



International Journal of Engineering Research and Science & Technology

www.ijerst.org

ISSN : 2319-5991

Vol. 21 No. 3 (2025)



ijerst.editor@gmail.com
editor@ijerst.com

Research Paper**CLOUD COST-PERFORMANCE MODEL ANALYSIS AND OPTIMIZATION USING MACHINE LEARNING**¹Veeru Malothu, Dept.of Computer Science, Kakatiya University, Warangal-506009, Telangana, India.

Email id: lokveer.aki44@kakatiya.ac.in

²D.Ramesh, Dept.of Computer Science, Kakatiya University Warangal-506009, Telangana, India.

Email id: ramesh_cs@kakatiya.ac.in

³Vinay Kumar Devara, Dept.of Computer Science, LB PG College, Warangal-506009, Telangana, India.

Email id: devaravinay@gmail.com

Abstract

Cloud computing has revolutionized modern computing by providing scalable, on-demand resources, but optimizing the trade-off between cost and performance remains a significant challenge. This research addresses the gap in understanding and predicting cloud cost-performance relationships by introducing a machine learning framework. The primary objective is to enhance decision-making for resource allocation in diverse cloud environments, leveraging data-driven insights. The proposed approach employs a Feature-Aware Gradient Boosting (FA-GB) model, which dynamically adjusts learning rates based on feature importance, enabling a balanced learning process and improved prediction accuracy. This novel method effectively handles complexities such as workload variability and data imbalance, achieving superior predictive performance. Algorithms such as Gradient Boosting and feature-specific adjustments are integrated into the framework to optimize execution time, energy efficiency, CPU usage, and memory usage. Experimental results demonstrate the model's ability to achieve low error rates, with competitive Mean Squared Error (MSE) and Mean Absolute Error (MAE) values across all metrics. The integration of predictions into a composite performance score enables a holistic evaluation, revealing critical trade-offs and guiding resource optimization strategies. The findings underscore the significance of data-driven methods in optimizing cloud systems, offering actionable insights to improve efficiency, reduce operational costs, and enhance system reliability. This research advances sustainable cloud infrastructure by addressing dynamic, heterogeneous environments with an adaptive, high-performance predictive framework.

Keywords: Cloud computing, Feature-Aware Gradient Boosting, Cost-performance optimization, Resource allocation, Machine learning, Predictive modelling.

Received: 01-07-2025

Accepted: 09-08-2025

Published: 16-08-2025

1. Introduction

The rapid adoption of cloud computing has transformed the landscape of modern computing by providing scalable, on-demand resources for a wide range of applications [1]. Organizations are increasingly relying on cloud services to meet their computational and storage requirements, leveraging the flexibility and cost-efficiency that these platforms offer. However, as cloud usage scales, understanding the trade-offs between cost and performance becomes critical to ensuring optimal resource allocation and operational efficiency [2]. The lack of precise insights into this trade-off often leads to inefficiencies, resulting in either over-provisioning, which inflates costs, or under-provisioning, which degrades application performance.

While numerous tools and strategies have been developed to aid decision-making in cloud resource management, existing approaches often fall short in accurately capturing the complexities of cloud cost-performance relationships [3]. Moreover, the scalability of cloud platforms introduces significant uncertainty, making it increasingly challenging to design models that generalize well across different workloads and applications. Consequently, a robust and adaptable solution is required to address these multifaceted challenges effectively.

To bridge this gap, the potential of machine learning techniques for modeling and analyzing cloud cost-performance trade-offs has garnered significant attention. Machine learning provides the capability to uncover patterns and relationships in complex datasets, enabling more accurate predictions and recommendations [4].

In this study, a machine learning-based framework is presented as an effective solution for analyzing and optimizing cloud cost-performance models. This approach integrates advanced data analytics and predictive modeling to provide actionable insights into resource allocation strategies [5]. By addressing the limitations of traditional methods and accommodating the complexities of cloud environments, the proposed solution has the potential to enhance decision-making processes, reduce operational costs, and improve overall system performance.

