

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

Alistair McLean, Roy Sterritt



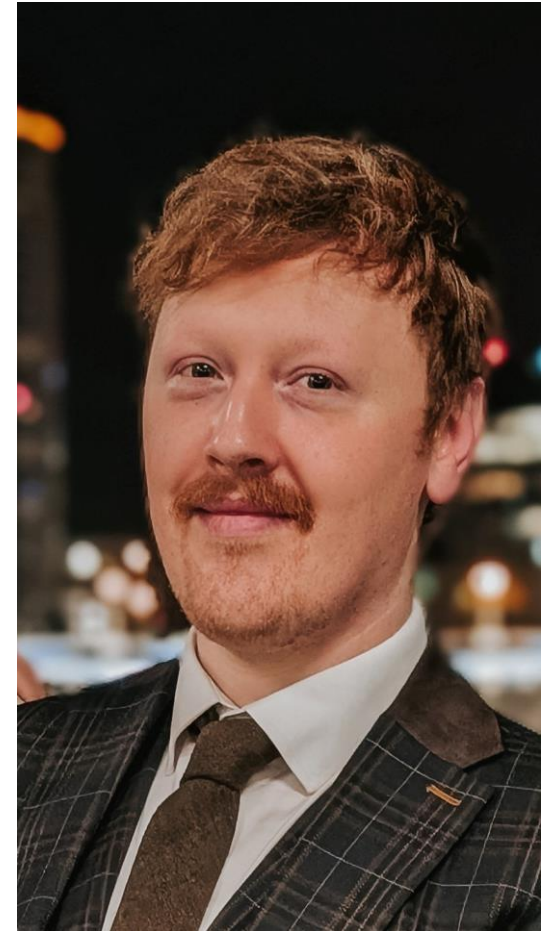
Presented by **Alistair McLean**
School of Computing
Faculty of Computing, Engineering
and the Built Environment
Ulster University
McLean-A13@ulster.ac.uk

IARIA ICAS 2025

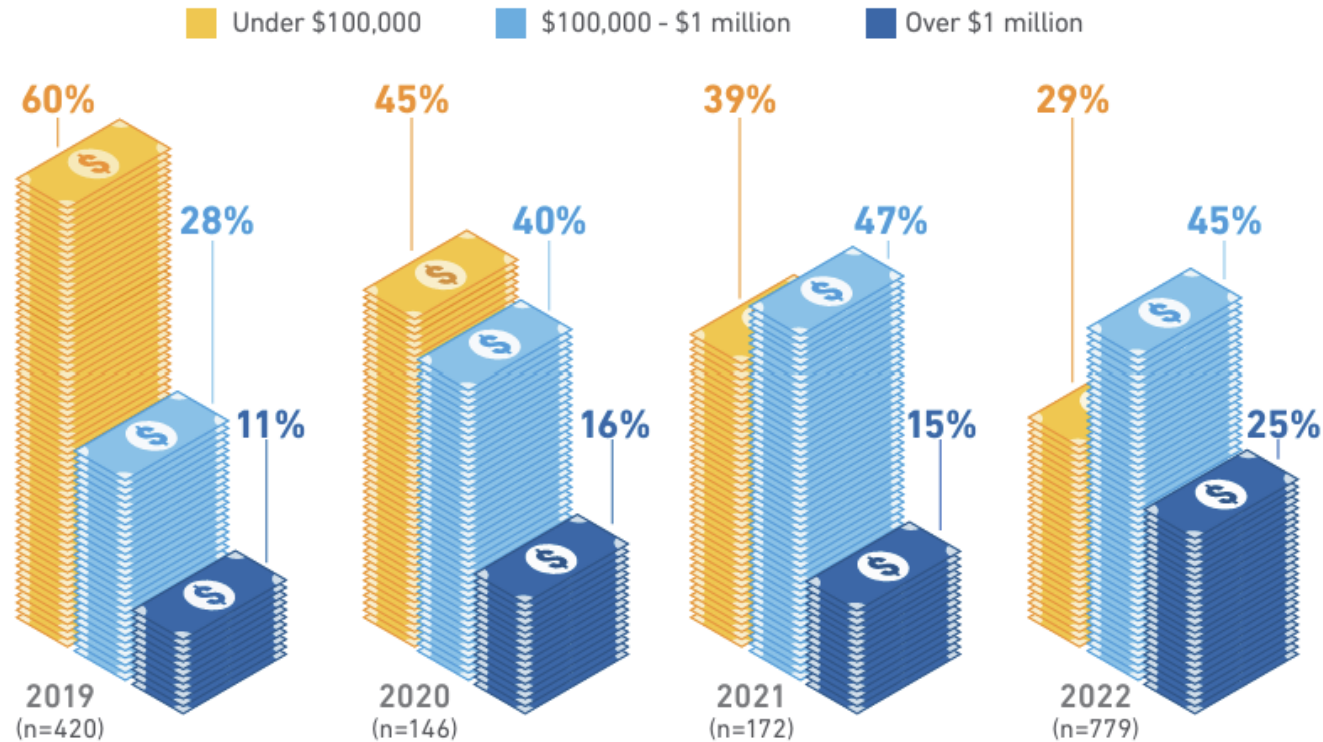
AIAC: AI and Autonomic Computing special session

