# signADAM: Learning Confidences for Deep Neural Networks

Dong Wang[*]
2947365183hk@gmail.com
School of Artificial Intelligence,
Xidian University

Yicheng Liu[*]
moooooore66@gmail.com
School of Electronic Engineering,
Xidian University

Wenwo Tang
tww1507019@126.com
School of Artificial Intelligence,
Xidian University

Fanhua Shang[†]
fhshang@xidian.edu.cn
School of Artificial Intelligence,
Xidian University

Hongying Liu
hyliu@xidian.edu.cn
School of Artificial Intelligence,
Xidian University

Qigong Sun
xd_qigongsun@163.com
School of Artificial Intelligence,
Xidian University

Licheng Jiao
lchjiao@mail.xidian.edu.cn
School of Artificial Intelligence,
Xidian University

## ABSTRACT

In this paper, we propose a new first-order gradient-based algorithm to train deep neural networks. We first introduce the sign operation of stochastic gradients (as in sign-based methods, e.g., SIGN-SGD) into ADAM, which is called as signADAM. Moreover, in order to make the rate of fitting each feature closer, we define a confidence function to distinguish different components of gradients and apply it to our algorithm. It can generate more sparse gradients than existing algorithms do. We call this new algorithm signADAM++. In particular, both our algorithms are easy to implement and can speed up training of various deep neural networks. The motivation of signADAM++ is preferably learning features from the most different samples by updating large and useful gradients regardless of useless information in stochastic gradients. We also establish theoretical convergence guarantees for our algorithms. Empirical results on various datasets and models show that our algorithms yield much better performance than many state-of-the-art algorithms including SIGN-SGD, SIGNUM and ADAM. We also analyze the performance from multiple perspectives including the loss landscape and develop an adaptive method to further improve generalization. The source code is available at https://github.com/DongWanginxdu/signADAM-Learn-by-Confidence.

## KEYWORDS

Deep Learning, Optimization, Computer Vision, Image Classification, Confidence Function

---

[*]Equal contribution.
[†]Corresponding author.

---

## 1 INTRODUCTION

Deep neural networks (DNNs) are widely used in various fields of machine learning [19] (e.g., speech recognition, visual object recognition and object detection) and achieve huge success in all these applications. Moreover, training DNNs is usually considered as a non-convex optimization problem. Stochastic gradient descent (SGD) is one of the most effective algorithms to train various DNNs. It is a method to minimize an objective function, which is parameterized by the model's parameters $\theta \in \mathbb{R}^d$ ($d$ is the number of the parameters in a model). SGD updates the parameters through the negative directions of the gradients of an objective function.

SGD was first introduced by [30]. In recent years, a lot of work so far has focused on the adaptive methods. For example, the first algorithm, which achieves satisfactory performance in these fields, is ADAGRAD [11], and ADAGRAD performs well in sparse settings. However, in the high dimensional setting, the vanishing gradient problem becomes serious in the training of DNNs. To alleviate this issue, several algorithms based on the adaptive methods have been developed, such as ADADELTA [41], RMSPROP [38] and ADAM [17]. ADAM is one of the most popular optimization algorithms in deep learning. Based on the exponentially decaying average of past gradients, ADAM derives from the moment estimation and uses the first and second moments to compute individual learning rates for different parameters. However, the adaptive methods not only require the first moment estimate but also need the second moment estimate, which means extra computation amount. Recently, SIGN-SGD and SIGNUM [5] have been proposed. They applied the sign operation to SGD and its momentum variant, and also achieved a competitive performance in several benchmarks.

Hinton et al. [14] indicated that in some computer vision tasks especially image classification, the gradients mainly come from the incorrect samples and we can improve convergence rate by repeating
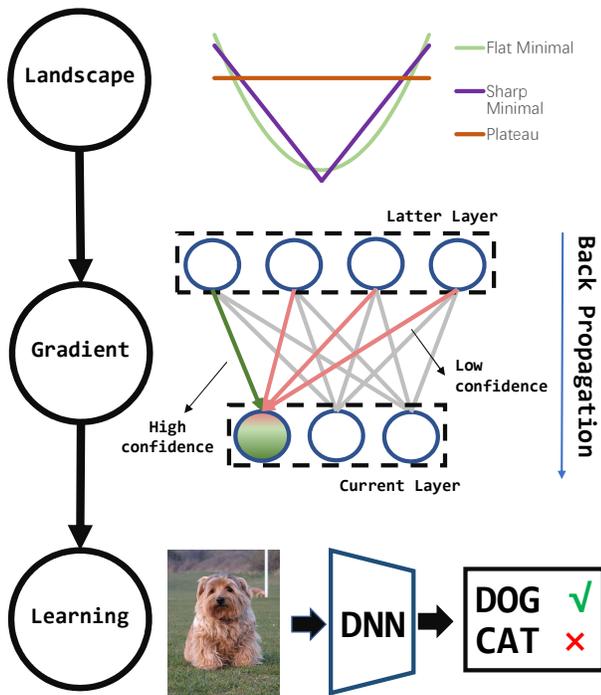
**Figure 1: The connection between the loss landscape and learning. Loss landscape reflects the degree of feature learning. The gradients in different loss landscapes are various, and thus we give them different confidences by our defined confidence function in the training stage to enhance feature learning. A small gradient means that the corresponding feature is easy to learn so it should get a small confidence. An unbalanced feature learning leads to a sharp minimal. If the models learn features too unequally which is caused by updating too much or stopping learning, the algorithm would convergence to a very sharp minimal around a plateau.**

sampling them. In other words, if we make the fitting rate of correct samples much slower than that of the incorrect samples, the network training will converge faster in fact. In practice, although adaptive methods have a faster convergence rate than SGD with momentum, SGD usually has a better performance on the test set than the others. We hold the view that the adaptive method gives the easy features a excessively fast fitting rate so that the models do not learn each useful feature of the dataset equally and reasonably. Because deep neural networks can not only learn the features that humans can easily understand, but also learn the features that are not easy to explain. So far, we have not been able to illustrate which features are really useful and which are useless. Inspired by the maximum entropy model [2], we tend to treat each feature equally. In terms of the difference of the fitting rate for each feature, the parameters which are used to extract and choose features should be not updated in a similar rate.

Motivated by all the above insights, the sign-based methods proposed by [5], and the work [3], we propose an efficient algorithm (called signADAM) instead of spending time in calculating the accurate moment information. We first integrate the sign operation to

ADAM and hope to combine the advantages of both sign-based methods and ADAM. Furthermore, we introduce the confidence function, which is used to distinguish different components of gradients, as shown in Figure 1. Therefore, we need a larger confidence for large components of gradients than the others. We multiply the gradients by its confidence element-wise. Considering that the second moment estimate in signADAM is close to a constant, we remove it and apply the confidence function to its gradients from mini-batch samples. Then we get a new algorithm called signADAM++.

Above all, we first put forward a sign-based method by applying the sign operation into ADAM. And by introducing the confidence function, signADAM is developed into signADAM++. It completely suits our new framework. In particular, our framework explained that the confidence function can help to distinguish which features are necessary and urgent to learn. Besides, this also meets biological principles. We also analyze the results by the loss landscape and reveal the relationship between feature learning and loss landscape.

The remainder of this paper is organized as follows. Section 2 discusses some recent advances in optimization algorithms for deep learning. In Section 3, we show our motivation by experimental results and propose a new optimization framework based on the confidence function. In Section 4, we reveal that our signADAM++ can be summarized into the framework. Empirically, our approach consistently achieves significantly better results than other methods on various datasets and models as shown in Section 6.

