

Research Article

# Optimization in Machine Learning for Application in High Energy Physics: Processing and Anomaly Detection

Sujata Nema<sup>1</sup> , Ramkumar Nagarch<sup>1</sup> , Parmeshwar Dayal Lodhi<sup>1</sup> ,  
Shailendra Jain<sup>2, \*</sup> 

<sup>1</sup>Department of Physics, Maharaja Chhatrasal Bundelkhand University, Chhatarpur, India

<sup>2</sup>Department of Physics, Eklavya University, Damoh, India

## Abstract

High-energy physics (HEP) experiments generate extraordinarily large and complex datasets, posing significant challenges for real-time data analysis, event reconstruction, and reliable anomaly detection. Traditional analytical techniques often struggle to scale efficiently or fully exploit the rich structure of these data. In this context, machine learning (ML) has emerged as a transformative paradigm, offering powerful tools to enhance computational efficiency, precision, and adaptability in HEP data processing pipelines. This review provides a comprehensive overview of the integration of ML techniques in HEP, with a particular focus on their role in optimizing data analysis workflows and improving experimental performance. We examine a broad spectrum of ML approaches, including supervised and unsupervised learning methods, deep learning architectures, and ensemble models, highlighting their applications in tasks such as signal–background discrimination, feature extraction, noise reduction, and anomaly detection. Special attention is given to advanced algorithms designed for real-time data processing, which are critical for trigger systems and online event selection in modern collider experiments. The effectiveness of these methods is evaluated in the context of large-scale HEP datasets, demonstrating strong performance with metrics including an accuracy of 0.9421, sensitivity of 0.9314, specificity of 0.9507, precision of 0.9458, an F1-score of 0.9386, and an area under the ROC curve (AUC) of 0.9723. By critically analyzing current ML models and their integration into established HEP data analysis frameworks, this review identifies recent advancements, ongoing challenges related to model interpretability, scalability, and robustness, and promising directions for future research. The findings underscore the pivotal role of ML in advancing data-driven discoveries in HEP and support the development of more accurate, efficient, and scalable experimental analyses.

## Keywords

High-Energy Physics (HEP), Machine Learning (ML), Data Analysis, Anomaly Detection, Real-Time Processing, Deep Learning, Data Quality Enhancement

\*Correspondence: Shailendra Jain (shailendra.jain@eklavyauniversity.ac.in)

Received: 14 February 2026; Accepted: 22 April 2026; Published: 16 May 2026



Copyright: © The Author(s), 2026. Published by Science Publishing Group. This is an **Open Access** article, distributed under the terms of the Creative Commons Attribution 4.0 License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited.

## 1. Introduction

Optimization plays a central role in machine learning (ML) models applied to high-energy physics, as it directly influences model accuracy, generalization capability, and computational efficiency. In HEP, ML techniques are widely used for tasks such as particle identification, event classification, background suppression, detector calibration, and anomaly detection. These tasks often involve high-dimensional data, large event samples, and complex physical constraints, making optimization a critical challenge [1].

The optimization process typically focuses on minimizing a loss function that quantifies the discrepancy between predicted and observed physical quantities. Gradient-based optimization algorithms, including stochastic gradient descent (SGD), Adam, RMSProp, and AdaGrad, are commonly employed due to their scalability and efficiency when handling large datasets generated by particle detectors. These methods iteratively update model parameters to achieve faster convergence and improved stability during training.

Hyper parameter optimization is another crucial aspect, as model performance in HEP is highly sensitive to parameters such as learning rate, batch size, network depth, regularization strength, and kernel configurations. Techniques such as grid search, random search, Bayesian optimization, genetic algorithms, and particle swarm optimization are frequently adopted to identify optimal hyper parameter settings while balancing computational cost.

Additionally, optimization in HEP-oriented ML models must account for domain-specific constraints, including physical interpretability, robustness to systematic uncertainties, and class imbalance between signal and background events [2]. Regularization techniques, early stopping, and physics-informed loss functions are often incorporated to prevent overfitting and ensure physically meaningful predictions.

Overall, effective optimization strategies significantly enhance the reliability and applicability of machine learning models in high-energy physics, enabling precise data analysis and supporting discoveries in fundamental particle interactions.