Alistair McLean

- MSc Artificial Intelligence Graduate from Ulster University (2024).
- Autonomic Computing interest from MSc AI module: COM760 Autonomic Computing and Robotics
see Alistair McLean, [Roy Sterritt](#), “Autonomic Computing in the Cloud: An Overview of Past, Present and Future Trends”, *The 2023 IARIA Annual Congress on Frontiers in Science, Technology, Services, and Applications: Technical Advances and Human Consequences*, Valencia, Spain, pp77-82. Available at [ThinkMind\(TM\) Digital Library](#)
- Currently working as an AI Engineer in the telecommunications industry.
- <https://www.linkedin.com/in/alistair-mclean/>



Problem Description



Uptime Institute Global Survey of IT and Data Centre Managers 2019-2022 [1].

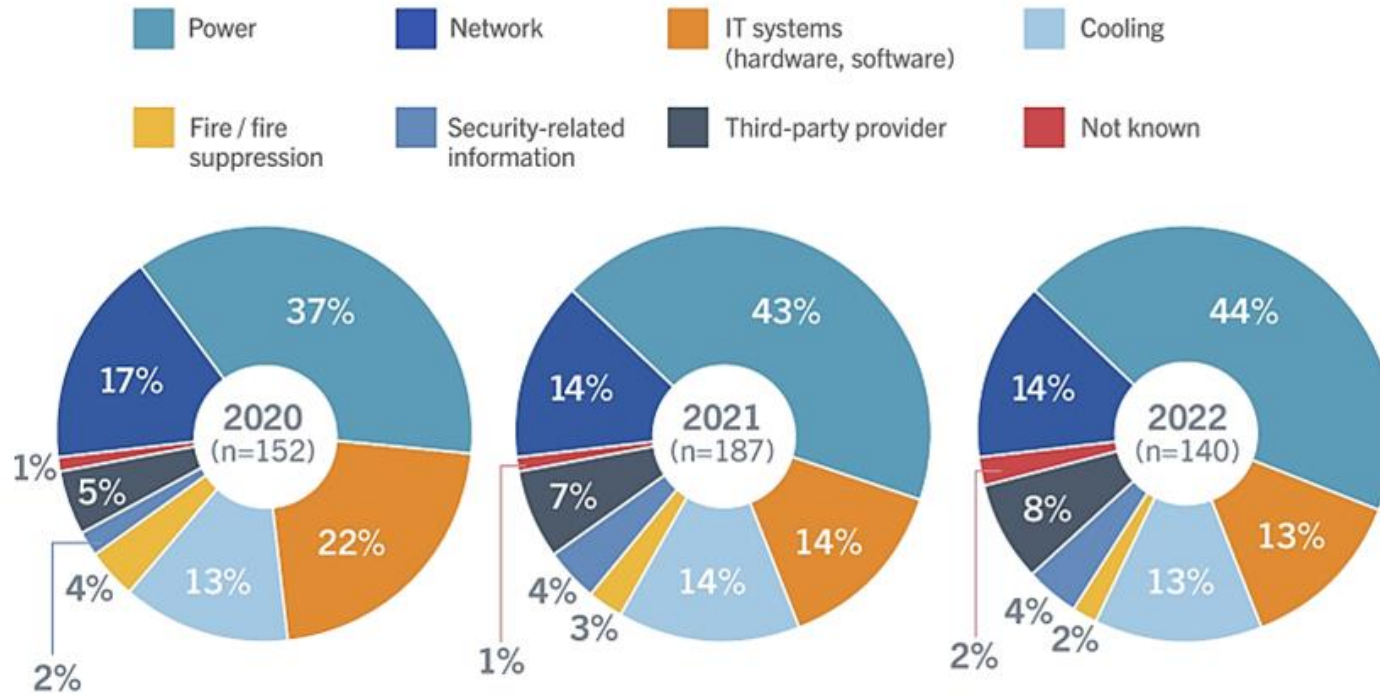
Unplanned downtime and IT system outages can cost organisations millions of dollars in lost revenue, loss of opportunity, and reduced reputation.