## 2 RELATED WORK

In this section, we mainly review the recently developed stochastic optimization techniques for deep learning.

### 2.1 Traditional Gradient-based Methods

There are several lines of research studies on the steepest gradient descent method. For training deep neural networks, stochastic gradient descent (SGD) is an efficient method. To solve the problem of navigating ravines [36], the momentum term [26] takes steps in relevant directions frequently. It dampens the oscillation impact along the short axis directions and makes contributions along the long axis directions [32].

Inspired by the DNNs' optimization development and theoretical advance, many adaptive algorithms have been proposed, such as ADAGRAD [11], ADADELTA [41], ADAM [17], Nadam [10] and AMSgrad [31]. They show great abilities to tackle the issues which have sparse data and non-stationary objectives. These methods attract great attention and they are successfully employed in several applications particularly popular in GANs and Q-learning. However, they also have been observed not to converge in some settings and circumstances. The linkages between adaptive methods especially ADAM and sign-based gradient descent approaches have begun to appear in recent work [3], which shows "sign" is a special case of ADAM under convex assumptions. They think ADAM combines two components: variance adaptation and taking sign, while the latter is a dominant one and they have proven it in greater details.

### 2.2 Quantization Methods

As both the model complexity and the amount of data increases, there are some prior works on quantization of models and optimization to
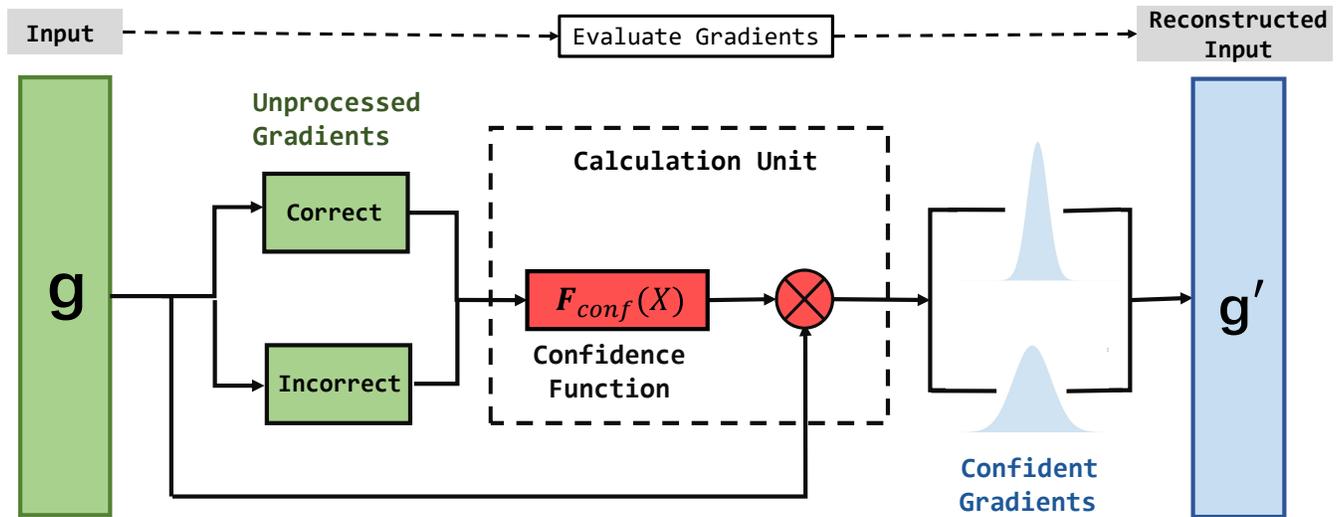
**Figure 2: Raw gradients have two parts coming from correct and incorrect samples, respectively. After being processed by the calculation unit, the gradients of correct samples have more zero components than those of the incorrect samples. The input "X" for the confidence function can be data (features), gradients, weights and so on. Because we think there is a close relationship between models, data and learning in deep neural networks. What we know about the optimization in deep learning is not only a non-convex problem, but also related to the model's structure property. In this paper, we show a kind of confidence functions with gradients as its input. So our confidence function is gradient-based.**

deal with large-scale problems. Inspired by delta-sigma modulation, Seide et al. [34] proposed a compression gradient method in SGD, which uses 1-bit to quantify those gradients during data exchange. This reduces the bandwidth requirement. Surprisingly, quantization can also improve the performance of ADAGRAD in the experiments, which gives us a higher accuracy and shorter training time compared with the common setting. De Sa et al. [7] developed a unified theoretical framework, which has proven that using low-precision arithmetic SGD can converge in a reasonable rate under convex assumptions. Moreover, in [1], they proposed a method called QSGD, which is applied in the non-convex case and allows a smooth trade-off between convergence and compression per iteration, and they also showed the effectiveness of gradient quantization. Another work [40] quantified gradients in a different way. It constrains gradients into three values $\{1, 0, -1\}$ to speed up distributed training. They proved that gradient clipping is necessary and effective.

The sign-based method was first introduced by [29]. They proposed an adaptive method (called RPROP) standing for "resilient propagation". According to the signs of gradients, RPROP adapts each element update magnitudes. Recently, with the theoretical development of sign-based optimization, [16] and [5] provided the proof of the sign-based method convergence in Polyak-Lojasiewicz condition and non-convex condition, respectively. Moreover, a more recent work [4] showed that the sign-based method outperforms existing quantization methods, for it does not need to ensure unbiasedness, while other methods requiring randomization.

## 2.3 What We Do

In this paper, we want to propose an efficient algorithm, which incorporates the sign operation into the adaptive gradient-based methods, e.g., ADAM. At a high level, we propose an efficient framework

of optimization algorithms, which evaluate gradients and update models depending on the confidence function setting the gradients to adaptive values in each iteration. This simple approach takes the advantages of a low cost complexity, strong performance guarantee and the adaptive property. For instance, as shown in Section 6, the experimental results reveal our algorithms perform even better than other adaptive and sign-based methods on a variety of both models (e.g., GoogLeNet [37], VGG-19 [35] and ResNet-18 [13]) and datasets (e.g., CIFAR-10 and CIFAR-100).

The main contributions of this paper can be summarized as follows:

(1) We first define a confidence function, which could be thought as augmented sampling and give a common optimization framework for training various DNNs, as shown in Figure 2.
(2) We then combine the sign operation and the momentum accelerated method with adaptive learning rates into non-convex optimization. The proposed method achieves a satisfactory result, as shown in Section 6.
(3) We provide our algorithms a theoretical guarantee to prove that they at least have the same convergence rate as SIGN-SGD in standard assumptions.
(4) We explain the experimental results from the perspective of the loss landscape and expose the relationship between loss landscape and feature learning.
(5) Experimental results show our motivation is reasonable and it suits the biological neural processes.

## 3 MOTIVATION

Backpropagation (BP) [33] is an essential tool to train various DNNs. Using this technique we could improve the accuracy by leveraging information from the training samples. Gradients come from models

and data. In this perspective, during updating, gradients are guiding models to extract the unknown features by learning from data. To learn all the features from the training set equally in a better way, we introduce the confidence function and show how it works in the experiment below.

## 3.1 How Our Confidence Function Works



**Figure 3: Remaining ratio changes as the confidence factor increases.**

**Set-up**   We run VGG-Net on CIFAR-10 and randomly choose a layer with many weights as samples. By using different confidence factors, we show the various remaining ratios for correct and incorrect samples, respectively. The definition of the remaining ratio as follows.

$$\text{Remaining Ratio} = \frac{\|grad\|_2}{\|grad'\|_2},$$

where $grad \in \mathbb{R}^d$ is the raw gradient, $grad'$ is the processed gradient, and $\| \cdot \|_2$ denotes the $L_2$-norm.

