



MÄLARDALEN UNIVERSITY  
SCHOOL OF INNOVATION, DESIGN AND ENGINEERING  
VÄSTERÅS, SWEDEN

---

Thesis for the Degree of Bachelor in Computer Science 15.0 credits

# **MACHINE LEARNING BASED PREDICTIVE DATA ANALYTICS FOR EMBEDDED TEST SYSTEMS**

Fayad Al Hanash  
*Fah20002@student.mdu.se*

Examiner: Mobyen Uddin Ahmed  
Mälardalen University, Västerås, Sweden

Supervisor: Shahina Begum, Md Alamgir Kabir & Md Rakibul Islam  
Mälardalen University, Västerås, Sweden

Company Supervisor: Daniel Bengtsson  
Nordic Engineering Partner, Stockholm, Sweden

29/09/2023

## Abstract

Organizations gather enormous amounts of data and analyze these data to extract insights that can be useful for them and help them to make better decisions. Predictive data analytics is a crucial subfield within data analytics that make accurate predictions. Predictive data analytics extracts insights from data by using machine learning algorithms. This thesis presents the supervised learning algorithm to perform predictive data analytics in Embedded Test System at the Nordic Engineering Partner company. Predictive Maintenance is a concept that is often used in manufacturing industries which refers to predicting asset failures before they occur. The machine learning algorithms used in this thesis are support vector machines, multi-layer perceptrons, random forests, and gradient boosting. Both binary and multi-class classifier have been provided to fit the models, and cross-validation, sampling techniques, and a confusion matrix have been provided to accurately measure their performance. In addition to accuracy, recall, precision, f1, kappa, mcc, and roc auc measurements are used as well. The prediction models that are fitted achieve high accuracy.

## Acknowledgements

I would like to thank my supervisors Shahina Begum, Md Alamgir Kabir & Md Rakibul Islam for their support and advice during the thesis. I would also like to thank Daniel Bengtsson for sharing the ETS data and knowledge that allowed me to work on the thesis. Finally, I would like to thank my family for their encouragement and support during my studies.

# Table of Contents

|   |           |
|---|-----------|
| <b>1. Introduction .....</b>                          | <b>6</b>  |
| <b>2. Background .....</b>                            | <b>6</b>  |
| 2.1. <i>Embedded Test System.....</i>                 | 6         |
| 2.2. <i>Machine Learning.....</i>                     | 7         |
| 2.2.1. <i>Supervised Learning .....</i>               | 7         |
| 2.2.2. <i>Machine Learning Algorithms.....</i>        | 7         |
| 2.3. <i>Predictive Data Analytics.....</i>            | 10        |
| 2.4. <i>Predictive Maintenance.....</i>               | 11        |
| 2.4.1. <i>Machine learning for PdM.....</i>           | 12        |
| 2.5. <i>PdM based ML for ETS.....</i>                 | 12        |
| <b>3. Related Work .....</b>                          | <b>12</b> |
| <b>4. Problem Formulation .....</b>                   | <b>13</b> |
| 4.1. <i>Research questions.....</i>                   | 13        |
| 4.2. <i>Limitations.....</i>                          | 13        |
| <b>5. Method .....</b>                                | <b>13</b> |
| 5.1. <i>Research Methodology.....</i>                 | 13        |
| <b>6. Ethical and Societal Considerations.....</b>    | <b>13</b> |
| <b>7. Experimental study and implementation .....</b> | <b>14</b> |
| 7.1. <i>Data collecting.....</i>                      | 15        |
| 7.2. <i>Data preparing .....</i>                      | 15        |
| 7.3. <i>Normalization.....</i>                        | 17        |
| 7.4. <i>Modeling.....</i>                             | 17        |
| 7.5. <i>Evaluation.....</i>                           | 18        |
| 7.5.1. <i>Cross Validation.....</i>                   | 18        |
| 7.5.2. <i>Data balancing .....</i>                    | 18        |
| 7.5.3. <i>Performance Measures .....</i>              | 19        |
| <b>8. Results.....</b>                                | <b>20</b> |
| 8.1. <i>Binary Classification.....</i>                | 20        |
| 8.1.1. <i>Imbalanced dataset .....</i>                | 20        |
| 8.1.2. <i>Balanced dataset.....</i>                   | 22        |
| 8.2. <i>Multi-class Classification.....</i>           | 23        |
| 8.2.1. <i>Imbalanced dataset .....</i>                | 23        |
| 8.2.2. <i>Balanced dataset.....</i>                   | 25        |

---

|   |           |
|---|-----------|
| <b>9. Discussion .....</b>                        | <b>26</b> |
| <b>10. Conclusions.....</b>                       | <b>28</b> |
| <b>11. Future Work.....</b>                       | <b>28</b> |
| <b>References.....</b>                            | <b>29</b> |
| <b>Appendices .....</b>                           | <b>32</b> |
| <i>Appendix A: Feature extraction .....</i>       | <i>32</i> |
| <i>Appendix B: Modeling &amp; Evaluation.....</i> | <i>40</i> |

## List of Figures

|  |    |
|--|----|
| Figure 1 Embedded Test System.....   | 7  |
| Figure 2 Support Vectors Classification [15].....  | 8  |
| Figure 3 Random Forest Classification [18] .....   | 9  |
| Figure 4 Multi-Layer Perceptron [19].....  | 9  |
| Figure 5 Gradient Boosting [24] .....  | 10 |
| Figure 6 A diagram of the CRISP-DM [1].....  | 11 |
| Figure 7 Experimental flow diagram.....  | 14 |
| Figure 8 Distribution of the classes .....   | 18 |
| Figure 9 Confusion matrix of the binary classifiers for the unbalanced dataset.....        | 21 |
| Figure 10 Roc curve for the binary classifiers (unbalanced).....                           | 21 |
| Figure 11 Confusion matrix of the binary classifiers for the balanced dataset.....         | 22 |
| Figure 12 Roc curve for the binary classifiers (balanced) .....                            | 23 |
| Figure 13 Confusion matrix of the multi-class classifiers for the unbalanced dataset ..... | 24 |
| Figure 14 Confusion matrix of the multi-class classifiers for the balanced dataset .....   | 25 |
| Figure 15 Distribution of the scaled features in the training set .....                    | 27 |
| Figure 16 Correlation matrix of the training dataset .....                                 | 28 |

## List of Tables

|  |    |
|--|----|
| Table 1 Details of test result features .....  | 15 |
| Table 2 Details of measurement result features .....   | 15 |
| Table 3 Details of extracted features .....  | 16 |
| Table 4 The selected features .....  | 17 |
| Table 5 Confusion matrix.....  | 19 |
| Table 6 Accuracy, Precision, Recall, F1, Kappa, MCC, and ROC_AUC measurements of the binary classifiers for the unbalanced dataset ..... | 21 |
| Table 7 Cross validation scores of the binary classifiers.....   | 22 |
| Table 8 Accuracy, Precision, Recall, F1, Kappa, MCC, ROC_AUC measurements of the binary classifiers for the balanced dataset .....       | 22 |
| Table 9 Accuracy, Precision, Recall, and F1 measurements of the multi-class SVM for the unbalanced dataset.....                          | 24 |
| Table 10 Accuracy, Precision, Recall, and F1 measurements of the multi-class RF for the unbalanced dataset.....                          | 24 |
| Table 11 Accuracy, Precision, Recall, and F1 measurements of the multi-class MLP for the unbalanced dataset.....                         | 24 |

Table 12 Accuracy, Precision, Recall, and F1 measurements of the multi-class GB for the unbalanced dataset.....25

Table 13 Cross validation scores of the multi-class classifiers .....26

Table 14 Accuracy, Precision, Recall, and F1 measurements of the multi-class SVM for the balanced dataset.....26

Table 15 Accuracy, Precision, Recall, and F1 measurements of the multi-class RF for the balanced dataset.....26

Table 16 Accuracy, Precision, Recall, and F1 measurements of the multi-class MLP for the balanced dataset.....26

Table 17 Accuracy, Precision, Recall, and F1 measurements of the multi-class GB for the balanced dataset.....26

## List of abbreviations

|                         |  |
|-------------------------|--|
| <i>AI</i>               | <i>Artificial Intelligence</i>                                 |
| <i>ANN</i>              | Artificial Neural Networks                                     |
| <i>CM</i>               | Confusion Matrix   |
| <i>CRISP-DM</i>         | Cross Industry Standard Process for Data Mining                |
| <i>CV</i>               | Cross Validation   |
| <i>EOL</i>              | End of Line  |
| <i>EQ</i>               | Measurement equal to the low limit                             |
| <i>ETS</i>              | Embedded Test System   |
| <i>FN</i>               | False Negative   |
| <i>FP</i>               | False Positive   |
| <i>FPR</i>              | False Positive Rate  |
| <i>GB</i>               | Gradient Boosting  |
| <i>GELE</i>             | Measurement $\geq$ low limit and measurement $\leq$ high limit |
| <i>GTLT</i>             | Measurement $>$ low limit and measurement $<$ high limit       |
| <i>ICT</i>              | In-Circuit Test  |
| <i>LINQ</i>             | Language-Integrated Query                                      |
| <i>MCC</i>              | Matthews Correlation Coefficient                               |
| <i>ML</i>               | Machine Learning   |
| <i>MLP</i>              | Multi-Layer Perceptron   |
| <i>NEP</i>              | Nordic Engineering Partner                                     |
| <i>NULL</i>             | No comparison  |
| <i>PCB</i>              | Printed Circuit Board  |
| <i>PCBA</i>             | Printed Circuit Board Assembly                                 |
| <i>PdM</i>              | Predictive Maintenance   |
| <i>PvM</i>              | Preventive Maintenance   |
| <i>R2F</i>              | Run-to-failure   |
| <i>RF</i>               | Random Forest  |
| <i>ROC</i>              | Receiver Operating Characteristic                              |
| <i>ROC AUC</i>          | Area Under the Receiver Operating Characteristic Curve         |
| <i>SMOTE</i>            | Synthetic Minority Oversampling Technique                      |
| <i>SVM</i>              | Support Vector Machines  |
| <i>TN</i>               | True Negative  |
| <i>TP</i>               | True Positive  |
| <i>TPR</i>              | True Positive Rate   |
| <i>V</i>                | Volt   |
| $\text{\AA}\mu\text{A}$ | Microampere  |

## 1. Introduction

Nowadays, organizations gather enormous amounts of data and analyze these data to extract insights that can be useful for them and help them to make better decisions [1]. Predictive data analytics is a crucial subfield within data analytics that make accurate predictions. Predictive data analytics extracts insights from data by using machine learning algorithms. There are four types of machine learning algorithms such as supervised learning, unsupervised learning, semi-supervised learning, and reinforcement [1]. This thesis presents the supervised learning algorithm to perform predicative data analytics in Embedded Test System (ETS) at the Nordic Engineering Partner company (NEP). The increasing availability of data and the increasing capability of hardware have encouraged industries to apply machine learning (ML) approaches in maintenance management domains [2]. The concept Predictive Maintenance (PdM) is often used in manufacturing industries which refers to predicting the asset failures before they occur [3]. There are many advantages to using predictive maintenance in production environments, but there are also challenges to be overcome [4]. The PdM increases productivity and efficiency and decreases the system faults and unplanned downtimes [4]. The necessity of integrating data from diverse systems and sources within a facility to gather accurate data for creating prediction models is a challenge that has to be overcome [4]. Additionally, the use of artificial intelligence (AI) raises several challenges such as “obtaining training data, dealing with dynamic operating environments, selecting the ML algorithm that better fits to a given scenario, and the necessity of context-aware information, such as operational conditions, and production environment” [4]. Furthermore, the use of ML for predictive data analytics faces several challenges as well, such as extracting and selecting the features from the database and fitting and developing the models to predict the failure of ETS. In this thesis, the implementation of a ML-based application is carried out using the experimental and Cross Industry Standard Process for Data Mining (CRISP-DM) methodologies [1].

## 2. Background

### 2.1. Embedded Test System

The Nordic Engineering Partner is a company that develops and assembles advanced In-Circuit Test (ICT) fixtures executing through End of Line (EOL) tests on Printed Circuit Board Assemblies (PCBAs), these fixtures are called Embedded Test System. ICT is a type of electrical test and used to ensure that the PCB is functioning correctly and identify any potential defects. ICT uses a bed-of-nails test fixture that makes electrical contact to every net on the PCB [5]. EOL testing is a crucial part in ensuring the quality and reliability of PCBAs and represents the final functional check of PCBA before it is shipped to the customers [6]. PCBA is the process of attaching electronic components to a substrate that is printed with a particular circuit. PCBA ensures components are electrically interconnected and allows them to communicate with each other [7]. The Embedded Test System is shown in Figure 1 and is an “all-in-one solution” for testing electronic products. The ETS consist of a test fixture, a computer with operators panel and tester specific boards such as power supply, input/output channels etc.<sup>1</sup> ETS fixtures are hardware and software tools that are designed to automate tests on the PCBA, such as functional tests, stress tests, and environmental tests. ETS fixtures improve the quality and reliability of the PCBA by increasing test efficiency, reducing human error, and improving test coverage and accuracy.

---

<sup>1</sup> Embedded Test System. <https://nepartner.se/products/embedded-test-system/>



Figure 1 Embedded Test System

## 2.2. Machine Learning

Arthur Samuel coined the term “Machine Learning” in 1959 in the context of letting the machine solve the game of checkers. Machine learning refers to a computer program that has the ability to learn in order to produce behavior that was not explicitly programmed by the programmer, and it can display behavior that the programmer might not even be aware of [8][9]. The program learns the behaviors based on the following factors: the consumed data by the program, a metric that measures the error or the deviation between the ideal behavior and current behavior, and a feedback mechanism that utilizes the measured error to guide the program to produce better behavior in future occurrences [8]. ML is a subset of artificial Intelligence [10]. In contrast to AI-Applications, ML involves learning of hidden patterns within the data and uses these patterns for prediction or classification [10]. Since making inferences from samples is the core task of ML, the ML builds statical models by using the theory of statistics [11].

### 2.2.1. Supervised Learning

There are four types of machine learning algorithms such as supervised learning, unsupervised learning, semi-supervised learning, and reinforcement [1]. Supervised learning simply refers to learning from examples. Supervised learning techniques aim to learn a model using a training set of input-output pairs and later use this model to make predictions (predicate outputs) on a test set. The inputs are also called descriptive features and the output is called a target feature or label [1][12][13].

Let  $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$  represent a set of training examples where:

- $x_i$  is a vector or tuple containing n numbers of descriptive features.
- $y_i$  is a target feature.
- $(x_i, y_i)$  is a pair of  $x_i$  and its corresponding  $y_i$ .

More specifically, the dataset represents a table in a relational database, and each row in the table is called a record or an example and represents such  $(x_i, y_i)$ . In addition, each column in the table is called an attribute or a feature and represent such  $x_i$  or  $y_i$ . The test set is represented in the same way as the training set, but target feature  $y_i$  is not known. By applying the supervised learning,  $y_i$  can be predicted.

### 2.2.2. Machine Learning Algorithms

ML provides the knowledge that intelligent machines need to maintain and sustain their functionality. Essentially, ML algorithms are embedded into machines and data flows to extract knowledge and feed it into the system for managing the processes more rapidly and efficiently [10]. Machine learning algorithms utilize the patterns and relationships existing in massive amounts of multidimensional data to automatically learn how to anticipate outcomes, classify data, and perform other tasks. They gain knowledge from a training set of examples, and then they apply that knowledge to any new input [14]. There are a lot of ML algorithms, but in this thesis, support vector machines (SVM), random forest (RF), Multilayer perceptron (MLP) and gradient boosting (GB) were chosen to perform the work. The model’s implementation was carried out using the Scikit library. Scikit is a library in Python that supports the implementation of those algorithms<sup>2</sup>.

---

<sup>2</sup> <https://scikit-learn.org/>



### 2.2.2.1. Support Vector Machines

The support vector machine algorithm is designed to solve binary classification problems. The idea is to use a hyperplane that separates two classes and maximizes the margin between the classes' closest vectors [8] [15]. As shown in Figure 2 below, the hyperplane is the middle of the margin, and the support vectors are the points on the boundary [15]. When it is difficult to find a hyperplane between the two classes, the kernel approach can be used to separate the classes by projecting the support vectors into a higher-dimensional space [15]. SVM can be used to solve multi-class classification problems by separating the classes into several binary classifiers and solving the corresponding problem. This can be achieved via two techniques:

- One-against-all:
 

The one-against-all technique splits the dataset into  $n$ -classes, then applies the SVM algorithm for each class against all other classes. This technique creates  $n$ -numbers of SVM [16]. For example, the ETS's dataset consists of four classes: passed, failed, error, and terminated. By applying this technique, four SVM binary classifications can be obtained, as follows:

  - Passed vs (failed, error, and terminated).
  - Failed vs (passed, error, and terminated).
  - Error vs (passed, failed, and terminated).
  - Terminated vs (passed, failed, and error).
  
- One-against-one:
 

The one-against-one technique splits the dataset into  $n$ -classes, then applies the SVM algorithm for one class against each other classes. This technique creates  $(n(n-1)/2)$  numbers of SVMs [16]. By applying this technique to the ETS's dataset, six SVM binary classifications can be obtained, as follows:

  - Passed vs failed.
  - Passed vs error.
  - Passed vs terminated.
  - Failed vs error.
  - Failed vs terminated.
  - Error vs terminated.

Scikit provides a great support for developing SVM<sup>3</sup>. Scikit uses the one-against-one technique to model SVM for multi-class strategies. The One-against-all matrix is constructed from the one-against-one matrix.

