# Using Gradient Descent to an Optimisation Algorithm that Uses the Optimal Value of Parameters (Coefficients) for a Differentiable Function

**Falah Amer Abdulazeez**
*University of Anbar, Ramadi city - Al Anbar Governorate, Iraq*
*falah.amer.azeez@uoanbar.edu.iq*
**Abdul Sttar Ismail**
*University of Anbar, Ramadi city - Al Anbar Governorate, Iraq*
*sttarwdaa2019@uoanbar.edu.iq*
**Rafid S. Abdulaziz**
*University of Anbar, Ramadi city - Al Anbar Governorate, Iraq*
*rafid.alhashimy@uoanbar.edu.iq*

| Article History | Abstract |
|---|---|
| | Deep neural networks (DNN) are commonly employed. Deep networks' many parameters require extensive training. Complex optimisers with multiple hyperparameters speed up network training and increase generalisation. Complex optimiser hyperparameter tuning is generally trial-and-error. In this study, we visually assess the distinct contributions of training samples to a parameter update. Adaptive stochastic gradient descent is a batch stochastic gradient descent variation for neural networks using ReLU in hidden layers (aSGD). It involves the mean effective gradient as the genuine slope for boundary changes, in contrast to earlier procedures. Experiments on MNIST show that aSGD speeds up DNN optimisation and improves accuracy without added hyperparameters. Experiments on synthetic datasets demonstrate it can locate redundant nodes, which helps model compression. |
| | |

## 1. Introduction

One effective method for local optimisation is known as gradient descent. It is a first-order approach that only involves the gradient and has several applications, including optimum control, video coding, robotics, and localisation. Nesterov has devised a quick gradient approach, which may be used for various applications, including model predictive control. One of the most well-known algorithms for optimisation, gradient descent, is also the method used most of the time while optimising neural networks. At the same time, every cutting-edge Deep Learning library offers implementations of several methods to optimise gradient descent (for example, the documentation for lasagne, Caffe, and keras).

Utilising a strategy known as slope plunge, plausible to limit a goal capability is defined by a model's boundaries to accomplish the ideal outcomes. In this strategy, the boundaries are refreshed in a way that is counter to where the slope of the goal capability takes as to the boundaries. The learning rate, represented by the letter $alpha$, concludes the greatness of the means that we need to do to come to a (nearby) negligible degree of execution. To put it another way, we continue descending

toward the surface's incline framed by the goal capability until we arrive at a valley. This continues until we arrive at the surface's lower part.

### 1.1 Optimisation Algorithm

Streamlining alludes to the most common way of limiting or boosting the worth of a goal capability defined by x. This system is alluded to as the advancement interaction. This try is alluded to as the objective of limiting the expense/loss capability J(w), defined by the model's boundaries w Rd, in the language of AI and profound learning. As such, the goal is to diminish how much cash is lost. With regards to depreciation, advancement calculations pursue achieving one of the accompanying objectives:

• Find where the goal capability is at its generally limited least wherever on the planet. If the goal capability is raised, implying that any nearby least is likewise a worldwide least, then this is the kind of thing that is plausible. If the goal capability isn't curved, then, at that point, this isn't possible.

• Sort out the worth of the goal capability that is the least practicable one inside its encompassing region and utilise that as your beginning stage. If the goal capability isn't curved, which is the situation in most issues, including profound learning, then, at that point, this is often the condition that emerges.

There are three basic classes for improvement calculations, which are as per the following:

• A streamlining strategy that is iterative and meets an appropriate arrangement regardless of the statement of the boundaries, for example, slope plunge when applied to calculated relapse.

• A streamlining methodology that doesn't repeat and settles for a solitary point in the issue space.

• A streamlining technique that is of the iterative kind in its way to deal with tackling issues, as brain organisations, and works on an assortment of issues that have non-raised cost works; this strategy is known as brain network improvement. As a result, introducing the boundaries is a critical stage during the time spent accomplishing lower mistake rates and accelerating the combination methodology.

### 1.2 Gradient Descent

In AI and profound learning, the technique known as Slope Drop is the one that is utilised most frequently for improvement. This specific strategy is for first-request streamlining. While executing the reports on the boundaries, it just thinks about the primary subsidiary. During every emphasis, we report on the boundaries to make them change on the contrary course of the slope of the goal capability J(w) regarding the boundaries. This is done because the slope shows which course the rising will be the most difficult. The learning rate decides the size of the step we take all through every cycle to move toward the nearby least and draw nearer to it. Consequently, we travel downward while adhering to the direction that the slope slows in order to reach the lowest point in the region.

Gradient descent is an optimisation procedure used to identify the values of the parameters (coefficients) of a function (f) that result in the lowest feasible cost function. This is done to find the optimal solution (cost).

When it is impossible to identify the parameters analytically (for instance, by utilising linear algebra), and they must instead be searched for via an optimisation process, gradient descent is the most efficient method.

The concept of gradient descent as a whole is founded on the notion of randomly initialising the variable w and then executing consecutive updates:

$$w_{k+1} \leftarrow w_k - a\nabla f(w_k) \tag{1}$$

Where the convergence is dependent not just on but also on the structural features of F itself.

### 1.3 Different Iterations of the Gradient Descent Algorithm

In AI and profound learning, the technique known as Slope Drop is the one that is utilised most frequently for improvement.

• Batch Gradient Descent;

• Mini-batch Gradient Descent;

• Stochastic Gradient Descent.

### 1.4 Differentiable Function

In math, a capability with a solitary variable is supposed to be differentiable if the subsidiary of the capability can be situated at every single point over the whole of the capability's space. All in all, a differentiable capability can be separated. The digression line to the diagram of a differentiable capability won't ever be in an upward position at any inside point given inside its space. This is because different capabilities can take on various qualities. In science, a differentiable capability is one that contains no breaks, cusps, or points. Be that as it may, few out of every odd constant capability can be separated, even though a differentiable capability should continuously be nonstop.

