# Hard Disk Drive Reliability: A Comparative Study of Supervised Machine Learning Algorithms for Predicting Drive Failure

Alistair McLean, Roy Sterritt

School of Computing and Mathematics
Ulster University
Belfast, Northern Ireland
mclean-a13@ulster.ac.uk | r.sterritt@ulster.ac.uk

*Abstract*— Unexpected downtime and IT system outages can cost organisations millions of dollars in lost revenue, loss of opportunity, and negatively impacted reputation. Third party cloud services and infrastructure are commonly used by individuals and organisations as it offers the ability to create highly scalable applications without the huge cost of purchasing and maintaining their own hardware facility. Consequently, cloud service providers are challenged with ensuring that their data centres are reliable, as they have shared responsibility for the applications deployed in them. One of the most common causes of IT system failure in data centres is failing Hard Disk Drives (HDDs). It is proposed that if data centres were able to accurately predict imminent HDD failures, then appropriate action could be taken to prevent potential outages. This paper investigates the relationship between Self-Monitoring, Analysis, and Reporting Technology (SMART) attributes and HDD failure, implementing supervised machine learning methods to predict drive failure at various prediction horizons. Random Forest and XGBoost classifiers are observed to achieve the best prediction performance, with the Area Under the Receiver Operating Characteristic Curve (AUROC) calculated at 0.9185±0.0066 and 0.9162±0.0066 respectively at the shortest prediction horizon (0-24 hours prior to failure). Reallocated sectors count (SMART 5), reported uncorrectable errors (SMART 187), current pending sector count (SMART 197), and uncorrectable sector count (SMART 198) were found to be the most important SMART attributes for HDD failure prediction.

*Keywords-hard disk drive; hdd reliability; machine learning; failure prediction.*

LIST OF ABBREVIATIONS

| Abbreviation | Definition |
|---|---|
| **AUROC** | Area Under the Receiver Operating Characteristic Curve |
| DT | Decision Tree |
| FAR | False Alarm Rate |
| FDR | False Discovery Rate |
| FPR | False Positive Rate |
| HDD | Hard Disk Drive |
| k-NN | K-Nearest Neighbour |
| LR | Logistic Regression |
| ML | Machine Learning |
| MLP | Multi-Layer Perceptron |
| RF | Random Forest |
| RUL | Remaining Useful Life |
| SLA | Service Level Agreement |
| SMART | Self-Monitoring, Analysis, and Reporting Technology |
| TPR | True Positive Rate |
| XGB | XGBoost |

## I. INTRODUCTION

Unexpected downtime or outages of IT systems can have major consequences for businesses and users. It is reported that a single outage can cost an organisation millions of dollars through the loss of revenue, loss of opportunities, and diminished reputation, and that this cost impact is increasing year on year [1]. Many organisations today use cloud computing as part of their products and services, reducing the need for purchasing and maintaining their own IT infrastructure while improving the scalability of their applications. As cloud adoption continues to grow, cloud service providers have shared responsibility for their users' applications and are tasked with providing highly available services and reliable IT infrastructure. Application downtime or system unavailability can be detrimental for both the cloud provider and the cloud user, with cloud providers liable to fines for breaches of Service Level Agreements (SLAs).

Data centre outages occur from time to time, which may result in application unavailability or system downtime. The most common causes of data centre outages are on-site power-related problems, such as generator or grid failures, followed by network problems and IT system failure [2]. With respect to IT system failures, which include hardware and software failures, Hard Disk Drives (HDDs) are believed to be one of the main offenders of causing problems. HDDs are one of the most replaced hardware components and one of the least reliable, with [3] reporting that 78% of faults or replacements are attributable to hard disks. Another investigation [4] of the data centres of a major internet service organisation observed that 82% of hardware failure tickets were attributable to HDDs. Therefore, data centres could potentially improve their reliability by monitoring the health of HDDs in their IT estate and take appropriate action before a drive failure occurs.

Self-Monitoring, Analysis, and Reporting Technology (SMART) was developed in 1995 and is commonly used by manufacturers today, providing measurements collected by sensors within HDDs to report on various indicators of reliability and drive health. SMART attributes are numbered from 1 to 255 giving raw and normalised values of each measurement. For example, SMART 9 reports the power-on hours, the total count of hours that the drive has been in a power-on state across its operational lifetime. SMART 190 and 194 give measurements of internal temperature within the HDD unit. SMART 240 records the total time in hours

that the read-write head has been operating to position itself across the surface of the disk, while SMART 241 and 242 record the lifetime data (in bytes) written to, and read from, the drive respectively. Other attributes count the number of errors, such as SMART 1, which records the rate of read errors between the read-write head and the disk surface.

