
AdaSmooth: An Adaptive Learning Rate Method based on Effective Ratio

Jun Lu¹

Abstract

It is well known that we need to choose the hyper-parameters in Momentum, AdaGrad, AdaDelta, and other alternative stochastic optimizers. While in many cases, the hyper-parameters are tuned tediously based on experience becoming more of an art than science. We present a novel per-dimension learning rate method for gradient descent called AdaSmooth. The method is insensitive to hyper-parameters thus it requires no manual tuning of the hyper-parameters like Momentum, AdaGrad, and AdaDelta methods. We show promising results compared to other methods on different convolutional neural networks, multi-layer perceptron, and alternative machine learning tasks. Empirical results demonstrate that AdaSmooth works well in practice and compares favorably to other stochastic optimization methods in neural networks.

1. Introduction

Over the years, stochastic gradient-based optimization has become a core method in many fields of science and engineering such as computer vision and automatic speech recognition processing (Krizhevsky et al., 2012; Hinton et al., 2012a; Graves et al., 2013). Stochastic gradient descent (SGD) and deep neural network (DNN) play a core role in training stochastic objective functions. When a new deep neural network is developed for a given task, some hyper-parameters related to the training of the network must be chosen heuristically. For each possible combination of structural hyper-parameters, a new network is typically trained from scratch and evaluated over and over again. While much progress has been made on hardware (e.g. Graphical Processing Units) and software (e.g. cuDNN) to speed up the training time of a single structure of a DNN, the exploration of a large set of possible structures remains very slow making the need of a stochastic optimizer that is insensitive to hyper-parameters.

¹Correspondence to: Jun Lu <jun.lu.locky@gmail.com>. Copyright 2022 by the author(s)/owner(s). April 2nd, 2022.

1.1. Gradient Descent

Gradient descent (GD) is one of the most popular algorithms to perform optimization and by far the most common way to optimize machine learning tasks. And this is particularly true for optimizing neural networks. The neural networks or machine learning in general find the set of parameters $\mathbf{x} \in \mathbb{R}^d$ in order to optimize an objective function $L(\mathbf{x})$. The gradient descent finds a sequence of parameters

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T, \quad (1)$$

such that when $T \rightarrow \infty$, the objective function $L(\mathbf{x}_T)$ achieves the optimal minimum value. At each iteration t , a step $\Delta \mathbf{x}_t$ is applied to change the parameters. Denoting the parameters at the t -th iteration as \mathbf{x}_t . Then the update rule becomes

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \Delta \mathbf{x}_t. \quad (2)$$

The most naive method of stochastic gradient descent is the vanilla update: the parameter moves in the opposite direction of the gradient which finds the steepest descent direction since the gradients are orthogonal to level curves (a.k.a., level surface, see Lemma 16.4 in Lu (2022b)):

$$\Delta \mathbf{x}_t = -\eta \mathbf{g}_t = -\eta \frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t} = -\eta \nabla L(\mathbf{x}_t), \quad (3)$$

where the positive value η is the learning rate and depends on specific problems, and $\mathbf{g}_t = \frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t} \in \mathbb{R}^d$ is the gradient of the parameters. The learning rate η controls how large of a step to take in the direction of negative gradient so that we can reach a (local) minimum. While if we follow the negative gradient of a single sample or a batch of samples iteratively, the local estimate of the direction can be obtained and is known as the stochastic gradient descent (SGD) (Robbins & Monro, 1951). In the SGD framework, the objective function is stochastic that is composed of a sum of subfunctions evaluated at different subsamples of the data.

For a small step-size, gradient descent makes a monotonic improvement at every iteration. Thus, it always converges, albeit to a local minimum. However, the speed of the vanilla GD method is usually slow, while it can take an exponential rate when the curvature condition is poor. While choosing higher than this rate may cause the procedure to diverge

in terms of the objective function. Determining a good learning rate (either global or per-dimension) becomes more of an art than science for many problems. Previous work has been done to alleviate the need for selecting a global learning rate (Zeiler, 2012), while it is still sensitive to other hyper-parameters.

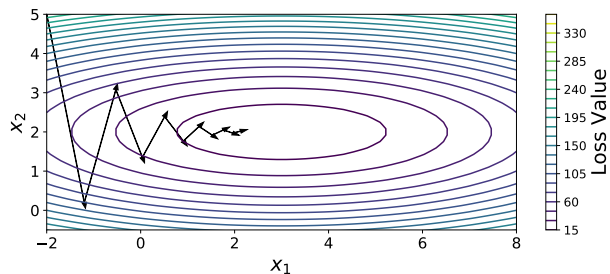
The main contribution of this paper is to propose a novel stochastic optimization method which is insensitive to different choices of hyper-parameters resulting in adaptive learning rates, and naturally performs a form of step size annealing. We propose the AdaSmooth algorithm to both increase optimization efficiency and out-of-sample accuracy. While previous works propose somewhat algorithms that are insensitive to the global learning rate (e.g., (Zeiler, 2012)), the methods are still sensitive to hyper-parameters that influence per-dimension learning rates. The proposed AdaSmooth, a method for efficient stochastic optimization that only requires first-order gradients and (accumulated) past update steps with little memory requirement, allows flexible and adaptive per-dimension learning rates. Meanwhile, the method is memory efficient and easy to implement. Our method is designed to combine the advantages of the following methods: AdaGrad (Duchi et al., 2011), which works well with sparse gradients, RMSProp (Hinton et al., 2012b), which works well in online and non-stationary settings, and AdaDelta (Zeiler, 2012), which is less sensitive in global learning rate.

2. Related Work

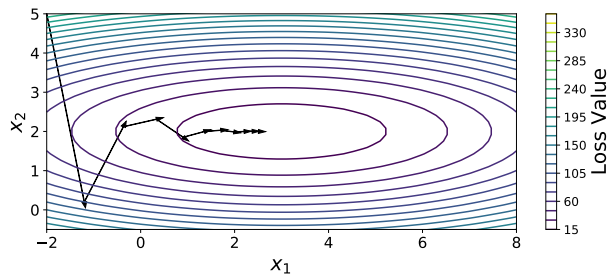
There are several variants of gradient descent to use heuristics for estimating a good learning rate at each iteration of the progress. These methods either attempt to accelerate learning when suitable or to slow down learning near a local minima (Zeiler, 2012; Kingma & Ba, 2014; Ruder, 2016).

2.1. Momentum

If the cost surface is not spherical, learning can be quite slow because the learning rate must be kept small to prevent divergence along the steep curvature directions (Rumelhart et al., 1986; Qian, 1999; Sutskever et al., 2013). The SGD with Momentum (that can be applied to full batch or mini-batch learning) attempts to use previous step to speed up learning when suitable such that it enjoys better converge rates on deep networks. The main idea behind the Momentum method is to speed up the learning along dimensions where the gradient consistently point in the same direction; and to slow the pace along dimensions in which the sign of the gradient continues to change. Figure 1(a) shows a set of updates for vanilla GD where we can find the update along dimension x_1 is consistent; and the move along dimension x_2 continues to change in a zigzag pattern. The GD with Momentum keeps track of past parameter updates with an



(a) Optimization without Momentum. A higher learning rate may result in larger parameter updates in dimension across the valley (direction of x_2) which could lead to oscillations back as forth across the valley.



(b) Optimization with Momentum. Though the gradients along the valley (direction of x_1) are much smaller than the gradients across the valley (direction of x_2), they are typically in the same direction and thus the momentum term accumulates to speed up movement, dampen oscillations and cause us to barrel through narrow valleys, small humps and (local) minima.

