

PID Controller-Based Stochastic Optimization Acceleration for Deep Neural Networks

Haoqian Wang^{ib}, Member, IEEE, Yi Luo, Wangpeng An^{ib}, Qingyun Sun,
Jun Xu^{ib}, and Lei Zhang^{ib}, Fellow, IEEE

Abstract—Deep neural networks (DNNs) are widely used and demonstrated their power in many applications, such as computer vision and pattern recognition. However, the training of these networks can be time consuming. Such a problem could be alleviated by using efficient optimizers. As one of the most commonly used optimizers, stochastic gradient descent-momentum (SGD-M) uses past and present gradients for parameter updates. However, in the process of network training, SGD-M may encounter some drawbacks, such as the overshoot phenomenon. This problem would slow the training convergence. To alleviate this problem and accelerate the convergence of DNN optimization, we propose a proportional-integral-derivative (PID) approach. Specifically, we investigate the intrinsic relationships between the PID-based controller and SGD-M first. We further propose a PID-based optimization algorithm to update the network parameters, where the past, current, and change of gradients are exploited. Consequently, our proposed PID-based optimization alleviates the overshoot problem suffered by SGD-M. When tested on popular DNN architectures, it also obtains up to 50% acceleration with competitive accuracy. Extensive experiments about computer vision and natural language processing demonstrate the effectiveness of our method on benchmark data sets, including CIFAR10, CIFAR100, Tiny-ImageNet, and PTB. We have released the code at <https://github.com/tensorboy/PIDOptimizer>.

Index Terms—Deep neural network (DNN), optimization, proportional-integral-derivative (PID) control, stochastic gradient descent (SGD)-momentum.

I. INTRODUCTION

BENEFITTING from the availability of a great number of data (e.g., ImageNet [1]) and the fast-growing power

of GPUs, deep neural networks (DNNs) success in a wide range of applications, such as computer vision and natural language processing. Despite the significant successes of DNNs, the training and inference of deep and wide DNNs are often computationally expensive, which may take several days or longer even with powerful GPUs. Many stochastic optimization algorithms are not only used in the field of machine learning [2] but also deep learning [3]. It is very important to explore how to boost the speed of training DNNs while maintaining performance. Furthermore, with a better optimization method, even computation limited hardware (e.g., IoT device) can save lots of time and memory usage. The accelerating methods of the computational time for DNNs can be divided into two parts, the speed-up of training and that of test. The methods in [4]–[6] aiming to speed up test process of DNNs often focus on not only the decomposition of layers but also the optimization solutions to the decomposition. Besides, there has been other streams on improving testing performance of DNNs, such as the FFT-based algorithms [7] and reduced parameters in deep nets [8]. As for the methods to speed up the training speed of DNNs, the key factor is the way to update the millions of parameters of a DNN. This process mainly depends on the optimizer, and the choice of the optimizer is also a key point of a model. Even with the same data set and architecture, different optimizers could result in very different training effects, due to different directions of the gradient descent, different optimizers may reach completely different local minimum [9].

The learning rate is another principal hyperparameter for DNN training [10]. Based on different strategies of choosing learning rates, DNN optimizers can be categorized into two groups: 1) Hand-tuned learning rate optimizers, such as stochastic gradient descent (SGD) [11], SGD-momentum (SGD-M) [12], Nesterov's momentum [12], and so on, and 2. Auto learning rate optimizers, such as AdaGrad [13], RMSProp [14], Adam [15], and so on.

The SGD-M method puts past and current gradients into consideration and then updates the network parameters. Although SGD-M performs well in most cases, it may encounter an overshoot phenomenon [16], which indicates the case where the weight exceeds its target value too much and fails to correct its update direction. Such an overshoot problem costs more resources (e.g., time and GPUs) to train a DNN and also hampers the convergence of SGD-M. Therefore, a more efficient DNN optimizer is eagerly desired to alleviate the overshoot problem and achieve better convergence.

Manuscript received April 1, 2019; revised September 20, 2019 and December 5, 2019; accepted December 25, 2019. Date of publication January 28, 2020; date of current version December 1, 2020. This work was supported in part by the National Natural Science Foundation of China (NSFC) fund under Grant 61571259, Grant 61831014, and Grant 61531014, and in part by the Shenzhen Science and Technology Project under Grant JCYJ20170817161916238, Grant JCYJ20180508152042002, and Grant GGF2017040714161462. (Corresponding author: Haoqian Wang.)

Haoqian Wang, Yi Luo, and Wangpeng An are with the Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen 518055, China, and also with the Shenzhen Institute of Future Media Technology, Shenzhen 518055, China (e-mail: wanghaoqian@tsinghua.edu.cn; vast2stars@gmail.com; anwangpeng@gmail.com).

Qingyun Sun is with the Department of Mathematics, Stanford University, Stanford, CA 94305 USA (e-mail: qysun@stanford.edu).

Jun Xu is with the College of Computer Science, Nankai University, Tianjin 300071, China (e-mail: nankaimathxujun@gmail.com).

Lei Zhang is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong, and also with the Artificial Intelligence Center, Alibaba DAMO Academy, Hangzhou 311121, China (e-mail: cslzhang@comp.polyu.edu.hk).

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2019.2963066

2162-237X © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

The similarity between optimization algorithms popularly employed in DNN training and classic control methods has been investigated in [17]. In automatic control systems, feedback control is essential. Proportional-integral-derivative (PID) controller is the most widely used feedback control mechanism due to its simplicity and functionality [18]. Most of industrial control system are based on PID [19], such as unmanned aerial vehicles [20], robotics [21], and autonomous vehicles [22]. PID control takes current error, change in error (differentiation of the error over time), and the past cumulative error (integral of the error over time) into account. Therefore, the difference between current and expected outputs will be minimized.

On the other hand, a few studies have been done on the connections between PID with DNN optimization. In this article, we investigate specific relationships analytically and mathematically toward this research line. We first clarify the intrinsic connection between the PID controller and stochastic optimization methods, including SGD, SGD-M, and Nesterov's momentum. Finally, we propose a PID-based optimization method for DNN training. Similar to SGD-M, our proposed PID optimizer also considers the past and current gradients for network update. The Laplace transform [23] is further introduced for hyperparameter initialization, which makes our method simple yet effective. Our major contributions to this article can be summarized in three folds.

- 1) By combining the error calculation in the feedback control system with network parameters' update, we reveal a potential relationship between DNN optimization and feedback system control. We also find that some optimizers (e.g., SGD-M) are special cases of PID control device.
- 2) We propose a PID-based DNN optimization approach by taking the past, current, and changing information of the gradient into consideration. The hyperparameter in our PID optimizer is initialized by classical Laplace transform.
- 3) We systematically experiment with our proposed PID optimizer on CIFAR10, CIFAR100, Tiny-ImageNet, and PTB data sets. The results show that the PID optimizer is faster than SGD-M in the DNN training process.

A preliminary version of this article was presented at a conference version [24]. In this article, we incorporate additional content in significant ways.

- 1) We evaluate the performance of our PID optimizer on the language modeling application by utilizing the character-level Penn Treebank (PTB-c) data set with an LSTM network.
- 2) The proposed PID optimizer is applied to GAN with the MNIST data set and shows the digital images generated by them separately to illustrate that our method is also applicable in GAN.
- 3) We update the conclusion that the proposed PID optimizer exceeds SGD-M in GANs and RNNs.

We organize the rest of this article as follows. Section II briefly surveys related works. Section III investigates the relationship between PID controller and DNN optimization algorithms. Section IV introduces the proposed PID approach for DNN optimization. Experimental results and detailed

analysis are reported in Section V. Section VI concludes this article.

II. RELATED WORKS

A. Classic Deep Neural Network Architectures

1) *CNN*: Convolutional neural networks (CNNs) [25] have recently achieved great successes in visual recognition tasks, including image classification [26], object detection [27]–[29], and scene parsing [30]. Recently, lots of deep CNN architectures, such as VGG, ResNet, and DenseNet, have been proposed to improve the performance of these tasks mentioned above. Network depth tends to improve network performance. However, the computational cost of these deep networks also increases significantly. Moreover, real-world systems may be affected by the high cost of these networks.