The aim of this paper is to investigate the relationship between SMART metrics and HDD failure, and to assess if Machine Learning (ML) methods can accurately predict imminent drive failure using reported SMART attribute measurements. Additionally, a comparative study will determine the best ML algorithm for application in this problem domain.

The remainder of the paper is organised as follows: Section II highlights existing work in the field, describing research approaches and their effectiveness for HDD failure prediction. Section III describes the dataset analysis and ML methods used in this paper. Section IV reports the results and performance of the ML implementation. Section V discusses the limitations of the work, highlighting improvements and opportunities for extension. Section VI summarises the research and concludes the paper.

## II. EXISTING WORK

"Autonomic systems are examples of accelerated AI automation. They are self-managing physical or software systems, performing domain-bounded tasks that exhibit three fundamental characteristics: autonomy, learning and agency. When traditional AI techniques aren't able to achieve business adaptability, flexibility and agility, autonomic systems can be successful in helping with implementation. Autonomic systems will take five to ten years until mainstream adoption but will be transformational to organizations" [5]. This work on supervised Machine Learning for HDD failure prediction fits generally with that Autonomic vision [6][7].

The ability to predict HDD failures would allow data centres to mitigate against potential outages by proactively replacing drives before they reach a state of failure. Unsurprisingly, there have been many works of research investigating and attempting to address this problem. Machine learning and probabilistic techniques are popular among researchers, applying traditional ML classification and regression methods, Bayesian networks, deep learning, or combining multiple methods with ensemble learning.

The approach taken in [8] uses SMART attributes to create HDD failure prediction models using classification and regression trees. Their experiments show the classification tree model was able to successfully predict 95% of failures with a False Alarm Rate of less than 0.1% when applied to a real-world data centre containing 25,792 HDDs. Additionally, they propose a regression tree model to evaluate the health status of the drives, where the probability of a fault occurring is predicted. Also using a tree-based model, [9] takes a binary classification approach to predict the health of HDDs in Meta's Tectonic storage fleet. The SMART metrics of 53,000 failed HDDs were used, alongside a random sample of non-failing drives, to categorise the HDDs as healthy or unhealthy at 1 and 30

days prior to failure. Their XGBoost classifier showed limited prediction performance, achieving low precision when applied to unseen data from a different time window. However, they report noticeable improvements when using the difference, or delta, between SMART measurements over time as opposed to using the singular measurements from a set prediction horizon.

Highlighting the limitation of using SMART attributes with their default thresholds to detect failing HDDs, [10] proposes a failure prediction method using a Bayesian network to provide Remaining Useful Life (RUL) estimates of drives. Using a subset of SMART attributes and their temporal trends, the proposed Bayesian Network for Failure prediction in HDDs (BNFH) was applied to a dataset containing 49,056 drives from Backblaze's data centres. Their evaluation showed the model outperformed standard reliability-based methods and other Bayesian network-based methods presented in [11]; and achieved similar relative accuracy to a Recurrent Neural Network presented in [12]. The work in [13] utilises ensemble learning to create a Combined Bayesian Network (CBN), where the learning results from four individual classifiers are combined to predict the remaining time before a drive fails. The individual classifiers used backpropagation neural networks, evolutionary neural networks, support vector machines, and classification tree methods. Experimental results indicate the CBN performs similarly to the classification tree model and outperforms the other models. However, the CBN has additional benefit over the classification tree model by indicating when the drive will fail, not just that it will fail.

Other research papers propose deep learning methods for HDD failure prediction. The work in [14] uses bidirectional LSTM models with multi-day lookback periods to learn the temporal progression of key health indicators present in SMART data. The proposed model achieved 96.4% accuracy in predicting HDD failure for a 15-day lookback period, outperforming a standard LSTM implementation. However, due to the inconsistency in SMART measurements recorded by different HDD manufacturers and models, the data used in this work only related to a single Seagate model (ST4000DM000) over the course of 9 months. Another deep learning approach, presented in [15], proposes a model based on Gated Recurrent Unit (GRU) neural networks and TimeGAN adversarial networks to analyse the temporal sequences of SMART attributes in HDDs, while addressing data imbalance issues. Their proposed approach achieved an average failure detection rate of 95% and a false alarm rate of 0.2%. This work also only applies to a single Seagate drive (ST6000DX000).

While existing work has achieved success in HDD failure prediction, it is not always clear which method performs best in this problem domain. The listed works in this section apply numerous ML algorithms to HDD SMART metrics using data centre drives. However, the data centres, drive manufacturers, drive models, and timeframes within the datasets will vary from paper to paper. Therefore, it is not necessarily viable to make direct comparisons. As such, the purpose of this paper is to apply multiple machine learning methods to the same dataset of operational data centre hard

drives. The comparative study will derive insight into the best performing methods for HDD failure prediction.

