**2024 Sem 1**

**CITS5508 Machine Learning**

**Assignment 2**

Mila Zhang (22756463)

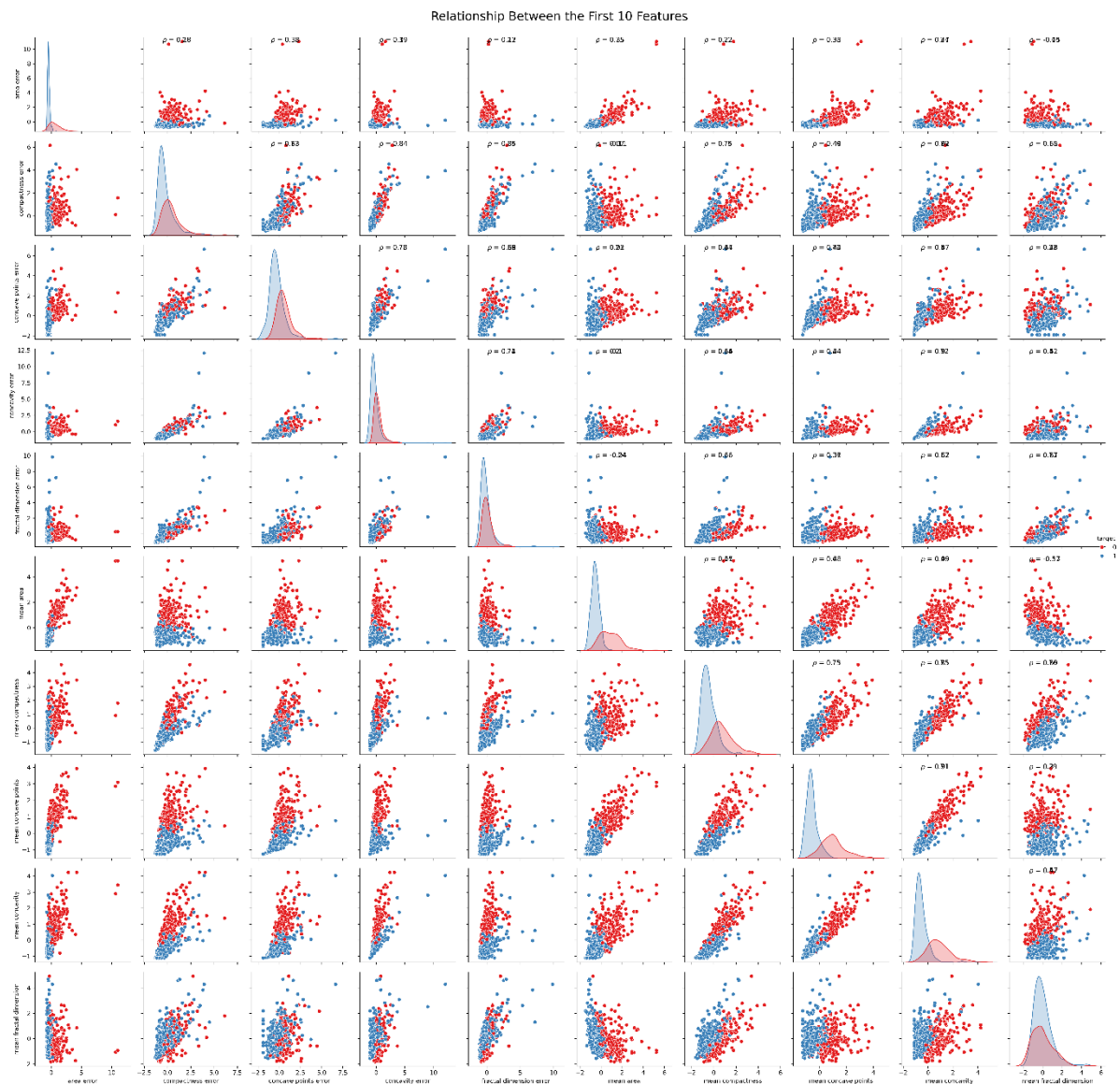Due date: Friday, 3 May, 8pm

# Table of Contents

# 1. Initial Inspections and Data Preprocessing
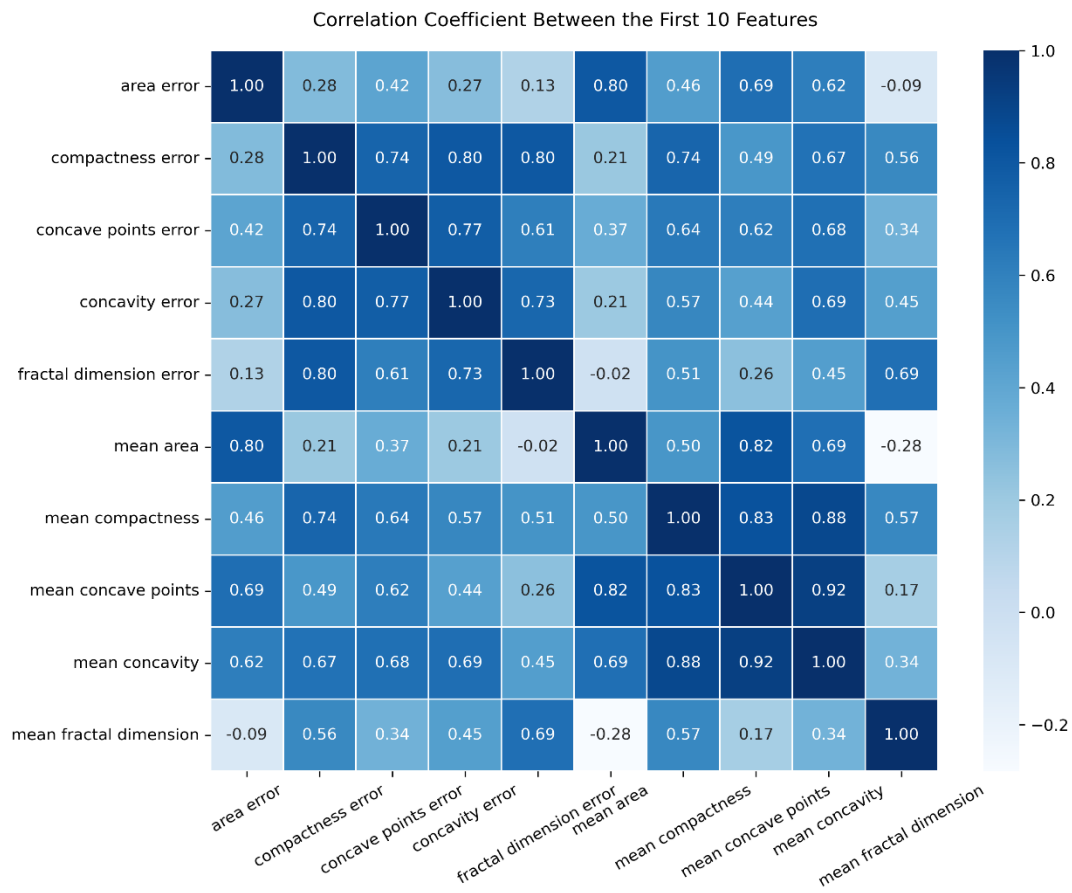
## 1.1 Task D1



Relationship Between the First 10 Features

## 1.2 Task D2

Some pairs of features are highly correlated and the distribution of data is almost linear, while most pairs of features do not exhibit multicollinearity. Some clusters of data points are well separated regarding different target labels, while some clusters are mixed up and cannot be separated well. Instances that are located far away from the clusters might be outliers. For features that are highly correlated, one of them in each pair can be removed because they have a strong relationship with each other and will essentially provide redundant information for splitting the data.

