



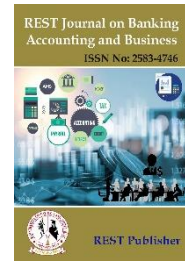
REST Journal on Banking, Accounting and Business

Vol: 2(2), June 2023

REST Publisher; ISSN: 2583 4746 (Online)

Website: <http://restpublisher.com/book-series/jbab/>

DOI: <https://doi.org/10.46632/jbab/2/2/7>



Unleashing the Power of Gradient Descent: Dive into the World of Optimization Algorithms

Vishal Triloknath Jaiswar

S.S.T College of Arts and Commerce, Mumbai, Maharashtra, India.

Corresponding Author Email: vishal.mit21009@sstcollege.edu.in

Abstract: While gradient descent optimization algorithms have gained immense popularity, they are often treated as mysterious black-box optimizers due to the scarcity of practical explanations about their strengths and weaknesses. This article endeavors to equip readers with intuitive insights into the behavior of various algorithms, enabling them to harness their potential. Throughout this comprehensive overview, we explore diverse variants of gradient descent, address challenges, introduce prominent optimization algorithms, delve into parallel and distributed architectures, and explore additional strategies to optimize gradient descent. Prepare to unravel the enigmatic realm of gradient descent and unleash its true power.

Keywords: *Gradient descent Variants, Unveiling the Obstacles, Unveiling the Arsenal.*

1. INTRODUCTION

Among the vast array of optimization algorithms, gradient descent reigns supreme as one of the most widely embraced approaches, particularly in the realm of neural network optimization. Curiously, even cutting-edge Deep Learning libraries such as lasagne's², caffe's³, and keras⁴ include a repertoire of gradient descent optimization algorithm implementations. However, these algorithms often remain shrouded in mystery, used as enigmatic black-box optimizers, with scarce practical explanations highlighting their merits and limitations. This article endeavors to demystify these optimization algorithms, bestowing the reader with invaluable intuitions that shed light on their behavior. In Section 2, we embark on a captivating exploration of the diverse variants of gradient descent. Delving deeper, Section 3 succinctly encapsulates the challenges encountered during training. Building upon this foundation, Section 4 unveils the most prominent optimization algorithms, uncovering their intrinsic motivations in addressing these challenges and unveiling the derivation of their ingenious update rules. Venturing further into the realm of optimization, Section 5 unveils a brief yet illuminating glimpse into the fascinating domain of parallel and distributed architectures for optimizing gradient descent. Finally, Section 6 delves into additional strategies that prove instrumental in fine-tuning the performance of gradient descent. At its core, gradient descent emerges as an indispensable technique for minimizing an objective function $J(\theta)$, wherein the model's parameters $\theta \in \mathbb{R}^d$ serve as vital components. This is achieved by iteratively adjusting the parameters in the opposite direction of the gradient $\nabla_{\theta} J(\theta)$ of the objective function with respect to the parameters. The learning rate η acts as the compass, dictating the magnitude of steps taken to traverse towards a (local) minimum. In essence, we follow the path guided by the slope of the objective function's terrain, descending into valleys until our destination is reached.⁵ Get ready to embark on an enthralling expedition into the inner workings of optimization, unraveling the enigmatic intricacies of gradient descent.

2. GRADIENT DESCENT VARIANTS

Unveiling the Trio: Exploring the Fascinating Variants of Gradient Descent Within the realm of gradient descent, we encounter three captivating variants, each distinguished by the amount of data employed to compute the gradient of the objective function. This trade-off between parameter update accuracy and computational efficiency has far-reaching implications.

Journeying through Batch Gradient Descent: The first variant, known as batch gradient descent or vanilla gradient descent, embraces the entirety of the training dataset when computing the gradient of the cost function with respect to the parameters θ : $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$ (1) This approach demands the computation of gradients for the entire

dataset, leading to potential sluggishness and intractability when working with memory-intensive datasets. Additionally, batch gradient descent lacks the capability to update the model online, inhibiting real-time adaptation. Implementation-wise, the code for batch gradient descent manifests as follows: for i in range (nb epochs): `params_grad = evaluate_gradient(loss_function, data, prams)` `prams = prams - learning_rate * params_grad` During a predetermined number of epochs, we calculate the gradient vector `params_grad` of the loss function for the entire dataset, with cutting-edge deep learning libraries offering efficient automatic differentiation. If manually deriving gradients, it is advisable to perform gradient checking. Parameters are subsequently updated in the direction of the gradients, with the learning rate determining the magnitude of the update. Batch gradient descent guarantees convergence to the global minimum for convex error surfaces and to a local minimum for non-convex surfaces.

