



**NILKAMAL SCHOOL OF MATHEMATICS,
APPLIED STATISTICS & ANALYTICS**

A Study On Analysing Failure Rate Of Storage Devices Using New Weibull Log Logistic Grey Forecasting Model

**A DISSERTATION SUBMITTED TO
SVKM'S NMIMS (DEEMED TO BE UNIVERSITY)
IN PARTIAL FULFILMENT FOR THE DEGREE OF
MASTERS OF SCIENCE
IN
STATISTICS AND DATA SCIENCE**

**BY
RITESH PATIL
DARSHANA PARLE
SEJAL MAHADIK
REETIK PRAJAPATI
NISHAD SANGALE**

**UNDER THE SUPERVISION OF
PROF. VAIBHAV VASUNDEKAR**

NILKAMAL SCHOOL OF MATHEMATICS, APPLIED STATISTICS AND ANALYTICS

SVKM's Narsee Monjee Institute of Management Studies

(Deemed-To-Be-University)

V.L. Mehta Rd, Vile Parle (West), Mumbai – 400056

NOV- 2024

INDEX

| | |
|-------------------------|----|
| ACKNOWLEDGEMENT | 3 |
| MOTIVATION | 4 |
| ABSTRACT..... | 5 |
| INTRODUCTION | 6 |
| OBJECTIVES..... | 8 |
| METHODOLOGY | 9 |
| RESULTS | 14 |
| CONCLUSIONS | 18 |
| LIMITATIONS | 19 |
| SCOPE..... | 20 |
| REFERENCES | 21 |
| ANNEXURE/ APPENDIX..... | 22 |

Acknowledgement

We wholeheartedly thank our project mentor, Prof. Vaibhav Vasundekar for providing valuable guidance and support in all stages of our project work, his meticulous insights helped us to improve the quality of the project we have been making. Our project greatly benefited from his statistical knowledge and experience which were useful at the critical stages of the project.

We would like to show our gratitude towards the SDS course coordinator Dr. Pradnya Khandeparkar and Prof. Rushina Singhi for providing us with this opportunity to work on this project. We extend our heartfelt gratitude to all the faculty members for their support, suggestions, and encouragement.

Finally, we extend our heartfelt gratitude to all the people who directly or indirectly are involved with our project.

Motivation

The motivation behind this topic, “A Study of Analysing Failure Rate of Storage Devices Using New Weibull Log Logistic Grey Forecasting Model” is to improve the accuracy and stability of predicting hard disk drive (HDD) failures. Current HDD failure prediction methods often suffer from high false positive rates and significant resource consumption, which affects their overall efficiency. This paper introduces a novel forecasting model based on the Weibull Log-Logistic distribution, which is designed to capture the complex nonlinear interactions between factors that lead to HDD failure. By addressing the limitations of existing statistical and machine learning models, the proposed model incorporates a new accumulation generation operator and multiple interaction effects between variables, allowing for more precise failure prediction. This improvement can significantly reduce the economic impact of false alarms in early warning systems, enhancing the reliability of server operations that depend on the continuous functionality of hard disk drives.

Abstract

The paper presents a new forecasting model to predict hard disk drive (HDD) failures using a novel Weibull Log-Logistic Grey forecasting approach. HDD failure prediction is critical for the reliability of servers and storage systems. However, existing methods have limitations, such as high false positive rates and significant resource usage, which lead to inefficient early warning systems. This paper aims to address these issues by developing a more accurate and stable model that can effectively predict HDD failures and reduce economic losses caused by false alarms.

The new model is based on the Weibull Log-Logistic mixture distribution, which helps to better capture the failure trends in HDD data by reducing volatility and fitting complex patterns in the data. The model also introduces a new accumulation generation operator to improve the model's flexibility in handling nonlinear data. This operator helps in reducing the influence of erratic changes in the data, making the predictions more stable.

Another key feature of the model is the inclusion of multiple interaction effects, which allows it to consider the combined impact of various independent factors that influence HDD failure, such as temperature, disk activity levels, and SMART (Self-Monitoring, Analysis, and Reporting Technology) attributes. These attributes, collected by HDDs to monitor their own health, are critical in identifying potential failures. By considering the interaction between these factors, the model enhances its predictive power and accuracy.

The given model is applied to real-world data from four different HDD types provided by BackBlaze, a company that tracks HDD failures. The results demonstrated that the new model consistently outperformed existing models in terms of prediction accuracy. The model achieved lower Mean Absolute Percentage Error (MAPE) values, indicating that it was better at forecasting HDD failures with fewer false positives and improved reliability.

The paper concludes that the new Weibull Log-Logistic Grey forecasting model can be a valuable tool for companies that rely on large-scale data storage systems. It provides a more accurate and resource-efficient method for predicting HDD failures, which could lead to significant cost savings by preventing data loss and minimizing the need for costly, unnecessary HDD replacements. This paper also suggests that future research could further refine the model and explore its application to other types of mechanical systems with similar failure characteristics.

Keywords: Weibull Log-Logistic Grey forecasting model, Hard disk drive (HDD) failure prediction, Multiple interaction effects, Accumulation generation operator, SMART (Self-Monitoring, Analysis, and Reporting Technology) attributes.

Introduction

HDDs are a critical component of modern data infrastructure, from personal computers to large-scale cloud storage systems used by corporations and organizations. As the amount of data stored on these systems continues to grow, the need for reliable and efficient prediction methods to anticipate HDD failures has become increasingly important.

The failure of an HDD can lead to significant problems, including data loss, system downtime, and the high cost of replacement and recovery. Therefore, the ability to predict HDD failures before they occur is crucial in avoiding such disruptions. Effective early warning systems help companies minimize the risk of data loss and ensure smoother operations.

However, predicting HDD failures is a complex task. Current predictive models and techniques, while useful, have notable limitations. Many existing methods have high false positive rates, meaning they often mistakenly predict failure when there is none. This results in unnecessary replacements of HDDs, which leads to extra costs and wasted resources. Moreover, the accuracy of current models can be unreliable, especially when dealing with large and complex datasets. As a result, there is a need for new approaches that can provide more accurate and efficient predictions.

This paper aims to tackle these issues by proposing a new forecasting model, called the Weibull Log-Logistic Grey forecasting model, which is designed to improve the accuracy of HDD failure predictions while reducing false positives. The model builds on the strengths of two statistical distributions – Weibull and Log-Logistic – which are widely used in failure analysis and reliability studies. By combining these distributions, the proposed model captures the complex patterns and trends in HDD failure data more effectively.

In addition to this, the model introduces a new accumulation generation operator. This operator plays a crucial role in stabilizing the model by making it more flexible in handling nonlinear data. This is particularly important in the context of HDD failure prediction, where the data is often subject to fluctuations and erratic changes. The accumulation generation operator helps smooth out these fluctuations, enabling the model to provide more consistent and stable predictions.

Another important aspect of the proposed model is its ability to consider multiple interaction effects. HDD failures are rarely caused by a single factor. Instead, they are typically the result of a combination of various independent factors such as temperature, usage intensity, and the health status of the drive as monitored through SMART (Self-Monitoring, Analysis, and Reporting Technology) attributes. SMART attributes are a set of diagnostic features that HDDs use to assess their own health, and they provide valuable

information that can be used to predict potential failures. By incorporating the interaction between these factors, the model increases its predictive power and accuracy.

To test the effectiveness of the new model, we applied it to real-world data obtained from BackBlaze, a company that tracks HDD failures. BackBlaze's dataset includes information from four different types of HDDs, making it a valuable resource for validating the performance of the forecasting model. The results of the study showed that the Weibull Log-Logistic Grey forecasting model consistently outperformed existing methods in terms of accuracy. The model was able to predict failures with lower Mean Absolute Percentage Error (MAPE) values, which means it produced fewer errors and false positives compared to other models.

In conclusion, the introduction sets the stage for the rest of the paper by explaining the importance of HDD failure prediction and the limitations of current methods. The given research paper argue that their new Weibull Log-Logistic Grey forecasting model can help overcome these limitations by providing a more accurate, stable, and resource-efficient solution. The model's ability to capture complex patterns in HDD failure data and its consideration of multiple interaction effects make it a valuable tool for companies that rely on large-scale data storage systems. By reducing false positives and improving prediction accuracy, the model can help prevent costly data loss and minimize unnecessary HDD replacements.