**Results**   When the confidence factor is set to a proper value, the remaining ratio of correct samples is obviously less than that of the incorrect samples. In addition, we find our confidence function weakens the components of the gradient from both correct and incorrect samples.

**Discussion**   Correct and incorrect samples have different features, which are easy and hard to learn, respectively. And correct samples have more easy features than incorrect samples. Our confidence function has a greater inhibition for the features, which are easy to learn. Also this makes the fitting rate of easy features faster than the hard ones, which lets the models learn the features equally. Of course, our confidence function has a similar effect as online hard sample mining [42]. Our improvement is given as follows:

- We first implement the augmented sampling in the optimization of deep learning by using our gradient-based confidence function.
- Most existing methods only deal with samples, while our method can deal with features. Because hard samples have features easy to learn and easy samples have features hard to learn.
- Our method can be more widely used in the deep learning community because it has a great performance on various datasets and models.

| | Confidence functions | Momentum |
|---|---|---|
| SIGN-SGD | $\frac{\delta}{|grad|}$ | × |
| SIGNUM | $\delta/|(1-\beta)\sum_i^k \beta^{t-i}grad_i|^*$ | ✓ |
| ADAM | $\delta/(1-\beta)\sum_{i=1}^k \beta^{k-i}grad_i \odot grad_i^*$ | ✓ |
| signADAM | $\delta/|grad|(1-\beta)\sum_{i=1}^k \beta^{k-i}{}^*$ | ✓ |
| signADAM++ | $F_{conf}(grad, \alpha) = \begin{cases} 0, & \text{if } grad \in [-\alpha, \alpha], \\ \frac{\delta}{|grad|}, & \text{otherwise} \end{cases}$ | ✓ |

**Table 1: Comparison of SIGN-SGD, SIGNUM, ADAM, our proposed signADAM and signADAM++ algorithms. Note that * means the above marked form is only approximation to the true confidence function because the algorithms first use the moving average. These methods are not applying the confidence function directly to the unprocessed gradients from the mini-batch samples, which is different from SIGN-SGD and our signADAM++ algorithm.**

## 3.2 Biological Perspective

For the view of neural processes, some neuroscience research [9] pointed out that cortical neurons are rarely active in their maximum saturation regime. And they suggest that neurons be encoded information in a sparse and distributed way [2]. Additionally, the research revealed that the brain uses a sparse representation: only around 1-4% of the neurons are active together at a fixed time [2, 20]. When a neuron is not fired, it does not change its state. Similarly, if we regard weights as neurons, neurons not activated represents weights not updated. We can implement this step by assigning small gradients a confidence with zero. In addition, the research from [21] shows that more independent genes may be the key life keep prosperity in the earth, while in deep learning, learning more distinguished features can make models more robust, which can weaken overfitting.

## 3.3 Our Framework Based on The Confidence Function

We define a new learning framework for training deep neural networks by using our confidence function. It is showed in Algorithm 1. We now consider the gradient descent as not only an optimization method but also a learning problem in deep neural networks. Our confidence function makes DNNs know how to learn features and which features to learn. So how to design a simple and efficient confidence function becomes attractive. Because such a new confidence function can help the model learn better features so that the model can get more satisfactory solutions. We also analyze some state-of-the-art optimization methods from the perspective of the confidence function. The analysis is showed in Table 1.

## 4 SIGNADAM++

In ADAM, the real learning rate is inversely proportional to a $L_2$-norm of their all current and past gradients. To extend the update rule from $L_2$-norm to $L_\infty$-norm, then the second moment estimate term becomes a recursive formula: $u_k = \max(\beta_2 u_{k-1}, |g_k|)$ in ADAMAX [17]. The update formula will have a sign term when $\beta_2 u_{k-1} \leqslant |g_k|$. Furthermore, the gradient descent using the sign can be viewed as the steepest descent with $L_\infty$-norm [6]. Balles and Hennig [3] proved that the ADAM's update rule is dominated by taking the sign of stochastic gradients. In addition, by setting the hyper-parameters $\beta_1$

---

**Algorithm 1** Generic framework based on our confidence function

---

**Input:** $\beta$ : Exponential decay rate for moment estimate.
**Initialize:** $\theta_0$ : Initial parameter vector,
  $\alpha$ : Initial confidence factor.
  **while** $k \neq$ max epoch **do**
    $k \leftarrow k + 1$;
    $\tilde{g}_k \leftarrow$ Stochastic Gradient($x_{k-1}$);
    $m_k = \beta m_{k-1} + F_{conf}(\alpha, \tilde{g}_k) \odot \tilde{g}_k$;
    $\theta_k \leftarrow \theta_{k-1} - m_k$;
  **end while**
**Output:** $\theta_k$.

---

and $\beta_2$ in ADAM to zero ($\beta_1, \beta_2 \to 0$), ADAM can be changed into SIGN-SGD [5].

In the optimization field, to combine the advantages of ADAM and SIGN-SGD, we first design an algorithm, called signADAM. Another idea of this algorithm is to reinforce the sign aspect in ADAM. Meanwhile, using "sign" can make communication and calculation effective. ADAM can be easily incorporated with "sign" merely by taking the "sign" of all stochastic gradients in ADAM.

---

**Algorithm 2** signADAM

---

**Input:** $\delta$ : Learning rate.
  $\beta_1, \beta_2 \in [0, 1)$ : Exponential decay rates for moment estimates,
  $\epsilon \leftarrow$ 1e-8.
  $\theta_0$ : Initial parameter vector.
  $m_0 \leftarrow 0, v_0 \leftarrow 0, k \leftarrow 0$.
  Current point $x_k$, current first and second moment $m_k, v_k$.
  **while** $k \neq$ max epoch **do**
    $k \leftarrow k + 1$;
    $\tilde{g}_k \leftarrow$ Stochastic Gradient($x_{k-1}$);
    $\tilde{g}_k \leftarrow sign(\tilde{g}_k)$;
    $m_k \leftarrow \beta_1 m_{k-1} + (1 - \beta_1)\tilde{g}_k$; //(update first moment estimate)
    $v_k \leftarrow \beta_2 v_{k-1} + (1 - \beta_2)\tilde{g}_k \odot \tilde{g}_k$; //(update second moment estimate)
    $\theta_k \leftarrow \theta_{k-1} - \delta \frac{m_k}{\sqrt{v_k} + \epsilon}$; //(update parameters)
  **end while**
**Output:** $\theta_k$.

---

For each dimension of the gradients, the second moment estimate term in our signADAM algorithm is a value determined by the iteration $k$. Since $[sign(\tilde{g}_k)]^2$ is a fixed number equal to 1, so the term $v_k$ can be rewritten as follows:

$$v_k = -\beta_2^{k-1} + 1.$$

As the global iteration $k$ increases, the magnitude of the difference between 1 and $v_k$ would become infinitesimal ($k \to \infty, v_k \to 1$). In numerical the effect of this term can be neglected. Thus the new parameter update rule is formulated as follows:

$$\theta_k \leftarrow \theta_{k-1} - \delta m_k.$$

Then we analyze ADAM and SIGN-SGD from the perspective of the confidence function. ADAM gives an approximate adaptive confidence to all gradients. In contrast, SIGN-SGD sets the confidence to 0 only when the gradients equal to 0 and gives other

gradients a real adaptive confidence. Bernstein et al. [5] pointed out that SIGN-SGD belongs to the family of the adaptive methods. This confirms our adaptive confidence on large gradients for both ADAM and SIGN-SGD. We have shown the confidence functions of all the algorithms in Table 1.