Figure 1. A 2-dimensional convex function $L(\mathbf{x}) = 2(x_1 - 3)^2 + 20(x_2 - 2)^2 + 5$ and $\frac{\partial L(\mathbf{x})}{\partial \mathbf{x}} = (4x_1 - 12, 8x_2 - 16)^\top$. Starting point to descent is $(-2, 5)^\top$.

exponential decay, and the update method has the following step:

$$\Delta \mathbf{x}_t = \rho \Delta \mathbf{x}_{t-1} - \eta \frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t}, \quad (4)$$

where the algorithm remembers the latest update and adds it to present update by multiplying a parameter ρ called *momentum parameter*. That is, the amount we change the parameter is proportional to the negative gradient plus the previous weight change; the added momentum term acts as both a smoother and an accelerator. The momentum parameter ρ works as a *decay constant* where $\Delta \mathbf{x}_1$ may have effect on $\Delta \mathbf{x}_{100}$; however, its effect is decayed by this decay constant. In practice, the momentum parameter ρ is usually set to be 0.9 by default. Momentum simulates the concept inertia in physics. It means that in each iteration, the update mechanism is not only related to the gradient descent, which refers to the dynamic term, but also maintains a component which is related to the direction of last update iteration, which refers to the momentum.

The momentum works extremely better in ravine-shaped loss curve. Ravine is an area, where the surface curves are much more steeply in one dimension than in another (see the

loss curve in Figure 1, i.e., a long narrow valley). Ravines are common near local minima in deep neural networks and vanilla SGD has trouble navigating them. As shown by the toy example in Figure 1(a), SGD will tend to oscillate across the narrow ravine since the negative gradient will point down one of the steep sides rather than along the ravine towards the optimum. Momentum helps accelerate gradients in the correct direction (Figure 1(b)).

2.2. AdaGrad

The learning rate annealing procedure modifies a single global learning rate that applies to all dimensions of the parameters (Smith, 2017). Duchi et al. (2011) proposed a method called AdaGrad where the learning rate is updated on a per-dimension basis. The learning rate for each parameter depends on the history of gradient updates of that parameter in a way such that parameters with a scarce history of updates are updated faster using a larger learning rate. In other words, parameters that have not been updated much in the past are more likely to have higher learning rates now. Denoting the element-wise vector multiplication between \mathbf{a} and \mathbf{b} by $\mathbf{a} \odot \mathbf{b}$, formally, the AdaGrad has the following update step:

$$\Delta \mathbf{x}_t = -\frac{\eta}{\sqrt{\sum_{\tau=1}^t \mathbf{g}_\tau^2 + \epsilon}} \odot \mathbf{g}_t, \quad (5)$$

where ϵ is a smoothing term to better condition the division, η is a global learning rate shared by all dimensions, \mathbf{g}_τ^2 indicates the element-wise square $\mathbf{g}_\tau \odot \mathbf{g}_\tau$, and the denominator computes the l_2 norm of sum of all previous squared gradients in a per-dimension fashion. Though the η is shared by all dimensions, each dimension has its own dynamic learning rate controlled by the l_2 norm of accumulated gradient magnitudes. Since this dynamic learning rate grows with the inverse of the accumulated gradient magnitudes, larger gradient magnitudes have smaller learning rates and smaller absolute values of gradients have larger learning rates. Therefore, the accumulated gradient in the denominator has the same effects as the learning rate annealing.

One pro of the AdaGrad method is that it partly eliminates the need to tune the learning rate controlled by the accumulated gradient magnitude. However, AdaGrad’s main weakness is its accumulation of the squared gradients in the denominator. Since every added term is positive, the accumulated sum keeps growing or exploding during every training step. This in turn causes the per-dimension learning rate to shrink and eventually decrease throughout training and become infinitesimally small, eventually falling to zero and stopping training any more. Moreover, since the magnitudes of gradients are factored out in AdaGrad, this method can be sensitive to the initialization of the parameters and the corresponding gradients. If the initial magnitudes of the gradients are large or infinitesimally huge, the per-dimension

learning rates will be low for the remainder of training. This can be partly combated by increasing the global learning rate, making the AdaGrad method sensitive to the choice of learning rate. Further, since AdaGrad assumes the parameter with fewer updates should favor a larger learning rate; and one with more movement should employ a smaller learning rate. This makes it consider only the information from squared gradients, or the absolute value of the gradients. And thus AdaGrad does not include information from total move (i.e., the sum of updates; not the sum of absolute updates).

To be more succinct, AdaGrad has the following main drawbacks: 1) the continual decay of learning rates throughout training; 2) the need for a manually selected global learning rate; 3) considering only the absolute value of gradients.

2.3. AdaDelta

AdaDelta is an extension of AdaGrad that overcomes the main weakness of AdaGrad (Zeiler, 2012). The original idea of AdaDelta is simple: it restricts the window of accumulated past gradients to some fixed size w rather than t (i.e., current time step). However, since storing w previous squared gradients is inefficient, the AdaDelta introduced in Zeiler (2012) implements this accumulation as an exponentially decaying average of the squared gradients. This is very similar to the idea of momentum term (or decay constant).

2.3.1. ADADELTA: FORM 1 (RMSPROP)

We first discuss the exact form of the window AdaGrad (we here call it AdaGradWin for short). Assume at time t this running average is $E[\mathbf{g}^2]_t$ then we compute:

$$E[\mathbf{g}^2]_t = \rho E[\mathbf{g}^2]_{t-1} + (1 - \rho) \mathbf{g}_t^2, \quad (6)$$

where ρ is a decay constant similar to that used in the momentum method and \mathbf{g}_t^2 indicates the element-wise square $\mathbf{g}_t \odot \mathbf{g}_t$.

As Eq (6) is just the root mean squared (RMS) error criterion of the gradients, we can replace it with the criterion shorthand. Let $\text{RMS}[\mathbf{g}]_t = \sqrt{E[\mathbf{g}^2]_t + \epsilon}$, where again a constant ϵ is added to better condition the denominator. Then the resulting step size can be obtained as follows:

$$\Delta \mathbf{x}_t = -\frac{\eta}{\text{RMS}[\mathbf{g}]_t} \odot \mathbf{g}_t, \quad (7)$$

where again \odot is the element-wise vector multiplication.

As aforementioned, the form in Eq (6) is originally from the exponential moving average (EMA). In the original form of EMA, $1 - \rho$ is also known as the smoothing constant (SC) where the SC can be written as $\frac{2}{N+1}$ and the period N can be thought of as the number of past values to do the moving

average calculation (Lu, 2022a):

$$\text{SC} = 1 - \rho = \frac{2}{N + 1}. \quad (8)$$

The above Eq (8) links different variables: the decay constant ρ , the smoothing constant (SC), and the period N . If $\rho = 0.9$, then $N = 19$. That is, roughly speaking, $E[\mathbf{g}^2]_t$ at iteration t is approximately equal to the moving average of past 19 squared gradients and the current one (i.e., moving average of 20 squared gradients totally). The relationship in Eq (8) though is not discussed in Zeiler (2012), it is important to decide the lower bound of the decay constant ρ . Typically, a time period of $N = 3$ or 7 is thought to be a relatively small frame making the lower bound of decay constant $\rho = 0.5$ or 0.75; when $N \rightarrow \infty$, the decay constant ρ approaches 1.