2) *GAN*: Goodfellow *et al.* [31] first proposed generative adversarial network (GAN), which consists of generative and adversarial networks. The generator tries to obtain very realistic outputs to foolish the discriminator, which would be optimized to distinguish between the real data and the generated outputs. GANs will be trained to generate synthetic data, mimicking genuine data distribution.

In machine learning, models can be classified into two categories: the generative model and the discriminative model. A discriminative network (denoted as D) can discriminate between two (or more) different classes of data, such as CNN trained for image classification. A generative network (denoted as G) can generate new data, which fits the distribution of the training data. For example, a trained Gaussian mixture model is able to generate new random data, which more-or-less fits the distribution of the training data.

GANs pose a challenging optimization problem due to the multiple loss functions, which must be optimized simultaneously. The optimization of GAN is conducted by two steps: 1) optimize the discriminative network while fixing the generative one and 2) optimize the generative network while fixing the discriminative network. Here, fixing a network means only allowing the network to pass forward and not perform back-propagation. These two steps are seamlessly alternating, updated, and dependent on each other for efficient optimization. After enough training cycles, the optimization objective $V(D, G)$ introduced in [31] will reach the situation, where the probability distribution of the generator exactly matches the true probability distribution of the training data. Meanwhile, the discriminator has the capability to distinguish the realistic data from the virtual generated ones. However, the perfect cooperation between the generator and the discriminator will fail occasionally. The whole system will reach the status of “model collapse,” indicating that the discriminator and the generator tend to produce the same outputs.

3) *LSTM*: Hochreiter *et al.* first proposed the long short term network, generally called LSTM, to obtain long-term dependence information from the network. As a type of recurrent neural network (RNN), LSTM has been widely used and obtained excellent success in many applications. LSTM is deliberately designed to avoid long-term dependence problems. Remember that long-term information is the default

behavior of LSTM in practice, rather than the ability to acquire at great cost. All RNNs have a chained form of repeating network modules. In the standard RNN, this repeating module often has a simple structure (e.g., “tanh” layer). The outputs of all LSTM cells are utilized to construct a new feature, where multinomial logistic regression is introduced to form the LSTM model.

One widely used way to evaluate RNN models is the adding task [32], [33], which takes two sequences of length T as input. By sampling in the range $(0, 1)$ uniformly, we form the first sequence. For another sequence, we set two entries as 1 and the rest as 0. The output is obtained by adding two entries in the first sequence. The positions of the entries are determined by the two entries of 1 from the second sequence.

B. Accelerating the Training/Test Process of DNNs

1) *Training Process Acceleration*: Since DNNs are mostly computationally intensive, Han *et al.* [34] proposed a deep compression method to reduce the storage requirement of DNNs by $35\times$ – $49\times$ without affecting the accuracy. Moreover, the compressed model has $3\times$ – $4\times$ layerwise speedup and $3\times$ – $7\times$ better energy efficiency. Unimportant connections are pruned. Weight sharing and Huffman coding are applied to quantize the network. This article mainly attempts to reduce the number of parameters of neural networks. Liu *et al.* [35] proposed the network slimming technique that can simultaneously reduce the model size, running-time memory, and computing operations. He *et al.* [36] proposed a new filter pruning strategy based on the geometric median to accelerate the training of deep CNNs. Dai *et al.* [37] proposed a synthesis tool to synthesize compact yet accurate DNNs. Du *et al.* [38] proposed a continuous growth and pruning (CGaP) scheme to minimize the redundancy from the beginning. Hubara *et al.* [39] introduced a method to train Quantized Neural Networks that reduce memory size and accesses during the forward pass. Kidambi *et al.* [40] presented an intuitive and easier-to-tune version of ASGD (please refer to Section IV) and showed that ASGD leads to faster convergence significantly with a comparable accuracy than SGD, heavy ball, and Nesterov’s momentum [12].

2) *Test Process Acceleration*: Denton *et al.* [4] proposed a method that compresses all convolutional layers. This is achieved by approximating proper low rank and then updating the upper layers until the prediction result is enhanced. Based on singular value decompositions (SVDs), this process consists of numerous tensor decomposition operations and filter clustering approaches to make use of similarities among learned features. Jaderberg *et al.* [5] introduced an easy-to-implement method that can significantly speed up pretrained CNNs with minimal modifications to existing frameworks. There can be a small associated loss in performance, but this is tunable to the desired accuracy level. Zhang *et al.* [6] first proposed a response reconstruction method, which introduces the nonlinear neurons and a low-rank constraint. Without the usage of SGD and based on generalized SVD (GSVD), a solution is developed for this nonlinear problem. Li *et al.* presented a method to prune filters with relatively low-weight

magnitudes to produce CNNs with reduced computation costs without introducing irregular sparsity [41].

C. Deep Learning Optimization

In the training of DNN [10], the learning rate is an essential hyperparameter. DNN optimizers can be categorized into two groups based on different strategies of setting the learning rate: 1) hand-tuned learning rate optimizers, such as SGD [11], SGD momentum [12], Nesterov’s momentum [12], and so on, and 2) autolearning rate optimizers, such as AdaGrad [13], RMSProp [14], Adam [15], and so on. Good results have been achieved on the CIFAR10, CIFAR100, ImageNet, PASCAL VOC, and MS COCO data sets. They were mostly obtained by residual neural networks [42]–[45] trained by using SGD-M. This article focuses on the improvement of the first category of optimizers. The introduction to these optimizers is as follows.

Classical momentum [25] is the first-ever variant of gradient descent involving the usage of a momentum parameter. In the objective across iterations, it accelerates gradient descent that collects a velocity vector in directions of continuous reduction.

SGD [11] is a widely used optimizer for DNN training. SGD is easy to apply, but the disadvantage of SGD is that it converges slowly and may oscillate at the saddle point. Moreover, how to choose the learning rate reasonably is a major difficulty of SGD.

SGD-M [12] is an optimization method that considers momentum. Compared with the original gradient descent step, the SGD-M introduces variables related to the previous step. It means that the parameter update direction is decided not only by the present gradient but also by the previously accumulated direction of the fall. This allows the parameters to change a little in the direction where gradient change frequently. On the contrary to this, SGD-M changes parameters a lot in the direction, where gradient change slowly.

Nesterov’s momentum [12] is another momentum optimization algorithm motivated by Nesterov’s accelerated gradient method [46]. Momentum is improved from the SGD algorithm so that each parameter update direction depends not only on the gradient of the current position but also on the direction of the last parameter update. In other words, Nesterov’s momentum essentially uses the second-order information of the objective (loss function), and therefore, it can accelerate the convergence better.

D. PID Controller

Traditionally, the PID controller has been used to control a feedback system [19] by exploiting the present, past, and future information of prediction error. The theoretical basis of the PID controller was first proposed by Maxwell [47] in his seminal article “On Governors.” The mathematical formulation was given by Minorsky [48]. In recent years, several advanced control algorithms have been proposed.