*1.4.1    Differentiable Functions and Their Associated Rules*

On the off chance that f and g are differentiable capabilities, the subsidiaries of their aggregate, contrast, item, and remainder might be acquired by applying specific principles to every one of their individual definitions. This is true whether or not f and g are both differentiable functions. The following are various differentiability formulae that may be used to locate the derivatives of a function that can be differentiated:

- $(f + g)' = f' + g'$
- $(f - g)' = f' - g'$
- $(fg)' = f'g + fg'$
- $(f/g)' = (f'g - fg')/f^2$

Example

Let's calculate the derivative of the function by applying the differentiability principles to the problem $f(x) = (2x+1)^3$

$df/dx = d(2x+1)^3/dx$

$= d(8x^3 + 12x^2 + 6x + 1)/dx$

$= 24x^2 + 24x + 6$

$= 6(2x+1)^2$

## 2.  Related Work

In [2], whether you are fostering an answer for an issue, all things considered, or composing code for a product application, streamlining ought to be your primary goal continuously. Notwithstanding, this goal might be achieved by carrying out one of the accessible streamlining procedures. The always well-known Gradient Descent (GD) streamlining calculations are habitually utilised as "black box" enhancers regarding the goal of issues including unconstrained improvement. A calculation that depends on slopes will, through every emphasis of the calculation, try to draw nearer to the minimiser or maximise the expense capability by using the data given by the inclination's goal capability.

In [10], a few unique kinds of AI are drawing near; for example, brain organisations, support vector machines, and transformative registering all incorporate the method of settling complex improvement issues. The essential focal points of the improvement hypothesis are (1) the standards that should be met for an extremum worth to exist (known as slope examination) and (2) the advancement of enhancement calculations and the investigation of their intermingling. This part centres around the hypothesis and methods of curved improvement, explicitly on the slope and sub-angle approaches in smooth and non-smooth raised advancements and confined arched streamlining.

In [5], this article aims to research the qualities of sign inclination drop calculations. The fact that pioneered these calculations makes these estimations contrast with conventional slope plummet techniques in that they focus on the indication of the angle as opposed to the actual inclination; the Rprop calculation is the one. Customary inclination plunge techniques centre around the actual angle. This article gives two unmistakable consequences of intermingling for neighbourhood advancement: the first is for ostensible frameworks liberated from vulnerability, and the second is for frameworks that incorporate vulnerabilities. Both of these ends are introduced for frameworks that incorporate vulnerabilities. To exhibit the viability of novel sign slope plunge calculations, for example, the polarity calculation DICHO, as far as the speed at which they meet, it is essential to apply these calculations to different settings and afterwards dissect the results of those applications. Since these calculations may, as a curiosity, meet truly towards various minima than the closest least of the beginning condition, they are great for use in worldwide improvement as a new metaheuristic approach. This is because worldwide advancement is one of the main applications of metaheuristics.

Along these lines, the methods for sign slope plummet are appropriate for use in worldwide improvement.

In [4], cell networks have an immense functional obligation that expects them to expand both the limit of their organisations and the inclusion they give. Changing the azimuths and slants of receiving wires to accomplish the most significant organisation inclusion is the goal of this task. The techniques presently being utilised, which are prevalently angle-free strategies, are being utilised to achieve this objective. We can exhibit that the angle vector is scanty, which permits us to demonstrate that the calculation of the slope can be finished all the more rapidly regardless of the enormous number of decision factors that incorporate azimuths and slants. This is because there is a limitation on the number of base stations arranged inside a particular distance of a specific example site. Thus, there are several requirements.

Furthermore, the SGD strategy requires a reasonable measure of process, as it is predicated on minimal expense assessments of the slopes. This makes the technique a suitable choice. As a result of this, it can be used in a proficient way when applied to organisations of a vast scope. Tests show that the proposed techniques perform very well when contrasted with meta-heuristic calculations regarding the close ideal arrangements they give and the productivity with which they register such responses. This is valid both as far as the close ideal arrangements they furnish and the productivity with which they process them. This discussion will likewise zero in on how versatile the suggested calculations are and the way that pertinent they are in different settings.

In [9], the article suggests a cross-breed calculation that might be utilised to look for the worldwide least of a multimodal capability. A hunt technique is partitioned into two phases: the primary stage is the chaotic optimisation algorithm (COA) in view of two times transporter waves for worldwide looking. The subsequent stage is the gradient descent algorithm (GDA) for precise neighbourhood looking. The turbulent elements are rejuvenated by involving a one-layered map in one of three distinct variations: the strategic guide, the cubic guide, or the site map. Three different testing capabilities are utilised. The proposed half-and-half calculation and the GDA working all alone were utilised to do a sum of 100 recreations, with every reenactment starting from a particular introductory point picked indiscriminately. These recreations were completed for every one of the test capabilities. It was discussed whether it was feasible to find the extremum with progress and accuracy and whether the calculations could combine while utilising the three detailed turbulent guides.

## 3. The Gradient of a Differentiable Function

### 3.1 Convex Functions

We want to examine GD (and SGD in the future) more comprehensively. We will continuously operate under the premise that convexity is true, even though these techniques are also used—and occasionally analysed—when this assumption is invalid. In the rest of what follows, f will signify the goal, and x or y will be its variables. Exceptions will be made for the instances (they do not stand anymore for a predictor or training variables).

Definition 1: (Convex function) A differentiable function: $\mathbb{R}^P \to \mathbb{R}$ is said convex if

$$f(y) \geq f(x) + \nabla f(x)^T (y - x), \forall x, y \in \mathbb{R}^P \tag{2}$$

If function f can be differentiated twice, then this is the same as needing $\nabla^2 f(x) \succeq 0, \forall x \in \mathbb{R}^P$ (here $\succeq$ represents the semidefinite partial ordering known as the Loewner order and is denoted by the letter $A \succeq B \iff A - B$ is positive and semidefinite). A definition of convexity that is more comprehensive would say that A $\forall x, y \in \mathbb{R}^P$ and $\alpha \in [0,1]$,

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y). \tag{3}$$

As a practice problem, demonstrate that if f can be differentiated, then this is the same as our definition. The following inequality is seen relatively often in the demonstrations utilising convexity.