However, we can find that the AdaGradWin still only considers the absolute value of gradients and a fixed number of past squared gradients is not flexible which can cause a small learning rate near (local) minima as we will discuss in the sequel.

RMSProp The AdaGradWin is actually the same as the RMSProp method developed independently by Geoff Hinton in Hinton et al. (2012b) both of which are stemming from the need to resolve AdaGrad’s radically diminishing per-dimension learning rates. Hinton et al. (2012b) suggests ρ to be set to 0.9 and the global learning rate η to be 0.001 by default.

2.3.2. ADADELTA: FORM 2

Zeiler (2012) shows the units of the step size shown above do not match (so as the vanilla SGD, the momentum, and the AdaGrad). To overcome this weakness, from the correctness of the second order method, the author considers to rearrange Hessian to determine the quantities involved. It is well known that, though the calculation of Hessian or approximation to the Hessian matrix is a tedious and computationally expensive task, its curvature information is useful for optimization, and the units in Newton’s method are well matched. Given the Hessian matrix \mathbf{H} , the update step in Newton’s method can be described as follows (Becker & Le Cun, 1988; Dauphin et al., 2014):

$$\Delta \mathbf{x}_t \propto -\mathbf{H}^{-1} \mathbf{g}_t \propto \frac{\frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t}}{\frac{\partial^2 L(\mathbf{x}_t)}{\partial \mathbf{x}^2}}. \quad (9)$$

This implies

$$\frac{1}{\frac{\partial^2 L(\mathbf{x}_t)}{\partial \mathbf{x}_t^2}} = \frac{\Delta \mathbf{x}_t}{\frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t}}, \quad (10)$$

i.e., the units of the Hessian matrix can be approximated by the right-hand side term of the above equation. Since the

RMSProp update in Eq (7) already involves $\text{RMS}[\mathbf{g}]_t$ in the denominator, i.e., the units of the gradients. Putting another unit of the order of $\Delta \mathbf{x}_t$ in the numerator can match the same order as Newton’s method. To do this, define another exponentially decaying average of the update steps:

$$\begin{aligned} \text{RMS}[\Delta \mathbf{x}]_t &= \sqrt{E[\Delta \mathbf{x}^2]_t} \\ &= \sqrt{\rho E[\Delta \mathbf{x}^2]_{t-1} + (1 - \rho) \Delta \mathbf{x}_t^2}. \end{aligned} \quad (11)$$

Since $\Delta \mathbf{x}_t$ for the current iteration is not known and the curvature can be assumed to be locally smoothed making it suitable to approximate $\text{RMS}[\Delta \mathbf{x}]_t$ by $\text{RMS}[\Delta \mathbf{x}]_{t-1}$. So we can use an estimation of $\frac{1}{\frac{\partial^2 L(\mathbf{x}_t)}{\partial \mathbf{x}_t^2}}$ to replace the computationally expensive \mathbf{H}^{-1} :

$$\frac{\Delta \mathbf{x}_t}{\frac{\partial L(\mathbf{x}_t)}{\partial \mathbf{x}_t}} \sim \frac{\text{RMS}[\Delta \mathbf{x}]_{t-1}}{\text{RMS}[\mathbf{g}]_t}. \quad (12)$$

This is an approximation to the diagonal Hessian using only RMS measures of \mathbf{g} and $\Delta \mathbf{x}$, and results in the update step whose units are matched:

$$\Delta \mathbf{x}_t = -\frac{\text{RMS}[\Delta \mathbf{x}]_{t-1}}{\text{RMS}[\mathbf{g}]_t} \odot \mathbf{g}_t. \quad (13)$$

The idea of AdaDelta from the second method overcomes the annoying choosing of learning rate. Similarly, in Kingma & Ba (2014), an exponentially decaying average is incorporated into the gradient information such that convergence in online convex setting is improved. Meanwhile, the diminishing problem was also attacked in second order methods (Schaul et al., 2013). The idea of averaging gradient or its alternative is not unique, previously in Moulines & Bach (2011), Polyak-Ruppert averaging has been shown to improve the convergence of vanilla SGD (Ruppert, 1988; Polyak & Juditsky, 1992). Other stochastic optimization methods, including the natural Newton method, AdaMax, Nadam, all set the step-size per-dimension by estimating curvature from first-order information (Le Roux & Fitzgibbon, 2010; Kingma & Ba, 2014; Dozat, 2016). The LAMB adopts layerwise normalization due to layerwise adaptivity (You et al., 2019) and we shall not go into the details.

3. AdaSmooth Method

In this section we will discuss the effective ratio based on previous updates in the stochastic optimization process and how to apply it to accomplish adaptive learning rates per-dimension via the flexible smoothing constant, hence the name AdaSmooth. The idea presented in the paper is derived from the first form of AdaDelta (Zeiler, 2012) in order to improve two main drawbacks of the method: 1) consider only the absolute value of the gradients rather than the total movement in each dimension; 2) the need for manually selected hyper-parameters.

3.1. Effective Ratio (ER)

Kaufman (2013; 1995) suggested replacing the smoothing constant in the EMA formula with a constant based on the *efficiency ratio* (ER). And the ER is shown to provide promising results for financial forecasting via classic quantitative strategies (Lu, 2022a) where the ER of the closing price is calculated to decide the trend of the asset. This indicator is designed to measure the *strength of a trend*, defined within a range from -1.0 to +1.0 where the larger magnitude indicates a larger upward or downward trend. Recently, Lu & Yi (2022) shows the ER can be utilized to reduce over-estimation and underestimation in time series forecasting. Given the window size M and a series $\{h_1, h_2, \dots, h_T\}$, it is calculated with a simple formula:

$$e_t = \frac{s_t}{n_t} = \frac{h_t - h_{t-M}}{\sum_{i=0}^{M-1} |h_{t-i} - h_{t-1-i}|} \quad (14)$$

$$= \frac{\text{Total move for a period}}{\text{Sum of absolute move for each bar}},$$

where e_t is the ER of the series at time t . At a strong trend (i.e., the input series is moving in a certain direction, either up or down) the ER will tend to 1 in absolute value; if there is no directed movement, it will be a little more than 0.

Instead of calculating the ER of the closing price of asset, we want to calculate the ER of the moving direction in the update methods for each parameter. And in the descent methods, we care more about how much each parameter moves apart from its initial point in each period, either move positively or negatively. So here we only consider the absolute value of the ER. To be specific, the ER for the parameters in the proposed method is calculated as follows:

$$e_t = \frac{s_t}{n_t} = \frac{|\mathbf{x}_t - \mathbf{x}_{t-M}|}{\sum_{i=0}^{M-1} |\mathbf{x}_{t-i} - \mathbf{x}_{t-1-i}|} \quad (15)$$

$$= \frac{|\sum_{i=0}^{M-1} \Delta \mathbf{x}_{t-1-i}|}{\sum_{i=0}^{M-1} |\Delta \mathbf{x}_{t-1-i}|},$$

where $e_t \in \mathbb{R}^d$ whose i -th element $e_{t,i}$ is in the range of $[0, 1]$ for all i in $[1, 2, \dots, d]$. A larger value of $e_{t,i}$ indicates the descent method in the i -th dimension is moving in a certain direction; while a smaller value approaching 0 means the parameter in the i -th dimension is moving in a zigzag pattern, interleaved by positive and negative movement. In practice, and in all of our experiments, the M is selected to be the batch index for each epoch. That is, $M = 1$ if the training is in the first batch of each epoch; and $M = M_{\max}$ if the training is in the last batch of the epoch where M_{\max} is the maximal number of batches per epoch. In other words, M ranges in $[1, M_{\max}]$ for each epoch. Therefore, the value of $e_{t,i}$ indicates the movement of the i -th parameter in the most recent epoch. Or even more aggressively, the window can range from 0 to the total number of batches