We define the difference between the actual output and the desired output as error $e(t)$. The PID controller calculates the error $e(t)$ in every step t , and then applies a correction $u(t)$ to the system as a function of the proportional (P), integral (I),

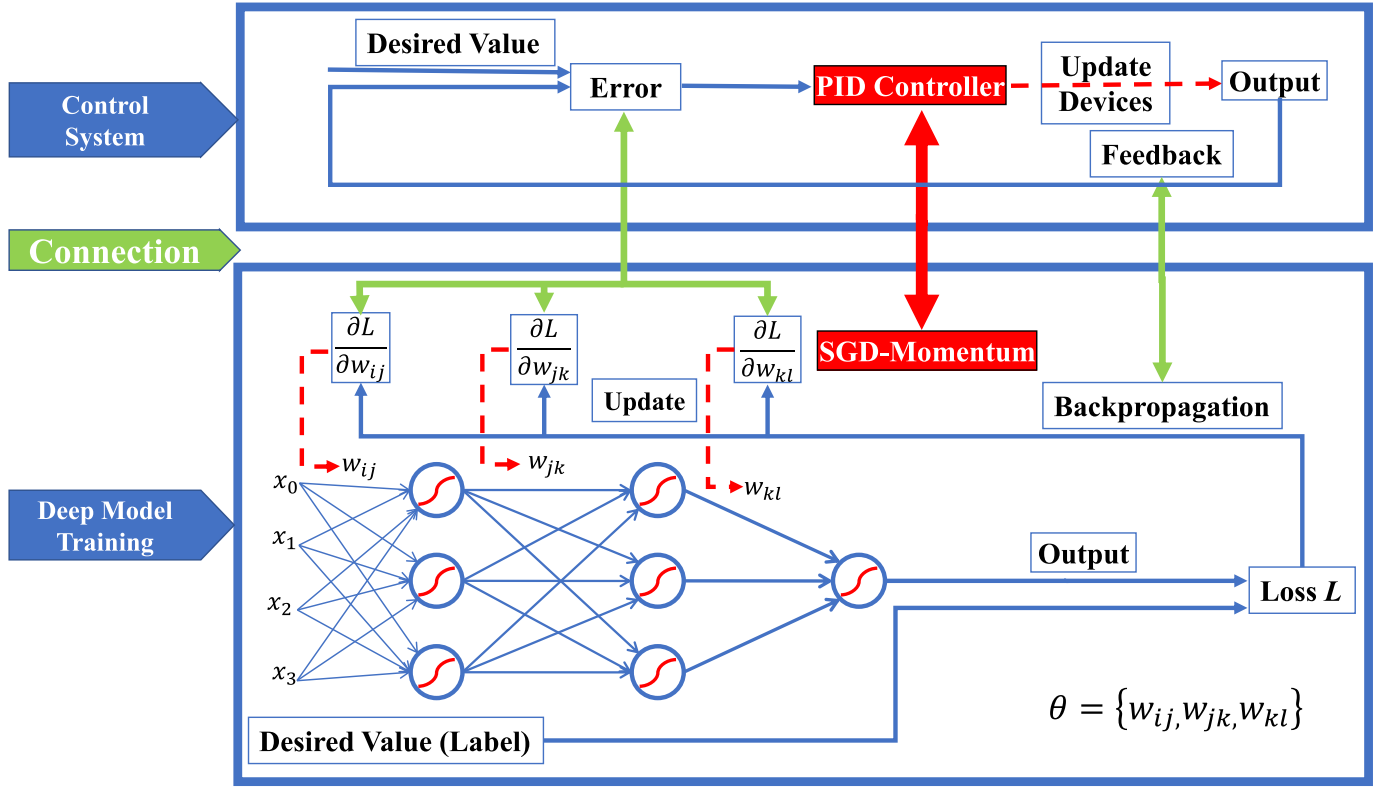


Fig. 1. Illustrations about the relationships between control system and deep model training. It also shows the connection between PID controller and SGD-M.

and derivative (*D*) terms of $e(t)$. Mathematically, the PID controller can be described as

$$u(t) = K_p e(t) + K_i \int_0^t e(t) dt + K_d \frac{d}{dt} e(t) \quad (1)$$

where K_p , K_i , and K_d correspond to the gain coefficients of the *P*, *I*, and *D* terms, respectively. The function of error $e(t)$ is the same as the gradient in optimization of deep learning. The coefficients K_p , K_i , and K_d reflect the contribution to the current correction to the current, past, and future errors, respectively.

According to our analyses, we find that PID control techniques can be more useful for the optimization of deep networks. The study presented in this article is one of the first investigations to apply PID as a new optimizer to the deep learning field. Our studies have succeeded in demonstrating the significant advantages of the proposed optimizer. With the inheritance of the advantages of the PID controller, the proposed optimizer performs well despite its simplicity.

III. PID AND DEEP NEURAL NETWORK OPTIMIZATION

We reveal the intrinsic relation between PID control and DNNs optimization. The intrinsic relation inspires us to explore new DNNs optimization methods. The core idea of this section is to regard the parameter update in the DNNs training process as using a PID controller in the system to reach equilibrium.

A. Overall Connections

At first, we summarize the training process for deep learning. DNNs need to map the input x to the output y through

parameters θ . To measure the gap between the DNN output and the desired output, the loss function L is introduced. Given some training data, we can calculate the loss function $L(\theta, X_{\text{train}})$. In order to minimize the loss function L , we find the derivative of the loss function L with respect to the parameter θ and update θ with the gradient descent method in most cases. DNNs gradually learn the complex relationship between input x and output y by constantly updating the parameters θ , which called DNN's training. The updating of θ is driven by the gradient of loss function until it is converged.

Then, the purpose of an automated control system is to evaluate the system status and make it to the desired status through a controller. In a feedback control system, the controller's action is affected by the system's output. The error $e(t)$ between the measured system status and desired status is taken into consideration so that the controller can make the system get close to the desired status.

More specifically, as shown in (1), PID controller estimates a control variable $u(t)$ by considering the current, past, and future (derivative) of the error $e(t)$.

From here, we can see that the error in the PID control system is related to the gradient in the DNN training process. The update of parameters during DNN training can be analogized to the adjustment of the system by the PID controller.

As can be seen from the discussion above, there is a high similarity between DNNs optimization and PID-based control system. Fig. 1 shows their flowchart, respectively, and we can see the similarity more intuitively. Based on the difference between the output and target, both of them change the system/network. The negative feedback process in

the PID controller is similar to the backpropagation in DNNs optimization.

One key difference is that the PID controller computes the update utilizing system error $e(t)$. However, the DNN optimizer decides the updates by considering gradient $\partial L/\partial \theta$. Let us regard the gradient $\partial L/\partial \theta$ as the incarnation of error $e(t)$. Then, the PID controller could be fully related to DNN optimization. In the following, we prove that SGD, SGD-M, and Nesterov's momentum all are special cases of PID controller.

B. Stochastic Gradient Descent

In DNN training, there are widely used optimizers, such as SGD and its variants. The parameter update rule of SGD from iteration t to $t+1$ is determined by

$$\theta_{t+1} = \theta_t - r \partial L_t / \partial \theta_t \quad (2)$$

where r is the learning rate. We now regard the gradient $\partial L_t / \partial \theta_t$ as error $e(t)$ in the PID control system. Comparing with the PID controller in (1), we find that SGD can be viewed as one type of P controller with $K_p = r$.

C. SGD-Momentum

SGD-M is faster than SGD to train a DNN because it can use history gradient. The rule of SGD-M updating parameter is given by

$$\begin{cases} V_{t+1} = \alpha V_t - r \partial L_t / \partial \theta_t \\ \theta_{t+1} = \theta_t + V_{t+1} \end{cases} \quad (3)$$

where V_t is a term that accumulates historical gradients. $\alpha \in (0, 1)$ is the factor that balances the past and current gradients. It is usually set to 0.9 [49]. Dividing two sides of the first formula of (3) by α^{t+1}

$$\frac{V_{t+1}}{\alpha^{t+1}} = \frac{V_t}{\alpha^t} - r \frac{\partial L_t / \partial \theta_t}{\alpha^{t+1}}. \quad (4)$$

By applying (4) from time $t+1$ to 1, we have

$$\begin{cases} \frac{V_{t+1}}{\alpha^{t+1}} - \frac{V_t}{\alpha^t} = -r \frac{\partial L_t / \partial \theta_t}{\alpha^{t+1}} \\ \frac{V_t}{\alpha^t} - \frac{V_{t-1}}{\alpha^{t-1}} = -r \frac{\partial L_{t-1} / \partial \theta_{t-1}}{\alpha^t} \\ \vdots \\ \frac{V_1}{\alpha^1} - \frac{V_0}{\alpha^0} = -r \frac{\partial L_0 / \partial \theta_0}{\alpha^1}. \end{cases} \quad (5)$$

