

Using Machine Learning for Analyzing Performance Metrics in Kubernetes

Diana Kutsa

Bachelor of Management in Ternopil National Economic University
Crystal Lake IL, USA

Abstract:

This article examines modern approaches to analyzing the performance metrics of containerized applications in Kubernetes using machine learning methods. The key classes of algorithms (regression models, classification and anomaly detection methods, reinforcement learning) are analyzed, along with their advantages and limitations. The architectural principles of integrating machine learning into the Kubernetes environment are presented, and typical application scenarios are proposed, including load forecasting, automatic scaling, and failure and anomaly detection. A review of existing frameworks (Kubeflow, Katib, etc.) and monitoring tools (Prometheus, Grafana) is conducted, highlighting their role in facilitating the implementation of MLOps pipelines within Kubernetes. Based on the findings, recommendations are formulated for selecting and combining machine learning methods depending on workload profiles. Additionally, future research directions are outlined, including deeper integration of reinforcement learning and predictive models to enhance the resilience of microservice architectures. This study is relevant for researchers and practitioners at the intersection of distributed computing, artificial intelligence, and DevOps who seek to integrate machine learning for automated monitoring, diagnostics, and performance optimization of Kubernetes clusters in scalable cloud infrastructures.

Keywords: Container Technologies, Kubernetes, Machine Learning, Performance Metrics, Container Orchestration, Resource Forecasting, Anomaly Detection, Reinforcement Learning, MLOps.

Introduction

Modern cloud technologies increasingly utilize containers and orchestration systems (Docker, Kubernetes, Nomad, etc.) to simplify the deployment and management of microservice applications [12]. As these systems scale and business applications become more complex, performance monitoring becomes critical, ranging from accurate resource load forecasting (CPU, memory, network resources) to identifying bottlenecks during peak loads [5, 10, 11]. Due to the dynamic nature of workloads and the high variability of container cluster configurations, traditional approaches (manual tuning, static scaling) are becoming less effective. In this context, the use of machine learning (ML) tools is becoming more relevant, as they not only enable real-time analysis of large streams of performance metrics but also allow for the development of predictive models for auto-tuning and performance optimization [1, 7].

Ahmad I. et al. [1] provide a comprehensive review of existing container scheduling techniques, emphasizing the challenge of optimal computational resource allocation under dynamically changing loads. Bandari V. [2] explores the application of artificial intelligence in automated container orchestration, predictive maintenance, security, and resource optimization. Additionally, Chiang R. C. [3] proposes a container placement strategy that accounts for resource competition in Docker Swarm using machine learning-based clustering algorithms. Zhong Z. et al. [12] contribute to the systematization of research by offering a taxonomy of machine learning-based container orchestration and outlining promising directions for future developments. Kusepova L. T. et al. [13] further analyze containerized application deployment processes using machine learning, aiming to assess the efficiency of existing methods and identify opportunities for optimization. The study by Yadav M. P., Rohit,

and Yadav D. K. [9] focuses on maintaining the resilience of containerized systems through predictive analysis, addressing the long-term stability of distributed computing environments.

Dartois J. E. et al. [5] investigate the application of machine learning algorithms for modeling SSD input/output performance in containerized virtualization environments. Ye K. and Ji Y. [10] focus on optimizing performance for big data processing applications in Docker containers, paying particular attention to fine-tuning key system metrics. Ye K. et al. [11] propose methods for fault injection and detection to enhance the fault tolerance of artificial intelligence in containerized cloud environments. Additional contributions to this field are made by Imdoukh M., Ahmad I., and Alfailakawi M. G. [7], who explore the development of an automatic scaling system for containerized applications.

Chowdary M. N. et al. [4] examine strategies for accelerating the deployment of machine learning models using MLOps practices, with the goal of optimizing integration, testing, and model delivery processes in production environments. Doukha R. et al. [6] address the deployment challenges of containerized deep learning models in cloud environments, aiming to eliminate bottlenecks associated with the computational intensity of such applications. Singh P. [8] provides a practical guide on deploying machine learning models in industrial production settings, emphasizing container technologies.