If we set an appropriate confidence factor in our confidence function (as shown in Figure 3), there will be different effects on the unprocessed gradients from both correct and incorrect samples. The confidence factor can be used to control the sparsity of the gradients. Considering the motivations in Section 3 and the framework based on our confidence function, we design a new confidence function and apply it into signADAM. The new algorithm is called signADAM++.

We use the following confidence function in the proposed signADAM++ algorithm:

$$F_{conf}(grad, \alpha) = \begin{cases} 0, & \text{if } grad \in [-\alpha, \alpha], \\ \frac{\delta}{|grad|}, & \text{otherwise,} \end{cases}$$

where $\alpha$ is the confidence factor.

From a computational point of view, choosing the sign operation is appealing for the following reasons.

- Tasks based on deep learning such as image classification usually have a large amount of training data (e.g., the popular image dataset, ImageNet, has 1.5 millions images). "Training" on such dataset may take a week even a month to shrink calculation in every step which can enhance the efficiency.

- The algorithm can induce sparsity of gradients by not updating some gradients. signADAM++ makes gradients sparse and leads models to learn distinguished information. Reddi et al. [27] indicated that it is important to update large and informative gradients. One of the ADAM's drawbacks is to ignore these gradients by using the exponential average algorithm [17]. Our algorithm can prevent this issue by giving small gradients a very low confidence (e.g., 0).

- Nowadays, learning from samples in deep learning is always stochastic due to the giant amount of data and models. In fact, the confidence of some gradients produced by loss functions and samples should be low. From the perspective of maximum entropy theory, each feature should have the equal right to make efforts in deep neural networks. This method will make the models work well.

- Although there are large gradients caused by the incorrect samples, large gradients may hurt models' generalization and learning ability found by [24] and [27]. signADAM++ addresses this issue by using moving average after applying confidence for unprocessed gradients and an adaptive confidence for some large gradients. These can give models a better performance.

## 5 CONVERGENCE ANALYSIS

In the training process of deep neural networks, given a finite length sequence of input and output pairs, the goal is to minimize the following empirical loss,

$$f(\theta) := \frac{1}{N} \sum_{i=1}^{N} f_i(\theta), \tag{1}$$

---

**Algorithm 3** signADAM++

---

**Input:** $\delta$ : Learning rate.

$\quad\quad \beta \in [0, 1)$ : Exponential decay rates for moment estimates.

$\quad\quad \alpha$ : Confidence factor.

**Initialize:** $\theta_0$ : Initial parameter vector, $m_0 \leftarrow 0$ (Initial moment vector), $k \leftarrow 0$ (Initial global step).

$\quad$ Current point $x_k$, and current momentum $m_k$.

$\quad$ **while** $k \neq$ max epoch **do**

$\quad\quad k \leftarrow k + 1;$

$\quad\quad \tilde{g}_k \leftarrow$ Stochastic Gradient($x_{k-1}$);

$\quad\quad$ **if** $|\tilde{g}_k| < \alpha$ **then**

$\quad\quad\quad \tilde{g}_k \leftarrow 0;$

$\quad\quad$ **else**

$\quad\quad\quad \tilde{g}_k \leftarrow sign(\tilde{g}_k);$

$\quad\quad$ **end if**

$\quad\quad m_k \leftarrow \beta m_{k-1} + (1 - \beta)\tilde{g}_k;$ //(update first moment estimate)

$\quad\quad \theta_k \leftarrow \theta_{k-1} - \delta m_k;$ //(update parameters)

$\quad$ **end while**

**Output:** $\theta_k$.

---

where the loss function $f_i(\theta)$ is discrepancy between the model output and ground truth. A neural network can be regarded as a function with the parameter $\theta \in \mathbb{R}^d$ and produces an output for each input $\xi^i$. We aim to analyze the objective function $f(\theta, \xi^i)$, which is a non-convex function of its parameters. Thus the optimization problem is given by

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} f(\theta). \tag{2}$$

To simplify the problem, we make several assumptions. It would not loss anything compared to typical SGD assumptions, since these assumptions are obtained from SGD.

ASSUMPTION 1 (GLOBAL MINIMAL). *For all parameters $\theta \in \mathbb{R}^d$, there exists a d-dimensional vector $\theta^*$ that satisfies the following inequality*

$$f(\theta) \geqslant f(\theta^*). \tag{3}$$

ASSUMPTION 2 (L-SMOOTH). *Let $g(x)$ denote the real gradient of the objective function $f(\cdot)$ evaluated at point $x$. Then $\forall x, y$, we require that for some non-negative constants $\vec{L} := [L_1, \cdots, L_d]$,*

$$\left| f(y) - [f(x) + g(x)^T(y - x)] \right| \leqslant \frac{1}{2} \sum_i L_i(y_i - x_i)^2. \tag{4}$$

*For the twice differentiable function $f(\cdot)$, this implies that $-diag(\vec{L}) \prec H \prec diag(\vec{L})$. This is related to the slightly weaker coordinate-wise Lipschitz condition used in the block coordinate descent literature[28].*

ASSUMPTION 3 (VARIANCE BOUND). *Since $x \in \mathbb{R}^d$ is stochastic, its homologous gradient $\tilde{g}$, which has coordinate bounded variance, is an independent unbiased estimate for the real gradient $g$.*

$$\mathbb{E}[\tilde{g}(x)] = g(x), \ \mathbb{E}[(\tilde{g}(x) - g(x))^2] \leqslant \sigma_i{}^2, \tag{5}$$

*for a vector containing non-negative constants $\vec{\sigma} := [\sigma_1, \cdots, \sigma_d]$*

Assumptions 2 and 3 are the standard assumptions of Lipschitz smoothness and total bounded variance, respectively. Under these assumptions and the framework of SIGN-SGD [5], we have the following results.

THEOREM 1 (NON-CONVEX CONVERGENCE RATE). *Suppose Assumptions 1, 2 and 3 hold. Let $\{\tilde{g}_k\}$ ($k = 0, \cdots, K$) be a sequence generated by Algorithm 1, where the learning rate is set to $\delta_k := 1/\sqrt{\|\vec{L}\|_1}K$ and the size of mini-batch as $n_k = K$. Let N be the cumulative number of stochastic gradient calls up to step K (i.e.,$N = O(K^2)$). Then we have*

$$\mathbb{E}\left[ \frac{1}{K} \sum_{k=0}^{K-1} \|g_k\|_1 \right]^2$$

$$\leqslant \frac{1}{\sqrt{N}} \left[ \sqrt{\|\vec{L}\|_1} \left( f(\theta_0) - f(\theta^*) + \frac{1}{2} + \frac{8\beta}{(1-\beta)^2} \right) + 2\|\vec{\sigma}\|_1 \right]^2.$$

To prove Theorem 1, we first present and prove the following lemmas. In the following lemma, we apply mathematical induction for exponential average of the gradients' signs.

LEMMA 1. *For all $m_k \in \mathbb{R}^d$, they satisfy the following inequality*

$$\left| m_{k,i} \right| \leqslant 1, \ \ i \in \{1, 2, \cdots, d\}, \tag{6}$$

*where $\left| m_{k,i} \right|$ is the exponential average of the gradients' sign.*

The detailed proofs of Theorem 1, Lemma 1 and the lemma below are provided in the APPENDIX. Next we give a lemma, which shows the exponential average of gradients' signs.