By adding the aforementioned equations together, we get

$$\frac{V_{t+1}}{\alpha^{t+1}} = \frac{V_0}{\alpha^0} - r \sum_{i=0}^t \frac{\partial L_i / \partial \theta_i}{\alpha^{i+1}}. \quad (6)$$

To make it more general, we set the initial condition $V_0 = 0$, and thus, the above-mentioned equation can be simplified as follows:

$$V_{t+1} = -r \sum_{i=0}^t \alpha^{t-i} \partial L_i / \partial \theta_{t-1}. \quad (7)$$

Put V_{t+1} into the second formula of (3), we have

$$\theta_{t+1} - \theta_t = -r \frac{\partial L_t}{\partial \theta_t} - r \sum_{i=0}^{t-1} \alpha^{t-i} \partial L_i / \partial \theta_i. \quad (8)$$

We could learn that parameter update process considers both the current gradient (P control) and the integral of past gradients (I control). If we assume $\alpha = 1$, we get the following equation:

$$\theta_{t+1} - \theta_t = -r \partial L_t / \partial \theta_t - r \sum_{i=0}^{t-1} \partial L_i / \partial \theta_i. \quad (9)$$

Comparing (9) with (1), we can see that SGD-M is a PI controller with $K_p = r$ and $K_i = r$. By using some mathematical skill [50], we simplify (3) by removing V_t . Then, (9) can be rewritten as

$$\theta_{t+1} = \theta_t - r \partial L_t / \partial \theta_t - r \sum_{i=0}^{t-1} \partial L_i / \partial \theta_i \alpha^{t-i}. \quad (10)$$

We can see it clear that the network parameter update depends on both current gradient $r \partial L_t / \partial \theta_t$ and the integral of past gradients $r \sum_{i=0}^{t-1} \partial L_i / \partial \theta_i \alpha^{t-i}$. It should be noted that the term I includes a decay factor α . Due to the huge number of training data, it is better to calculate the gradient based on mini-batch of training data. Therefore, the gradients behave in a stochastic manner. The purpose of the introduction of decay term α is to keep the gradients away from current value so that it can alleviate noise. In all, based on the analyses, we can view SGD-M as a PI controller.

D. Nesterov's Momentum

Momentum is improved from the SGD algorithm and it considers the second-order information of the objective (loss function), and therefore, it can accelerate the convergence better. We set the update rule as

$$\begin{cases} V_{t+1} = \alpha V_t - r \partial L_t / \partial (\theta_t + \alpha V_t) \\ \theta_{t+1} = \theta_t + V_{t+1}. \end{cases} \quad (11)$$

By using a variable transform $\hat{\theta}_t = \theta_t + \alpha V_t$ and formulating the update rule with respect to $\hat{\theta}$, we have

$$\begin{cases} V_{t+1} = \alpha V_t - r \partial L_t / \partial \hat{\theta}_t \\ \hat{\theta}_{t+1} = \hat{\theta}_t + (1 + \alpha) V_{t+1} - \alpha V_t. \end{cases} \quad (12)$$

Similar to the derivation process in (4)–(6) of SGD-M, we have

$$V_{t+1} = -r \left(\sum_{i=1}^t (\alpha^{t-i} \partial L_i / \partial \hat{\theta}_i) \right). \quad (13)$$

With (13), (11) can be rewritten as

$$\hat{\theta}_{t+1} - \hat{\theta}_t = -r(1 + \alpha) \partial L_t / \partial \hat{\theta}_t - \alpha r \left(\sum_{i=1}^{t-1} (\alpha^{t-i} \partial L_i / \partial \hat{\theta}_i) \right). \quad (14)$$

We could conclude that the network parameter update considers the current gradient (P control) and the integral of past gradients (I control). If we assume $\alpha = 1$, then

$$\hat{\theta}_{t+1} - \hat{\theta}_t = -2r(\partial L_t / \partial \hat{\theta}_t) - r \left(\sum_{i=0}^{t-1} (\partial L_i / \partial \hat{\theta}_i) \right). \quad (15)$$

Comparing (15) with (1), we can prove that Nesterov's momentum is a PI controller with $K_p = 2r$ and $K_i = r$. Furthermore, compared with SGD-M, Nesterov's momentum would utilize the current gradient and integral of past gradients to update the network parameters, but achieves larger gain coefficient K_p .

IV. PID-BASED DNN OPTIMIZATION

A. Overshoot Problem of SGD-Momentum

We can learn it from (10) and (14) that the Momentum optimizer will accumulate history gradients to accelerate. But on the other hand, the updating of parameters may be in the wrong path if the history gradients lag the update of parameters. According to the definition "the maximum peak value of the response curve measured from the desired response of the system" in discrete-time control systems [16], this phenomenon is named as overshoot. Specifically, it can be written as

$$\text{Overshoot} = \frac{\theta_{\max} - \theta^*}{\theta^*} \quad (16)$$

where θ_{\max} and θ^* are the maximum and optimum values of the weight, respectively.

The overshoot problem's test benchmark is the first function of De Jong [51] due to its smooth, unimodal, and symmetric characteristics. The function can be written as

$$f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2 \quad (17)$$

whose search domain is $-10 \leq x_i \leq 10, i = 1, 2$. For this function $\mathbf{x}^* = (0, 0)$, $f(\mathbf{x}^*) = 0$, we can pursue a global minimum rather than a local one.

To build a simple PID optimizer, we introduce a derivative term of gradient based on SGD-M

$$\text{PID} = \text{Momentum} + K_d(\partial f(x)/\partial x_c - \partial f(x)/\partial x_{c-1}) \quad (18)$$

where c is the present iteration index for x . With different choices of K_d in (18), we show the results of simulation in Fig. 2, where the loss-contour map is represented as the background. The redder, the bigger the loss function value is. In contrast, the bluer, the smaller the loss function value is. The x -axis and y -axis denote x_1 and x_2 , respectively. Both x_1 and x_2 are initialized to -10 . We use red and yellow lines to show the path of PID and SGD-M, respectively. It is obvious that SGD-M optimizer suffers from overshoot problem. By increasing K_d gradually (0.1, 0.5, and 0.93, respectively), our PID optimizer uses more "future" error so that it can largely alleviate the overshoot problem.

B. PID Optimizer for DNN

We are motivated by the simple example in Section IV-A and seek a PID optimizer to boost the convergence of DNN training. From (10), SGD-M can be viewed as a PI controller, which takes current and past gradient information actually. Fig. 2 shows that the PID controller introduces a derivative term of the gradient to use future information. Then, the overshoot problem can be alleviated obviously.

On the other hand, it is very easy to introduce noise when computing of gradients, because the training is often conducted in a mini-batch manner. We also try to estimate the average moving of the derivative part. Our proposed PID optimizer updates network parameter θ in iteration $(t + 1)$ by

$$\begin{cases} V_{t+1} = \alpha V_t - r \partial L_t / \partial \theta_t \\ D_{t+1} = \alpha D_t + (1 - \alpha)(\partial L_t / \partial \theta_t - \partial L_{t-1} / \partial \theta_{t-1}) \\ \theta_{t+1} = \theta_t + V_{t+1} + K_d D_{t+1}. \end{cases} \quad (19)$$

We could learn from (19) that a hyperparameter K_d is introduced in the proposed PID optimizer. We initialize K_d by introducing the Laplace transform [23] theory and the Ziegler-Nichols [52] tuning method.

C. Initialization of Hyperparameter K_d

The Laplace transform converts the function of real variable t to a function of complex variable s . The most common usage is to convert time to frequency. Denote the Laplace transformation of $f(t)$ as $F(s)$. There is

$$F(s) = \int_0^\infty e^{-st} f(t) dt, \quad \text{for } s > 0. \quad (20)$$

In general, it is easier to solve $F(s)$ than $f(t)$, which can be reconstructed from $F(s)$ with the inverse Laplace transform

$$f(t) = \frac{1}{2\pi i} \lim_{T \rightarrow \infty} \int_{\gamma - iT}^{\gamma + iT} e^{st} F(s) ds \quad (21)$$

where i is the unit of imagery part and γ is a real number.