Uptime Institute reports:

- This **cost is increasing** year on year
- In 2022, **70%** of downtime incidents **cost more than \$100,000**
- **25%** of incidents **cost more than \$1 million**

[1] A. Lawrence and L. Simon, "Annual Outage Analysis 2023," Uptime Institute, New York, NY 10174, 2023. Accessed: Jul. 23, 2024. [Online]. Available: <https://uptimeinstitute.com/resources/research-and-reports/annual-outage-analysis-2023>

Problem Description



Primary cause of significant site outages - Uptime Institute Global Survey of IT and Data Centre Managers 2020-2022 [1].

Top Causes of Data Centre Outages:

- Power
- Network
- **IT Systems**
- Cooling

[1] J. Davis, D. Bizo, A. Lawrence, O. Rogers, and M. Smolaks, "Global Data Center Survey 2022," Uptime Institute, New York, NY 10174, 2022. Accessed: Jul. 23, 2024. [Online]. Available: <https://uptimeinstitute.com/resources/research-and-reports/uptime-institute-global-data-center-survey-results-2022>

Problem Description

Causes of IT system failure in cloud computing infrastructures:

[1]: Study on data centres containing more than 100,000 servers

- **Hard Disk Drives (HDDs) are the most replaced components and one of the least reliable**
- **78% of faults or replacements were attributed to hard disks**
- **8% of servers can expect 1 hardware incident in a given year – but value is higher for machines with many HDDs**

[2]: Study on data centres containing hundreds of thousands of servers over 4 years

- **82% of component failures were related to HDDs**

[1] 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.

[2] 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.

Problem Statement

“Data centres could improve their reliability with effective monitoring of the health of HDDs...

...and the ability to predict imminent HDD failure would facilitate preventative action to mitigate against outages”

HDD Monitoring

Self-Monitoring, Analysis, and Reporting Technology (S.M.A.R.T)

- Collects measurements from sensors within the HDD unit to report on various indicators of health and reliability
- Used by many HDD manufacturers today
- Numbered 1-255
- Raw value
- Normalised value

Attribute	Definition
SMART 1	Read Error Rate
SMART 3	Spin Up Time
SMART 4	Start/Stop Count
SMART 5	Reallocated Sectors Count
SMART 7	Seek Error Rate
SMART 9	Power-On Hours
SMART 10	Spin Retry Count
SMART 12	Power Cycle Count
SMART 184	End-to-End error / IOEDC
SMART 187	Reported Uncorrectable Errors
SMART 188	Command Timeout
SMART 189	High Fly Writes
SMART 190	Temperature Difference
SMART 191	G-sense Error Rate
SMART 192	Power-off Retract Count
SMART 193	Load Cycle Count
SMART 194	Temperature
SMART 197	Current Pending Sector Count
SMART 198	Uncorrectable Sector Count
SMART 199	UltraDMA CRC Error Count
SMART 240	Head Flying Hours
SMART 241	Total LBAs Written
SMART 242	Total LBAs Read

Existing Work

Paper	Methodology	Research Output
[1]	Classification and Regression Trees	Classifier successfully predicted 95% of failures with False Alarm Rate of less than 0.1%
[2]	XGBoost Classification	XGBoost achieved low precision but they found that using the difference in SMART measurements over time as features led to better results
[3]	Bayesian Network	Predicts Remaining Useful Life (RUL) estimates. Showed good performance, achieving similar relative accuracy to a Recurrent Neural Network in other existing work
[4]	Combined Bayesian Network	Ensemble learning approach using learning results from 4 individual classifiers to create a combined Bayesian network. Performed similarly to a classification tree model but had addition benefit of providing estimated time before failure.
[5]	Bidirectional LSTM	Achieved 96.4% accuracy in predicting HDD failure for a 15-day lookback period
[6]	GRU Neural Network	Achieved 95% failure detection rate and 0.2% false alarm rate

[1] 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.

[2] 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.

[3] 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.

[4] 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.

[5] 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.

[6] 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.