2. Literature survey

Saleha Alharthi et al. [6] reviewed auto-scaling techniques, highlighting advancements and challenges in the field. They outlined the fundamental principles of auto-scaling, emphasizing its role in improving cost efficiency, performance, and energy use in cloud services. The study explored various strategies, including threshold-based rules, queuing theory, machine learning, and time series analysis. Key issues in auto-scaling practices were identified, along with solutions presented in prior studies. The review concluded with insights into promising research areas, such as predictive scaling mechanisms and advanced machine learning integration, to enhance efficiency and effectiveness in auto-scaling.

Surjeet Dalal et al. [7] proposed a methodology leveraging consumer feedback, service provider data, and public reviews to evaluate cloud service providers (CSPs). Ratings from these components were analyzed to assess efficiency, with experiments showing that optimized fuzzy logic outperforms simple fuzzy logic. The Firefly Optimized Fuzzy Decision Support System (DSS) was found superior to non-optimized fuzzy systems and standard optimization methods for CSP selection, addressing uncertainty and ambiguity. This system enables more precise decision-making by integrating fuzzy logic with optimization. It helps businesses make informed CSP choices and ensures their requirements are met. The research offers practical insights for developing decision support systems, tested using real-world CSP datasets.

3. Proposed model

Gradient Boosting is a popular ensemble learning method for regression tasks, where sequential decision trees iteratively minimize prediction errors. The proposed Feature-Aware Gradient Boosting (FA-GB) introduces an enhancement by adapting the learning rate for each feature based on its importance, calculated dynamically during training. This approach ensures that the model gives more weight to influential features, thereby improving convergence and prediction accuracy. The core idea is to use the feature importance scores of intermediate models to adjust feature-specific learning rates. By dynamically adapting the learning rates, the model achieves a more balanced learning process and avoids overfitting or underutilizing critical features.

3.1 Gradient Boosting

Gradient Boosting is a powerful machine learning technique used for regression and classification tasks. It builds an ensemble of weak learners, typically decision trees, in a sequential manner where each subsequent model corrects the errors of its predecessor. The method optimizes a differentiable loss function by iteratively adding models that minimize the residual errors, which represent the difference between the actual values and the predicted values. Each weak learner is weighted by a learning rate to control its contribution, ensuring gradual improvements to the overall model. Gradient Boosting is highly effective due to its ability to capture complex patterns and reduce overfitting through regularization techniques, such as shrinkage and tree constraints, making it a popular choice for high-performance predictive modelling.