LEMMA 2. $\frac{1-\beta}{1-\beta^k} \sum_{t=0}^{K-1} \beta^t [\|g_{k-t} - g_k\|_1] \leq 4 \frac{\sqrt{\|\vec{L}\|_1}}{\sqrt{K}} \frac{\beta}{(1-\beta)^2}.$

# 6 EXPERIMENTS

In this section, we discuss our empirical evaluation of signADAM++ and other related algorithms such as ADAM [17], SIGN-SGD and SIGNUM [5]. We study the classical deep learning problems (e.g., image classification) and conduct several experiments for different state-of-the-art DNNs such as GoogLeNet [37], VGG-19 [35], ResNet-18 [13] and SeNet-18 [15]. Using the torch-rnn library, we also train and test an LSTM network.
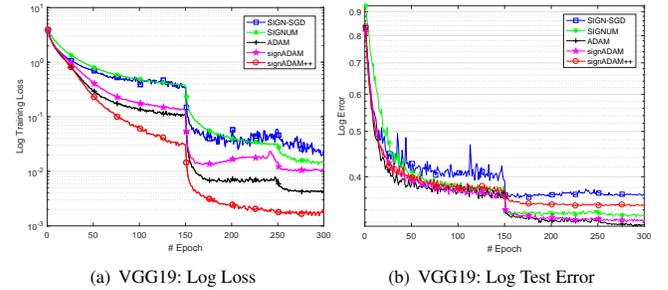


(a) VGG19: Log Loss $\quad\quad$ (b) VGG19: Log Test Error

**Figure 4: Comparison of the training loss and test error of all the algorithms on CIFAR-100.**

## 6.1 Setup

We compare the proposed algorithms (i.e., signADAM and signADAM++) with SIGN-SGD [5], SIGNUM [5] and ADAM [17]. We run all the experiments on the three standard image datasets:

MNIST, CIFAR-10 and CIFAR-100, all of which are divided into two parts: the training set and test set.

In the experiments on CIFAR-10 and CIFAR-100, all the training images are resized to $32 \times 32 \times 3$. These images are then randomly cropped to $32 \times 32 \times 3$, with padding equals to 4. They are also randomly flipped and finally normalized. For the test images, they are only normalized and no other data augmentation tricks are used in the training and test processes. In the experiments on MNIST, the training and test images are both normalized. We record the objective loss on the training set and test error on the test set and evaluated algorithms in the following deep learning models.

- LeNet consists of two sets of convolutional and average pooling layers, followed by a flattening layer, then two fully-connected layers and finally a softmax classifier.
- VGG-19 uses 16 convolution layers and 3 fully connected layers.
- GoogLeNet consists of 22 layers, using the Inception module to make itself deeper than VGG.
- ResNet-18 increases depth of the network results in less extra parameters and reduces the effect of gradient vanishing problems. ResNets achieve the improvement by adding simple skip connections.
- SENet-18 is an architectural using dynamic channel-wise feature recalibration, which can improve the representational power of a network.
- Two layers of LSTM units for the classification of MNIST.

We first test LeNet and LSTM on the MNIST dataset, which can be viewed as shallow neural networks. Then to verify our algorithms for training DNNs, we evaluate the performance of all the algorithms on the stat-of-the-art DNNs (e.g., VGG-19, GoogLeNet, ResNet-18 and SeNet-18). In addition, to prove that our algorithms work well on different datasets, we also design experiments on the CIFAR-100 dataset by running VGG-19. For all experiments, we have used optimal hyperparameters for all algorithms and applied learning rate schedulers for all algorithms. We also implement the $L_2$ normalization proposed by [22] for all methods. And the hyperparameters are given as follows. In experiments on MNIST, we use 0.001 as the learning rate of SIGN-SGD, SIGNUM and ADAM. On CIFAR-10 and CIFAR-100, we use 0.0001 as their learning rate. In all the experiments, we use the recommended initial values for $\beta_1$ and $\beta_2$ (i.e., $\beta_1 = 0.9$, $\beta_2 = 0.999$) and set the regularization coefficient to 5e-4.

## 6.2 Experimental Results

Figures 4, 6 and 7 show the experimental results of all the algorithms. From all the results, we make the following observations.

### 6.2.1 *Sparsity*.

***Figure 5 shows the sparsity of the confidence gradients***. These parameters are sampled from a randomly chosen layer of VGG-19 on the CIFAR-10 dataset. We find that the gradients are indeed sparse when applying our confidence function into the unprocessed gradients. It fits the biological neural process as mentioned in Section 3. By observing the training loss on CIFAR-10, our signADAM++ enjoys faster convergence than the other state-of-the-art algorithms. Besides, the effect of our confidence factor is similar as the the
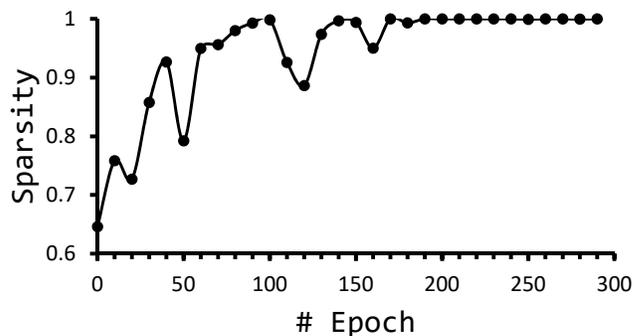


**Figure 5: Sparsity of confidence gradients.**

hard threshold methods in [8], which is used to tackle with Wavelet noises. The effectiveness of sparsity can be also confirmed by ReLU activations, as shown in [12, 18, 25].

Our method is to regard optimization as a better learning strategy for models. We decrease the randomness created by the stochastic mini-batch, which can increase the convergence rate. As to the feature learning, we tend to make the models learn all the features equally [39].

### 6.2.2 *Training Loss*.

***signADAM++ has the fastest convergence speed***. The training loss of signADAM++ is distinguished from that of SIGN-SGD, SIGNUM, ADAM and signADAM during the entire experiments. signADAM++ makes rapid initial progress, and achieves lower loss than other algorithms. There is a big reduction at epoch 150 after we apply learning-rate-decay to all the algorithms. It means that signADAM++ can reach our expectation in a shorter time. We infer that our signADAM++ has removed some randomness so that the parameters can be updated in main and correct directions. This can speed up the training of deep neural networks. In other words, thanks to our reasonable confidence function, the significant samples and their corresponding features have been learned by neural networks preferentially.