Research Contribution

Comparative study to determine the best performing machine learning methods for HDD failure prediction

- Use a common dataset of operational data centre HDDs
- Investigate the relationships between SMART metrics and HDD failure
 - Highlight the most important SMART metrics that indicate drive health status

Dataset

Source:  **Backblaze** [1]

- Cloud storage provider
- Currently has data centers in:
 - Sacramento, California
 - Stockton, California
 - Phoenix, Arizona
 - Reston, Virginia
 - Amsterdam, Netherlands
- Report daily SMART metrics collected from HDDs in their data centres

Dataset

- **Date** (yyyy-mm-dd)
- **Serial Number** – The manufacturer-assigned serial number of the drive.
- **Model** – The manufacturer-assigned model number of the drive.
- **Capacity** – The drive capacity in bytes.
- **Failure** – Contains a “0” if the drive is OK. Contains a “1” if this is the last day the drive was operational before failing.
- **SMART measurements** (raw and normalised values as reported by the given drive)
 - 2013-2014 SMART Stats – 80 columns of data, that are the Raw and Normalized values for **40 different SMART stats** as reported by the given drive.
 - 2015-2017 SMART Stats – 90 columns of data, that are the Raw and Normalized values for **45 different SMART stats** as reported by the given drive.
 - 2018 (Q1) SMART Stats – 100 columns of data, that are the Raw and Normalized values for **50 different SMART stats** as reported by the given drive
 - 2018 (Q2) SMART Stats – 104 columns of data, that are the Raw and Normalized values for **52 different SMART stats** as reported by the given drive.
 - 2018 (Q4) SMART Stats – 124 columns of data, that are the Raw and Normalized values for **62 different SMART stats** as reported by the given drive.

Example:

date	serial_number	model	capacity_bytes	failure	smart_1_normalized	smart_1_raw	smart_n...
26/12/2023	ZA180RV5	ST8000NM0055	8.00156E+12	1	81	116937464	...
23/12/2023	ZLW17S54	ST14000NM001G	1.40005E+13	1	82	163135856	...
23/12/2023	ZHZ4WMSE	ST12000NM0008	1.20001E+13	1	79	83946168	...

Data Exploration

For 10-year period from 01/01/2014 to 31/12/2023:

- Over 454 million rows
- 388,485 HDDs (unique serial numbers)
- 21,356 HDD failures

- **Seagate is most prevalent manufacturer of HDDs in the dataset (49.78%)**

- **Seagate HDDs account for 75.75% of failures in the dataset**

- **8.37% of all Seagate HDDs experienced failure**

Manufacturer	Total HDDs	% of HDDs in Dataset
Seagate	193,378	49.78
Toshiba	86,785	22.34
HGST	53,913	13.88
WDC	39,741	10.23
Hitachi	13,138	3.38
Other	1,531	0.39

Manufacturer	Total HDDs	Total Failures	% of All Failures	Failure Rate (%)
Seagate	193,378	16,177	75.75	8.37
Toshiba	86,785	2,223	10.41	2.56
HGST	53,913	1,651	7.73	3.06
WDC	39,741	682	3.19	1.72
Hitachi	13,138	467	2.19	3.55
Other	1,531	156	0.73	10.12

Data Exploration

Inconsistent use of SMART metrics between manufacturers:

Seagate accounts for 75.75% of all HDD failures

Selection of top failing Seagate models:

- ST4000DM000
- ST12000NM0007
- ST8000NM0055
- ST3000DM001
- ST12000NM0008
- ST8000DM002
- ST14000NM001G

Manufacturer	Total HDDs	Total Failures	% of All Failures	Failure Rate (%)
Seagate	193,378	16,177	75.75	8.37
Toshiba	86,785	2,223	10.41	2.56
HGST	53,913	1,651	7.73	3.06
WDC	39,741	682	3.19	1.72
Hitachi	13,138	467	2.19	3.55
Other	1,531	156	0.73	10.12