Embracing Stochastic Gradient Descent: In stark contrast, stochastic gradient descent (SGD) opts for a different approach by performing parameter updates for each training example ($x(i)$) and label ($y(i)$): $\theta = \theta - \eta \cdot \nabla \theta J(\theta; x(i); y(i))$ (2) Batch gradient descent incurs redundant computations when dealing with large datasets, as gradients for similar examples are recomputed before each parameter update. SGD eradicates this redundancy by updating parameters one example at a time, rendering it significantly faster and enabling online learning. However, SGD's frequent updates introduce high variance, causing the objective function to fluctuate intensely, as depicted in Figure 1. While batch gradient descent converges to the minimum of the parameter space, SGD's fluctuation grants it the ability to leap towards new and potentially superior local minima. Nonetheless, this very characteristic complicates the quest for the exact minimum, as SGD tends to overshoot. Notwithstanding, research has demonstrated that gradually decreasing the learning rate aligns SGD's convergence behavior with that of batch gradient descent, allowing it to almost certainly converge to a local or global minimum for non-convex and convex optimization, respectively. The code snippet below showcases the implementation of SGD, with the training data shuffled at each epoch, as elucidated in Section 6.1. for i in range (nb_epochs): `np.random.Shuffle(data)` for example in `data`: `params_grad = evaluate_gradient(loss_function, example, prams)` `prams = prams - learning_rate * params_grad` Prepare to traverse the captivating landscapes of gradient descent, where these three variants offer unique journeys towards optimization excellence.

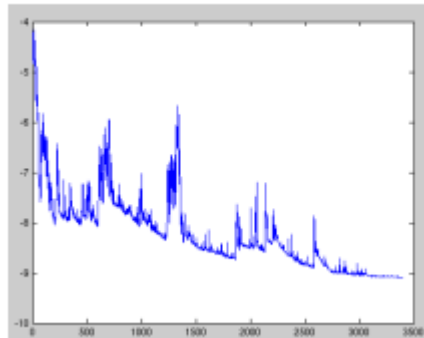


FIGURE 1. SGD fluctuation (Source: Wikipedia)

Embarking on a Harmonious Journey: Mini-Batch Gradient Descent: Finally, we encounter the convergence of two worlds, as mini-batch gradient descent emerges to seize the best of both realms. By performing updates for every mini-batch consisting of n training examples, this variant exhibits the following elegant equation: $\theta = \theta - \eta \cdot \nabla \theta J(\theta; x(i:i+n); y(i:i+n))$ (3) This approach delivers a twofold advantage: a) it diminishes the variance of parameter updates, fostering more stable convergence; and b) it capitalizes on highly optimized matrix operations prevalent in state-of-the-art deep learning libraries, resulting in efficient computation of gradients with respect to a mini-batch. Mini-batch sizes commonly range between 50 and 256, although specific applications may warrant variations. Notably, when working with neural networks, mini-batch gradient descent typically assumes the role of the algorithm of choice, with the term SGD often employed even when mini-batches are utilized. Please note that, for the sake of simplicity, the parameters $x(i:i+n)$ and $y(i:i+n)$ are omitted in subsequent modifications of SGD throughout this discourse. In the realm of code, a transformative shift takes place. Instead of iterating over individual examples, we now embark on a journey through mini-batches, each comprising 50 instances: for i in range(nb_epochs): `np.random.shuffle(data)` for batch in `get_batches(data, batch_size=50)`: `params_grad = evaluate_gradient(loss_function, batch, prams)` `prams = prams - learning_rate * params_grad` Prepare to be captivated as mini-batch gradient descent illuminates the path towards optimization excellence, harmonizing the strengths of its predecessors into a formidable force.

3. UNVEILING THE OBSTACLES

While vanilla mini-batch gradient descent presents itself as a promising contender, it brings forth a set of challenges that necessitate careful consideration: Selecting an optimal learning rate proves to be a formidable task. A learning rate that is too minuscule yields agonizingly slow convergence, whereas a rate that is too substantial can impede convergence and cause the loss function to oscillate around the minimum or even diverge. Learning rate schedules [18] emerge as a potential solution, aiming to adapt the learning rate throughout training. Techniques such as annealing involve reducing the learning rate according to a predetermined schedule or when the objective change between epochs falls below a threshold. However, these fixed schedules and thresholds fail to adapt to the unique characteristics of each dataset [4]. Furthermore, a single learning rate governs all parameter updates, which may not be suitable for scenarios where data is sparse and features exhibit varying frequencies. In such instances, it may be beneficial to apply more substantial updates to infrequent features, rather than treating all features equally. One of the paramount challenges in the quest to minimize highly non-convex error functions prevalent in neural networks is the perils of becoming ensnared within suboptimal local minima. Surprisingly, Dauphin et al. [5] argue that the true difficulty lies not in local minima but rather in saddle points—points where one dimension ascends while another descends. These saddle points are often encircled by a plateau of comparable error, rendering it incredibly arduous for SGD to break free, as the gradient approaches zero across all dimensions. Embark on this intrepid journey as we navigate through the hurdles that arise in the realm of mini-batch gradient descent, uncovering ingenious solutions to overcome these challenges and forge ahead toward optimization excellence.