## III. METHODOLOGY

The dataset used in this paper was obtained from Backblaze [16], which records the daily SMART metrics of HDDs within their data centres. Each row in the dataset contains the date, serial number, model, capacity (in bytes), and raw and normalised values for various SMART metrics reported by each HDD. Additionally, a failure column records a binary value indicating if the drive is functional (0), or if the entry represents the last operational day before the drive failed (1). This paper uses the reported data across a 10-year period, from 1$^{st}$ January 2014 to 31$^{st}$ December 2023.

### A. Data Exploration

Initial analysis of the dataset was conducted to gain a better understanding of the size and complexity of data, the HDD models that are reported on, and the prevalence of drive failures. Over the 10-year period, the dataset contained over 450 million rows, reporting the daily SMART metrics of 388,485 unique hard disk drives (identified by their serial number). In that time, 21,356 rows indicated that an HDD had failed. It appears that the data centre organisation preferred to replace failed drives, rather than repair them, as only a very small proportion of drives (0.005%) were seen to fail more than once. Table I shows that most of the drives (94.5%) did not experience failure, and nearly all failed drives only failed once.

Using the model number provided in the dataset, it was possible to analyse the failure rate with respect to drive manufacturer and drive model. Seagate models accounted for the largest proportion of drives in the data centre throughout the years. Hitachi models also made up a large proportion in the earlier years, but their presence almost entirely disappeared by 2018. Other drive manufacturers present in the data include Toshiba, HGST, and WDC. Of the 193,378 Seagate drives, 16,177 resulted in failure accounting for 75.75% of all failures in the dataset and indicates an 8.37% failure rate for all Seagate HDDs. Table II shows the top 10 models with the highest number of failures, indicating the model failure rate and their proportional contribution to the overall failures present in the data. Models prefixed with 'ST' are Seagate drives and Table II shows that 7 of the top 10 most failing drives belong to this manufacturer.

One of the known issues with SMART attributes is that manufacturers do not always use them equally, as mentioned in [14]. The same SMART attribute may be used to report

TABLE I. DRIVE FAILURE COUNTS ACROSS ALL HDDS IN THE DATASET SHOWING PROPORTION OF DRIVES WITH MULTIPLE FAILURES

| Number of Failures | Number of Drives | % of HDDs | % of Failures |
|---|---|---|---|
| 0 | 367,147 | 94.51 | - |
| 1 | 21,320 | 5.49 | 99.92 |
| 2 | 18 | 0.005 | 0.08 |

TABLE II. TOP 10 HDD MODELS WITH THE HIGHEST NUMBER OF FAILURES

| Model | Total HDDs | Total Failures | % of All Failures | Model Failure % |
|---|---|---|---|---|
| ST4000DM000 | 36,983 | 5,602 | 26.23 | 15.15 |
| ST12000NM0007 | 38,838 | 2,106 | 9.86 | 5.42 |
| ST8000NM0055 | 15,680 | 1,718 | 8.04 | 10.96 |
| ST3000DM001 | 4,354 | 1,454 | 6.81 | 33.39 |
| ST12000NM0008 | 20,836 | 1,349 | 6.32 | 6.47 |
| MG07ACA14TA | 39,292 | 1,173 | 5.49 | 2.99 |
| ST8000DM002 | 10,305 | 1,037 | 4.86 | 10.06 |
| HUH721212ALN604 | 11,166 | 600 | 2.81 | 5.37 |
| HMS5C4040BLE640 | 16,349 | 426 | 1.99 | 2.61 |
| ST14000NM001G | 11,154 | 418 | 1.96 | 3.75 |

different measurements by different manufacturers. This would make it difficult to train a machine learning model and therefore, for the purposes of this paper, HDDs belonging to a single manufacturer will be used for failure prediction. As Seagate drives are the most prevalent model of HDD in this dataset, and account for the most failures, the Seagate models from Table II were selected. These are the top failing drives and include the following models: ST4000DM000, ST12000NM0007, ST8000NM0055, ST3000DM001, ST12000NM0008, ST8000DM002, and ST14000NM001G.

### B. Features

Analysis of the data quality measured the prevalence of null or missing values to determine which SMART columns in the dataset could be used as features for machine learning. The HDD model, capacity (in bytes), and the raw and normalised SMART measurements were used for training and evaluating the ML classifiers.

