

# A Comprehensive Overview of Gradient Descent and its Optimization Algorithms

Atharva Tapkir<sup>1</sup>

Student, Computer Engineering, Sinhgad College of Engineering, Pune, India<sup>1</sup>

**Abstract:** This study explores machine learning gradient-based optimization algorithms, highlighting the critical importance of gradient descent and investigating adaptive strategies to improve its performance. The fundamental technique of optimization is gradient descent, although balancing convergence speed and accuracy can be difficult due to gradient descent's reliance on fixed learning rates. The study explores a variety of adaptive learning techniques, including as drop, decay, cyclic learning, and adaptive learning. These techniques are designed to modify learning rates in real-time during optimization, hence affecting stability and convergence. Additionally, the research delves into momentum-based methods like Adam, RMSProp, AdaGrad, and AdaDelta, clarifying their use in reducing the difficulties associated with traditional gradient descent. The study also clarifies gradient clipping methods, addressing the problem of exploding gradients and offering solutions to stabilize and enhance machine learning models. The goal of this thorough investigation is to provide practitioners with a sophisticated grasp of optimization techniques so they may guide machine learning models toward effectiveness, precision, and robustness in a variety of application domains.

**Keywords:** Gradient Descent, Optimizations, Learning Rate, Adam, Neural Networks

## I. INTRODUCTION

In the realm of machine learning, the quest for optimizing algorithms to achieve rapid convergence and precise parameter estimation is perpetual. At the core of this pursuit lies gradient descent, a foundational algorithm pivotal in optimizing objective functions within machine learning models. Gradient descent is the bedrock of optimization in machine learning, leveraging the gradient of an objective function to iteratively update model parameters. Its significance lies in its ability to navigate the multidimensional parameter space, steering the model towards optimal solutions by minimizing the objective function. However, traditional gradient descent is not without limitations. Its dependence on learning rates presents challenges in striking a balance between convergence speed and accuracy. The need for adaptive optimization algorithms arises from these shortcomings, aiming to fine-tune the learning process and mitigate the complexities associated with conventional gradient descent.

This paper seeks to unravel the nuances of optimization strategies encompassing various adaptive learning methods that enhance gradient-based algorithms. The focus will extend to elucidating decay, drop, cyclic learning, and adaptive learning approaches, each tailored to dynamically adjust learning rates during the optimization process. The objective is to elucidate their impact on convergence speed, model stability, and the ability to navigate complex optimization landscapes within machine learning tasks.

Additionally, the research will delve into momentum-based approaches, such as Adam, RMSProp, AdaGrad, and AdaDelta, each designed to address specific challenges encountered in traditional gradient descent. These algorithms incorporate momentum, adaptive learning rates, and variance reduction techniques to expedite convergence, tackle oscillations, and navigate saddle points effectively.

Furthermore, the paper will shed light on the significance of gradient clipping techniques, addressing issues like exploding gradients common in deep learning architectures. By exploring these techniques, the research aims to equip practitioners with insights into mitigating challenges associated with gradient magnitudes, fostering stable and efficient training of machine learning models.

In essence, this paper serves as a technical exploration into the realm of gradient-based optimization techniques and their adaptive counterparts. By dissecting their mathematical formulations, technical underpinnings, and empirical effects, the research endeavours to offer a comprehensive understanding of these algorithms' mechanisms, enabling practitioners to navigate and harness the optimization landscape for more efficient and robust machine learning models across various domains.

## II. GRADIENT DESCENT

Tasks in machine learning can be expressed as the problem of optimizing an objective function  $f(\theta)$  defined over some domain  $\theta \in \lambda$ . The goal in this case is to find the minimizer  $\theta^* = \arg \min_{\theta \in \lambda} f(\theta)$ . Although any technique that can minimize this objective function can be used, gradient descent—a technique that produces a series of updates—is the typical method for differentiable functions.

$$\theta_{t+1} = \theta_t - \alpha_t \nabla f(\theta_t)$$

We only use the beginning point arbitrarily to assess the performance. We will calculate the derivative, or slope, from that initial point. Using a tangent line, we can then determine how steep the slope is. The weights and bias changes to the parameters will be based on the slope. The slope will be steeper at the beginning, but it should progressively get less steep as new parameters are generated, until it reaches the point of convergence, which is the lowest point on the curve.