4. UNVEILING THE ARSENAL: GRADIENT DESCENT OPTIMIZATION ALGORITHMS

In the realm of Deep Learning, the community has embraced various algorithms to combat the challenges discussed earlier. However, we will focus solely on algorithms that are feasible to compute in practical scenarios for high-dimensional datasets, excluding infeasible options like second-order methods such as Newton's method [7].

The Momentum Revolution: When navigating through treacherous ravines, where the surface sharply curves in one dimension while remaining relatively flat in another [20], standard SGD encounters significant difficulties. It oscillates along the slopes of these ravines, making hesitant progress toward local optima, as depicted in Figure 2a.



FIGURE 2. Source: Genevieve B. Orr

Enter Momentum [17], a method that injects a dose of acceleration into SGD, mitigating oscillations and boosting progress in the desired direction, as illustrated in Figure 2b. Momentum achieves this by incorporating a fraction γ of the update vector from the previous time step into the current update vector [8]: $v_t = \gamma v_{t-1} + \eta \nabla \theta J(\theta)$ $\theta = \theta - v_t$ (4) Typically, the momentum term γ is set to 0.9 or a similar value. In essence, when utilizing momentum, we unleash a rolling ball on a hill. As the ball descends, it accumulates momentum, gradually increasing its velocity (unless air resistance intervenes, i.e., $\gamma < 1$). A similar phenomenon occurs with our parameter updates: the momentum term amplifies for dimensions where gradients align and dampens updates for dimensions where gradients change directions. The outcome? Expedited convergence and reduced oscillation, propelling us towards optimization success.

The Intelligent Ball: Nesterov Accelerated Gradient: Merely rolling down a hill without foresight leaves much to be desired. What if we could imbue our ball with intelligence, enabling it to anticipate and slow down before encountering an uphill slope? Nesterov Accelerated Gradient (NAG) [14] grants our momentum term this remarkable prescience. We know that we employ the momentum term, γv_{t-1} , to propel our parameters, θ . Consequently, computing $\theta - \gamma v_{t-1}$ approximates the forthcoming parameter position (lacking the gradient for a full update), offering a glimpse into our parameters' trajectory. Now, we can effectively peer into the future by

calculating the gradient not with respect to our current parameters, θ , but with respect to the estimated future position of our parameters: $v_t = \gamma v_{t-1} + \eta \nabla \theta J(\theta - \gamma v_{t-1})$ $\theta = \theta - v_t$

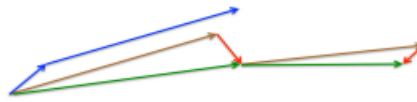


FIGURE 3. Nesterov update (Source: G. Hinton s lecture 6c)

Once again, we assign a momentum term, γ , with a value around 0.9. While Momentum first computes the current gradient (small blue vector in Figure 3) and then takes a substantial leap in the direction of the updated accumulated gradient (big blue vector), NAG takes a daring leap towards the previous accumulated gradient's direction (brown vector), measures the gradient, and then adjusts its course accordingly (green vector). This anticipatory update prevents excessive speed and enhances responsiveness, significantly boosting the performance of Recurrent Neural Networks (RNNs) across various tasks [2]. With the ability to adapt our updates to the error function's slope and accelerate SGD, we yearn for the capability to tailor updates to individual parameters, allowing for larger or smaller adjustments based on their significance.

Embracing Parameter Diversity: Adagrad: Adagrad [8] presents an algorithm for gradient-based optimization that tackles this very challenge: it dynamically adjusts the learning rate for each parameter, enabling more substantial updates for infrequent ones and smaller updates for frequent ones. Consequently, Adagrad proves particularly effective when handling sparse data. At Google, Dean et al. [6] discovered that Adagrad significantly bolstered the resilience of SGD, successfully employing it to train expansive neural networks capable of recognizing cats in YouTube videos [10]. Furthermore, Pennington et al. [16] harnessed Adagrad to train GloVe word embeddings, where infrequent words demanded more pronounced updates than their frequent counterparts. Previously, we performed simultaneous updates for all parameters θ , employing a uniform learning rate η for each parameter θ_i . However, Adagrad departs from this approach by assigning a distinct learning rate to every parameter θ_i at each time step t . To offer a concise illustration, let's denote $g_{t,i}$ as the gradient of the objective function with respect to parameter θ_i at time step t : $g_{t,i} = \nabla_{\theta_i} J(\theta_{t,i})$ (6) Consequently, the SGD update for each parameter θ_i becomes: $\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$ (7) In the update rule of Adagrad, the general learning rate η at time step t for parameter θ_i is adjusted based on the accumulated past gradients computed for θ_i : $\theta_{t+1,i} = \theta_{t,i} - \eta \sqrt{G_{t,i}} \cdot g_{t,i}$ (8) Here, $G_t \in \mathbb{R}^{d \times d}$ represents a diagonal matrix where each diagonal element $G_{t,ii}$ is the sum of the squares of the gradients with respect to θ_i up to time step t . Additionally, ϵ is a smoothing term that prevents division by zero (typically on the order of $1e-8$). Intriguingly, omitting the square root operation severely compromises the algorithm's performance. By conducting an element-wise matrix-vector multiplication between G_t and g_t , we can now leverage vectorization to streamline our implementation: $\theta_{t+1} = \theta_t - \eta \sqrt{G_t + \epsilon} \cdot g_t$ (9) One of the primary advantages of Adagrad is its ability to eliminate the manual tuning of the learning rate. Most implementations utilize a default value of 0.01 and maintain it throughout the process. Nevertheless, Adagrad does have a drawback: the accumulation of squared gradients in the denominator. As each added term is positive, the cumulative sum grows incessantly during training. Consequently, the learning rate diminishes to infinitesimal levels, rendering the algorithm incapable of acquiring further knowledge. The subsequent algorithms aim to address this limitation.