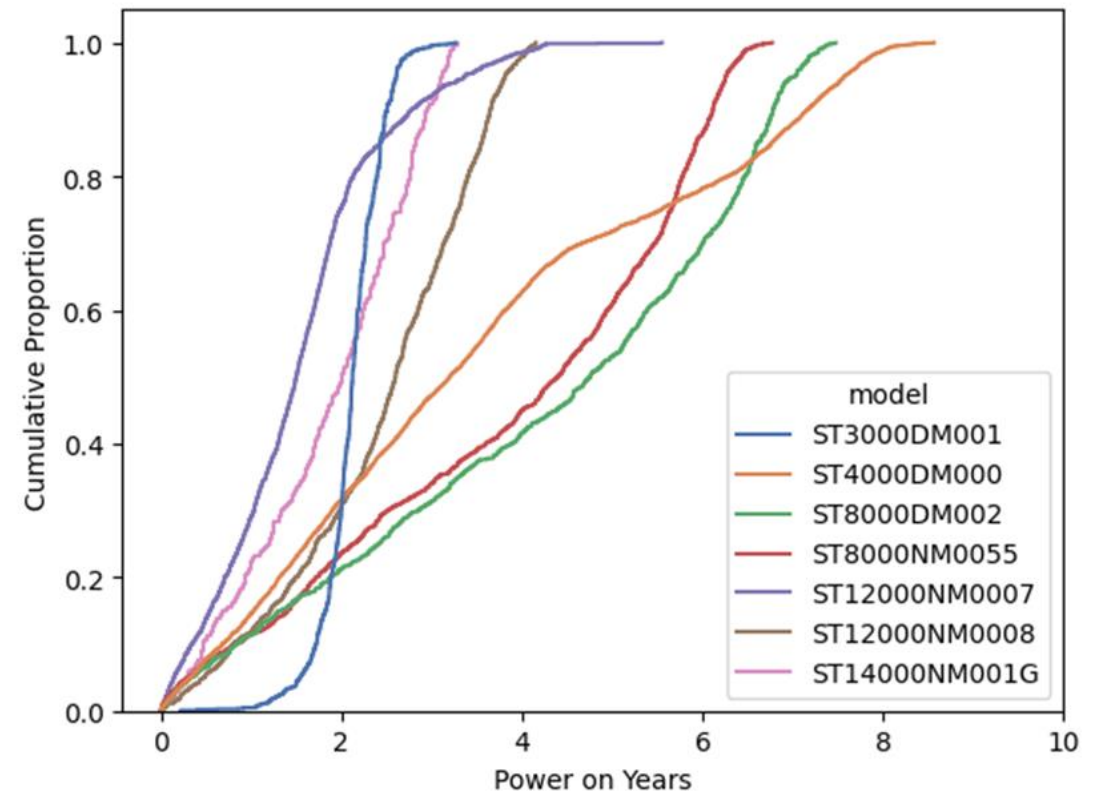
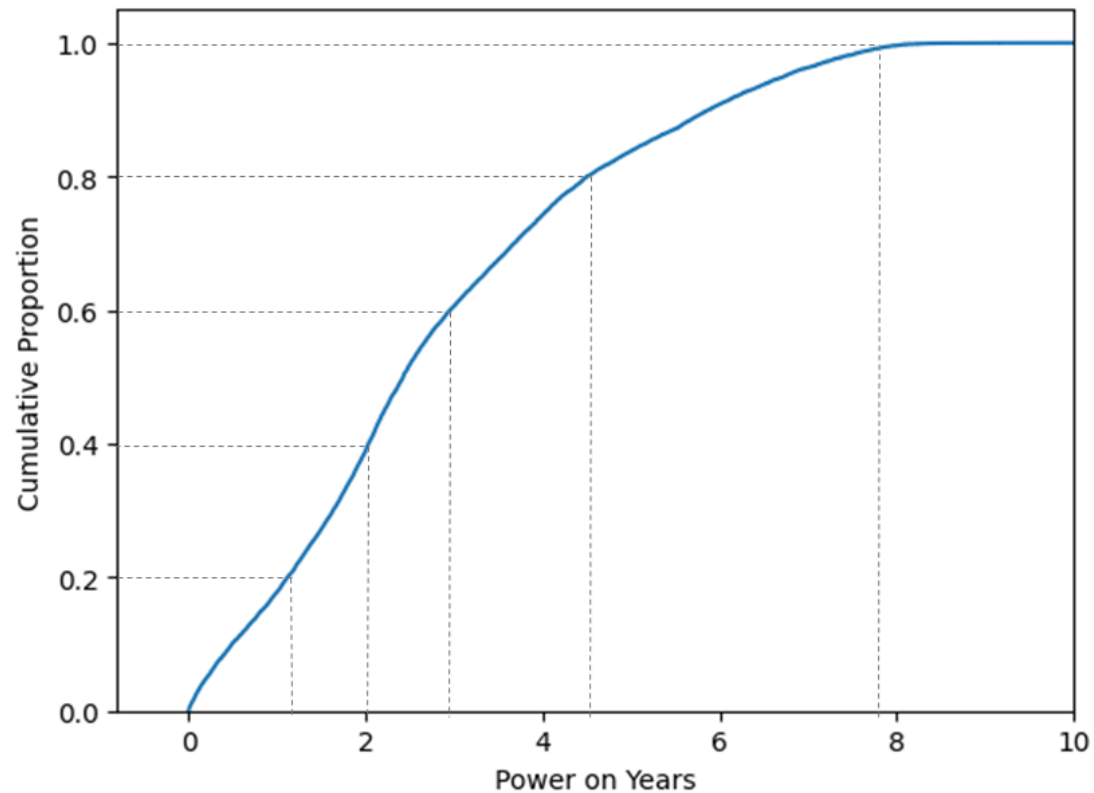
Model	Total HDDs	Total Failures	% of All Failures	Model Failure %
ST4000DM000	36,983	5,602	26.23	15.15
ST12000NM0007	38,838	2,106	9.86	5.42
ST8000NM0055	15,680	1,718	8.04	10.96
ST3000DM001	4,354	1,454	6.81	33.39
ST12000NM0008	20,836	1,349	6.32	6.47
TOSHIBA MG07ACA14TA	39,292	1,173	5.49	2.99
ST8000DM002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