The objective of gradient descent is to minimize the cost function, or the error between the anticipated and actual values of  $y$ . This is similar to the process of determining the linear regression line of best fit. It needs two data points to accomplish this: a direction and a learning rate. Future iterations' partial derivative computations are determined by these elements, enabling the process to progressively approach the local or global minimum (also known as the point of convergence).

The number of steps required to reach the minimum is known as the learning rate, sometimes known as the alpha or step size. This is usually a small value that is updated and assessed in accordance with the cost function's behavior. Larger steps are produced by high learning rates, but there is a chance of exceeding the minimum. Low learning rates, on the other hand, have tiny step sizes. Although it offers the benefit of greater precision, the quantity of iterations reduces overall efficiency because it requires more calculations and time to reach the minimum. The difference, or inaccuracy, between actual and anticipated value at a given position is measured by the cost (or loss) function. By giving the machine learning model feedback, this increases the model's effectiveness by enabling it to modify its parameters in order to reduce error and locate the local or global minimum. Until the cost function approaches or reaches zero, it iterates constantly in the direction of the steepest descent, also known as the negative gradient. The model will then cease to learn at this time. Furthermore, there is a small distinction between the phrases "cost function" and "loss function," despite the fact that they are often used interchangeably.

The steepest-descent method, which updates parameters using the gradient of the loss function, is the most often used technique for parameter learning in neural networks. When steps of finite size are taken into account, the steepest-gradient approach does not always lead in the optimal direction of progress, so occasionally it will behave unexpectedly. Only when looking at tiny steps can the direction with the sharpest descent be considered the best one. Occasionally, a minor adjustment to the parameters can turn a steepest-descent direction into an ascending direction. Numerous course modifications are therefore required. Whenever the steepest-descent vector travels along a direction of high curvature in the loss function, oscillation and zigzagging become a very common problem.

## III. LEARNING RATE

In the realm of machine learning, parameters can be categorized into two main types: machine learnable parameters and hyper-parameters. Machine learnable parameters are inherently learned or estimated by algorithms during training on a given dataset. On the other hand, hyper-parameters are values specifically assigned by machine learning engineers or data scientists. These hyper-parameters play a crucial role in governing how algorithms learn and in fine-tuning the model's performance. The learning rate specifically controls the pace at which an algorithm adjusts parameter estimates or acquires the values of these parameters during the learning process. As training progresses, altering and decreasing the learning rate gradually allows for smaller steps towards convergence. This adjustment prevents oscillations around the optimal values and aids in fine-tuning the model, enhancing its ability to find the best parameters for improved performance. The following are some effective strategies:

### A. *Decaying Learning Rate*

Because it puts the analyst in a difficult situation, a constant learning rate is undesirable. This is the conundrum. The algorithm will take too long to reach an optimal solution if a lower learning rate is applied early on. However, if the high learning rate is sustained, the algorithm will oscillate around the point for an extended period of time or diverge in an unstable manner. On the other hand, a big initial learning rate will initially enable the algorithm to approach a good

solution very close. It is not optimal to maintain a constant learning rate in either scenario. To automatically obtain the required learning-rate adjustment and avoid these issues, one can let the learning rate decay  $\alpha_t$  be expressed in terms of the initial decay rate  $\alpha_0$  and epoch  $t$  as follows:

Exponential Decay:

$$\alpha_t = \alpha_0 \exp(-k.t)$$

Inverse Decay:

$$\alpha_t = \frac{\alpha_0}{1 + k.t}$$

The decay rate is determined by the parameter  $k$ . A different strategy is to use step decay, where every few epochs the learning rate is lowered by a specific factor. One possible approach is to multiply the learning rate every five epochs by 0.5. Reducing the learning rate whenever the loss on a held-out portion of the training data set stops improving is a common strategy. The analyst may even use an implementation where the learning rate can be manually adjusted based on progress, and in certain cases, even supervise the learning process. Although it ignores many other troubling issues, this kind of approach can be applied to basic gradient descent implementations.