Redefining the Learning Rate: Adelta: Adadelta [22] emerges as an extension of Adagrad, aiming to mitigate its aggressive and monotonically decreasing learning rate. Unlike Adagrad, which accumulates all past squared gradients, Adadelta limits the accumulation window to a fixed size, denoted as w . Instead of storing w previous squared gradients, Adadelta employs a decaying average of all past squared gradients as the sum of gradients. The running average $E[g^2]_t$ at time step t is computed as a fraction γ of the previous average combined with $(1 - \gamma)$ times the current gradient squared: $E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g^2_t$ Here, γ assumes a value similar to the momentum term, typically around 0.9. To enhance clarity, let's redefine our vanilla SGD update in terms of the parameter update vector $\Delta\theta_t$: $\Delta\theta_t = -\eta \cdot g_t$ $\theta_{t+1} = \theta_t + \Delta\theta_t$ The parameter update vector derived for Adagrad can now be expressed as: $\Delta\theta_t = -\eta \sqrt{G_t + \epsilon} \cdot g_t$ To align with Adadelta, we replace the diagonal matrix G_t with the decaying average of past squared gradients $E[g^2]_t$: $\Delta\theta_t = -\eta \sqrt{E[g^2]_t + \epsilon} \cdot g_t$ As the denominator corresponds to the root mean squared (RMS) error criterion of the gradient, we can utilize a shorthand for the criterion: $\Delta\theta_t = -\eta \text{RMS}[g]_t \cdot g_t$ The authors highlight that the units in this update, as well as in SGD, Momentum, or Adagrad, do not align with the parameter units. To address this, they introduce another exponentially decaying average, not of squared gradients, but of squared parameter updates: $E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta^2_t$ Consequently, the root mean squared error of parameter updates becomes: $\text{RMS}[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$ As the value of $\text{RMS}[\Delta\theta]_t$ is

unknown, we approximate it with the RMS of parameter updates until the previous time step. By replacing the learning rate η in the previous update rule with $\text{RMS}[\Delta\theta]_{t-1}$, we arrive at the Adadelta update rule: $\Delta\theta_t = -\text{RMS}[\Delta\theta]_{t-1} \text{RMS}[g]_t$ $\theta_{t+1} = \theta_t + \Delta\theta_t$ Adadelta eliminates the need for a default learning rate, as it has been seamlessly integrated into the update rule.

Unveiling an Adaptive Learning Approach: RM Sprop: RMSprop, an adaptive learning rate technique, was introduced by Geoff Hinton during Lecture 6e of his Coursera Class12. It emerged independently around the same time as Adadelta, both aiming to address the diminishing learning rates experienced in Adagrad. Remarkably, RMSprop's first update vector is identical to the one derived for Adadelta: $E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g^2_t$ $\theta_{t+1} = \theta_t - \eta \text{p} E[g^2]_t + g_t$ Similar to Adadelta, RMSprop also divides the learning rate by an exponentially decaying average of squared gradients. Hinton suggests setting γ to 0.9, while a recommended default value for the learning rate η is 0.001.