TABLE III. SMART ATTRIBUTE FEATURES AND THEIR SPEARMAN RANK CORRELATION WITH HDD FAILURE

| ID | Attribute Name | Null % | Correlation with Failure |
|---|---|---|---|
| 1 | Read Error Rate | 0.39 | -0.001 |
| 3 | Spin Up Time | 1.32 | - |
| 4 | Start/Stop Count | 1.32 | 0.1015 |
| 5 | Reallocated Sectors Count | 0.38 | **0.5352** |
| 7 | Seek Error Rate | 1.32 | 0.0584 |
| 9 | Power-On Hours | 0.38 | 0.0314 |
| 10 | Spin Retry Count | 1.32 | - |
| 12 | Power Cycle Count | 1.32 | 0.0959 |
| 187 | Reported Uncorrectable Errors | 1.32 | **0.6114** |
| 188 | Command Timeout | 1.32 | 0.1378 |
| 190 | Temperature Difference | 1.32 | 0.0429 |
| 192 | Power-Off Retract Count | 1.32 | 0.0455 |
| 193 | Load Cycle Count | 1.32 | 0.0448 |
| 194 | Temperature | 0.38 | 0.0429 |
| 197 | Current Pending Sector Count | 0.38 | **0.5056** |
| 198 | Uncorrectable Sector Count | 1.32 | **0.5056** |
| 199 | UltraDMA CRC Error Count | 1.32 | 0.0705 |
| 240 | Head Flying Hours | 1.32 | -0.002 |
| 241 | Total LBAs Written | 1.33 | 0.0368 |
| 242 | Total LBAs Read | 1.33 | 0.0482 |

Furthermore, the Spearman rank correlation was measured to indicate which SMART metrics were more associated with drive failure. Table III shows the list of SMART metrics used and their Spearman rank correlation coefficients.

### C. Machine Learning Implementation

The failure prediction in this paper was treated using a classification approach to determine if an HDD will fail or not at specific prediction horizons, or lookahead days. The SMART metrics were collected at 0, 1, 2, and 7 days prior to failure occurrences of the selected Seagate models described previously. As most drives in the dataset do not fail (94.5%), the datasets are highly imbalanced with vastly more examples of non-failing drives than failed ones. As imbalanced data can introduce bias to machine learning models, the non-failing class was under-sampled. A random sample of non-failing drives was collected to create balanced, unbiased training and testing datasets, where each Seagate model had equal representation of failure and non-failure. The models were trained using 80% of the balanced datasets, reserving 20% for testing on unseen data. Any categorical fields, such as the HDD model and capacity, were converted to numerical values using ordinal and one-hot encoders.

The machine learning classifiers implemented in this paper include Random Forests (RF), XGBoost (XGB), Decision Trees (DT), Neural Networks (Multi-Layer Perceptron or MLP), k-Nearest Neighbour (k-NN), and Logistic Regression (LR). These methods were selected as they have shown good performance in other existing works of research. Appropriate hyperparameters were selected for each model using Bayesian optimisation, grid search, or random search with 5-fold cross-validation to measure the combination of parameters that achieved the best mean performance.

Feature importance was assessed for each model, collected if available from the classifier, or measured using permutation. Permutation calculates the decrease in model performance as a result of randomly altering the values of each feature after the model has been trained. If the model performance is not greatly affected by permutations of a feature, then it is assumed that the model does not consider that feature important. Conversely, if the model's performance reduces then the feature is considered important, with larger performance reductions implying a relatively more important feature.

### D. Model Evaluation

The performance of each model was evaluated by generating failure predictions using the test dataset, and by comparing these predictions to the true failure status of the drives. A confusion matrix of the test predictions allowed for calculation of several evaluation metrics using the True Positive (TP), False Positive (FP), True Negative (TN) and False Negative (FN) values, as shown in Fig. 1.

The accuracy and failure detection rate, or True Positive Rate (TPR), were calculated for each classifier at each lookahead window. Accuracy measures the percentage of

| | | Predicted Label | |
|---|---|---|---|
| | | 0 (Not Failed) | 1 (Failed) |
| True Label | 0 (Not Failed) | **TN** | **FP** |
| | 1 (Failed) | **FN** | **TP** |

Figure 1. Confusion matrix used to evaluate binary classifiers.

correct predictions made by the model and the TPR measures the proportion of failed drives that were correctly predicted as failing by the classifier. The accuracy and TPR calculations are as follows:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \qquad TPR = \frac{TP}{TP + FN}$$