Proposition 1: (Jensen's inequality) Iff: $\mathbb{R}^P \to \mathbb{R}$ is convex and μ is a probability calculated om $\mathbb{R}^P$

$$f(\int x d\mu(x)) \leq \int f(x) d\mu(x) \tag{4}$$

To put it another way: "the average picture is smaller than the average image of the images".

Proof Let $x* = R\ xd\mu(x)$. By the definition of convexity, we have $f(x) \geq f(x*) + \nabla f(x*)^T(x - x*) \forall x \in \mathbb{R}^P$ After integrating, Jensen's inequality follows, noting that $f + \nabla f(x*)^T(x - x*)d\mu(x) = 0$. The following stability characteristics are met by the class of convex functions (exercise):

- If $(f_j)_{j \in [m]}$ are convex and $(a_j)_{j \in [m]}$ are nonnegative, then $\sum_{j=1}^m a_j f_j$ is convex

- If $f: \mathbb{R}^P \to \mathbb{R}$ is convex and $A: \mathbb{R}^{P'} \to \mathbb{R}^P$ is linear then $f \circ A: \mathbb{R}^{P'} \to \mathbb{R}$ is convex

Example. If the loss in the second variable is convex, then problems of the type Equation (1) are convex, ( ) fixes linear in w, and $\Omega$ is convex.

It is also essential to highlight the following characteristic (immediately from the definition).

Proposition 2: Assume that: $\mathbb{R}^P \to \mathbb{R}$ is both differentiable and convex. Then $x* \in \mathbb{R}^P$ is a global minimiser of f iff

$$\nabla f(x*) = 0 \tag{5}$$

Example 1: (Low-rank matrix completion) Let $M \in \mathbb{R}^{n \times d}$ be a data matrix, where low rank is presumed to exist "$\min\{n, d\}^3$ Consider sampling a portion of the M coefficients, represented by $S \in 1, ..., n \times 1, ...,$ the next step is to rebuild the matrix M from the samples, with the hope that the low-rank characteristic will allow us to accomplish this job with substantially fewer samples than $n \times d$.

Instead of putting the issue into a matrix $W \in \mathbb{R}^{n \times d}$ We take into consideration the so-called Burer-Monteiro technique with explicit low-rank restrictions in which the issues is resolved over two smaller matrices $\in \mathbb{R}^{n \times r}$ and $V \in \mathbb{R}^{d \times r}$ which, in terms of the issue factors, may be relatively affordable $((n + d)r$ versus $nd)$. As a consequence of this, we will examine:

$$\underset{U \in \mathbb{R}^{n \times r}, V \in \mathbb{R}^{d \times r}}{\text{minimize}} \frac{1}{2} \sum_{(i,j) \in S} \left([UV^T]_{ij} - M_{ij}\right)^2 \tag{6}$$

The equation does not satisfy the convexity condition in the variables U and V. Nevertheless, it is easy to prove that there are no bogus local minima and that all saddle points are strict if sufficient assumptions are made about the data the number of samples. These assumptions include the following.

## 4. Analysis Of Gd for Strongly Convex and Smooth Functions Based On Differentiable Functions

Definition 2 (Strong convexity) One may say that a differentiable function f is highly convex if there is $a > 0$, iff

$$f(\mathcal{Y}) \geq f(x) + \nabla f(x)^T(y - x) + \frac{a}{2} \|x - y\|_2^2, \forall x, y \in \mathbb{R}^P \tag{7}$$

This equals twice differentiable functions $\nabla^2 f \succeq \alpha Id$ This characteristic means that f permits a singular minimiser x*, denoted by $\nabla f(x*) = 0$. Additionally, it ensures that the gradient is substantial when a point is far from optimality:

Lemma 1 It holds if f is differentiable and highly convex with minimiser x*

$$\|\nabla f(x)\|_2^2 \geq 2a\left(f(x) - f(x*)\right), \forall x \in \mathbb{R}^P \tag{8}$$

Proof: In Definition 2, the right-hand side is highly convex in y and reduced by $\tilde{\mathcal{Y}} = x - \frac{1}{a}\nabla f(x)$. When we take $y = x*$ on the left-hand side and plug this value into the bound, we obtain

$$f(x*) \geq f(x) - \frac{1}{a}\|\nabla f(x)\|_2^2 + \frac{1}{2a}\|\nabla f(x)\|_2^2 = f(x) - \frac{1}{2a}\|\nabla f(x)\|_2^2 \tag{9}$$

The conclusion is then rearranged.

Definition 3 (Smoothness): If f is a differentiable function and is β-smooth, then:

$$\left|f(y) - f(x) - \nabla f(x)^T(y - x)\right| \leq \frac{\beta}{2}\|x - \mathcal{Y}\|^2, \forall x, y \in \mathbb{R}^P \tag{10}$$

In other words, this is the same as f having a β-Lipschitz gradient $\|\nabla f(x) - \nabla f(\mathcal{Y})\|_2^2 \leq \|x - y\|_2^2, \forall x, y \in \mathbb{R}^P$ For twice differentiable functions, this is equal to $-\beta Id \preceq \nabla^2 f \preceq \beta Id$.

In the following theorem, we demonstrate that the gradient descent algorithm converges exponentially for problems of this kind.

Theorem 3: Take it for granted that f is both β-smooth and substantially convex. Selecting $\eta_t = 1/\beta$ the iterates $(\quad t x)_{t \geq 0}$ of GD on f satisfy

$$f(x_t) - f(x*) \leq \exp(-t\beta/a)\big(f(x_0) - f(x*)\big) \tag{11}$$

Proof The following is a descending quality that we have due to smoothness, with $\eta_t = 1/\beta$,

$$f(x_{t+1}) = f(x_t - \nabla f(x_t)/\beta) \leq f(x_t) - \|\nabla f(x_t)\|_2^2/\beta + \frac{1}{2\beta}\|\nabla f(x_t)\|_2^2$$
(12)

Rearranging, we obtain:

$$f(x_{t+1}) - f(x*) \leq \big(f(x_t) - f(x*)\big) - \frac{1}{2\beta}\|\nabla f(x_t)\|_2^2$$
(13)

By using Lemma 1, it can be shown that:

$$f(x_{t+1}) - f(x*) \leq (1 - \alpha/\beta)\big(f(x_t) - f(x*)\big) \leq \exp(-\alpha/\beta)\big(f(x_t) - f(x*)\big) \tag{14}$$

In conclusion, we will use a recursion.