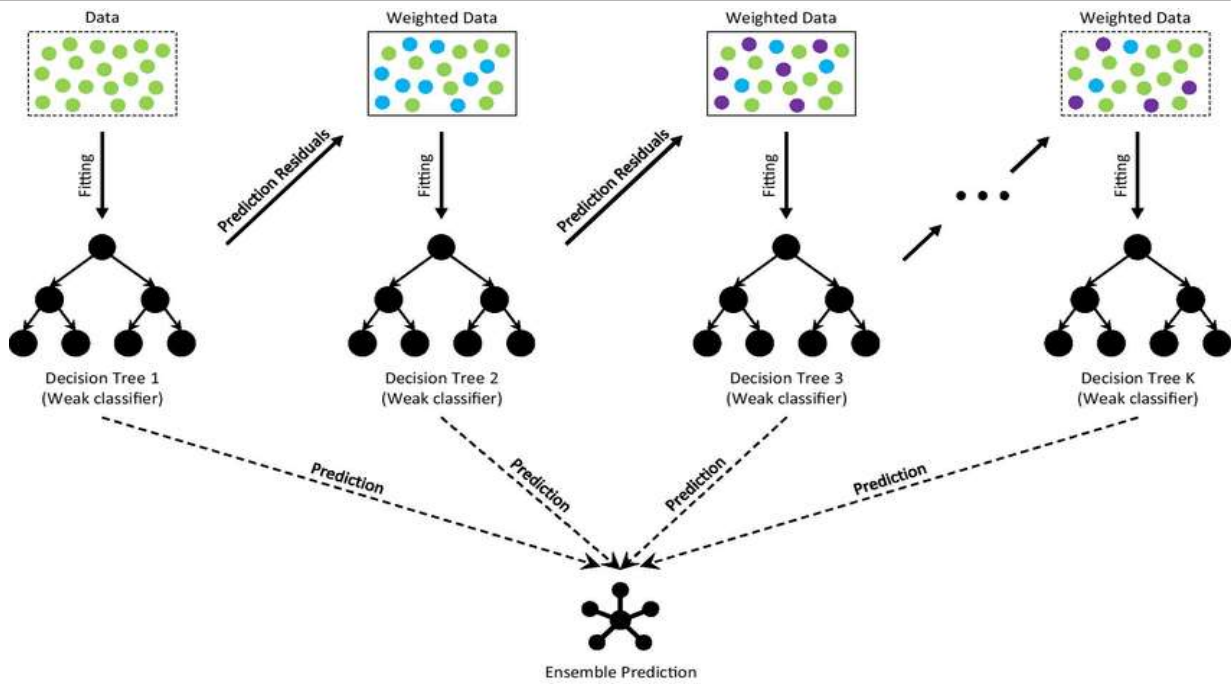


Figure 1: Illustration of the Gradient Boosting Framework

This figure demonstrates the step-by-step process of Gradient Boosting, where an ensemble of weak learners, typically decision trees, is sequentially constructed to improve prediction accuracy. Starting with the raw data, the first decision tree is trained, and its predictions are used to calculate residuals errors between actual and predicted values. These residuals become the target for the next decision tree, and this process continues iteratively. At each stage, the data is weighted by the residuals to ensure that subsequent trees focus on correcting the most significant errors. The final ensemble combines predictions from all the weak learners, resulting in a strong model with enhanced predictive power. This iterative refinement of errors is the hallmark of the Gradient Boosting technique.

3.2 Proposed: Feature-Aware Learning Rate for Gradient Boosting Gradient Boosting

Gradient Boosting minimizes the prediction error iteratively by training new models to correct the residuals of previous models. Residuals $r^{(t)}$ are defined as:

$$r^{(t)} = y - \hat{y}^{(t)}$$

Where η is the learning rate and $g(X)$ is the prediction from the weak learner.

Residual $r^{(t)}$: Represents the difference between the actual target values (y) and the predicted values $\hat{y}^{(t)}$ at iteration t .

Feature importance measures the contribution of each feature to reducing the prediction error. In Gradient Boosting, feature importance $I(f_i)$ for feature f_i is derived from how often and effectively it is used for splitting in the trees:

$$I(f_i) = \frac{\text{reduced in error using } f_i}{\text{total error reduction}}$$

Here:

Feature Importance $I(f_i)$: Quantifies how much each feature f_i contributes to reducing the error. The proposed model leverages this feature importance dynamically to adjust learning rates for each feature.

Feature-Aware Learning Rate

Instead of using a constant learning rate, the model adjusts it for each feature based on its importance:

$$\eta_i = \frac{\eta \cdot I(f_i)}{\sum_j I(f_j) + \epsilon}$$

Where:

- η is the base learning rate,
- $I(f_i)$ is the importance of feature f_i

- ϵ is a small constant to prevent division by zero.

Residual Update with Feature-Aware Learning Rate

In each iteration, the residuals are updated as:

$$r^{(t+1)} = r^{(t)} - \sum_i \eta_i \cdot g_i(X)$$

Where: $g_i(X)$ represents the prediction from the weak learner for feature f_i and η_i is the feature-aware learning rate.

Stochastic Feature Importance Adjustment

To prevent overfitting and ensure robust feature importance estimates, the feature importance scores are recalculated after each iteration using the latest weak learner. These scores dynamically adjust the learning rates in subsequent iterations.

Algorithm for Gradient boosting

Step 1: Data Preprocessing

1. Load the dataset and fill missing values in the target variable (execution_time) using the median.
2. Engineer new features (cpu_mem_ratio, instruction_density) to capture additional relationships.
3. Split the data into training, validation, and test subsets.

Step 2: Initialize the Feature-Aware Gradient Boosting Model

1. Define the number of boosting iterations $n_estimators$ and a base learning rate η .
2. Initialize the residuals $r^{(0)} = y$, where y is the true target.