Thus, a comprehensive approach to analyzing Kubernetes performance metrics using machine learning remains insufficiently explored, particularly in the context of integrating resource forecasting, anomaly detection, and dynamic infrastructure optimization within a unified framework. Furthermore, there is no established methodology for matching various ML algorithms (regression, clustering, reinforcement learning) to specific metrics and usage scenarios (e.g., business-critical services versus high-load but less critical microservices).

Based on this identified gap, this study aims to:

- Assess the potential and limitations of key machine learning algorithms for analyzing Kubernetes performance metrics.

- Develop an approach for integrating predictive and diagnostic models into the orchestration framework (via HPA, Kubeflow, etc.).

- Evaluate the effectiveness of the proposed methods under real or near-real test loads, demonstrating resource cost reduction and increased cluster stability.

The scientific novelty lies in proposing a comprehensive application of multiple ML algorithm classes (regression, clustering, RL) to manage diverse Kubernetes performance metrics.

The research hypothesis posits that the use of specialized machine learning models for analyzing key Kubernetes performance metrics will not only improve forecasting accuracy and prevent cluster overloads but also reduce overall response time to anomalies and failures.

Research Results

Regression methods are most commonly used for predicting key resource utilization metrics such as CPU, memory, and network bandwidth [10, 11].

1. Linear regression is a simple algorithm that provides a clear understanding of the relationship between factors (pod configuration, traffic intensity) and target metrics (response time, CPU usage). However, its accuracy decreases under highly dynamic workloads, as it fails to account for nonlinear and complex dependencies.

2. Support Vector Regression (SVR) is widely applied when data exhibits complex dependencies and sample size is relatively small. It models nonlinear relationships between resource utilization in Kubernetes and external factors (workload types, temporal patterns) more effectively [5, 11].

3. Random Forest and Gradient Boosting are commonly used for analyzing large volumes of telemetry data, allowing the consideration of multiple features (number of replicas, image versions, resource limits). For high interpretability requirements (e.g., in critical systems), additional interpretation methods such as Shapley values are necessary [7].

4. Time series models, including ARIMA, LSTM, and Prophet, are applicable for proactive scaling when predicting metric behavior several steps ahead is essential [3]. LSTM (Long Short-Term Memory) effectively

handles complex seasonal patterns and sudden load spikes, which are characteristic of microservice architectures.

5. Classification methods (e.g., logistic regression, Naïve Bayes, neural networks) and clustering techniques (e.g., k-means, DBSCAN) are used for anomaly detection and workload segmentation based on resource consumption intensity [6, 12].

6. Workload clustering segments applications and containers into clusters based on similar CPU/RAM consumption profiles, simplifying pod placement planning within a cluster [12]. K-means identifies groups of pods

with similar peak profiles, while DBSCAN helps ignore noisy data points with atypical behavior.

For real-time adaptation requirements, reinforcement learning (RL) and related decision-making methods are recommended [2, 9]. The model (agent) selects actions (pod scaling, resource reallocation) based on the environment's state (performance metrics) and receives rewards for achieving target objectives (latency minimization, SLA compliance) [4].

To systematize knowledge on the applicability of various ML approaches in Kubernetes, Table 1 is presented.

Table 1. Comparative overview of machine learning methods for analyzing kubernetes performance metrics [2-3, 5-7, 10-12]

