

Health Levels Modeling for SSD Failure Prediction

Gustavo W. M. Valença¹, Francisco L. F. Pereira¹,
Felipe T. Brito¹, Victor A. E. de Farias¹, Javam C. Machado¹

¹Laboratório de Sistemas e Bancos de Dados (LSBD)
Departamento de Computação
Universidade Federal do Ceará, Brazil

{gustavo.valenca, lucas.falcao}@lsbd.ufc.br,
{felipe.timbo, victor.farias, javam.machado}@lsbd.ufc.br

Abstract. *The increasing adoption of solid-state drives (SSDs) due to their high performance and reliability has made failure prediction crucial for ensuring data integrity and availability. Self-monitoring, Analysis, and Reporting Technology (SMART) is a system for drives that periodically reports various operational parameters that facilitate early detection of potential issues. Although many studies have used SMART attributes for approaching this matter – as a binary problem – we test new ways of predicting SSD failures, considering multiple health levels. In this paper, we first use feature selection for selecting the best SMART attributes as learning features. Then, we test the selected features on several classification models and two different prediction horizons of one month and one year ahead of the failure. The preliminary results effectively validate our approaches to address that problem, mainly in the smaller prediction horizon with non-linear models.*

1. Introduction

Solid State Drives (SSDs) are well known for their low latency capabilities and reliability over HDDs. Therefore, they have been widely used as a storage medium in modern data centers and households [Maneas et al. 2020]. Hence, the study of failure prediction in SSDs is a critical task to better maintain large storage reliability since it can be used to replace drives and avoid data loss and other associated costs.

The drive manufacturers implement the SMART (Self-Monitoring, Analysis, and Reporting Technology) [Ottem and Plummer 1995], a built-in system for HDDs and SSDs that gathers data on the state of the drive daily, providing a time series that can be used for failure prediction. Although it can change for each manufacturer, commonly gathered data generally includes Raw Read Error Rate, Power On Hours, Power Cycles, Device Temperature, Total LBAs Written, and Total LBAs Read. It is helpful since it provides several indicators that can be analyzed to check drive reliability based on thresholds. Despite all these preceding efforts, the exclusive use of SMART attributes to predict failure is not enough [Murray et al. 2005].

Machine Learning techniques have shown to be a good tool to determine whether a given drive is failing or not [Pereira et al. 2022, Xu et al. 2021] by using the SMART attributes as inputs to models. However, as the literature is more extensive on failure prediction in Hard Disk Drives, more study is necessary for SSDs.

In this paper, we focus on providing a different view on SSD failure prediction and present exploratory results. Although we are still using SMART attributes like other papers on SSDs [Chen et al. 2022], we are the first to study the modeling of the Remaining Useful Life (RUL) of the drive to six health levels [dos Santos Lima et al. 2017], representing time spans, instead of the common classification of healthy versus unhealthy. We test two horizons of prediction (1 year and 1 month). Also, we use a feature selection method to select the best learning features to increase prediction accuracy and reduce model complexity. Our results show that when non-linearity and temporal aspects are considered together, predicting RUL with health levels is less challenging, especially in a more tight time period, giving outcomes of up to around 90% accuracy depending on the selected manufacturer model.

2. Failure Prediction of Storage Devices

Studies have tackled the problem with various approaches to address the failure prediction task. One approach is to deal with the problem as a regression over the Remaining Useful Life [Lima et al. 2018], where the model tries to predict the number of days left the disk has before it fails. Another approach is to create classes over the RUL to perform a classification task [Lima et al. 2021]. A common modeling of the classes is to use the last month of the drive's life as the class, indicating the failure (unhealthy state) and all the rest as the healthy state. Also, it is noticeable a constant effort of researchers to perform feature selection, as not all SMART attributes can contribute to failure prediction [Felix et al. 2023].

Although HDD SMART data have been more extensively addressed for both regression and classification of the RUL, SSD SMART data have also been leveraged in other works, such as being used to investigate its correlation to failures [Han et al. 2021, Lu et al. 2022] and to find the best learning features for failure prediction in the classification approach of the last month versus the rest. [Xu et al. 2021]. Even though the usual 2-class classification problem for SSDs is suitable for a lot of cases, it may be of one's interest to increase the granularity of the failure prediction to better understand the state of the drives, as described in [dos Santos Lima et al. 2017] for HDDs. To the best of our knowledge, no study has ever performed this type of classification with SSD SMART data.