In addition to accuracy and TPR, the False Alarm Rate (FAR) was also measured, which calculates the percentage of drives that were incorrectly labelled as failing. Two measurements were used to determine the FAR: the False Positive Rate (FPR), which is the proportion of non-failing drives in the test dataset that were incorrectly labelled as failing; and the False Discovery Rate (FDR), which measures the proportion of predicted drives labelled as failing that are incorrect. The FAR calculations are as follows:

$$FPR = \frac{FP}{FP + TN} \qquad FDR = \frac{FP}{FP + TP}$$

Another measure of performance used to evaluate the models is the Area Under the Receiver Operating Characteristic (AUROC) curve. The ROC curve plots the TPR against the FPR and the area under the curve gives a measure of the model's prediction performance. An AUROC value of 1 represents a perfect classifier, and a value of 0.5 represents the performance obtained by a random classifier. While accuracy is commonly used and important for evaluating the likely real-world benefit of the prediction model, the AUROC represents the goodness of the model. Using the prediction probabilities of belonging to a particular class, rather than the resulting binary label of the classification, the AUROC can give a better indication of model performance and is useful for comparing the performance between different models. Consequently, the AUROC was measured for each classification model, and for each lookahead window, using 5-fold cross-validation of the test datasets to measure the mean prediction score.

## IV. RESULTS

The model prediction performance is shown in Table IV, comparing the mean AUROC scores of each of the models and the standard deviation from the 5-fold cross-validation evaluation on test data predictions. The random forest classifier achieved the highest AUROC score of 0.9185±0.0066 at a lookahead window of 0 days. This was followed very closely by the XGBoost classifier, which achieved an AUROC score of 0.9162±0.0066 for the same lookahead window. In all cases, as the prediction horizon increased, the model performance decreased. The worst performing classifier was logistic regression, achieving an

TABLE IV.  AUROC AND STANDARD DEVIATION OF FAILURE PREDICTION CLASSIFIERS USING 5-FOLD CROSS-VALIDATION

| Method | Lookahead Days (N) | | | |
|---|---|---|---|---|
| | 0 | 1 | 2 | 7 |
| Random Forest | 0.9185±0.0066 | 0.8976±0.0142 | 0.8830±0.0092 | 0.8653±0.0068 |
| XGBoost | 0.9162±0.0066 | 0.8954±0.0126 | 0.8841±0.0083 | 0.8653±0.0071 |
| Decision Tree | 0.8818±0.0086 | 0.8648±0.0084 | 0.8477±0.0132 | 0.8293±0.0053 |
| Neural Network | 0.8721±0.0105 | 0.8526±0.0132 | 0.8517±0.0131 | 0.8254±0.0133 |
| k-NN | 0.8617±0.0121 | 0.8414±0.0111 | 0.8482±0.0150 | 0.8176±0.0088 |
| Logistic Regression | 0.8484±0.0117 | 0.8166±0.0135 | 0.8192±0.0117 | 0.7871±0.0099 |

AUROC score of 0.8484±0.0117 at the shortest prediction horizon. The AUROC scores are generally higher in this paper compared to [17], but the prediction performance ranking of classification methods agrees with their findings.

The accuracy of the models generally follows the same trend and rankings as the AUROC scores, as shown in Table V. As accuracy uses the predicted label of the HDDs, and not the prediction probabilities associated with each class, they are lower than the AUROC as expected. Again, random forest and XGBoost performed the best at a lookahead window of 0 days with accuracies of 0.862 and 0.864 respectively.

Although the AUROC and accuracy scores are important evaluators of prediction performance, it is likely that a real-world application would place importance on how well the classifiers predicted failing drives. Therefore, the failure detection rate, also known as the True Positive Rate (TPR), was calculated for each model, and at each prediction horizon, as shown in Table VI. The XGBoost classifier achieved the best failure detection rate at most prediction horizons, successfully identifying 77.6% and 74.8% of failing drives with lookahead windows of 0 and 1 day respectively. It was able to successfully predict the imminent failure of 70.7% of drives 7 days in advance, better than any

TABLE V.  ACCURACY OF FAILURE PREDICTION CLASSIFIERS ON TESTING DATASETS AT EACH LOOKAHEAD WINDOW (N)

| | N | RF | XGB | DT | MLP | k-NN | LR |
|---|---|---|---|---|---|---|---|
| Accuracy | 0 | 0.862 | 0.864 | 0.854 | 0.810 | 0.801 | 0.778 |
| | 1 | 0.841 | 0.844 | 0.832 | 0.793 | 0.788 | 0.753 |
| | 2 | 0.822 | 0.822 | 0.813 | 0.787 | 0.786 | 0.752 |
| | 7 | 0.800 | 0.804 | 0.790 | 0.753 | 0.753 | 0.720 |

of the other classifiers. Additionally, any model used in a real-world scenario would require false alarms to be minimised to gain the confidence of users. Also shown in Table VI, the FPR and FDR for each model at each prediction horizon was measured. Random forest and XGBoost show the best ratios between TPR and FAR. At the shortest prediction horizon, the random forest classifier achieves 76.7% TPR with 4.1% FPR, while XGBoost achieves 77.6% TPR with 4.5% FPR.