## 1.3 Task D3

Correlation Coefficient Between the First 10 Features
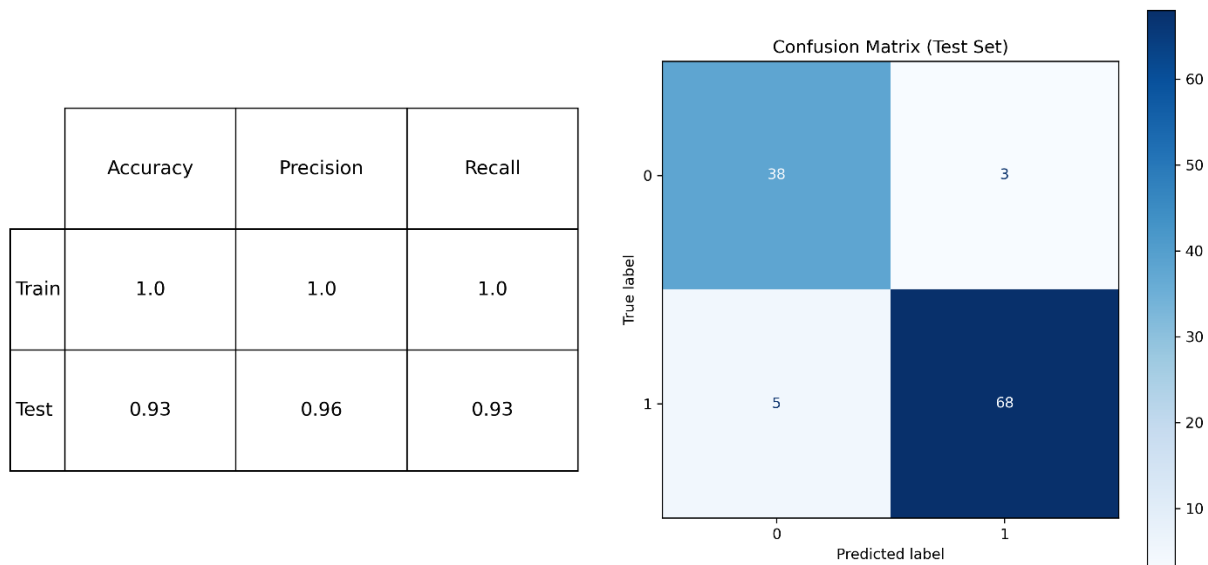


## 1.4 Task D4

The correlation coefficients support the previous observations. The pairs of features that exhibit a linear distribution have high correlation coefficient values above 0.88, corresponding to the blocks with darker shades in the correlation matrix.

# 2. Decision Tree Models with Default Hyperparameters

## 2.1 Task D6

Decision Tree with Default Hyperparameters

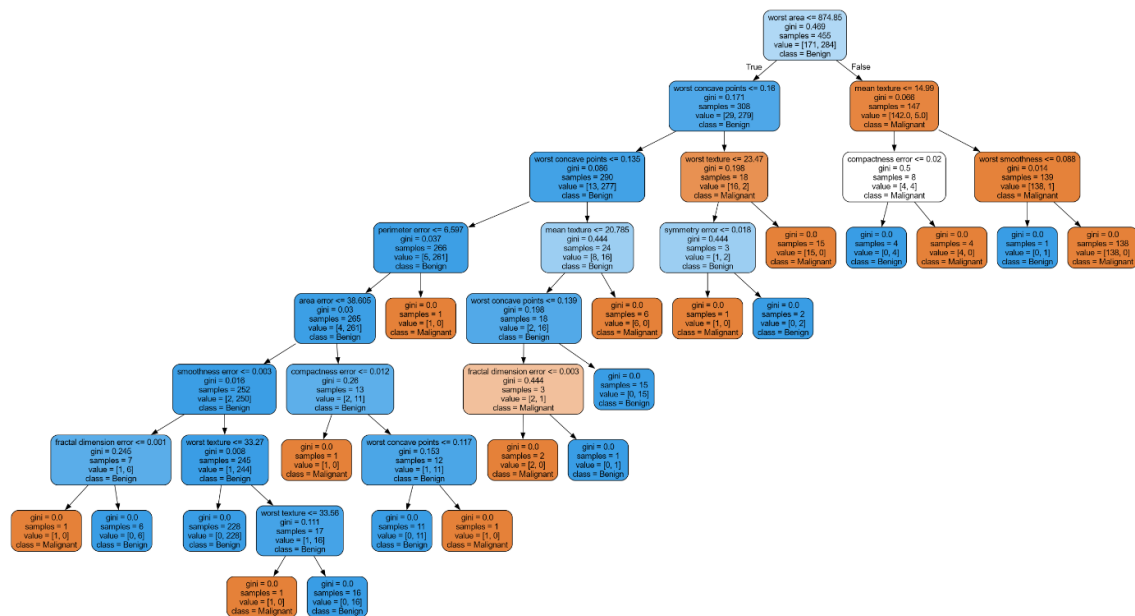|       | Accuracy | Precision | Recall |
|-------|----------|-----------|--------|
| Train | 1.0      | 1.0       | 1.0    |
| Test  | 0.93     | 0.96      | 0.93   |



