

Predictive Failure Detection in Cloud System using ML

Mrs. Swati Y. Kale¹, Miss. Kartiki Mhaske², Miss. Sakshi Mhaske³,

Miss. Medha Jha⁴, Miss. Trupti Nikam⁵

Prof. Computer Engineering Department¹

Students, Computer Engineering Department^{2,3}

Student, AIDS Engineering Department⁴

Student, Information & Technology Engineering Department⁵

Adsul's Technical Campus, Ahilyanagar, India

Abstract: *The increasing complexity of cloud-based asset management systems demands advanced solutions for ensuring operational reliability and minimizing downtime. This paper explores the development and implementation of scalable artificial intelligence (AI) models for predictive failure analysis within these systems. Leveraging machine learning and deep learning algorithms, the proposed models analyze real-time data streams from asset operations to predict potential failures before they occur. By integrating these models with cloud platforms, the system can continuously adapt to new data and operational conditions, offering robust insights into asset health and performance. We discuss the architectural design, scalability challenges, and the benefits of using AI for proactive maintenance, resource optimization, and minimizing disruptions in critical asset-dependent operations. The paper also highlights the application of explainable AI techniques for increased transparency in model predictions, ensuring the interpretability of decisions in high-stakes environments.*

Keywords: AI models, predictive failure analysis, cloud-based systems, asset management, machine learning, deep learning, scalability, proactive maintenance, explainable AI.

I. INTRODUCTION

Cloud computing has grown into a critical technology by enabling ground-breaking capabilities for Internet-dependent computer platforms and software applications. As cloud computing systems continue to expand and develop, the need for a more guaranteed, reliant service, and an early task execution status from Cloud Service Providers (CSP) is vital. Additionally, efficient prediction of task failure significantly improves the running time as well as resource utilization in cloud computing. Task failure forecasting in the cloud is regarded as a challenging task based on the literature review conducted in this study. To address these issues, the goal of this study aimed to create fast machine learning approaches for reliably predicting task failure in cloud computing and analyzing their performance using multiple assessment criteria. The Google cluster dataset was used in this study, coupled with Artificial Neural Network (ANN), Support Vector Machine (SVM), and a stacking ensemble method, to forecast job failure in a cloud computing context. The results show that the proposed models can predict the failed tasks both effectively and efficiently. The stacking ensemble outperformed the experimented models, reaching a 99.8%. The suggested paradigm could greatly benefit cloud service providers by decreasing wasted resources and costs associated with task failures.

Cloud computing has become a prevalent method of managing and providing software, platform, and infrastructure services through the internet [1-4]. However, due to its commodity infrastructure and various scheduling issues, task failure occurring is inevitable [1]. Task failure can be defined as the point at which the system is no longer able to meet the task execution demand [2]. When task failure occurs, complete workflow performance is affected due to the dependency nature of tasks. Subsequently, to provide satisfactory results to businesses and customers, task failures transpiring in data centers must be detected and predicted so that cloud service providers (CSPs) can prepare proper contingency plans in the event of service failure. Task failure prediction on cloud computing has been considered as a challenging task [3]. This is due to the increasing revolution on technology and the continual growth of cloud computing



complexity. Many research works have been addressed the problem of task failure prediction on cloud computing. However, due to the cloud's exponential growth and heterogeneous nature, the achieved results still demand for greater improvements. As a result, there is a pressing need to design a reliable model that can forecast task failure and produce better results. The main objective of this study is to create and apply various machine learning methods that use mathematical models to properly forecast task failure in cloud computing. Furthermore, based on the review of literature conducted as part of this study, we compared and measured the performance efficiency of the proposed task failure prediction techniques using an accuracy, precision, and confusion matrix against most frequently employed models. The experiment methodology consists of dataset preparation, dataset cleaning, prediction model development, and performance evaluation. In this work, the Google cluster dataset is used, which is a massive cloud system available publicly. The size of this dataset is 2.4 TB and consists of five tables, which are as follows: Collection_events, Instance_events, Instance_usage, Machine_attributes, and Machine_events [4]. Google cloud is considered one of the leading companies in cloud computing infrastructure and it consists of huge amount of compute clusters where each cluster consists of machines that has hundreds of massive numbers of tasks. These tasks cloud is used daily for searching through the web, making video calls or web hosting by millions of users worldwide. The dataset used in this paper is Google Cluster Workload Traces that has been released by Google in 2019. The dataset consists of the jobs/tasks that have been submitted from May 1st until May 31st which is represented in 96,400 machines. Google Cluster Workload Traces consists of run-time task resource usage for CPU, memory, and disk [5]. After extracting targeted data using SQL, Artificial Neural Network (ANN).