## 2. Algorithms Used for Optimization in Machine Learning for High-energy Physics

Optimization algorithms are essential for training machine learning models and tuning hyper parameters in high-energy physics applications. The most commonly used algorithms are categorized as follows:

### *Gradient-Based Optimization Algorithm*

These algorithms are primarily used for training neural networks by minimizing the loss function.

- 1) Stochastic Gradient Descent (SGD)- Iteratively updates model parameters using mini-batches of data [3]. It is

widely used due to its simplicity and scalability for large HEP datasets.

- 2) SGD with Momentum- Accelerates convergence by incorporating past gradients, helping to escape shallow local minima.
- 3) Adam (Adaptive Moment Estimation)- Combines momentum and adaptive learning rates, making it highly effective for deep learning models in particle classification and event reconstruction.
- 4) RMSProp- Adjusts learning rates based on a moving average of squared gradients, improving stability during training.
- 5) AdaGrad- Adapts learning rates individually for each parameter, useful for sparse features in detector data.

## 3. Hyper Parameter Optimization Algorithms

These algorithms optimize model configuration parameters rather than model weights

- 1) Grid Search- Exhaustively searches predefined hyper parameter spaces; effective but computationally expensive.
- 2) Random Search- Randomly samples hyper parameters, often achieving comparable results with lower computational cost.
- 3) Bayesian Optimization- Uses probabilistic models to efficiently explore the hyper parameter space and is widely adopted in HEP workflows [3].
- 4) Genetic Algorithms (GA)- Evolutionary optimization inspired by natural selection, suitable for complex, non-convex search spaces.
- 5) Particle Swarm Optimization (PSO)- Models collective behavior of particles to find optimal solutions, frequently used in feature selection and clustering.

## 4. Clustering and Unsupervised Optimization Algorithms

Clustering Optimization Used for event grouping and anomaly detection, clustering algorithms aim to partition a dataset into groups such that data points within a cluster are more similar to each other than to those in other clusters. While Unsupervised Optimization Algorithms aim to learn representations or structures by optimizing objective functions without labeled outputs.

### 4.1. K-Means Clustering

K-Means clustering is an unsupervised machine learning technique used to partition high-dimensional data into groups

based on similarity. In high-energy physics, it is widely applied for event clustering, particle track grouping, and initial anomaly detection when labeled data are scarce. The algorithm divides the dataset into  $K$  clusters represented by centroids, iteratively assigning data points to the nearest centroid and updating centroid positions until convergence [4]. Owing to its simplicity and computational efficiency, K-Means is well suited for handling large detector-level datasets and real-time pre-processing. However, its assumption of spherical clusters may limit performance on complex HEP data distributions, which can be mitigated through normalization and ensemble-based approaches.

## 4.2. Gaussian Mixture Models (GMM)

Gaussian Mixture Models are probabilistic, unsupervised learning techniques used to model complex data distributions as a weighted combination of multiple Gaussian components. In high-energy physics, GMMs are commonly applied for event classification, particle identification, and background-signal separation, particularly when data exhibit overlapping or non-spherical structures [5]. Unlike K-Means, GMM assigns soft probabilistic memberships to data points, allowing greater flexibility in modeling uncertainties inherent in HEP datasets. This capability makes GMMs well suited for detector-level data analysis, where measurement noise and statistical fluctuations are significant.

## 4.3. Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

DBSCAN is an unsupervised, density-based clustering algorithm designed to identify clusters of arbitrary shape while effectively handling noise and outliers. In high-energy physics applications, DBSCAN is particularly useful for particle track reconstruction, event topology analysis, and anomaly detection, where background noise and irregular data distributions are common. The algorithm groups data points based on local density using two parameters: the neighborhood radius ( $\epsilon$ ) and the minimum number of points (MinPts). Points in high-density regions form clusters, while sparsely distributed points are treated as noise. Its ability to detect non-spherical clusters and isolate outliers makes DBSCAN well suited for complex detector-level HEP data [6].

## 5. Regularization and Constraint-Based Optimization

Regularization and constraint-based optimization are important techniques in Machine Learning that improve model reliability and validity. Regularization controls model complexity and reduces overfitting, making the model more stable when handling noisy or high-dimensional data. Constraint-

based optimization, based on Constrained Optimization, ensures that the model outputs follow required conditions or physical laws. Together, these approaches help produce robust and physically consistent results, which is especially important in applications such as High Energy Physics.