- We necessarily have $\alpha \leq \beta$. The ratio $\kappa := \beta/\alpha$ is called the condition number.
- If we make the single assumption that the function is smooth and convex (but not highly convex), then GD with a constant step-size of $\eta = 1/\beta$ converges when there is a minimiser present as well, but at a slower pace in O(1/t)
- When choosing the step size, all required is an upper limit $\beta$ on the smoothness constant of (if it is overestimated, the convergence rate will decrease somewhat).

Example 2: regularised logistic regression Take, for example, an assignment involving classification with $\mathcal{Y} \in \{-1, +1\}$, the logistic loss $\ell(\mathcal{Y}, z) = \log(1 + e^{-\mathcal{Y}z})$ a linear model $f_w(x) = x^T w$ and regularization $\lambda\|w\|_2^2$. The goal of empirical risk reduction is

$$F(w) = \frac{1}{n}\sum_{i=1}^{n}\log\left(1 + e^{-\mathcal{Y}_i x_i^T w}\right) + \lambda\|w\|_2^2 \tag{15}$$

This function can be differentiated and has a convex shape. Because of the regularisation term, it has a convexity of at least $2\lambda$- substantial degrees. Its gradient is:

$$\nabla F(w) = \frac{1}{n}\sum_{i=1}^{n}\frac{-\mathcal{Y}_i x_i}{1 + e^{-\mathcal{Y}_i x_i^T w}} + 2\lambda w$$
(16)

and its Hessian $\nabla^2 F(w) = (\partial_{ij}F(w))_{i,j=1}^{d}$ is

$$\nabla^2 F(w) = \frac{1}{n}x_i x_i^T \frac{e^{-\mathcal{Y}_i x_i^T w}}{(1 + e^{-\mathcal{Y}_i x_i^T w})^2} + 2\lambda$$
(17)

Thus, F is $\beta$-smooth with $\beta = (1/n)\sum_{i=1}^{n}\|x_i\|_2^2 + 2\lambda$. Therefore, the condition number, which is what ultimately decides the convergence speed, is $\kappa = \beta/\alpha = 1 + (1/(2\lambda n))\sum_{i=1}^{n}\|x_i\|_2^2$. It was discovered that the regularisation, initially used to cut down on estimating error, actually helped optimisation.

## 5. Complex Gradient Descent Algorithm (CGD) Based on Differentiable Function

We start out by introducing a few concepts and notations. For $z = [z_1, \ldots, z_n]^T \in \mathbb{C}^n$, its conjugate is denoted by $\bar{z} = [\bar{z}_1, \ldots, \bar{z}_n]^T \in \mathbb{C}^n$, Write z = x+iy, r = [x, y], and $c = [z, \bar{z}]$ Hence, a complex function f(z): $\mathbb{C}^n \to \mathbb{C}$, with some creative use of notation, may be written in the following many forms:

$$f(z) = f(z, \bar{z}) = f(c) = f(x, y) = f(r) \tag{18}$$

As is customary, the gradient operator will be defined by $\frac{\partial}{\partial z} = \left[\frac{\partial}{\partial z_1}, \frac{\partial}{\partial z_2}, \ldots, \frac{\partial}{\partial}\right]$ the conjugate gradient operator by $\frac{\partial}{\partial \bar{z}} = \left[\frac{\partial}{\partial \bar{z}_1}, \frac{\partial}{\partial \bar{z}_2}, \ldots, \frac{\partial}{\partial}\right]$ and a differentiable function's gradient$\left(^{-}\right)$, by $\nabla_z f = \left(\frac{\partial f}{\partial}\right)$ where $(\cdot)$ H identifies the Hermitian transposition as being performed. Additionally, we indicate by $\langle a\ b\rangle = a^H b$ the result of multiplying two complex n-vectors together, $b \in \mathbb{C}^n$

We're assuming $\in \mathbb{C}^n$ n is a local minimal-value point of a real-valued function g(z). Let a0, which is located in the same area as a, be the first estimate for a. A is discovered via a gradient descent approach as the limit of the series of (ak):

$$a_{K+1} = a_k - t_k \nabla g(a_k), \quad k = 0,1,2 \ldots,$$

(19)

Where $\nabla g$ in (2) is Lipschitz constant and continuous$> 0$, i.e.,

$$\|\nabla g(a) - \nabla g(b)\| \leq L\|a - b\|$$ (20)

We use backtracking line search in this, where a fixed$\beta$, $0 < \beta < 1$, is employed for formulating $t_k$ by $t_k = \beta t_{k-1}, t_1 = 1$. When this condition is met, the iteration will come to an end.

$$g(a_k - t_k \nabla g(a_k)) > g(a_k) - \frac{t_k}{2} \|\nabla g(a_k)\|^2$$ (21)

The approach described above is known as the Complex Gradient Descent Method (CGD). It is also possible to use the convergence theorems of the real gradient descent techniques with CGD.