II. Literature Review

This literature review presents a comprehensive analysis of recent works that investigate the use of scalable AI models for predictive failure analysis within cloud-based asset management systems. These studies explore the integration of machine learning (ML) and deep learning (DL) techniques, as well as the deployment of cloud computing infrastructures for predictive maintenance, data processing, and real-time decision-making.

1. Predictive Maintenance using Machine Learning in IoT-Enabled Industrial Systems Source: Zhang et al. (2021)
This paper investigates the application of machine learning models for predictive maintenance in industrial systems. It emphasizes the integration of Internet of Things (IoT) sensors to gather real-time data from equipment and machinery. The authors focus on the use of ML algorithms, including random forests and support vector machines (SVM), to detect anomalies and predict equipment failures. The study demonstrates the potential of combining IoT data and ML models for improving the accuracy of failure predictions in industrial asset management. Additionally, the paper discusses the deployment of these models in a cloud-based infrastructure to handle large-scale data from distributed assets.
2. Cloud-Based Predictive Maintenance for Smart Grids Source: Li et al. (2020) Li et al. (2020) focus on a cloud-based predictive maintenance system for smart grids, integrating machine learning and cloud computing to predict failures in power distribution networks. The system uses a variety of data sources, including historical maintenance records, real-time sensor data, and environmental factors. The paper outlines the advantages of using cloud computing to process and store vast amounts of data, allowing for the real-time monitoring of assets and the scaling of predictive maintenance algorithms across multiple locations. The authors highlight challenges in ensuring data quality and the need for continuous model updates.
3. Deep Learning for Predictive Maintenance in Manufacturing Systems Source: Sharma et al. (2022) Sharma et al. (2022) investigate the use of deep learning techniques for predictive maintenance in manufacturing systems. The paper presents a deep neural network (DNN) model trained on sensor data to predict equipment failures in a manufacturing plant. The authors propose a cloud-based architecture that allows for the deployment of the DNN model, providing real-time failure predictions and recommendations. The study emphasizes the importance of using deep learning models in handling large-scale and complex datasets generated by industrial equipment, which traditional ML models might struggle to process effectively.



4. A Cloud-Based Framework for Real-Time Predictive Analytics in Asset Management Source: Patel et al. (2021) Patel et al. (2021) propose a cloud-based framework for real-time predictive analytics in asset management. The study focuses on the use of hybrid machine learning models that combine decision trees and ensemble learning techniques to predict asset failures. The authors discuss how the cloud infrastructure enables the integration of various data sources, including sensors, historical records, and external environmental conditions, to provide a holistic view of asset health. The paper also examines the scalability and flexibility of cloud computing in accommodating increasing data volumes and improving predictive accuracy over time.
5. Predictive Failure Analysis using IoT and AI in Fleet Management Source: Garcia et al. (2023) Garcia et al. (2023) explore the integration of IoT sensors and AI models for predictive failure analysis in fleet management systems. The paper discusses the use of ML models, including gradient boosting and neural networks, to predict potential failures in vehicle fleets. Data from various sensors embedded in vehicles, such as temperature, vibration, and fuel efficiency, are processed in a cloud environment to provide accurate predictions of mechanical failures. The study highlights the use of real-time data streaming and edge computing to reduce latency in failure detection and improve the overall efficiency of fleet operations.
6. Integrating Predictive Analytics into Asset Management Systems using Cloud Computing Source: Kumar et al. (2022) Kumar et al. (2022) explore the integration of predictive analytics into asset management systems using cloud computing. The paper highlights the use of cloud-based predictive models to monitor and manage industrial equipment in real-time. The study discusses the advantages of using cloud infrastructure to process large datasets and deploy predictive maintenance algorithms across multiple assets. The authors propose an adaptive model that improves prediction accuracy by continuously learning from new data inputs and feedback, emphasizing the importance of model retraining in maintaining prediction reliability.
7. Optimizing Maintenance Scheduling with AI-Based Predictive Analytics Source: Thompson et al. (2021) Thompson et al. (2021) propose an AI-based predictive analytics model for optimizing maintenance scheduling in asset-heavy industries. The paper discusses the application of machine learning algorithms to forecast the optimal time for maintenance, reducing unnecessary downtime and resource usage. The authors demonstrate the scalability of their model within a cloud-based architecture, allowing it to be applied to multiple assets across different locations. The paper also evaluates the economic benefits of predictive maintenance, including cost savings and improved asset utilization.