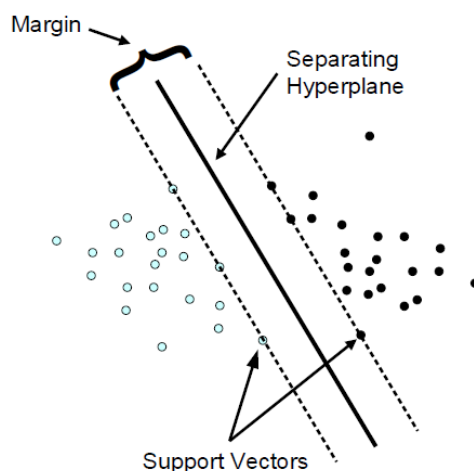
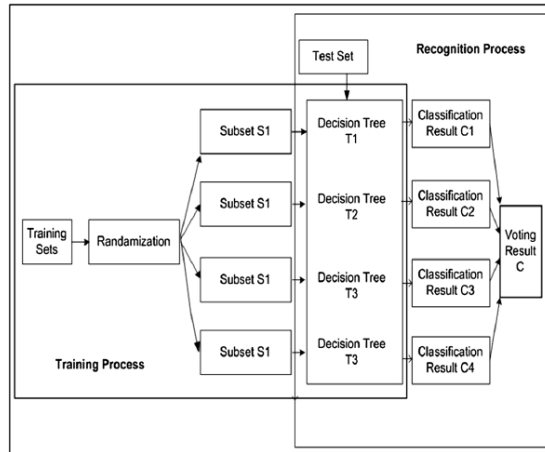


Figure 2 Support Vectors Classification [15]

<sup>3</sup> <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>

**2.2.2.2. Random Forest**

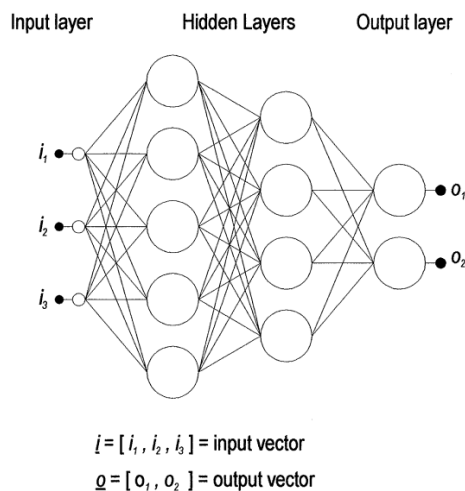
Random forest is a tree-based ensemble such that each tree depends on a set of random variables. RF can be used for classification purposes and is quite fast in training and predicting [17]. The approach, which creates many decision tree classifiers on various samples of the dataset while training the model, predicts the target feature on each classifier under the testing phase, then predicts the final target feature by voting for the majority of predicted target features. Figure 3 shows the process of random forest classification [18].



**Figure 3 Random Forest Classification [18]**

**2.2.2.3. Multi-Layer Perceptron**

Multi-Layer Perceptron is one of the most popular artificial neural networks (ANN) for solving classification problems. MLP is a feedforward network that consists of multiple layers of perceptrons, the input layer, the output layer, and hidden layers between the input and output layers. Figure 4 illustrates the MLP [19]. The size of hidden layers doesn't depend on the input and output. Each layer has several nodes or neurons that connect all layers with each other. MLP is described as fully connected, in which each layer must be connected to every node in the previous layer and every node in the next layer [8] [19] [20]. MLP is a supervised type and uses the backpropagation algorithm to train the model. The backpropagation algorithm works as follows: randomly initialize network weights; propagate the input through the network to get the output; calculate the error between the target output and the actual output; propagate the error back through the network; update the weights to reduce the error; repeat the steps (without initializing the weights again) until the error is as small as possible [19].



**Figure 4 Multi-Layer Perceptron [19]**

#### 2.2.2.4. Gradient Boosting

Boosting is an effective technique introduced in the domain of machine learning and statistics for solving classification problems [21] [22]. The concept behind boosting is the combination of simple classifiers that a weak or base learner obtains to improve prediction performance instead of using a simple classifier alone. A weak learner is a learning algorithm that produces classifiers with better performance than random guessing [17]. Boosting is a type of supervised learning that aims to find an optimal function of variables to predict the output by reducing the loss function [21]. One of the earliest boosting algorithms suggested is AdaBoost, which was developed to solve binary classification problems and then extended to solve multi-class classification as well. AdaBoost is an algorithm of type gradient descent, thus stating boosting as a statistical estimation and numerical optimization. Gradient boosting is an extension of the AdaBoost algorithm utilizing a statistical framework [8] [23]. Gradient boosting uses decision trees as weak or base learners [23]. Figure 5 illustrates the GB [24].

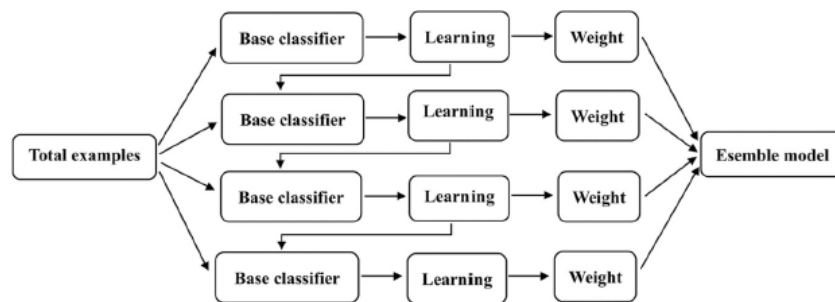


Figure 5 Gradient Boosting [24]

### 2.3. Predictive Data Analytics

Predictive data analytics is the technique of creating and utilizing models that make accurate predictions. Predictive data analytics analyses and extracts insights from current and previous data by using AI, ML, and statistics [1][11]. CRISP-DM is one of the most widely used methodologies for predictive data analytics. The CRISP-DM is preferred by data analytics practitioners because it is non-proprietary, explicitly includes both application-focused and technical perspectives in the data analytics process, and is neutral regarding application, industry, and tools [1][25]. Figure 6 illustrates the CRISP-DM. The phases of CRISP-DM are as follows:

- Business Understanding: addressing and determining business problems such as identifying the goal of building a prediction model and initializing and developing a plan.
- Data Understanding: understanding the available data sources and what they contain. In addition, addressing the data requirements such as the collection, description, exploration, and quality verification of the data.
- Data Preparation: preparing the data such as selecting the desired data, cleaning it, and organizing it into a specific structure.
- Modeling: building a predictive model by using different ML algorithms and selecting the best one.
- Evaluation: evaluating and testing the model to ensure that it fits the purpose (makes accurate predictions).
- Deployment: deploying and integrating the evaluated model into the process within an organization.

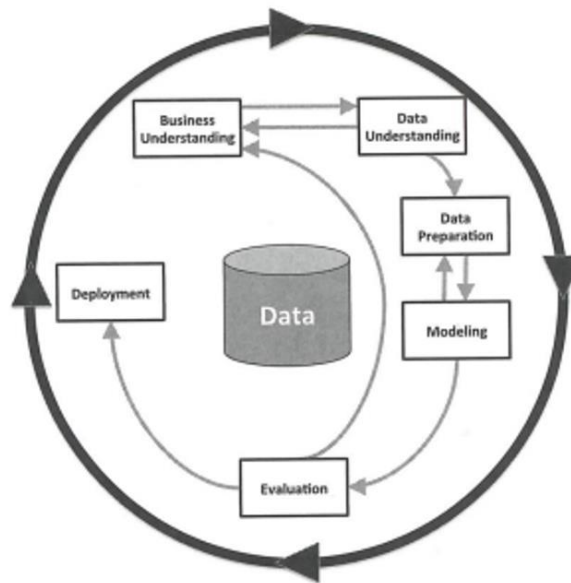


Figure 6 A diagram of the CRISP-DM [1]

## 2.4. Predictive Maintenance

Due to the increased availability of data, the industries started to apply a new approach called predictive maintenance to improve their products. This concept refers to the idea that each production component can inform when it will fail [26] [27]. According to a common estimate [3], downtime costs almost every company between 5% and 20% of its production, and maintenance costs can be enormous. By reducing downtime, industry can become more efficient [3]. In addition, industries can provide efficient solutions for managing the maintenance and improving quality by applying machine learning techniques. In order of rising complexity and efficiency, the maintenance management approaches have been categorized as follow:

- Run-to-failure (R2F): is reactive maintenance, which the maintenance is performed after the occurrence of failures and issues. This technique for maintenance is the easiest and the least efficient because the maintenance actions will be applied after the breakdown occurs and the cost will be higher than planned actions [26][27].
- Preventive maintenance (PvM): The maintenance is performed before the occurrence of failures by using a planned schedule that is “based on time or process iterations”. Failures are usually prevented using this technique, also known as scheduled maintenance. This technique may increase costs by unnecessarily maintaining a component that still works and has not failed [26][27].
- Predicative maintenance (PdM): Also known as “condition-based maintenance” or “on-line monitoring” or “risk-based maintenance” [28]. The maintenance is performed before the occurrence of failures by utilizing and monitoring the health status of the equipment. PdM applies statistical and machine learning techniques to automatically collect and analyze equipment data to predict failures [26][27]. With the help of the integrated sensors, PdM can reduce breakdowns and costs, increase efficiency, address different issues related to equipment replacement, and identify the underlying cause of the failure [28]. Both PvM and PdM may act in the same way to schedule the maintenance activities before failure occurrences, but the PdM approach is based on data gathered from sensors and analysis algorithms [28].

#### 2.4.1. Machine learning for PdM

ML-based PdM categorizes into two classes: supervised and unsupervised. In supervised class, the modeling dataset contains information about failure occurrences, while in unsupervised class, the modeling dataset contains information about the logistic and/or process but doesn't contain information about maintenance [27]. The availability of maintenance data depends on the type of existing maintenance management approach. However, in the run-to-fail approach, the supervised technique is suitable due to the availability of maintenance data such as the production process between two failure occurrences. On the other hand, in the preventive maintenance approach, the unsupervised technique is suitable because the failure maintenance is performed before actual failures occur, which leads to a lack of full availability of the maintenance data [27]. Due to the availability of datasets provided by NEP, the supervised approach is considered in this thesis. Two classes of supervised problems can exist depending on the type of outcome: a classification problem if the outcome is considered to have categorical values and a regression problem if the outcome is considered to have continuous values [27].

#### 2.5. PdM based ML for ETS

The NEP is trying to apply the ML algorithms within ETS to develop a system that can monitor the electrical performance of the test fixtures during testing and identify any changes in the test results that may indicate a problem with the bed-of-nails pins. For example, if the electrical resistance of the pins increases over time, it may indicate that the pins are becoming worn or contaminated and may need to be cleaned or replaced. This thesis can provide the essential steps to achieve the proposed system since the limitations of this thesis are limited to only classification purposes. To fully achieve the proposed system, the thesis can be extended for regression purposes as well. Furthermore, this thesis presents suitable ML algorithms that can be applied. Moreover, this thesis can provide a tool that is ready to extract the features, tune the parameters, fit ML models, evaluate them and plot the results. The tool can predict and classify the status of new ETS tests as well.

### 3. Related Work

The increasing availability of data and the increasing capability of hardware have recently encouraged industries to apply ML approaches in maintenance management domains. The application of ML approaches in maintenance management has made PdM a popular topic in manufacturing research [28][2]. PdM based ML applications have been a potential tool in preventing and predicting equipment failures. The performance of PdM applications is tied to the selection of the ML algorithm [29]. The research presented by Franco and Figueiredo [26], develops a PdM system in a mobile edge computing architecture to combine the low cost of embedded systems and the advantages of security with the needs of PdM. The research presented shows that the developed system provides an accuracy of 100%, whereas the ML algorithm was RF. According to Dey et al. [30], the performance of classification tasks that work online will be required by intelligent embedded systems in the future. SVM is one of the different classification techniques that has received a lot of interest. SVM appears to be a better option than traditional neural networks and provides "remarkable results both in classification and regression applications" [30]. Additionally, Chigurupati, Thibaux and Lassar [14] illustrate and propose that the SVM is a crucial algorithm for achieving very accurate prediction. The study developed by Castanheira et al. [31] shows that MLP is used to classify the incipient errors in power transformers. The provided procedure was evaluated using data collected from chromatographic tests of the transformers. The accuracy of MLP is between 75 and 90%. Gkamas et al. [32] proposed a ML framework based on a MLP classifier that can solve classification problems, analyze big data in a relative amount of time, is suitable for low-resource embedded systems, and is a utility for crucial event detection and industrial maintenance. The proposed framework can even be trained online. According to Carvalho et al. [29], research suggests a fault detection system that detects real-time faults on the hard disk drive. The suggested system is divided into two phases: batch training, where RF is used to train models, and real-time prediction, where the estimations are performed using data collected from the end-user device. However, the accuracy of this system is 85%. Furthermore, Carvalho et al. argue that RF is the most used and compared machine learning algorithm in predictive maintenance applications and the key reason is that RF is based on decision trees, which generate a huge number of observations that can be used in prediction. According to Liulys [22], recent developments in statistical and computer science domains have led to the manufacturing industry focusing on increasing the data from industrial

repositories in a sustainable manner. An exciting development is machine learning, where the GB is one of the most popular machine learning algorithms for analyzing big datasets.

## 4. Problem Formulation

The increasing complexity of ETS has made it difficult for engineers to ensure their reliability and performance. To address this issue, ML has been proposed as a potential solution for predicative data analytics in ETS. However, the adoption of ML-based predictive analytics in ETS is limited due to the lack of a systematic approach to implement ML-based predictive analytics in ETS and the lack of relevant research. A PdM-based ML approach allows engineers to ensure the performance of the ETS to maintain and replace pins frequently to prevent and avoid unexpected pin failures. This helps in reducing costs and unexpected failures, thus increasing the efficiency and reliability of ETS. This work aims to develop a systematic approach to implementing ML-based predictive analytics in ETS that can facilitate NEP's engineers' ability to predict the potential failures of ETS.

### 4.1. Research questions

1. How to extract and select the key Features from ETS dataset?
2. How to implement ML model with optimize the parameters to predict ETS failures?
3. How to develop a framework to be evaluated using real-world ETS dataset?

### 4.2. Limitations

The limitation of this thesis is that the dataset used in the implementation is limited to one year of data collection, which may not be enough to study all possible ETS tests. Additionally, this thesis is limited to the classification and prediction of the ETS results.

## 5. Method

### 5.1. Research Methodology

The main objective of the thesis was to apply ML algorithms to real-world ETS data to perform predictive analysis on the data that could facilitate NEP's engineers' ability to predict the potential failures on ETS. However, the developed approach could extract the features, prepare them for modeling, tuning the parameters and train the models and evaluate them. An experimental approach was carried out to address the research questions. The methodological steps were applied as follows:

- A brief understanding of the concepts, technologies, and processes of both ETS and PCBA was made such as understanding the test fixture, electronic components, test objects, and interactions between them.
- A literature study was carried out to review the existing literature and research in the domains of predictive data analytics, predictive maintenance, and embedded systems as well.
- A current state-of-the-art study within the machine learning domain was also carried out to identify ML algorithms that fit the application.
- An experimental implementation of the work was conducted as follows:
  - Collecting and preparing the data.
  - Medeling different ML-algorithms.
  - Evaluating and testing the performance of the built models.
- The conclusion was drawn and discussed.

## 6. Ethical and Societal Considerations

There are ethical considerations to consider while conducting the thesis. The data that is provided by NEP is confidential and considered to be used exclusively at Mälardalen University campuses and only for research purposes. A non-disclosure agreement has been signed with the NEP, which includes the

data sharing agreement and publication regulation. The data will be deleted as soon as the work is finished.

## 7. Experimental study and implementation

The experimental work was carried out following the CRISP-DM approach [1]. Figure 7 illustrates the core implementation of the work. The implementation consists of three phases: the data preparation phase, the modeling phase, and the evaluation phase. In the data preparation phase, feature extraction and normalization of the features have been carried out; more details provided in sections 7.1, 7.2, and 7.3. In the modeling phase, the steps are carried out as follows: 1. splitting the data set into an 80 percent train dataset and a 20 percent test dataset; 2. utilizing the train dataset to train the model; 3. selecting the type of classifier whether a binary classifier or a multi-class classifier; 4. selecting the classifier, such as SVM, RF, MLP, or GB; 5. selecting whether the dataset should be imbalanced or balanced. If the selected dataset is balanced, then resampling and cross validating the features, otherwise training the model with an imbalanced dataset. More detail about the modeling phase is described in section 7.4. In the evaluating phase, the evaluation of the trained model is carried out by using the test dataset to predict the outcomes and comparing them with the actual outcomes. More details are described in section 7.5.

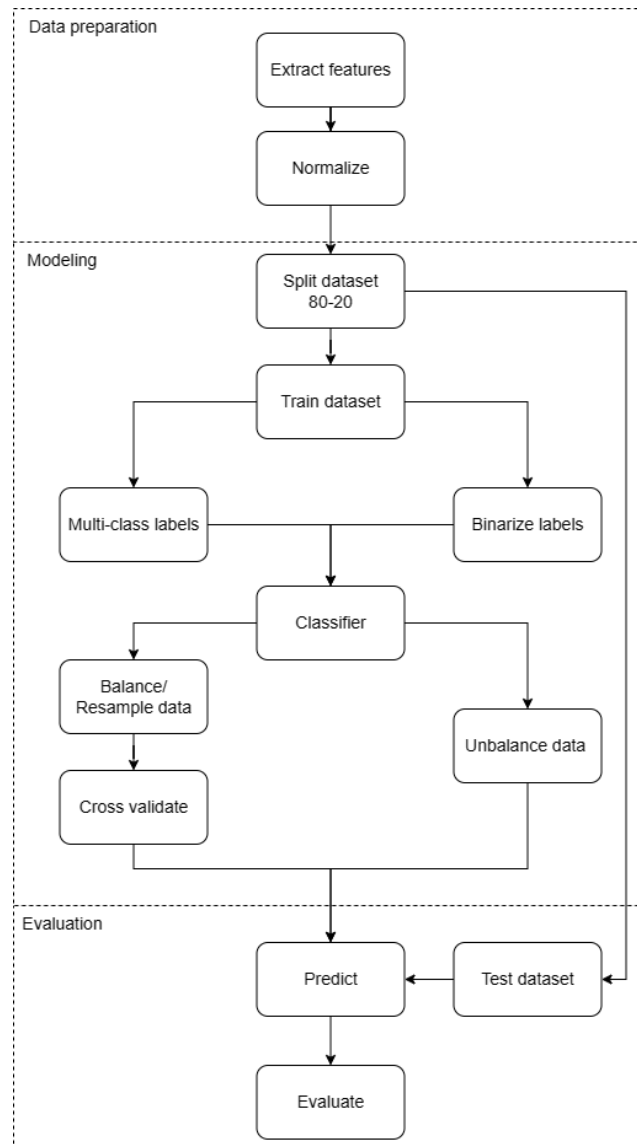


Figure 7 Experimental flow diagram

## 7.1. Data collecting

The dataset gathered from NEP and compiled from ETS tests that contains details about the EOL tests of PCBA. The dataset consists of two csv files: test result and measurement result. The test result provides information on complete tests, namely, each record in the test file is a complete test sequence of an EOL test of a single PCBA, while the measurement test provides in detail how a test sequence is processed. As shown in Tables 1 and 2, the test result has 4 features and contains 42 992 records, while the measurement has 11 features and contains 1 473 741 records.