ML Method	Analysis Tasks	Advantages	Limitations
Regression (including Linear, SVR, Random Forest)	CPU, RAM, and network metric forecasting; workload planning	+ Provides a good approximation of dependencies + Simple interpretation (linear models)	- May lack accuracy for complex nonlinear dependencies - Requires careful hyperparameter tuning
Clustering (k-means, DBSCAN)	Grouping pods with similar resource consumption patterns	+ Simplifies workload planning + Enables workload segment visualization	- Does not provide direct time-based or volume-based predictions - Requires feature selection and cluster number definition (k-means)
Anomaly Detection (One-Class SVM, Isolation Forest)	Detecting sudden deviations from normal behavior	+ Identifies hidden issues + Reduces failure response time	- Possible false positives - Requires balanced training on normal data
Reinforcement Learning (RL)	Automatic scaling, dynamic resource allocation	+ Adapts in real-time + Self-learning across different scenarios	- Complex implementation and training - Requires environment modeling and safe testing
Time Series Models (ARIMA, LSTM)	Short- and long-term load trend forecasting	+ Accounts for seasonality and recurring patterns + Works well with large historical datasets	- Potential forecasting lag - Requires periodic model retraining

The integration of machine learning methods into the Kubernetes orchestration pipeline relies on several tools, as illustrated in Figure 1.

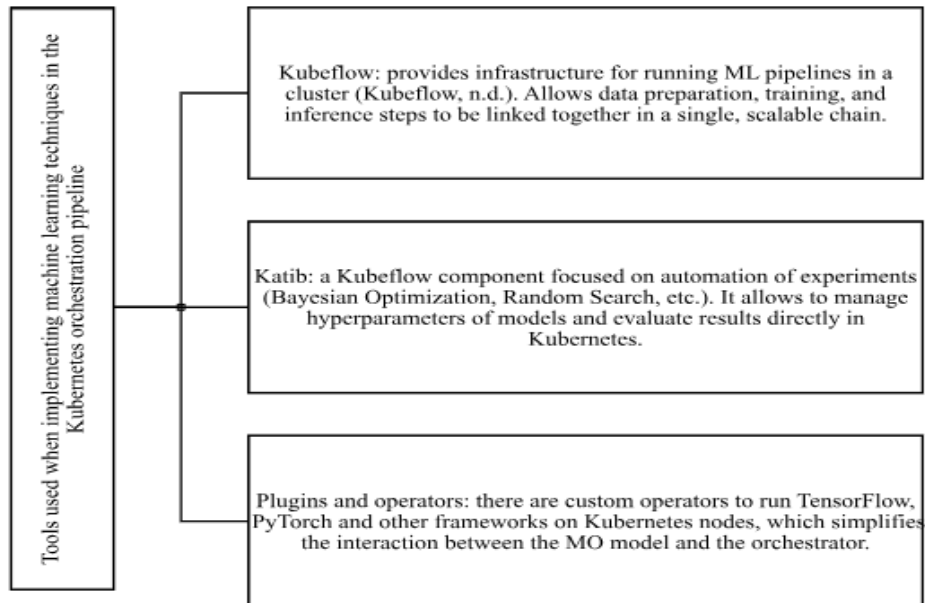


Fig.1. Tools Used in the Implementation of Machine Learning Methods in the Kubernetes Orchestration Pipeline [2, 8].

Thus, the machine learning methods applied for analyzing Kubernetes metrics cover a wide range of tasks, from predicting future states and failure prevention to adaptive resource management. The key criterion for selecting a specific algorithm is the nature of the metrics (linear or highly nonlinear dependencies), the required level of accuracy, and the system's ability to respond rapidly to external changes. The following sections will explore the practical aspects of implementing these methods in real Kubernetes infrastructures and present experimental results comparing their performance.

Using machine learning for analyzing performance metrics in Kubernetes

This section examines the practical aspects of implementing machine learning (ML) algorithms in the Kubernetes orchestration framework to analyze performance metrics and optimize cluster operations. Based on the previously reviewed methods and sources [7, 10, 12], a general approach to integrating ML into Kubernetes monitoring, scaling, and auto-configuration processes is described.

1. **Multilevel Data and Computation Organization.** Modern containerized environments utilize multiple sources of metrics,

including Prometheus, cAdvisor, and kube-state-metrics [12]. For ML applications, it is essential to centralize data collection while ensuring high availability and consistency [3]. Metric storage solutions such as InfluxDB and TimescaleDB enable scalable storage and analysis of historical data, which is critical for training models [11].