Unleashing the Power of Adaptive Moment Estimation: Adam, short for Adaptive Moment Estimation, is a remarkable technique [10] that computes adaptive learning rates for individual parameters. It goes beyond the capabilities of Adadelta and RMSprop by incorporating both an exponentially decaying average of past squared gradients (v_t) and an exponentially decaying average of past gradients (m_t), similar to momentum: $m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t$ $v_t = \beta_2 v_{t-1} + (1 - \beta_2)g^2_t$ (19) The parameters m_t and v_t represent estimations of the first moment (mean) and second moment (uncentered variance) of the gradients, respectively, hence the name of the method. Since m_t and v_t are initialized as vectors of zeros, the authors of Adam acknowledge that they tend to be biased towards zero, particularly during the initial time steps and when the decay rates (β_1 and β_2) are small. To address these biases, bias-corrected estimates of the first and second moments are computed: $\hat{m}_t = m_t / (1 - \beta_1^t)$ $\hat{v}_t = v_t / (1 - \beta_2^t)$ (20) These bias-corrected estimates are then utilized to update the parameters, following a similar approach as Ad delta and RM Sprop. This leads to the Adam update rule: $\theta_{t+1} = \theta_t - \eta / \sqrt{\hat{v}_t + \epsilon} \cdot \hat{m}_t$ (21) The authors propose default values of 0.9 for β_1 , 0.999 for β_2 , and 10^{-8} for ϵ . Empirical evidence demonstrates the effectiveness of Adam in practical scenarios, showcasing its superiority compared to other adaptive learning-method algorithms.

Breaking Boundaries with Ada Max: In the pursuit of further optimization, the Adam update rule introduced a factor, v_t , which scaled the gradient in relation to the 2 norm of past gradients (via the v_{t-1} term) and the current gradient's 2 norm, $|g_t|^2$: $v_t = \beta_2 v_{t-1} + (1 - \beta_2)|g_t|^2$ (22) To extend this update to the p norm, denoted by β_p^2 , we can introduce a generalization. However, norms with large p values often suffer from numerical instability, making 1 and 2 norms more prevalent in practical applications. Interestingly, ∞ norm exhibits remarkable stability. Thus, Ada Max [10] was proposed, demonstrating convergence of v_t with ∞ norm to a more stable value. To differentiate from Adam, the authors employ the variable u_t , which represents the infinity norm-constrained v_t : $u_t = \beta_\infty^2 v_{t-1} + (1 - \beta_\infty^2)|g_t|_\infty = \max(\beta_2 \cdot v_{t-1}, |g_t|)$ (24)By incorporating this into the Adam update equation and replacing $\sqrt{\hat{v}_t + \epsilon}$ with u_t , we unveil the Ada Max update rule: $\theta_{t+1} = \theta_t - \eta / u_t \cdot \hat{m}_t$ (25) Notably, as u_t utilizes the max operation, it avoids the bias towards zero observed in m_t and v_t of Adam. Consequently, a bias correction for u_t is unnecessary. For optimal results, recommended default values are $\eta = 0.002$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. Ada Max dares to transcend conventional boundaries and unlocks a new realm of optimization possibilities.

Embracing Momentum's Evolution: Nadam As we have previously explored, Adam is a fusion of RMSprop and momentum, where RMSprop contributes the exponentially decaying average of past squared gradients, v_t , and momentum factors in the exponentially decaying average of past gradients, m_t . Moreover, we discovered that Nesterov accelerated gradient (NAG) surpasses the effectiveness of vanilla momentum. In a quest to combine the strengths of Adam and NAG, Nadam (Nesterov-accelerated Adaptive Moment Estimation) [7] was conceived. To incorporate NAG into Adam, a modification of its momentum term, m_t , is required. Let us recall the momentum update rule using our current notation: $g_t = \nabla\theta_t J(\theta_t)$ $m_t = \gamma m_{t-1} + \eta g_t$ $\theta_{t+1} = \theta_t - m_t$ This equation once again highlights that momentum involves taking a step in the direction of the previous momentum vector and another step in the direction of the current gradient. NAG allows for a more accurate gradient step by updating the parameters with the momentum step before computing the gradient. By modifying the gradient g_t , we arrive at NAG: $g_t = \nabla\theta_t J(\theta_t - \gamma m_{t-1})$ $m_t = \gamma m_{t-1} + \eta g_t$ $\theta_{t+1} = \theta_t - m_t$ Dozat proposes a modification to NAG as follows: Rather than applying the momentum step twice, once to update the gradient g_t and a second time to update the parameters θ_{t+1} , we now directly apply the look-ahead momentum vector to update the current parameters: $g_t = \nabla\theta_t J(\theta_t)$ $m_t = \gamma m_{t-1} + \eta g_t$ $\theta_{t+1} = \theta_t - (\gamma m_t + \eta g_t)$ Notice that instead of using the previous momentum vector m_{t-1} as in Equation 27, we employ the current momentum vector m_t for looking ahead. To incorporate Nesterov momentum into Adam, we can similarly replace the previous momentum vector with the current momentum vector. Let's revisit the Adam update rule (note that we do not need to modify \hat{v}_t): $m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t$ $\hat{m}_t = m_t / (1 - \beta_1^t)$ $\theta_{t+1} = \theta_t - \eta / \sqrt{\hat{v}_t + \epsilon} \cdot \hat{m}_t$ Expanding the second equation using the definitions of \hat{m}_t and m_t , we obtain: $\theta_{t+1} = \theta_t - \eta / \sqrt{\hat{v}_t + \epsilon} \cdot (\beta_1 m_{t-1} / (1 - \beta_1^{t-1}) + (1 - \beta_1)g_t / (1 - \beta_1^t))$ (31) Notably, $\beta_1 m_{t-1} / (1 - \beta_1^{t-1})$ represents the