**Table 1 Details of test result features**

| <i>Header</i>         | <i>Description</i>                                       |
|-----------------------|--|
| <i>Test result ID</i> | The ID of a test sequence                                |
| <i>Timestamp</i>      | The execution time and data for each sequence            |
| <i>Execution time</i> | The execution time of a test sequence                    |
| <i>status</i>         | The status of a test (passed, failed, error, terminated) |

**Table 2 Details of measurement result features**

| <i>Header</i>                | <i>Description</i>   |
|------------------------------|--|
| <i>Measurement result ID</i> | The measurement ID for each test in a test sequence          |
| <i>Test result ID</i>        | The ID of a test sequence                                    |
| <i>Name</i>                  | The name of each test in a test sequence                     |
| <i>Timestamp</i>             | The execution time and data for each test in a test sequence |
| <i>Execution time</i>        | The execution time of a test sequence                        |
| <i>Numeric time</i>          | The value of a measurement test                              |
| <i>Numeric low limit</i>     | The lowest value accepted                                    |
| <i>Numeric high limit</i>    | The highest value accepted                                   |
| <i>Comparator</i>            | The comparator of a test                                     |
| <i>Unit</i>                  | The unit of the measured test                                |
| <i>Status</i>                | The status of a test (passed, failed, error, terminated)     |

## 7.2. Data preparing

Data preparation and feature extraction were carried out using the .NET<sup>4</sup> framework and the LINQ<sup>5</sup> library. .NET is a free open-source cross-platform to build modern applications. LINQ stands for language-integrated query and refers to a set of technologies based on the direct integration of query capabilities into the C# language. With LINQ, the dataset can be easily filtered, ordered, and grouped with a minimum of code. The code for feature extraction is provided in Appendix A and on GitHub<sup>6</sup> as well. The dataset was prepared into a format that could be used for training models. At the beginning, the two files have been merged into one file by including all features from the measurement file and replacing the status feature from the test file based on the test ID. Then a new dataset (format) was extracted, as shown in Table 3, some features were taken from the measurement file, and new features were added as follows:

- Grouping the tests where they had the same test ID.
- Calculating the minimum, maximum, average, median, and standard deviation of each execution time, numeric data, numeric low limit, and numeric high limit.
- Counting the frequency of comparators and units. The comparator and unit have different abbreviations, some of them are GELE, EQ, GTLT, LOG, V, ÅµA, and NULL.
- Counting the frequency of test names and measurement IDs.
- The missing values (NaN) have been replaced in the dataset.

<sup>4</sup> <https://dotnet.microsoft.com/en-us/>

<sup>5</sup> <https://learn.microsoft.com/en-us/dotnet/csharp/linq/>

<sup>6</sup> <https://github.com/FayadHanash/FeatureExtractionETS>



**Table 3 Details of extracted features**

| <i>Header</i>                         | <i>Description</i>   |
|---------------------------------------|--|
| <i>TestResultID</i>                   | The ID of a test sequence  |
| <i>Execution Time Total</i>           | The execution time for a test sequence                                   |
| <i>Execution Time Min</i>             | The minimum execution time for a test sequence                           |
| <i>Execution Time Max</i>             | The maximum execution time for a test sequence                           |
| <i>Execution Time Avg</i>             | The average execution time for a test sequence                           |
| <i>Execution Time Median</i>          | The median execution time for a test sequence                            |
| <i>Execution Time StdDev</i>          | The standard deviation execution time for a test sequence                |
| <i>Numeric Data Min</i>               | The minimum numeric data for a test sequence                             |
| <i>Numeric Data Max</i>               | The maximum numeric data for a test sequence                             |
| <i>Numeric Data Avg</i>               | The average numeric data for a test sequence                             |
| <i>Numeric Data Median</i>            | The median numeric data for a test sequence                              |
| <i>Numeric Data StdDev</i>            | The standard deviation numeric data for a test sequence                  |
| <i>Numeric Data Low Limit Min</i>     | The minimum numeric data low limit for a test sequence                   |
| <i>Numeric Data Low Limit Max</i>     | The maximum numeric data low limit for a test sequence                   |
| <i>Numeric Data Low Limit Avg</i>     | The average numeric data low limit for a test sequence                   |
| <i>Numeric Data Low Limit Median</i>  | The median numeric data low limit for a test sequence                    |
| <i>Numeric Data Low Limit StdDev</i>  | The standard deviation numeric data low limit for a test sequence        |
| <i>Numeric Data High Limit Min</i>    | The minimum numeric data high limit for a test sequence                  |
| <i>Numeric Data High Limit Max</i>    | The maximum numeric data high limit for a test sequence                  |
| <i>Numeric Data High Limit Avg</i>    | The average numeric data high limit for a test sequence                  |
| <i>Numeric Data High Limit Median</i> | The median numeric data high limit for a test sequence                   |
| <i>Numeric Data High Limit StdDev</i> | The standard deviation numeric data high limit for a test sequence       |
| <i>Name Frequency</i>                 | The frequency of the name for each test within a test sequence           |
| <i>Measurement ID Frequency</i>       | The frequency of the measurement ID for each test within a test sequence |
| <i>GELE Frequency</i>                 | The frequency of GELE for a test sequence                                |
| <i>EQ Frequency</i>                   | The frequency of EOF for a test sequence                                 |
| <i>LOG Frequency</i>                  | The frequency of LOG for a test sequence                                 |
| <i>GTLT Frequency</i>                 | The frequency of GTLT for a test sequence                                |
| <i>V Frequency</i>                    | The frequency of V for a test sequence                                   |
| <i>ÂµA Frequency</i>                  | The frequency of ÂµA for a test sequence                                 |
| <i>NULL Frequency</i>                 | The frequency of NULL for a test sequence                                |
| <i>Status</i>                         | The status for a test sequence   |

After exploring all extracted features, the most varied features were selected to train the model. The feature variety improves the performance of ML models and avoids data overfitting. As shown in Table 4, the selected features are the total, maximum, average, and standard deviation of the execution time and the minimum, average, median, and standard deviation of the numeric data. In addition, the status feature is selected since it is the target feature. The name of each test in a test sequence is not considered while only its frequency is considered, due to the fact that the names are the same and frequently repeated in each sequence, which does not make any sense for modeling's purpose to have them in the dataset. Furthermore, the measurement ID is not considered in the dataset since the measurement ID is frequently increased in each test. Moreover, the timestamp feature for each test in a sequence test is not considered because the execution time total already represent the total time for a test sequence. Finally, all other features that are presented in Table 3 and are not presented in Table 4 have been ignored since they do not vary and are almost the same in each test sequence. The dataset has been divided into two sets: one for the training phase, which includes 80% of records, and the second for the testing phase, which includes 20% of records.

**Table 4** The selected features

| <i>Header</i>                | <i>Description</i>  |
|------------------------------|---|
| <i>Execution Time Total</i>  | The execution time for a test sequence                    |
| <i>Execution Time Max</i>    | The maximum execution time for a test sequence            |
| <i>Execution Time Avg</i>    | The average execution time for a test sequence            |
| <i>Execution Time StdDev</i> | The standard deviation execution time for a test sequence |
| <i>Numeric Data Min</i>      | The minimum numeric data for a test sequence              |
| <i>Numeric Data Avg</i>      | The average numeric data for a test sequence              |
| <i>Numeric Data Median</i>   | The median numeric data for a test sequence               |
| <i>Numeric Data StdDev</i>   | The standard deviation numeric data for a test sequence   |
| <i>Status</i>                | The status for a test sequence                            |

### 7.3. Normalization

Normalization is a technique used to scale the values in a given range. It is helpful for classification and predication purposes and facilitates the modeling process for machine learning algorithms. Min-Max normalization, also known as range normalization, is a normalization technique that fits the original data into a specific range. The Min-Max normalization can be calculated by using the following equation:

$$a' = \frac{a_i - \min(a)}{\max(a) - \min(a)} * (high - low) + low \quad (1)$$

Where  $a'$  is the normalized value,  $a$  is the actual value,  $\min$  is the minimum value of  $a$ ,  $\max$  is the maximum value of  $a$ ,  $low$  and  $high$  are the boundaries of the desired range. [1] [33]. In the extracted features, the feature ranges were extremely high. The normalization phase was carried out using the Min-Max technique, and the features were scaled into the range [0,1].

### 7.4. Modeling

The modeling implementation was carried out using Python and libraries such as Scikit<sup>7</sup>, Imbalanced-learn<sup>8</sup>, Matplotlib<sup>9</sup>, NumPy<sup>10</sup>, Pandas<sup>11</sup> and Seaborn<sup>12</sup>. Scikit is a comprehensive library for predictive data analysis and machine learning. It provides a wide range of tools and algorithms for classification, model evaluation, etc. Imbalanced-learn is a library that provides re-sampling techniques. Matplotlib is a plotting and visualization library. Pandas is a library for data analysis and manipulation. Seaborn is a visualization library based on Matplotlib. The code for normalization, modeling, and evaluating is provided in Appendix B and on GitHub<sup>13</sup> as well. The final dataset contains 42 844 records, and after counting insight into the target set, the class that is “passed” has 39855 records, the class that is “failed” has 378, the class that is “error” has 2518 records, and the class that is “terminated” has 93 records. Since the class “passed” is significantly higher than the other classes, the modeling phase was trained in two different phases: the binary classification phase and the multi-class classification phase. In the binary classification phase, the target set has been binarized into two classes: the “passed” class and the “rest” class, while in the multi-class classification, all the classes have been kept. As illustrated in Figure 7, the developed framework allows to select either binary classification or multi-class classification, then select a classifier, and after that, select whether the dataset should be imbalanced or balanced. If the selected dataset is balanced, then the features will be resampled and cross-validated. Then, the modeling process continues fitting until it is done and evaluates the performance of the modeled classifier. Parameter

<sup>7</sup> <https://scikit-learn.org/>

<sup>8</sup> <https://imbalanced-learn.org/>

<sup>9</sup> <https://matplotlib.org/>

<sup>10</sup> <https://numpy.org/>

<sup>11</sup> <https://pandas.pydata.org/>

<sup>12</sup> <https://seaborn.pydata.org/>

<sup>13</sup> <https://github.com/FayadHanash/MachineLearningBasedPredictiveDataAnalytics>

tuning is carried out using the GridSearchCV<sup>14</sup> approach, which is a method provided by Scikit that allows to search for the best parameters that fit the model.

## 7.5. Evaluation

In the evaluation phase, the combination of cross validation, sampling techniques, and confusion matrix were applied to the dataset for determining the classifier performance.

### 7.5.1. Cross Validation

Cross validation (CV) is a way to ensure how a model performs on unseen data. It is a procedure for evaluating robust models. The repeated stratified k-fold cross validation has been used to validate the balanced dataset. The repeated stratified k-fold cross validation is a version of k-fold cross validation that uses stratified random sampling to create the folds. The dataset was split into five folds and repeated three times; at each repeat, the training set was split into one testing fold and four training folds. This method can yield correct values if it is applied efficiently [34].

### 7.5.2. Data balancing

A dataset is imbalanced where the distribution of classes is significantly skewed, which leads to the classification of the majority classes being more accurate better than the minority classes or the minority classes being ignored at all. A way to address this issue is by re-sampling the dataset. Over sampling technique duplicates the records in the minority classes and replicates them in the dataset. Under sampling technique delates the records from the majority classes. Synthetic Minority Oversampling Technique (SMOTE) is a technique that creates synthetic examples to over-sample the minority classes [35]. By combining the SOMTE and the under-sample, a classifier can perform better [35]. In this experiment, the data balancing was carried out by pipelining the SMOTE and Random Under Sampler techniques. Figure 8 shows the distribution of the classes before and after balancing the dataset in both binary and multi-class cases.

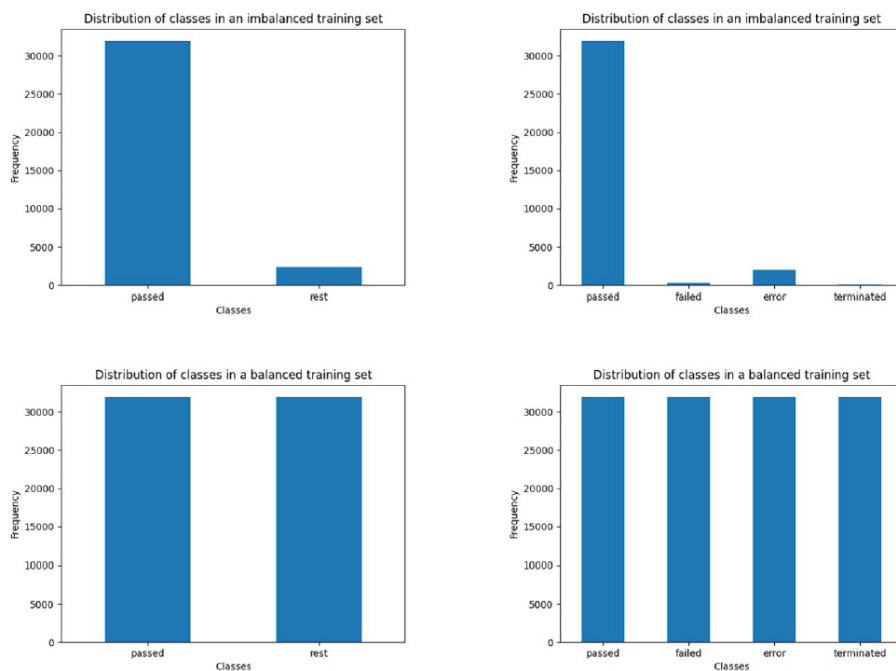


Figure 8 Distribution of the classes

<sup>14</sup> [https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)

### 7.5.3. Performance Measures

A confusion matrix (CM) was carried out to evaluate the performance of the algorithms. The CM as shown in Table 5 is a matrix containing the predicted class as columns and the actual class as rows.

**Table 5 Confusion matrix**

|                        | <i>Predicted Negative</i> | <i>Predicted Positive</i> |
|------------------------|---------------------------|---------------------------|
| <i>Actual Negative</i> | True Negatives            | False Positives           |
| <i>Actual Positive</i> | False Negatives           | True Positives            |

True Negatives (TN): represent the number of correctly classified negative examples. False Positives (FP): represent the number of falsely classified positive examples. False Negatives (FN): represent the number of falsely classified negative examples. True Positives (TP): represent the number of correctly classified positive examples [35]. Accuracy represents the proportion of correct prediction examples of the total prediction examples and measures as the ratio between sum of correct classification and total number of classifications. The accuracy can be calculated by using the following equation:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (2)$$

Recall represents the proportion of correctly predicted positive examples of all actual positive examples and measures as the ratio between of the true positive examples and the sum of true positive and false negative examples. The recall can be calculated by using the following equation:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

Precision represents the proportion of correctly predicted positive examples of all predicted positive examples and measures as the ratio of the true positive examples and the sum of true positive and false negative examples. The precision can be calculated by using the following equation:

$$Precision = \frac{TP}{TP + FP} \quad (4)$$

F1-score represents the harmonic mean of recall and precision. The f1-score can be calculated by using the following equation:

$$F1 = \frac{2 * precision * recall}{precision + recall} \quad (5)$$

The Kappa coefficient estimates the degree of agreement between a pair of raters. The kappa, also known as Cohen's kappa, was introduced by Jacob Cohen. The Kappa can be used as performance measure for both binary and multi-class classification. The value of kappa is a range between [-1, +1], where -1 indicates to wrong predication and +1 perfect predication [36]. The kappa can be calculated by using the following equation:

$$Kappa = \frac{2 * (TP * TN - FP * FN)}{(TP + FP) * (FP + TN) + (TP + FN) * (FN + TN)} \quad (6)$$

Matthews Correlation Coefficient (MCC), also known as the phi ( $\phi$ ) coefficient in 2 \* 2 confusion matrices [36]. MCC can be used to evaluate the performance of binary and multi-class classifications. The value of MCC is a range between [-1, +1], where -1 indicates to wrong predication and +1 perfect predication [36]. The MCC can be calculated by using the following equation:

$$MCC = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP) * (TP + FN) * (TN + FP) * (TN * FN)}} \quad (7)$$

The Receiver Operating Characteristic curve (ROC curve) is a plot that separates two classes in a binary classifier, such as a form of threshold. The ROC curve summarizes the classification performance between the true positive rate (TPR) and false positive rate (FPR) [35]. TPR and FPR are used to plot the ROC curve. The area under the ROC curve, often known as the ROC AUC, is a metric provided by the ROC curve, and it is a measurement that is considered the most essential property of the classifier. The larger the area, the better the performance of the classifier. TPR represents the same thing as Recall, also known as Sensitivity. FPR, also known as Specificity, represents the proportion of incorrectly predicted positive examples of all actual negative examples and measures as the ratio of the false positive examples and the sum of false positive and true negative examples [8]. The TPR and FPR can be calculated by using the following equations:

$$TPR = \frac{TP}{TP + FN} \quad (8)$$

$$FPR = \frac{FP}{FP + TN} \quad (9)$$

## 8. Results

The result of the experimental work is divided into two sections: binary classification and multi-class classification, and each section is divided into two subsections as well: imbalanced dataset and balanced dataset. The measurements that are used in the binary classifications to perform the evaluation are accuracy, precision, recall, f1-score, kappa, mcc, roc auc. In contrast, accuracy, precision, recall, and f1-score are used for the multi-class classification. Furthermore, the results of the CV are estimated only for the balanced dataset.