Objectives

1. To test a new forecasting model to accurately predict hard disk drive (HDD) failures, improving the reliability of data storage systems.
2. To capture and fit the trend of the hard disk failure data flexibly and reduce the volatility.
3. To better handle the multiple interaction effects in hard disk drive failure data and reduce the economic losses caused by the false alarms in early warning system.
4. To study the grey model with non-linear and interactive features to forecast the failure trend of hard disk drives.
5. Implement the model to minimize incorrect failure predictions, reducing unnecessary HDD replacements and associated costs.
6. To apply the model to real-world HDD failure data for practical validation.

Methodology

The methodology follows several stages, beginning with the theoretical foundation from the Weibull and Log-Logistic distributions, proceeding with model development, parameter estimation, and evaluation on real-world data sets.

1. Data Preprocessing

Data Cleaning: The raw dataset from Backblaze is preprocessed to remove irrelevant and redundant data entries, ensuring only critical data attributes related to hard drive failures are retained.

SMART Attribute Selection: Based on prior studies and correlations identified within SMART data, specific attributes were selected for the model, including:

- **Reallocated Sectors Count** (SMART5),
- **Reported Uncorrectable Errors** (SMART187),
- **Command Timeout** (SMART188),
- **Current Pending Sector Count** (SMART197),
- **Uncorrectable Sector Count** (SMART198).

The selected SMART attributes are normalized to maintain consistency across different data scales, ensuring model stability and improving prediction accuracy.

2. Model Formulation: Weibull Log-Logistic Grey Forecasting Model (WLLGM)

- **Distribution Choice:** The Weibull Log-Logistic (WLL) distribution is chosen to handle failure data due to its flexibility in modeling time-to-failure data and accommodating non-linear behavior. This distribution combines the Weibull and Log-Logistic distributions to capture the variability and trend in the failure data more effectively than traditional distributions.

1. Weibull Distribution-

- **Weibull Distribution Function:**

$$f_1(t; \psi_1) = \frac{\beta_1}{\alpha_1} \left(\frac{t}{\alpha_1} \right)^{\beta_1-1} e^{-\left(\frac{t}{\alpha_1} \right)^{\beta_1}}$$

2. Log-Logistic Distribution-

- **Log-Logistic Distribution Function:**

$$f_2(t; \psi_2) = \frac{\beta_2}{\alpha_2} \left(\frac{t}{\alpha_2} \right)^{\beta_2-1} \left(1 + \left(\frac{t}{\alpha_2} \right)^{\beta_2} \right)^{-2}$$

3. Weibull Log-Logistic Mixture Distribution-

- **Weibull Log-Logistic Mixture Distribution:**

$$f(t; \psi) = pf_1(t; \psi_1) + (1 - p)f_2(t; \psi_2)$$

- **Accumulation Generation Operator:** To account for time dependencies and trends in failure data, an accumulation generation operator based on the WLL mixture distribution is introduced. This operator applies a weighted accumulation across data points, reducing volatility and smoothing the data to enhance the model's adaptability to the nonlinear characteristics of failure patterns in hard drives.

$$x^{(1)}(k) = \sum_{i=1}^k h(k-i+1, \psi)x^{(0)}(i)$$

- **Model Structure:** The proposed model, WLLGM(1, N), is based on the GM(1, N) grey model but extends it to include multiple interaction effects among SMART attributes. This model incorporates both direct effects and interaction effects between SMART attributes, capturing the compound influence of these variables on hard drive failure probabilities.

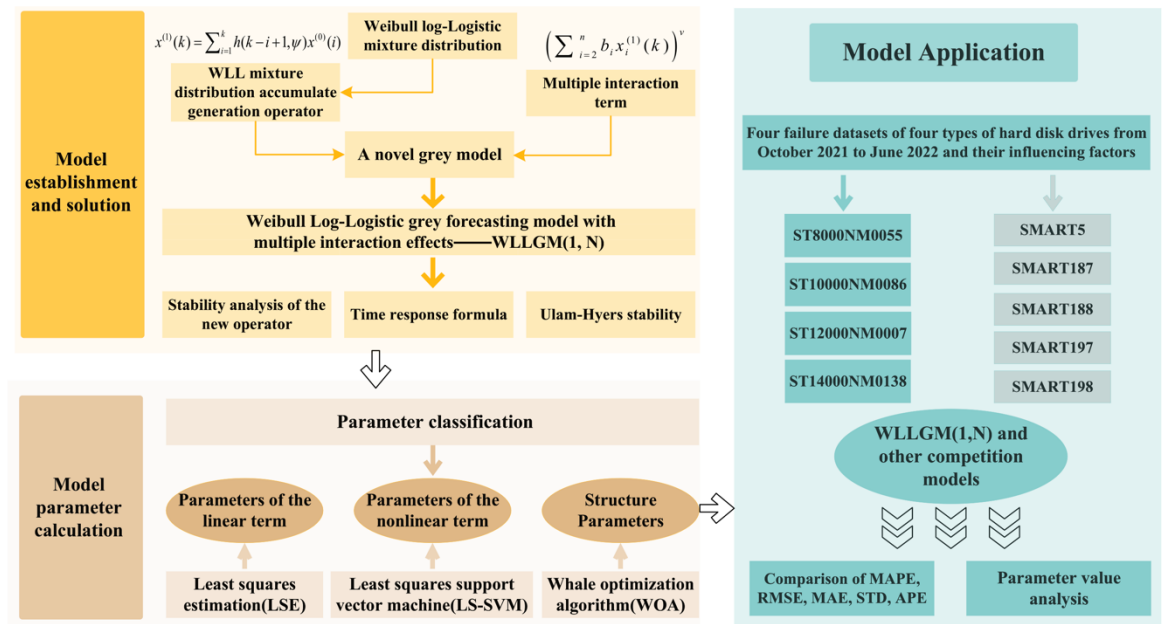


Fig. 1. Full-text framework diagram.

3. Parameter Estimation and Optimization

Parameter estimation for the WLLGM model is executed through a multi-step approach, applying both traditional and advanced optimization techniques to handle the linear, nonlinear, and complex interaction parameters within the model.

- **Linear Parameters:** Parameters representing linear relationships are estimated using the Least Squares Method.
- **Nonlinear Parameters:** Parameters that capture the nonlinear relationships among SMART attributes are estimated using the Least Squares Support Vector Machine (LS-SVM). LS-SVM offers advantages in handling small sample sizes and reducing sensitivity to outliers.
- **Structural Parameters:** Control parameters, such as the exponent for interaction terms and scaling factors, require meta-heuristic optimization due to the complex search space.
- **Optimization Using Whale Optimization Algorithm (WOA):**
The Whale Optimization Algorithm (WOA) is employed to optimize structure parameters like the power exponent and the control coefficient in the WLLGM model.
WOA is chosen for its ability to efficiently handle complex, non-linear optimization problems by simulating the social behavior of humpback whales. This algorithm iteratively refines the parameter values to minimize the Mean Absolute Percentage Error (MAPE) of the model's predictions.

4. Multiple Interaction Effects Modeling

To capture the compound effect of SMART attributes on failure predictions, the WLLGM model introduces multiple interaction terms. This involves:

- **Defining Interaction Terms:** The model considers interaction terms across selected SMART attributes to capture the non-linear dependencies and compounding effects that influence hard drive failure. For instance, increased values in Reallocated Sectors (SMART5) and Reported Uncorrectable Errors (SMART187) are compounded, creating a heightened failure risk.