By using the Laplace transform, we can first transform our PID optimizer into its Laplace transformed functions of s and then simplify the algebra. After obtaining the transformation $F(s)$, we can achieve the desired solution $f(t)$ with the inverse transform.

We initialize a parameter of a node in the DNN model as a scalar θ_0 . After enough times of updates, the optimal value of θ^* can be obtained. We simplify the parameter update in DNN optimization as a one-step response (from θ_0 to θ^*) in the control system. We introduce the Laplace transform to set K_d and denote the time domain change of weight θ as $\theta(t)$.

The Laplace transform of θ^* is θ^*/s [53]. We denote by $\theta(t)$ the weight at iteration t . The Laplace transform of $\theta(t)$ is denoted as $\theta(s)$, and that of error $e(t)$ as $E(s)$

$$E(s) = \frac{\theta^*}{s} - \theta(s).$$

The Laplace transform of PID [53] is

$$U(s) = \left(K_p + K_i \frac{1}{s} + K_d s \right) E(s). \quad (22)$$

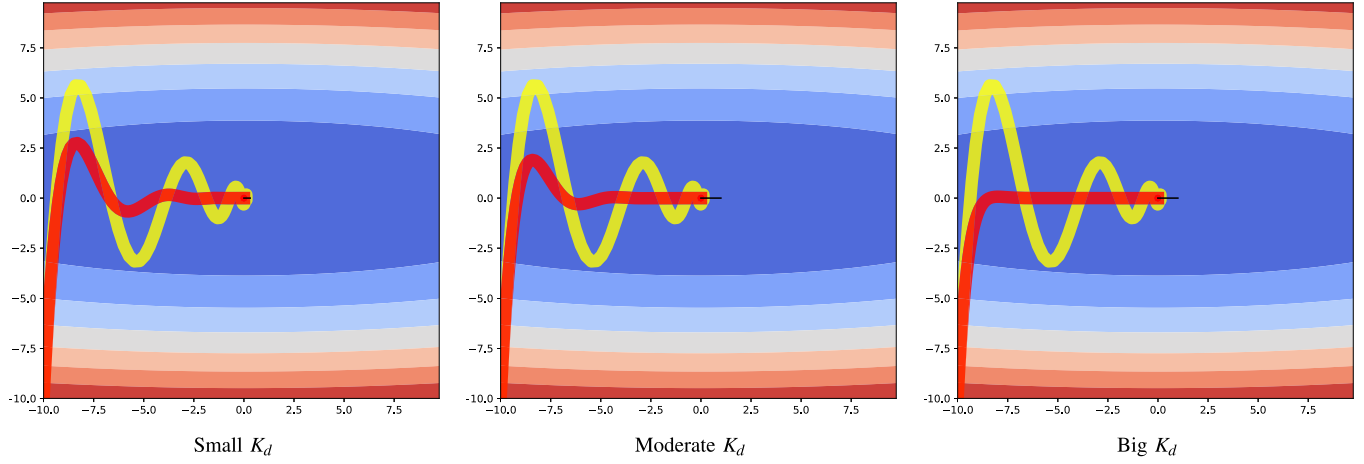


Fig. 2. Overshoot problem of momentum with different values of K_d . The red and yellow lines indicate the results obtained by PID and SGD-M, respectively.

In our case, $u(t)$ corresponds to the update of $\theta(t)$. Therefore, we replace $U(s)$ with $\theta(s)$ and with $E(s) = (\theta^*/s) - \theta(s)$. Equation (22) can be rewritten as

$$\theta(s) = \left(K_p + K_i \frac{1}{s} + K_d s \right) \left(\frac{\theta^*}{s} - \theta(s) \right). \quad (23)$$

With this form, it is easy to derive a standard closed-loop transfer function [54] as

$$\frac{\theta^*}{s} - \theta(s) = \frac{1}{K_d} \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \quad (24)$$

where

$$\begin{cases} \frac{K_p + 1}{K_d} = 2\zeta\omega_n \\ \frac{K_i}{K_d} = \omega_n^2. \end{cases} \quad (25)$$

Equation (24) can be rewritten as

$$\frac{\theta^*}{s} - \theta(s) = \frac{(s + \zeta\omega_n) + \frac{\zeta}{\sqrt{1-\zeta^2}}\omega_n\sqrt{1-\zeta^2}}{(s + \zeta\omega_n)^2 + \omega_n^2(1-\zeta^2)}. \quad (26)$$

We can get the time (iteration) domain form of $\theta(s)$ by using the inverse Laplace transform table [53] and the initial condition of θ (θ_0)

$$\theta(t) = \theta^* - \frac{(\theta^* - \theta_0) \sin(\omega_n \sqrt{1-\zeta^2} t + \arccos(\zeta))}{e^{\zeta\omega_n t} \sqrt{1-\zeta^2}} \quad (27)$$

and

$$\begin{cases} (K_p + 1)/K_d = 2\zeta\omega_n \\ K_i/K_d = \omega_n^2 \end{cases} \quad (28)$$

where ζ is damping ratio and ω_n is natural frequency of the system. The evolution process of a weight as an example of $\theta(t)$ is shown in Fig. 3. From (28), we get

$$K_i = \frac{(K_p + 1)^2}{4K_d\zeta}. \quad (29)$$

From (29), we know that K_i is a monotonically decreasing function of ζ . Based on the definition of overshoot in (16),

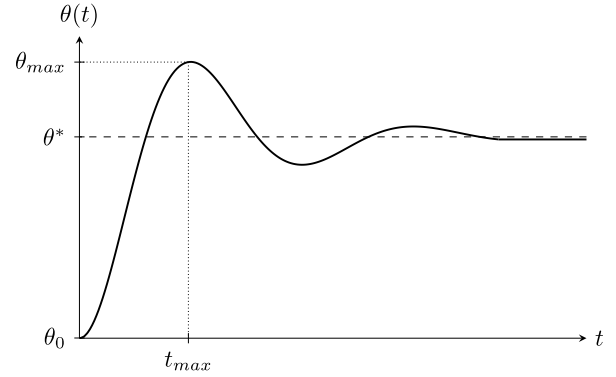


Fig. 3. Evolution process of the weight $\theta(t)$ for PID optimizer.

it is obvious that ζ is monotonically decreasing with overshoot. Then, K_i is a monotonically increasing function of overshoot. In a word, the more history error (integral part), the more overshoot problem in the system. This is a good explanation of why SGD-M overshoots its target and needs more training time.

By differentiating $\theta(t)$ w.r.t. time t , and let

$$\frac{d\theta(t)}{dt} = 0.$$

We have the peak time of the weight as

$$t_{\max} = \frac{\pi}{\omega_n \sqrt{1-\zeta^2}}. \quad (30)$$

Put t_{\max} into (27), we have θ_{\max} , and put θ_{\max} into (16), we have

$$\text{Overshoot} = \frac{\theta(t_{\max}) - \theta^*}{\theta^*} = e^{\frac{-\zeta\pi}{\sqrt{1-\zeta^2}}}. \quad (31)$$

We could learn from (27) that the term $\sin(\omega_n(1-\zeta^2)^{1/2}t + \arccos(\zeta))$ brings periodically oscillation change to the weight, which is no more than 1. The term $e^{-\zeta\omega_n t}$ mainly controls the convergence rate. It should be noted that the value of hyperparameter K_d in calculating the derivate

$$e^{-\zeta\omega_n} = e^{-\frac{K_p+1}{2K_d}}. \quad (32)$$

Based on the above-mentioned analyses, we know that the training of DNN can be accelerated by using large derivate. But on the other hand, if K_d is too large, the system will be fragile. After some experiments, we set K_d based on the Ziegler–Nichols optimum setting rule [52].