2. **Computational Resources for ML.** Deploying ML models within a Kubernetes cluster (via Kubeflow, Katib) helps reduce latency in data transmission [2, 13]. For large datasets, GPU/TPU-powered nodes (Nomad, Prefect) or specialized distributed training frameworks (Spark on Kubernetes, Ray on Kubernetes) are recommended [4].

3. **Continuous MLOps Cycle.** This process includes data collection, model training, accuracy evaluation, and production deployment. Integration with CI/CD pipelines ensures automatic model updates when new workload patterns emerge [8]. Maintaining model versions (model versioning) and conducting regular A/B testing is essential to compare ML-based scaling with baseline approaches (e.g., standard Horizontal Pod Autoscaler, HPA) [12].

Table 2 summarizes the stages of implementing ML-based performance metric analysis in Kubernetes, including their objectives and key actions.

Table 2. The stages of implementing machine learning for analyzing Kubernetes performance metrics [2, 7, 12, 13]

Stage	Objective	Key Actions
1. Data Collection and Storage	Establish a unified database for metrics and logs	- Integrate Prometheus, cAdvisor - Configure long-term storage (InfluxDB, TimescaleDB) - Define a storage schema
2. Data Preprocessing	Remove noise, handle missing values, normalize features	- Clean and aggregate time-series data - Remove anomalous data points - Extract key features
3. Model Selection and Tuning	Determine the most suitable ML method for the scenario	- Compare regression, clustering, and RL approaches - Optimize hyperparameters (Katib, Grid Search)
4. Training and Validation	Obtain a trained model with high accuracy	- Split data into train/test sets - Perform K-fold cross-validation - Evaluate using MAE, RMSE, F1-score, etc.
5. Deployment (Inference)	Automate model inference in Kubernetes	- Implement microservice inference (Flask, FastAPI) - Containerize the model (Docker image) - Set up CI/CD (Jenkins, GitLab)
6. Monitoring and Iterations	Continuously track model performance and adjust when failures occur	- Log prediction results - Automatically roll back to baseline scaling policy in case of failures

Thus, achieving sustainable success requires not only training the model but also ensuring a complete cycle from data collection to continuous improvement.

Practical examples and implementation experience:

- Auto-scaling based on LSTM. Imdoukh M., Ahmad I., Alfaiakawi M. G. [7] demonstrated that LSTM neural networks provide accurate workload predictions, allowing timely deployment of additional pods and preventing performance degradation.

- Anomaly detection using One-Class SVM. Doukha R. et al. [6] successfully detected anomalous memory usage patterns in containerized deep learning applications within a Kubernetes experimental environment.

Despite the high efficiency of machine learning in analyzing Kubernetes performance metrics, several limitations exist:

1. Dependence on data quality. If the collected metrics are not representative or contain excessive noise, even the most advanced

algorithms may fail to produce reliable results [10].

2. High computational resource requirements. Deep neural networks (LSTM, CNN) and reinforcement learning (RL) algorithms may require specialized GPU nodes.

3. Debugging complexity. Automated solutions must include fallback mechanisms to revert to baseline Kubernetes settings to prevent uncontrolled failures [1].

Future developments are expected to focus on tighter integration of MLOps tools (Kubeflow, Katib, MLflow) with Kubernetes monitoring systems, as well as the advancement of hybrid approaches that combine predictive models with reinforcement learning methods [12, 13]. This will enable the development of more accurate and adaptive strategies for scaling and resource planning.

Thus, the application of machine learning in Kubernetes extends beyond simple metric analysis. It enables proactive resource optimization, enhances fault tolerance, and allows for flexible responses to unexpected

workload changes. The next section will present conclusions and recommendations on a comprehensive approach to container orchestration using machine learning, along with prospects for further research.

Conclusion

The conducted study demonstrates that the combination of container technologies and machine learning algorithms can significantly improve performance management in Kubernetes. The analysis of metrics using regression models, anomaly detection methods, and reinforcement learning enables rapid responses to dynamic workload changes in the cluster, forecasting application resource consumption, and proactive scaling.