### B. *Scheduled Drop*

The drop method involves a specified proportional reduction in the learning rate at a predetermined frequency, as opposed to the decay method's monotonous learning rate drop. The equation below displays the formula used to calculate for a specific epoch:

$$\alpha_n = \alpha_0 \times D^{\frac{n}{p}}$$

The initial learning rate ( $\alpha_0$ ), the epoch/iteration number ( $n$ ), a hyper-parameter ( $D$ ) that indicates how much the learning rate must decrease, and another hyper-parameter ( $p$ ) that indicates the frequency of learning rate drops based on epochs are all included in the equation above. The drawback of both the decay and drop approaches is that they do not assess whether or not lowering the learning rate is necessary. In both approaches, regardless of the complexity of the cost function minimization, the learning rate falls.

### C. *Adaptive Learning*

With this method, the cost function's gradient value determines how quickly or slowly the learning rate changes. The learning rate will increase with decreasing gradient value and decrease with increasing gradient value. Therefore, for steeper and shallower regions of the cost function curve, respectively, the learning accelerates and decelerates. The equation below displays the formula utilized in this method.

$$\alpha_n = \frac{\alpha_0}{\sqrt{S_n}}$$

The momentum factor, ' $S_n$ ' is determined by using the following equation to the previous equation, where ' $\alpha_0$ ' is the initial learning rate. The number of epochs or iterations is ' $n$ '.

$$S_n = \gamma S_{n-1} + (1 - \gamma) \frac{\partial CF}{\partial \beta} ]_n$$

The hyperparameter " $\gamma$ " in the equation above usually has a value between 0.7 and 0.9. Keep in mind that the momentum factor  $S_n$  in the previous equation is the exponential weighted average of the gradients. Therefore, in order to calculate the momentum component, not only the value of the present gradient but also the values of gradients from earlier epochs are taken into account.

The momentum factor " $S_n$ " is larger and the gradient is large when the cost function curve is steep. As a result, the learning rate is lower. The learning rate is higher, the gradient is tiny, the momentum factor " $S_n$ " is likewise small, and the cost function curve is shallow. The gradient adapted learning rate strategy eliminates the drawbacks of the decay and drop approaches by considering the gradient of the cost function when deciding how to raise or lower the learning rate. One well-liked technique for training deep neural networks is stochastic gradient descent.

D. Cyclic Learning

This method involves cyclically varying the learning rate between a base rate and a maximum rate. At a set frequency, the learning rate fluctuates in a triangle shape between the maximum and base rates. It has been observed that alternative forms, including parabolic or sinusoidal, also provide comparable outcomes. The value of "step size" can be changed to change the frequency of variation. This method's formula is displayed below.

for  $E > S$ :

$$\alpha_E = \alpha_{E-1} + (\alpha_{max} - \alpha_{base})X(-1)^{\left(\frac{E}{S}+1\right)}$$

for  $E \leq S$ :

$$\alpha_E = \alpha_{max} - \frac{\alpha_{max} - \alpha_{base}}{S} X (S - E)$$

In the formulae above, "max" and "base" stand for the maximum and base learning rates, respectively, and E represents the learning rates for a specific epoch. There is a step size, S. Finding the ideal base and maximal learning rates is a crucial first step in making this strategy effective. The "LR range test" method is used to find these; it involves training the model for a few epochs while allowing the learning rate to vary linearly from a modest initial value. The accuracy of the model is then captured for various learning rates and plotted. Determine the two learning rate values—1) the point at which accuracy starts to rise, and 2) the point at which accuracy starts to fall or fluctuate—from the plot. The base learning rate is represented by the first point, while the maximal learning rate is represented by the second.