In our approach, we put $g = E$ in (2), where $E(a) = E(f, A)$ is the energy function of f at A. The energy function of f at A is denoted by $E(a)$. In addition, to ensure that the new n-tuple is produced after each iteration step, ak is still in $\mathbb{D}^n$, $t_k$, tk in (3) and must satisfy

$$a_k + t_k \nabla E(a_k) \in \mathcal{N}_{a_k} \cap \mathbb{D}^n,$$ (22)

$$\mathcal{N}_{a_k} = \left\{ z; \left\| (a_k)_{j-z_j} \right\| < r, 1 \leq j \leq n \right\}$$ (23)

CGD must have an efficient formulation of the gradient $-\nabla E (= \nabla A)$. The evidence points to the fact that:

$$e_a(z) = \frac{\sqrt{1-|a|^2}}{1-\overline{a}z}, \quad a \in \mathbb{D},$$ (24)

And let $P_\ell, l = 1, \ldots, n$be the permutation of the set that serves as the index $\{1, 2, \cdots, n\}$ such that $P_\ell(n) = \ell$. In consideration of the existing analytic function$f \in H^2$We determine n functions using inductive reasoning$fP_{\ell(j)}, 1 \leq j \leq n$, by

$$f_{P_\ell(1)}(z) = f(z),$$
$$f_{P_\ell(j)}(z) = \frac{1 - z\overline{a}_{P_\ell(j-1)}}{z - a_{P_\ell(j-1)}} \left( f_{P_\ell(j-1)}(z) - \langle f_{P_\ell(j-1)}, e_{a_{P_\ell(j-1)}} \rangle e_{a_{P_\ell(j-1)}}(z) \right).$$

(25)

It was proved that all $fP_{\ell(j)}(z)$ is analytical in the letter $\mathbb{D}$. The energy function $E(a)$ may be represented in n distinct ways because it is not affected by the permutation.:

$$E(a) = \sum_{j=1}^{n} \left( 1 - \left| a_{P_{\ell(j)}} \right|^2 \right) \left| fP_{\ell(j)} \left( a_{P_{\ell(j)}} \right) \right|^2, \ell = 1, .., n,$$ (26)

Note that the variable$a_\ell = aP_{\ell(n)}$only appears in the last component of the total when it (5). Since$fP_{\ell(n)}$is analytic, we have$\frac{\overline{fP_{\ell(n)}}}{\partial z_\ell} = 0$, so that

$$\frac{\partial z^\ell}{\partial(-E(a))} = \overline{f}_{P_\ell(n)}(a^\ell)\left(a^\ell \overline{f}_{P_\ell(n)}(a^\ell) - (1 - |a^\ell|_2)\overline{f}'_{P_\ell(n)}(a^\ell)\right), \quad \ell = 1, \cdots, n.$$

(27)

Where$f'P_{\ell(n)}$may be determined using

$$f'_{P_\ell(n)}(z) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{f_{P_\ell(n)}(e^{i\theta})e^{-i\theta}}{(1 - ze^{-i\theta})^2} d\theta.$$

(28)

Remark 1. The permutation P is not unique. In practice, we select$P_\ell = P^\ell$, where P is the $1-$ shift permutation: $P(1, 2, \cdots, n) = (2, \cdots, n, 1)$.

Remark 2. In fact, in practice, we set$\beta = \frac{1}{\|f\|^2}$ if $\|f\|^2 > 1$; or $\beta = \frac{1}{2^k \|f\|^2}$ if $\|f\|^2 \leq 1$, where K is an integer such that$2^k \|f\|^2 > 1$Such a setting of $\beta$ in this paper is just to ensure that$a_k \in \mathbb{D}^n$at every stage. However, we must point out that this option is not the best available one $\beta$. It has come to our attention that the specifications , which are governed mainly by $\beta$, significantly impact the precision and efficiency of Algorithm 1 (CGD), as in Figure 1.

```
Require: a, f, β, neighbor size r, and the tolerance ε.
1: Compute ∇E(a) using f and a.
2: while ‖∇E(a)‖² > ε do
3: Find s₁ > 0, s₂ > 0 such that ‖a + s₁∇E(a)‖∞ = 1 and ‖s₂∇E(a)‖∞ = r. Set s = min (s₁, s₂).
4: Compute c = a + s∇E(a).
5: while E(c) < E(a) + s/2 ‖∇E(a))‖² do
6: Update s: s = βs
7: Update c: c = a + s∇E(a).
8: end while
9: Update a: a = c.
10: Re-compute ∇E(a)
11: end while
12: Set the output: b = a
Ensure: b
```

*Figure 1. CGD: An Intricate Gradient Descent Search Technique for Finding the Optimal Tuple*

## 6. Methodology

*6.1 Gradient Descent*

While attempting to distinguish a neighbourhood least of a differentiable capability, the enhancement calculation known as "inclination plummet" is much of the time utilised. In the field of AI, enhancing different parameters is habitually utilised. Regarding AI, enhancement difficulties can be considered endeavours to limit an expense capability. The name given to this sort of capability is the goal capability. By more than once adjusting the boundaries of the bearing inverse to that of the angle (otherwise called the vector that incorporates the fractional subsidiaries as a whole).

AI calculations habitually break down the goal capability they use as an assumption for the per-tests misfortune capability (e.g., mean squared mistake loss and cross-entropy loss). Complete delegate works out.

$$L(x, y, \theta) = \frac{1}{M} \sum_{i=1}^{M} l\left(x^{(i)}, y^{(i)}, \theta\right)$$

(29)

$$g = \frac{1}{M} \sum_{i=1}^{M} \nabla_\theta l\left(x^{(i)}, y^{(i)}, \theta\right)$$

(30)

Where $\nabla_\theta l$ the gradients of the particular sample and g are the mean gradient of the whole loss, Equation (2) has an O computational cost (M). One calculation will need significant time if the training set is enormous.

As per the goal capability, we can slowly bring down the worth of the objective capability to move toward the base worth in the close by region. The answer for the restrictively high handling costs was proposed to be the stochastic slope plummet, curtailed as SGD, as shown in Figure 2. The basic idea driving SGD is that extended misfortunes all through the whole preparation set (that is, L(x, y, θ) in Equation (1)) can be approximately approximated by utilising a more modest batch of preparing tests. The preliminary calculation exhibits how to assess the anticipated misfortune and the angles of that misfortune. m is fixed, in spite of the way that the size of the preparation set M might increment. Subsequently, the expense of refreshing every boundary doesn't change straightforwardly depending upon how much information is in the preparation set, M, while the processing cost of SGD is O (m).

