# A Comprehensive Assessment of Optimisation Methods for Machine Learning Instruction

[1]Prakruti D. Dave, [2]Ankit J. Solanki, [3*]Hiren S. Lekhadiya
[1,3] Assistant Professor, School of Engineering, P P Savani University
[2] Teaching Assistant, School of Engineering, P P Savani University

**Abstract: An extensive review of optimisation methods for training machine learning (ML) models a crucial branch of artificial intelligence is provided in this article. ML uses statistical techniques to allow systems to learn from experience and become more intelligent without the need for explicit programming. The study emphasises the value of optimisation in machine learning, emphasising how it can be used to improve training efficiency and generalisation by modifying model parameters to minimise loss functions. Numerous optimisation techniques are examined, such as Constraint-based techniques, Gradient Descent Variants, Adaptive Learning Rate Techniques, Second-Order Optimisation Techniques, and Bayesian Optimisation. Every segment delves into the fundamentals, uses, and advantages of these methods, highlighting their significance in addressing issues like overfitting, scalability, and computational effectiveness. The purpose of this page is to help practitioners, academics, and enthusiasts navigate the wide range of optimisation approaches designed for various machine learning algorithms and applications.**

**Keywords: Optimization, Adaptive, Second-Order Optimization, Constraint-based Methods, Overfitting, Scalability, Artificial Intelligence.**

## 1. INTRODUCTION

Optimization methods are fundamental to the efficient training of machine learning models. Artificial intelligence (AI) includes machine learning (ML), which gives systems the ability to learn from experience and get better on their own without explicit programming. By employing statistical techniques, ML enables machines to enhance their performance on tasks through exposure to increasing amounts of data over time. Three main categories exist for machine learning: supervised, unsupervised, and reinforcement learning techniques [1 - 2]. The goal of supervised learning algorithms is to map inputs to outputs based on input-output pairings by learning from labelled data. Image classification is a common example, in which the algorithm is trained to identify objects from labelled images. Conversely, unsupervised learning algorithms use unlabelled data to find patterns and structures in the data. This type is exemplified by clustering, where the algorithm groups comparable data points without predetermined labels. Reinforcement learning algorithms pick up new skills by interacting with their surroundings in order to maximise cumulative rewards or accomplish goals [3]. When teaching a computer programme to play chess, for example, the algorithm must learn the best moves via trial and error and be rewarded for employing winning tactics.

Optimisation in machine learning refers to changing model parameters in order to minimise or maximise a loss function. Finding the optimal model parameters that produce the lowest loss function value is the main goal. The difference between the model's expected and actual outputs is quantified by a loss function, which evaluates the model's performance. The objective of training is to reduce this "loss," honing the model to increase forecast precision. By iteratively modifying parameters to lessen the difference between forecasts and actual results, optimisation makes sure that the model improves its predictive power over time and becomes more capable of generalising to new data.

Algorithms that modify model parameters in order to minimise a loss function are used in the training of machine learning models. In this process of adjustment, optimisation techniques are essential because they guarantee that models learn from data in an efficient manner [4-6]. Millions or even billions of parameters in sophisticated models are present in many machine learning problems. These models can be trained within appropriate time and computational resource restrictions thanks to effective optimisation techniques. In order to minimise the possibility of overfitting—a situation in which a model performs

well on training data but badly on test data—effective optimisation approaches are crucial for creating models that generalise effectively to unseen data [7–8]. Scalable optimisation techniques are more and more important as datasets get bigger and models get more intricate. By using these methods, models are guaranteed to be computationally feasible even when dealing with large volumes of data [9]. Many models and algorithms fall under the umbrella of machine learning, and each one needs unique optimisation strategies based on its unique properties. Optimisation techniques offer the adaptability to modify and refine these algorithms to fit particular issues.

## 2. VARIANTS OF GRADIENT DESCENT

Gradient Descent updates parameters by moving them in the direction opposite to the gradient of the loss function. Stochastic Gradient Descent (SGD) accelerates computation by using a subset of data to compute the gradient, though it introduces more noise compared to standard gradient descent. Mini-Batch Gradient Descent strikes a balance between full dataset processing and SGD, updating parameters with small random batches of data. Momentum enhances stability by incorporating a fraction of the previous update vector into the current one, aiding navigation through regions of high curvature. These Gradient Descent Variants are versatile optimization techniques applicable across diverse domains such as machine learning, deep learning, and optimization theory.