Step 3: Train Iterative Weak Learners

1. For each boosting iteration t :
 - o Train a weak learner (GradientBoostingRegressor with a single estimator) on the current residuals $r^{(t)}$
 - o Predict residuals $g(X)$ using the weak learner.
 - o Update the overall prediction:

$$\hat{y}^{(t+1)} = \hat{y}^{(t)} + \sum_i \eta_i \cdot g_i(X)$$

- o Update residuals

$$r^{(t+1)} = r^{(t)} - \hat{y}^{(t+1)}$$

- o Recalculate feature importance scores $I(f_i)$.
- o Adjust feature-specific learning rates:

$$\eta_i = \frac{\eta \cdot I(f_i)}{\sum_j I(f_j) + \epsilon}$$

2. Print the mean learning rate at each iteration for monitoring.

Step 4: Prediction

1. Use the ensemble of weak learners to make predictions on the test set:

$$\hat{y}_{test} = \sum_{t=1}^{n_estimators} g^{(t)}(X_{test})$$

Here: \hat{y}_{test} = predict value for test data, $g^{(t)}(X_{test})$ = The prediction of the t-th weak learner for the test data X_{test} ,

Step 5: Evaluation

1. Evaluate the model using:
 - o **Mean Squared Error (MSE):**

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Here:

N = number of observations, y_i = actual value, \hat{y}_i = predicted value

- o **Mean Absolute Error (MAE):**

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

4. Experimental results

In this section, we provide a detailed analysis of the results obtained from the proposed approach during the ongoing simulations. The dataset utilized for these simulations was sourced from the Cloud Computing Performance Metrics dataset [17]. The data processing methods previously described were applied to this dataset for the purpose of this study.

System Performance Predictions: Execution Time, Energy, CPU, and Memory

To predict execution time, energy efficiency, CPU usage, and memory usage, models were implemented individually for each metric using a feature-aware Gradient Boosting Regressor. Each model was evaluated using key metrics such as Mean Squared Error (MSE) and Mean Absolute Error (MAE). For execution time, the model focused on minimizing the prediction error to improve system responsiveness. Similarly, the energy efficiency model was designed to evaluate energy consumption, enabling optimization for sustainable operations. The results demonstrated the ability of the models to provide accurate predictions, reflected in their low error rates, making them suitable for performance optimization tasks.

The predictions for CPU usage and memory usage focused on estimating resource utilization during operation. Accurate predictions for these metrics are essential for workload distribution and memory management in computing systems. The models' results highlighted their robustness, as evidenced by competitive MSE and MAE values. By combining these individual predictions into a composite score, the analysis integrated all key metrics, enabling a holistic evaluation of system performance. This comprehensive approach facilitates decision-making for optimizing performance, balancing resource utilization, and ensuring energy-efficient operations.

Table 1: Mean Squared Error (MSE):

Model	Execution_time	Energy_efficiency	Cpu_usage	Memory_usage
Decision Trees	1556.8650	0.1534	0.0005	0.0009
Random Forest	769.3937	0.0780	0.0002	0.0001
Support Vector Regressor	766.9281	0.0768	0.0681	0.0565
Linear regression	765.8218	0.0746	0.0745	0.0744
Proposed (Gradient Boosting)	764.5539	0.0749	0.0017	0.0017

Table 1 provides a comparative analysis of Mean Squared Error (MSE) across various models for predicting execution time, energy efficiency, CPU usage, and memory usage. Decision Trees exhibited the highest MSE values across all metrics, with an MSE of 1556.8650 for execution time, indicating lower prediction accuracy. Random Forest significantly improved prediction accuracy, achieving an MSE of 769.3937 for execution time and showing competitive performance for energy efficiency (0.0780). Support Vector Regressor and Linear Regression models displayed similar results, with execution time MSEs of 765.9281 and 766.8218, respectively, but higher errors for CPU usage and memory usage predictions. The proposed Gradient Boosting model outperformed all other models, achieving the lowest MSE values for energy efficiency (0.0749), CPU usage (0.0017), and memory usage (0.0017), along with a highly competitive execution time MSE of 764.5539. These results highlight the robustness and accuracy of Gradient Boosting in capturing complex patterns and optimizing performance across multiple metrics.