Furthermore, the techniques that include more than one yet less than the preparation can all be viewed as tests that have recently been alluded to as small-scale or little-batch stochastic strategies. It is presently more considered normal to utilise the term test. At present everyday practice, these strategies are alluded to as stochastic techniques.

---
**Algorithm 1:** Stochastic gradient descent

**Input:** learning rate, $\eta$, epoch, $T$, initial parameters, $\theta$, the neural network, $f$, loss function, $L$, the training set, $\left\{ \left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right) \ldots, \left(x^{(M)}, y^{(M)}\right) \right\}$.

1 **for** $t = 1$ to $T$ **do**
2     Sample $m$ samples from training set as a mini-batch ;
3     Estimate expectation: $L(x, y, \theta_t) = \frac{1}{m} \sum_{i=1}^{m} l\left(x^{(i)}, y^{(i)}, \theta_t\right)$;
4     Estimate gradient: $g = \nabla_\theta L(x, y, \theta_t)$ ;
5     Parameter update $\theta_{t+1} = \theta_t - \eta g$ ;
6 **end**

---

*Figure 2. Stochastic Gradient Descent*

### 6.2 The Importance of Samples in the Process of Updating Parameters

It has been found that the streamlining difficulties can be considered, including the decrease of the assumption for the misfortune capability related to every individual preparation test taken overall. SGD utilises a small assortment of guides to estimate this assumption. The normal of the multitude of individual misfortunes comprises the complete misfortune cases in this batch; subsequently, the slope of this complete misfortune rises to the mean of the angles of the singular misfortunes, the singular sufferings and misfortunes. Then again, the assumption that was gotten from the can be seen by checking out at line 3 in Algorithm 1.

An assembling of tests Because of the way that the data remembered for each example is novel, the consequences of each. There is a lot of variety across the model all's boundaries. To explain, this alludes to both the fractional subordinate's misfortune on a few qualities of a solitary example and the incomplete subsidiaries of per-test misfortune. There is a wide assortment of ways that one boundary could prompt misfortune. Like this, various fractional subsidiaries can be found in irrelevant angles. If we use ReLU as an actuation capability, it will bring about 0 like clockwork. To appropriately represent this point and guarantee appreciation, we will utilise a model from reality.

### 6.2.1 Notations

The "concentric circle" information is utilised as a twofold characterisation model in Figure 1a to feature each example's different jobs in the boundary update. It is a twofold order issue with classes as a result and organises as the information. The external red focuses are considered a "negative class," while the internal blue focuses are viewed as a "positive class." A completely associated network with one info layer, two hidden layers, and one result layer is utilised to prepare the information. The ReLU is utilised to enact neurons in subcutaneous layers. For simplicity of clarification, we make notational shows. L1, L2, L3, and L4 separately represent the four levels in Figure 3.
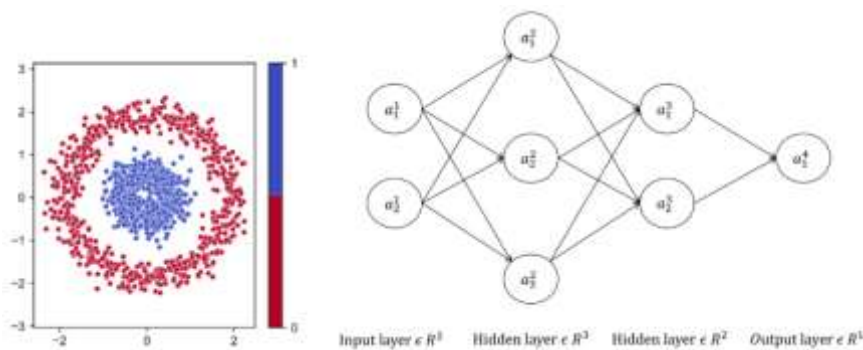


*Figure 3. A Dataset Called "Concentric Circle" with a Network Structure*

### 6.2.2 Dummy Node

It is important to note that specific nodes, such as the $a\_^3$, are muted for the input of every sample. We define it as follows: A node with the property that for every sample, $\sigma(z_l) < 0$, is an actual dummy node. On the other hand, a node is referred to as a pseudo dummy node if it is activated for any $\sigma(z_l) > 0$ and all samples, i.e., for any sample. They are both referred to as dummy nodes.

### 6.3 Gradient Descent with Sample-Based Adaptive Batch Size

The anticipated misfortunes over the whole preparation set are assessed utilising a small batch of preparing tests, as demonstrated in lines 3 and 4 of Calculation 1. The slope is then gotten from the mean of every slope of the singular examples. The commitment of each example to a similar boundary update is, be that as it may, unmistakable in light of the perception results. Since their "mistakes" never again spread in the regressive engendering and the slope disappears, specific examples don't add to the boundary update because of the property of actuation capability ReLU. The inclination gauge of SGD strays accordingly. A few examples in the smaller-than-normal batch are by and by viewed as genuine examples for computing the assumptions, even though they don't add to boundary refreshing. Accordingly, the anticipated angle over the preparation set is underrated, and the batch size (m in Algorithm 1 line 3) is exaggerated.

Accordingly, we recommend a methodology called test-based adaptive batch size slope plunge (aSGD), which independently computes the misfortune and inclination for each model boundary.

## 7. Result and Discussion

Three experiments are carried out in this section to show the properties of the suggested methodology. The network topology and training data used in Explore 1 are the same as those in Section 2. To study the impact of ASD on the training process, we kept track of all parameter changes during the model training procedure. Explore 2: To look at the effect of USD and SGD on the brain organisation and preparing results, we embraced organisation geography, utilising two engineered datasets, and counted the step, exactness, and misfortune toward the finish of preparing. In Explore 3, we tried the viability of aSGD and SGD utilising the MNIST transcribed digit data set. We ran the examination multiple times and counted the outcomes.