bias-corrected estimate of the momentum vector from the previous time step. Consequently, we can replace it with \hat{m}_{t-1} : $\theta_{t+1} = \theta_t - \eta \sqrt{\hat{v}_t} + (\beta_1 \hat{m}_{t-1} + (1 - \beta_1) g_t) / (1 - \beta_1^t)$ (32) This equation bears a striking resemblance to our expanded momentum term in Equation 27. By incorporating Nesterov momentum as we did in Equation 29, we simply replace the bias-corrected estimate of the momentum vector from the previous time step.

Illuminating Algorithmic Behaviors: In these captivating visuals (Figure 4a and Figure 4b), we gain valuable intuitions into the optimization behaviors exhibited by the presented optimization algorithms [13].

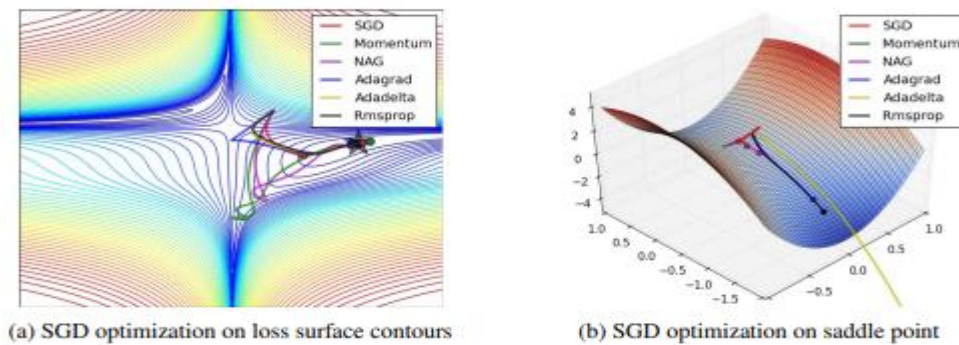


FIGURE 4. Source and full animations: Alec Radford

Figure 4a unveils the paths traversed by these algorithms on the contours of a loss surface, specifically the Beale function. Originating from the same starting point, they embark on distinct trajectories towards the minimum. Notably, Adagrad, Adadelta, and RMSprop promptly set off in the right direction, converging with similar swiftness. Meanwhile, Momentum and NAG encounter detours, akin to a ball rolling off course down a hill. Nevertheless, NAG swiftly corrects its path, leveraging its enhanced responsiveness by looking ahead, ultimately steering towards the minimum. Turning our attention to Figure 4b, we explore the algorithms' behavior when confronted with a saddle point—a point where one dimension exhibits a positive slope while the other dimension displays a negative slope—posing a challenge for SGD, as mentioned earlier. Here we observe that SGD, Momentum, and NAG grapple with breaking symmetry, although the latter two eventually triumph in escaping the saddle point. On the other hand, Adagrad, RMSprop, and Adadelta promptly descend down the negative slope, with Adadelta leading the charge. It becomes evident that adaptive learning-rate methods, namely Adagrad, Adadelta, RMSprop, and Adam, prove most adept in such scenarios, showcasing superior convergence capabilities. These captivating visualizations provide compelling evidence that adaptive learning-rate methods offer optimal convergence and effectively navigate complex optimization landscapes. *Deciding on the Optimal Optimizer: A Wise Choice:* Now, the question arises: Which optimizer should you employ? Consider the nature of your input data. If it leans towards sparsity, then it is highly likely that you will achieve optimal results by employing one of the adaptive learning-rate methods. An additional advantage lies in the fact that you won't need to fine-tune the learning rate, as the default value often yields excellent outcomes. To summarize, RMSprop emerges as an extension of Adagrad, addressing its issue of rapidly diminishing learning rates. Adadelta shares similarities with RMSprop, but instead of using the root mean square (RMS) of parameter updates in its numerator update rule, it utilizes the RMS of gradients. Finally, Adam incorporates bias-correction and momentum into the framework of RMSprop. Thus far, RMSprop, Adadelta, and Adam showcase comparable performance in similar scenarios. Kingma et al. [10] shed light on Adam's slight advantage over RMSprop, particularly in later stages of optimization when gradients become sparser, thanks to its bias-correction mechanism. Consequently, Adam could potentially be the optimal choice overall. Interestingly, recent studies reveal a preference for vanilla SGD (Stochastic Gradient Descent) without momentum and a straightforward learning rate annealing schedule. While SGD often manages to find a minimum, it may require significantly more time compared to some of the other optimizers. It heavily relies on a robust initialization and annealing schedule and may encounter challenges in escaping saddle points instead of local minima. Hence, if your priority lies in fast convergence and you are training a deep or complex neural network, selecting one of the adaptive learning rate methods would be wise. The key to making an informed decision lies in understanding your data's characteristics and aligning them with the strengths and capabilities of the various optimization algorithms at your disposal