According to the Ziegler–Nichols rule, the ideal setup of K_d should be one third of the oscillation period (T), which means $K_d = (1/3)T$. From (27), we can get $T = (2\pi/(\omega_n(1-\zeta^2)^{1/2}))$. If we make a simplification that α in momentum is equal to 1, then $K_i = K_d = r$. Combined with (28), K_d will have a closed form solution

$$K_d = 0.25r + 0.5 + \left(1 + \frac{16}{9}\pi^2\right) / r. \quad (33)$$

For real-world cases, where different DNNs are applied train on different data sets, we would first start with this ideal setting of K_d and change it slightly latter.

V. EXPERIMENTAL RESULTS

We introduce four commonly used data sets for the experiments. Then, we compare our proposed optimizer with other optimizers by using CNN and LSTM on four commonly used data sets. Specifically, we first train a multilayer perceptron (MLP) on the MNIST data set to demonstrate the advantages of the PID optimizer. We then train CNNs on the CIFAR data sets to show that our PID optimizer achieves competitive accuracy compared with other optimizers, but it has a faster training speed. Further studies are carried out to prove that our PID optimizer also performs well on a larger data set. Based on the Tiny-ImageNet data set [55], we carry out a series of experiments. The results indicate that it is applicable for our PID optimizer to be extended to modern networks. Our proposed PID optimizer is set to use all hyperparameters that are detailed for SGD-M. The initial learning rate and learning rate schedule vary with different experiments.

A. Data Sets

1) *MNIST Data Set*: The MNIST data set [56] of handwritten numbers from 0 to 9. Being a subset of a larger data set NIST, MNIST consists of 60 000 training data, and 10 000 test ones. The digits have been size-normalized and centered in a fixed-size image of 28×28 pixels. With the usage of anti-aliasing techniques, the preprocessed images contain gray levels.

2) *CIFAR Data Sets*: The CIFAR10 data set [57] has 60 000 RGB color images, the shape of which is 32×32 . There are 10 classes, each of which includes 6000 images. For training and testing, 50 000 and 10 000 images are used, respectively. Similar to CIFAR10, the CIFAR100 data set [57] consists of 100 classes with 600 images for each class. From each class, 500 and 100 images are extracted for training and testing, respectively. The 100 classes in CIFAR100 [57] are further arranged into 20 superclasses. We performed random crops, horizontal flips, and padded four pixels around each side on the original image for data augmentation.

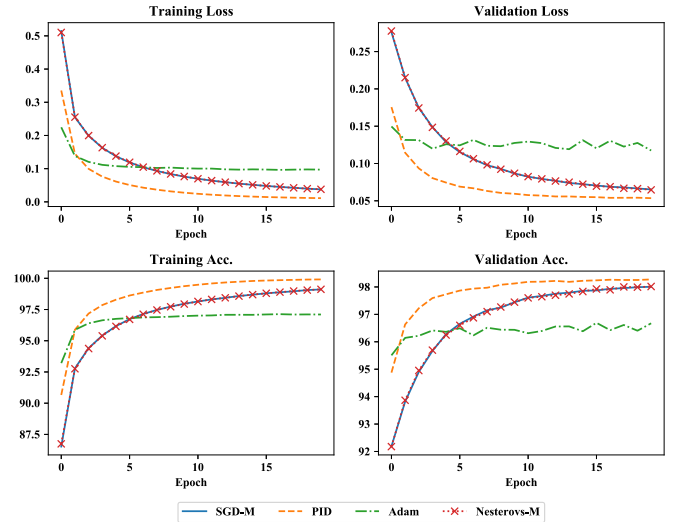


Fig. 4. Comparison between PID and other optimizers on the MNIST data set for 20 epochs. Top row: the curves of training loss and validation loss. PID optimizer obtains lower losses and converges faster. Bottom row: the curves of training accuracy and validation accuracy. PID optimizer performs much better than others for both training and test accuracies.

3) *Tiny ImageNet Data Set*: There are 200 classes in the Tiny-Imagenet [55] data set. Each class contains 500, 50, and 50 images for training, validation, and testing, respectively. Tiny-ImageNet is harder to be classed correctly than the CIFAR data set. It is not only because of a larger number of classes but also the relevant objects that need to be classified usually occupy little pixels of the whole image.

4) *PTB Data Set*: The Penn Treebank data set, known as the PTB data set, is widely used in machine learning of natural language processing research. The PTB data set has 2499 stories, which come from a three-year Wall Street Journal (WSJ) collection of 98 732 stories.

B. Results of CNNs

1) *Results of MLP on MNIST Data Set*: To compare the proposed PID optimizer with SGD-M [12], we first carry out a series of experiments. We use the MNIST data set to train a basic network, MLP. There are 1000 hidden nodes in the hidden layer. ReLU acts as a nonlinearity layer in the MLP network. We place the softmax layer on the top. The training batch size is 128 for 20 epochs. After running the experiments for ten times, we obtain the average results. Fig. 4 shows comparisons among the four methods in terms of training statistics. The Adam data set performs well in the early stages of training, but overall it could be very unstable and slower than the PID optimizer. As Fig. 4 shows, the PID optimizer performs faster convergence than other optimizers. Furthermore, the PID optimizer achieves lower loss and higher accuracy in both training and validation phases. In addition, it has stronger generalization ability on the test data set. It can be seen from Fig. 5 that the standard deviation of the PID optimizer during training is minimal, which proves its training stability. The accuracy is 98% in the PID optimizer and 97.5% in SGD-M.

TABLE I
COMPARISONS BETWEEN PID AND SGD-M OPTIMIZERS IN TERMS OF TEST ERRORS AND TRAINING EPOCHS.
WE REPORT THE RESULTS BASED ON CIFAR10 AND CIFAR100

Model	Depth-k	Params (M)	Runs	CIAFR10	Epochs	CIFAR100	Epochs
-	-	-	-	PID/SGD-M	PID/SGD-M	PID/SGD-M	PID/SGD-M
Resnet [42]	110	1.7	5	6.23 /6.43	239 /281	24.95 /25.16	237 /293
	1202	10.2	5	7.81 /7.93	230 /293	27.93/ 27.82	251 /296
PreActResNet [44]	164	1.7	5	5.23 /5.46	230 /271	24.17 /24.33	241 /282
	8-64	34.43	10	3.65/ 3.43	221 /294	17.46 /17.77	232 /291
ResNeXt29 [58]	16-64	68.16	10	3.42 /3.58	209 /289	17.11 /17.31	229 /283
	16-8	11	10	4.42 /4.81	213 /290	21.93 /22.07	229 /283
WRN [45]	28-20	36.5	10	4.27/ 4.17	208 /290	20.21 /20.50	221 /295
	100-12	0.8	10	3.83 /4.30	196 /291	19.97 /20.20	213 /294
DenseNet [43]	190-40	25.6	10	3.11 /3.32	194 /293	16.95 /17.17	208 /297

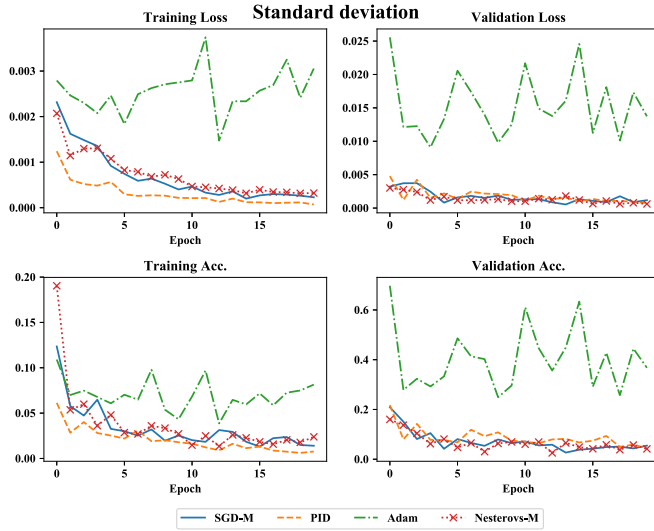


Fig. 5. PID versus other optimizers on the MNIST data set for 20 epochs. The standard deviation of ten runs. Top row: the curves of training loss and validation loss. Bottom row: the curves of training accuracy and validation accuracy.