3. Methodology

Our main goal is to predict SSD failures with SMART data by classifying the drive's Remaining Useful Life in health levels, a higher granularity than the usually employed approach of 2-classes. We can divide our work into three main steps.

Data Preprocessing. For each SSD model, two prediction horizons are considered (360 and 30 days), labeling the day intervals in six health levels as depicted in Figure 1. Features with more than 99% of missing data are completely deleted.

Feature Selection. To evaluate the importance of SMART attributes in the SSD model dataset, four feature selection approaches that were presented in [Xu et al. 2021] are considered, covering a diversity of methods: Pearson correlation, Spearman correlation, Random forest, and XGBoost. After passing all features to all four methods, a ranking is obtained in every approach, and the first columns that together makeup at least

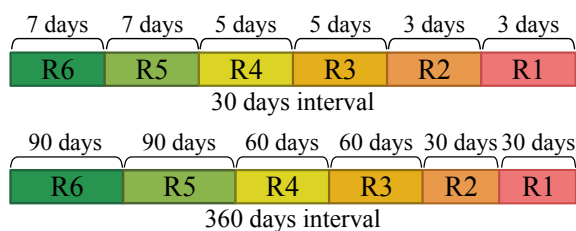


Figure 1. RUL interval setup used for the two prediction horizons.

70% of the total importance are picked as learning features for the given method and used in the prediction experiments.

Health Level Prediction. After the preceding steps are finished, the last task is to evaluate the classification. As an exploratory study, we selected three Machine Learning models as predictors based on their characteristics: Logistic Regression as a linear model, Random Forest as a nonlinear model, and Long Short-Term Memory network (LSTM) as a nonlinear model that considers the temporal aspect of the data. As each SSD manufacturer model can report different SMART features and may degrade differently, we train different classification models for each of them. Also, we train one model for each prediction horizon of 1 year and 1 month. For the experiments, we train and test failure prediction considering all columns as learning features, and considering only the ones of each feature selection approach.

To assess the performance of the failure prediction models, we are taking into account classic classification metrics, such as the Accuracy of the classification predictions, the Precision, and the Recall. Both Precision and Recall are calculated using the macro aggregation of the multiclass case. It is important to state that Precision is a very important metric for this type of application. This happens because, as the Precision measures how good the model does not label as positive a sample that is negative, it helps prevent the costs of replacing a healthy disk. However, as there are cases where high Precision does not represent the desired behavior, especially in the multiclass case, the F0.5 score is a good balanced metric that also considers the recall but puts a bigger weight on the Precision value.

4. Experiments and Results

This section provides a detailed analysis of the experiments conducted to evaluate the effectiveness of the proposed methodology.

4.1. Dataset

We use the data released by [Xu et al. 2021], as it was the most complete public dataset containing SMART time series. The dataset includes SMART logs and the timestamp of the failed drives from six different drive manufacturer models, comprising almost 500K SSDs, spanned over two years. Manufacturers are represented by MA, MB, and MC, following a number that refers to one of the two drive models included by each manufacturer (e.g. MA1, MA2). Every drive’s data is a daily SMART data time series that is potentially irregular, considering that SMART attributes are not collected daily for all drives. After data preprocessing, SMART features can make up to 42 columns, including raw and normalized values for each SMART. Altogether, around 16K of drives had failed

within the two-year period, with total percentages by SSD model as follows: 8.40% for MA1, 5.42% for MA2, 11.08% for MB1, 3.70% for MB2, 64.46% for MC1, and 6.94% for MC2. All six drive models are SATA SSDs.

4.2. Experimental procedure

In this study, we are training Machine Learning models to classify SSD Smart Data into different health levels. Therefore, we are splitting the dataset using 70% of the time series for training and 30% for testing, considering each time series as a data point instead of each observed day. The same training and testing split was used for the Feature Selection and the Health Level Prediction steps. Also, from the six models at hand in the dataset, we selected the three with more failure examples: MA1 (1,370), MB1 (1,807), and MC1 (10,510).