### 8.1. Binary Classification

#### 8.1.1. Imbalanced dataset

Figure 9 presents the CM of the SVM, RF, MLP, and GB classifiers, where the dataset was unbalanced and binarized to 'pass' and 'rest' classes. Table 6 shows the measurement as follows: The accuracy and precision of the classifiers are almost 100% where the data is unbalanced. The recall, F1-score, kappa, and MCC measurements are highest for the RF classifier and lowest for the SVM classifier. The ROC AUC measurement is highest for RF and GB at 100%, as plotted in Figure 10.

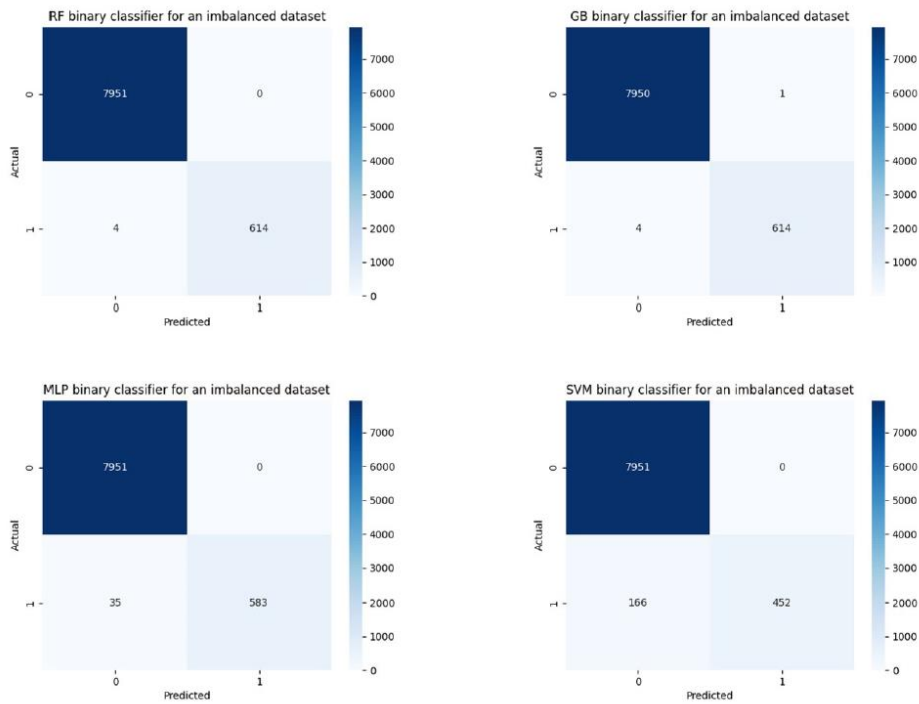


Figure 9 Confusion matrix of the binary classifiers for the unbalanced dataset

Table 6 Accuracy, Precision, Recall, F1, Kappa, MCC, and ROC\_AUC measurements of the binary classifiers for the unbalanced dataset

| Classifier Algorithm | Accuracy | Precision | Recall | F1-Score | Kappa | MCC   | ROC_AUC |
|----------------------|----------|-----------|--------|----------|-------|-------|---------|
| SVM                  | 98.06    | 1         | 0.73   | 0.84     | 0.83  | 0.85  | 0.97    |
| RF                   | 99.95    | 1         | 0.99   | 0.98     | 0.98  | 0.98  | 1       |
| MLP                  | 99.59    | 1         | 0.94   | 0.97     | 0.97  | 0.97  | 0.98    |
| GB                   | 99.94    | 0.998     | 0.994  | 0.996    | 0.996 | 0.996 | 1       |

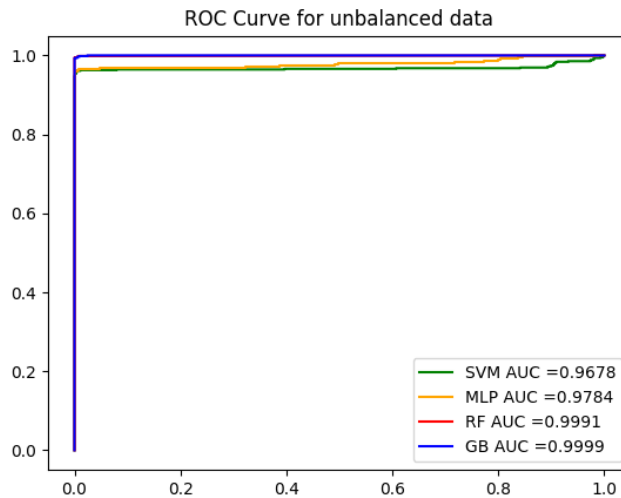
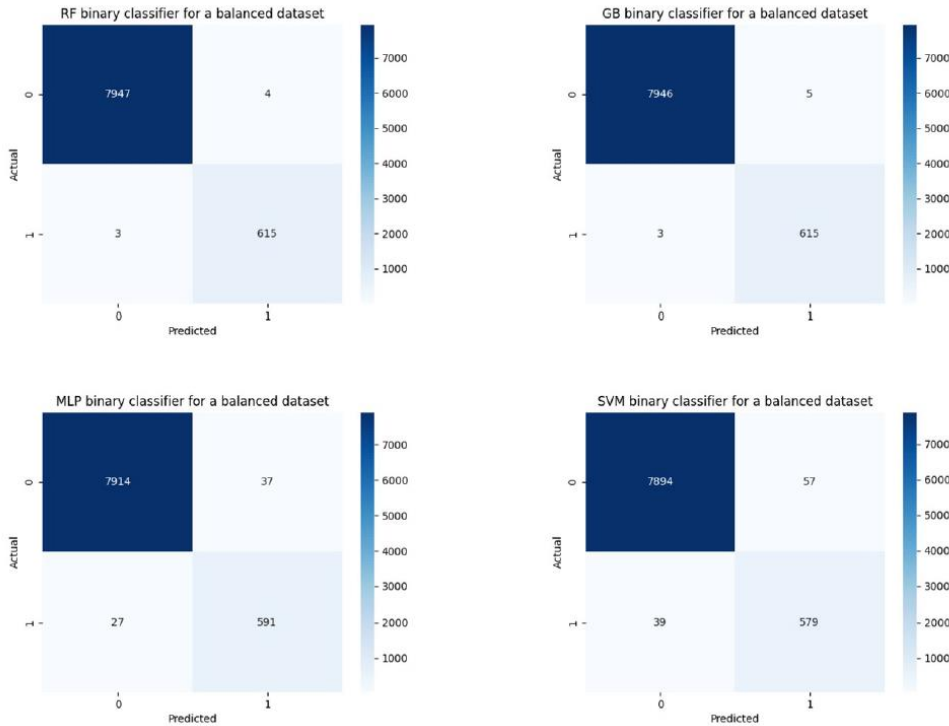


Figure 10 Roc curve for the binary classifiers (unbalanced)

**8.1.2. Balanced dataset**

Figure 11 presents the CM of the SVM, RF, MLP, and GB classifiers, where the dataset was balanced and binarized to ‘pass’ and ‘rest’ classes. Table 8 shows that almost all the measurements are equal. The ROC AUC measurement is highest for RF and GB at 100%, as plotted in Figure 12. Table 7 shows the average of the CV scores that are in the range [97%, 100%], and the best classifier was RF under training the model.



**Figure 11 Confusion matrix of the binary classifiers for the balanced dataset**

**Table 7 Cross validation scores of the binary classifiers**

| <i>Classifier Algorithm</i> | <i>Cross validation score</i> |
|-----------------------------|-------------------------------|
| <i>SVM</i>                  | 97.07                         |
| <i>RF</i>                   | 99.94                         |
| <i>MLP</i>                  | 98.44                         |
| <i>GB</i>                   | 99.90                         |

**Table 8 Accuracy, Precision, Recall, F1, Kappa, MCC, ROC\_AUC measurements of the binary classifiers for the balanced dataset**

| <i>Classifier Algorithm</i> | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> | <i>Kappa</i> | <i>MCC</i> | <i>ROC_AUC</i> |
|-----------------------------|-----------------|------------------|---------------|-----------------|--------------|------------|----------------|
| <i>SVM</i>                  | 98.88           | 0.91             | 0.937         | 0.923           | 0.917        | 0.918      | 0.971          |
| <i>RF</i>                   | 99.92           | 0.994            | 0.995         | 0.994           | 0.994        | 0.994      | 1              |
| <i>MLP</i>                  | 99.25           | 0.94             | 0.956         | 0.949           | 0.945        | 0.945      | 0.998          |
| <i>GB</i>                   | 99.91           | 0.992            | 0.995         | 0.994           | 0.993        | 0.993      | 1              |

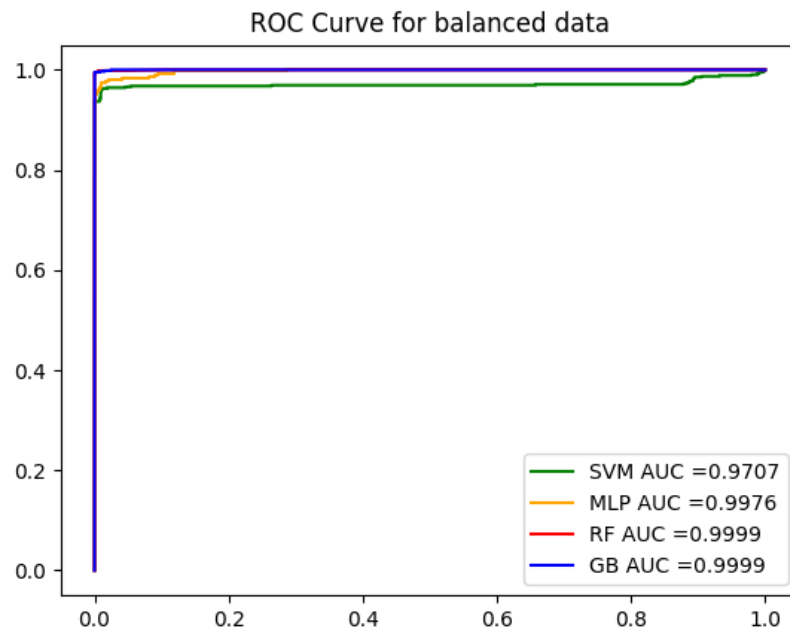


Figure 12 Roc curve for the binary classifiers (balanced)

## 8.2. Multi-class Classification

### 8.2.1. Imbalanced dataset

Figure 13 presents the CM of the SVM, RF, MLP, and GB classifiers where the dataset was unbalanced. The f1-score is focused here because the f1-score is between recall and precision. Tables 9, 10, 11, and 12 show the accuracy, precision, recall, and f1-score of the SVM, RF, MLP, and GB classifiers, respectively. Table 9 shows that the SVM classifier classified the 'passed' class almost correctly while failing to classify the 'terminated' class, which classified one sample correctly of the total terminated 15 samples. The 'failed' and 'error' classes were classified with 0.7 and 0.9, respectively. Tables 10, 11, and 12 show that the RF, MLP, and GB classifiers classified almost all the classes correctly.



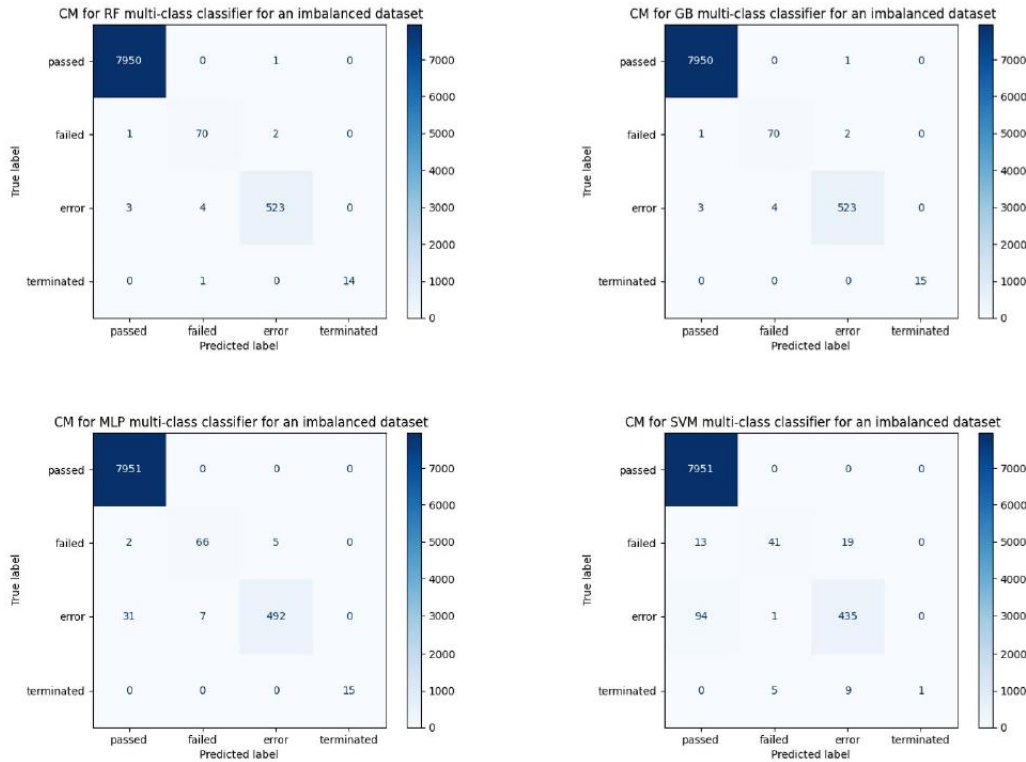


Figure 13 Confusion matrix of the multi-class classifiers for the unbalanced dataset

Table 9 Accuracy, Precision, Recall, and F1 measurements of the multi-class SVM for the unbalanced dataset

| Class      | Accuracy | Precision | Recall | F1-Score |
|------------|----------|-----------|--------|----------|
| Passed     | 98.75    | 0.987     | 1      | 0.993    |
| Failed     | 99.56    | 0.872     | 0.562  | 0.683    |
| Error      | 98.56    | 0.94      | 0.82   | 0.876    |
| Terminated | 99.84    | 1         | 0.0667 | 0.125    |

Table 10 Accuracy, Precision, Recall, and F1 measurements of the multi-class RF for the unbalanced dataset

| Class      | Accuracy | Precision | Recall | F1-Score |
|------------|----------|-----------|--------|----------|
| Passed     | 99.94    | 0.999     | 1      | 1        |
| Failed     | 99.91    | 0.93      | 0.959  | 0.946    |
| Error      | 99.88    | 0.994     | 0.987  | 0.991    |
| Terminated | 99.99    | 1         | 0.933  | 0.966    |

Table 11 Accuracy, Precision, Recall, and F1 measurements of the multi-class MLP for the unbalanced dataset

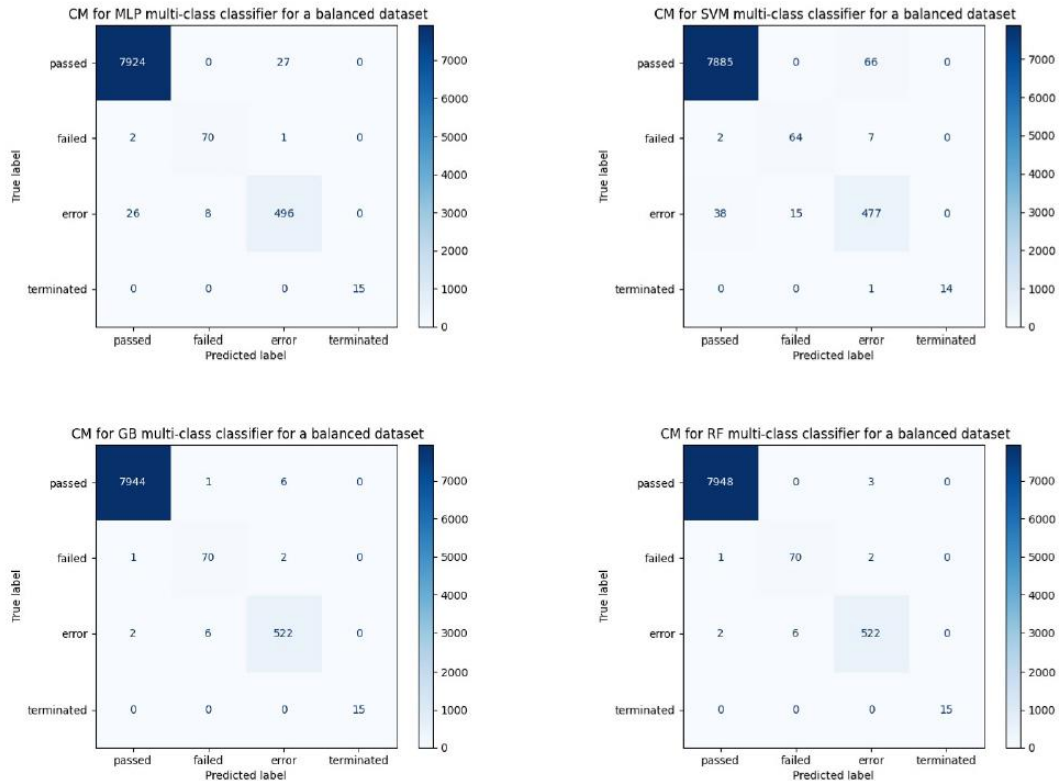
| Class      | Accuracy | Precision | Recall | F1-Score |
|------------|----------|-----------|--------|----------|
| Passed     | 99.61    | 0.996     | 1      | 0.998    |
| Failed     | 99.84    | 0.904     | 0.904  | 0.904    |
| Error      | 99.50    | 0.99      | 0.928  | 0.958    |
| Terminated | 100      | 1         | 1      | 1        |

**Table 12 Accuracy, Precision, Recall, and F1 measurements of the multi-class GB for the unbalanced dataset**

| <i>Class</i>      | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> |
|-------------------|-----------------|------------------|---------------|-----------------|
| <i>Passed</i>     | 99.94           | 0.999            | 1             | 1               |
| <i>Failed</i>     | 99.92           | 0.946            | 0.959         | 0.952           |
| <i>Error</i>      | 99.88           | 0.994            | 0.987         | 0.991           |
| <i>Terminated</i> | 100             | 1                | 1             | 1               |

**8.2.2. Balanced dataset**

Figure 14 presents the CM of the SVM, RF, MLP, and GB classifiers where the dataset was balanced. The f1-score is focused here also. Tables 14, 15, 16, and 17 show the accuracy, precision, recall, and f1-score of the SVM, RF, MLP, and GB classifiers, respectively. Table 15 shows the improvement of the SVM classifier with balanced data over unbalanced data, which this time misclassified only one sample of the ‘terminated’ class. While still the RF, MLP, and GB correctly classified all the classes, with a nuance improvement compared to the unbalanced data, as shown in tables 15, 16, and 17, respectively. Table 13 shows the average of the CV scores that are in the range [95%, 100%], and the best classifier was RF under training the model.