IV. MOMENTUM-BASED LEARNING

Momentum-based approaches understand that zigzagging is caused by strongly opposing actions that cancel each other out and lessen the actual magnitude of the steps taken in the right (long-term) direction. An example of this scenario is illustrated in Figure 1.

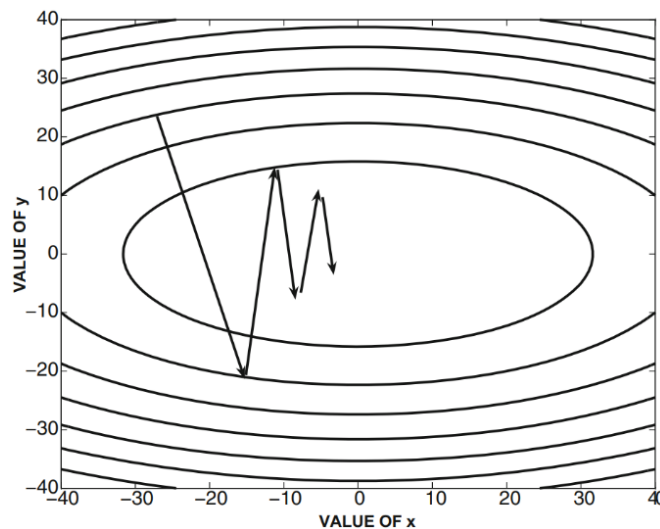


Fig. 1 Loss function is elliptical in bowl

$$L = x^2 + 4y^2$$

It is possible that even trying to make the step bigger will result in less movement in the right direction and push the existing solution even farther away from the ideal one. From this vantage point, it makes far more sense to proceed in the direction of the most recent "averaged" step to smooth out the zigzagging.

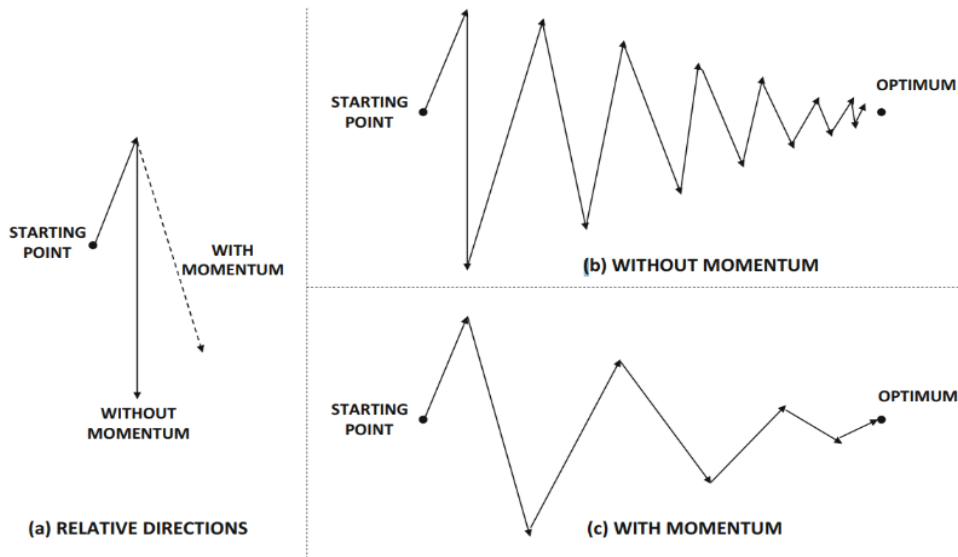


Fig. 2 Effect of momentum in smoothing zigzag updates

Examine a situation where gradient descent is being performed with respect to the parameter vector  $\bar{W}$  in order to comprehend this point. The following are the normal updates for gradient-descent with regard to loss function L (defined over a mini-batch of instances):

$$\bar{v} \leftarrow -\alpha \frac{\partial L}{\partial \bar{W}}; \bar{W} \leftarrow \bar{W} + \bar{v}$$