The feature importance ranking for each classifier is shown in Table VII, obtained from the model where available, or estimated with feature permutation. Features with smaller values of ranking order indicate more importance to the classification model. Missing values indicate that the feature was not present in the top 5 most important features for that classifier. SMART 187, reported uncorrectable errors, is a key indicator for HDD failure and is the most important feature for almost all classifiers.

TABLE VII.  MOST COMMON IMPORTANT FEATURES AMONGST CLASSIFIERS INDICATED BY FEATURE RANKING ORDER

| | Feature Ranking Order of Importance if Present in Top 5 Most Important Features | | | | | |
|---|---|---|---|---|---|---|
| | RF | DT | XGB | MLP | k-NN | LR |
| SMART 5 | 2 | 3 | 4 | 3 | 1 | 2 |
| SMART 187 | 1 | 1 | 1 | 1 | 4 | 1 |
| SMART 197 | 3 | 2 | 3 | 4 | 2 | - |
| SMART 198 | 4 | - | 2 | 2 | 3 | - |
| SMART 240 | - | 5 | - | - | - | - |
| SMART 241 | 5 | - | - | - | - | - |
| SMART 242 | - | 4 | - | 5 | - | - |

TABLE VI.  TRUE POSITIVE RATE (TRP), FALSE POSITIVE RATE (FPR) AND FALSE DISCOVERY RATE (FDR) OF CLASSIFIERS AT EACH LOOKAHEAD WINDOW (N)

| | N | RF | XGB | DT | MLP | k-NN | LR |
|---|---|---|---|---|---|---|---|
| TPR | 0 | 0.767 | 0.776 | 0.759 | 0.761 | 0.682 | 0.599 |
| | 1 | 0.738 | 0.748 | 0.746 | 0.728 | 0.669 | 0.566 |
| | 2 | 0.707 | 0.717 | 0.695 | 0.726 | 0.681 | 0.582 |
| | 7 | 0.689 | 0.707 | 0.701 | 0.689 | 0.636 | 0.507 |
| FPR | 0 | 0.041 | 0.045 | 0.049 | 0.139 | 0.077 | 0.040 |
| | 1 | 0.052 | 0.058 | 0.078 | 0.141 | 0.089 | 0.056 |
| | 2 | 0.061 | 0.072 | 0.066 | 0.150 | 0.106 | 0.074 |
| | 7 | 0.085 | 0.095 | 0.118 | 0.182 | 0.130 | 0.065 |
| FDR | 0 | 0.049 | 0.054 | 0.060 | 0.153 | 0.100 | 0.062 |
| | 1 | 0.064 | 0.070 | 0.092 | 0.160 | 0.116 | 0.088 |
| | 2 | 0.078 | 0.090 | 0.086 | 0.168 | 0.132 | 0.111 |
| | 7 | 0.107 | 0.116 | 0.140 | 0.208 | 0.169 | 0.114 |

Comparing these results with Table III, the four features with the highest correlation coefficient are the top four most important features in Table VII.

## V. DISCUSSION

The classification approach adopted in this paper achieved relatively high levels of prediction performance. It highlighted the most accurate machine learning classifiers for failure prediction using a common dataset of HDDs from an operational data centre. It proved that SMART metrics can be used as an indication of imminent failure, with some more useful than others. However, there are some limitations and weaknesses that need to be highlighted.

As HDDs are constructed of both mechanical and electrical components, there are a number of potential reasons for their failure. Over time the subcomponents within the hard drive unit can degrade, causing problems with reading and writing data, and eventually leading to failure. In this scenario, where the drives fail slowly with accumulated usage and workload, SMART metrics may be a good indicator of HDD health. However, hard disk drives are susceptible to external factors, such as physical disturbances from knocks and vibrations, water damage, and power-related problems, including voltage spikes. It is therefore unlikely that SMART metrics would indicate imminent failure for quick-failing drives due to external factors. With the best TPR of 77.6% achieved in this paper, at least 22.4% of failed drives were not detected. Without knowing the root cause of failure, it could be possible that those drives did not contain any indication of failure in their SMART metrics.