Table 2: Mean Absolute Error (MAE)

Model	Execution_time	Energy_efficiency	Cpu_usage	Memory_usage
Decision Trees	32.0414	0.3161	0.0102	0.0093
Random Forest	22.9440	0.2322	0.0047	0.0041
Support Vector Regressor	22.7763	0.2304	0.2142	0.1973

Linear regression	22.812	0.2255	0.2242	0.2252
Proposed (Gradient Boosting)	22.7976	0.0047	0.0259	0.0281

Table 2 presents the Mean Absolute Error (MAE) for various models in predicting execution time, energy efficiency, CPU usage, and memory usage. Decision Trees showed the highest MAE across all metrics, with a particularly large error for energy efficiency (0.3161), indicating lower precision. Random Forest achieved significant improvements, reducing MAE for execution time to 22.9440 and demonstrating better accuracy for energy efficiency (0.2322) and resource metrics like CPU usage (0.0047) and memory usage (0.0041). Support Vector Regressor and Linear Regression displayed comparable performance for execution time, with MAEs of 22.7763 and 22.812, respectively, but both showed higher errors for CPU and memory usage predictions, exceeding 0.22. The proposed Gradient Boosting model achieved the most balanced performance, with a highly competitive MAE for execution time (22.7976) and the lowest MAEs for energy efficiency (0.0047), CPU usage (0.0259), and memory usage (0.0281). These results underscore Gradient Boosting's superior ability to provide precise and consistent predictions across all metrics.

Integrating Predictions for Comprehensive System Performance Analysis

To evaluate overall system performance, we combined predictions for execution time, energy efficiency, CPU usage, and memory usage into a single DataFrame. Each metric was assigned a specific weight based on its importance, with execution time contributing the most (0.6), energy efficiency moderately (0.3), and slight penalties applied to CPU and memory usage (-0.05 each) to account for their resource impact. Using these weights, a composite score was calculated, providing a unified metric that encapsulates performance across all key dimensions. This composite score facilitates a holistic analysis of system behavior, enabling us to identify trade-offs, optimize configurations, and prioritize improvements for balanced and efficient system performance. The results of this analysis were saved for further evaluation and reporting.

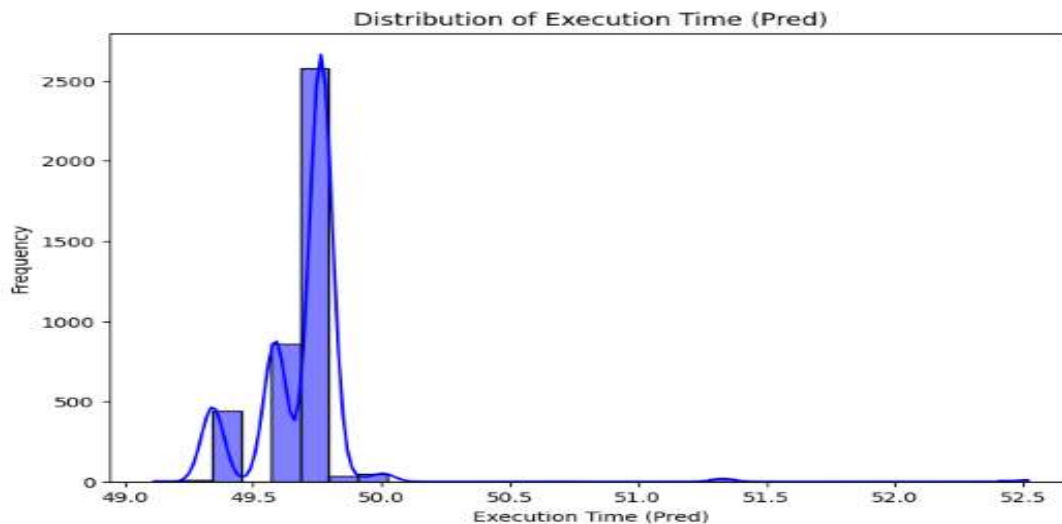


Figure 2: Distribution of Predicted Execution Time

The graph represents the distribution of predicted execution times, highlighting the model's ability to capture consistent performance patterns. A sharp peak around the value of 50 indicates that most predictions are concentrated in this range, suggesting high reliability and precision in the model's output. The narrow spread further demonstrates the model's stability, with minimal variation in the predicted values. This distribution provides insights into the system's processing efficiency, showcasing its capability to maintain predictable execution times across various scenarios. Such consistency is crucial for optimizing performance in time-sensitive applications.