2) *Results on CIFAR Data Sets*: In order to fully test our proposed PID optimizer, we compare it with the SGD-M optimizer on recent leading DNN models (ResNet [42], PreActResNet [44], ResNeXt29 [58], WRN [45], and DenseNet [43]). The details are shown in Table I, where the second column lists the depth of networks and k . The value of k in ResNeXt29, WRN, and DenseNet represents cardinality, widening factor, and growth rate, respectively. The third column lists the number of parameters. The fourth column shows the update numbers to calculate the mean test error. The next four columns show the average test error and the numbers of the epoch, when they accomplish the test errors firstly (the minimum number of epoch to reach the best accuracy).

The following conclusion can be given from Table I. First, compared with SGD-M, our PID optimizer obtains lower test errors for all architectures (except for ResNet with depth 1202) based on results from CIFAR10 and CIFAR100 data sets. Second, for the training epochs needed to reach the best results, the PID optimizer needs less training than SGD-M.

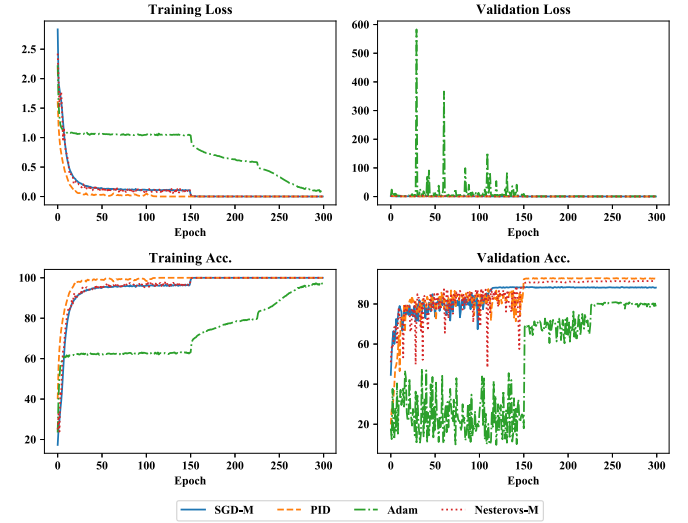


Fig. 6. Comparison among PID and other optimizers on the CIFAR10 data set by using DenseNet 190-40. Top row: the curves of training and validation loss. PID optimizer obtains lower losses and behaves more stable. Bottom row: the curves of training accuracy and validation accuracy. PID optimizer performs slightly better than SGD-M for both training and test accuracies.

Specifically, compared with SGD-M, our proposed PID optimizer achieves 35% and up to 50% acceleration on average. This reveals that the gradient descent's direction acts a very important role, which can be utilized to alleviate the overshoot problem and contribute to faster convergence for the training of DNNs. In Fig. 6, we further present more training statistics on CIFAR10 to compare PID and SGD-M optimizers. For the backbone DenseNet 190-40 [43], we set its network depth as 190 and growth rate as 40. Based on the experiments, we can obviously conclude that our PID optimizer achieves faster converges than SGD-M. More important, in both the training and validation phases, the PID optimizer obtains lower loss and higher accuracy.

3) *Results on Tiny-ImageNet*: We also apply our proposed PID optimizer on the Tiny-ImageNet data set with DenseNet 100-12 architecture to indicate its effectiveness. The initial learning rate of four optimizers is 0.1. The decreasing schedule is set to 50% and 75% of training epochs. The batch size is 500. In Fig. 7, we show the curves of training loss and accuracy and validation loss and accuracy over the change of

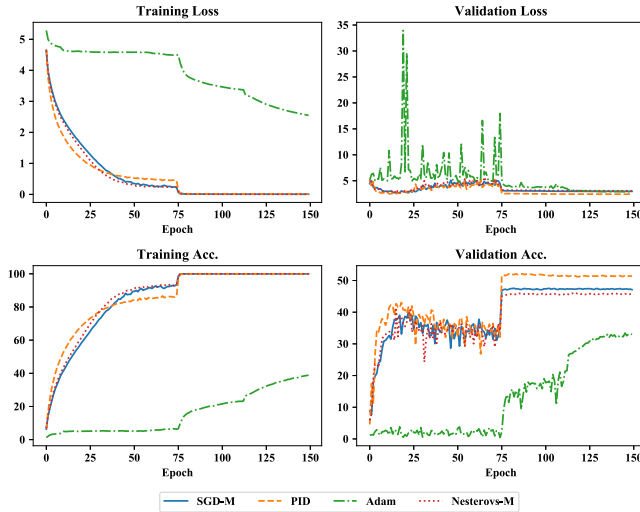


Fig. 7. Comparison among PID and other optimizers on the Tiny-imagenet data set with DenseNet 100–12 backbone. Top row: the curves of training and validation loss. We can see PID optimizer obtains lower training and validation losses. Bottom row: the curves of training accuracy and validation accuracy. PID optimizer achieves the best performance.

epochs for the four optimizers. Similar to results tested on the CIFAR data sets, the proposed PID optimizer not only converges faster but also obtains better performance. These results prove the generalization ability of our proposed PID optimizer.

C. Results of GANs

During the training of GANs, both G and D need to be trained. We train them both in an alternating manner. Each of their objectives can be expressed as a loss function that we can optimize via gradient descent. Therefore, we train G for a couple of steps, then train D for a couple of steps, then give G the chance to improve itself, and so on. The result is that the generator and the discriminator get better at their objectives in terms so that the generator can fool the most sophisticated discriminator finally. In practice, this method ends up with generative neural nets that are good at producing new data.

In the experiments, we use deep convolutional GANs (DCGANs) to test our proposed PID optimizer. The discriminator of this DCGAN consists of two convolutional layers (with ReLU function and max pooling) and two fully connected layers. The generator of this DCGAN consists of a fully connected layer (with batch normalization and ReLU function) and three convolutional layers. The binary cross entropy is used as a loss function. The learning rate is initialized to 0.0003 for all optimizers. The qualitative results of PID are shown in Fig. 8(b), and the SGD-M results are demonstrated in Fig. 8(a). From Fig. 8, we could find that the generated images with the PID optimizer are more realistic than these with the SGD-M optimizer.

D. Results of RNNs

In this experiment, we employ a simple LSTM that only has one layer with 100 hidden units. Mean squared error

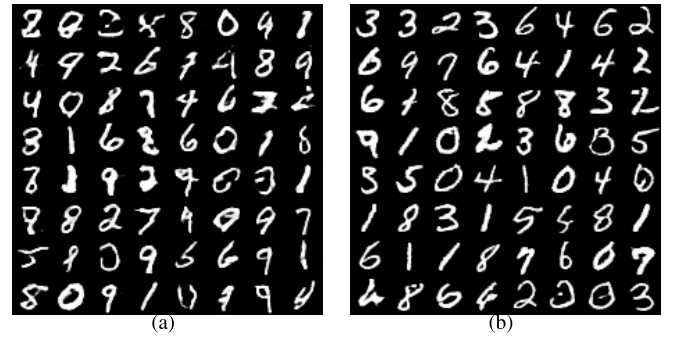


Fig. 8. PID versus SGD-M for generating images through GANs on MNIST data set. (a) Generated images from SGD-M. (b) Generated images from PID.