### 7.1 aSGD's Properties

In the second section, we trained the same network using USD and SGD. We discovered two features of an aSGD by contrasting the parameter change under the two optimisers.

#### 7.1.1 Speeding Up the Training Process

When comparing the parameter change curves for the two optimisers, such $w\_{11}^3$, $w\_{12}^3, w\_{13}^3$ and $b\_1^3$, it is clear that aSGD has a higher range and rate of parameter change than SGD. This is because ASD adjusts the parameter values using the mean effective gradient. The mean viable slope will be more noteworthy in a solitary step when contrasted with the slope in SGD. This is because the adaptive batch size rectifies the ReLU capability's slanted assessment of the general angle in SGD. It permits the mean viable angle to assess the usual inclinations across the preparation set. The preparation interaction is accelerated, and the slope plummet strategy turns out to be more productive.

Although the adaptive learning rate and batch size are comparable, their basic standards are fundamentally divergent. Adagrad, Adadelta, RMSprop, and Adam are instances of adaptive learning rate calculations that total the slope across various time moves toward altering the learning rate progressively. The model will unite all the more rapidly because of every boundary getting the right angle. The examples utilised in an aSGD make a proper slope. The slope of every boundary changes because of the distinctions among tests and elements with different recurrences. The mean compelling inclination is collected from many examples each time the boundaries are changed. The adaptive batch size shows this transformation. As such, the adaptive learning rate utilises the distinction achieved by the boundary's change during preparation, while the adaptive batch size utilises the boundary's distinction achieved by tests in a batch. Be that as it may, the ReLU's remarkable actuation attributes must make sense of this differentiation. At the point when the straight initiation esteem is under 0, the ReLU will explicitly bring 0 back. Thus, these examples do not affect the boundary update. ReLU misses the mark on the kind of property, as Tanh.

### 7.2 Artificial Datasets

We employ synthetic data for training and testing with a more redundant network structure to look at the characteristics of aSGD and dummy nodes in greater detail. On fictitious data, we evaluated how well aSGD and SGD performed. The experiment was run 10,000 times, and after the training, we recorded the step, accuracy, and loss. To create a probability distribution map, we counted the number and areas of phony hubs in the organisation. The particular plan is as follows: we train two kinds of manufactured information, "concentric circles" and "twisting lines," with SGD and aSGD, individually, using a fully connected network. The experiments are divided into four categories. The ratio of the samples from the training and test datasets, which follow the same pattern of data, is 7:3.

While Groups 2 and 4 train using SGD, Groups 1 and 3 train with aSGD for "concentric circle" and "spiral line," respectively. The ultimately linked network has five hidden levels with 10 nodes in each. Other nodes are ReLU, except the output node's sigmoid activation function. Before each training session, parameters are initialised randomly. Ten thousand training batches are run. Ten thousand batches or a minor loss causes us to discontinue training.

Groups 1 and 2 learn at a rate of 0.01, while Groups 3 and 4 learn at a rate of 0.05. The experiment was performed 10,000 times, with 128 samples being used for training in the four groups.

### 7.3 MNIST Handwritten Digit Dataset

We devise the accompanying analysis to explore the qualification in grouping execution between the SGD and aSGD. The information assortment comprises manually written digits from the MNIST. Each picture is a 28 by 28 dark scale picture, with 60,000 preparation photographs and 10,000 test pictures. The ten digits are partitioned into ten classes. The framework utilises a completely associated network with two hidden layers, each with 784, 30, 20, and 10 hubs. While different layers utilise ReLU, the result layer involves Softmax as its actuation capability. We train 500 ages at different learning rates with a batch size of 64. In this analysis, overfitting isn't forestalled by early halting. Thus, we can assess how well the two analysers act in different situations, like underfitting and overfitting. Multiple times of investigations are directed.

Because of an all the more unequivocally assessed mean compelling slope, it tends to be seen that, for similar learning rates, the paces of exactness improvement and misfortune decline in aSGD are much higher than those in SGD. In any case, raising SGD's learning rate might prompt a practically identical cycle rate. Subsequently, we, bit by bit, raise the learning rate while remaining within adequate reach and watch as the two enhancers change from underfitting to overfitting. While SGD can't accomplish both, it tends to be seen that at the learning pace of 0.03 in Figure 8b, SGD combines to decent test misfortune esteem. However, the preparation time could be shorter. At the learning pace of 0.01 in Figure 8b, the test loss of aSGD diminishes quickly, and the most reduced test misfortune esteem that can be accomplished in lower than SGD. The test misfortune quickly declines at the learning pace of 0.1 in SGD but can't arrive at good misfortune esteem.

Both in preparation and testing, aSGD performs better concerning exactness. While contrasting the two analysers' exhibitions, aSGD outflanks SGD in two key regions: preparing time effectiveness and precision. The expectation to learn and adapt support in SGD, notwithstanding, is considerably more noteworthy than it is in SGD. The test misfortune bend's cushion fills in size in the late preparation stages at the learning pace of 0.01 in aSGD. This may be because the adaptive batch size is generally lower or equivalent to the batch size for the hyper boundaries, which has a similar impact as accelerating learning. Subsequently, aSGD is shakier than SGD, particularly in preparing advances. Accordingly, more methods ought to be remembered for ASD, including early quit, learning rate timetables, dropout, and batch standardisation.

Tables 1 and 2 outline the mean and standard deviation of the most minimal test misfortune from 100 examinations and the connected test precision, train exactness, and train misfortune. As displayed in Table 1, when prepared with indistinguishable hyper boundaries, aSGD outflanks SGD in preparation and testing. The two analysers' learning rates are brought down to inside their appropriate reaches in Table 2. At a learning pace of 0.005, an SGD gets the most reduced test misfortune, while SGD is 0.03. Intense sort shows higher scores. Even though SGD's misfortune is fairly more modest than Sgd's, it is critical to note that SGD's insignificant misfortune arrived after 500 ages. However, aSGD's base misfortune is reached at 100 ages.