Confusion Matrix (Test Set)

## 2.2 Task D7

The classifier, when configured with default hyperparameters is overfitting. This overfitting is evidenced by the significant discrepancy in performance between the training and test datasets (e.g. the accuracy of the training set is 1.0 while that of the test set is just 0.93).

Decision trees inherently operate by recursively splitting the data until each branch terminates in a leaf node, where either all data in the node belong to a single class or no further splits can reduce Gini impurity.

Relying on default hyperparameters, this method optimises performance exclusively for the training dataset and does not necessarily generalise well on unseen data. A fully-grown decision tree can perfectly classify all training data points, capturing all the noise and anomalies in the training data. These anomalies might not represent the actual trends or patterns in the overall data population, leading to poorer generalisation on new unseen test data.

## 2.3 Task D8



## 2.4 Task D9

There are 9 levels in this decision tree model. This diagram helped to confirm the classifier has an overfitting issue. According to the diagram, it is observed that all samples in each leaf node belong to the same class, and each leaf node has a Gini impurity of zero. Additionally, the size of leaf nodes is very small. This is an interpretable model because the cause of each decision is demonstrated clearly in the tree and thus provides a transparent and straightforward nature of the decision-making process.

# 2.5 Task D10

Accuracy Scores Over Different Random States



Precision Scores Over Different Random States



Recall Scores Over Different Random States

According to the figures, the scores range from 0.91 to 1.00, demonstrating high consistency with some variation among the different models.

The observed variability across different runs with different seeds results from their inherent sensitivity to the specific makeup of the training data used. Decision trees can produce different structures and therefore slightly different performance metrics depending on how the nodes are split at each decision point in the tree.

Despite this variability, the models demonstrate a high degree of consistency in performance, suggesting that the decision tree has captured the underlying patterns in the data effectively across all seeds.

## 2.6 Task D11



Accuracy Scores Over Different Splits



Precision Scores Over Different Splits



Recall Scores Over Different Splits

The accuracy and precision scores over different split ratios exhibit an overall tendency of increase as the training set ratio gets larger. This performance behaves as expected because the model is better trained when fed more training data generally. However, the recall scores increase then drop then increase again which is not as optimal as I expected.
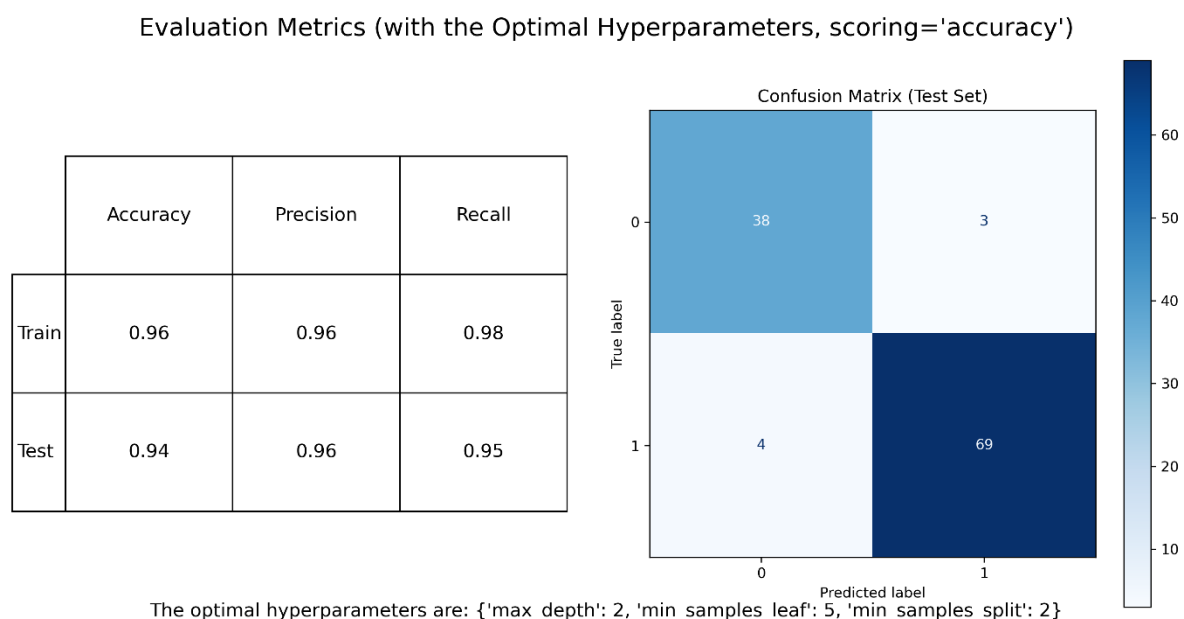
# 3. Decision Tree Models with Optimal Hyperparameters