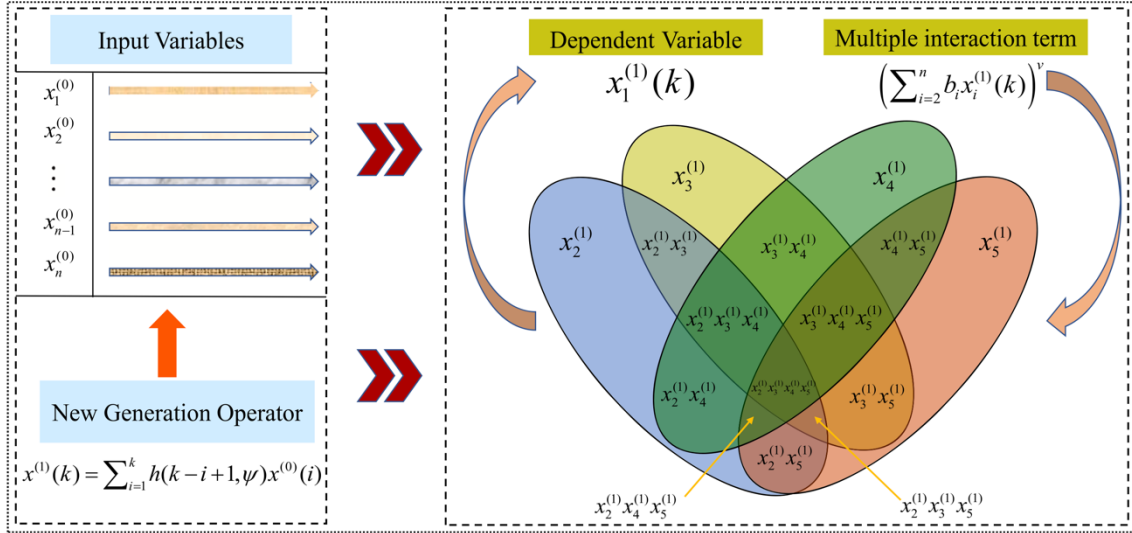


Fig. 2. The multiple interaction mechanism of the WLLGM(1, N) model.

5. Model Evaluation

Performance Metrics: These metrics are essential in assessing how closely the model's predictions align with the actual HDD failure rates. The model's effectiveness is assessed using several evaluation metrics:

- **Mean Absolute Percentage Error (MAPE)**- quantifying prediction accuracy
- **Root Mean Square Error (RMSE)** - assessing residuals' spread
- **Mean Absolute Error (MAE)** - assessing residuals' spread
- **Standard Deviation of Absolute Percentage Error (STD).**
- **Comparative Analysis:** The WLLGM model's performance is compared with Multiple Linear Regression (MLR), The WLLGM model is found to outperform MLR, particularly in predictive accuracy and robustness.

$$\text{MAPE} = \frac{1}{n} \sum_{s=1}^n \frac{|\hat{x}_1^{(0)}(s) - x_1^{(0)}(s)|}{x_1^{(0)}(s)} \times 100\%,$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{s=1}^n [\hat{x}_1^{(0)}(s) - x_1^{(0)}(s)]^2},$$

$$\text{MAE} = \frac{1}{n} \sum_{s=1}^n |\hat{x}_1^{(0)}(s) - x_1^{(0)}(s)|,$$

$$\text{STD} = \sqrt{\frac{1}{n} \sum_{s=1}^n \left(\frac{|\hat{x}_1^{(0)}(s) - x_1^{(0)}(s)|}{x_1^{(0)}(s)} - \text{MAPE} \right)^2},$$

$$\text{APE} = \frac{|\hat{x}_1^{(0)}(s) - x_1^{(0)}(s)|}{x_1^{(0)}(s)} \times 100\%.$$

The methodology described above outlines the process of developing and validating a new Weibull Log-Logistic Grey Forecasting Model (WLLGM) to predict HDD failures. By incorporating both Weibull and Log-Logistic distributions, along with sophisticated parameter estimation techniques, the model effectively captures the nonlinear trends and interactions in HDD failure data, significantly improving prediction accuracy and reliability.

Results

Table 1

Experiment results of WLLGM and competition models for hard disk drive failure.

| Model - ST4000DM000 | | | |
|---------------------|--------|---------|---------|
| Month | Actual | WLLGM | MLR |
| Fitting | | | |
| February | 22 | 20.1690 | 24.0485 |
| March | 26 | 25.3026 | 22.8306 |
| April | 20 | 20.7971 | 17.2586 |
| May | 26 | 26.1530 | 30.2511 |
| Predicted | | | |
| June | 22 | 21.8419 | 23.0721 |

| Model - ST14000NM0138 | | | |
|-----------------------|--------|----------|----------|
| Month | Actual | WLLGM | MLR |
| Fitting | | | |
| February | 112 | 110.3756 | 115.5808 |
| March | 93 | 92.3876 | 98.2647 |
| April | 86 | 88.2275 | 90.0256 |
| May | 89 | 91.6438 | 85.2450 |
| Predicted | | | |
| June | 113 | 111.5675 | 118.8821 |

| Model - ST14000NM001G | | | |
|-----------------------|--------|---------|---------|
| Month | Actual | WLLGM | MLR |
| Fitting | | | |
| February | 12 | 12.7933 | 16.2408 |
| March | 13 | 12.5304 | 15.2672 |
| April | 16 | 18.3272 | 14.0256 |
| May | 19 | 20.3961 | 21.0256 |
| Predicted | | | |
| June | 13 | 12.7079 | 14.0256 |

Table 1 presents the results obtained by the WLLGM(1, N) model against the competition model MLR, where the actual data represents the total number of failures of hard disk drives per month. it can be found that that for the hard disk drives ST4000DM000 ,ST14000NM0138 ,ST14000NM001G , the fitting and prediction values of WLLGM are closest to the actual values, highlighting better fitting with less deviation.

Table 2**Evaluation Metrics of models for fitting and prediction.**

| Model - ST4000DM000 | | |
|---------------------|--------|--------|
| Criterion | WLLGM | MLR |
| Fitting | | |
| MAPE(%) | 0.4474 | 1.0425 |
| RMSE | 0.1422 | 0.3439 |
| STD% | 0.1179 | 0.3734 |
| MAE | 0.082 | 0.0303 |
| Prediction | | |
| MAPE(%) | 0.5567 | 2.3328 |
| RMSE | 0.7689 | 0.3728 |
| STD% | 0.5104 | 0.4212 |
| MAE | 0.6119 | 0.0312 |

| Model - ST14000NM0138 | | |
|-----------------------|--------|--------|
| Criterion | WLLGM | MLR |
| Fitting | | |
| MAPE(%) | 0.5564 | 0.8745 |
| RMSE | 0.2887 | 0.3988 |
| STD% | 0.2373 | 0.4284 |
| MAE | 0.1958 | 0.0533 |
| Prediction | | |
| MAPE(%) | 0.3425 | 4.615 |
| RMSE | 0.6275 | 0.4748 |
| STD% | 0.3699 | 0.5005 |
| MAE | 0.4308 | 0.0594 |

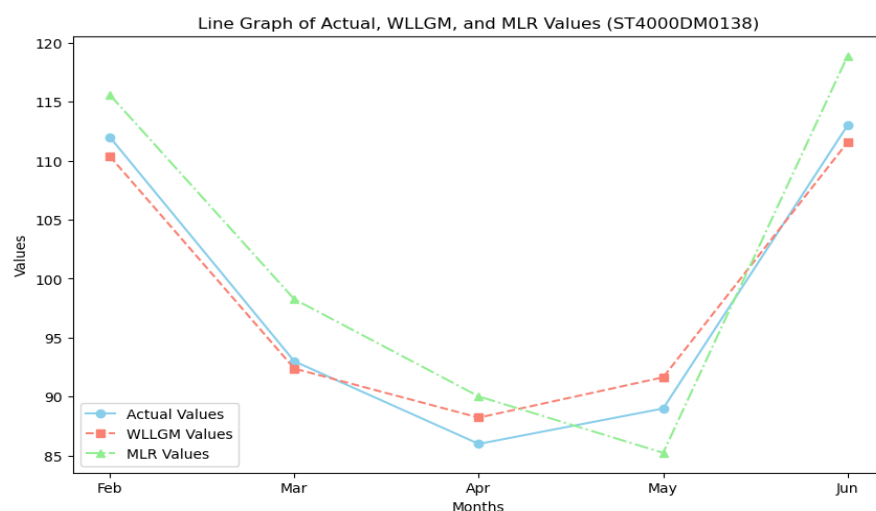
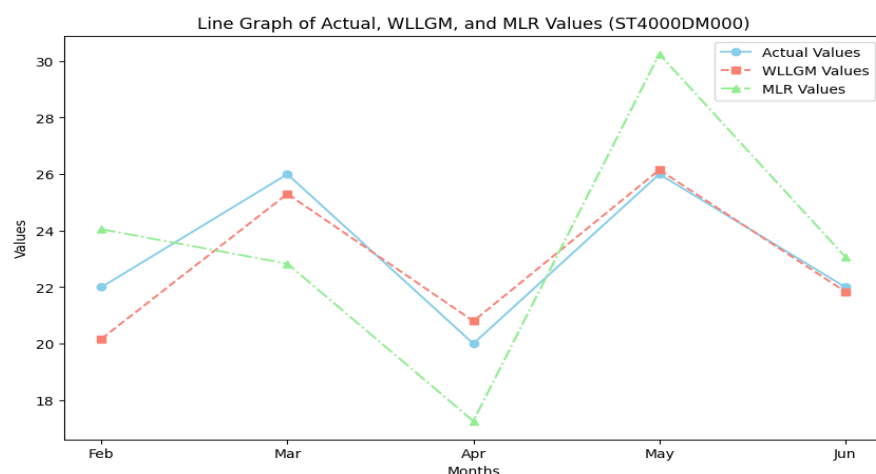
| Model - ST14000NM001G | | |
|-----------------------|--------|--------|
| Criterion | WLLGM | MLR |
| Fitting | | |
| MAPE(%) | 0.2002 | 1.2036 |
| RMSE | 0.1890 | 0.3624 |
| STD% | 0.0580 | 0.1254 |
| MAE | 0.1027 | 0.2145 |
| Prediction | | |
| MAPE(%) | 0.6077 | 4.1824 |
| RMSE | 0.5737 | 0.3214 |
| STD% | 0.3664 | 0.2823 |
| MAE | 0.557 | 0.7523 |