Gradient-based methods are pivotal for minimizing loss functions, optimizing model parameters, and configuring objectives effectively in complex optimization landscapes. Training deep learning models like multi-layer perceptron's (MLPs) relies heavily on gradient-based optimization to minimize loss functions. Similarly, Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) employ various forms of gradient descent to adjust weights and biases during training. Figure 1 [10] depicts the difference between noiseless and noisy SGD training. Using a single example introduces fluctuations, resulting in winding and slower convergence paths during iterations. Figure 2 [11] illustrates vertical oscillations in the gradient's direction. A straightforward approach to mitigate this issue is to stabilize the gradient horizontally while minimizing vertical fluctuations.
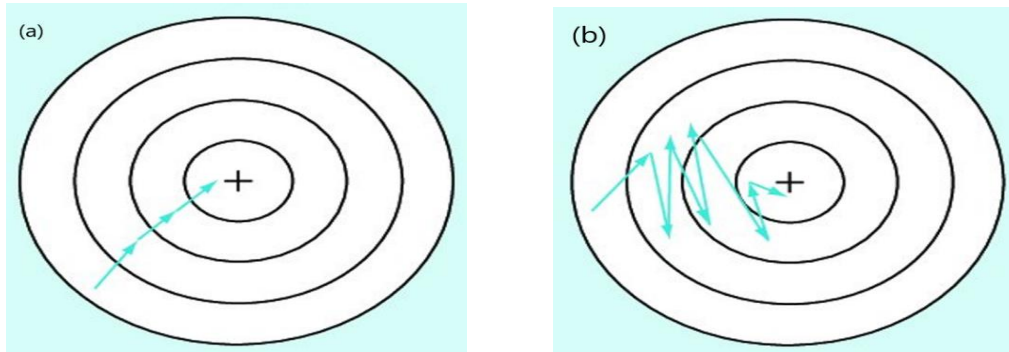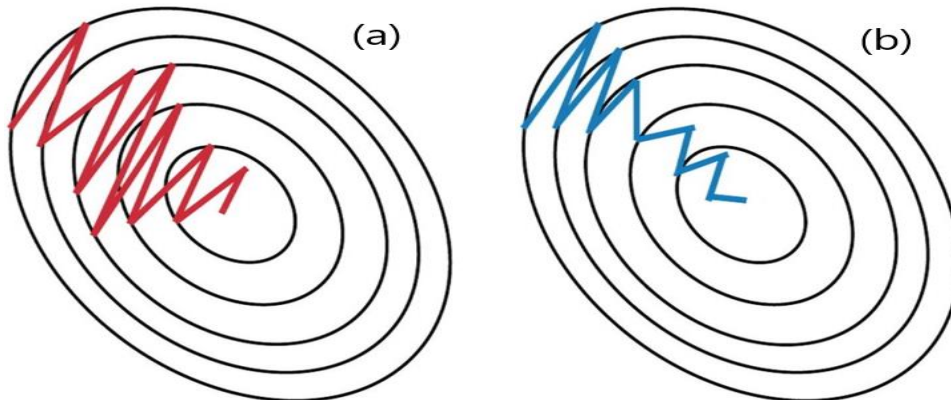


Figure 1. Comparison of (a) GD and (b) SGD



Figure 2. Comparison of SGD algorithms (a) without and (b) with momentum

## 3. METHODS FOR ADAPTIVE LEARNING RATES

Adaptive learning rate techniques, such as RMSprop and Adam, are essential for machine learning model optimisation because they dynamically modify all model parameter learning rates according to their historical gradients. An exponentially decreasing average of squared gradients is used to divide the learning rate in RMSprop, an Adagrad modification. Adam (Adaptive Moment Estimation) adapts learning rates by using both the first and second moments of gradients, combining ideas from Momentum and RMSprop.

These methods accelerate optimization algorithms by dynamically adjusting learning rates, ensuring efficient progress in scenarios where the loss function landscape varies across dimensions or iterations. Fixed learning rates in standard algorithms can lead to oscillations or divergence, particularly in parameter spaces with steep gradients. Adaptive methods mitigate these issues by scaling the learning rate based on gradient nature and history. Adaptive learning rate methods provide resilience against challenges such as saddle points, gradient issues, and non-stationary objectives. By adjusting learning rates adaptively, these methods navigate complex optimization landscapes more effectively.

Traditional optimization approaches often require manual tuning of hyperparameters like the learning rate, a task that can be cumbersome and time-consuming. Adaptive learning rate methods alleviate this burden by automatically adjusting hyperparameters based on optimization problem characteristics, thereby reducing the need for extensive manual tuning.

## 4. OPTIMIZATION METHODS THAT UTILIZE SECOND-ORDER

Newton's Method guides the optimisation process, which is a computationally demanding effort, especially with huge datasets, by using the second derivative, also referred to as the Hessian matrix [12]. Without explicitly calculating the second derivative, Quasi-Newton Methods [13] use gradient information to approximate the Hessian matrix and speed up convergence. As the name implies, second-order derivatives—specifically, the Hessian matrix—are incorporated into the optimisation process by means of second-order optimisation techniques. On the other hand, first-order techniques such as gradient descent only use the objective function's first derivative, or gradient.