## 3.1 Task D12

Evaluation Metrics (with the Optimal Hyperparameters, scoring='accuracy')

|       | Accuracy | Precision | Recall |
|-------|----------|-----------|--------|
| Train | 0.96     | 0.96      | 0.98   |
| Test  | 0.94     | 0.96      | 0.95   |

Confusion Matrix (Test Set)

The optimal hyperparameters are: {'max_depth': 2, 'min_samples_leaf': 5, 'min_samples_split': 2}

## 3.2 Task D13

After fine-tuning the hyperparameters, the issue of overfitting is addressed and the model has better generalisation. This aligns with my expectations.

The scores of the training set are no longer 1.0 as obtained in D6, and the scores of the fine-tuned model on the test set are generally higher than that of the model without fine-tuning. Now the discrepancies in scores between the training set and the test set are smaller compared to those of the model without fine-tuning. These have shown that fine-tuning

reduced the issue of overfitting in D6. However, after fine-tuning, the performance on test set is lower than the model without fine-tuning.

## 3.3 Task D14

Confusion Matrices Over Different Scoring Options (Test Set)



| Optimal Hyperparameter\ Scoring Method | max_depth | min_samples_leaf | min_samples_split |
|---|---|---|---|
| accuracy | 4 | 5 | 2 |
| precision | 5 | 2 | 2 |
| recall | 3 | 5 | 2 |

When using different scoring methods, the optimal hyperparameters are different, especially "max_depth" and "min_samples_leaf". When hyperparameters are tuned to optimise different scoring metrics, they adjust the model to emphasise different aspects of the data.

The confusion matrices resemble each other which indicate similar model performance on the test set. This suggests that the models' predictions are robust across various decision boundaries imposed by the hyperparameters. This similarity can indicate that the model is consistently identifying true positives, true negatives, false positives, and false negatives, regardless of the hyperparameter settings prioritised during tuning.

# 4. Decision Tree Models with a Reduced Feature Set

## 4.1 Task D15

Feature Importances (Descending Order)

| | |
|---|---|
| worst area | 0.84 |
| worst concave points | 0.13 |
| mean smoothness | 0.02 |
| area error | 0.0 |
| compactness error | 0.0 |
| concave points error | 0.0 |
| concavity error | 0.0 |
| fractal dimension error | 0.0 |
| mean area | 0.0 |
| mean compactness | 0.0 |
| mean concave points | 0.0 |
| mean concavity | 0.0 |
| mean fractal dimension | 0.0 |
| mean symmetry | 0.0 |
| mean texture | 0.0 |
| perimeter error | 0.0 |
| smoothness error | 0.0 |
| symmetry error | 0.0 |
| texture error | 0.0 |
| worst compactness | 0.0 |
| worst concavity | 0.0 |
| worst fractal dimension | 0.0 |
| worst smoothness | 0.0 |
| worst symmetry | 0.0 |
| worst texture | 0.0 |

## 4.2 Task D16

Retained features:  ['mean smoothness', 'worst area', 'worst concave points']

Removed features:  ['area error', 'compactness error', 'concave points error', 'concavity error', 'fractal dimension error', 'mean area', 'mean compactness', 'mean concave points', 'mean concavity', 'mean fractal dimension', 'mean symmetry', 'mean texture', 'perimeter

error', 'smoothness error', 'symmetry error', 'texture error', 'worst compactness', 'worst concavity', 'worst fractal dimension', 'worst smoothness', 'worst symmetry', 'worst texture']

Total feature importance value after dimension reduction step:  1.0.