Table 2 presents the result of the model evaluation matrix, we can conclude that, for each dataset, the WLLGM consistently outperforms MLR across metrics:

WLLGM has a significantly lower MAPE, reflecting its accuracy in both fitted and predicted values. It's lower RMSE across datasets shows it more effectively minimizes error magnitude. It also shows lower variability in prediction errors, which enhances its robustness across datasets. WLLGM typically had a lower MAE, indicating that its errors are consistently smaller on average compared to MLR, aligning with WLLGM's focus on minimizing prediction errors.

Hence, the Weibull Log-Logistic Grey Forecasting Model (WLLGM) demonstrates superior accuracy and consistency over Multiple Linear Regression (MLR) for HDD failure prediction. The WLLGM's lower MAPE, RMSE, STD, and MAE values imply it can better adapt to the complex and non-linear interactions in HDD failure data.

The visualization results of each model.



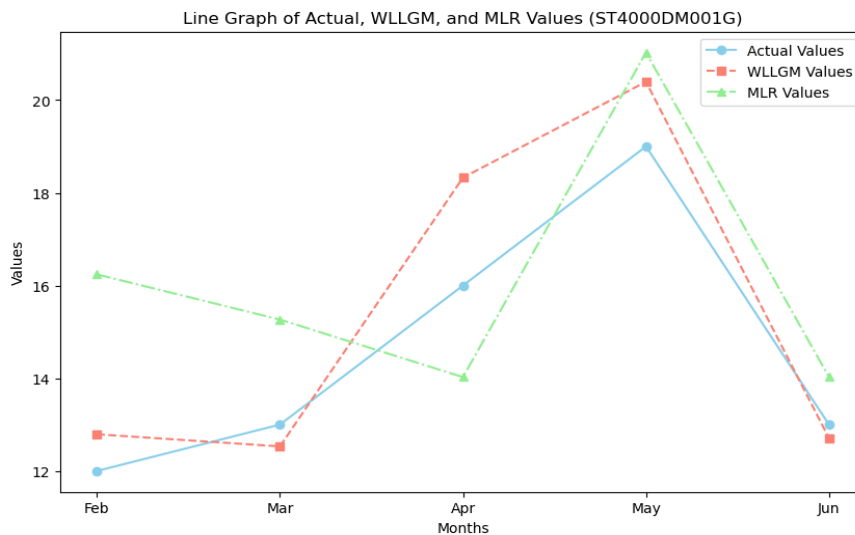


Table 3

The results of parameters of WLLGM.

| Parameter | ST4000DM000 | ST14000NM0138 | ST14000NM001G |
|-----------|-----------------|-----------------|-----------------|
| v | 1.3953904 | 1.90212157 | 0.11779333 |
| u | 1.20138130 | 0.86652997 | 0.64586261 |
| a | -0.00348812 | -0.00049662 | 0.00600237 |
| b2 | 4.04635767e-01 | 4.55979758e-01 | 2.84539493e-01 |
| b3 | 1.23023154e-04 | 5.02074157e-06 | 1.27294058e-05 |
| b4 | -1.30001196e-01 | 3.43576572e-06 | -1.30049882e-01 |
| b5 | -2.25882886e-06 | -3.39541682e-06 | 2.18662073e-04 |
| b6 | -2.25882888e-06 | -3.39541674e-06 | 2.18662073e-04 |
| b7 | 0.50247033 | 0.49206755 | 0.49895195 |

Table 3 presents the values of the parameters of the proposed model for forecasting the failure of hard disk drives. Compared to ST4000DM000 and ST14000NM001G, which have stronger nonlinear relationships, ST14000NM0138's failure patterns are simpler due to its minimal nonlinearity, as indicated by the nonlinear degree (v) parameter. The consistency of the model is shown by the development coefficient (a), where a negative value suggests high model alignment with observed failure patterns and a positive but minor value represents dependable prediction performance. Lastly, the influence coefficients (bi) show that the most important feature for HDD failures is SMART5 (Reallocated Sectors Count), highlighting the significance of keeping an eye on this particular property.

Conclusions

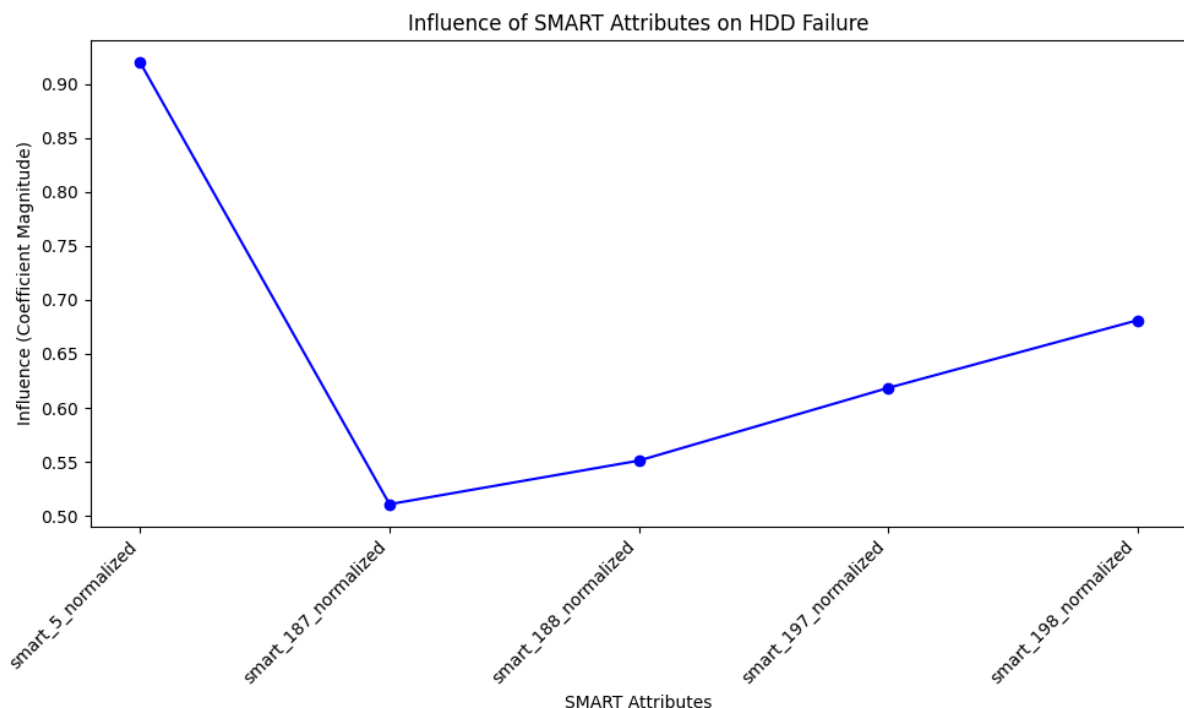
1.A novel Weibull Log-Logistic grey forecasting model with multiple interaction effects is proposed in this paper to accurately and stably predict hard disk drive failure.