5. REIMAGINING PARALLELIZATION AND DISTRIBUTION OF SGD: ADVANCEMENTS UNVEILED

Considering the widespread adoption of large-scale data solutions and the accessibility of low-cost clusters, it is evident that distributing SGD to further expedite its progress is a logical choice. SGD, in its inherent form, operates sequentially—step by step, inching closer to the minimum. While it delivers satisfactory convergence, it can prove sluggish, especially when dealing with vast datasets. In contrast, asynchronous execution of SGD accelerates the process but often suffers from suboptimal communication between workers, resulting in compromised convergence. Additionally, we can also parallelize SGD on a single machine without necessitating a massive computing cluster. In this section, we explore the algorithms and architectures proposed to optimize the parallelized and distributed versions of SGD. *Introducing Hogwild!*: In their work, Niu et al. [15] introduce a parallel update scheme called Hogwild!, enabling simultaneous SGD updates on CPUs. The scheme allows processors to access shared memory without the need for parameter locking. This approach is effective for sparse input data, where each update affects only a fraction of the parameters. The authors demonstrate that Hogwild! achieves near-optimal convergence rates in such cases, as the likelihood of processors overwriting vital information is minimal. *Unveiling Downpour SGD*: Downpour SGD, an asynchronous variant of SGD, was employed by Dean et al. [6] in their Dist Belief framework (predecessor to Tensor Flow) at Google. It involves running multiple replicas of a model in parallel, each processing a subset of the training data. These replicas communicate their updates to a parameter server, distributed across multiple machines. Each machine is responsible for storing and updating a fraction of the model's parameters. However, since the replicas do not engage in direct communication or share weights and updates, their parameters are constantly at risk of diverging, impeding convergence. *Delay-Tolerant Algorithms for SGD*: McMahan and Streeter [12] extend Ada Grad to the parallel setting by introducing delay-tolerant algorithms that adapt not only to past gradients but also to the delays in updates. This approach has proven to be effective in practical scenarios. *The Power of Tensor Flow*: Tensor Flow [1], Google's recently open-sourced framework for implementing and deploying large-scale machine learning models, builds upon their experience with Dist Belief. It is already extensively employed internally for computations on a diverse array of mobile devices and large-scale distributed systems. The distributed version, released in April 2016, leverages a computation graph split into subgraphs for each device, with communication facilitated through Send/Receive node pairs. *Elastic Averaging SGD*: In the pursuit of innovation, Zhang et al. [23] propose Elastic Averaging SGD (EASGD), which connects the parameters of asynchronous SGD workers with an elastic force represented by a center variable stored in the parameter server. This design enables local variables to deviate further from the center variable, fostering enhanced exploration of the parameter space. Empirical evidence demonstrates that this heightened exploration capacity leads to improved performance by uncovering new local optima. These advancements in parallelization and distribution of SGD showcase the ongoing efforts to optimize its performance, enabling faster convergence and wider applicability in diverse computing environments.

6. UNVEILING ADDITIONAL STRATEGIES FOR ENHANCING SGD PERFORMANCE

In this section, we explore supplementary strategies that can be employed in conjunction with the aforementioned algorithms to further optimize the performance of SGD. For a comprehensive overview of additional effective techniques, we recommend referring to [11]. *Shuffling and Curriculum Learning*: Typically, it is desirable to present training examples to the model in a randomized order to prevent optimization algorithm bias. Consequently, shuffling the training data after each epoch is often a prudent choice. However, in certain scenarios where the objective is to tackle progressively challenging problems, presenting the examples in a meaningful order can actually improve performance and convergence. This approach, known as Curriculum Learning [3], establishes a structured order for the examples. Zaremba and Sutskever [21] successfully employed a combined or mixed strategy in training LSTMs for evaluating simple programs, outperforming the naive approach that sorts examples solely based on increasing difficulty. *Harnessing the Power of Batch Normalization*: To facilitate effective learning, it is customary to normalize the initial parameter values by initializing them with zero mean and unit variance. However, as training progresses and parameter updates vary, this normalization is lost, leading to slower training and amplified changes as the network deepens. Batch normalization [9] restores these normalizations for each mini-batch and enables backpropagation of changes through the operation. By incorporating normalization into the model architecture, higher learning rates can be utilized, and less emphasis needs to be placed on initialization parameters. Additionally, batch normalization serves as a regularizer, reducing or even eliminating the need for Dropout. *Employing Early Stopping*: As emphasized by Geoff Hinton, early stopping represents a "beautiful free lunch" [16]. It is crucial to continuously monitor the error on a validation set during training and exercise patience when observing insufficient improvement in validation error. Stopping the training process at an appropriate stage can prevent unnecessary computations and resource consumption. *Embracing Gradient Noise*: Neelakantan et al. [13] introduce the concept of adding noise following a Gaussian distribution $N(0, \sigma^2_t)$ to each gradient update. This is expressed by the equation: $gt_i = gt_i + N(0, \sigma^2_t)$ (34)