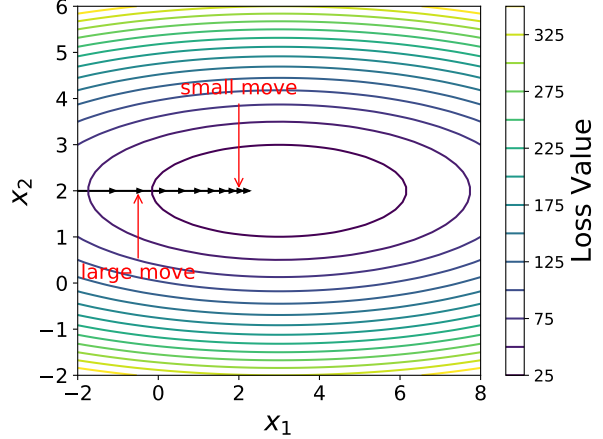


Figure 2. Demonstration of how the effective ratio works. Stochastic optimization tends to move a large step when it is far from the (local) minima; and a relatively small step when it is close to the (local) minima.

seen during the whole training progress. The adoption of the adaptive window size M rather than a fixed one has a benefit that we do not need to keep the past $M + 1$ steps $\{\mathbf{x}_{t-M}, \mathbf{x}_{t-M+1}, \dots, \mathbf{x}_t\}$ to calculate the signal and noise vectors $\{s_t, n_t\}$ in Eq (15) since they can be obtained in an accumulated fashion.

3.2. AdaSmooth

If the ER in magnitude of each parameter is small (approaching 0), the movement in this dimension is zigzag, the proposed AdaSmooth method tends to use a long period average as the scaling constant to slow down the movement in that dimension. When the absolute ER per-dimension is large (tend to 1), the path in that dimension is moving in a certain direction (not zigzag), and the learning actually is happening and the descent is moving in a correct direction where the learning rate should be assigned to a relatively large value for that dimension. Thus the AdaSmooth tends to choose a small period which leads to a small compensation in the denominator; since the gradients in the closer periods are small when it's near the (local) minima. A particular example is shown in Figure 2, where the descent is moving in a certain direction, and the gradient in the near periods is small; if we choose a larger period to compensate for the denominator, the descent will be slower due to the large factored denominator. In short, we want a smaller period to calculate the exponential average of the squared gradients in Eq (6) if the update is moving in a certain direction without a zigzag pattern; while when the parameter is updated in a zigzag basis, the period for the exponential average should be larger.

The obtained value of ER is used in the exponential smooth-

ing formula. Now, what we want to go further is to set the time period N discussed in Eq (8) to be a smaller value when the ER tends to 1 in absolute value; or a larger value when the ER moves towards 0. When N is small, SC is known as a “fast SC”; otherwise, SC is known as a “slow SC”.

For example, let the small time period be $N_1 = 3$, and the large time period be $N_2 = 199$. The smoothing ratio for the fast movement must be as for EMA with period N_1 (“fast SC” = $\frac{2}{N_1+1} = 0.5$), and for the period of no trend EMA period must be equal to N_2 (“slow SC” = $\frac{2}{N_2+1} = 0.01$). Thus the new changing smoothing constant is introduced, called the “scaled smoothing constant” (SSC), denoted by a vector $\mathbf{c}_t \in \mathbb{R}^d$:

$$\mathbf{c}_t = (\text{fast SC} - \text{slow SC}) \times \mathbf{e}_t + \text{slow SC}. \quad (16)$$

By Eq (8), we can define the *fast decay constant* $\rho_1 = 1 - \frac{2}{N_1+1}$, and the *slow decay constant* $\rho_2 = 1 - \frac{2}{N_2+1}$. Then the scaled smoothing constant vector can be obtained by:

$$\mathbf{c}_t = (\rho_2 - \rho_1) \times \mathbf{e}_t + (1 - \rho_2), \quad (17)$$

where the smaller \mathbf{e}_t , the smaller \mathbf{c}_t . For a more efficient influence of the obtained smoothing constant on the averaging period, Kaufman recommended squaring it. The final calculation formula then follows:

$$E[\mathbf{g}^2]_t = \mathbf{c}_t^2 \odot \mathbf{g}_t^2 + (1 - \mathbf{c}_t^2) \odot E[\mathbf{g}^2]_{t-1}. \quad (18)$$

We notice that $N_1 = 3$ is a small period to calculate the average (i.e., $\rho_1 = 1 - \frac{2}{N_1+1} = 0.5$) such that the EMA sequence will be noisy if N_1 is less than 3. Therefore, the minimal value of ρ_1 in practice is set to be larger than 0.5 by default. While $N_2 = 199$ is a large period to compute the average (i.e., $\rho_2 = 1 - \frac{2}{N_2+1} = 0.99$) such that the EMA sequence almost depends only on the previous value leading to the default value of ρ_2 no larger than 0.99. Experiment study will show that the AdaSmooth update will be insensitive to the hyper-parameters in the sequel. We also carefully notice that when $\rho_1 = \rho_2$, the AdaSmooth algorithm recovers to the RMSProp algorithm with decay constant $\rho = 1 - (1 - \rho_2)^2$ since we square it in Eq (18). After developing the AdaSmooth method, we realize the main idea behind it is similar to that of SGD with Momentum: to speed up (compensate less in the denominator) the learning along dimensions where the gradient consistently points in the same direction; and to slow the pace (compensate more in the denominator) along dimensions in which the sign of the gradient continues to change.

Empirical evidence shows the ER used in simple moving average with a fixed windows size w can also reflect the trend of the series/movement in quantitative strategies (Lu, 2022a). However, this again needs to store w previous squared gradients in the AdaSmooth case, making it inefficient and we shall not adopt this extension.

3.3. AdaSmoothDelta

Notice the ER can also be applied to the AdaDelta setting:

$$\Delta \mathbf{x}_t = -\frac{\sqrt{E[\Delta \mathbf{x}^2]_t}}{\sqrt{E[\mathbf{g}^2]_t + \epsilon}} \odot \mathbf{g}_t, \quad (19)$$

where

$$E[\mathbf{g}^2]_t = \mathbf{c}_t^2 \odot \mathbf{g}_t^2 + (1 - \mathbf{c}_t^2) \odot E[\mathbf{g}^2]_{t-1}, \quad (20)$$

and

$$E[\Delta \mathbf{x}^2]_t = (1 - \mathbf{c}_t^2) \odot \Delta \mathbf{x}_t^2 + \mathbf{c}_t^2 \odot E[\Delta \mathbf{x}^2]_{t-1}, \quad (21)$$

in which case the difference in $E[\Delta \mathbf{x}^2]_t$ is to choose a larger period when the ER is small. This is reasonable in the sense that $E[\Delta \mathbf{x}^2]_t$ appears in the numerator while $E[\mathbf{g}^2]_t$ is in the denominator of Eq (19) making their compensation towards different directions. Or even, a fixed decay constant can be applied for $E[\Delta \mathbf{x}^2]_t$:

$$E[\Delta \mathbf{x}^2]_t = (1 - \rho_2) \Delta \mathbf{x}_t^2 + \rho_2 E[\Delta \mathbf{x}^2]_{t-1}, \quad (22)$$

The AdaSmoothDelta optimizer introduced above further alleviates the need for a hand specified global learning rate which is set to $\eta = 1$ from the Hessian context. However, due to the adaptive smoothing constants in Eq (20) and (21), the $E[\mathbf{g}^2]_t$ and $E[\Delta \mathbf{x}^2]_t$ are less locally smooth making it less insensitive to the global learning rate than the AdaDelta method. Therefore, a smaller global learning rate, e.g., $\eta = 0.5$ is favored in AdaSmoothDelta. The full procedure for computing AdaSmooth is then formulated in Algorithm 1.