The presented examples and literature review confirm that tasks such as container placement planning, SLA monitoring, and overload prevention are handled more effectively through a comprehensive approach, where historical data analysis is combined with real-time configuration adaptation. Developed MLOps tools (Kubeflow, Katib) and orchestration mechanisms (Prometheus, Horizontal Pod Autoscaler, Vertical Pod Autoscaler) enable the automation of the full model training and inference cycle, as well as seamless integration of these models into standard DevOps processes.

The scientific novelty of this study lies in the combination of multiple ML methods (regression, clustering, RL) within a unified Kubernetes performance management framework, producing a synergistic effect compared to traditional autoscalers. From a practical perspective, this approach reduces system downtime risks, enhances service quality, and optimizes cloud infrastructure costs.

Future research directions include expanding the range of analyzed metrics (network latency, I/O operations, log files), developing hybrid methods (combining reinforcement learning with time-series forecasting), and improving MLOps practices to ensure continuous model improvement and rapid adaptation to changing cluster operating conditions.

References

- 1 Ahmad I. et al. Container scheduling techniques: A survey and assessment //Journal of King Saud University-Computer and Information Sciences. – 2022. – Vol. 34 (7). – pp. 3934-3947.
- 2 Bandari V. A comprehensive review of AI applications in Automated Container Orchestration, Predictive maintenance, security and compliance, resource optimization, and continuous Deployment and Testing //International Journal of Intelligent Automation and Computing. – 2021. – Vol. 4 (1). – pp. 1-19.
- 3 Chiang R. C. Contention-aware container placement strategy for docker swarm with machine learning based clustering algorithms //Cluster Computing. – 2023. – Vol. 26 (1). – pp. 13-23.
- 4 Chowdary M. N. et al. Accelerating the Machine Learning Model Deployment using MLOps //Journal of Physics: Conference Series. – IOP Publishing, 2022. – Vol. 2327 (1). – pp. 1-10.
- 5 Dartois J. E. et al. Investigating machine learning algorithms for modeling ssd i/o performance for container-based virtualization //IEEE transactions on cloud computing. – 2019. – Vol. 9 (3). – pp. 1103-1116.
- 6 Doukha R. et al. Deployment of containerized deep learning applications in the cloud //2020 5th International Conference on Cloud Computing and Artificial Intelligence: Technologies and Applications (CloudTech). – IEEE. - 2020. – pp. 1-6.
- 7 Imdoukh M., Ahmad I., Alfailakawi M. G. Machine learning-based auto-scaling for containerized applications //Neural Computing and Applications. – 2020. – Vol. 32 (13). – pp. 9745-9760.
- 8 Singh P. Deploy machine learning models to production //Cham, Switzerland: Springer. – 2021. – pp.7-29.
- 9 Yadav M. P., Rohit, Yadav D. K. Maintaining container sustainability through machine learning //Cluster Computing. – 2021. – Vol. 24 (4). – pp. 3725-3750.
- 10 Ye K., Ji Y. Performance tuning and modeling for big data applications in docker containers //2017 International Conference on

Networking, Architecture, and Storage (NAS). – IEEE. - 2017. – pp. 1-6.

11 Ye K. et al. Fault injection and detection for artificial intelligence applications in container-based clouds //Cloud Computing–CLOUD 2018: 11th International Conference, Held as Part of the Services Conference Federation, SCF 2018, Seattle, WA, USA, June 25–30, 2018, Proceedings 11. – Springer International Publishing. - 2018. – pp. 112-127.

12 Zhong Z. et al. Machine learning-based orchestration of containers: A taxonomy and future directions //ACM Computing Surveys (CSUR). – 2022. – Vol. 54 (10). – pp. 1-35.

13 Kusepova L. T. et al. Deploying Container applications using machine learning: overview and analysis //Bulletin of Kazatk. – 2024. – Vol. 2 (23). – pp. 94-102.