***signADAM performs like SIGNUM***. We notice that the curves of signADAM and SIGNUM are close especially in the experiments on CIFAR-10 and MNIST. There are two following main reasons:

- signADAM and SIGNUM both have moving average and sign operation. Although moving average indeed can help accelerate in the relevant directions and dampen oscillations, some small oscillations cannot be completely removed because its momentum coefficient is not adaptive. These small oscillations are increased by the sign operation. Similarly, signADAM strengthens the small oscillations by using the sign operation and dampens oscillations by using the moving average.
- Considering that the high precision of the gradients, the square of the sign for a gradient is close to a constant as the iteration increases. So we can merge the second order estimate into the learning rate. Our experimental results indeed confirm this statement.
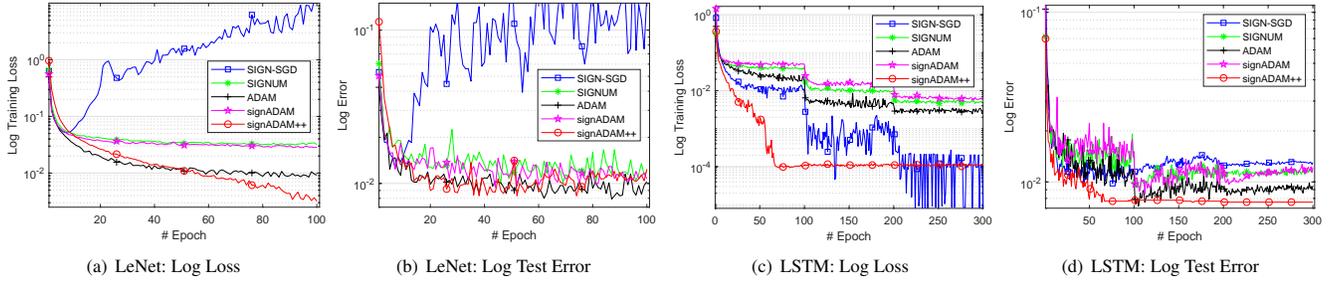
(a) LeNet: Log Loss      (b) LeNet: Log Test Error      (c) LSTM: Log Loss      (d) LSTM: Log Test Error

**Figure 6: Comparison of the logarithm training loss and test error of all the algorithms on MNIST.**



(a) VGG19: Log Loss      (b) VGG19: Log Test Error      (c) GoogLeNet: Log Loss      (d) GoogLeNet: Log Test Error

(e) ResNet18: Log Loss      (f) ResNet18: Log Test Error      (g) SENet18: Log Loss      (h) SENet18: Log Test Error
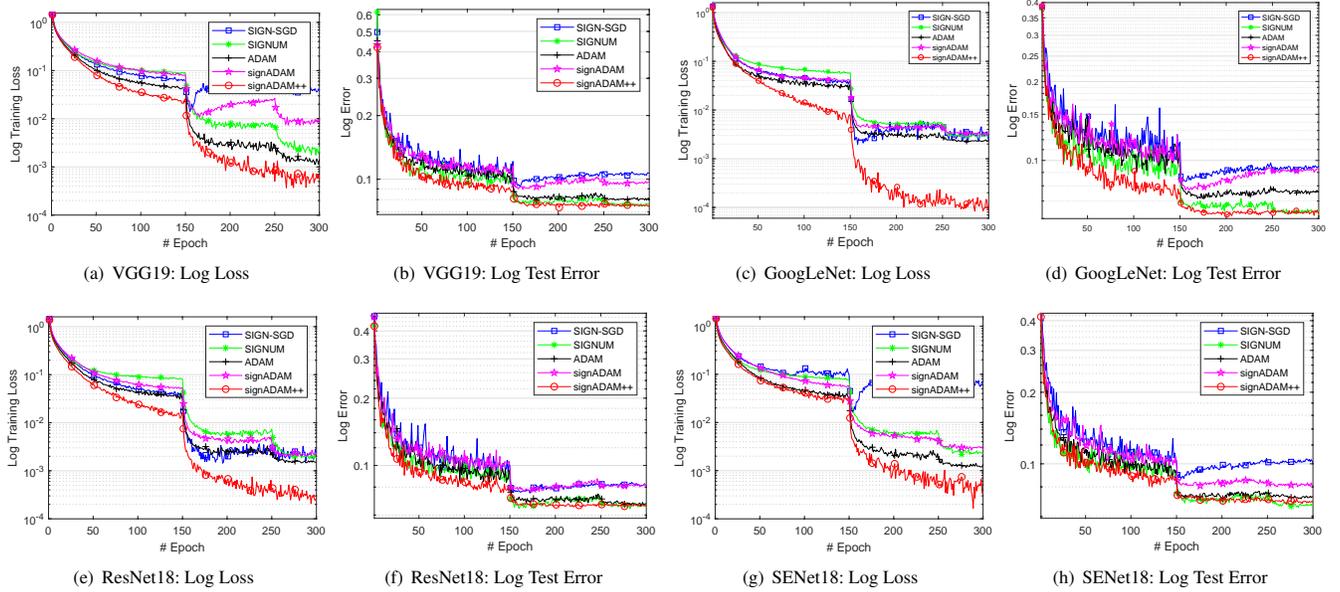
**Figure 7: Comparison of the logarithm training loss and test error of all the algorithms on CIFAR-10.**

***SIGN-SGD has large oscillations***. SIGN-SGD updates parameters by using the signs of its gradient. This means that SIGN-SGD does not distinguish the components of gradients. It makes neural networks learn some samples and features so fast that SIGN-SGD finds a local minimal with higher loss.

### 6.2.3 *Test Error*.

***signADAM++ has a competitive optimization rate.*** On the test set, we find that signADAM++ also has a satisfactory convergence rate. Table 2 shows the numbers of epochs for all the algorithms when the error reaches a given tolerance. In the experiments of MNIST and CIFAR-10, signADAM++ performs well. However, signADAM++ performs slightly worse on CIFAR-100. We give the following explanation: Because there are 100 classes in CIFAR-100, each class in CIFAR-100 has obviously fewer samples than that in CIFAR-10. So we do not have enough samples for each class to learn. Thanks to the adaptive property of ADAM and SIGNUM, they can both continue learning some features from samples regardless the confidence of information. But for signADAM++, because we apply a very low confidence (e.g., 0) into some information, this makes some features hard even impossible to be learned so specific

as CIFAR-10. In other words, due to the lack of data, features on CIFAR-100 are not treated as equally as on CIFAR-10. This result in that the optimal point is located around a sharp minimal, which is usually considered to have a worse generalization than a flat minimal. The adaptive methods can have a large step size in some directions to skip out of the sharp minimal points. And it learns some features relatively sufficient. So the confidence factor decides which features to learn. To keep the models learning features from samples all the time and finding some flatter minimal points, we develop an adaptive method to solve the problem in the next subsection.

### 6.3   Our Adaptive and Weight-decay Technique

Considering that the "stop learning" issue above, we improve our algorithm and introduce the adaptive confidence function. In particular, we set the confidence factor in the confidence function to a dynamic value, which is determined by the standard deviation of the gradients. So the new confidence function is defined as follows.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\sqrt{var(grad)},$$
$$\alpha_t = \beta_2 m_t,$$

| | MNIST (0.01) | | CIFAR-10 (0.1) | | | | CIFAR-100 (0.4) |
|---|---|---|---|---|---|---|---|
| | LeNet | LSTM | VGG-19 | GoogLeNet | RseNet-18 | SENet-18 | VGG-19 |
| SIGN-SGD | - | 44 | 151 | 31 | 49 | 83 | 91 |
| SIGNUM | 46 | 101 | 72 | 39 | 76 | 71 | 48 |
| ADAM | **20** | 58 | 123 | **23** | 41 | 45 | **31** |
| signADAM | 44 | 101 | 151 | 32 | 57 | 69 | 39 |
| signADAM++ | **20** | **21** | **50** | 24 | **37** | **41** | 46 |

**Table 2: The number of epochs for all algorithms to reach a constant error. For MNIST, the error tolerance is 0.01, while it is 0.1 and 0.4 for CIFAR-10 and CIFAR-100, respectively. Note that "-" means that the algorithm does not achieve the expected error in the entire run.**

| Methods | Accuracy (%) | # Epoch |
|---|---|---|
| SIGN-SGD | 63.96% | 166 |
| SIGNUM | 67.12% | 291 |
| ADAM | 68.62% | 288 |
| signADAM | 67.97% | 227 |
| signADAM++ | 65.38% | 285 |
| ADAMW | 69.11% | 291 |
| signADAM++-AW | **72.31**% | 278 |

**Table 3: Comparison of the accuracies and needed epochs of all the algorithms on CIFAR-100. It is clear that the accuracy of signADAM++ is significantly better than those of the other algorithms on CIFAR-100.**

where $\beta_1$ and $\beta_2$ are the hyper-parameters in the moving average, and $\alpha$ is the confidence factor. By introducing this new adaptive technique, there are always some parameters being updated. These parameters are used to extract and choose the most distinguished features. This keeps the models learning from the samples. Besides, for the objective function with the traditional $L_2$-norm, there exists some components from the regularization term, which can disturb the sampling gradients. Our goal is to make models continue learning the features equally. So we do not need the traditional $L_2$-norm. Instead, we replace the $L_2$-norm using weight-decay [23]. This insight is called the "AW" method. Table 3 shows the highest accuracy, which is achieved by each algorithm.