The decision to approach the research as a binary classification problem means that the prediction only has the option to label a drive with a failure status of failing or non-failing. However, the health of HDDs may indicate that the drive is at low, moderate, or high risk of failing, in which case a multi-class classification approach might be better suited. Using the same logic, it may be desirable to predict the probability of drive failure. In the case of the binary classification approach, if a drive is predicted to have a 51% chance of failing it would be labelled as a failing drive (assuming a 50% threshold). However, data centre maintainers may dismiss that drive as a low risk if they were presented with the probability of failing, whereas they would be forced to investigate the drive if the binary classification prediction was presenting its impending failure. Hence, the AUROC was used in this project to better evaluate the classifiers' prediction performance using the prediction probabilities of belonging to a particular class.

A well-known issue with SMART metrics is that the data they report isn't always consistent between various manufacturers and drive models. Some attribute fields may be used to record a particular measurement by one manufacturer, but a completely different measurement by another. And the format or scale used may not be consistent even when reporting the same measurement. Therefore, to mitigate against this, the drive models used in the ML implementation of this work only considered drives of a single manufacturer, Seagate. Consequently, it is not guaranteed that the prediction classifiers would generalise well for predicting failures or other manufacturers' drives.

The prediction horizons examined in this paper use the SMART attribute measurements from 0, 1, 2, and 7 days prior to HDD failure. The prediction performance improved as the lookahead days decreased, with 0 days achieving the best AUROC, accuracy, TPR, and FAR rates. The 0-day lookahead window means that the drives failed sometime in the next 0-24 hours. In a real-world application it is likely that this prediction window is too short, not allowing for enough time to act. Increasing the window would decrease the prediction accuracy, potentially reducing the likelihood for users to trust the classification output.

The features used to train and evaluate the ML classifiers consisted only of the raw and normalised SMART measurements with scaling applied. As indicated in [9], the rate of change of SMART measurements over a given period can provide additional features that potentially improve the prediction performance of HDD failure prediction models. The work presented in this paper would benefit by extending to include temporal disparities of SMART measurements as features for machine learning.

Further work may include extending the scope to compare the prediction performance of the classification methods implemented here with other methodologies, such as time series prediction, and analysing the AUROC, accuracy, TPR, and FAR at varying prediction horizons.

## VI. CONCLUSSION

In this paper, the SMART attributes of operational HDDs in a large data centre were analysed with respect to drive failure. SMART attributes 5, 187, 197, and 198 (reallocated sectors count, reported uncorrectable errors, current pending sector count, and uncorrectable sector count) were observed to have the highest correlation with HDD failure.

A subset of the SMART attribute measurements, reported daily by the data centre HDDs, was used to create machine learning classifiers for drive failure prediction. The ML classification models implemented in this work include Random Forest, XGBoost, Decision Tree, Neural Network (Multi-Layer Perceptron), k-Nearest Neighbour, and Logistic Regression methods. The SMART metrics were collected at 0, 1, 2, and 7 days prior to drive failure to evaluate the prediction performance at multiple prediction horizons. It was found that as the prediction horizon decreases, the performance of the failure prediction increased for all classifiers.

Random Forest and XGBoost classifiers achieved the best results, with 86% prediction accuracy and 4-5% False Alarm Rate (FAR) at the shortest prediction horizon. The failure detection rate ranged from 67% when making predictions 7 days prior to HDD failure, to 77% when using the SMART measurements recorded in the last 24 hours before failure. The AUROC was calculated to make better comparisons between the classifiers, which again showed Random Forest and XGBoost as the best performing, with AUROC scores of 0.9185±0.0066 and 0.9162±0.0066 respectively at the shortest prediction horizon.

The relative feature importances of the ML models were obtained, either directly from the classifier or estimated using feature permutation. It was found that, in all but one of the classifiers, SMART 187 was regarded as the most important metric for predicting HDD failure. The top four most important features across all classifiers were those with the highest Spearman rank correlation coefficient relating to failure as described above (SMART 5, 187, 197, and 198).

The classification models generated in this work could benefit from future advancements, and enhanced feature engineering would likely improve the performance of the models. Additionally, SMART attributes do not account for many of the external factors that can affect the health of HDDs, such as physical disturbances. Therefore, it may be valuable to extend the work of this paper by considering other relevant datasets alongside SMART data for predicting failure. For example, force sensor data may indicate knocks or jolts to the HDD, and [18] has shown success in using machine learning to classify force signals and determine if a collision occurred. Other future work may include incorporating the classification models with autonomic computing, where the failure predictions can inform autonomic actions. Such actions may involve pre-emptively backing up data to another storage device to mitigate the risk of data loss.

## REFERENCES

[1] A. Lawrence and L. Simon, "Annual Outage Analysis 2023," Uptime Institute, New York, NY 10174, 2023. [Online]. Available: https://uptimeinstitute.com/resources/research-and-reports/annual-outage-analysis-2023. [retrieved: Feb. 2025].