4. Experiments

To evaluate the strategy and demonstrate the main advantages of the proposed AdaSmooth method, we conduct experiments with different machine learning models; and different data sets including real handwritten digit classification task, MNIST (LeCun, 1998)¹, and Census Income² data sets are used. In all scenarios, same parameter initialization is adopted when training with different stochastic optimization algorithms. We compare the results in terms of convergence speed and generalization. In a wide range of scenarios across various models, AdaSmooth improves optimization rates, and leads to out-of-sample performances that are as good or better than existing stochastic optimization algorithms in terms of loss and accuracy.

4.1. Experiment: Multi-Layer Perceptron

Multi-layer perceptrons (MLP, a.k.a., multi-layer neural networks) are powerful tools for solving machine learning

¹It has a training set of 60,000 examples, and a test set of 10,000 examples.

²Census income data has 48842 number of samples where 70% of them are used as training set in our case: <https://archive.ics.uci.edu/ml/datasets/Census+Income>.

Algorithm 1 Computing AdaSmooth at iteration t : the proposed AdaSmooth algorithm. All operations on vectors are element-wise. Good default settings for the tested tasks are $\rho_1 = 0.5, \rho_2 = 0.99, \epsilon = 1e-6, \eta = 0.001$; see Section 3.2 or Eq (8) for a detailed discussion on the explanation of the decay constants' default values. Empirical study in Section 4 shows that the AdaSmooth algorithm is not sensitive to the hyperparameter ρ_2 , while $\rho_1 = 0.5$ is relatively a lower bound in this setting. The AdaSmoothDelta iteration can be calculated in a similar way.

- 1: **Input:** initial parameter \mathbf{x}_1 , Constant ϵ ;
- 2: **Input:** global learning rate η , by default $\eta = 0.001$;
- 3: **Input:** fast decay constant ρ_1 , slow decay constant ρ_2 ;
- 4: **Input:** assert $\rho_2 > \rho_1$, by default $\rho_1 = 0.5, \rho_2 = 0.99$;
- 5: **for** $t = 1 : T$ **do**
- 6: Compute gradient $\mathbf{g}_t = \nabla L(\mathbf{x}_t)$;
- 7: Compute ER $e_t = \frac{|\mathbf{x}_t - \mathbf{x}_{t-M}|}{\sum_{i=0}^{M-1} |\Delta \mathbf{x}_{t-1-i}|}$;
- 8: Compute smoothing $\mathbf{c}_t = (\rho_2 - \rho_1) \times e_t + (1 - \rho_2)$;
- 9: Compute normalization term:

$$E[\mathbf{g}^2]_t = \mathbf{c}_t^2 \odot \mathbf{g}_t^2 + (1 - \mathbf{c}_t^2) \odot E[\mathbf{g}^2]_{t-1};$$

- 10: Compute step $\Delta \mathbf{x}_t = -\frac{\eta}{\sqrt{E[\mathbf{g}^2]_t + \epsilon}} \odot \mathbf{g}_t$;
- 11: Apply update $\mathbf{x}_t = \mathbf{x}_{t-1} + \Delta \mathbf{x}_t$;
- 12: **end for**
- 13: **Return:** resulting parameters \mathbf{x}_t , and the loss $L(\mathbf{x}_t)$.

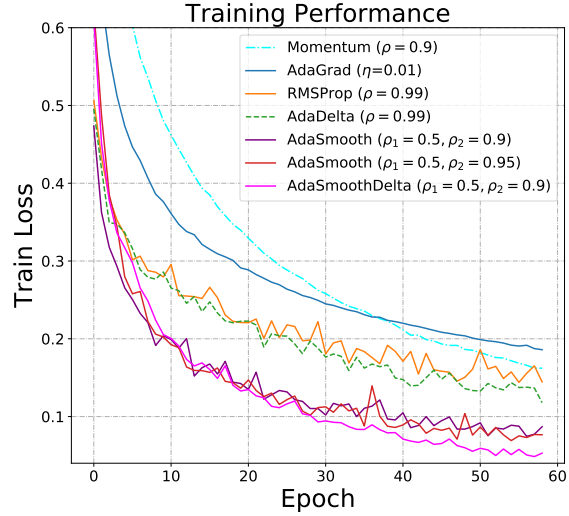
Method	MNIST	Census
Momentum ($\rho = 0.9$)	98.64%	85.65%
AdaGrad ($\eta=0.01$)	98.55%	86.02%
RMSProp ($\rho = 0.99$)	99.15%	85.90%
AdaDelta ($\rho = 0.99$)	99.15%	86.89%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.9$)	99.34%	86.94%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.95$)	99.45%	87.10%
AdaSmDel. ($\rho_1 = 0.5, \rho_2 = 0.9$)	99.60%	86.86%

Table 1. **MLP:** Best in-sample evaluation in training accuracy (%). AdaSmDel is short for AdaSmoothDelta.

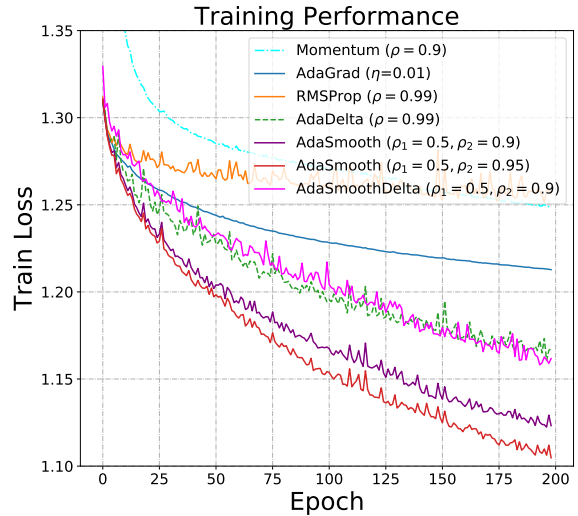
Method	MNIST	Census
Momentum ($\rho = 0.9$)	94.38%	83.13%
AdaGrad ($\eta=0.01$)	96.21%	84.40%
RMSProp ($\rho = 0.99$)	97.14%	84.43%
AdaDelta ($\rho = 0.99$)	97.06%	84.41%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.9$)	97.26%	84.46%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.95$)	97.34%	84.48%
AdaSmDel. ($\rho_1 = 0.5, \rho_2 = 0.9$)	97.24%	84.51%

Table 2. **MLP:** Best out-of-sample evaluation in test accuracy for the first 5 epochs. AdaSmDel is short for AdaSmoothDelta.

tasks finding internal linear and nonlinear features behind



(a) MNIST training Loss



(b) Census Income training loss

Figure 3. **MLP:** Comparison of descent methods on MNIST digit and Census Income data sets for 60 and 200 epochs with MLP.

