 = sns.boxplot(data=od, y="Process temperature [K]", x="Failure
Type", ax=axes[0, 1])
g3 = sns.boxplot(data=od, y="Rotational speed [rpm]", x="Failure
Type", ax=axes[1, 0])
g4 = sns.boxplot(data=od, y="Torque [Nm]", x="Failure Type",
ax=axes[1, 1])
g5 = sns.boxplot(data=od, y="Tool wear [min]", x="Failure Type",
ax=axes[2, 0])
# Rotate x-axis labels
for ax in axes.flatten():
    ax.tick_params(axis='x', labelrotation=45)
g1.set_title("Air Temperature")
```

```
g2.set_title("Process Temperature")
g3.set_title("Rotational speed")
g4.set_title("Torque")
g5.set_title("Tool wear")
plt.tight_layout()
plt.show()
```

**MULTICLASS CLASSIFICATION:**
```
from sklearn.model_selection import train_test_split
X_multi = od.drop("Failure Type", axis=1)
y_multi = od["Failure Type"]
X_multi_train, X_multi_test, y_multi_train, y_multi_test =
train_test_split(X_multi, y_multi, test_size=0.2, random_state=42)
X_multi_train.shape, X_multi_test.shape, y_multi_train.shape,
y_multi_test.shape
from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder()
# Use this
y_multi_train = ohe.fit_transform(y_multi_train.values.reshape(-1,
1))
y_multi_test = ohe.transform(y_multi_test.values.reshape(-1, 1))
y_multi_train = pd.DataFrame(y_multi_train.toarray(),
columns=ohe.categories_)
y_multi_test = pd.DataFrame(y_multi_test.toarray(),
columns=ohe.categories_)
y_multi_train.head(10)
y_multi_train.describe()
y_multi_train[y_multi_train > 0.0].count()
```

**SCALING:**
```
X_multi_train.describe()
from sklearn.preprocessing import StandardScaler
to_scale = ["Air temperature [K]", "Process temperature [K]",
"Rotational speed [rpm]", "Torque [Nm]", "Tool wear [min]"]
```

```python
sc = StandardScaler()
# Instead of this
sc.fit(X_multi_train[to_scale])
X_multi_train[to_scale] = sc.transform(X_multi_train[to_scale])
X_multi_test[to_scale] = sc.transform(X_multi_test[to_scale])
X_multi_train.describe()
```

**MODEL TRAINING:**

```python
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score
def fit_model(model, X_train, X_test, y_train, y_test):
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred, average="micro")
    recall = recall_score(y_test, y_pred, average="micro")
    f1 = f1_score(y_test, y_pred, average="micro")
    return accuracy, precision, recall, f1
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
knn = KNeighborsClassifier()
dtc = DecisionTreeClassifier(random_state=21)
rfc = RandomForestClassifier(random_state=21)
result = {}
for model,name in zip([knn, dtc, rfc],
                      ["KNN", "Decision Tree", "Random Forest"]):
    result[name] = fit_model(model, X_multi_train, X_multi_test,
y_multi_train, y_multi_test)
```

**COMPARISION:**

```python
df_result = pd.DataFrame(np.array(list(result.values())),
                         columns= ["ACCURACY", "PRECISION",
"RECALL", "F1-SCORE"],
                         index= result.keys())
```

```python
df_result.index.name = "Model"    # name the index of the result1
dataframe as 'Model'
df_result
from sklearn.model_selection import GridSearchCV

grid_params = {
    "max_depth": [2, 4, 6, 8, 10, 12, 15, 20],
    "min_samples_split": [2, 3, 4, 5],
    "min_samples_leaf": [2, 3, 4, 5, 6]
}
gs = GridSearchCV(dtc, grid_params, cv=5)
print(fit_model(gs, X_multi_train, X_multi_test, y_multi_train,
y_multi_test))
print(f"Best model: {gs.best_estimator_}")
from sklearn.tree import DecisionTreeRegressor
final_model = DecisionTreeRegressor(max_depth=10, random_state=21)
final_model.fit(X_multi_train, y_multi_train)
from sklearn.ensemble import BaggingClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
# Assuming you have your X and y data ready
X_train, X_test, y_train, y_test = train_test_split(X_multi,
y_multi, test_size=0.2, random_state=42)
# Instantiate the base classifiers
knn = KNeighborsClassifier()
dtc = DecisionTreeClassifier(random_state=21)
rfc = RandomForestClassifier(random_state=21)
# Create the BaggingClassifier
bagging_knn = BaggingClassifier(knn, n_estimators=10,
random_state=42)
bagging_dtc = BaggingClassifier(dtc, n_estimators=10,
random_state=42)
```

```python
bagging_rfc = BaggingClassifier(rfc, n_estimators=10,
random_state=42)
# Train the BaggingClassifier
bagging_knn.fit(X_train, y_train)
bagging_dtc.fit(X_train, y_train)
bagging_rfc.fit(X_train, y_train)
# Make predictions
y_pred_knn = bagging_knn.predict(X_test)
y_pred_dtc = bagging_dtc.predict(X_test)
y_pred_rfc = bagging_rfc.predict(X_test)
# Evaluate the BaggingClassifiers
accuracy_knn = accuracy_score(y_test, y_pred_knn)
accuracy_dtc = accuracy_score(y_test, y_pred_dtc)
accuracy_rfc = accuracy_score(y_test, y_pred_rfc)
print("Bagging KNN Accuracy:", accuracy_knn)
print("Bagging Decision Tree Accuracy:", accuracy_dtc)
print("Bagging Random Forest Accuracy:", accuracy_rfc)
```

**PREDICTION:**
```python
import warnings
warnings.filterwarnings("ignore", category=UserWarning)
def predict_failure(air_temperature, process_temperature,
rotational_speed, torque, tool_wear, model):
    # Create a feature vector
    user_input = [[air_temperature, process_temperature,
rotational_speed, torque, tool_wear]]
    # Make prediction using the provided model
    prediction = model.predict(user_input)
    return prediction
# Get input from the user
air_temperature = float(input("Enter Air temperature [K]: "))
process_temperature = float(input("Enter Process temperature [K]:
"))
rotational_speed = float(input("Enter Rotational speed [rpm]: "))
```

```python
torque = float(input("Enter Torque [Nm]: "))
tool_wear = float(input("Enter Tool wear [min]: "))
# Predict using the bagging_dtc model
predicted_failure_dtc = predict_failure(air_temperature,
process_temperature, rotational_speed, torque, tool_wear,
bagging_dtc)
print("Predicted Failure Type :",predicted_failure_dtc[0])
```