[2] J. Davis, D. Bizo, A. Lawrence, O. Rogers, and M. Smolaks, "Global Data Center Survey 2022," Uptime Institute, New York, NY 10174, 2022. [Online]. Available: https://uptimeinstitute.com/resources/research-and-reports/uptime-institute-global-data-center-survey-results-2022. [retrieved: Feb. 2025].

[3] K. V. Vishwanath and N. Nagappan, "Characterizing Cloud Computing Hardware Reliability," in Proceedings of the 1st ACM Symposium on Cloud Computing, Indiana, IN, USA, 2010, pp. 193–204.

[4] G. Wang, L. Zhang, and W. Xu, "What Can We Learn from Four Years of Data Center Hardware Failures?," in 2017 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), Denver, CO, USA, 2017, pp. 25-36.

[5] L. Perri, "What's New in the 2022 Gartner Hype Cycle for Emerging Technologies," Gartner, 2022. [Online]. Available: https://www.gartner.com/en/articles/what-s-new-in-the-2022-gartner-hype-cycle-for-emerging-technologies. [retrieved: Feb. 2025].

[6] J. O. Kephart and D. M. Chess, "The vision of autonomic computing", Computer, vol. 36, no. 1, pp. 41-50, Jan. 2003, doi: 10.1109/MC.2003.1160055.

[7] R. Sterritt, (2005), "Autonomic computing". Innovations Syst Softw Eng: A NASA Journal, 1, 79–88 (2005). doi:10.1007/s11334-005-0001-5

[8] J. Li et al., "Hard Drive Failure Prediction Using Classification and Regression Trees," in 2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks, Atlanta, GA, USA, 2014, pp. 383-394.

[9] Z. Miller, O. Medaiyese, M. Ravi, A. Beatty, and F. Lin, "Hard Disk Drive Failure Analysis and Prediction: An Industry View," in 2023 53rd Annual IEEE/IFIP International Conference on Dependable Systems and Networks - Supplemental Volume (DSN-S), Porto, Portugal, 2023, pp. 21-27.

[10] I. C. Chaves, M. R. P. de Paula, L. G. M. Leite, J. P. P. Gomes, and J. C. Machado, "Hard Disk Drive Failure Prediction Method Based On A Bayesian Network," in 2018 International Joint Conference on Neural Networks (IJCNN), Rio de Janeiro, Brazil, 2018, pp. 1-7.

[11] I. C. Chaves, M. R. P. de Paula, L. G. M. Leite, L. P. Queiroz, J. P. P. Gomes, and J. C. Machado, "BaNHFaP: A Bayesian Network Based Failure Prediction Approach for Hard Disk Drives," 2016 5th Brazilian Conference on Intelligent Systems (BRACIS), Recife, Brazil, 2016, pp. 427-432.

[12] C. Xu, G. Wang, X. Liu, D. Guo, and T. -Y. Liu, "Health Status Assessment and Failure Prediction for Hard Drives with Recurrent Neural Networks," in IEEE Transactions on Computers, vol. 65, issue. 11, pp. 3502-3508, 1 Nov. 2016.

[13] S. Pang, Y. Jia, R. Stones, G. Wang, and X. Liu, "A combined Bayesian network method for predicting drive failure times from SMART attributes," in 2016 International Joint Conference on Neural Networks (IJCNN), Vancouver, BC, Canada, 2016, pp. 4850-4856.

[14] A. Coursey, G. Nath, S. Prabhu, and S. Sengupta, "Remaining Useful Life Estimation of Hard Disk Drives using Bidirectional LSTM Networks," in 2021 IEEE International Conference on Big Data (Big Data), Orlando, FL, USA, 2021, pp. 4832-4841

[15] Q. Hai, S. Zhang, C. Liu, and G. Han, "Hard Disk Drive Failure Prediction Based on GRU Neural Network," in 2022 IEEE/CIC International Conference on Communications in China (ICCC), Sanshui, Foshan, China, 2022, pp. 696-701.

[16] Backblaze. Hard Drive Data and Stats [Online]. Available: https://www.backblaze.com/cloud-storage/resources/hard-drive-test-data. [retrieved: Jan. 2024].

[17] R. Pinciroli, L. Yang, J. Alter, and E. Smirni, "Lifespan and Failures of SSDs and HDDs: Similarities, Differences, and Prediction Models," in IEEE Transactions on Dependable and Secure Computing, vol. 20, issue. 1, pp. 256-272, 2023.

[18] A.-N. Sharkawy, A. Ma'arif, Furizal, R. Sekhar, and P. Shah, "A comprehensive pattern recognition neural network for collision classification using force sensor signals," Robotics, vol. 12, issue. 5, p.124, 2023.