the model inputs and outputs. We adopt the simplest MLP structure: an input layer, a hidden layer, and an output layer. We notice that rectified linear unit (Relu) outperforms Tanh, Sigmoid, and other nonlinear units in practice making it the default nonlinear function in our structures. Since dropout has become a core tool in training neural networks (Srivastava et al., 2014), we adopt 50% dropout noise to the network architecture during training to prevent overfitting. To be more concrete, the detailed architecture for each fully connected layer is described by

$F(\langle num\ outputs \rangle : \langle activation\ function \rangle)$; and for a dropout layer is described by $DP(\langle rate \rangle)$. Then the network structure we use can be described as follows:

$$F(128:Relu)DP(0.5)F(num\ of\ classes:Softmax). \quad (23)$$

All methods are trained on mini-batches of 64 images per batch for 60 or 200 epochs through the training set. Setting the hyper-parameter to $\epsilon = 1e - 6$. If not especially mentioned, the global learning rates are set to $\eta = 0.001$ in all scenarios. While a relatively large learning rate ($\eta = 0.01$) is used for AdaGrad method since its accumulated decaying effect; learning rate for the AdaDelta method is set to 1 as suggested by Zeiler (2012) and for the AdaSmoothDelta method is set to 0.5 as discussed in Section 3.3 while we will show in the sequel the learning rate for AdaSmoothDelta will not influence the result significantly. In Figure 3(a) and 3(b) we compare SGD with Momentum, AdaGrad, RMSProp, AdaDelta, AdaSmooth and AdaSmoothDelta in optimizing the training set losses for MNIST and Census Income data sets respectively. The SGD with Momentum method does the worst in this case. AdaSmooth performs slightly better than AdaGrad and RMSProp in the MNIST case and much better than the latter in the Census Income case. AdaSmooth shows fast convergence from the initial epochs while continuing to reduce the training losses in both the two experiments. We here show two sets of slow decay constant for AdaSmooth, i.e., “ $\rho_2 = 0.9$ ” and “ $\rho_2 = 0.95$ ”. Since we square the scaled smoothing constant in Eq (18), when $\rho_1 = \rho_2 = 0.9$, the AdaSmooth recovers to RMSProp with $\rho = 0.99$ (so as the AdaSmoothDelta and AdaDelta case). In all cases, the AdaSmooth results perform better while there is almost no difference between the results of AdaSmooth with various hyper-parameters in the MLP model. Table 1 shows the best training set accuracy for different algorithms, indicating the superiority of AdaSmooth. While we notice the best test set accuracy for various algorithms are very close; we only report the best ones for the first 5 epochs in Table 2. In all scenarios, the AdaSmooth method converges slightly faster than other optimization methods in terms of the test accuracy.

4.2. Experiment: Convolutional Neural Networks

Convolutional neural networks (CNN) are powerful models with non-convex objective functions. CNN with several layers of convolution, pooling and nonlinear units have mostly demonstrated remarkable success in computer vision tasks, e.g., face identification, traffic sign detection, and medical picture segmentation (Krizhevsky et al., 2012) and speech recognition (Hinton et al., 2012a; Graves et al., 2013); which is partly from the local connectivity of the convolutional layers, and the rotational and shift invariance due to the pooling layers. To evaluate the strategy and demonstrate the main advantages of the proposed AdaSmooth method, a real

handwritten digit classification task, MNIST is used. For comparison with Zeiler (2012)’s method, we train with Relu nonlinearities and 2 convolutional layers in the front of the structure, followed by two fully connected layers. Again, the dropout with 50% noise is adopted in the network to prevent from overfitting.

To be more concrete, the detailed architecture for each convolutional layer is described by $C(\langle kernel\ size \rangle : \langle num\ outputs \rangle : \langle activation\ function \rangle)$; for each fully connected layer is described by $F(\langle num\ outputs \rangle : \langle activation\ function \rangle)$; for a max pooling layer is described by $MP(\langle kernel\ size \rangle : \langle stride\ number \rangle)$; and for a dropout layer is described by $DP(\langle rate \rangle)$. Then the network structure we use can be described as follows:

$$C(5:10:Relu)MP(2:2)C(5:20:Relu)MP(2:2) \\ - DP(0.5)F(50:Relu)DP(0.5)F(10:Softmax). \quad (24)$$

All methods are trained on mini-batches of 64 images per batch for 50 epochs through the training set. Setting the hyper-parameters to $\epsilon = 1e - 6$ and $\rho_1 = 0.5, \rho_2 = 0.9$ or 0.99 , i.e., the number of periods chosen find the exponential moving average are between 3 and 19, or between 3 and 199 iterations as discussed in Section 2.3.1 and 3. The AdaSmooth with $\rho_2 = 0.9$ or 0.99 is a wide range of upper bound on the decay constant from this context; and we shall see the AdaSmooth results are not sensitive to these choices. If not specially described, a small learning rate $\eta = 0.001$ for the convolutional networks is used in our experiments when applying stochastic descent. While again, the learning rate for the AdaDelta method is set to 1 as suggested by Zeiler (2012) and for the AdaSmoothDelta method is set to 0.5 as discussed in Section 3.3.

In Figure 4(a), we compare SGD with Momentum, AdaGrad, RMSProp, AdaDelta, AdaSmooth, and AdaSmoothDelta in optimizing the training set loss (negative log likelihood). The RMSProp ($\rho = 0.9$) does the worst in this case, whereas tuning the decay constant to $\rho = 0.99$ can significantly improve performance making the RMSProp sensitive to the hyper-parameter (hence the learning rates per-dimension). Further, though the training loss of AdaDelta ($\rho = 0.9$) decreases fastest in the first 3 epochs, the performance becomes poor at the end of the training, with an average loss larger than 0.15; while the overall performance of AdaDelta ($\rho = 0.99$) works better than the former, however, its overall accuracy is still worse than AdaSmooth. This also reveals the same drawback of AdaGrad for the AdaDelta method, i.e., they are sensitive to the choices of hyper-parameters.

In order to evaluate whether the AdaSmooth actually can find the compensation we want, we also explore a random selection for the number of periods/iterations N , termed as *AdaSmooth(Random)* in the sequel, which selects the

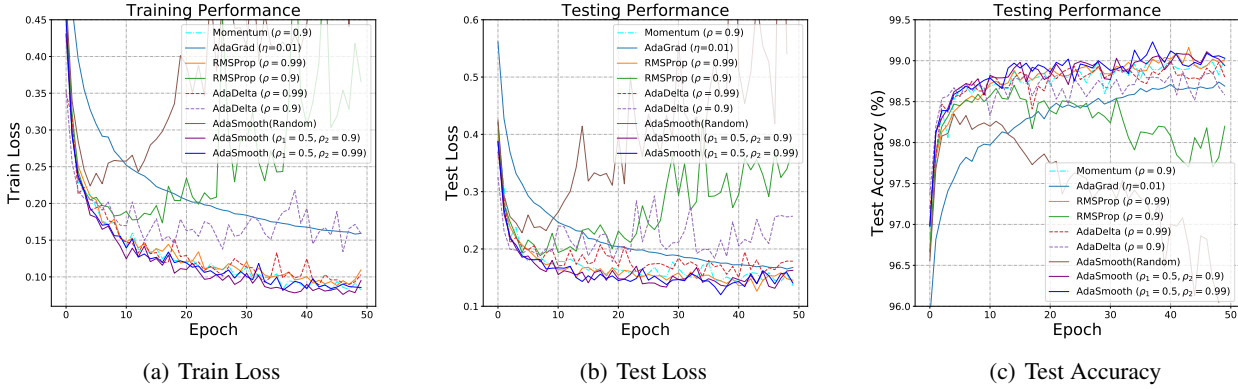


Figure 4. CNN: Comparison of descent methods on MNIST digit data set for 50 epochs with CNN. Though not significantly better than RMSProp or SGD with Momentum, AdaSmooth converges slightly faster and obtains better out-of-sample accuracy.