Three classification models are trained for the failure prediction step: Logistic Regression, Random Forest, and LSTM. For the Random Forest, the number of estimators is set to 100, and the maximum depth is set to 13. For the LSTM model, the hidden state of the recurrent cell is set to 32, and its output is passed through two fully connected neural network layers with outputs of size 16 and 6 and ReLU and Softmax activation functions, respectively. Also, since the data can be irregularly sampled, we are adding a time interval vector [Che et al. 2018] exclusively to the LSTM experiments. This new attribute contains the information of the number of days before the last SMART sampling for each row.

Except for the LSTM model, due to its high execution time, the experiments were performed three times, and the reported value is the average. All experiments were run using Scikit-Learn version 1.3.0 and PyTorch version 2.3.1

4.3. Results

Method	30 Days			360 Days		
	MA1	MB1	MC1	MA1	MB1	MC1
No Feature Selection (No FS)	40	38	42	40	38	42
Pearson	6	12	12	4	11	5
Spearman	5	7	9	7	10	8
Random Forest (RF)	4	4	3	5	5	4
XGBoost	10	11	13	13	9	10

Table 1. Number of features chosen for each Feature Selection method.

The number of selected attributes in the Feature Selection step is presented in Table 1, and the results for the health level prediction step are presented in Tables 2, 3, and 4 for SSD models MA1, MB1, and MC1, respectively. The tables show the results for the three classification models: Logistic Regression (LR), Random Forest (RF), and Long Short-Term Memory network (LSTM), for each of the Feature Selection methods: No Feature Selection (No FS), Pearson correlation, Spearman correlation, Random Forest (RF) importance, and XGBoost importance. The reported classification metrics are Accuracy (A), Precision (P), Recall (R), and F0.5 score. The best values for each metric are marked in bold font for each classification model, and the overall better is marked with an underline.

Method	30 Days				360 Days			
	A	P	R	F0.5	A	P	R	F0.5
LR (No FS)	23.7%	47.7%	18.3%	36.1%	32.2%	43.2%	25.4%	37.9%
LR (Pearson)	23.9%	62.2%	18.4%	42.2%	27.1%	56.0%	20.8%	41.8%
LR (Spearman)	23.5%	62.2%	18.2%	41.9%	27.5%	42.7%	21.3%	35.5%
LR (RF)	23.6%	74.6%	17.1%	44.6%	27.2%	37.3%	20.8%	32.2%
LR (XGBoost)	23.9%	47.9%	18.4%	36.3%	32.2%	43.3%	25.19%	37.8%
RF (No FS)	37.8%	42.5%	35.7%	41.0%	44.9%	48.3%	42.1%	46.9%
RF (Pearson)	32.1%	37.3%	28.1%	35.0%	41.8%	41.5%	40.6%	41.3%
RF (Spearman)	39.1%	39.0%	38.3%	38.9%	42.1%	42.3%	40.1%	41.8%
RF (RF)	38.3%	38.2%	36.4%	37.8%	47.9%	48.6%	46.7%	48.2%
RF (XGBoost)	41.0%	42.9%	39.8%	42.2%	49.8%	51.4%	47.6%	50.6%
LSTM (No FS)	70.0%	62.9%	66.4%	63.6%	43.4%	43.5%	44.3%	43.7%
LSTM (Pearson)	69.9%	67.3%	68.2%	67.5%	38.6%	41.3%	35.8%	40.1%
LSTM (Spearman)	78.6%	80.5%	76.8%	79.7%	43.7%	43.9%	43.4%	43.8%
LSTM (RF)	72.6%	74.2%	70.36%	73.3%	42.8%	43.5%	41.3%	43.0%
LSTM (XGBoost)	73.0%	72.9%	72.5%	72.8%	43.1%	43.6%	43.3%	43.6%

Table 2. Results for the MA1 manufacturer model.

In general, the results improve as the different classifiers consider non-linearity and the temporal aspect of the data, getting results as low as 24% for the linear Logistic Regression to results as good as 78% for the recurrent model on the model MA1. Also, the biggest prediction horizon of 360 days had overall results considerably worse than the 30-day horizon (e.g., 59% versus 83.1% on the model MB1). This can happen because, unlike HDDs, SSDs may show significant indications of failure when it is closer to the actual failure. However, further investigation is needed. To the LSTM model, this difference was significantly bigger. This can also be associated with the chosen structure of the network, where a more in-depth study of its hyperparameters, especially the size of the hidden state, can make the network deal better with long-term time-series predictions.