**Figure 14 Confusion matrix of the multi-class classifiers for the balanced dataset**

**Table 13 Cross validation scores of the multi-class classifiers**

| <i>Classifier Algorithm</i> | <i>Cross validation score</i> |
|-----------------------------|-------------------------------|
| <i>SVM</i>                  | 94.63                         |
| <i>RF</i>                   | 99.97                         |
| <i>MLP</i>                  | 98.81                         |
| <i>GB</i>                   | 99.93                         |

**Table 14 Accuracy, Precision, Recall, and F1 measurements of the multi-class SVM for the balanced dataset**

| <i>Class</i>      | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> |
|-------------------|-----------------|------------------|---------------|-----------------|
| <i>Passed</i>     | 98.76           | 0.995            | 0.992         | 0.993           |
| <i>Failed</i>     | 99.72           | 0.81             | 0.877         | 0.842           |
| <i>Error</i>      | 98.52           | 0.866            | 0.9           | 0.883           |
| <i>Terminated</i> | 99.99           | 1                | 0.93          | 0.966           |

**Table 15 Accuracy, Precision, Recall, and F1 measurements of the multi-class RF for the balanced dataset**

| <i>Class</i>      | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> |
|-------------------|-----------------|------------------|---------------|-----------------|
| <i>Passed</i>     | 99.93           | 1                | 1             | 1               |
| <i>Failed</i>     | 99.89           | 0.921            | 0.959         | 0.94            |
| <i>Error</i>      | 99.85           | 0.991            | 0.985         | 0.988           |
| <i>Terminated</i> | 100             | 1                | 1             | 1               |

**Table 16 Accuracy, Precision, Recall, and F1 measurements of the multi-class MLP for the balanced dataset**

| <i>Class</i>      | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> |
|-------------------|-----------------|------------------|---------------|-----------------|
| <i>Passed</i>     | 99.36           | 0.996            | 0.997         | 0.997           |
| <i>Failed</i>     | 99.87           | 0.897            | 0.959         | 0.927           |
| <i>Error</i>      | 99.28           | 0.947            | 0.936         | 0.941           |
| <i>Terminated</i> | 1               | 1                | 1             | 1               |

**Table 17 Accuracy, Precision, Recall, and F1 measurements of the multi-class GB for the balanced dataset**

| <i>Class</i>      | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1-Score</i> |
|-------------------|-----------------|------------------|---------------|-----------------|
| <i>Passed</i>     | 99.88           | 1                | 0.999         | 0.999           |
| <i>Failed</i>     | 99.88           | 0.909            | 0.959         | 0.933           |
| <i>Error</i>      | 99.81           | 0.985            | 0.985         | 0.985           |
| <i>Terminated</i> | 100             | 1                | 1             | 1               |

## 9. Discussion

The measurement of the experiment is very high and accurate, which indicates doubts about it being unbelievable and not trustable. The SVM seems to be more acceptable than the RF, MLP, and GB, even though the RF, MLP, and GB are perfect. In contrast, it is also good to mention that when training the SVM, where the regularization factor was set to 10, the training phase was done for the binary SVM, where the data was balanced and unbalanced, and for the multi-class SVM, where the data was only unbalanced. The modeling time for multi-class SVM where the data was balanced, has taken 5 days without finishing the modeling phase, which leads to remodeling the SVM, but this time the regularization factor was set to 5 and 4 for binary SVM and multi-class SVM, respectively. The measurements reflect that the training is being overfitted despite the normalization and sampling techniques. One of the reasons why the measurements are high is the data itself, because the distribution of the features is not varied and the classes in each feature are discriminated, which leads to the classifier's bias. Figure 15 presents the distribution of the scaled features in the training set. Figure 16 shows the correlation matrix of features in the training set. The matrix shows that feature 7 is highly correlated to feature 5 with a value of 0.93 and to feature 1 with a value of 0.81. Feature 3 is highly

correlated to features 1 and 2, with values 0.82 and 0.82, respectively. In addition, as early mentioned in the limitation section, that the dataset is only for one-year collection, which means that after exploring the dataset and reviewing the document provided by NEP, not all the components are included in the dataset, especially the comparator and unit components. Regarding the research questions:

1. How to extract and select the key Features from ETS dataset?  
The answer to this question is explained in the data collection and preparation sections, section 7.1 and 7.2, respectively.
2. How to implement ML model with optimize the parameters to predict ETS failures?  
The answer to this question is explained in the modeling section, section 7.4.
3. How to develop a framework to be evaluated using real-world ETS dataset?  
The answer to this question is explained in the evaluation section, section 7.5.

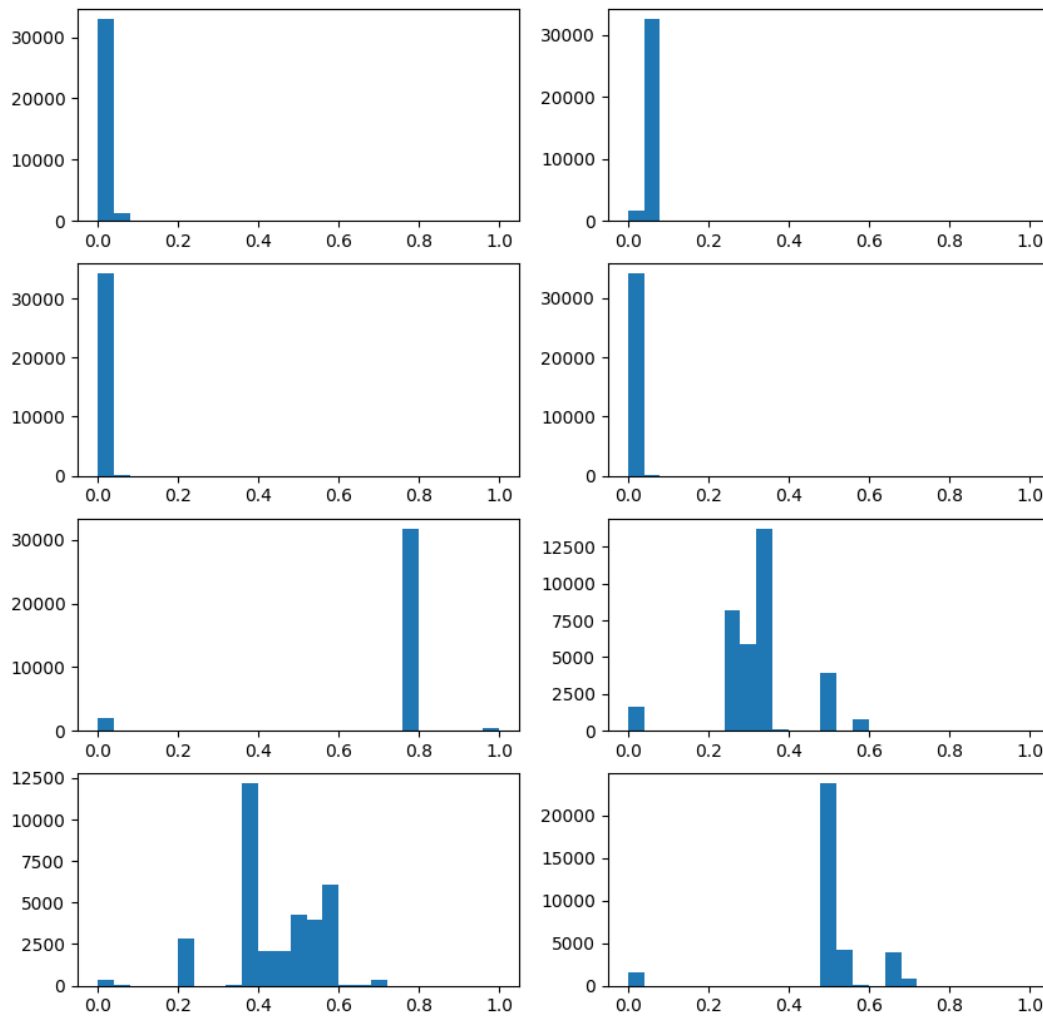


Figure 15 Distribution of the scaled features in the training set.

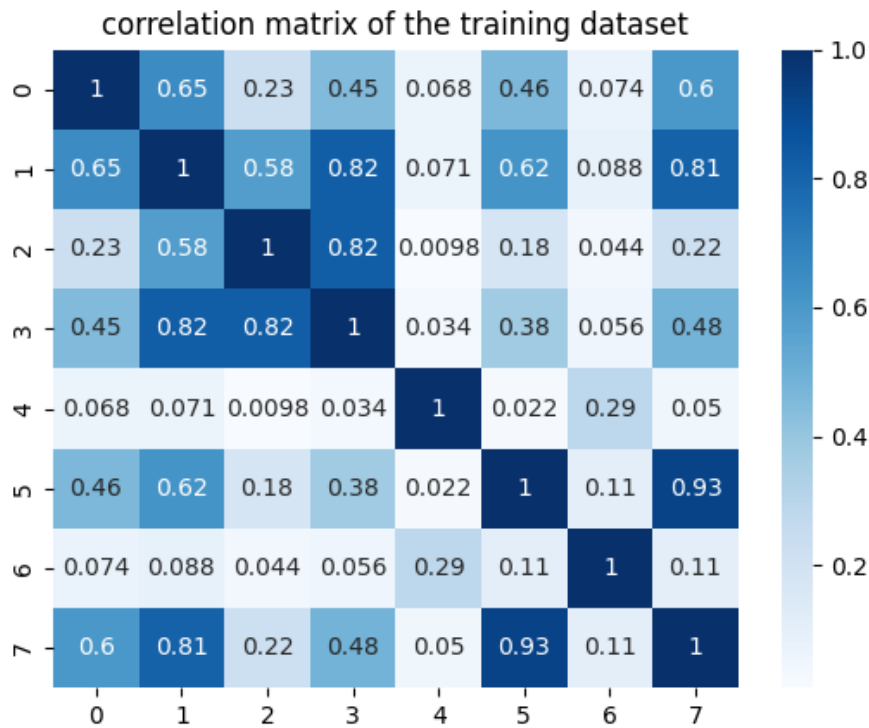


Figure 16 Correlation matrix of the training dataset

## 10. Conclusions

In this thesis, supervised machine learning algorithms were used to predict and classify the data that was collected from embedded test system. The machine learning algorithms were support vector machines, multi-layer perceptrons, random forests, and gradient boosting. Predictive maintenance is often used in manufacturing industries to predict failures before they occur. Predictive maintenance increases productivity and efficiency and decreases system failures and unplanned downtimes. Predictive data analytics is a crucial subfield within data analytics that make accurate predictions. Predictive data analytics extracts insights from data by using machine learning algorithms. The use of machine learning for predictive data analytics faces several challenges, such as extracting and selecting the features from the database and fitting and developing the models to predict the failure of ETS. The CRISP-DM was used in the experiment. Both binary and multi-class classifiers have been provided to fit the models, and cross-validation, sampling techniques, and a confusion matrix have been provided to accurately measure their performance. In addition to accuracy, recall, precision, f1, kappa, mcc, and roc auc measurements are also used. The random forest and gradient boosting performed best in all cases in this experiment. The trained models were supposed to be overfitted and biased, despite all the sampling techniques and normalization.

## 11. Future Work

In this experiment, only one single classifier is used to train the model. For future work, multiple different classifiers can be used at the same time to train the model by using the model ensemble technique, such as boosting or bagging techniques. The model ensemble technique is a technique that allows to combine a set of classification algorithms for building the models and then aggregates the prediction results of these models. Furthermore, the experimental work can also be extended to regression rather than just classification tasks.

## References

- [1] J. D. Kelleher, B. M. Namee & A. D'Arcy, *Fundamentals of machine learning for predictive data analytics: algorithms, worked examples, and case studies*. Cambridge, MA, England: The MIT Press, 2015.
- [2] G. A. Susto, J. Wan, S. Pampuri, M. Zanon, A. B. Johnston, P. G. O'Hara & S. McLoone, *An Adaptive Machine Learning Decision System for Flexible Predictive Maintenance*. 2014 IEEE International Conference on Automation Science and Engineering (CASE). New Taipei, Taiwan. 2014. pp. 806-811. [online]. Available: <https://doi.org/10.1109/CoASE.2014.6899418>. Accessed: September. 27, 2023.
- [3] A. Kane, A. Kore, A. Khandale, S. Nigade & P. P. Joshi, *Predictive Maintenance using Machine Learning*. arXiv.org. 2022. [online]. Available: <https://doi.org/10.48550/arXiv.2205.09402>. Accessed: September. 27, 2023.
- [4] J. Dalzochio, R. Kunst, E. Pignaton, A. Binotto, S. Sanyal, J. Favilla & J. Barbosa, *Machine learning and reasoning for predictive maintenance in industry 4.0: Current status and challenges*. Elsevier. 2020. [online]. Available: <https://doi.org/10.1016/j.compind.2020.103298>. Accessed: September. 27, 2023.
- [5] A. J. Albee, *The Evolution of ICT: PCB Technologies, Test Philosophies, and Manufacturing Business Models Are Driving In-Circuit Test Evolution and Innovations*. IPC APEX EXPO. Conf. 2013. Available: [https://www.circuitinsight.com/pdf/evolution\\_ict\\_changing\\_pcb\\_technologies\\_ipc.pdf](https://www.circuitinsight.com/pdf/evolution_ict_changing_pcb_technologies_ipc.pdf). Accessed: September. 27, 2023.
- [6] V. Hirsch, P. Reimann, O. Kirn & B. Mitschang, *Analytical Approach to Support Fault Diagnosis and Quality Control in End-Of-Line Testing*. Proceda CIRP. 2018. [online]. Available: <https://www.sciencedirect.com/science/article/pii/S2212827118301240>. Accessed: September. 27, 2023.
- [7] A. Maria & K. Srihari, *A Review of Knowledge-Based Systems in Printed Circuit Board Assembly*. The international journal of advanced manufacturing technology. 1992. [online]. Available: <https://link.springer.com/article/10.1007/BF03500680>. Accessed: September. 27, 2023.
- [8] A. V. Joshi, *Machine Learning and Artificial Intelligence*. Second Edition Cham. Switzerland: Springer. 2023. [online]. Available: <https://doi.org/10.1007/978-3-031-12282-8> Accessed: September. 27, 2023.
- [9] K. Malhotra & Y. Kumar, *Challenges to implement Machine Learning in Embedded Systems*. 2020 2nd International Conference on Advances in Computing, Communication Control and Networking (ICACCCN). Greater Noida. India. 2020. pp. 477-481. [online]. Available: <https://doi.org/10.1109/ICACCCN51052.2020.9362893>. Accessed: September. 27, 2023.
- [10] M. W. Berry, A. Mohamed & B. W. Yap, *Supervised and Unsupervised Learning for Data Science*. Cham. Switzerland: Springer. 2020. [online]. Available: [https://doi.org/10.1007/978-3-030-22475-2\\_1](https://doi.org/10.1007/978-3-030-22475-2_1). Accessed: September. 27, 2023.
- [11] E. Alpaydin, *Introduction to Machine Learning*. Third Edition. The MIT Press. 2014.
- [12] C. Elkan, *Predictive analytics and data mining*. San Diego. University of California. 2013. [online]. Available: [https://cseweb.ucsd.edu/classes/sp15/cse190-c/files/elkan\\_dm.pdf](https://cseweb.ucsd.edu/classes/sp15/cse190-c/files/elkan_dm.pdf). Accessed: September. 27, 2023.
- [13] S. J. Russell & P. Norvig, *Artificial Intelligence: A Modern Approach*. Harlow. England: Pearson Education. 2022.
- [14] A. Chigurupati, R. Thibaux & N. Lassar, *Predicting Hardware Failure Using Machine Learning*. 2016 Annual Reliability and Maintainability Symposium (RAMS). 2016. [online]. Available: <https://doi.org/10.1109/RAMS.2016.7448033>. Accessed: September. 27, 2023.
- [15] D. Mayer, *Support Vector Machines: The Interface to libsvm in package e1071*. FH Technikum Wien. Austria. 2015. [online]. Available: <https://vps.fmvz.usp.br/CRAN/web/packages/e1071/vignettes/svmdoc.pdf>. Accessed: September. 27, 2023.
- [16] C. Stoean & R. Stoean, *Support Vector Machines and Evolutionary Algorithms for Classification Single or Together?.* Springer. 2014. pp. 23. [online]. Available: <https://doi.org/10.1007/978-3-319-06941-8>. Accessed: September. 27, 2023.
- [17] C. Zhang & Y. Ma, *Ensemble Machine Learning Methods and Applications*. Springer. 2012. pp. 35-163. [online]. Available: <https://doi.org/10.1007/978-1-4419-9326-7>. Accessed: September. 27, 2023.