2.The model's ability to account for multiple interaction effects between SMART attributes, enhances its predictive power. This multi-factor approach allows for a more comprehensive understanding of the causes behind HDD failures.

3.The new parameter estimation method is effective for the proposed model. The model was validated by taking three different varieties of hard disk drives as datasets, which shows superior prediction performance and robustness compared to the MLR model and can better characterize hard disk drive failure data.

4.The model effectively reduces false positive rates, which is a critical improvement over traditional HDD failure prediction models. This leads to fewer unnecessary HDD replacements, saving resources and costs.

5.Among the five SMART attributes that are most associated with hard disk drive failure, the reported uncorrectable errors, command timeout, current pending sector count, and uncorrectable sector count have the same influence on the failure of hard disk drives. The reallocated sectors count has the greatest impact.



Limitations

1. The model was validated using data from a specific source (BackBlaze) with particular HDD types. The performance of the model may vary when applied to different datasets or HDD brands, limiting its generalizability to all types of hard disk drives.
2. The model focuses on predicting failures based on historical data and certain factors (e.g., SMART attributes). It may not capture other relevant factors, such as manufacturing defects, power fluctuations, or external environmental factors, which can also lead to HDD failures.
3. The model's real-time prediction capability has not been fully explored. Implementing it in systems requiring live, continuous monitoring might pose challenges, such as data processing speed and resource allocation.
4. The reliance on SMART attributes assumes that these indicators are always accurate and sufficient for predicting failures. However, in some cases, SMART data may not capture all failure modes or might fail to provide early warning.
5. The current model is tailored for HDD failure prediction. Its application to other types of data storage technologies, such as solid-state drives (SSDs) or hybrid systems, is not explored, limiting its wider application.
6. While the model shows promising results, further refinement is necessary, particularly in addressing more complex nonlinearities and optimizing the balance between prediction accuracy and computational efficiency.

Scope

1. The current model is designed specifically for HDDs. Future research could extend the model to predict failures in other storage technologies such as solid-state drives (SSDs), hybrid storage systems, or cloud-based storage, where failure mechanisms and data characteristics differ.
2. Exploring the potential of integrating the Weibull Log-Logistic Grey forecasting model with machine learning algorithms could lead to hybrid models that combine statistical methods with advanced predictive analytics. This approach might further enhance prediction accuracy and adaptability to different datasets.
3. Since the Weibull and Log-Logistic distributions are commonly used in reliability analysis, future research could explore applying the proposed model to other mechanical systems with similar failure characteristics, such as turbines, vehicles, or medical devices, broadening its utility beyond HDDs.

References

- [1] G.F. Hughes, J.F. Murray, K. Kreutz-Delgado, et al., Improved disk-drive failure warnings, *Trans. Reliab.* 51 (2002) 350–357, <https://doi.org/10.1109/tr.2002.802886>.
- [2] E. Pinheiro, W.D. Weber, L.A. Barroso, Usenix, Failure trends in a large disk drive population, in: *Proceedings of the 5th USENIX Conference on File and Storage Technologies*, Usenix Assoc, San Jose, CA, 2007, p. 2. <https://dl.acm.org/doi/10.5555/1267903.1267905>.
- [3] S. Sankar, M. Shaw, K. Vaid, et al., Datacenter scale evaluation of the impact of temperature on hard disk drive failures, *ACM Trans. Storage* 9 (2013) 1–24, <https://doi.org/10.1145/2491472.2491475>.
- [4] B.D. Strom, S. Lee, G.W. Tyndall, et al., Hard disk drive reliability modeling and failure prediction, *Trans. Magn.* 43 (2007) 3676–3684, <https://doi.org/10.1109/tmag.2007.902969>.
- [5] L.P. Queiroz, F.C.M. Rodrigues, J.P.P. Gomes, et al., A fault detection method for hard disk drives based on mixture of Gaussians and nonparametric statistics, *Trans. Ind. Inform.* 13 (2017) 542–550, <https://doi.org/10.1109/tii.2016.2619180>.
- [6] R. Chen, X. Xiao, M. Gao, et al., A novel mixed frequency sampling discrete grey model for forecasting hard disk drive failure, *ISA Trans.* (2024), <https://doi.org/10.1016/j.isatra.2024.02.023>.
- [7] J. Zhao, Y.Z. He, H.M. Liu, et al., Disk failure early warning based on the characteristics of customized SMART, in: *Proceedings of the 19th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems*, IEEE, Electr Network, 2020, pp. 1282–1288, <https://doi.org/10.1109/ITherm45881.2020.9190324>.
- [8] V. Tomer, V. Sharma, S. Gupta, et al., Hard disk drive failure prediction using SMART attribute, *Mater. Today: Proc.* 46 (2021) 11258–11262, <https://doi.org/10.1016/j.matpr.2021.03.229>.
- [9] R. Jiang, D. Murthy, Mixture of Weibull distributions—Parametric characterization of failure rate function, *Appl. Stoch. Models Data. Anal.* 14 (1998) 47–65, [https://doi.org/10.1002/\(SICI\)1099-0747\(199803\)14:1<47::AID-ASM306>3.0.CO;2-E](https://doi.org/10.1002/(SICI)1099-0747(199803)14:1<47::AID-ASM306>3.0.CO;2-E).

Annexure/ Appendix

LOAD THE DATA AND SPLIT INTO TRAINING AND TESTING SETS WHERE X ARE THE INDEPENDENT SMART ATTRIBUTES I.E SMART_198_NORMALIZED,SMART_197_NORMALIZED,SMART_188_NORMALIZED,SMART_187_NORMALIZED AND THE DEPENDENT VARIABLE IS SMART_5_NORMALIZED WHICH AFFECTS THE FAILURE RATE.

''''''

IMPORTING THE REQUIRED LIBRARIES

!PIP INSTALL OPENPYXL

IMPORT NUMPY AS NP

IMPORT PANDAS AS PD

FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT

FROM SCIPY.SIGNAL IMPORT CONVOLVE

FOR GOOGLE COLAB:

FROM GOOGLE.COLAB IMPORT FILES

UPLOADED = FILES.UPLOAD()

GET THE FILENAME

FILENAME = LIST(UPLOADED.KEYS())[0]

READ THE EXCEL FILE USING PANDAS

DATA = PD.READ_EXCEL(FILENAME)

SELECT THE SMART ATTRIBUTES AS FEATURES

X = DATA[['SMART_5_NORMALIZED', 'SMART_187_NORMALIZED', 'SMART_188_NORMALIZED',
 'SMART_197_NORMALIZED', 'SMART_198_NORMALIZED']]

USING 'SMART_5_NORMALIZED' AS TARGET VARIABLE

Y = DATA['SMART_5_NORMALIZED']

SPLIT DATA INTO TRAINING AND TEST SETS

X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = TRAIN_TEST_SPLIT(X, Y, TEST_SIZE=0.2, RANDOM_STATE=42)

OUTPUT SHAPES OF THE TRAINING AND TEST SETS

PRINT("TRAINING SET SHAPE:", X_TRAIN.SHAPE, Y_TRAIN.SHAPE)

PRINT("TESTING SET SHAPE:", X_TEST.SHAPE, Y_TEST.SHAPE)