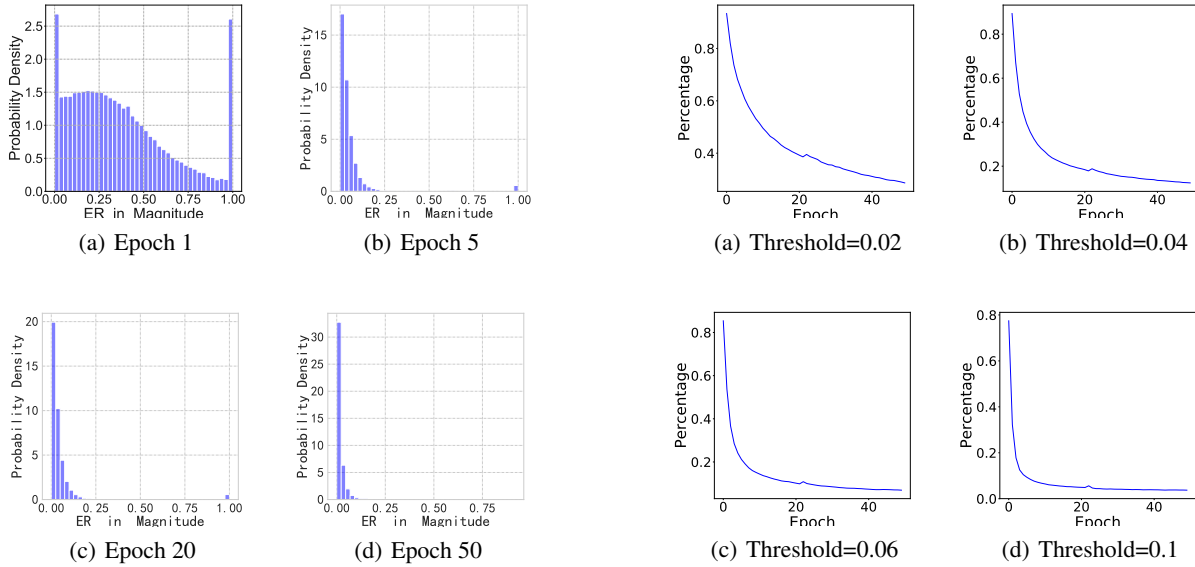


Figure 5. CNN: The distribution of ERs for all parameters in different epochs when training MNIST data for AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.99$) model.

Figure 6. CNN: Percentage of ERs when they are larger than a threshold for different epochs in AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.99$) model.

decay constant between ρ_1 and ρ_2 randomly during different batches. The result shows the AdaSmooth(Random) cannot find the adaptive learning rates per-dimension in the right way showing our method actually finds the correct choices between the lower and upper bound of the decay constant.

The analysis of different methods in test loss and test accuracy are consistent with that of training loss as shown in Figure 4(b) and Figure 4(c). We further save the effective ratios in magnitude as stated in Eq (15) during each epoch; the distributions of them are shown in Figure 5. In the first few epochs, the ERs in magnitude have more weights in

large values, while in the last few epochs, the probability decreases into a stationary distribution; in other words, the ER distributions in the last few epochs are similar. In Figure 6, we put a threshold on the ER values, the percentages of ERs larger than the threshold are plotted; when the training is in progress, the ERs tend to approach smaller values indicating the algorithm goes into stationary points, the (local) minima. While Figure 7 shows the absolute values of update step Δx_t for different epochs; when the training is in the later stage, the parameters are closer to the (local) minima, indicating a smaller update step and favoring a relatively smaller period (resulting from both the step and the direc-

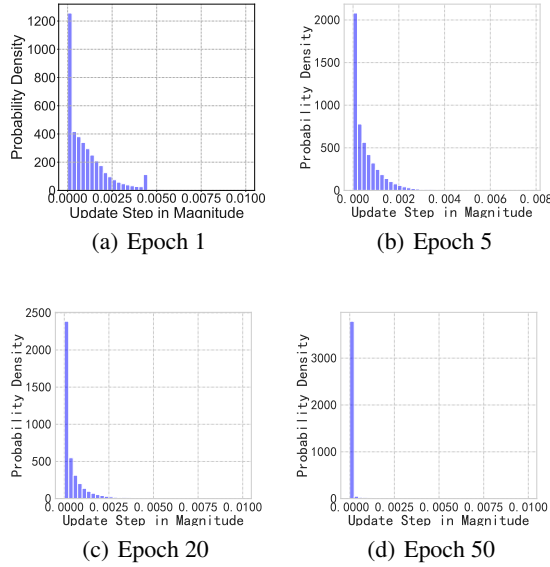
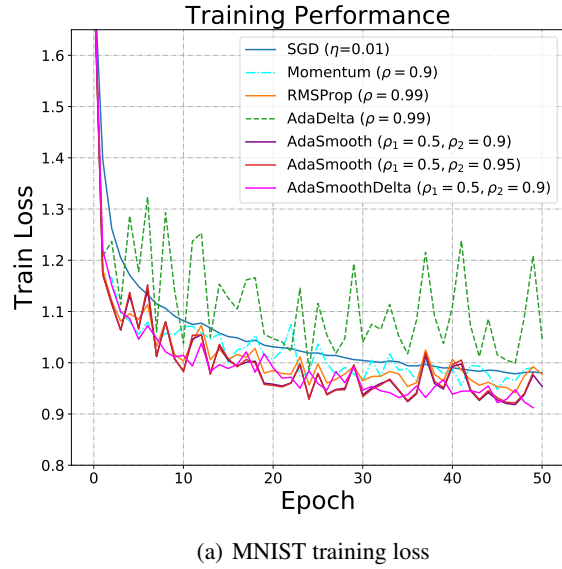


Figure 7. CNN: The distribution of step Δx_t in magnitude for different epochs in AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.99$) model.

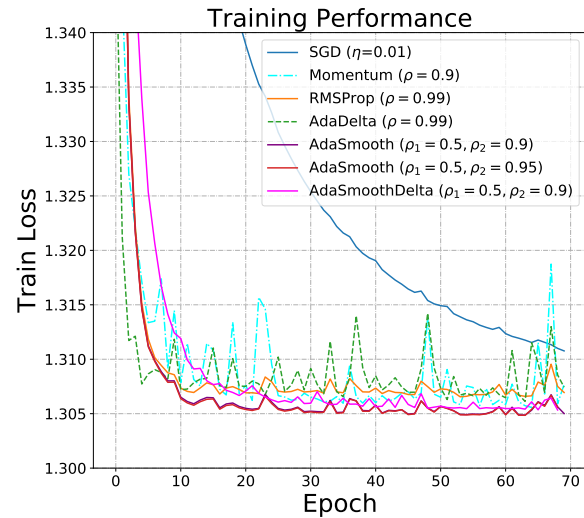
tion of the movement) of exponential moving average in Eq (18), which in turn compensates relatively less for the learning rate per-dimension and makes the update faster. In Table 3, we further show test accuracy for AdaSmooth and AdaSmoothDelta with various hyper-parameters after 10 epochs in which case different choices of the parameters do not significantly alter performance.