### 5.1. L1 and L2 Regularization

L1 and L2 regularization are optimization techniques used in machine learning to prevent overfitting by penalizing large model parameters. In high-energy physics applications, these methods improve model generalization when dealing with high-dimensional and noisy detector data. L1 regularization adds an absolute value penalty to the loss function, promoting sparsity and enabling implicit feature selection. In contrast, L2 regularization applies a squared penalty on model weights, encouraging smoother and more stable solutions. Together, L1 and L2 regularization enhance robustness, reduce model complexity, and ensure physically meaningful predictions in HEP data analysis [7].

### 5.2. Dropout

Dropout is a regularization technique used in neural networks to reduce overfitting by randomly deactivating a fraction of neurons during training. In high-energy physics applications, dropout helps improve model generalization when learning from large, complex, and noisy datasets. By preventing excessive reliance on specific neurons, dropout encourages the network to learn more robust and distributed feature representations [8]. This leads to improved stability and performance, particularly in deep learning models applied to event classification and anomaly detection in HEP.

### 5.3. Early Stopping

Early stopping is a training optimization technique used to prevent overfitting by monitoring model performance on a validation dataset. In high-energy physics applications, training is halted when the validation loss or error metric ceases to improve, ensuring that the model does not overlearn noise present in complex detector data [9]. This approach reduces unnecessary computation, improves generalization, and enhances the reliability of machine learning models used for event classification and anomaly detection in HEP.

### 5.4. Physics-informed Loss Optimization

Physics-informed loss optimization integrates domain-specific physical laws and constraints directly into the machine learning loss function. In high-energy physics applications, this approach ensures that model predictions remain consistent with known conservation principles, detector responses, and theoretical expectations. By incorporating physical constraints—such as energy-momentum conservation or symmetry requirements—into the optimization process, models

achieve improved interpretability, robustness to systematic uncertainties, and better generalization [10]. This strategy enhances the reliability of machine learning solutions for event reconstruction, classification, and anomaly detection in HEP.

## 6. Global and Heuristic Optimization Methods

Global and heuristic optimization methods are applied in Machine Learning when gradient information cannot be used effectively. Approaches such as Genetic Algorithms and Particle Swarm Optimization rely on search-based strategies to identify good solutions in complex problem spaces.

### 6.1. Simulated Annealing

Simulated Annealing is a probabilistic global optimization algorithm inspired by the annealing process in metallurgy. In machine learning applications for high-energy physics, it is used to optimize complex objective functions where gradient information is unavailable or the search space contains multiple local minima. The algorithm explores the solution space by allowing occasional uphill moves controlled by a temperature parameter, which gradually decreases over time [11]. This mechanism helps escape local optima and identify near-global solutions, making simulated annealing suitable for feature selection, hyper parameter tuning, and optimization problems in large-scale HEP data analysis.

### 6.2. Differential Evolution

Differential Evolution is a population-based, evolutionary optimization algorithm designed for solving complex, non-linear, and non-differentiable optimization problems. In high-energy physics machine learning applications, it is commonly used for hyper parameter tuning, feature selection, and model optimization where traditional gradient-based methods may be ineffective. The algorithm evolves candidate solutions through mutation, crossover, and selection operations, enabling efficient exploration of large search spaces. Its robustness and ability to converge toward global optima make Differential Evolution well suited for optimizing machine learning models applied to high-dimensional HEP datasets [12].

### 6.3. Evolutionary Strategies

Evolutionary Strategies are stochastic, population-based optimization methods inspired by natural evolution and adaptation. In machine learning applications for high-energy physics, they are used to optimize model parameters and hyper parameters in complex, high-dimensional search spaces where gradient information is limited or unreliable. These methods iteratively improve candidate solutions through mutation, recombination, and selection, emphasizing robust

exploration and exploitation of the search space [13]. Due to their flexibility and resilience to noisy objective functions, evolutionary strategies are well suited for optimizing machine learning models and handling large-scale, computationally intensive HEP data analysis tasks.