Compared to first-order approaches, second-order methods frequently reach optimal solutions in fewer iterations, especially when working with poorly conditioned functions. When curvature fluctuates greatly, they work particularly well for extremely nonlinear or ill-conditioned functions [14]. By revealing the local curvature of the optimisation landscape, these techniques help algorithms find local minima or maxima more quickly. This feature comes in handy when dealing with intricate optimisation issues where the first-order gradient might not provide much insight. Insufficient curvature information might cause first-order techniques to oscillate or behave erratically in certain situations. Second-order approaches provide more steady convergence behaviour by accounting for curvature.

For optimization problems where the Hessian matrix [15] provides crucial curvature and condition information about the objective function, second-order methods excel in adaptively adjusting step sizes to enhance convergence.

## 5. METHODS OF CONSTRAINTS-BASED

By using gradient descent, Projected Gradient Descent updates parameters while guaranteeing that, in the event that certain requirements are not met, the new values are restricted to a workable range. Evolutionary Algorithms, such as Genetic Algorithms, simulate natural selection by employing mechanisms like mutation, crossover, and selection to evolve solutions for optimization problems. Particle Swarm Optimization (PSO) manages a population of candidate solutions (particles) that navigate the search space based on their individual best positions and the swarm's best-known position.

Constraint-based methods in machine learning and optimization seek optimal solutions while adhering to specific constraints. These techniques enforce predefined boundaries or conditions to ensure that solutions meet criteria essential to the problem's domain. Many real-world problems, like resource allocation or scheduling, impose constraints related to

time, capacity, or resource availability. Constraint-based methods guarantee that solutions are both optimal and feasible within these constraints. Certain optimization problems, particularly in engineering design, involve multiple constraints related to materials, safety, and cost. Constraint-based methods enable the identification of optimal solutions that satisfy all specified constraints, balancing optimality with practical implementation in real-world scenarios.

Integrating constraints helps mitigate risks associated with decisions by ensuring solutions adhere to safety, operational, or regulatory requirements. In complex optimization scenarios with multiple objectives, constraints provide clear guidelines that structure decision-making, ensuring coherent and aligned solutions with the problem's overarching goals.

## 6.    BAYESIAN OPTIMIZATION

Bayesian Optimization is a technique for optimizing expensive and black-box functions using probabilistic models. Unlike traditional methods that rely on derivatives or gradients, Bayesian Optimization constructs a probabilistic model of the objective function. This model guides the selection of the next point to evaluate [16].

Fundamentally, Bayesian Optimisation describes the behaviour of the goal function using probabilistic models, most commonly Gaussian Processes. These models calculate the function's value at as-yet-undiscovered places and express the degree of uncertainty in those estimates. Iteratively choosing points based on the existing model, evaluating the objective function, updating the model with new observations, and repeating the process until convergence or a stopping requirement is satisfied are all part of the sequential character of Bayesian Optimisation. A key strength of Bayesian Optimization lies in its ability to balance exploration and exploitation. It efficiently explores uncertain regions of the parameter space while exploiting promising regions to avoid getting trapped in local optima.

In practical applications, evaluating the objective function such as assessing the performance of machine learning models or conducting physical experiments— can be time-consuming or costly. Bayesian Optimization aims to minimize the number of evaluations needed to find the optimal solution, making it particularly suitable for scenarios where function evaluations are expensive.

When the mathematical form of the objective function is unknown or too complex for analytical methods, Bayesian Optimization provides a robust framework. Traditional optimization methods like gradient descent may struggle with non-convex or multimodal functions, often converging to local optima. Bayesian Optimization's strategy of balancing exploration and exploitation enables efficient navigation of complex landscapes to seek the global optimum. In machine learning and deep learning, where model performance hinges on hyperparameters (e.g., learning rate, regularization parameters), Bayesian Optimization offers an efficient approach to tuning these parameters. This method enhances model performance without resorting to exhaustive grid or random search methods.

## 7.    ADVANTAGES AND DRAWBACKS OF OPTIMIZATION TECHNIQUES

Gradient descent variants are widely favoured for their simplicity in implementation, applicability across diverse problem domains, and computational efficiency. Despite these advantages, they can converge slowly, especially in high-dimensional spaces, and are susceptible to local minima. Adam serves as an example of adaptive learning rate techniques, which dynamically modify learning rates for each parameter to promote quicker convergence and better results in deep learning tasks. But they could be unstable and hyperparameter-sensitive, so careful adjustment is required. Second-order optimization methods like Newton's Method utilize curvature information to achieve rapid convergence, particularly in regions with pronounced curvature. Yet, they are computationally demanding due to the need for computing and inverting the Hessian matrix, and they may encounter numerical instability issues.