- [18] J. Hemanth, X. Fernando, P. Lafata & Z. Baig, *International Conference on Intelligent Data Communication Technologies and Internet of Things (ICICI) 2018*. Springer. 2019. pp. 760. [online]. Available: <https://doi.org/10.1007/978-3-030-03146-6>. Accessed: September. 27, 2023.
- [19] M. W. Gardner & S. R. Dorling, *Artificial Neural Networks (The Multilayer Perceptron) – A review of applications in the atmospheric sciences*. Atmospheric Environment Elsevier. 1998. pp. 2627-2636. [online]. Available: [https://doi.org/10.1016/S1352-2310\(97\)00447-0](https://doi.org/10.1016/S1352-2310(97)00447-0). Accessed: September. 27, 2023.
- [20] M. M. Ali, B. K. Paul, K. Ahmed, F. M. Bauji, J. M. W. Quinn & M. A. Moni, *Heart disease prediction using supervised machine learning algorithms: Performance analysis and comparison*. *Computers in Biology and Medicine*. Elsevier. 2021. [online]. Available: <https://doi.org/10.1016/j.compbiomed.2021.104672>. Accessed: September. 27, 2023.
- [21] K. Li, S. Yao, Z. Zhang, B. Cao, C. M. Wilson, D. Kalos, P. f. Kuan, R. Zhu & X. Wang, *Efficient gradient boosting for prognostic biomarker discovery*. *Bioinformatics*. 2022. pp. 1631-1638. [online]. Available: <https://doi.org/10.1093/bioinformatics/btab869>. Accessed: September. 27, 2023.
- [22] K. Liulys, *Machine Learning Application in Predictive Maintenance*. 2019 Open Conference of Electrical, Electronic and Information Sciences. Vilnius. Lithuania. 2019. pp. 1-4. [online]. Available: <https://doi.org/10.1109/eStream.2019.8732146>. Accessed: September. 27, 2023.
- [23] G. Biau, B. Cadre & L. Rouviere, *Accelerated gradient boostin*. *Machine Learning 2019*. Springer. 2019. [online]. Available: <https://doi.org/10.1007/s10994-019-05787-1>. Accessed: September. 27, 2023.
- [24] X. W. Liu, Z. L. Long, W. Zhang & L. M. Yang, *Key feature space for predicting the glass-forming ability of amorphous alloys revealed by gradient boosted decision trees model*. *Journal of Alloys and Compounds*. Elsevier. 2022. [online]. Available: <https://doi.org/10.1016/j.jallcom.2021.163606>. Accessed: September. 27, 2023.
- [25] D. L. Olson & D. Delen, *Advanced Data Mining*. Springer. 2008. [online]. Available: <https://doi.org/10.1007/978-3-540-76917-0>. Accessed: September. 27, 2023.
- [26] I. T. Franco & R. M. d. Figueiredo, *Predictive Maintenance: An Embedded System Approach*. *Journal of Control, Automation and Electrical Systems*. 2023. [online]. Available: <https://doi.org/10.1007/s40313-022-00949-4>. Accessed: September. 27, 2023.
- [27] G. A. Susto, A. Schirru & S. Pampuri, *Machine Learning for Predictive Maintenance: A Multiple Classifier Approach*. *IEEE Transactions on Industrial Informatics*. 2015. [online]. Available: <https://doi.org/10.1109/TII.2014.2349359>. Accessed: September. 27, 2023.
- [28] M. Calabrese, M. Cimmino, M. Manfrin, F. Fiume, D. Kapetis, M. Mengoni, S. Ceccacci, E. Frontoni, M. Paolanti, A. Carrotta & G. Toscano, *An Event Based Machine Learning Framework for Predictive Maintenance in Industry 4.0*. *IEEE/ASME International Conference on Mechatronic and Embedded Systems and Applications*. Anaheim. California. USA. 2019. [online]. Available: <https://doi.org/10.1115/DETC2019-97917>. Accessed: September. 27, 2023.
- [29] T. P. Carvalho, F. A. A. M. N. Soares, R. Vita, R. da P. Francisco, J. P. Basto & S. G. S. Alcala, *A systematic literature review of machine learning methods applied to predictive maintenance*. *Computers & Industrial Engineering*. Elsevier. 2019. [online]. Available: <https://doi.org/10.1016/j.cie.2019.106024>. Accessed: September. 27, 2023.
- [30] S. Dey, M. Kedia, N. Agarwal & A. Basu, *Embedded Support Vector Machine: Architectural Enhancements and Evaluation*. 20th International Conference on VLSI Design held jointly with 6th International Conference on Embedded Systems (VLSID'07). Bangalore. India. 2007. pp. 685-690. [online]. Available: <https://doi.org/10.1109/VLSID.2007.73>. Accessed: September. 27, 2023.
- [31] L. G. Castanheira, J. A. d. Vasconcelos, A. J. R. Reis, P. H. V. Magalhaes & S. A. L. d. Silva, *Application of Neural Networks in the Classification of Incipient Faults in Power Transformers: A Study of Case*. The 2011 International Joint Conference on Neural Networks. San Jose. CA. USA. 2011. pp. 3099-3104. [online]. Available: <https://doi.org/10.1109/IJCNN.2011.6033631>. Accessed: September. 27, 2023.
- [32] T. Gkamas, S. Kontogiannis, V. Karaiskos, C. Pikridas & I. A. Karolos, *Proposed Cloud-assisted Machine Learning Classification Process implemented on Industrial Systems: Application to Critical Events Detection and Industrial Maintenance*. 2022 5<sup>th</sup> World Symposium on Communication Engineering (WSCE). Nogoya. Japan. 2022. pp. 95-99. [online]. Available: <https://doi.org/10.1109/WSCE56210.2022.9916043>. Accessed: September. 27, 2023.
- [33] S. G. K. Patro & K. K. Sahu, *Normalization: A Preprocessing Stage*. arXiv preprint. 2015. [online]. Available: <https://doi.org/10.48550/arXiv.1503.06462>. Accessed: September. 27, 2023.
- [34] D. V. Souza, J. X. Santos, H. C. Vieira, T. L. Naide, S. Nisgoski & L. E. S. Oliveira, *An automatic recognition system of Brazilian flora species based on textural features of macroscopic images of wood*.

Wood Science and Technology. 2020. [online]. Available: <https://doi.org/10.1007/s00226-020-01196-z>. Accessed: September. 27, 2023.

[35] N. V. Chawla, K. W. Bowyer, L. O. Hall & W. P. Kegelmeyer, *SMOTE: Synthetic Minority Over-sampling Technique*. Journal of Artificial Intelligence Research. 2002. [online]. Available: <https://doi.org/10.48550/arXiv.1106.1813>. Accessed: September. 27, 2023.

[36] D. Chicco, M. J. Warrens & G. Jurman, *The Matthews Correlation Coefficient (MCC) is More Informative Than Cohen's Kappa and Brier Score in Binary Classification Assessment*. IEEE. 2021. [online]. Available: <https://doi.org/10.1109/ACCESS.2021.3084050>. Accessed: September. 27, 2023.



## Appendices

### Appendix A: Feature extraction

```

1 // Extract.cs
2 // Author: Fayad Al Hanash
3 namespace ExtractFeatures;
4 /// <summary>
5 /// A class to extract features from the measurement results and test results.
6 /// </summary>
7 2 references
8 public class Extract
9 {
10     /// <summary>
11     /// Method that extracts the features
12     /// Calls the ReadFile method to get the measurements and test result arrays.
13     /// Calls GetMeasurementResultFeatures and GetTestResultFeatures methods to create list of
14     /// MeasurementResultFeatures and TestResultFeatures objects
15     /// Calls the GetMergedList method to merge the lists into one list of MeasurementResultFeatures objects
16     /// Groups the measurement result features based on the test result ID
17     /// Creates a new list of ExtractedFeatures object
18     /// Iterates through each group to calculate the measurements and frequencies
19     /// Calls the WriteFile method to save the features in the file
20     /// </summary>
21     1 reference
22     public bool StartExtracting(string measurementPath, string testPath, string desPath)
23     {
24         (string[] measurementArray, string[] testArray) = FileUtility.ReadFile(measurementPath, testPath);
25         if (measurementArray == null) return false;
26         List<MeasurementResultFeatures> measurementResultList = GetMeasurementResultFeatures(measurementArray);
27         List<TestResultFeatures> testResultList = GetTestResultFeatures(testArray);
28         List<MeasurementResultFeatures> mergerdList = GetMergedList(measurementResultList, testResultList);
29         var groups = mergerdList.GroupBy(b => b.TestResultID);
30         List<ExtractedFeatures> featureAfterList = new List<ExtractedFeatures>();
31         foreach (var g in groups)
32         {
33             var ExtQuery = from b in g let execTime = b.ExecutionTime select execTime;
34             var NumericDataQuery = from b in g let numericData = b.NumericData select numericData;
35             var NumericLowLimitQuery = from b in g let numericLowLimit = b.NumericLowLimit select numericLowLimit;
36             var NumericHighLimitQuery = from b in g let numericHighLimit = b.NumericHighLimit select numericHighLimit;
37             var measurementIDQuery = from b in g let measurementID = b.MeasurementResultID select measurementID;
38
39             var nameQuery = from n in g let name = n.Name select name;
40             var comparator = from b in g let comp = b.Comparator select comp;
41             var unit = from b in g let u = b.Unit select u;
42             var GELE = from c in comparator where c.GELE == 1 select c.GELE;
43             var EQ = from c in comparator where c.EQ == 1 select c.EQ;
44             var LOG = from c in comparator where c.LOG == 1 select c.LOG;
45             var GTLT = from c in comparator where c.GTLT == 1 select c.GTLT;
46             var V = from u in unit where u.V == 1 select u.V;
47             var μA = from u in unit where u.μA == 1 select u.μA;
48             var NULL = from u in unit where u.NULL == 1 select u.NULL;
49
50             float extMin = ExtQuery.Min();
51             float numDMin = NumericDataQuery.Min();
52             float numLMin = NumericLowLimitQuery.Min();
53             float numHMin = NumericHighLimitQuery.Min();
54
55             ExtractedFeatures after = new ExtractedFeatures()
56             {
57                 TestResultID = g.Key,
58                 MeasurementTestFrequency = measurementIDQuery.Count(),
59                 NameFrequency = nameQuery.Count(),
60                 ExecTimeTotal = g.ElementAt(0).ExecutionTimeAll,
61                 ExecTimeMes = new Measurement
62                 {

```

```

61         Avg = ExtQuery.Average(),
62         Max = ExtQuery.Max(),
63         Min = ExtQuery.OrderBy(x => x != 0 && extMin != 0).First(),
64         Median = Helper.MEDIAN(ExtQuery),
65         StdDev = Helper.STDDEV(ExtQuery)
66     },
67     NumericDataMes = new Measurement
68     {
69         Avg = NumericDataQuery.Average(),
70         Max = NumericDataQuery.Max(),
71         Min = NumericDataQuery.OrderBy(x => x != 0 && numDMin != 0).First(),
72         Median = Helper.MEDIAN(NumericDataQuery),
73         StdDev = Helper.STDDEV(NumericDataQuery)
74     },
75     NumericLowLimitMes = new Measurement
76     {
77         Avg = NumericLowLimitQuery.Average(),
78         Max = NumericLowLimitQuery.Max(),
79         Min = NumericLowLimitQuery.OrderBy(x => x != 0 && numLMin != 0).First(),
80         Median = Helper.MEDIAN(NumericLowLimitQuery),
81         StdDev = Helper.STDDEV(NumericLowLimitQuery)
82     },
83     NumericHighLimitMes = new Measurement
84     {
85         Avg = NumericHighLimitQuery.Average(),
86         Max = NumericHighLimitQuery.Max(),
87         Min = NumericHighLimitQuery.OrderBy(x => x != 0 && numHMin != 0).First(),
88         Median = Helper.MEDIAN(NumericHighLimitQuery),
89         StdDev = Helper.STDDEV(NumericHighLimitQuery)
90     },
91     GELEFreq = GELE.Count(),
92     EQFreq = EQ.Count(),
93     LOQFreq = LOG.Count(),
94     GTLTFreq = GTLT.Count(),
95     VFreq = V.Count(),
96     μAFreq = μA.Count(),
97     NULLFreq = NULL.Count(),
98     Status = g.ElementAt(0).Status
99 };
100
101     featureAfterList.Add(after);
102 }
103     return FileUtility.WriteFile(featureAfterList, desPath);
104 }
105
106 <summary>
107 <summary> Method that returns a list of MeasurementResultFeatures object
108 </summary>
109 <summary>
110 <summary> 1 reference
111 <summary> List<MeasurementResultFeatures> GetMeasurementResultFeatures(string[] arr) =>
112 <summary>
113 <summary>
114 <summary>
115 <summary>
116 <summary>
117 <summary>
118 <summary>
119 <summary>
120 <summary>
121 <summary>
122 <summary>
123 <summary>
124 <summary>
125 <summary>
126 <summary>
127 <summary>
128 <summary>
129 <summary>

```

```

130         .ToList();
131
132     /// <summary>
133     /// Method that returns a list of TestResultFeatures object
134     /// </summary>
135     1 reference
136     List<TestResultFeatures> GetTestResultFeatures(string[] arr) =>
137     {
138         arr
139         .Skip(1)
140         .Select(line =>
141         {
142             string[] col = line.Split(',');
143             return new TestResultFeatures
144             {
145                 TestResultID = int.Parse(col[0]),
146                 Timestamp = col[1].Substring(0, col[1].LastIndexOf(':')),
147                 ExecutionTime = float.Parse(col[2]),
148                 Status = Helper.GetStatus(col[3])
149             };
150         })
151         .ToList();
152
153     /// <summary>
154     /// Method that merges list of MeasurementResultFeatures and TestResultFeatures lists and returns a new list
155     /// of MeasurementResultFeatures
156     /// </summary>
157     1 reference
158     List<MeasurementResultFeatures> GetMergedList(List<MeasurementResultFeatures> mRList, List<TestResultFeatures> tRList)
159     {
160         var list = (from mR in mRList
161                     join tR in tRList
162                     on mR.TestResultID equals tR.TestResultID
163                     into temp
164                     from tR in temp.DefaultIfEmpty()
165                     select new MeasurementResultFeatures
166                     {
167                         MeasurementResultID = mR.MeasurementResultID,
168                         TestResultID = mR.TestResultID,
169                         Name = mR.Name,
170                         TimeStamp = mR.TimeStamp,
171                         ExecutionTime = mR.ExecutionTime,
172                         ExecutionTimeAll = tR.ExecutionTime,
173                         NumericData = mR.NumericData,
174                         NumericLowLimit = mR.NumericLowLimit,
175                         NumericHighLimit = mR.NumericHighLimit,
176                         Comparator = mR.Comparator,
177                         Unit = mR.Unit,
178                         Status = tR.Status
179                     }).ToList();
180
181         return list;
182     }

```

## Extract.cs

```

1 // TestResultFeatures.cs
2 namespace ExtractFeatures;
3 /// <summary>
4 /// A test result features class represents the features of the test results.
5 /// </summary>
6 4 references
7 public class TestResultFeatures
8 {
9     2 references
10    public int TestResultID { get; set; }
11    1 reference
12    public string Timestamp { get; set; }
13    2 references
14    public float ExecutionTime { get; set; }
15    2 references
16    public int Status { get; set; }
17 }

```

## TestResultFeatures.cs

```
1 // MeasurementResultFeatures.cs
2 namespace ExtractFeatures;
3 /// <summary>
4 /// A measurement result features class represents the features of the measurement results.
5 /// </summary>
6 public class MeasurementResultFeatures
7 {
8     public int MeasurementResultID { get; set; }
9     public int TestResultID { get; set; }
10    public string Name { get; set; }
11    public string TimeStamp { get; set; }
12    public float ExecutionTime { get; set; }
13    public float ExecutionTimeAll { get; set; }
14    public float NumericData { get; set; }
15    public float NumericLowLimit { get; set; }
16    public float NumericHighLimit { get; set; }
17    public Comparator Comparator { get; set; }
18    public Unit Unit { get; set; }
19    public int Status { get; set; }
20 }
```

#### MeasurementResultFeatures.cs

```
1 // Comparators.cs
2 namespace ExtractFeatures;
3 /// <summary>
4 /// A comparator record represents the comparators used to evaluate a test.
5 /// </summary>
6 public record Comparator()
7 {
8     public int GELE { get; set; }
9     public int EQ { get; set; }
10    public int LOG { get; set; }
11    public int GTLT { get; set; }
12
13    public override string ToString() => $"{GELE},{EQ},{LOG},{GTLT}";
14 }
```

#### Comparators.cs

```

1 // Unit.cs
2 namespace ExtractFeatures;
3 4 references /// <summary>
4 /// A unit record represents the unit of the measured values.
5 /// </summary>
6 6 references
7 public record Unit
8 {
9     4 references
10    public int V { get; set; }
11    4 references
12    public int  $\mu$ A { get; set; }
13    4 references
14    public int NULL { get; set; }
15    0 references
16    public override string ToString() => $"{V},{ $\mu$ A},{NULL}";
17 }

```

Unit.cs

```

1 // Measurement.cs
2 namespace ExtractFeatures;
3 4 references /// <summary>
4 /// A measurement class represents the calculated Min, Max, Mean, Median,
5 /// and StdDev of the numerical and execution time measurements.
6 /// </summary>
7 8 references
8 public class Measurement
9 {
10    6 references
11    public float Min { get; set; }
12    6 references
13    public float Max { get; set; }
14    8 references
15    public float Avg { get; set; }
16    7 references
17    public float Median { get; set; }
18    8 references
19    public float StdDev { get; set; }
20
21    0 references
22    public override string ToString() =>
23    |   $"{Min},{Max},{Avg},{Median},{StdDev}";
24    1 reference
25    public string ExecToString() => $"{Max},{Avg},{StdDev}";
26    1 reference
27    public string NumDataToString() => $"{Min},{Avg},{Median},{StdDev}";
28    1 reference
29    public string NumDLowToString() => $"{Avg},{Median},{StdDev}";
30    0 references
31    public string NumDHighToString() => NumDLowToString();
32 }