## 4.3 Task D17

### Comparing Fine-tuned Decision Trees' Performance

**Using All Features**

|  | Accuracy | Precision | Recall |
|---|---|---|---|
| Train | 0.96 | 0.96 | 0.98 |
| Test | 0.94 | 0.96 | 0.95 |

**Using Selected Features**

|  | Accuracy | Precision | Recall |
|---|---|---|---|
| Train | 0.96 | 0.96 | 0.98 |
| Test | 0.94 | 0.96 | 0.95 |

### Confusion Matrices of Fine-tuned Desicion Trees (Test Set)



Using All Features



Using Selected Features

## 4.4 Task D18

I did not repeat the cross-validation process to find the optimal hyperparameters when using the reduced set of features. The only variable is the feature set (whether reduced or not). This enables comparison of using the full feature set and the reduced feature set.
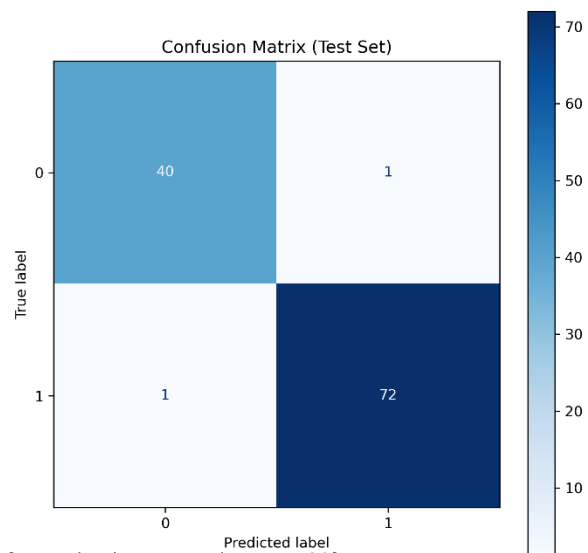
Reducing the number of features does not have any impact on the model performance. According to D15, those reduced features all have a feature importance of 0, which means these features do not reduce impurity nor contribute to the prediction in each split.

# 5. Random Forest Models

## 5.1 Task D19

Random Forest (with the Optimal Hyperparameters, scoring='accuracy')

|       | Accuracy | Precision | Recall |
|-------|----------|-----------|--------|
| Train | 0.99     | 0.98      | 1.0    |
| Test  | 0.98     | 0.99      | 0.99   |

Confusion Matrix (Test Set)

|             | Predicted 0 | Predicted 1 |
|-------------|-------------|-------------|
| True 0      | 40          | 1           |
| True 1      | 1           | 72          |

The optimal hyperparameters are: {'max_depth': 4, 'n_estimators': 20}

## 5.2 Task D20

The performance of the Random Forest Model is higher than the one obtained in D12. The scores of accuracy, precision, and recall are all higher than performance obtained in D12. The FN and FP misclassifications are fewer than that of model in D12. This result aligns with my expectations because Random Forests are an ensemble method that relies on the collective decision-making of multiple decision trees to improve accuracy and robustness. It effectively overcomes some of the limitations of individual decision trees and has higher overall performance.

## 5.3 Task D21

I do not think these models are good enough and can be trusted to be used in real life. These models were trained based on a small set of data points. And through inspection in D1 we can observe anomalies and outliers existing in the dataset. These affects the prediction accuracy and its generalisation capability for unseen data.

A more complex model might not be necessary due to the high recall, accuracy, precision scores and good performance of confusion matrix, especially for the Random Forest model in D19.

I do not think using a machine learning algorithm for this task at this stage is a good idea. This decision process should not be automated. Identifying breast cancer is a crucial task and we do not want to misclassify any case. It is more suitable for being used to assist the manual diagnosis of breast cancer. We should not rely solely on these machine learning models.

The dataset used is a small dataset with less than 600 instances which is small. Also, there are anomalies and outliers as observed in D1. It is better to collect more data and perform data cleaning before moving to the modelling stage.