Method	30 Days				360 Days			
	A	P	R	F0.5	A	P	R	F0.5
LR (No FS)	23.3%	28.6%	16.8%	25.1%	30.7%	42.1%	27.7%	38.1%
LR (Pearson)	23.6%	55.8%	16.9%	38.3%	28.3%	49.2%	23.6%	40.4%
LR (Spearman)	23.7%	36.7%	16.9%	29.8%	28.9%	43.9%	26.1%	38.6%
LR (RF)	23.2%	74.4%	16.6%	43.8%	22.6%	52.5%	18.7%	38.6%
LR (XGBoost)	23.2%	57.7%	16.6%	38.6%	26.1%	44.3%	21.5%	36.5%
RF (No FS)	43.1%	64.2%	40.7%	57.5%	46.4%	51.7%	44.9%	50.2%
RF (Pearson)	40.1%	51.8%	37.9%	48.3%	45.2%	48.6%	44.1%	47.6%
RF (Spearman)	42.0%	52.9%	39.3%	49.5%	39.7%	43.8%	38.2%	42.6%
RF (RF)	40.6%	55.2%	36.0%	49.9%	46.9%	49.5%	45.9%	48.7%
RF (XGBoost)	44.7%	62.9%	42.4%	57.3%	46.9%	50.1%	45.9%	49.2%
LSTM (No FS)	80.1%	78.3%	79.9%	78.6%	59.3%	69.8%	55.5%	66.4%
LSTM (Pearson)	70.1%	78.3%	66.9%	75.8%	43.4%	45.9%	41.2%	44.8%
LSTM (Spearman)	81.6%	80.6%	80.5%	80.6%	51.0%	62.7%	48.4%	59.2%
LSTM (RF)	82.3%	80.4%	81.7%	80.7%	57.2%	67.5%	53.8%	64.2%
LSTM (XGBoost)	83.1%	80.7%	82.5%	81.1%	55.8%	67.7%	53.7%	64.3%

Table 3. Results for the MB1 manufacturer model.

Additionally, with a few exceptions for Feature Selection, the nonlinear methods

of importance measure performed better than the linear methods. This is another evidence that the problem of failure prediction for SSD is nonlinear. Furthermore, the Random Forest Feature Selection method consistently chose fewer features and still achieved the best results for the MC1 model in the 30-day prediction horizon.

Besides that, an interesting result is that the better reported Precision and F0.5 scores are consistently from the same classification models with better prediction accuracies. This is ideal for this type of application, as explained in Section 3.

Method	30 Days				360 Days			
	A	P	R	F0.5	A	P	R	F0.5
LR (No FS)	24.2%	61.4%	16.8%	40.1%	30.7%	44.3%	23.9%	37.9%
LR (Pearson)	24.3%	61.4%	16.8%	40.1%	30.5%	39.8%	23.6%	35.0%
LR (Spearman)	24.4%	74.8%	16.9%	44.4%	30.8%	43.6%	23.9%	37.5%
LR (RF)	24.1%	62.2%	16.7%	40.3%	29.9%	56.4%	23.1%	43.8%
LR (XGBoost)	24.3%	61.3%	16.8%	40.1%	30.9%	40.8%	24.0%	35.8%
RF (No FS)	28.7%	49.5%	22.0%	39.6%	42.7%	49.5%	37.0%	46.4%
RF (Pearson)	27.7%	43.4%	20.8%	35.6%	41.4%	41.3%	36.7%	40.3%
RF (Spearman)	27.9%	43.6%	21.4%	36.1%	39.9%	40.6%	34.6%	39.2%
RF (RF)	30.2%	48.4%	23.1%	39.7%	45.8%	45.4%	41.4%	44.6%
RF (XGBoost)	30.4%	52.5%	24.0%	42.5%	46.3%	48.5%	41.3%	46.9%
LSTM (No FS)	92.3%	89.6%	90.3%	89.8%	46.9%	51.4%	47.3%	50.5%
LSTM (Pearson)	85.7%	82.1%	84.7%	82.6%	60.9%	66.1%	61.6%	65.1%
LSTM (Spearman)	83.2%	79.3%	81.2%	79.7%	55.0%	61.3%	55.8%	60.1%
LSTM (RF)	95.1%	93.3%	93.6%	93.4%	49.8%	55.8%	50.0%	54.5%
LSTM (XGBoost)	89.7%	85.9%	87.8%	86.3%	46.3%	49.0%	45.8%	48.3%

Table 4. Results for the MC1 manufacturer model.