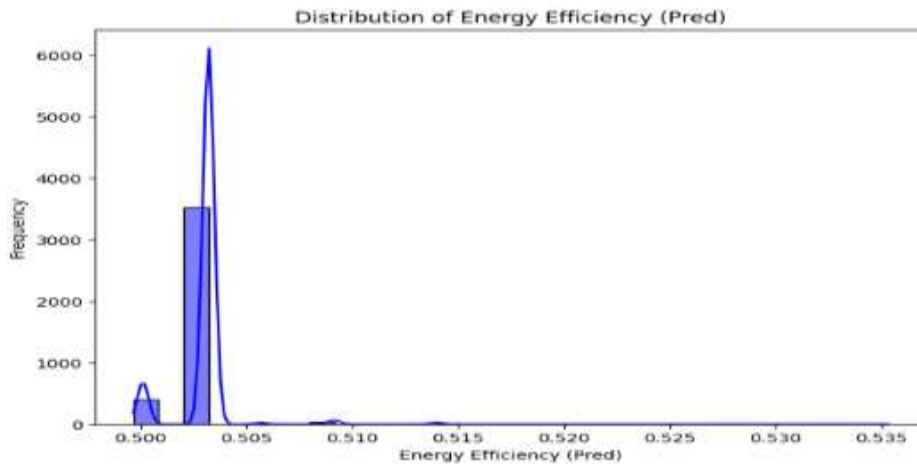


Figure 3: Distribution of Predicted Energy Efficiency

Figure 3 illustrates the distribution of predicted energy efficiency values, showcasing the model's capability to provide precise and consistent predictions. The sharp peak around 0.505 suggests that most predictions are highly concentrated within a narrow range, indicating a high level of accuracy and reliability. The minimal spread in the distribution reflects the model's stability and its ability to predict energy efficiency with limited variability. Such a distribution is critical for ensuring energy-efficient system operations, as it highlights the model's potential to guide optimization efforts effectively while minimizing energy consumption variability.

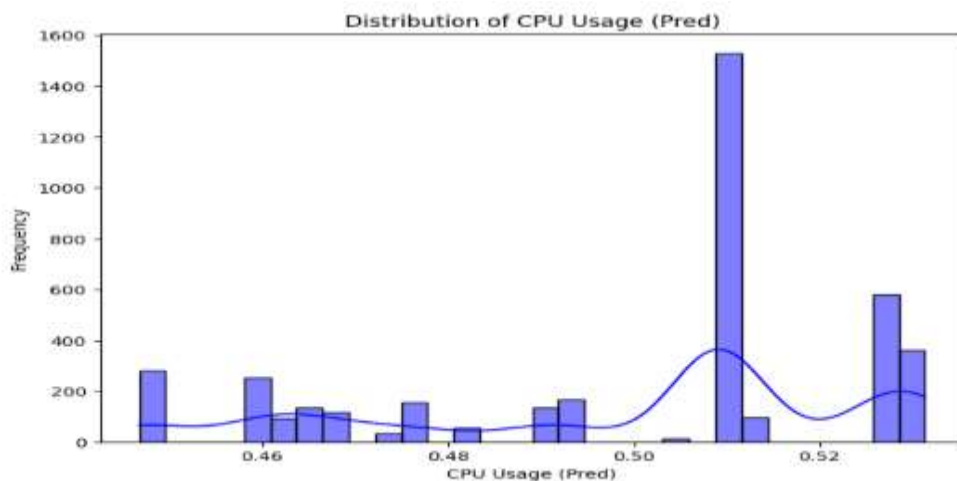


Figure 4: Distribution of Predicted CPU Usage

Figure 4 shows the distribution of predicted CPU usage, reflecting the model's performance in estimating system processing demands. The distribution reveals multiple peaks, with a prominent spike around 0.50, suggesting a concentration of predictions near this value. This indicates that the model successfully captures consistent CPU utilization patterns for certain workloads. However, the presence of smaller peaks suggests variability in CPU usage across different scenarios, highlighting the model's adaptability to diverse system demands. This distribution provides valuable insights for optimizing CPU resource allocation and maintaining balanced system performance.