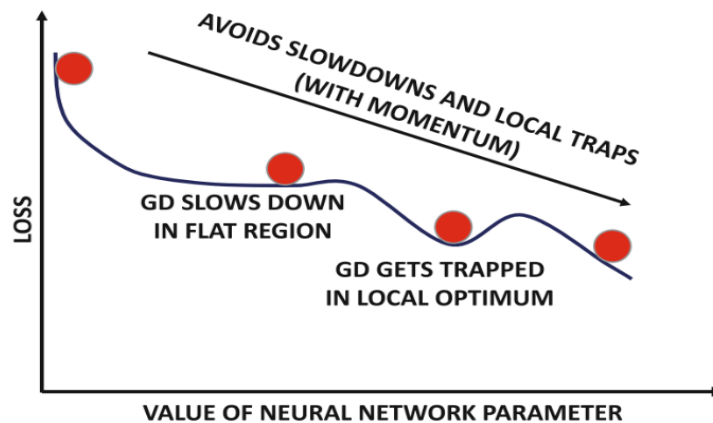


Fig. 3 Effect of momentum in navigating complex loss surfaces

The learning rate is represented by  $\alpha$ , and the annotation "GD" denotes pure gradient descent without momentum. In flat areas of the loss surface, momentum aids in maintaining optimization speed and preventing local optima.

Momentum-based descent accelerates learning by favoring consistent directions toward the optimal solution while minimizing ineffective oscillations. This approach prioritizes steady directions across multiple steps, allowing for larger strides in the right direction without causing issues in other directions. This acceleration is evident in Figure 2(a), showcasing increased gradient components in the correct direction. Figures 2(b) and (c) demonstrate how momentum-driven updates lead to quicker convergence to the optimal solution with fewer iterations. Momentum often causes a slight overshoot in the direction of velocity increase, akin to a marble rolling down a bowl. This overshooting effect actually aids the momentum-based strategy by accelerating progress toward the ideal solution.

Despite the overshoot, this approach outperforms by reaching the optimal solution quicker, compensating for any temporary deviation. Additionally, this overshooting tendency is beneficial as it helps avoid getting trapped in local optima.

## V. PARAMETER-SPECIFIC LEARNING RATES

The fundamental concept of the momentum methods discussed in the previous section is to take advantage of the gradient direction consistency of specific parameters to accelerate the updates. A more explicit way to accomplish this goal would be to have various learning rates for various factors. The concept is that whereas parameters with small partial derivatives tend to be more consistent yet move in the same direction, parameters with big partial derivatives frequently oscillate and zigzag. The delta-bar-delta method was one of the first approaches to be proposed in this direction. This method monitors whether each partial derivative's sign changes or remains constant. A partial derivative is a good indicator that the direction is accurate if its sign remains constant. The partial derivative increases in that direction in such a scenario. Alternatively, if the partial derivative's sign is always flipped, the partial derivative will decrease. Nevertheless, because the faults in stochastic gradient descent can amplify, this type of technique is intended for gradient descent rather than stochastic gradient descent. As a result, several techniques that can function effectively even when the mini-batch approach is applied have been put forth.

### A. *AdaGrad*

Over the course of the AdaGrad algorithm, the total squared magnitude of the partial derivative with respect to each parameter is recorded. Though the absolute value will rise with the number of epochs due to successive aggregation, the square-root of this value is proportionate to the parameter's root-mean-square slope.