5. Conclusion

In this study, we presented a novel approach to SSD failure prediction by classifying the Remaining Useful Life (RUL) of the drives into six health levels, providing a more granular view than the traditional healthy versus unhealthy classification. Our methodology involved data preprocessing to handle SMART attributes, feature selection using various methods, and testing multiple machine learning models, including Logistic Regression, Random Forest, and Long Short-Term Memory (LSTM) networks, across two prediction horizons of one month and one year. Future work includes exploring more sophisticated neural network architectures such as Transformer models, as long with an extensive hyperparameter optimization to provide better performance in handling temporal data, capturing complex patterns in SMART attributes, and improving overall prediction accuracy.

Acknowledgment

This research was partially funded by Lenovo, as part of its R&D investment under Brazilian Informatics Law. Additional funding was provided by CNPq/Brazil under grant number 316729/2021-3.

References

- Che, Z., Purushotham, S., Cho, K., Sontag, D., and Liu, Y. (2018). Recurrent neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):6085.
- Chen, L., Zhu, Z., Li, A., Mashhadi, N., Frickey, R., Ye, J., and Guo, X. (2022). Ssd drive failure prediction on alibaba data center using machine learning. In *2022 IEEE International Memory Workshop (IMW)*, pages 1–4. IEEE.
- dos Santos Lima, F. D., Amaral, G. M. R., de Moura Leite, L. G., Gomes, J. P. P., and de Castro Machado, J. (2017). Predicting failures in hard drives with lstm networks. In *2017 Brazilian Conference on Intelligent Systems (BRACIS)*, pages 222–227. IEEE.
- Felix, G. L., Pereira, F. L., Praciano, F. D., Gomes, J. P., and Machado, J. C. (2023). Feature selection for remaining useful life prediction in hard disk drives with missing data. In *Anais Estendidos do XXXVIII Simpósio Brasileiro de Bancos de Dados*, pages 57–63. SBC.
- Han, S., Lee, P. P., Xu, F., Liu, Y., He, C., and Liu, J. (2021). An in-depth study of correlated failures in production ssd-based data centers. In *19th USENIX Conference on File and Storage Technologies (FAST 21)*. USENIX Association.
- Lima, F. D. S., Pereira, F. L. F., Chaves, I. C., Gomes, J. P. P., and Machado, J. C. (2018). Evaluation of recurrent neural networks for hard disk drives failure prediction. In *2018 7th Brazilian Conference on Intelligent Systems (BRACIS)*, pages 85–90. IEEE.
- Lima, F. D. S., Pereira, F. L. F., Chaves, I. C., Machado, J. C., and Gomes, J. P. P. (2021). Predicting the health degree of hard disk drives with asymmetric and ordinal deep neural models. *IEEE Transactions on Computers*, 70(2):188–198.
- Lu, R., Xu, E., Zhang, Y., Zhu, Z., Wang, M., Zhu, Z., Xue, G., Li, M., and Wu, J. (2022). NVMe SSD failures in the field: the Fail-Stop and the Fail-Slow. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*, pages 1005–1020, Carlsbad, CA. USENIX Association.
- Maneas, S., Mahdavian, K., Emami, T., and Schroeder, B. (2020). A study of {SSD} reliability in large scale enterprise storage deployments. In *18th USENIX Conference on File and Storage Technologies (FAST 20)*, pages 137–149.
- Murray, J. F., Hughes, G. F., Kreutz-Delgado, K., and Schuurmans, D. (2005). Machine learning methods for predicting failures in hard drives: A multiple-instance application. *Journal of Machine Learning Research*, 6(5).
- Ottum, E. and Plummer, J. (1995). Playing it smart: The emergence of reliability prediction technology. Technical report, Technical report, Seagate Technology Paper.
- Pereira, F. L. F., Bucar, R. C., Brito, F. T., Gomes, J. P. P., and Machado, J. C. (2022). Predicting failures in hdds with deep nn and irregularly-sampled data. In *Brazilian Conference on Intelligent Systems*, pages 196–209. Springer.
- Xu, F., Han, S., Lee, P. P., Liu, Y., He, C., and Liu, J. (2021). General feature selection for failure prediction in large-scale ssd deployment. In *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pages 263–270. IEEE.