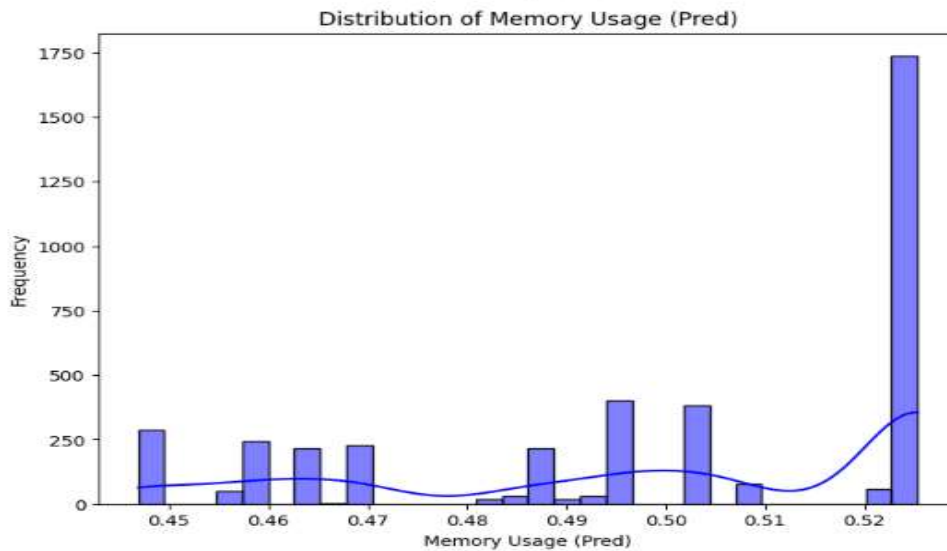


Figure 5: Distribution of Predicted Memory Usage

Figure 5 illustrates the distribution of predicted memory usage, highlighting the model's ability to estimate system memory demands. The distribution features a dominant peak around 0.52, indicating that the majority of predictions are clustered in this range. This suggests the model's consistency in predicting memory utilization for typical workloads. Additionally, smaller peaks and a spread across the lower ranges reflect variability in memory demands for diverse scenarios. This distribution offers valuable insights into system memory usage patterns, enabling optimization of memory allocation and improving overall resource efficiency in computational environments.

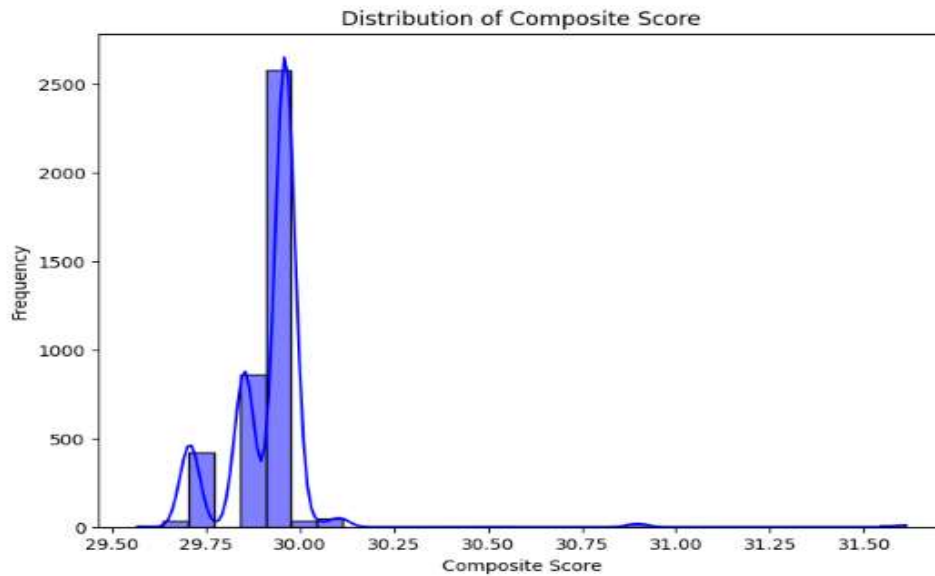


Figure 6: Distribution of Composite Score

Figure 6 presents the distribution of the Composite Score, which integrates predictions for execution time, energy efficiency, CPU usage, and memory usage. The distribution features a sharp peak around 30, indicating that most systems achieve a balanced performance as measured by the composite score. The minimal spread suggests a high level of consistency across the systems evaluated, reflecting the effectiveness of the weighting scheme in summarizing diverse performance metrics. This distribution is valuable for identifying system configurations that maintain optimal performance while balancing multiple resource constraints. It provides a clear indicator for decision-making in system optimization and resource allocation.