```

Measurement.cs

```

1  // ExtractedFeatures.cs
2  namespace ExtractFeatures;
3  /// <summary>
4  /// A feature class represents the extracted features (features after extraction).
5  /// </summary>
6  5 references
7  public class ExtractedFeatures
8  {
9      1 reference
10     public int TestResultID { get; set; }
11     1 reference
12     public int MeasurementTestFrequency { get; set; }
13     1 reference
14     public int NameFrequency { get; set; }
15     2 references
16     public float ExecTimeTotal { get; set; }
17     2 references
18     public Measurement ExecTimeMes { get; set; }
19     2 references
20     public Measurement NumericDataMes { get; set; }
21     1 reference
22     public Measurement NumericLowLimitMes { get; set; }
23     1 reference
24     public Measurement NumericHighLimitMes { get; set; }
25     1 reference
26     public int GELEFreq { get; set; }
27     1 reference
28     public int EQFreq { get; set; }
29     1 reference
30     public int LOQFreq { get; set; }
31     1 reference
32     public int GTLTFreq { get; set; }
33     1 reference
34     public int VFreq { get; set; }
35     1 reference
36     public int μAFreq { get; set; }
37     1 reference
38     public int NULLFreq { get; set; }
39     2 references
40     public int Status { get; set; }
41
42     25 v /*public override string ToString() => $"{TestResultID},{MeasurementTestFrequency}," +
43         26     $"{NameFrequency},{ExecTimeTotal},{ExecTimeMes},{NumericDataMes},{NumericLowLimitMes}," +
44         27     $"{NumericHighLimitMes},{GELEFreq},{EQFreq},{LOQFreq},{GTLTFreq},{VFreq},{μAFreq}, " +
45         28     $"{NULLFreq},{(int)Status}";*/
46
47     29 v public override string ToString() => $"{ExecTimeTotal},{ExecTimeMes.ExecToString()}" +
48         30     $"{NumericDataMes.NumDataToString()},{(int)Status}";
49 }

```

ExtractedFeatures.cs

```

1 // Helper.cs
2 namespace ExtractFeatures;
3
4 13 references
5 public class Helper
6 {
7     /*public static string ExtractedFeaturesNames = "TestResultID,MeasurementTestFrequency,NameFrequency," +
8     "ExecTimeTotal,ExecMin,ExecMax,ExecAvg,ExecMedian,ExecStdDev,NumDataMin,NumDataMax,NumDataAvg," +
9     "NumDataMedian,NumDataStdDev,NumDataLowLimMin,NumDataLowLimMax,NumDataLowLimAvg,NumDataLowLimMedian," +
10    "NumDataLowLimStdDev,NumDataHighLimMin,NumDataHighLimMax,NumDataHighLimAvg,NumDataHighLimMedian," +
11    "NumDataHighLimStdDev,GELEFreq,EQFreq,LOQFreq,GTLTFreq,VFreq,μAFreq,NULLFreq,Status";*/
12    1 reference
13    public static string ExtractedFeaturesNames = "ExecTimeTotal,ExecMax,ExecAvg,ExecStdDev,NumDataMin," +
14    "NumDataAvg,NumDataMedian,NumDataStdDev,Status";
15    0 references
16    public static string OriginalFeaturesNames = "measurement_result_id,test_result_id,name,timestamp," +
17    "execution_time,numeric_data,numeric_low_limit,numeric_high_limit,comparator,unit,status";
18
19    /// <summary>
20    /// Method that returns a new instance of comparator and sets the value of the comparator based on the string value.
21    /// </summary>
22    1 reference
23    public static Func<string, Comparator> GetComparator = x => x switch
24    {
25        "GELE" => new() { GELE = 1 },
26        "EQ" => new() { EQ = 1 },
27        "LOG" => new() { LOG = 1 },
28        "GTLT" => new() { GTLT = 1 },
29        _ => new(),
30    };
31
32    /// <summary>
33    /// Method that returns a new instance of unit and sets the value of the unit based on the string value.
34    /// </summary>
35    1 reference
36    public static Func<string, Unit> GetUnit = x => x switch
37    {
38        "V" => new Unit { V = 1 },
39        "μA" => new Unit { μA = 1 },
40        "NULL" => new Unit { NULL = 1 },
41        _ => new Unit()
42    };
43
44    /// <summary>
45    /// Method that returns the status of a test as an integer based on the string value.
46    /// </summary>
47    2 references
48    public static int GetStatus(string status)
49    {
50        return status switch
51        {
52            "passed" => 0,
53            "failed" => 1,
54            "error" => 2,
55            "terminated" => 3,
56            _ => 4,
57        };
58    }
59
60    /// <summary>
61    /// Method that returns the status of a test as a string based on the integer value.
62    /// </summary>

```

```

58 0 references
   public static string GetStatus(int status)
59  {
60     return status switch
61     {
62         0 => "passed",
63         1 => "failed",
64         2 => "error",
65         3 => "terminated",
66         _ => "unknown",
67     };
68 }
69
70 /// <summary>
71 /// Method that calculates the standard deviation.
72 /// </summary>
73 4 references
   public static float STDDEV(IEnumerable<float> s)
74  {
75     if (s.Count() > 0)
76     {
77         float x = (float) Math.Sqrt(s.Sum(d => Math.Pow(d - s.Average(), 2)) / (s.Count() - 1));
78         if (double.IsNaN(x))
79             return 0;
80         else
81             return x;
82     }
83     else return 0;
84 }
85
86 /// <summary>
87 /// Method that calculates the median
88 /// </summary>
89 4 references
   public static float MEDIAN(IEnumerable<float> s) =>
90     s.Count() > 0 ? s.OrderBy(x => x).ElementAt((int)Math.Ceiling((double)(s.Count() - 1) / 2)) : 0;
91
92 }

```

### Helper.cs

```

1 // FileUtility.cs
2 namespace ExtractFeatures;
3
4 2 references
   public class FileUtility
5  {
6     /// <summary>
7     /// Method that reads the measurement results and test results from the csv files,
8     /// returns a tuple of string arrays containing the measurement results and test results.
9     /// </summary>
10 1 reference
   public static (string[], string[]) ReadFile(string measuremetPath, string testPath)
11  {
12     try
13     {
14         if (!File.Exists(measuremetPath))
15             throw new FileNotFoundException("File not found", measuremetPath);
16
17         if (testPath == null)
18         {
19             return (File.ReadAllLines(measuremetPath), null);
20         }
21         else
22         {
23             if (!File.Exists(testPath))
24                 throw new FileNotFoundException("File not found", testPath);
25             return (File.ReadAllLines(measuremetPath), File.ReadAllLines(testPath));
26         }
27     }
28     catch (Exception e)
29     {
30         throw new Exception(e.ToString());
31     }
32     return (null, null);
33 }
34 }

```



```
36 | /// <summary>  
37 | /// Method that writes the extracted features to a csv file.  
38 | /// </summary>  
   | 1 reference  
39 | public static bool WriteFile(List<ExtractedFeatures> afterExtracts, string path)  
40 | {  
41 |     try  
42 |     {  
43 |         using (StreamWriter sw = new StreamWriter(path))  
44 |         {  
45 |             sw.WriteLine(Helper.ExtractedFeaturesNames);  
46 |             afterExtracts.ForEach(x => sw.WriteLine(x.ToString()));  
47 |             return true;  
48 |         }  
49 |     }  
50 |     catch (Exception e)  
51 |     {  
52 |         throw new Exception(e.ToString());  
53 |     }  
54 |     return false;  
55 | }  
56 | }  
57 |
```

FileUtility.cs

## Appendix B: Modeling & Evaluation

```
1 | # main.py  
2 | # Author: Fayad Al Hanash  
3 | # Date: 2023-08-30  
4 |   
5 | # This file is the main file for the project  
6 | # import libraries  
7 | import numpy as np  
8 | from helper import Helper  
9 | import sklearn.model_selection as ms  
10 | from file_utility import FileUtility  
11 | from parameters_tunning import ParametersTunning  
12 | from binary_classification import BinaryClassification  
13 | from multi_classification import MultiClassification  
14 |   
15 | # method that runs the binary classification for the given algorithm and dataset  
16 | def binary_classification( alg, is_blanced, is_modeled, x_train, x_test, y_train, y_test):  
17 |     cl = BinaryClassification(alg,is_blanced, x_train, x_test, y_train, y_test)  
18 |     if is_modeled:  
19 |         mod = FileUtility.load_model(f'{alg}_bin_{is_blanced}_balancing')  
20 |         cl.set_model(mod)  
21 |     else:  
22 |         cl.fit()  
23 |         FileUtility.save_model(cl.get_model(), f'{alg}_bin_{is_blanced}_balancing')  
24 |     cl.display_matrix(f"{alg} bin for {is_blanced} dataset")  
25 |     cl.plot()  
26 |
```

```

27
28 # method that runs the multi classification for the given algorithm and dataset
29 def multi_classification( alg, is_blanced, is_modeled, x_train, x_test, y_train, y_test):
30     cl = MultiClassification(alg,is_blanced, x_train, x_test, y_train, y_test)
31     if is_modeled:
32         mod = FileUtility.load_model(f'{alg}_mul_{is_blanced}_balancing')
33         cl.set_model(mod)
34     else:
35         cl.fit()
36         FileUtility.save_model(cl.get_model(), f'{alg}_mul_{is_blanced}_balancing')
37     cl.display_confusion_matrix(f'{alg} multi for {is_blanced} dataset")
38     cl.plot()
39
40
41 # method that runs the parameters tuning for the given algorithm and dataset
42 def tune_parameters(alg,is_balanced,cl_type, x_train, y_train):
43     pt = ParametersTunning(alg, x_train, y_train)
44     para = pt.tune_parameters(is_balanced)
45     FileUtility.save_text_json(para, f'{alg}_{cl_type}_{is_balanced}_balancing.txt')
46
47 # method that runs the print the cross validation for an already saved cross validation
48 def print_cross_validation(alg,is_balanced,cl_type):
49     cr = FileUtility.load_model(f'{alg}_{cl_type}_{is_balanced}_balancing_cross')
50     score = np.mean(cr['test_score'])
51     print(f'Mean Accuracy: {score:.2%}')
52
53
54 # methid that runs the main method
55 def main():
56
57     # read the data from the file
58     xData, yData = FileUtility.read_file('data/extracted.csv')
59     # split the data into train and test
60     x_train, x_test, y_train, y_test = ms.train_test_split(xData, yData, train_size=0.80,
61     test_size=0.20, random_state=0)
62     # scale the data
63     x_train, x_test = Helper.scale_min_max(x_train, x_test)
64
65
66     algorithm = "svm"
67     is_blanced = False
68     cl_type = 'bin'
69     is_modeled = True
70
71     if cl_type == 'bin':
72         binary_classification(algorithm, is_blanced, is_modeled, x_train, x_test, y_train, y_test)
73     else:
74         multi_classification(algorithm, is_blanced, is_modeled, x_train, x_test, y_train, y_test)
75
76
77 # run the main method
78 if __name__ == "__main__":
79     main()

```

main.py

```
1 | # model.py
2 | # import libraries
3 | from numpy import mean
4 | from imblearn.pipeline import Pipeline
5 | from imblearn.over_sampling import SMOTE
6 | from imblearn.under_sampling import RandomUnderSampler
7 | from sklearn.model_selection import cross_validate , RepeatedStratifiedKFold
8 |
9 | # The class Model which is used to build the model
10 | class Model:
11 |     # method that initializes the model
12 |     def __init__(self, x_train, x_test, y_train, y_test):
13 |         self.model = None
14 |         self.algorithm = None
15 |         self.score = None
16 |         self.x_train = x_train
17 |         self.x_test = x_test
18 |         self.y_train = y_train
19 |         self.y_test = y_test
20 |         self.y_predict = None
21 |         self.y_predict_proba = None
22 |         self.cross_validation = None
23 |
24 |     # method that fits the model
25 |     def fit(self):
26 |         pass
27 |
28 |     # method that fits the model with unbalanced data
29 |     def fit_unbalanced_data(self):
30 |         self.model = self.algorithm.fit(self.x_train, self.y_train)
31 |         self.score = self.model.score(self.x_test, self.y_test)
32 |         self.y_predict = self.model.predict(self.x_test)
33 |         self.y_predict_proba = self.model.predict_proba(self.x_test)
34 |
```



```
1 | # binary_classification.py
2 | # Importing the libraries
3 | from plot import Plot
4 | from sklearn import svm
5 | from model import Model
6 | from helper import Helper
7 | from sklearn.neural_network import MLPClassifier
8 | from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier
9 | from sklearn.metrics import confusion_matrix, accuracy_score, precision_score
10 | from sklearn.metrics import recall_score, f1_score, roc_auc_score, matthews_corrcoef, cohen_kappa_score, roc_curve
11 |
12 | # The class BinaryClassification which is used to build the binary classification model
13 | # inherits the Model class
14 | class BinaryClassification(Model):
15 |     # method that initializes the binary classification model
16 |     # calls the super class constructor
17 |     # calls the binarize method from the helper class to binarize the target features
18 |     def __init__(self, alg, is_balanced, x_train, x_test, y_train, y_test):
19 |         self.is_balanced = is_balanced
20 |         super().__init__(x_train, x_test, Helper.binarize(y_train), Helper.binarize(y_test))
21 |         self.set_algorithm(alg)
22 |
23 |     # method that initializes the algorithm
24 |     def set_algorithm(self, alg):
25 |         if alg == "mlp":
26 |             if self.is_balanced == False:
27 |                 self.algorithm = MLPClassifier(activation='relu', hidden_layer_sizes=(200,),
28 |                                                 learning_rate='constant', solver='adam')
29 |             else:
30 |                 self.algorithm = MLPClassifier(activation='relu', hidden_layer_sizes=(50,50),
31 |                                                 learning_rate='adaptive', solver='adam', max_iter=300)
32 |         elif alg == "gb":
33 |             if self.is_balanced == False:
34 |                 self.algorithm = GradientBoostingClassifier(n_estimators=300, criterion='squared_error',
35 |                                                           loss='deviance', max_features='sqrt')
36 |             else:
37 |                 self.algorithm = GradientBoostingClassifier(n_estimators=300, criterion='friedman_mse',
38 |                                                           loss='log_loss', max_features='log2')
39 |         elif alg == "svm":
40 |             self.algorithm = svm.SVC(C=5, degree=3, gamma='scale', kernel='poly', probability=True)
41 |         elif alg == "rf":
42 |             if self.is_balanced == False:
43 |                 self.algorithm = RandomForestClassifier(n_estimators=100, criterion='log_loss', max_features='sqrt')
44 |             else:
45 |                 self.algorithm = RandomForestClassifier(n_estimators=200, criterion='gini', max_features='sqrt')
46 |         else:
47 |             raise ValueError("Invalid algorithm")
48 |
```

```

49     # method that fits the model
50     def fit(self):
51         if self.is_balanced == False:
52             self.fit_unbalanced_data()
53         else:
54             self.fit_balanced_data()
55
56     # method that returns the confusion matrix
57     def get_confusion_matrix(self):
58         return confusion_matrix(self.y_test, self.y_predict)
59
60     # method that returns the accuracy score
61     def get_accuracy(self):
62         return accuracy_score(self.y_test, self.y_predict)
63
64     # method that returns the precision score
65     def get_precision(self):
66         return precision_score(self.y_test, self.y_predict)
67
68     # method that returns the recall score
69     def get_recall(self):
70         return recall_score(self.y_test, self.y_predict)
71
72     # method that returns the f1 score
73     def get_f1_score(self):
74         return f1_score(self.y_test, self.y_predict)
75
76     # method that returns the roc auc score
77     def get_roc_auc_score(self):
78         return roc_auc_score(self.y_test, self.model.predict_proba(self.x_test)[:,-1])
79
80     # method that returns the roc curve
81     def get_roc_curve(self):
82         return roc_curve(self.y_test, self.model.predict_proba(self.x_test)[:,-1])
83
84     # method that returns the matthews corrcoeff score
85     def get_matthews_corrcoef(self):
86         return matthews_corrcoef(self.y_test, self.y_predict)
87
88     # method that returns the cohens kappa score
89     def get_cohen_kappa_score(self):
90         return cohen_kappa_score(self.y_test, self.y_predict)
91
92     # method that displays the measures
93     def display_matrix(self, title):
94         print("*****")
95         print(f"Confusion Matrix: {title}\n\n{self.get_confusion_matrix()}\n")
96         print(f"Accuracy: \t{self.get_accuracy():.2%}")
97         print(f"Precision: \t{self.get_precision():.3}")
98         print(f"Recall: \t{self.get_recall():.3}")
99         print(f"F1 Score: \t{self.get_f1_score():.3}")
100        print(f"Kappa: \t{self.get_cohen_kappa_score():.3}")
101        print(f"ROC AUC: \t{self.get_roc_auc_score():.3}")
102        print(f"MCC: \t{self.get_matthews_corrcoef():.3}")
103        print("*****")
104
105     # method that plots the confusion matrix and roc curve
106     def plot(self):
107         Plot.plot_confusion_matrix_binary(self.get_confusion_matrix(), f"Confusion Matrix binary
108                                         {self.algorithm.__class__.__name__}")
109         Plot.plot_roc_curve(self.get_roc_curve()[0], self.get_roc_curve()[1], self.get_roc_auc_score(),
110                            f"ROC Curve binary {self.algorithm.__class__.__name__}")
111