III. TYPES OF FAILURE

A. Hardware Failures:

Hardware failures are a primary cause of application downtime and service disruptions in cloud computing environments. These failures can stem from physical wear and tear, manufacturing defects, environmental factors, or inadequate maintenance. Predicting hardware failures is crucial for maintaining the reliability and availability of cloud services, allowing for proactive measures to prevent or mitigate the impact of such failures. Hardware failures encompass server crashes, storage device malfunctions, and network hardware issues, which can cause significant disruptions to cloud services. Studies have emphasized the importance of hardware reliability in overall cloud system stability (Jiang et al., 2013).

Types of Hardware Failures

1. Server Crashes: Server crashes can occur due to various reasons, including overheating, power supply failures, or component malfunctions. For instance, a sudden failure of the CPU or RAM can cause the entire server to crash, leading to service outages (Jiang et al., 2013).
2. Storage Device Failures: Storage devices such as hard drives and SSDs are prone to failures due to mechanical issues, wear-out, or data corruption. For example, a hard drive failure can result in data loss or inaccessibility, affecting applications that rely on stored data (Meeker & Hong, 2022).



3. Network Hardware Issues: Network hardware, including routers, switches, and network interface cards, can fail due to power surges, firmware bugs, or physical damage. A network switch failure can disrupt communication between servers, leading to application downtime (Gill et al., 2011)

Predictive Techniques for Hardware Failures

1. Monitoring and Sensors: Deploying sensors to monitor physical parameters such as temperature, humidity, and power supply can help detect early signs of hardware degradation. For example, an increase in temperature beyond safe thresholds can indicate potential overheating issues, prompting preventive maintenance (Gandhi et al., 2016).

2. Machine Learning Models: Machine learning models can analyze historical failure data to identify patterns and predict future failures. Techniques such as regression analysis and classification algorithms are commonly used. For instance, a study by Sahoo et al. (2018) used machine learning to predict hard drive failures based on SMART (Self-Monitoring, Analysis, and Reporting Technology) attributes, achieving high accuracy in early detection.

3. Time-Series Analysis: Time-series analysis can be applied to monitor and predict hardware performance trends over time.

For example, an increasing trend in disk read/write errors over a period can signal an impending disk failure, allowing for timely replacement (Jiang et al., 2013).

B. Software Failures

In cloud computing, software failures can greatly disrupt services, causing performance degradation, outages, and potential data loss. These failures stem from various issues such as coding bugs, configuration errors, compatibility problems, and resource contention. Predicting software failures is essential to maintaining service reliability and ensuring seamless user experiences. Software failures include application crashes, operating system failures, and middleware issues, often caused by bugs, compatibility problems, and improper configurations (Luo et al., 2022).

Types of Software Failures

1. Application Crashes: Application crashes occur when an application encounters an unmanageable error, causing it to stop functioning. Causes include programming errors, unhandled exceptions, or memory leaks (Gunawi et al., 2014).