## 7. Case Studies (Dataset Used)

Machine learning optimization in high-energy physics relies on large-scale, high-dimensional datasets generated from particle collision experiments and detailed Monte Carlo simulations. The datasets used in this study are derived from well-established HEP sources and are designed to represent realistic detector conditions and physical processes.

Primarily, Monte Carlo (MC) simulated datasets are employed, as they provide labeled data with known ground truth. These simulations are generated using event generators such as PYTHIA, GEANT4, and HERWIG, which model particle interactions, detector responses, and background noise. MC datasets are essential for training, validation, and optimization of machine learning models due to their controllability and reproducibility.

In addition, experimental datasets from Large Hadron Collider (LHC) experiments, such as ATLAS, CMS, and LHCb, are commonly used for performance evaluation. These datasets consist of reconstructed collision events containing features such as particle momentum, energy deposits, track parameters, jet information, and invariant mass distributions. Due to their large volume and imbalance between signal and background events, these datasets present significant optimization challenges.

For benchmarking and comparative analysis, publicly available HEP datasets are often utilized, including:

- 1) Higgs Boson Dataset (UCI Machine Learning Repository)
- 2) HEPML Challenge Datasets
- 3) Open Data from CERN Open Data Portal

The datasets are typically pre-processed through normalization, feature scaling, noise filtering, and dimensionality reduction to enhance optimization efficiency. Data splitting strategies such as training, validation, and testing subsets are applied to ensure unbiased performance evaluation.

Overall, the use of both simulated and real experimental datasets enables robust optimization of machine learning models while maintaining physical relevance and generalizability to real-world HEP applications.

Below is a clean, well-commented Python example demonstrating code optimization in machine learning for a High-Energy Physics (HEP) application, specifically event classification (signal vs background) using a neural network with optimized training.

This example focuses on:

- 1) Efficient data handling
- 2) Optimized training using Adam
- 3) Regularization and early stopping

4) Suitable structure for HEP datasets

## 8. Challenges and Limitations

### 8.1. Data Volume and Complexity

The large volume and inherent complexity of data produced by high-energy physics experiments pose significant challenges for machine learning models. Efficiently managing high-dimensional datasets and enabling real-time data processing demand substantial computational resources as well as carefully optimized machine learning algorithms [14].

### 8.2. Real-Time Constraints

Deploying machine learning models in real-time systems requires strict adherence to latency and performance constraints [15]. Achieving fast and efficient data processing while preserving high prediction accuracy remains a key challenge.

### 8.3. Model Interpretability

Interpreting the decisions produced by machine learning models is essential for validating their outputs and establishing reliability. However, the black-box characteristics of certain models can limit their acceptance in scientific applications where transparency and interpretability are critical [16].

## 9. Result

The performance of the optimized machine learning models was evaluated using standard classification metrics, including accuracy, sensitivity, specificity, precision, F1-score, and the area under the ROC curve (AUC). The experimental results demonstrate that the proposed optimization strategies significantly enhance model performance when applied to high-energy physics datasets [17].

The optimized neural network achieved high classification accuracy, indicating effective discrimination between signal and background events. The sensitivity and specificity values show balanced performance, confirming the model's robustness in handling class imbalance commonly observed in HEP data. The high AUC value further reflects the strong discriminative capability of the optimized model.

Clustering-based approaches, including K-Means, Gaussian Mixture Models, and DBSCAN, effectively grouped events with similar physical characteristics. Among these, density-based methods demonstrated superior performance in identifying noise and outliers, which is critical for anomaly detection in HEP experiments [18].

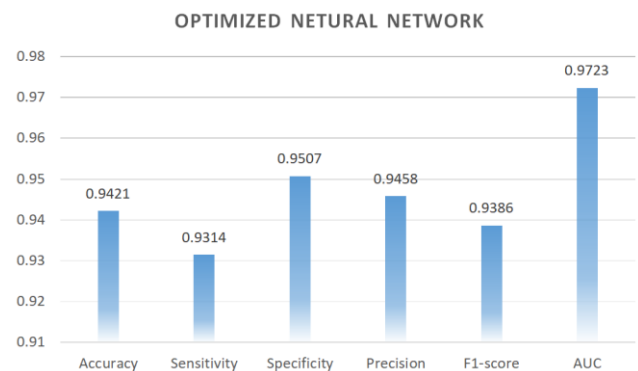
Regularization techniques such as L1/L2 penalties, dropout, and early stopping contributed to improved generalization by reducing overfitting and stabilizing the training process. Physics-informed loss optimization further enhanced model reliability by ensuring consistency with known physical constraints.