## 7 CONCLUSIONS

In this paper, we first proposed an efficient sign-based gradient descent optimization method for non-convex objective functions. Our method combines the low computation complexity advantage of SIGN-SGD and the adaptive property of ADAM. In particular, our algorithms incorporated the sign operation into ADAM. Besides, we defined the confidence function to make the model learn features equally. The motivation is supported by experimental results. signADAM++ will not be constrained by the depth of neural networks and has some obvious advantages over some state-of-the-art algorithms including SIGN-SGD, SIGNUM and ADAM. Some intuitive analysis shows that signADAM++ has much faster convergence than ADAM, which is a popular optimization algorithm for deep learning. We analyzed the convergence properties of our algorithms and the experimental results also confirmed our analysis. Furthermore, signADAM++ suits our new optimization framework based on the confidence function. In conclusion, we found that signADAM++ is a robust algorithm for training various deep neural networks. We

also proposed a framework on the confidence function and reveal the relationship between the loss landscape and feature learning.

## REFERENCES

[1] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. 2017. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *NeurIPS*. 1709–1720.
[2] David Attwell and Simon B Laughlin. 2001. An energy budget for signaling in the grey matter of the brain. *Journal of Cerebral Blood Flow & Metabolism* 21, 10 (2001), 1133–1145.
[3] Lukas Balles and Philipp Hennig. 2017. Dissecting adam: The sign, magnitude and variance of stochastic gradients. *arXiv preprint arXiv:1705.07774* (2017).
[4] Jeremy Bernstein, Kamyar Azizzadenesheli, Yu-Xiang Wang, and Anima Anandkumar. 2018. Convergence rate of sign stochastic gradient descent for non-convex functions. (2018).
[5] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Anima Anandkumar. 2018. signSGD: Compressed optimisation for non-convex problems. *arXiv preprint arXiv:1802.04434* (2018).
[6] Stephen Boyd and Lieven Vandenberghe. 2004. *Convex optimization*. Cambridge university press.
[7] Christopher M De Sa, Ce Zhang, Kunle Olukotun, and Christopher Ré. 2015. Taming the wild: A unified analysis of hogwild-style algorithms. In *NeurIPS*. 2674–2682.
[8] David L Donoho and Jain M Johnstone. 1994. Ideal spatial adaptation by wavelet shrinkage. *biometrika* 81, 3 (1994), 425–455.
[9] Rodney J Douglas and Kevan AC Martin. 2007. Recurrent neuronal circuits in the neocortex. *Current biology* 17, 13 (2007), R496–R500.
[10] Timothy Dozat. 2016. Incorporating Nesterov Momentum into Adam. In *ICLR Workshop*. 2013–2016.
[11] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12, Jul (2011), 2121–2159.
[12] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Deep sparse rectifier neural networks. In *AIStats*. 315–323.
[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
[14] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
[15] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.

[16] Hamed Karimi, Julie Nutini, and Mark Schmidt. 2016. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *ECML-PKDD*. Springer, 795–811.

[17] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[18] John Langford, Lihong Li, and Tong Zhang. 2009. Sparse online learning via truncated gradient. *JMLR* 10, Mar (2009), 777–801.

[19] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436.

[20] Peter Lennie. 2003. The cost of cortical computation. *Current biology* 13, 6 (2003), 493–497.

[21] Adi Livnat, Christos Papadimitriou, Jonathan Dushoff, and Marcus W Feldman. 2008. A mixability theory for the role of sex in evolution. *Proceedings of the National Academy of Sciences* 105, 50 (2008), 19803–19808.

[22] Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101* (2017).

[23] Ilya Loshchilov and Frank Hutter. 2017. Fixing Weight Decay Regularization in Adam. *CoRR* abs/1711.05101 (2017). arXiv:1711.05101 http://arxiv.org/abs/1711.05101

[24] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. 2019. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843* (2019).

[25] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *ICML*. 807–814.

[26] Ning Qian. 1999. On the momentum term in gradient descent learning algorithms. *Neural networks* 12, 1 (1999), 145–151.

[27] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. 2018. On the convergence of adam and beyond. (2018).

[28] Peter Richtárik and Martin Takáč. 2014. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming* 144, 1-2 (2014), 1–38.

[29] Martin Riedmiller and Heinrich Braun. 1993. A direct adaptive method for faster backpropagation learning: The RPROP algorithm. In *Proceedings of the IEEE international conference on neural networks*, Vol. 1993. San Francisco, 586–591.

[30] Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. *The annals of mathematical statistics* (1951), 400–407.

[31] Sebastian Ruder. 2016. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747* (2016).

[32] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. 1985. *Learning internal representations by error propagation*. Technical Report. California Univ San Diego La Jolla Inst for Cognitive Science.

[33] David E Rumelhart, Geoffrey E Hinton, Ronald J Williams, et al. 1988. Learning representations by back-propagating errors. *Cognitive modeling* 5, 3 (1988), 1.

[34] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 2014. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *Interspeech*.

[35] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[36] Richard Sutton. 1986. Two problems with back propagation and other steepest descent learning procedures for networks. In *Proceedings of the Eighth Annual Conference of the Cognitive Science Society, 1986*. 823–832.

[37] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1–9.

[38] T. Tieleman and G. Hinton. 2012. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning.

[39] Naftali Tishby and Noga Zaslavsky. 2015. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*. IEEE, 1–5.

[40] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. 2017. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In *NeurIPS*. 1509–1519.

[41] Matthew D Zeiler. 2012. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701* (2012).

[42] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. 2016. Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks. *IEEE Signal Processing Letters* 23, 10 (Oct 2016), 1499–1503. https://doi.org/10.1109/LSP.2016.2603342

# APPENDIX A: PROOF OF LEMMA 1

PROOF.

$$\forall i, \quad m_{1,i} = 0, \quad m_{2,i} \in \{\pm(1-\beta), 0\}.$$

Since $\beta \in (0, 1)$, then $|m_{2,i}| \leqslant 1$. Assume $\forall i, |m_{k,i}| \leqslant 1$, then we have

$$\forall i, |m_{k+1,i}| = |\beta \cdot m_{k,i} + (1-\beta)sign(\tilde{g}_{k,i})|$$
$$\leqslant \beta \cdot |m_{k,i}| + (1-\beta) \cdot |sign(\tilde{g}_{k,i})|$$
$$\leqslant \beta + 1 - \beta = 1.$$