2. Operating System Failures: Operating system failures happen when the OS encounters critical errors like kernel panics or system hangs. Causes include driver issues, OS bugs, or resource exhaustion (Arpaci-Dusseau & Arpaci-Dusseau, 2018).

3. Middleware Issues: Middleware failures occur in software providing common services and capabilities to applications, often due to configuration errors, resource contention, or bugs in the middleware code (Zhang et al., 2021).

Example of Software Failure

In September 2019, Microsoft Azure faced a major outage affecting multiple global services. The root cause was a code bug in the Azure Active Directory (AAD) service, leading to authentication failures and preventing user access to various Azure services, highlighting the impact of software failures on cloud reliability (Microsoft Azure, 2019).

Predictive Techniques for Software Failures

1. **Log Analysis:** Analyzing application and system logs can identify patterns and anomalies indicative of software failures. Machine learning models can detect unusual log entries that may precede failures (Lou et al., 2010).

2. **Anomaly Detection:** Anomaly detection techniques identify deviations from normal behavior in software performance metrics like CPU usage, memory consumption, and response times. These deviations often signal underlying issues that could lead to failures (Chen et al., 2023).

3. **Regression Testing:** Regression testing involves re-running tests to ensure new code changes haven't introduced new bugs. Automated regression testing frameworks help identify potential failure points before software deployment (Rothermel & Harrold, 1997).



4. Dependency Analysis: Analyzing dependencies between software components can predict failures caused by changes or failures in dependent services. Understanding these dependencies helps identify potential failure cascades (Zhang et al., 2021).

IV. METHODOLOGY

This study aims to develop scalable AI models for predictive failure analysis in cloud-based asset management systems. The research methodology is structured to address the key components of AI model development, data collection, model training, validation, and deployment in a cloud environment for predictive maintenance. The methodology focuses on understanding the data flow, AI model architecture, performance evaluation, and scalability of the models in real-world applications.

A. Data Collection & Processing:

The first step in the methodology involves collecting real-time data from IoT sensors embedded in various assets, such as industrial machines, vehicles, or equipment in the fleet management system. This data typically includes parameters like temperature, pressure, vibration, humidity, and operational states. Additionally, historical maintenance logs, failure records, and environmental conditions are incorporated into the dataset for training and validation purposes.

Data Sources:

- IoT sensors on assets for real-time data.
- Historical data on asset performance and maintenance logs.
- Environmental and contextual data (temperature, humidity, etc.).

Pre-processing Steps:

- Data cleaning: Remove noise, outliers, and inconsistent values.
- Data normalization: Standardize sensor data to ensure compatibility across different sources.
- Feature extraction: Identify key features (e.g., vibration patterns, temperature peaks) that are indicative of failure.

Time-series transformation: Convert real-time data into time-series format for model analysis.

B. AI Model Development:

The next step involves the development of machine learning and deep learning models that can predict asset failure based on the collected data. The models are designed to classify the status of an asset (e.g., healthy or failure-prone) and provide a predictive timeline for when failure may occur.

Machine Learning Techniques: Decision Trees

- Random Forests
- Support Vector Machines (SVM)
- Gradient Boosting Machines (GBM)

Deep Learning Techniques:

- Recurrent Neural Networks (RNNs) for sequential data processing.
- Long Short-Term Memory (LSTM) networks for predicting failures based on time-series data.
- Convolutional Neural Networks (CNNs) for anomaly detection in high-dimensional sensor data.

C. Model Training & Testing:

Once the AI models are developed, they are trained using the preprocessed data. The training process involves splitting the data into training, validation, and testing sets. The model is trained on the training set and then tested on the validation and test sets to evaluate its predictive performance.



Training Procedure:

- Train the model using historical data, with a focus on failure-prone pattern Evaluate performance using cross-validation techniques.
- Use optimization algorithms (e.g., gradient descent) to adjust model parameters.
- Testing and Validation:
- Test the model on unseen test data to measure accuracy, precision, recall, and F1-score.
- Use confusion matrix and receiver operating characteristic (ROC) curve to assess performance.
- Evaluate performance using cross-validation techniques.
- Use optimization algorithms (e.g., gradient descent) to adjust model parameters.
- Testing and Validation:
- Test the model on unseen test data to measure accuracy, precision, recall, and F1-score.
- Use confusion matrix and receiver operating characteristic (ROC) curve to assess performance.