*Table 1. Mean and Variance of MNIST Model Performance at the Same Learning Rates*

| Method | Learning Rate | Train Accuracy | Train Loss | Test Accuracy | Test Loss |
|--------|--------------|----------------|------------|---------------|-----------|
| aSGD | 0.001 | 0.9464 ± 0.0055 | 0.1774 ± 0.0194 | 0.9372 ± 0.0044 | 0.2253 ± 0.0172 |
| SGD | 0.001 | 0.8137 ± 0.0438 | 0.5859 ± 0.1174 | 0.8186 ± 0.0430 | 0.5783 ± 0.1189 |
| aSGD | 0.003 | 0.9645 ± 0.0033 | 0.1189 ± 0.0117 | 0.9488 ± 0.0027 | 0.1965 ± 0.0126 |
| SGD | 0.003 | 0.8874 ± 0.0234 | 0.3787 ± 0.0774 | 0.8876 ± 0.0231 | 0.3834 ± 0.0778 |
| aSGD | 0.01 | 0.9652 ± 0.0039 | 0.1168 ± 0.0135 | 0.9498 ± 0.0034 | 0.1922 ± 0.0164 |
| SGD | 0.01 | 0.9376 ± 0.0063 | 0.2131 ± 0.0252 | 0.9317 ± 0.0062 | 0.2423 ± 0.0257 |

*Table 2. Mean and Variance of MNIST-trained Model Performance*

| Method | Learning Rate | Train Accuracy | Train Loss | Test Accuracy | Test Loss |
|--------|--------------|----------------|------------|---------------|-----------|
| aSGD | 0.001 | 0.9464 ± 0.0055 | 0.1774 ± 0.0194 | 0.9372 ± 0.0044 | 0.2253 ± 0.0172 |
| | 0.003 | 0.9645 ± 0.0033 | 0.1189 ± 0.0117 | 0.9488 ± 0.0027 | 0.1965 ± 0.0126 |
| | 0.005 | 0.9652 ± 0.0039 | 0.1168 ± 0.0135 | 0.9498 ± 0.0034 | 0.1922 ± 0.0164 |
| | 0.01 | 0.9662 ± 0.0032 | 0.1135 ± 0.0102 | 0.9493 ± 0.0029 | 0.1937 ± 0.0112 |
| SGD | 0.01 | 0.9376 ± 0.0063 | 0.2131 ± 0.0252 | 0.9317 ± 0.0062 | 0.2423 ± 0.0257 |
| | 0.03 | 0.9615 ± 0.0036 | 0.1311 ± 0.0129 | 0.9484 ± 0.0034 | 0.1984 ± 0.0176 |
| | 0.05 | 0.9625 ± 0.0033 | 0.1276 ± 0.0111 | 0.9484 ± 0.0026 | 0.2023 ± 0.0141 |
| | 0.07 | 0.9611 ± 0.0043 | 0.1326 ± 0.0145 | 0.9474 ± 0.0036 | 0.2074 ± 0.0164 |
| | 0.1 | 0.9599 ± 0.0039 | 0.1374 ± 0.0137 | 0.9465 ± 0.0029 | 0.2099 ± 0.0126 |

## 8. Conclusion

By visualising error backpropagation, we evaluated the function of samples in parameter update. Sparse data and different-frequency characteristics can cause SGD to estimate the predicted gradient incorrectly. ReLU's property increased biassed estimation. aSGD, an adaptive optimisation approach based on samples, was proposed. aSGD'sadaptive batch size permitted it to appraise the normal inclination more precisely than SGD. This superior slope drop strategy. aSGD performed well on manufactured and written hand-digit datasets. aSGD changed pointless boundaries into genuine faker hubs, recognising repetitive hubs and decreased models. aSGD is straightforward, however. Later in training, it might be unstable and must be employed alongside other optimisation strategies. ReLU currently amplifies sample-based parameter adjustments. aSGD only works with ReLU neural networks. Deep learning uses sigmoid, tanh, and step activation functions. Differentiating parameter changes owing to samples under these activation functions and generating novel adaptive algorithms require more research.

## References

[1] P. V. Hai, and J. A. Rosenfeld, Joel, "The gradient descent method from the perspective of fractional calculus.," *Mathematical Methods in the Applied Sciences*, vol. 44, no.7, pp.5520-5547, 2021.

[2] S. H. Haji, and A. M. Abdulazeez, "Comparison Of Optimisation Techniques Based On Gradient Descent Algorithm: A Review," *Palarch's Journal Of Archaeology Of Egypt/Egyptology*, vol.18, no.4, pp.2715-2743, 2021.

[3] J. D. Lee, M. Simchowitz, M. I. Jordan, and B.Recht, "Gradient descent only converges to minimisers: Non-isolated critical points and invariant regions," arXiv preprint, arXiv:1605.00405. 2016.

[4] Y. Liu, W. Huangfu, H. Zhang, K. Long, "An Efficient Stochastic Gradient Descent Algorithm to Maximise the Coverage of Cellular Networks," *IEEE Transactions on Wireless Communications*, vol.18, no.7, pp. 3424-3436, 2019.

[5] E. Moulay, V. Léchappé and F. Plestan, "Properties of the sign gradient descent algorithms," *Information Sciences*. vol.492, pp.29-39.

[6] S.Bubeck. "Convex optimisation: Algorithms and complexity. Foundations and Trends R in Machine Learning," vol.8, no.3-4, pp.231-357, 2015.

[7] S. Ruder, "An overview of gradient descent optimisation algorithms," arXiv preprint, arXiv:1609.04747, 2016.

[8] S. Chatterjee, "Convergence Of Gradient Descent For Deep Neural Networks," arXiv:2203.16462v3, 2022

[9] D. Tsankova, and L. Svetla, "Global Optimisation Algorithm Based on One-Dimensional Chaotic Maps and Gradient Descent Technique," *Information Technologies and Control*. vol.15. pp.17-22, 2017.

[10] X. D. Zhang, "Gradient and Optimization," In *A Matrix Algebra Approach to Artificial Intelligence*. Singapore: Springer, 2020, pp.89-155.