''''''DEFINING THE PDF OF EACH DISTRIBUTION SUCH AS WEIBULL AND LOG LOGISTIC AND THEN USING IT TO DEFINE THE WEIBULL LOG LOGISTIC MODEL AND THEN CALCULATE THE WLL ACCUMULATION GENERATION OPERATOR.'''''

```

FROM SCIPY.SIGNAL IMPORT CONVOLVE

DEF WEIBULL_COMPONENT(T, ALPHA1, BETA1):

    """WEIBULL PROBABILITY DENSITY FUNCTION (PDF) """

    RETURN (BETA1 / ALPHA1) * (T / ALPHA1)**(BETA1 - 1) * NP.EXP(-(T / ALPHA1)**BETA1)

DEF LOG_LOGISTIC_COMPONENT(T, ALPHA2, BETA2):

    """LOG-LOGISTIC PROBABILITY DENSITY FUNCTION (PDF) """

    RETURN (BETA2 / ALPHA2) * (T / ALPHA2)**(BETA2 - 1) / ((1 + (T / ALPHA2)**BETA2)**2)

DEF WLL_ACCUMULATION_OPERATOR_VECTORIZED(DATA, P=0.3, ALPHA1=1.5, ALPHA2=0.5,
BETA1=2.0, BETA2=2.5):

    """APPLIES THE WEIBULL LOG-LOGISTIC (WLL) ACCUMULATION GENERATION OPERATOR ON A DATA
SEQUENCE WITHOUT USING FOR LOOPS. """

    N = LEN(DATA)

    T = NP.ARANGE(1, N + 1) # TIME STEPS FROM 1 TO N

    # COMPUTE H(T) FOR ALL T USING VECTORIZED OPERATIONS

    H = P * WEIBULL_COMPONENT(T, ALPHA1, BETA1) + (1 - P) * LOG_LOGISTIC_COMPONENT(T,
ALPHA2, BETA2)

    # PERFORM CONVOLUTION BETWEEN DATA AND H(T)

    X1 = NP.CONVOLVE(DATA, H, MODE='FULL')[:N]

    RETURN X1

# ENSURE X_TRAIN CONTAINS ONLY NUMERIC DATA

X_TRAIN_NUMERIC = X_TRAIN.APPLY(PD.TO_NUMERIC, ERRORS='COERCE').FILLNA(0)

# APPLY WLL ACCUMULATION GENERATION OPERATOR TO EACH SMART ATTRIBUTE IN THE
TRAINING SET

X_TRAIN_WLL = X_TRAIN_NUMERIC.APPLY(LAMBDA COL:
WLL_ACCUMULATION_OPERATOR_VECTORIZED(COL.VALUES), AXIS=0)

# DISPLAY THE TRANSFORMED TRAINING SET

PRINT("TRANSFORMED TRAINING SET WITH WLL OPERATOR APPLIED:\n", X_TRAIN_WLL.HEAD())

"""TO CALCULATE THE HAZARD FUNCTION I.E THE VALUE OF FAILURE RATE OF WLLGM MODEL"""

P=0.3 #MIXING PROPORTION

```

```

ALPHA1=1.5

ALPHA2=0.5

BETA1=2.0

BETA2=2.5

N = LEN(DATA)

T = NP.ARANGE(1, N + 1)

H = P * WEIBULL_COMPONENT(T, ALPHA1, BETA1) + (1 - P) * LOG_LOGISTIC_COMPONENT(T,
ALPHA2, BETA2)

MEAN_H=NP.MEAN(H)

DEF RELIABILITY_FUNCTION(T, ALPHA1 , BETA1 ,ALPHA2, BETA2 ,P):

    R1= NP.EXP(-(T / ALPHA1)**BETA1)

    R2=1/ (1 + (T / ALPHA2)**BETA2)

    RETURN P*R1+(1-P)*R2

RELIABILITY=RELIABILITY_FUNCTION(T, ALPHA1 , BETA1 ,ALPHA2, BETA2 ,P)

DEF FAILURE_RATE(T):

    F1 = WEIBULL_COMPONENT(T, ALPHA1, BETA1)

    F2 = LOG_LOGISTIC_COMPONENT(T, ALPHA2, BETA2)

    RETURN (P * F1 + (1 - P) * F2) / RELIABILITY

PRINT(FAILURE_RATE(T))

# CALCULATE X(1)_I SEQUENCE USING THE ACCUMULATION OPERATOR

DEF COMPUTE_X1_SEQUENCE(XO_SEQUENCE):

    N = LEN(XO_SEQUENCE)

    X1_SEQUENCE = NP.ZEROS(N)

    FOR K IN RANGE(N):

        X1_SEQUENCE[K] = SUM(FAILURE_RATE(K - J + 1) * XO_SEQUENCE[J] FOR J IN RANGE(K + 1))

    RETURN X1_SEQUENCE

# GENERATE X(1)_I FOR EACH ORIGINAL X(O)_I SEQUENCE

X_OTHER_SEQUENCES = [COMPUTE_X1_SEQUENCE(XO) FOR XO IN X_ORIGINAL_SEQUENCES]

""""USING THE WHALE OPTIMIZATION TECHNIQUE TO CALCULATE THE META-HEURISTIC PARAMTERS
AND CALCULATING THE BEST MAPE I.E THE FITTED MAPE FOR THE MODEL.""""

IMPORT RANDOM

# DEFINE THE OBJECTIVE FUNCTION (MAPE)

DEF CALCULATE_MAPE(Y_TRUE, Y_PRED):

    """"CALCULATE MAPE BETWEEN ACTUAL AND PREDICTED VALUES.""""

```



```

    RETURN NP.MEAN(NP.ABS((Y_TRUE - Y_PRED) / Y_TRUE)) * 100

# FUNCTION TO CALCULATE WLLGM PREDICTIONS, CONSIDERING SMART ATTRIBUTE
# COEFFICIENTS

DEF WLLGM_PREDICT(X, MU, GAMMA, V, SMART_COEFFICIENTS):

    RETURN NP.DOT(X, SMART_COEFFICIENTS) * MU + GAMMA - V

# WHALE OPTIMIZATION ALGORITHM

DEF WHALE_OPTIMIZATION_ALGORITHM(X_TRAIN, Y_TRAIN, SEARCH_SPACE, ITERATIONS=20,
    WHALES_COUNT=10):

    # DEFINE PARAMETERS AND INITIALIZE WHALES RANDOMLY WITHIN THE SEARCH SPACE

    WHALES = [NP.ARRAY([RANDOM.UNIFORM(LOW, HIGH) FOR LOW, HIGH IN SEARCH_SPACE]) FOR _
    IN RANGE(WHALES_COUNT)]

    BEST_WHALE = NONE

    BEST_MAPE = FLOAT("INF")

    # INITIALIZE RANDOM SMART ATTRIBUTE COEFFICIENTS

    SMART_COEFFICIENTS = NP.RANDOM.RAND(X_TRAIN.SHAPE[1]) # ONE COEFFICIENT PER SMART
    ATTRIBUTE

    FOR ITERATION IN RANGE(ITERATIONS):

        FOR WHALE IN WHALES:

            MU, GAMMA, V = WHALE

            Y_PRED = WLLGM_PREDICT(X_TRAIN, MU, GAMMA, V, SMART_COEFFICIENTS)

            MAPE = CALCULATE_MAPE(Y_TRAIN, Y_PRED)

        # UPDATE BEST WHALE

        IF MAPE < BEST_MAPE:

            BEST_MAPE = MAPE

            BEST_WHALE = WHALE

    # WHALE POSITION UPDATE BASED ON ENCIRCLING BEHAVIOR AND RANDOM EXPLORATION

    FOR I IN RANGE(WHALES_COUNT):

        IF RANDOM.RANDOM() < 0.5:

            WHALES[I] = BEST_WHALE + (RANDOM.RANDOM() - 0.5) * 2 * BEST_WHALE # ENCIRCLE

        ELSE:

            WHALES[I] = NP.ARRAY([RANDOM.UNIFORM(LOW, HIGH) FOR LOW, HIGH IN
            SEARCH_SPACE]) # EXPLORE

    PRINT(F"ITERATION {ITERATION + 1}/{ITERATIONS}, BEST MAPE: {BEST_MAPE}")

    RETURN BEST_WHALE, BEST_MAPE, SMART_COEFFICIENTS

# DEFINE SEARCH SPACE FOR EACH PARAMETER