Data Exploration

Time to Failure Analysis

- SMART 9 (Power-On Hours)

100% of failures occur within 8 years of operation; 80% within 4.5 years; 60% within 3 years.



Features

SMART metrics with low percentage of nulls or missing values were selected as features for machine learning.

Other features include:

- Model
- Capacity (bytes)

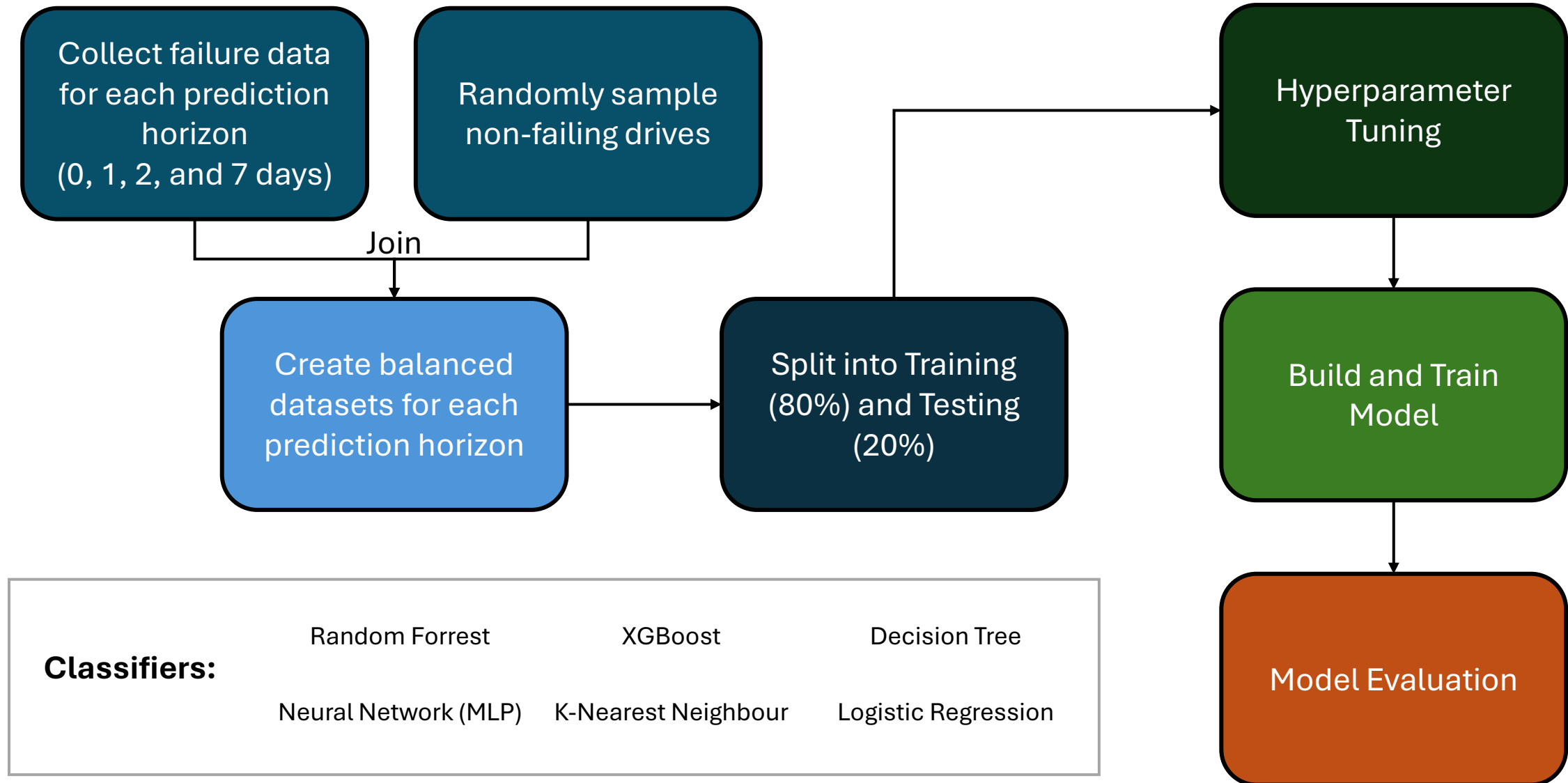
Spearman rank correlation was calculated to measure the association with HDD failure status.