Let  $A_i$  be the aggregate value for the  $i^{th}$  parameter. Therefore, in each iteration, the following update is performed:

$$A_i \leftarrow A_i + \left(\frac{\partial L}{\partial w_i}\right)^2; \forall i$$

The following is the update for the  $i^{th}$  parameter  $w_i$ :

$$w_i \leftarrow w_i - \frac{\alpha}{\sqrt{A_i}} \left(\frac{\partial L}{\partial w_i}\right); \forall i$$

Scaling the derivative inversely with  $A_i$  is a kind of "signal-to-noise" normalization because  $A_i$  only measures the historical magnitude of the gradient rather than its sign; it encourages faster relative movements along gently sloping directions with consistent sign of the gradient. If the gradient component along the  $e^{i^{th}}$  direction keeps wildly fluctuating between +100 and -100, this type of magnitude-centric normalization will penalize that component far more than another gradient component that consistently takes on the value in the vicinity of 0.1 (but with a consistent sign).

The fundamental issue with the strategy, though, is that absolute motions along all components would tend to slow down over time. The slowdown can be attributed to the fact that  $A_i$  represents the total value of all partial derivatives throughout history. The aggregate scaling factors' reliance on old history, which might eventually grow stale, is another issue.

### B. *RMSProp*

Similar to AdaGrad, the RMSProp technique employs the absolute magnitude  $A_i$  of the gradients to do "signal-to-noise" normalization. To estimate  $A_i$ , however, exponential averaging is used rather than just adding the squared gradients. The progress is not prematurely hindered by a continuously growing scaling factor  $A_i$ , since averaging is used to normalize rather than aggregating values. Using a decay factor  $\rho \in (0, 1)$ , the main idea is to weight the squared partial derivatives occurring  $t$  updates ago by  $\rho^t$ . The running estimate is initialized to 0. Early iterations experience some (undesired) bias as a result, but this eventually goes away. Therefore, if  $A_i$  is the exponentially averaged value of the  $i^{th}$  parameter  $w_i$ , we have the following way of updating  $A_i$ :

$$A_i \leftarrow \rho A_i + (1 - \rho) \left(\frac{\partial L}{\partial w_i}\right)^2; \forall i$$



The gradient of each parameter is normalized by taking its square root. Next, the (global) learning rate  $\alpha$  is updated using the following formula:

$$w_i \leftarrow w_i - \frac{\alpha}{\sqrt{A_i}} \left( \frac{\partial L}{\partial w_i} \right); \forall i$$

An additional benefit of RMSProp in comparison to AdaGrad is the exponential decay of the significance of old, or stale, gradients over time. The disadvantage of RMSProp is that because it is initialized to 0, the running estimate  $A_i$  of the second-order moment is biased in the early iterations.

#### C. *RMSProp with Nesterov Momentum*

Nesterov momentum takes the future position of the parameters into account to maximize the gradient descent. In addition to calculating the gradient at the current position, it also does it at an adjusted position that accounts for movement induced by momentum. In this hybrid technique, the learning rates for each parameter are adjusted using RMSProp based on the squared gradients, and Nesterov momentum improves the optimization process by taking momentum into account to predict future parameter updates. Combining RMSProp with Nesterov momentum leverages each of their unique advantages: Nesterov momentum provides more informed parameter updates, while RMSProp offers adjustable learning rates.