This completes the proof.  □

# APPENDIX B: PROOF OF LEMMA 2

PROOF. Using Lemma F.3 in SIGNSGD [5], then we can have the following conclusion: Under Assumption 2, for any $s \in R^d$, $(|s_i \leq 1|, i = 1, 2, \cdots, d)$, any x $\in R^d$ and any $\epsilon \leq \delta$, then

$$\|g(x + \epsilon s) - g(x)\|_1 \leq 2\epsilon \|\vec{L}\|_1.$$

We also have

$$\|g_{k-t} - g_k\|_1 \leq \sum_{l=0}^{t-1} \|g_{k-l} - g_{k-l-1}\|_1$$

$$\leq 2\|\vec{L}\|_1 \sum_{l=0}^{t-1} \delta_{k-l-1} \leq \sum_{l=0}^{t-1} \frac{2\|\vec{L}\|_1 \delta}{\sqrt{k-l}}$$

$$\leq 2\|\vec{L}\|_1 \delta \int_{k-t}^{k} \frac{dx}{\sqrt{x}} = 4\|\vec{L}\|_1 \delta(\sqrt{k} - \sqrt{k-t})$$

$$\leq 4\|\vec{L}\|_1 \delta \sqrt{k}\left(1 - \sqrt{1 - \frac{t}{k}}\right)$$

$$\leq 4\|\vec{L}\|_1 \delta \frac{t}{\sqrt{k}}.$$

Let $\delta = \frac{1}{\sqrt{\|\vec{L}\|_1}}$, then

$$\frac{1-\beta}{1-\beta^k} \sum_{t=0}^{K-1} \beta^t [\|g_{k-t} - g_k\|_1] \leq \frac{1-\beta}{1-\beta} \sum_{t=0}^{\infty} \beta^t [\|g_{k-t} - g_k\|_1]$$

$$\leq 4\frac{\sqrt{\|\vec{L}\|_1}}{\sqrt{k}} \frac{\beta}{(1-\beta)^2}.$$

This completes the proof.  □

# APPENDIX C: PROOF OF THEOREM 1

PROOF. Suppose Assumptions 1, 2 and 3 hold. Following the proof non-convex SGD, we modify Assumption 2 as follows.

$$f_{k+1} - f_k \leqslant g_k^T(x_{k+1} - x_k) + \sum_{i=1}^{d} \frac{L_i}{2}(x_{k+1} - x_k)_i^2$$

$$= -\delta_k g_k^T m_k + \delta_k^2 \sum_{i=1}^{d} \frac{L_i}{2} m_{k,i}^2.$$

Then we apply Lemma1 (i.e., $|m_{k,i}| \leqslant 1$) into the above inequality and can get

$$f_{k+1} - f_k \leqslant -\delta_k g_k^T m_k + \delta_k^2 \sum_{i=1}^{d} \frac{L_i}{2}.$$

Now we calculate the mathematical expectation of both sides of the above inequality.

$$\mathbb{E}[f_{k+1} - f_k | x_k] \leqslant -\delta_k \mathbb{E}[g_k{}^T m_k] + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$= -\delta_k \mathbb{E}[g_k{}^T \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t sign(\tilde{g}_{k-t})] + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$= -\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \mathbb{E} \left\{ \beta^t \left[ \|g_k\|_1 - 2 \sum_{i=1}^{d} |g_{k,i}| \mathbb{I}(sign(g_{k,i} \neq sign(\tilde{g}_{k-t,i})) \right] \right\} + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$= -\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \left( \beta^t \left[ \|g_k\|_1 - 2 \sum_{i=1}^{d} |g_{k,i}| \mathbb{P}(sign(g_{k,i} \neq sign(\tilde{g}_{k-t,i})) \right] \right) + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1.$$

$$\mathbb{E}[f_{k+1} - f_k | x_k] \leqslant -\delta_k \|g_k\|_1 + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t \sum_{i=1}^{d} |g_{k,i}| \mathbb{P}[|\tilde{g}_{k-t,i} - g_{k,i}| \geqslant g_{k,i}] + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$\leqslant -\delta_k \|g_k\|_1 + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t \sum_{i=1}^{d} \mathbb{E}[|\tilde{g}_{k-t,i} - g_{k,i}|] + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$\leqslant -\delta_k \|g_k\|_1 + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t \sum_{i=1}^{d} \left( \mathbb{E}[|\tilde{g}_{k-t,i} - g_{k-t,i}|] + \mathbb{E}[|g_{k-t,i} - g_{k,i}|] \right) + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$\leqslant -\delta_k \|g_k\|_1 + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t \sum_{i=1}^{d} \left( \sqrt{\mathbb{E}[|\tilde{g}_{k-t,i} - g_{k-t,i}|^2]} + \mathbb{E}[|g_{k-t,i} - g_{k,i}|] \right) + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$\leqslant -\delta_k \|g_k\|_1 + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t \frac{\|\vec{\sigma}\|_1}{\sqrt{n_k}} + 2\delta_k \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t [\|g_{k-t} - g_k\|_1] + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1.$$

Note that $\mathbb{I}(\cdot) = 1$ when $(\cdot)$ holds, otherwise $\mathbb{I}(\cdot) = 0$. $\mathbb{P}(\cdot)$ stands for the probability of $(\cdot)$.

Using Lemma 2, we can obtain

$$\mathbb{E}[f_{k+1} - f_k | x_k] \leqslant -\delta_k \|g_k\|_1 + 2\delta_k \left( \frac{\|\vec{\sigma}\|_1}{\sqrt{n_k}} + \frac{1-\beta}{1-\beta^k} \sum_{t=0}^{k-1} \beta^t [\|g_{k-t} - g_k\|_1] \right) + \frac{\delta_k^2}{2} \left\| \vec{L} \right\|_1$$

$$\leqslant -\frac{1}{\sqrt{\|\vec{L}\|_1 K}} \|g_k\|_1 + \frac{2\|\vec{\sigma}\|_1}{\sqrt{\|\vec{L}\|_1 K}} + \frac{1}{K} \frac{8\beta}{(1-\beta)^2} + \frac{1}{2K}.$$

Since

$$f(\theta_0) - f(\theta^*) \geq f(\theta_0) - \mathbb{E}[f_K]$$

$$= \mathbb{E} \left[ \sum_{k=0}^{K-1} f_k - f_{k+1} \right]$$

$$\geq \mathbb{E} \sum_{k=0}^{K-1} \left[ \frac{1}{\sqrt{\|\vec{L}\|_1 K}} \|g_k\|_1 - \frac{1}{2\sqrt{\|\vec{L}\|_1 K}} \left( 4\|\vec{\sigma}\|_1 + \sqrt{\|\vec{L}\|_1} + \frac{16\beta\sqrt{\|\vec{L}\|_1}}{(1-\beta)^2} \right) \right]$$

$$= \sqrt{\frac{K}{\|\vec{L}\|_1}} \mathbb{E} \left[ \frac{1}{K} \sum_{k=0}^{K-1} \|g_k\|_1 \right] - \frac{1}{2\sqrt{\|\vec{L}\|_1}} \left( 4\|\vec{\sigma}\|_1 + \sqrt{\|\vec{L}\|_1} + \frac{16\beta\sqrt{\|\vec{L}\|_1}}{(1-\beta)^2} \right),$$

then

$$\mathbb{E} \left[ \frac{1}{K} \sum_{k=0}^{K-1} \|g_k\|_1 \right] \leqslant \frac{1}{\sqrt{K}} \left[ \sqrt{\|\vec{L}\|_1} \left( f(\theta_0) - f(\theta^*) + \frac{1}{2} + \frac{8\beta}{(1-\beta)^2} \right) + 2\|\vec{\sigma}\|_1 \right].$$

To square both sides of the above inequality, and we complete the proof.                    □