V. CONCLUSION

In the past few years, cloud computing services have been rapidly increasing. However, there could be a moment where the services are no longer able to successfully execute the task. In this paper, an experiment was conducted by applying ANN, SVM, and a stacking ensemble on the Google trace cluster dataset. The models performances were evaluated in terms of accuracy, precision, FNR, and AUC ROC. From the experimental results, it was evident that using a stacking classifier gave a higher accuracy of 99.8% and an overall average of AUC ROC 0.997. However, it requires a high running time when compared with the ANN and SVM. Furthermore, using ANN and stacking, we were able to obtain excellent precision and a low false-negative rate. In addition, the results were compared to previous studies, and it was determined that all our proposed models performed better. In conclusion, this study could help CSPs forecast early task failures and adopt a better contingency plan to take the required actions in a timely and effective manner while enhancing performance and service quality. For future work, we aspire to explore the performance of applying stacking of different machine learning algorithms with less run time speed to detect failures in cloud environments efficiently in terms of speed and accuracy, both of which should be further improved. Furthermore, we are also looking at predicting the task failures in cloud-based applications in real-time through providing instantaneous response and assisting in higher level of service availability.

VI. FUTURE SCOPE

The future scope of scalable AI models for predictive failure analysis in cloud-based asset management systems is vast, with several opportunities for further research and development. As industries continue to adopt more sophisticated technologies and data-driven strategies, the potential for improving asset management systems through AI and cloud computing will only expand. Several key areas for future exploration include the following:

1. Integration of Edge Computing and AI While cloud computing offers significant advantages in terms of scalability and data processing, edge computing holds immense potential for further enhancing predictive maintenance systems. Edge computing involves processing data closer to the source (i.e., at the asset level), reducing latency and improving realtime decision-making. Future research could explore hybrid architectures that combine both edge and cloud computing to ensure faster predictions and reduced reliance on cloud resources for time-sensitive applications, such as autonomous vehicles or industrial robotics.
2. Explainable AI for Predictive Maintenance One of the primary challenges with deep learning models, such as LSTM, is their "black box" nature, where it is difficult to understand how the model makes predictions. In asset management, where decisions based on AI predictions can have significant operational and financial implications, there is a need for more transparent models. The future of AI in asset management lies in developing explainable AI (XAI) techniques that provide insights into how predictions are made, offering greater transparency and trust for decision-makers. Researchers



can explore methods like SHAP (Shapley Additive Explanations) and LIME (Local Interpretable Model-agnostic Explanations) to improve the interpretability of complex models.

3. Real-Time Learning and Adaptation AI models for predictive failure analysis should be able to learn and adapt in real-time as new data streams in. Future advancements in reinforcement learning and online learning techniques could allow models to continuously adjust their parameters without requiring manual intervention. This would enable a system that is truly autonomous, capable of learning from evolving operational conditions and improving its predictions over time. These capabilities would be particularly useful in dynamic environments such as manufacturing plants or fleet management systems, where asset conditions may change frequently.

4. Multi-Model and Ensemble Approaches As the complexity of asset management systems increases, combining multiple AI models or ensemble learning techniques could lead to more robust and accurate failure predictions. Future work can explore the integration of various machine learning and deep learning algorithms to create hybrid models that combine the strengths of each approach. For instance, combining timeseries forecasting models with anomaly detection models could result in more comprehensive and precise predictions, providing a deeper understanding of potential failure points and preventive actions.

ACKNOWLEDGMENT

It gives us great pleasure in presenting the paper on “Predictive Failure Detection in cloud System using ML”. We would like to take this opportunity to thank our guide, Prof. Swati Y. Kale, Professor, Computer Department, Adsul’s technical Campus, Ahilyanagar, for giving us all the help and guidance we needed. We are grateful to her for hers kind support, and valuable suggestions were very helpful.