The variance is annealed according to the following schedule: $\sigma^2_t = \eta / (1 + t)^\gamma$ (35) The inclusion of such noise enhances the robustness of networks against poor initialization and proves particularly beneficial for training deep and complex networks. It is postulated that the introduced noise provides the model with additional opportunities to escape and discover new local minima, which are more prevalent in deeper models. These additional strategies expand the repertoire of techniques available to optimize SGD, enabling practitioners to further enhance its performance in various settings.

7. CONCLUSION: UNVEILING THE POWER OF OPTIMIZING SGD

In this comprehensive article, we embarked on a journey through the realm of stochastic gradient descent (SGD), exploring its three variants. Among them, mini-batch gradient descent emerged as the favored choice within the machine learning community. Delving deeper into the optimization of SGD, we examined a plethora of algorithms that are widely employed to maximize its potential. From the stalwart Momentum and Nesterov accelerated gradient to the adaptive marvels of Adagrad, Adadelta, RM Sprop, Adam, Ada Max, and Nadam, we covered the spectrum of optimization techniques. Furthermore, we explored various algorithms tailored specifically for asynchronous SGD, enabling parallelization and distribution of computations. To elevate the performance of SGD even further, we ventured into additional strategies that prove invaluable. We unraveled the benefits of shuffling and curriculum learning, where intelligently ordering training examples can enhance performance and convergence. The transformative power of batch normalization, a technique that restores parameter normalizations during training, was unveiled. Additionally, we highlighted the importance of early stopping, a valuable approach for monitoring and halting training when validation error fails to improve sufficiently. By assimilating these diverse strategies, ranging from algorithmic choices to auxiliary techniques, practitioners can unlock the full potential of SGD, revolutionizing the optimization landscape for a broad array of applications. In summary, this article serves as a comprehensive guide, equipping readers with an arsenal of tools and knowledge to unleash the true power of SGD optimization.

REFERENCES

- [1]. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2015). TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems.
- [2]. Bengio, Y., Boulanger-Lewandowski, N., & Pascanu, R. (2012). Advances in Optimizing Recurrent Networks.
- [3]. Bengio, Y., Louradour, J., Collobert, R., & Weston, J. (2009). Curriculum Learning. Proceedings of the 26th Annual International Conference on Machine Learning.
- [4]. Darken, C., Chang, J., & Moody, J. (1992). Learning Rate Schedules for Faster Stochastic Gradient Search. Neural Networks for Signal Processing II Proceedings.
- [5]. Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., & Bengio, Y. (2014). Identifying and Attacking the Saddle Point Problem in High-Dimensional Non-Convex Optimization.
- [6]. Dean, J., Corrado, G. S., Monga, R., Chen, K., Devin, M., Le, Q. V., ... Ng, A. Y. (2012). Large-Scale Distributed Deep Networks. Neural Information Processing Systems.
- [7]. Dozat, T. (2016). Incorporating Nesterov Momentum into Adam. ICLR Workshop.
- [8]. Duchi, J., Hazan, E., & Singer, Y. (2011). Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. Journal of Machine Learning Research.
- [9]. Ioffe, S., & Szegedy, C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.
- [10]. Kingma, D. P., & Ba, J. L. (2015). Adam: A Method for Stochastic Optimization. International Conference on Learning Representations.
- [11]. LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. R. (1998). Efficient BackProp. Neural Networks: Tricks of the Trade.
- [12]. McMahan, H. B., & Streeter, M. (2014). Delay-Tolerant Algorithms for Asynchronous Distributed Online Learning. Advances in Neural Information Processing Systems.
- [13]. Neelakantan, A., Vilnis, L., Le, Q. V., Sutskever, I., Kaiser, L., Kurach, K., & Martens, J. (2015). Adding Gradient Noise Improves Learning for Very Deep Networks.
- [14]. Nesterov, Y. (n.d.). A Method for Unconstrained Convex Minimization Problem with the Rate of Convergence $O(1/k^2)$.
- [15]. Niu, F., Recht, B., Christopher, R., & Wright, S. J. (2011). Hogwild!: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent.
- [16]. Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global Vectors for Word Representation. Conference on Empirical Methods in Natural Language Processing.
- [17]. Qian, N. (1999). On the Momentum Term in Gradient Descent Learning Algorithms. Neural Networks.
- [18]. Robbins, H., & Monro, S. (1951). A Stochastic Approximation Method