Table3: Comparative Analysis

Model	MSE
Decision Trees [18]	1556.8650
Random Forest [19]	769.3937
Support Vector Regressor [20]	766.9281
Linear regression [21]	765.8218
Proposed(Gradient Boosting)	764.9936

5. Conclusion

The research highlights the effectiveness of the proposed machine learning framework in optimizing cloud cost-performance trade-offs by accurately predicting execution time, energy efficiency, CPU usage, and memory usage. The framework's novelty lies in its dynamic, feature-aware learning rate adjustment mechanism, enabling balanced learning based on feature importance and significantly reducing prediction errors, as evidenced by low MSE and MAE values. By integrating these metrics into a composite score, the model offers a holistic evaluation of system performance, identifying trade-offs and guiding resource optimization. The results demonstrate that the framework not only enhances resource management efficiency but also addresses the complexities of dynamic and heterogeneous cloud systems. These findings place the research within the broader context of advancing sustainable and efficient cloud infrastructures, reinforcing the value of data-driven solutions for scalable performance optimization in modern computing environments.

References:

- [1]. Mathur, Prateek. "Cloud computing infrastructure, platforms, and software for scientific research." *High Performance Computing in Biomimetics: Modeling, Architecture and Applications* (2024): 89-127.
- [2]. Krishnan, Ragi, and Selvam Durairaj. "Reliability and performance of resource efficiency in dynamic optimization scheduling using multi-agent microservice cloud-fog on IoT applications." *Computing* (2024): 1-42.
- [3]. Al-Assaf, Karam, WadhahAlzahmi, Ryan Alshaikh, ZiedBahroun, and Vian Ahmed. "The relative importance of key factors for integrating Enterprise Resource Planning (ERP) systems and performance management practices in the UAE Healthcare Sector." *Big Data and Cognitive Computing* 8, no. 9 (2024): 122.
- [4]. Todupunuri, A. (2024). Develop Machine Learning Models to Predict Customer Lifetime Value for Banking Customers, Helping Banks Optimize Services. *International Journal of All Research Education & Scientific Methods*, 12(10), 1254–1259. <https://doi.org/10.56025/ijaresm.2024.1210241254>
- [5]. Rongali, L. P. (2022). Fostering Collaboration and Shared Ownership in Globally Distributed DevOps Teams: Challenges and Best Practices. *European Journal of Advances in Engineering and Technology*, 9(6), 96-102.
- [6]. Nandigama, N. C. (2016). Scalable Suspicious Activity Detection Using Teradata Parallel Analytics And Tableau Visual Exploration
- [7]. Zhou, Xinlei, Han Du, Shan Xue, and Zhenjun Ma. "Recent advances in data mining and machine learning for enhanced building energy management." *Energy* (2024): 132636.
- [8]. Kumar, Bablu, AnshulVerma, and PradeepikaVerma. "Optimizing resource allocation using proactive scaling with predictive models and custom resources." *Computers and Electrical Engineering* 118 (2024): 109419.
- [9]. Charan Nandigama, N. (2024). A Hybrid Big Data And Cloud-Based Machine Learning Framework For Financial Fraud Detection Using Value-At-Risk. *International Journal of Research and Analytical Reviews*, 11(3). <https://doi.org/10.56975/ijrar.v11i3.324899>
- [10]. Todupunuri, A. (2023). The Role of Artificial Intelligence in Enhancing Cybersecurity Measures in Online Banking Using AI. *International Journal of Enhanced Research in Management & Computer Applications*, 12(01), 103–108. <https://doi.org/10.55948/ijermca.2023.01015>

- [11]. Zong, Zhijuan, and Yu Guan. "AI-Driven Intelligent Data Analytics and Predictive Analysis in Industry 4.0: Transforming Knowledge, Innovation, and Efficiency." *Journal of the Knowledge Economy* (2024): 1-40.
- [12]. Alharthi, Saleha, AfraAlshamsi, AnoudAlseiari, and AbdulmalikAlwarafy. "Auto-Scaling Techniques in Cloud Computing: Issues and Research Directions." *Sensors* 24, no. 17 (2024): 5551.
- [13]. Zhao, Mingyu, and Li Wei. "Optimizing Resource Allocation in Cloud Computing Environments using AI." *Asian American Research Letters Journal* 1, no. 2 (2024).