```

binary\_classification.py

```

1  # multi_classification.py
2  # Importing the libraries
3  import numpy as np
4  from model import Model
5  from helper import Helper
6  from plot import Plot
7  from sklearn import svm
8  from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier
9  from sklearn.neural_network import MLPClassifier
10 from sklearn.metrics import multilabel_confusion_matrix, classification_report
11
12 # The class MultiClassification which is used to build the multi classification model
13 # inherits the Model class
14 class MultiClassification(Model):
15
16     # method that initializes the multi classification model
17     # calls the super class constructor
18     def __init__(self, alg, is_blanced, x_train, x_test, y_train, y_test):
19         super().__init__(x_train, x_test, y_train, y_test)
20         self.is_blanced = is_blanced
21         self.set_algorithm(alg)
22         self.num_classes = len(np.unique(y_test))
23
24     # method that initializes the algorithm
25     def set_algorithm(self, alg):
26         if alg == "mlp":
27             if self.is_blanced == False:
28                 self.algorithm = MLPClassifier(activation='relu',hidden_layer_sizes=(50,50),
29                                                 learning_rate='adaptive',solver='adam',max_iter=200)
30             else:
31                 self.algorithm = MLPClassifier(activation='relu',hidden_layer_sizes=(200,),
32                                                 learning_rate='adaptive',solver='adam')
33         elif alg == "gb":
34             if self.is_blanced == False:
35                 self.algorithm = GradientBoostingClassifier(n_estimators=300,criterion='squared_error',
36                                                           loss='log_loss',max_features='sqrt')
37             else:
38                 self.algorithm = GradientBoostingClassifier(n_estimators=300,criterion='friedman_mse',
39                                                           loss='deviance',max_features='sqrt')
40         elif alg == "svm":
41             if self.is_blanced == False:
42                 self.algorithm = svm.SVC(C=5,degree=4,gamma='scale',kernel='poly' , probability=True)
43             else:
44                 self.algorithm = svm.SVC(C=3.5,degree=3, gamma='scale',kernel='poly', probability=True)
45         elif alg == "rf":
46             if self.is_blanced == False:
47                 self.algorithm = RandomForestClassifier(n_estimators=200,criterion='log_loss',max_features='sqrt')
48             else:
49                 self.algorithm = RandomForestClassifier(n_estimators=300,criterion='log_loss',max_features='sqrt')
50         else:
51             raise ValueError("Invalid algorithm")
52

```

```

53 # method that fits the model
54 def fit(self):
55     if self.is_balanced == False:
56         self.fit_unbalanced_data()
57     else:
58         self.fit_balanced_data()
59
60 # method that returns the multilabel confusion matrix
61 def get_multilabel_confusion_matrix(self):
62     return multilabel_confusion_matrix(self.y_test, self.y_predict)
63
64 def get_classification_report(self):
65     return classification_report(self.y_test, self.y_predict,
66                                 target_names=["passed", "failed", "error", "terminated"])
67
68 # method that returns the accuracy, precision, recall and f1 score
69 def get_accuracy_recall_precision_f1(self, arr):
70     tn = arr[0][0]
71     fp = arr[0][1]
72     fn = arr[1][0]
73     tp = arr[1][1]
74
75     precision = 0
76     recall = 0
77     accuracy = 0
78     f1 = 0
79     # ignore the divide by zero warning
80     np.seterr(divide='ignore', invalid='ignore')
81
82     if (tp+tn+fp+fn == 0):
83         accuracy = 0
84         #raise ZeroDivisionError("Accuracy is undefined when TP+TN+FP+FN=0")
85     else:
86         accuracy = (tp+tn)/(tp+tn+fp+fn)
87     if (tp == 0):
88         precision = 0
89     elif (tp+fp == 0):
90         precision = 0
91         #raise ZeroDivisionError("Precision is undefined when TP+FP=0")
92     #elif (fp == 0):
93     #    precision = "precision is not estimated when FP=0, the value is 1"
94     else:
95         precision = tp/(tp+fp)
96     if (tp + fn == 0):
97         recall = 0
98         #raise ZeroDivisionError("Recall is undefined when TP+FN=0")
99     #elif (fn == 0):
100     #    recall = "recall is not estimated when FN=0, the value is 1"
101     else:
102         recall = tp/(tp+fn)
103
104     f1 = 2 * (precision * recall) / (precision + recall)
105
106     return accuracy, precision, recall, f1
107
108 # method that returns the kappa and mcc
109 def get_kappa_mcc(self, arr):
110     tn = arr[0][0]
111     fp = arr[0][1]
112     fn = arr[1][0]
113     tp = arr[1][1]
114

```



```

115     # ignore the divide by zero warning
116     np.seterr(divide='ignore', invalid='ignore')
117
118     kappa_divisor = (tp+fp)*(fp+tn)+(tp+fn)*(fn+tn)
119     if (kappa_divisor == 0):
120         kappa = 0
121         #raise ZeroDivisionError("Kappa is undefined when divisor=0")
122     else:
123         kappa = (tp+tn)/kappa_divisor
124     mcc_divisor = np.sqrt((tp+fp)*(tp+fn)*(tn+fp)*(tn+fn))
125     if (mcc_divisor == 0):
126         mcc = 0
127         #raise ZeroDivisionError("MCC is undefined when divisor=0")
128     else:
129         mcc = (tp*tn-fp*fn)/mcc_divisor
130     return kappa,mcc
131
132 # method that displays the confusion matrix
133 def display_confusion_matrix(self,title):
134     print("*****")
135     cm = self.get_multilabel_confusion_matrix()
136     print(f"Confusion Matrix for {title}:\n",cm)
137     for i in range(self.num_classes):
138         print(f"Class: ({Helper.get_status(i)})")
139         accuracy, precision, recall, f1 = self.get_accuracy_recall_precision_f1(cm[i])
140         kappa, mcc = self.get_kappa_mcc(cm[i])
141         print(f"Accuracy:\t{accuracy:.2%}\nPrecision:\t{precision:.3}\nRecall:\t\t{
142             recall:.3}\nF1 Score:\t{f1:.3}\nKappa:\t\t{kappa:.3}\nMCC:\t\t{mcc:.3}")
143         #f"Cohen's Kappa Score:\t{kappa:.3}\nMatthews Correlation Coefficient:\t{mcc:.3}")
144     print("*****")
145
146 # method that plots the confusion matrix
147 def plot(self):
148     Plot.plot_confusion_matrix_multi_class(self.y_predict,self.y_test,
149         f"Confusion Matrix multi-class {self.algorithm.__class__.__name__}")
150
151 # method that displays the classification report
152 def display_report(self):
153     print("*****")
154     print("Classification Report:\n",self.get_classification_report())
155     print("*****")
156

```

multi\_classification.py

```

1 | # plot.py
2 | # The class Plot which is used to plot the data
3 | # import the libraries
4 | import numpy as np
5 | import pandas as pd
6 | import seaborn as sns
7 | from helper import Helper
8 | from matplotlib import pyplot as plt
9 |
10 | class Plot:
11 |
12 |     # method that plots the confusion matrix
13 |     @staticmethod
14 |     def plot_confusion_matrix_binary(cm,title):
15 |         cm = cm
16 |         sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
17 |         plt.title(title)
18 |         plt.xlabel('Predicted')
19 |         plt.ylabel('Actual')
20 |         plt.show()
21 |
22 |     # method that plots the roc curve
23 |     @staticmethod
24 |     def plot_roc_curve(fpr, tpr, roc_auc_score, title):
25 |         fpr = fpr
26 |         tpr = tpr
27 |         plt.plot(fpr, tpr, label=f'ROC Curve: {roc_auc_score:.3}')
28 |         plt.plot([0,1],[0,1], 'k--', label='Random Guess')
29 |         plt.xlabel('False Positive Rate')
30 |         plt.ylabel('True Positive Rate')
31 |         plt.legend(loc='lower right')
32 |         plt.title(title)
33 |         plt.show()
34 |
35 |     # method that plots the confusion matrix for multi-class classification
36 |     @staticmethod
37 |     def plot_confusion_matrix_multi_class(y_pred, y_test, title):
38 |         fig, ax = plt.subplots(figsize=(7,5))
39 |         from sklearn.metrics import ConfusionMatrixDisplay
40 |         ConfusionMatrixDisplay.from_predictions(y_test, y_pred, ax=ax, cmap='Blues')
41 |         ax.xaxis.set_ticklabels(Helper.get_the_classes())
42 |         ax.yaxis.set_ticklabels(Helper.get_the_classes())
43 |         _=ax.set_title(title)
44 |         plt.show()
45 |
46 |     # method that plots the correlation matrix
47 |     @staticmethod
48 |     def plot_correlation_matrix(data, title):
49 |         corr = np.corrcoef(data, rowvar=False)
50 |         sns.heatmap(corr, annot=True, cmap='Blues')
51 |         plt.title(title)
52 |         plt.show()
53 |
54 |     # method that plots the distribution of the classes both binary and multi-class classification
55 |     @staticmethod
56 |     def plot_class_distribution(x,str):
57 |         class_count = pd.DataFrame(x).value_counts(sort = False)
58 |         class_count.plot(kind = 'bar', rot=0)
59 |         plt.title(f"{str}")
60 |         plt.xlabel("Classes")
61 |         plt.ylabel("Frequency")
62 |         if len(class_count) <= 2 :
63 |             plt.xticks(range(2), Helper.get_the_binary_classes())
64 |         else:
65 |             plt.xticks(range(4), Helper.get_the_classes())
66 |         plt.show()
67 |
68 |

```

```

69 # method that plots the histogram of each feature
70 @staticmethod
71 def plot_one_feature_histogram(data, title, xlabel, ylabel, bins=10):
72     plt.barh(data,np.arange(len(data)))
73     #plt.hist(data, bins=bins)
74     plt.title(title)
75     plt.xlabel(xlabel)
76     plt.ylabel(ylabel)
77     plt.show()
78
79 # method that plots the histogram of all features
80 @staticmethod
81 def plot_features_histogram(data):
82     num_features = data.shape[1]
83     num_rows = (num_features + 1) // 2
84     fig, axes = plt.subplots(num_rows, 2, figsize=(10, 6 * num_rows))
85     for i in range(num_features):
86         ax = axes[i // 2, i % 2]
87         ax.hist(data.iloc[:, i], bins=25)
88         #ax.set_title(f"feature {i}")
89         #ax.title.set_position([.5, 1.05])
90         #ax.set_xlabel(xlabel)
91         #ax.set_ylabel(ylabel)
92     for i in range(num_features, num_rows * 2):
93         fig.delaxes(axes.flatten()[i])
94     plt.show()
95

```

plot.py

```

1 | # helper.py
2 # A helper class for the main program
3 class Helper:
4
5     # method that binarizes the data
6     @staticmethod
7     def binarize(y_data):
8         for i in range(len(y_data)):
9             if y_data[i] == 0:
10                y_data[i] = 0
11            else:
12                y_data[i] = 1
13        return y_data
14
15
16 # method that returns the classes of the target feature as a list
17 @staticmethod
18 def get_the_classes():
19     return ['passed', 'failed', 'error', 'terminated']
20
21 # method that returns the binary classes of the target feature as a list
22 @staticmethod
23 def get_the_binary_classes():
24     return ['passed', 'rest']
25
26
27 # method that normalizes the data using the min max scaler
28 # returns the normalized train and test data
29 @staticmethod
30 def scale_min_max(x_train, x_test):
31     from sklearn.preprocessing import MinMaxScaler
32     scaler = MinMaxScaler(feature_range=(0, 1))
33     s_train = scaler.fit_transform(x_train)
34     s_test = scaler.transform(x_test)
35     return s_train, s_test

```

```

36
37     # method that returns the status of the target feature
38     @staticmethod
39     def get_status(i):
40         if i ==0:
41             x = "passed"
42         elif i ==1:
43             x = "failed"
44         elif i ==2:
45             x = "error"
46         elif i ==3:
47             x = "terminated"
48         else:
49             x = "other"
50         return x

```

helper.py

```

1 | # file_utility.py
2 | # FileUtility class to read and write files
3 | import json
4 | import pickle
5 |
6 | class FileUtility:
7 |
8 |     # method that reads file and returns x_data and y_data
9 |     @staticmethod
10 |    def read_file(filename):
11 |        try:
12 |            x_data = []
13 |            y_data = []
14 |            with open(filename, 'r') as f:
15 |                lines = f.readlines()
16 |                for line in lines[1:lines.__len__():]:
17 |                    line = line.strip().split(',')
18 |                    x_data.append([float(line[0]), float(line[1]), float(line[2]),float(line[3]),
19 |                                   float(line[4]),float(line[5]),float(line[6]),float(line[7])])
20 |                    y_data.append(int(line[8]))
21 |            return x_data, y_data
22 |        except IOError as e:
23 |            print("read method error", e.strerror)
24 |
25 |    # method that writes x_data and y_data to file
26 |    @staticmethod
27 |    def write_to_file(x_train, x_test, filename):
28 |        try:
29 |            with open(filename, 'w') as f:
30 |                for i in range(len(x_train)):
31 |                    for j in range(len(x_train[i])):
32 |                        f.write(str(x_train[i][j]))
33 |                        f.write(',')
34 |                    #f.write(str(y_train[i]))
35 |                    f.write('\n')
36 |                for i in range(len(x_test)):
37 |                    for j in range(len(x_test[i])):
38 |                        f.write(str(x_test[i][j]))
39 |                        f.write(',')
40 |                    #f.write(str(y_test[i]))
41 |                    f.write('\n')
42 |        except IOError as e:
43 |            print("write method error", e.strerror)
44 |

```

```
45
46 # method that saves the model
47 @staticmethod
48 def save_model(model, filename):
49     try:
50         with open(f"modeled/{filename}", 'wb') as f:
51             pickle.dump(model, f)
52     except IOError as e:
53         print("save error", e.strerror)
54
55 # method that loads the model
56 @staticmethod
57 def load_model(filename):
58     try:
59         with open(f"modeled/{filename}", 'rb') as f:
60             return pickle.load(f)
61     except IOError as e:
62         print("load error", e.strerror)
63
64 # method that saves texts as json format
65 @staticmethod
66 def save_text_json(data, filename):
67     try:
68         with open(f"parameters/{filename}", 'w') as f:
69             json.dump(data, f)
70     except IOError as e:
71         print("save json error", e.strerror)
```

file\_utility.py

```

1 | # parameters_tunning.py
2 | # import libraries
3 | from sklearn import svm
4 | from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier
5 | from sklearn.neural_network import MLPClassifier
6 | from sklearn.model_selection import GridSearchCV
7 | from sklearn.model_selection import RepeatedStratifiedKFold
8 | from imblearn.over_sampling import SMOTE
9 | from imblearn.under_sampling import RandomUnderSampler
10 | from imblearn.pipeline import Pipeline
11 |
12 | # The class ParametersTunning which is used to tune the parameters of the model
13 | class ParametersTunning:
14 |     # method that initializes the parameters tunning
15 |     def __init__(self, alg_name, x_train, y_train):
16 |         self.x_train = x_train
17 |         self.y_train = y_train
18 |         self.algorithm = None
19 |         self.parameters = {}
20 |         self.set_algorithm(alg_name)
21 |
22 |     # method that initializes the algorithm and the parameters
23 |     def set_algorithm(self, alg):
24 |         if alg == "mlp":
25 |             self.algorithm = MLPClassifier()
26 |             self.parameters = {'hidden_layer_sizes': [(50,),(100,),(200,),(50,50)], 'activation':
27 |                               ['relu','tanh','sigmoid'], 'learning_rate': ['constant','adaptive'],
28 |                               'solver': ['sgd', 'adam'],'max_iter':[200,300,400,500]}
29 |
30 |         elif alg == "gb":
31 |             self.algorithm = GradientBoostingClassifier()
32 |             self.parameters = { 'n_estimators':[100,200,300],'criterion' :
33 |                               ['friedman_mse','squared_error'], 'max_features': ['sqrt', 'log2','auto'],
34 |                               'loss': ['log_loss', 'deviance','exponential']}
35 |
36 |         elif alg == "svm":
37 |             self.algorithm = svm.SVC()
38 |             self.parameters = { 'C': [0.1,1,1.5,2,2.5,3,4,4.5,5], 'kernel': [ 'rbf','poly'],
39 |                               'degree' : [3,4],'gamma': ['scale', 'auto']}
40 |
41 |         elif alg == "rf":
42 |             self.algorithm = RandomForestClassifier()
43 |             self.parameters = { 'n_estimators':[100,200,300],'class_weight' :
44 |                               ['balanced_subsample','balanced'], 'max_features': ['sqrt', 'log2','None'],
45 |                               'criterion' : ['gini','entropy','log_loss']}
46 |         else:
47 |             raise ValueError("Invalid algorithm")

```

```
46 # method that tunes the parameters of the model and returns the best parameters
47 def tune_parameters(self, is_balanced):
48     if is_balanced == False:
49         return self.tune_parameters_unbalanced()
50     else:
51         return self.tune_parameters_balanced()
52
53 # method that tunes the parameters of the model with unbalanced data and returns the best parameter
54 def tune_parameters_unbalanced(self):
55     cv = RepeatedStratifiedKFold(n_splits=5, n_repeats=3, random_state=1)
56     #grid_search = GridSearchCV(estimator=self.algorithm, param_grid=self.parameters, cv=cv)
57     grid_search = GridSearchCV(estimator=self.algorithm, param_grid=self.parameters)
58     grid_search.fit(self.x_train, self.y_train)
59     print(grid_search.best_params_)
60     return grid_search.best_params_
61
62
63 # method that tunes the parameters of the model with balanced data and returns the best parameters
64 def tune_parameters_balanced(self):
65     oversample = SMOTE()
66     undersample = RandomUnderSampler()
67     steps = [('o', oversample), ('u', undersample)]
68     pipeline = Pipeline(steps=steps)
69     x_train_res, y_train_res = pipeline.fit_resample(self.x_train, self.y_train)
70     cv = RepeatedStratifiedKFold(n_splits=5, n_repeats=3, random_state=1)
71     grid_search = GridSearchCV(estimator=self.algorithm, param_grid=self.parameters, cv=cv)
72     grid_search.fit(x_train_res, y_train_res)
73     print(grid_search.best_params_)
74     return grid_search.best_params_
75
```

parameters\_tunning.py