REFERENCES

1. Jampani, Sridhar, Aravind Ayyagari, Kodamasimham Krishna, Punit Goel, Akshun Chhapola, and Arpit Jain. (2020). Crossplatform Data Synchronization in SAP Projects. *International Journal of Research and Analytical Reviews (IJRAR)*, 7(2):875. Retrieved from www.ijrar.org.
2. Gudavalli, S., Tangudu, A., Kumar, R., Ayyagari, A., Singh, S. P., & Goel, P. (2020). AI-driven customer insight models in healthcare. *International Journal of Research and Analytical Reviews (IJRAR)*, 7(2). <https://www.ijrar.org>
3. Gudavalli, S., Ravi, V. K., Musunuri, A., Murthy, P., Goel, O., Jain, A., & Kumar, L. (2020). Cloud cost optimization techniques in data engineering. *International Journal of Research and Analytical Reviews*, 7(2), April 2020. <https://www.ijrar.org>
4. Sridhar Jampani, Aravindsundeeep Musunuri, Pranav Murthy, Om Goel, Prof. (Dr.) Arpit Jain, Dr. Lalit Kumar. (2021). Optimizing Cloud Migration for SAP-based Systems. *Iconic Research And Engineering Journals*, Volume 5 Issue 5, Pages 306- 327.
5. Gudavalli, Sunil, Vijay Bhasker Reddy Bhimanapati, Pronoy Chopra, Aravind Ayyagari, Prof. (Dr.) Punit Goel, and Prof. (Dr.) Arpit Jain. (2021). Advanced Data Engineering for MultiNode Inventory Systems. *International Journal of Computer Science and Engineering (IJCSE)*, 10(2):95–116.
6. Gudavalli, Sunil, Chandrasekhara Mokkalapati, Dr. Umababu Chinta, Niharika Singh, Om Goel, and Aravind Ayyagari. (2021). Sustainable Data Engineering Practices for Cloud Migration. *Iconic Research And Engineering Journals*, Volume 5 Issue 5, 269- 287
7. Ravi, Vamsee Krishna, Chandrasekhara Mokkalapati, Umababu Chinta, Aravind Ayyagari, Om Goel, and Akshun Chhapola. (2021). Cloud Migration Strategies for Financial Services. *International Journal of Computer Science and Engineering*, 10(2):117–142.
8. Vamsee Krishna Ravi, Abhishek Tangudu, Ravi Kumar, Dr. Priya Pandey, Aravind Ayyagari, and Prof. (Dr) Punit Goel. (2021). Real-time Analytics in Cloud-based Data Solutions. *Iconic Research And Engineering Journals*, Volume 5 Issue 5, 288-305.



9. Ravi, V. K., Jampani, S., Gudavalli, S., Goel, P. K., Chhapola, A., & Shrivastav, A. (2022). Cloud-native DevOps practices for SAP deployment. *International Journal of Research in Modern Engineering and Emerging Technology (IJRMEET)*, 10(6). ISSN: 2320- 6586.
10. Gudavalli, Sunil, Srikanthudu Avancha, Amit Mangal, S. P. Singh, Aravind Ayyagari, and A. Renuka. (2022). Predictive Analytics in Client Information Insight Projects. *International Journal of Applied Mathematics & Statistical Sciences (IJAMSS)*, 11(2):373–394.
11. Gudavalli, Sunil, Bipin Gajbhiye, Swetha Singiri, Om Goel, Arpit Jain, and Niharika Singh. (2022). Data Integration Techniques for Income Taxation Systems. *International Journal of General Engineering and Technology (IJGET)*, 11(1):191–212.
12. Gudavalli, Sunil, Aravind Ayyagari, Kodamasimham Krishna, Punit Goel, Akshun Chhapola, and Arpit Jain. (2022). Inventory Forecasting Models Using Big Data Technologies. *International Research Journal of Modernization in Engineering Technology and Science*, 4(2). <https://www.doi.org/10.56726/IRJMETS19207>