Constraint-based approaches are invaluable in a variety of applications because they impose limits on parameter values to guarantee model stability and interpretability. But they could require specific optimisation methods and be computationally demanding, especially for complicated constraints. Bayesian optimization models the objective function using a probabilistic surrogate, effectively exploring

the parameter space to enhance sample efficiency. Nonetheless, its effectiveness can depend on prior knowledge or assumptions about the objective function and may diminish in high-dimensional spaces.

## 8. CONCLUSION

Optimization methods in machine learning offer a spectrum of strengths and limitations tailored to diverse problem scenarios and objectives. Gradient descent variants are favoured for their simplicity, wide applicability, and efficiency, yet they face challenges like slow convergence and susceptibility to local minima in high-dimensional spaces. Adaptive learning rate methods such as Adam promise quicker convergence and improved performance in deep learning but require meticulous hyperparameter tuning to manage instability. Second-order methods like Newton's Method leverage curvature insights for rapid convergence but at the cost of computational intensity and numerical stability concerns. Regularization methods effectively curb overfitting yet introduce bias and demand fine-tuning. Constraint-based approaches ensure stability and interpretability under defined constraints, though they may need specialized techniques for intricate rules. Bayesian optimization enhances efficiency through probabilistic modeling but hinges on accurate prior knowledge, posing limitations in high-dimensional contexts. Each method addresses specific challenges from parameter optimization to resilience against overfitting and navigation of complex landscapes empowering practitioners to tailor techniques according to data, model characteristics, and optimization goals, thereby enhancing the reliability and efficacy of machine learning applications.

## REFERENCES

[1] M. M. Taye, "Understanding of Machine Learning with Deep Learning: Architectures, Workflow, Applications and Future Directions," Computers, vol. 12, no. 5, p. 91, 2023.

[2] A. Haleem, M. Javaid, M. A. Qadri, R. P. Singh, R. Suman, "Artificial intelligence (AI) applications for marketing: A literature-based study," International Journal of Intelligent Networks, vol. 3, pp. 119-132, 2022.

[3] R. Pugliese, S. Regondi, R. Marini, "Machine learning-based approach: global trends, research directions, and regulatory standpoints," Data Science and Management, vol. 4, pp. 19-29, 2021.

[4] M. Gupta, K. Rajnish, V. Bhattacharjee, "Impact of Parameter Tuning for Optimizing Deep Neural Network Models for Predicting Software Faults," Scientific Programming, vol. 2021, pp. 1-17, 2021.

[5] M. J. Bianco et al., "Machine learning in acoustics: Theory and applications," The Journal of the Acoustical Society of America, vol. 146, no. 5, pp. 3590-3628, 2019.

[6] T. Q. Bao and C. Gutiérrez, "Ekeland variational principles for vector equilibrium problems," Optimization, vol. 73, no. 1, pp. 29-62, 2024.

[7] S. Ali, "Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence," Information Fusion, vol. 99, p. 101805, 2023.

[8] S. Raschka, J. Patterson, C. Nolet, "Machine Learning in Python: Main Developments and Technology Trends in Data Science, Machine Learning, and Artificial Intelligence," Information, vol. 11, no. 4, p. 193, 2020.

[9] U. Sivarajah, M. M. Kamal, Z. Irani, V. Weerakkody, "Critical analysis of Big Data challenges and analytical methods," Journal of Business Research, vol. 70, pp. 263-286, 2017.

[10] Y. Tian, Y. Zhang, H. Zhang, "Recent Advances in Stochastic Gradient Descent in Deep Learning," Mathematics, vol. 11, no. 3, p. 682, 2023.

[11] L. Bottou, O. Bousquet, "The tradeoffs of large scale learning," Advances in Neural Information Processing Systems 20 (NIPS 2007), pp. 161-168, 2008.

[12] K. Dassios, "Analytic Loss Minimization: Theoretical Framework of a Second Order Optimization Method," Symmetry, vol. 11, no. 2, p. 136, 2019.

[13] J. Z. Zhang, N.Y. Deng, L. H. Chen, "New Quasi-Newton Equation and Related Methods for Unconstrained Optimization," Journal of Optimization Theory and Applications, vol. 102, pp. 147-167, 1999.

[14] T. Guo, Y. Liu, C. Han, "An Overview of Stochastic Quasi-Newton Methods for Large-Scale Machine Learning," Journal of the

Operations Research Society of China, vol. 11, pp. 245-275, 2023.

[15] J. M. Bofill, "Updated Hessian matrix and the restricted step method for locating transition structures," Journal of Computational Chemistry, vol. 15, no. 1, pp. 1-11, 1994.

[16] X. Wang, Y. Jin, S. Schmitt, M. Olhofer, "Recent Advances in Bayesian Optimization," ACM Computing Surveys, vol. 55, no. 13s, pp. 1-36, 2023.