#### D. *AdaDelta*

The AdaDelta approach computes the update as a function of incremental updates from earlier iterations, doing away with the requirement for a global learning parameter. This update is similar to that used by RMSProp. Examine the RMSProp update, which is repeated below:

$$w_i \leftarrow w_i - \frac{\alpha}{\sqrt{A_i}} \left( \frac{\partial L}{\partial w_i} \right); \forall i$$

here,  $\frac{\alpha}{\sqrt{A_i}} \left( \frac{\partial L}{\partial w_i} \right) \sim \Delta w_i$

We will demonstrate how  $\alpha$  is changed to a value that is based on earlier incremental adjustments. The increment in the value of  $w_i$  is represented by the value of  $\Delta w_i$  in each update. We maintain an exponentially smoothed value  $\delta_i$  of the values of  $\Delta w_i$  from earlier iterations using the same decay parameter  $\rho$ , just like with the exponentially smoothed gradients  $A_i$ .

$$\delta_i \leftarrow \rho \delta_i + (1 - \rho)(\Delta w_i)^2; \forall i$$

Since the value of  $\Delta w_i$  is not yet known, the value of  $\delta_i$  for a given iteration can only be calculated using the iterations that came before it. However, in the current iteration,  $A_i$  can also be calculated using the partial derivative. There is a slight distinction in the calculation of  $A_i$  and  $\delta_i$ . As a result, the AdaDelta update is can be obtained by replacing:

$$\Delta w_i \sim \sqrt{\frac{\delta_i}{A_i}} \left( \frac{\partial L}{\partial w_i} \right)$$

Notably, this update does not include a single  $\alpha$  parameter for the learning rate.

#### E. *Adam*

The Adam approach incorporates momentum into the update by exponentially smoothing the first-order gradient in addition to using a similar "signal-to-noise" normalization as AdaGrad and RMSProp. Additionally, it directly resolves the bias that arises from starting a smoothed value's running estimate at zero, an impractical starting point for exponential smoothing. Let  $A_i$  be the exponentially averaged value of the  $i^{th}$  parameter  $w_i$ , just as in the case of RMSProp. This value is updated using the decay parameter  $\rho \in (0, 1)$  in the same manner as RMSProp:

$$A_i \leftarrow \rho A_i + (1 - \rho) \left( \frac{\partial L}{\partial w_i} \right)^2; \forall i$$

$F_i$  represents the  $i^{th}$  component of the gradient, which is maintained at the same time as an exponentially smoothed value of the gradient. This smoothing is done using an alternative decay parameter,  $\rho_f$ :

$$F_i = \rho_f F_i + (1 - \rho_f) \left( \frac{\partial L}{\partial w_i} \right); \forall i$$

This particular kind of gradient smoothing with  $\rho_f$  is an adaptation of the momentum method. Next, at the  $t^{th}$  iteration, the following update is applied at learning rate  $\alpha_t$ :

$$w_i \leftarrow w_i - \frac{\alpha_t}{\sqrt{A_i}} (F_i); \forall i$$

Compared to the RMSProp algorithm, there are two main distinctions. To include momentum, the gradient is first substituted by its exponentially smoothed value. Secondly, the learning rate  $\alpha_t$ , which is defined as follows, is now dependent on the iteration index  $t$ .

$$\alpha_t = \alpha \left( \frac{\sqrt{1 - \rho^t}}{1 - \rho_f^t} \right)$$

In actuality, the learning rate adjustment is essentially a bias correction factor that is used to correct for the two exponential smoothing processes' unrealistic initialization. This adjustment is especially significant in the initial repetitions. Early iterations are biased because  $F_i$  and  $A_i$  are both initialized at 0. Notably, since  $\rho, \rho_f \in (0, 1)$ , each of  $\rho^t$  and  $\rho_f^t$  converges to 0 for large  $t$ . Consequently, initialization bias correction factor converges to 1, and  $\alpha_t$  converges to  $\alpha$ . The original Adam study suggests that the default values of  $\rho_f$  and  $\rho$  are 0.9 and 0.999, respectively. Because it combines most of the benefits of other algorithms and frequently outperforms the best of the competition, the Adam algorithm is incredibly popular.