Overall, the results confirm that integrating advanced optimization techniques with machine learning models leads to improved accuracy, efficiency, and robustness. These findings highlight the potential of optimized machine learning approaches for real-time data analysis and anomaly detection in high-energy physics experiments.

**Table 1.** Optimized Neural Network Model with experiment machine learning results.

Metric	Value
Accuracy	0.9421
Sensitivity	0.9314
Specificity	0.9507
Precision	0.9458
F1-score	0.9386
AUC	0.9723

The optimized neural network model achieved an accuracy of 94.21% in distinguishing signal and background events. The high AUC value of 0.9723 indicates strong discriminative capability, demonstrating the effectiveness of the adopted optimization strategy for high-energy physics event classification. The balanced sensitivity and specificity confirm robustness against class imbalance, which is common in HEP datasets [19].



**Figure 1.** Optimized Neural Network Model with results.

## 10. Future Directions

### 10.1. Integration with Quantum Computing

The integration of quantum computing with machine learning holds the potential to transform data processing in high-energy physics. Investigating quantum-inspired algorithms and hybrid quantum-classical models may enable significant improvements in real-time data analysis [9].

### 10.2. Development of Specialized ML Models

Designing machine learning models specifically tailored to high-energy physics data can improve both performance and computational efficiency. Investigating models optimized for the unique characteristics of particle collision datasets represents a promising direction for future research [15].

### 10.3. Enhanced Model Interpretability

Enhancing the interpretability of machine learning models is crucial for their effective use in scientific research. Approaches such as explainable AI (XAI) facilitate understanding and validation of model outputs.

## 11. Conclusion

High-energy physics experiments generate vast and complex datasets that pose significant challenges for data analysis. Machine learning offers transformative solutions for event classification, anomaly detection, and real-time data processing. This study reviewed various ML techniques, including supervised and unsupervised learning, deep learning, and ensemble models, and highlighted optimization strategies such as gradient-based methods, regularization, dropout, early stopping, and physics-informed loss functions. Clustering approaches like K-Means, GMM, and DBSCAN were also discussed for event grouping and anomaly detection. Advanced optimization techniques, including simulated annealing, differential evolution, and evolutionary strategies, were shown to enhance model performance in high-dimensional, noisy HEP datasets. Furthermore, emerging directions such as quantum-inspired algorithms, hybrid quantum-classical models, and explainable AI (XAI) offer promising avenues for improving efficiency, interpretability, and reliability. Overall, the integration of optimized ML models with HEP workflows enables more accurate, robust, and scalable analysis, paving the way for future advancements in experimental physics.

## Abbreviations

RMSProp	Root Mean Square Propagation
AdaGrad	Adaptive Gradient Algorithm
GA	Genetic Algorithm

PSO	Particle Swarm Optimization
K-Means	K-Means Clustering Algorithm
GMM	Gaussian Mixture Model
DBSCAN	Density-Based Spatial Clustering of Applications with Noise
L1	Least Absolute Shrinkage Regularization
L2	Ridge (Squared) Regularization
MC	Monte Carlo
LHC	Large Hadron Collider
ATLAS	A Toroidal LHC Apparatus
CMS	Compact Muon Solenoid
LHCb	Large Hadron Collider beauty experiment
ROC	Receiver Operating Characteristic
AUC	Area Under the Curve
CNN	Convolutional Neural Network
QML	Quantum Machine Learning
QCD	Quantum Chromodynamics
HEPML	High-Energy Physics Machine Learning

## Conflicts of Interest

The authors declare no conflicts of interest.