```

```

SEARCH_SPACE = [(0.01, 1), # RANGE FOR MU

                (0.01, 5), # RANGE FOR GAMMA

                (0.01, 2)] # RANGE FOR V

# RUN WOA TO FIND THE OPTIMAL PARAMETERS

BEST_PARAMS, BEST_MAPE, SMART_COEFFICIENTS =
WHALE_OPTIMIZATION_ALGORITHM(X_TRAIN_WLL, Y_TRAIN, SEARCH_SPACE)

PRINT("OPTIMAL PARAMETERS FOUND BY WOA:")

PRINT("MU:", BEST_PARAMS[0])

PRINT("GAMMA:", BEST_PARAMS[1])

PRINT("V:", BEST_PARAMS[2])

MODEL=WLLGM_PREDICT(X,BEST_PARAMS[0] ,BEST_PARAMS[1] , BEST_PARAMS[2],
SMART_COEFFICIENTS)

PRINT(NP.MEAN(MODEL))

""""USING THE LEAST SQUARE SUPPORT VECTOR MACHINE (LS-SVM) TO MODEL THE NON LINEAR
PARAMETERS AND TO CALCULATE THE TEST MAPE I.E FOR PREDICTED VALUES.""""

FROM SKLEARN.SVM IMPORT SVR

FROM SKLEARN.PREPROCESSING IMPORT STANDARDSCALER

# PARAMETERS OBTAINED FROM WOA

MU_OPT, GAMMA_OPT, V_OPT = BEST_PARAMS

# DEFINE THE LS-SVM

SCALER = STANDARDSCALER() # NORMALIZE THE DATA TO IMPROVE SVM PERFORMANCE
X_TRAIN_SCALED = SCALER.FIT_TRANSFORM(X_TRAIN_WLL) # SCALE TRANSFORMED DATA

# INSTANTIATE AND TRAIN THE SUPPORT VECTOR REGRESSOR (SVR) AS LS-SVM SUBSTITUTE
LS_SVM = SVR(KERNEL='LINEAR', C=1.0) # LINEAR KERNEL FOR LINEAR RELATIONSHIP
LS_SVM.FIT(X_TRAIN_SCALED, Y_TRAIN)

# OUTPUT THE SUPPORT VECTOR COEFFICIENTS (EQUIVALENT TO THE WEIGHTS BI IN THE PAPER)
SUPPORT_VECTOR_COEFFFS = LS_SVM.COEF_

PRINT("ESTIMATED COEFFICIENTS (B_I) FOR EACH FEATURE:", SUPPORT_VECTOR_COEFFFS)

# CALCULATE PREDICTIONS ON THE TEST SET FOR EVALUATION

X_TEST_SCALED = SCALER.TRANSFORM(X_TEST) # SCALE THE TEST SET WITH THE TRAINING
SCALER

Y_PRED = LS_SVM.PREDICT(X_TEST_SCALED)

# EVALUATE MODEL PERFORMANCE ON THE TEST SET

TEST_MAPE = CALCULATE_MAPE(Y_TEST, Y_PRED)

PRINT("TEST MAPE:", TEST_MAPE)

```

""""TO COMPUTE THE WLLGM VALUE OF THE MODEL""""

```
X1_SEQUENCE = NP.ARRAY(X_TRAIN_WLL)

X_OTHER_SEQUENCES = [NP.ARRAY([X1_SEQUENCE])]

# DEFINE Z(1)(K) BASED ON X1_SEQUENCE

DEF COMPUTE_Z1(K, X1_SEQUENCE):

    RETURN 0.5 * (X1_SEQUENCE[K-1] + X1_SEQUENCE[K])

# COMPUTE THE WLLGM(1, N) VALUE

DEF COMPUTE_WLLGM_1_N(K, X1_SEQUENCE, X_OTHER_SEQUENCES, A, GAMMA, V, B_I,
B_N_PLUS_1):

    # COMPUTE THE SUM TERM (SUM_I=2^N B_I * X(1)_I(K))^V

    SUM_TERM = SUM(B_I[I] * X_OTHER_SEQUENCES[I][K] FOR I IN RANGE(LEN(B_I)))

    INTERACTION_TERM = (SUM_TERM) ** V

# COMPUTE Z(1)(K)

    Z1_K = COMPUTE_Z1(K, X1_SEQUENCE)

# APPLY EQUATION (14)

    LEFT_TERM = X1_SEQUENCE[K] - X1_SEQUENCE[K-1] + A * Z1_K

    RIGHT_TERM = GAMMA * INTERACTION_TERM + B_N_PLUS_1

# CALCULATE WLLGM(1, N) AT INDEX K

    WLLGM_1_N_VALUE = LEFT_TERM - RIGHT_TERM

    RETURN WLLGM_1_N_VALUE

K = 2

WLLGM_VALUE = COMPUTE_WLLGM_1_N(K, X1_SEQUENCE, X_OTHER_SEQUENCES, A=0.5, GAMMA,
V, B_I, B_N_PLUS_1)

PRINT(F"WLLGM(1, N) VALUE AT K={K}: {WLLGM_VALUE}")
```

""""LASSO REGRESSION IS USED WHEN THE MAPE IS LARGER AND WE REDUCE IT BY TAKING THE CROSS VALIDATION TERMS.""""

```
# APPLY THE WLL ACCUMULATION GENERATION OPERATOR TO THE TEST SET AS WE DID WITH THE
TRAINING SET

X_TEST_WLL = X_TEST.APPLY(LAMBDA COL:
WLL_ACCUMULATION_OPERATOR_VECTORIZED(COL.VALUES), AXIS=0)

# NOW, SCALE THE TRANSFORMED TEST SET WITH THE SAME SCALER USED FOR THE TRAINING SET

X_TEST_SCALED = SCALER.TRANSFORM(X_TEST_WLL)

# RE-RUN LASSOCV WITH THE TRANSFORMED TEST SET
```

```

FROM SKLEARN.LINEAR_MODEL IMPORT LASSOCV

FROM SKLEARN.MODEL_SELECTION IMPORT CROSS_VAL_SCORE

# APPLY LASSOCV FOR REGULARIZATION TO MINIMIZE OVERFITTING

LASSO_MODEL = LASSOCV(ALPHAS=NP.LOGSPACE(-3, 3, 10), CV=5) # ADJUST ALPHA RANGE IF
NEEDED

LASSO_MODEL.FIT(X_TRAIN_SCALED, Y_TRAIN) # TRAIN WITH SCALED WLL-TRANSFORMED DATA

# EVALUATE WITH CROSS-VALIDATION

CV_MAPE_SCORES_LASSO = -CROSS_VAL_SCORE(LASSO_MODEL, X_TRAIN_SCALED, Y_TRAIN, CV=5,
SCORING='NEG_MEAN_ABSOLUTE_PERCENTAGE_ERROR')

PRINT("CROSS-VALIDATED MAPE WITH LASSO:", NP.MEAN(CV_MAPE_SCORES_LASSO))

# PREDICTIONS AND EVALUATION ON THE TEST SET WITH CONSISTENT SCALING

Y_PRED_LASSO = LASSO_MODEL.PREDICT(X_TEST_SCALED)

TEST_MAPE_LASSO = CALCULATE_MAPE(Y_TEST, Y_PRED_LASSO)

PRINT("TEST MAPE AFTER REGULARIZATION WITH LASSO:", TEST_MAPE_LASSO)

""""USED TO CALCULATE THE EVALUATION METRICS IN CASE OF LASSO REGRESSION""""

FROM SKLEARN.METRICS IMPORT MEAN_SQUARED_ERROR, MEAN_ABSOLUTE_ERROR

IMPORT NUMPY AS NP

# CALCULATE RMSE

RMSE = NP.SQRT(MEAN_SQUARED_ERROR(Y_TEST, Y_PRED_LASSO))

# CALCULATE MAE

MAE = MEAN_ABSOLUTE_ERROR(Y_TEST, Y_PRED_LASSO)

# CALCULATE STANDARD DEVIATION OF ABSOLUTE PERCENTAGE ERROR

APE = NP.ABS((Y_TEST - Y_PRED_LASSO) / Y_TEST) * 100

STD_APE = NP.STD(APE)

# DISPLAY METRICS

PRINT("EVALUATION METRICS ON TEST SET:")

PRINT(F"MAPE: {TEST_MAPE_LASSO}%")

PRINT(F"RMSE: {RMSE}")

PRINT(F"MAE: {MAE}")

PRINT(F"STD OF APE: {STD_APE}")

""""USED TO CALCULATE THE EVALUATION METRICS SUCH AS RMSE,MAE,STD FOR THE MODEL.""""

FROM SKLEARN.METRICS IMPORT MEAN_SQUARED_ERROR, MEAN_ABSOLUTE_ERROR