## VI. GRADIENT CLIPPING

A method for handling situations when the partial derivatives in several directions have wildly disparate magnitudes is called gradient clipping. Some gradient clipping techniques aim to make the various components of the partial derivatives more equal, which is a concept related to adaptive learning rates. But rather than using the gradients' historical values, the clipping is just based on their current values. The most popular types of gradient clipping are two:

### A. Value-based clipping

A minimum and maximum threshold are established on the gradient values in value-based clipping. A partial derivative is set to the minimal threshold if it is less than the minimum. The maximum threshold is applied to all partial derivatives that exceed the maximum.

### B. Norm-based clipping:

In this instance, the L2-norm of the entire vector normalizes the gradient vector as a whole. It should be noted that the relative magnitudes of the updates along different directions remain unchanged by this form of clipping. The effects of the two types of clipping, however, are extremely comparable for neural networks (like recurrent neural networks) that exchange parameters across layers. By clipping, the values can be more effectively conditioned, resulting in updates that are generally comparable amongst mini-batches. Consequently, it would stop an abnormal gradient explosion in a specific mini-batch from having a significant impact on the solution.

In general, gradient clipping has much fewer impacts than many other techniques. But it works especially well to prevent the exploding gradient issue that arises with recurrent neural networks.



## VII. CONCLUSION

The research outlined in this paper delves deeply into the optimization landscape of machine learning through the lens of gradient-based algorithms, shedding light on the intricate relationship between learning rates, convergence, and parameter updates. The comprehensive exploration of optimization strategies, particularly focusing on gradient descent and its variants, provides a nuanced understanding of how these techniques influence the efficiency and effectiveness of machine learning models. By dissecting fundamental concepts like learning rates and adaptive learning methods including decay, drop, cyclic learning, and adaptive learning, this study underscores their pivotal roles in navigating the trade-offs between convergence speed and accuracy. The examination of momentum-based approaches, such as Adam, RMSProp, AdaGrad, and AdaDelta, underscores their significance in addressing challenges like oscillations and slow convergence, highlighting their versatility and effectiveness across various scenarios.

This paper not only elucidates the theoretical underpinnings and mathematical foundations of these optimization strategies but also emphasizes their practical implications in real-world machine learning applications. The discussion on gradient clipping methods provides valuable insights into handling issues like exploding gradients, contributing to stable training and improved model performance, particularly in recurrent neural networks. The study accentuates the significance of selecting appropriate optimization techniques and fine-tuning learning rates, offering a roadmap for practitioners to navigate the complex landscape of optimization strategies based on dataset characteristics and model architectures. Overall, the insights gleaned from this research equip both researchers and practitioners with a deeper understanding and a diverse toolkit to optimize machine learning models, paving the way for more efficient, accurate, and robust algorithms in diverse application domains.

## REFERENCES

- [1]. S. Ruder, "An overview of gradient descent optimization algorithms," arXiv preprint arXiv:1609.04747, 2016.
- [2]. John Duchi, Elad Hazan, and Yoram Singer. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [3]. Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M. W., Pfau, D., Schaul, T., and de Freitas, N. Learning to learn by gradient descent by gradient descent. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 29*, pp. 3981–3989. Curran Associates, Inc., 2016.
- [4]. Kingma, D. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [5]. J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [6]. N. N. Schraudolph. Local gain adaptation in stochastic gradient descent. In *International Conference on Artificial Neural Networks*, volume 2, pages 569–574, 1999.
- [7]. F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4(1):1–106, 2012.
- [8]. C. Darken, J. Chang, and J. Moody, "Learning rate schedules for faster stochastic gradient search," in *Neural Networks for Signal Processing [1992] II*, Proceedings of the 1992 IEEE-SP Workshop. IEEE, 1992, pp. 3–12.
- [9]. Arvind Neelakantan, Luke Vilnis, Quoc V. Le, Ilya Sutskever, Lukasz Kaiser, Karol Kurach, and James Martens. Adding Gradient Noise Improves Learning for Very Deep Networks. pages 1–11, 2015.
- [10]. Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks: the official journal of the International Neural Network Society*, 12(1):145–151, 1999.