Method	MNIST
AdaGrad ($\eta=0.01$)	96.82%
AdaGrad ($\eta=0.001$)	89.11%
RMSProp ($\rho = 0.99$)	97.82%
RMSProp ($\rho = 0.9$)	97.88%
AdaDelta ($\rho = 0.99$)	97.83%
AdaDelta ($\rho = 0.9$)	98.20%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.9$)	98.13%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.95$)	98.12%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.99$)	98.12%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.9, \eta = 0.5$)	98.86%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.95, \eta = 0.5$)	98.91%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.99, \eta = 0.5$)	98.78%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.99, \eta = 0.6$)	98.66%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.99, \eta = 0.7$)	98.66%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.99, \eta = 0.8$)	98.58%

Table 3. CNN: Best out-of-sample evaluation in test accuracy for AdaSmooth and AdaSmoothDelta with various hyper-parameters after 10 epochs. AdaSmDel is short for AdaSmoothDelta.



(a) MNIST training loss



(b) Census Income training loss

Figure 8. **Logistic regression:** Comparison of descent methods on MNIST digit and Census Income data sets for 50 and 70 epochs with logistic regression. Notice that the training curve of AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.9$) is close to ($\rho_1 = 0.5, \rho_2 = 0.95$) in this case for both data sets.

4.3. Experiment: Logistic Regression

Logistic regression is an analytic technique for multivariate modeling of categorical dependent variables and has a well-studied convex objective, making it suitable for comparison of different optimizers without worrying about going into local minima (Menard, 2002). To evaluate, again if not mentioned explicitly, the global learning rates are set to

$\eta = 0.001$ in all scenarios. While a relatively large learning rate ($\eta = 0.01$) is used for AdaGrad method since its accumulated decaying effect; learning rate for the AdaDelta method is set to 1 as suggested by Zeiler (2012) and for the AdaSmoothDelta method is set to 0.5 as discussed in Section 3.3. The loss curves for training processes are shown in Figure 8(a) and 8(b), where we compare SGD, SGD with Momentum, RMSProp, AdaDelta, AdaSmooth, and AdaSmoothDelta in optimizing the training set losses for MNIST and Census Income data sets respectively. The unaltered SGD method does the worst in this case. AdaSmooth performs slightly better than RMSProp in the MNIST case and much better than the latter in the Census Income case. AdaSmooth matches the fast convergence of AdaDelta (in which case AdaSmooth converges slightly slower than AdaDelta in the first few epochs), while AdaSmooth continues to reduce the training loss, converging to the best performance in these models. As aforementioned, when $\rho_1 = \rho_2 = 0.9$, the AdaSmooth recovers to RMSProp with $\rho = 0.99$ (so as the AdaSmoothDelta and AdaDelta case). In all cases, the AdaSmooth results perform better while the difference between various hyper-parameters for AdaSmooth is not significant; there are almost no differences in AdaSmooth results with different hyper-parameter settings, indicating its insensitivity to hyper-parameters. In Table 4, we compare the training set accuracy of various algorithms; AdaSmooth works best in the MNIST case, while the superiority in Census Income data is not significant. Similar results as the MLP scenario can be observed in the test accuracy and we shall not give the details for simplicity.

Method	MNIST	Census
SGD ($\eta=0.01$)	93.29%	84.84%
Momentum ($\rho = 0.9$)	93.39%	84.94%
RMSProp ($\rho = 0.99$)	93.70%	84.94%
AdaDelta ($\rho = 0.99$)	93.48%	84.94%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.9$)	93.74%	84.92%
AdaSmooth ($\rho_1 = 0.5, \rho_2 = 0.95$)	93.71%	84.94%
AdaSmoDel. ($\rho_1 = 0.5, \rho_2 = 0.9$)	93.66%	84.97%

Table 4. **Logistic regression:** Best in-sample evaluation in training accuracy. AdaSmoDel is short for AdaSmoothDelta.

5. Conclusion

The aim of this paper is to solve the hyper-parameter tuning in the gradient-based optimization methods for machine learning problems with large data sets and high-dimensional parameter spaces. We propose a simple and computationally efficient algorithm that requires little memory and is easy to implement for gradient-based optimization of stochastic objective functions. Overall, we show that AdaSmooth is a versatile algorithm that scales to large-scale high-dimensional machine learning problems and AdaSmoothDelta is an al-

gorithm insensitive to both the global learning rate and the hyper-parameters with caution to set the global learning rate smaller than 1.

References

- Becker, Sue and Le Cun, Yann. Improving the convergence of back-propagation learning with. 1988.
- Dauphin, Yann N, Pascanu, Razvan, Gulcehre, Caglar, Cho, Kyunghyun, Ganguli, Surya, and Bengio, Yoshua. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in neural information processing systems*, 27, 2014.
- Dozat, Timothy. Incorporating nesterov momentum into adam. 2016.
- Duchi, John, Hazan, Elad, and Singer, Yoram. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Graves, Alex, Mohamed, Abdel-rahman, and Hinton, Geoffrey. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6645–6649. Ieee, 2013.
- Hinton, Geoffrey, Deng, Li, Yu, Dong, Dahl, George E, Mohamed, Abdel-rahman, Jaitly, Navdeep, Senior, Andrew, Vanhoucke, Vincent, Nguyen, Patrick, Sainath, Tara N, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97, 2012a.
- Hinton, Geoffrey, Srivastava, Nitish, and Swersky, Kevin. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012b.
- Kaufman, Perry J. Smarter trading, 1995.
- Kaufman, Perry J. *Trading Systems and Methods*, + Website, volume 591. John Wiley & Sons, 2013.
- Kingma, Diederik P and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- Le Roux, Nicolas and Fitzgibbon, Andrew W. A fast natural newton method. In *ICML*, 2010.
- LeCun, Yann. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Lu, Jun. Exploring classic quantitative strategies. *arXiv preprint arXiv:2202.11309*, 2022a.
- Lu, Jun. Matrix decomposition and applications. *arXiv preprint arXiv:2201.00145*, 2022b.
- Lu, Jun and Yi, Shao. Reducing overestimating and underestimating volatility via the augmented blending-ARCH model. *arXiv preprint arXiv:2203.12456*, 2022.
- Menard, Scott. *Applied logistic regression analysis*, volume 106. Sage, 2002.

- Moulines, Eric and Bach, Francis. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24, 2011.
- Polyak, Boris T and Juditsky, Anatoli B. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- Qian, Ning. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- Robbins, Herbert and Monro, Sutton. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Ruder, Sebastian. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- Rumelhart, David E, Hinton, Geoffrey E, and Williams, Ronald J. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- Ruppert, David. Efficient estimations from a slowly convergent robbins-monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- Schaul, Tom, Zhang, Sixin, and LeCun, Yann. No more pesky learning rates. In *International conference on machine learning*, pp. 343–351. PMLR, 2013.
- Smith, Leslie N. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pp. 464–472. IEEE, 2017.
- Srivastava, Nitish, Hinton, Geoffrey, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Sutskever, Ilya, Martens, James, Dahl, George, and Hinton, Geoffrey. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pp. 1139–1147. PMLR, 2013.
- You, Yang, Li, Jing, Reddi, Sashank, Hseu, Jonathan, Kumar, Sanjiv, Bhojanapalli, Srinadh, Song, Xiaodan, Demmel, James, Keutzer, Kurt, and Hsieh, Cho-Jui. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Zeiler, Matthew D. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.