Top four SMART attributes correlated with HDD failure:

- SMART 5
- SMART 187
- SMART 197
- SMART 198

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

Machine Learning Implementation



Machine Learning Evaluation

Confusion Matrix:

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

TP = True Positives, TN = True Negatives, FP = False Positives, FN = False Negatives

Accuracy:

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

True Positive Rate (TPR):

$$TPR = \frac{TP}{TP + FN}$$

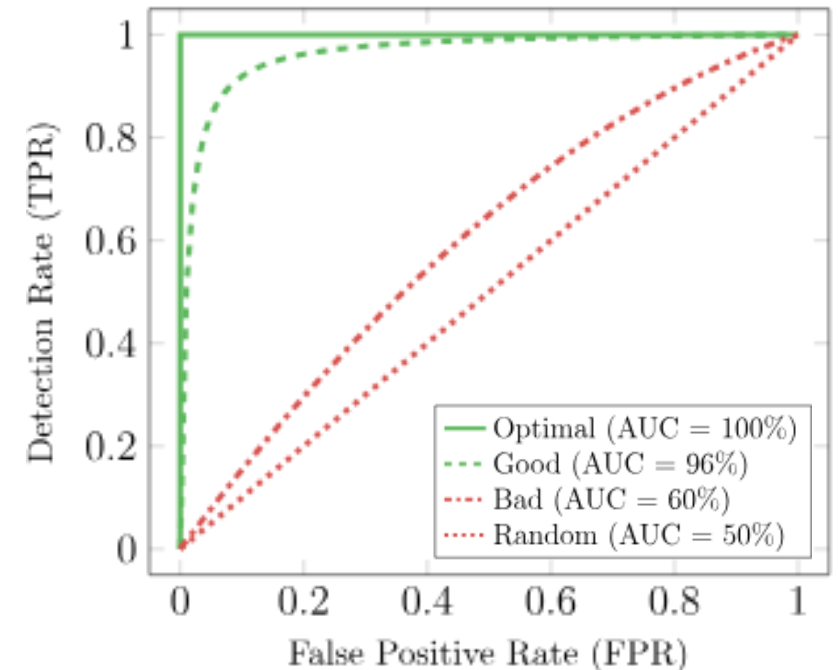
False Positive Rate (FPR):

$$FPR = \frac{FP}{FP + TN}$$

False Discovery Rate (FDR):

$$FDR = \frac{FP}{FP + TP}$$

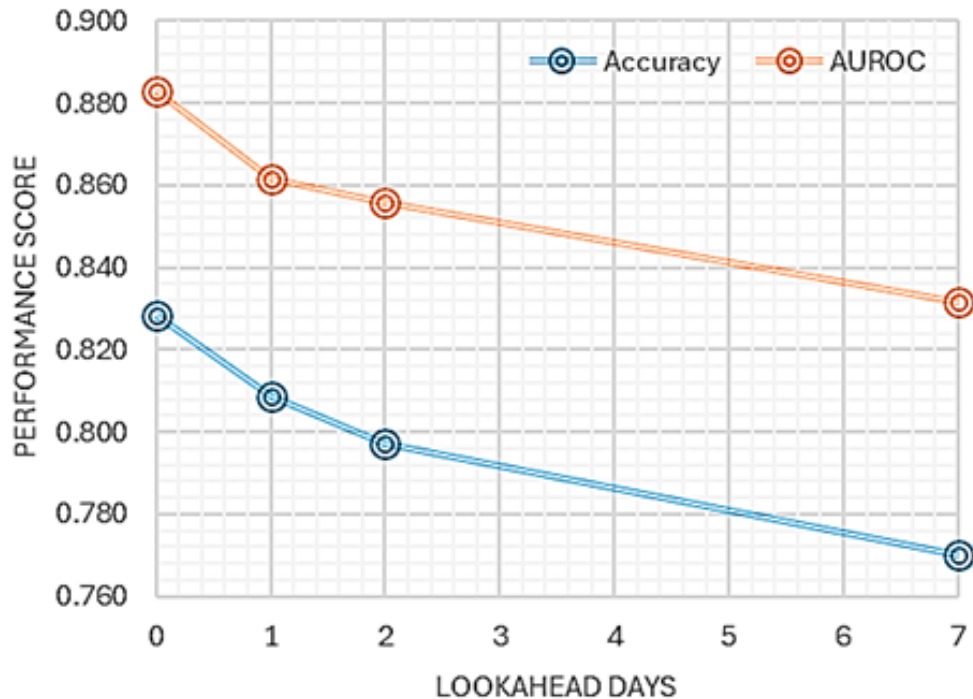
AUROC [1]:



[1] SecuML. *Detection Performance Metrics* [Online]. Available: https://anssi-fr.github.io/SecuML/miscellaneous.detection_perf.html (accessed 2024, Sept. 10).

Machine Learning Results

- As prediction horizon increased, accuracy and AUROC decreased
- Random Forest and XGBoost achieved the best results
- Logistic Regression performed the worst



Accuracy:

	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

AUROC:

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

Machine Learning Results

- **Best TPR was achieved by XGBoost (77.6%) with FPR of 4.5% at N = 0.**
- **Random Forest performed similarly with TPR of 76.7% and FPR of 4.1% at N = 0.**
- **XGBoost also achieved best result at N = 7, with TPR of 70.7% and FPR of 9.5%.**

	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

Machine Learning Results

- **Top 4 most important features for each classifier:**
 - SMART 187
 - SMART 5
 - SMART 197
 - SMART 198
- **These SMART attributes also have the highest Spearman rank correlation with HDD failure**

	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	-	-

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

Critical Evaluation

- **Limitations of SMART**
 - HDDs are susceptible to external factors that can cause failure
 - No root cause details
 - TPR of 77.6% means 22.4% of failing drives weren't detected
- **Binary Classification**
 - Is multi-class more representative?
 - Probability of failure may be more relevant
- **HDD Manufacturers and Models**
 - Not guaranteed to use SMART attributes consistently
 - ML implementation only applied to a single manufacturer
- **Prediction Horizon**
 - Only used set prediction horizons of 0, 1, 2, and 7 days
 - 0 days (0-24hrs) achieved best results – but is this too short for real-world application?
 - Longer prediction horizon will reduce prediction performance
- **Feature Engineering**
 - Limited feature engineering in this work
 - Temporal disparities in SMART measurements could be used as features to potentially achieve better prediction performance

Summary

- The following SMART attributes were identified as the most important indicators of imminent HDD failure:
 - SMART 5 (Reallocated Sectors Count)
 - SMART 187 (Reported Uncorrectable Errors)
 - SMART 197 (Current Pending Sector Count)
 - SMART 198 (Uncorrectable Sector Count)
- Prediction performance improves as the prediction horizon decreases
- Random Forest and XGBoost classifiers achieve the best results with:
 - 86% accuracy, 77% Failure Detection Rate, and 4.5% False Alarm Rate at shortest prediction horizon (0-24 hours prior to failure)
 - Random Forest: AUROC of 0.9185 ± 0.0066 at $N = 0$
 - XGBoost: AUROC of 0.9162 ± 0.0066 at $N = 0$



Faculty of Computing,
Engineering and the
Built Environment

MSc Artificial Intelligence

ONLY MSC ARTIFICIAL
INTELLIGENCE COURSE
IN NORTHERN IRELAND

ulster.ac.uk

Acknowledgements

Alistair McLean is sponsored by his employer to undertake the MSc in Artificial Intelligence at Ulster University in parttime mode.

This paper was produced during Ulster University's MSc in Artificial Intelligence, as part of the research project/dissertation (COM 748).

Conference funding provided by the School of Computing as acknowledgement of Alistair's project/dissertation being one of the top marks in 2024.

Thank You



Presented by **Alistair McLean**
School of Computing
Faculty of Computing, Engineering
and the Built Environment
Ulster University
McLean-A13@ulster.ac.uk
IARIA ICAS 2025

AIAC: AI and Autonomic Computing special session