```

```

import numpy as np

# Calculate RMSE
rmse = np.sqrt(mean_squared_error(y_test, y_pred))

# Calculate MAE
mae = mean_absolute_error(y_test, y_pred)

# Calculate Standard Deviation of Absolute Percentage Error
ape = np.abs((y_test - y_pred) / y_test) * 100
std_ape = np.std(ape)

# Display Metrics
print("Evaluation Metrics on Test Set:")
print(f"MAPE: {test_mape}%")
print(f"RMSE: {rmse}")
print(f"MAE: {mae}")
print(f"STD of APE: {std_ape}")

"""To determine the smart attribute that affects the failure rate of the model the
most by comparing all the smart attributes."""

# Determine the most influential smart attribute
influence_strength = np.abs(smart_coefficients) # Absolute values of the
coefficients

# Get the index of the most influential smart attribute
most_influential_index = np.argmax(influence_strength)

# Print the smart attribute with the greatest influence
print(f"The most influential smart attribute is smart attribute
{most_influential_index + 1}")

# Display the influence of all smart attributes
print("Smart attribute influences:", influence_strength)

# Optional: Plot the influence of smart attributes
import matplotlib.pyplot as plt

smart_attribute_names = ['SMART_5_NORMALIZED', 'SMART_187_NORMALIZED',
'SMART_188_NORMALIZED', 'SMART_197_NORMALIZED', 'SMART_198_NORMALIZED']

plt.bar(smart_attribute_names, influence_strength)

plt.xlabel('Smart Attributes')
plt.ylabel('Influence (Coefficient Magnitude)')
plt.title('Influence of Smart Attributes on HDD Failure')
plt.show()

```

```

IMPORT PANDAS AS PD

# SELECT RELEVANT SMART ATTRIBUTE COLUMNS

SMART_COLUMNS = ['SMART_5_RAW', 'SMART_187_RAW', 'SMART_188_RAW', 'SMART_197_RAW',
'SMART_198_RAW']

SMART_DATA = FILTERED_DATA[SMART_COLUMNS]

# APPLY CUMULATIVE SUM TO EACH SMART ATTRIBUTE COLUMN

ACCUMULATED_DATA = SMART_DATA.CUMSUM()

# OPTIONALLY, VISUALIZE TO CHECK THE TREND

IMPORT MATPLOTLIB.PYLOT AS PLT

ACCUMULATED_DATA.PLOT(TITLE=F'ACCUMULATED DATA FOR {DATA}', FIGSIZE=(10, 6))

PLT.XLABEL("DATE")

PLT.YLABEL("ACCUMULATED VALUES")

PLT.SHOW()

"""TO COMPARE THE WLLGM MODEL WITH OTHER MODEL WE USE THE BELOW MLR MODEL"""

#IMPORTING THE LIBRARIES

IMPORT PANDAS AS PD

FROM GOOGLE.COLAB IMPORT FILES

FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT

FROM SKLEARN.LINEAR_MODEL IMPORT LINEARREGRESSION

FROM SKLEARN.METRICS IMPORT R2_SCORE, MEAN_ABSOLUTE_ERROR, MEAN_SQUARED_ERROR

IMPORT NUMPY AS NP

"""HERE WE CHECK FOR NULL OR MISSING VALUES AND DROP IF ANY. THEN WE SET THE
INDEPENDENT AND DEPENDENT VARIABLES AND SPLIT THE DATA INTO TEST AND TRAIN
DATASETS."""

# REMOVE ROWS WITH MISSING VALUES IF ANY

DATA = DATA.DROPNA()

# SELECT SMART ATTRIBUTES AS INDEPENDENT VARIABLES (BASED ON THE RESEARCH PAPER)

INDEPENDENT_VARS = ['SMART_198_NORMALIZED', 'SMART_188_NORMALIZED',
'SMART_187_NORMALIZED', 'SMART_197_NORMALIZED']

DEPENDENT_VAR = 'SMART_5_NORMALIZED'

X = DATA[INDEPENDENT_VARS]

Y = DATA[DEPENDENT_VAR]

# SPLIT THE DATA INTO TRAINING AND TESTING SETS

```

```

X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = TRAIN_TEST_SPLIT(X, Y, TEST_SIZE=0.3, RANDOM_STATE=42)

# APPLY MULTIPLE LINEAR REGRESSION TO ESTABLISH THE LINEAR RELATION

MODEL = LINEARREGRESSION()

MODEL.FIT(X_TRAIN, Y_TRAIN)

# STEP 7: RETRIEVE THE ESTIMATED PARAMETERS (COEFFICIENTS AND INTERCEPT)

COEFFICIENTS = MODEL.COEF_

INTERCEPT = MODEL.INTERCEPT_

PRINT("INTERCEPT:", INTERCEPT)

PRINT("COEFFICIENTS:", COEFFICIENTS)

# DISPLAY THE RESULTS

PRINT("\nESTIMATED PARAMETERS (MLR):")

FOR VAR, COEF IN ZIP(INDEPENDENT_VARS, COEFFICIENTS):

    PRINT(F"{VAR}: {COEF:.4F}")

""""CALCULATING THE EVALUATION METRICS FOR BOTH FITTED AND PREDICTED VALUES.""""

FROM SKLEARN.METRICS IMPORT MEAN_ABSOLUTE_ERROR, MEAN_SQUARED_ERROR

IMPORT NUMPY AS NP

# FUNCTION TO CALCULATE APE (ABSOLUTE PERCENTAGE ERROR)

DEF ABSOLUTE_PERCENTAGE_ERROR(Y_TRUE, Y_PRED):

    Y_TRUE, Y_PRED = NP.ARRAY(Y_TRUE), NP.ARRAY(Y_PRED)

    RETURN NP.ABS((Y_TRUE - Y_PRED) / Y_TRUE) * 100

# FUNCTION TO CALCULATE RMSE

DEF ROOT_MEAN_SQUARED_ERROR(Y_TRUE, Y_PRED):

    RETURN NP.SQRT(MEAN_SQUARED_ERROR(Y_TRUE, Y_PRED))

# FUNCTION TO CALCULATE STD OF ABSOLUTE PERCENTAGE ERROR

DEF STD_APE(Y_TRUE, Y_PRED):

    APE = ABSOLUTE_PERCENTAGE_ERROR(Y_TRUE, Y_PRED)

    RETURN NP.STD(APE)

# STEP 6: GET FITTED AND PREDICTED VALUES

FITTED_VALUES = MODEL.PREDICT(X_TRAIN) # FITTED VALUES FOR THE TRAINING SET

PREDICTED_VALUES = MODEL.PREDICT(X_TEST) # PREDICTED VALUES FOR THE TEST SET

# CALCULATE METRICS FOR FITTED VALUES (TRAINING SET)

MAPE_FITTED = NP.MEAN(ABSOLUTE_PERCENTAGE_ERROR(Y_TRAIN, FITTED_VALUES))

RMSE_FITTED = ROOT_MEAN_SQUARED_ERROR(Y_TRAIN, FITTED_VALUES)

```

```

STD_APE_FITTED = STD_APE(Y_TRAIN, FITTED_VALUES)

MAE_FITTED = MEAN_ABSOLUTE_ERROR(Y_TRAIN, FITTED_VALUES)

# CALCULATE METRICS FOR PREDICTED VALUES (TEST SET)

MAPE_PREDICTED = NP.MEAN(ABSOLUTE_PERCENTAGE_ERROR(Y_TEST, PREDICTED_VALUES))

RMSE_PREDICTED = ROOT_MEAN_SQUARED_ERROR(Y_TEST, PREDICTED_VALUES)

STD_APE_PREDICTED = STD_APE(Y_TEST, PREDICTED_VALUES)

MAE_PREDICTED = MEAN_ABSOLUTE_ERROR(Y_TEST, PREDICTED_VALUES)

# DISPLAY RESULTS

PRINT("METRICS FOR FITTED VALUES (TRAINING SET):")

PRINT(F"MAPE: {MAPE_FITTED:.2F}%")

PRINT(F"RMSE: {RMSE_FITTED:.2F}")

PRINT(F"STD OF APE: {STD_APE_FITTED:.2F}")

PRINT(F"MAE: {MAE_FITTED:.2F}")

PRINT("\nMETRICS FOR PREDICTED VALUES (TEST SET):")

PRINT(F"MAPE: {MAPE_PREDICTED:.2F}%")

PRINT(F"RMSE: {RMSE_PREDICTED:.2F}")

PRINT(F"STD OF APE: {STD_APE_PREDICTED:.2F}")

PRINT(F"MAE: {MAE_PREDICTED:.2F}")

```