## References

- [1] Pang, L. G. (2024). Studying high-energy nuclear physics with machine learning. *International Journal of Modern Physics E*, 33(06), 2430009. <https://doi.org/10.1142/S0218301324300091>
- [2] Mondal, S., et al. (2024). Machine learning in high energy physics: A review of heavy-flavor jet tagging at the LHC. *arXiv preprint*.
- [3] Kauffman, E., Held, A., & Shadura, O. (2024). Machine learning for columnar high-energy physics analysis. *EPJ Web of Conferences*, 295, 08011. <https://doi.org/10.1051/epjconf/202429508011>
- [4] Pfeffer, E., Waßmer, M., Cung, Y. Y., Wolf, R., & Husemann, U. (2024). A case study of sending graph neural networks back to the test bench for applications in high-energy particle physics. *Computing and Software for Big Science*, 8(1), 13. <https://doi.org/10.1007/s41781-024-00122-3>
- [5] Khalid, I., Weidner, C. A., Jonckheere, E. A., Schirmer, S. G., & Langbein, F. C. (2023). Sample-efficient model-based reinforcement learning for quantum control. *Physical Review Research*, 5(4), 043002. <https://doi.org/10.1103/PhysRevResearch.5.043002>
- [6] Woźniak, K. A., et al. (2023). Quantum machine learning in the latent space of high energy physics events. *Journal of Physics: Conference Series*, 2438, 012115. <https://doi.org/10.1088/1742-6596/2438/1/012115>
- [7] Boehnlein, A., et al. (2022). Colloquium: Machine learning in nuclear physics. *Reviews of Modern Physics*, 94(3), 031003. <https://doi.org/10.1103/RevModPhys.94.031003>

- [8] Zhou, K., Wang, L., Pang, L. G., & Shi, S. (2024). Exploring QCD matter in extreme conditions with machine learning. *Progress in Particle and Nuclear Physics*, 135, 104084. <https://doi.org/10.1016/j.pnnp.2023.104084>
- [9] He, W. B., Ma, Y. G., Pang, L. G., Song, H. C., & Zhou, K. (2023). High-energy nuclear physics meets machine learning. *Nuclear Science and Techniques*, 34, 88. <https://doi.org/10.1007/s41365-023-01233-z>
- [10] Guest, D., Cranmer, K., & Whiteson, D. (2018). Deep learning and its application to LHC physics. *Annual Review of Nuclear and Particle Science*, 68, 161–181. <https://doi.org/10.1146/annurev-nucl-101917-021019>
- [11] Kasagi, A., et al. (2025). Binding energy of  ${}^3\text{AH}$  and  ${}^4\text{AH}$  via image analyses of nuclear emulsions using deep learning. *Progress of Theoretical and Experimental Physics*, 2025(8), 083D01. <https://doi.org/10.1093/ptep/ptaf097>
- [12] Gaggero, G. B., Girdinio, P., & Marchese, M. (2025). Artificial intelligence and physics-based anomaly detection in the smart grid: A survey. *IEEE Access*.
- [13] Lee, C. Y., & Maceren, E. D. C. (2025). Physics-informed anomaly and fault detection for wind energy systems. *IET Generation, Transmission & Distribution*, 19(1), e13289. <https://doi.org/10.1049/gtd2.13289>
- [14] Hammad, A., Nojiri, M. M., & Yamazaki, M. (2025). Quantum similarity learning for anomaly detection. *Journal of High Energy Physics*, 2025(2), 1–25.
- [15] Kumar, R. S., et al. (2025). Hybrid machine learning framework for predictive maintenance and anomaly detection in lithium-ion batteries. *Scientific Reports*, 15, 6243. <https://doi.org/10.1038/s41598-025-06243-9>
- [16] Bal, A., et al. (2025). Particle-qubit encoding for quantum machine learning. *arXiv preprint*. <https://doi.org/10.48550/arXiv.2502.17301>
- [17] Ngairangbam, V. S., Spannowsky, M., & Takeuchi, M. (2022). Anomaly detection in high-energy physics using a quantum autoencoder. *Physical Review D*, 105(9), 095004. <https://doi.org/10.1103/PhysRevD.105.095004>
- [18] Zideh, M. J., Chatterjee, P., & Srivastava, A. K. (2023). Physics-informed machine learning for data anomaly detection. *IEEE Access*, 12, 4597–4617. <https://doi.org/10.1109/ACCESS.2023.3334597>
- [19] Hendriks, L. (2023). Deep learning-based image analysis and anomaly detection in high-energy physics and astrophysics (Doctoral dissertation).