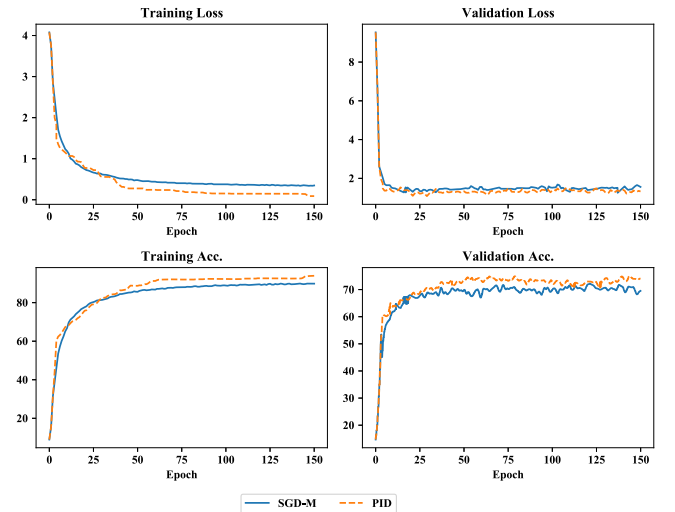


Fig. 9. Comparison between PID and SGD-M for the adding task of RNN. Top row: the curves of training and validation loss. PID optimizer achieves lower training and validation losses than SGD-M. Bottom row: the curves of training and validation accuracy. Our PID optimizer performs better in both training and test performance.

is used as the objective function for the adding problem. The initial learning rate is set to 0.002 for the SGD-M and PID optimizers. The learning rate is reduced by a factor of 10 every 20 000 training steps. We randomly generate all the training and testing data throughout the whole experiment. The results are shown in Fig. 9. The LSTM model with SGD-M has troubles in convergence. However, our proposed PID optimizer can reach to a small error with very fast convergence. It indicates that our proposed PID optimizer could effectively train LSTM.

Results on PTB Data Set: In this subsection, we evaluate the PTB-c data set to evaluate our proposed PID optimizer. We follow similar experimental settings as in [59]. Specifically, we apply the framewise batch normalization [60] and set batch size as 128. The learning rate is initially set to 0.0002 and decreases by $10\times$ when the validation performance no longer improves. We also introduce dropout [61] by using the dropping probability of 0.25 and 0.3. There is no overlapping in the sequences, whose length is set as $T = 50$ for both training and testing. Then, we train networks with the PID and SGD-M optimizers. The results are shown

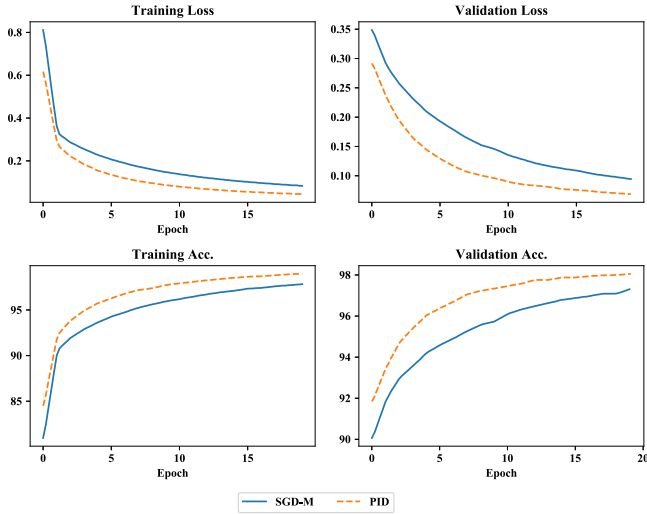


Fig. 10. Comparison between PID and SGD-M to train LSTM on PTB data set. Top row: the curves of training and validation loss. PID optimizer helps to achieve smaller training and validation losses. Bottom row: the curves of training and validation accuracy. PID optimizer helps to achieve higher training and validation performance.

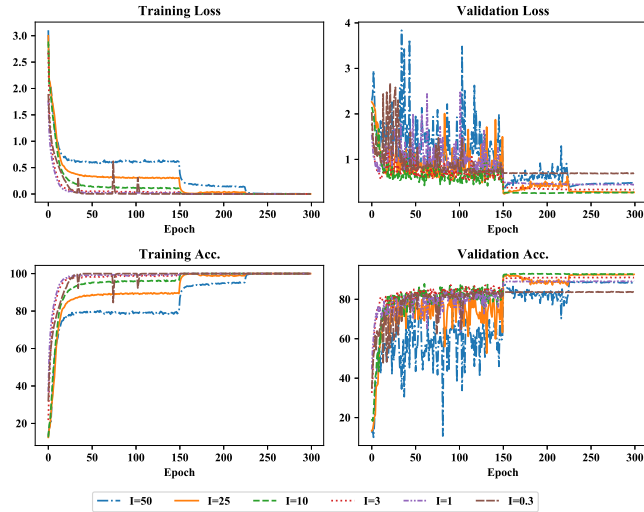


Fig. 11. Comparison among PID controllers with different K_i on the CIFAR10 data set by using DenseNet 100–12. K_d is fixed to 10. Top row: the curves of training and validation loss. Bottom row: the curves of training and validation accuracy. Within a certain range, larger K_i achieves better validation accuracy.

in Fig. 10. Comparing with the SGD-M, we can see that our proposed PID optimizer achieves better performance on the LSTM model.

E. Results of Different K_i and K_d

We also perform an ablation study on the hyperparameters of the PID controller. The experiments are run on the CIFAR10 data set with DenseNet 100–12. The initial learning rate is 0.1, and it is reduced by 10 in the 150 and 225 epochs.

The first group of experiments investigates the variation of training and verification statistics with K_i while K_d is fixed. Fig. 11 demonstrates six PID controllers whose K_d is 10. In the training, the performance of all controllers differs from each other at an early stage, but eventually, they can reach the same level. In validation, the controller with $K_i = 10$ achieves

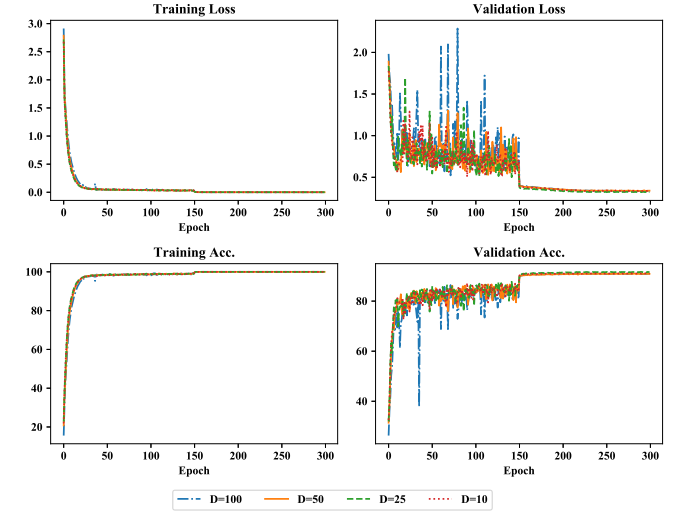


Fig. 12. Comparison among PID controllers with different K_d on the CIFAR10 data set by using DenseNet 100–12. K_i is fixed to 3. Top row: the curves of training and validation loss. Bottom row: the curves of training and validation accuracy.

the lowest loss and highest validation accuracy. We also repeat this experiment with $K_d = 10, 25, 50$, and 100, respectively, and results are highly similar to Fig. 11. One interesting phenomenon is that the larger the K_i , the more affected by the decreasing schedule.

Then, we change the research object to K_d . The settings of the second group of experiments are kept the same as previous experiments, but K_i is fixed. Fig. 12 shows that their performance is highly consistent. It is also shown that the larger K_d , the more unstable the validation performance. The reasons may be that large K_d leads to more change of optimization path.

As can be seen from these experiments, K_i is more important than K_d in these specific tasks (CIFAR10 with DenseNet100–12). K_i not only affects the speed of convergence but also affects the accuracy of verification.

VI. CONCLUSION

Motivated by the outstanding performance of the PID controller in the field of automatic control, we reveal the connections between the PID controller and stochastic optimizers and their variants. Then, we propose a new PID optimizer used in DNN training. The proposed PID optimizer reduces the overshoot phenomenon of SGD-M and accelerates the training process of DNNs by combining the present, the past, and the change information of gradients to update parameters. Our experiments on both image recognition tasks with the MNIST, CIFAR, and Tiny-ImageNet data sets and LSTM tasks with PTB data set validate that the proposed PID optimizer is 30% to 50% faster than SGD-M while obtaining lower error rate. We will continue to study the relationship among optimal hyperparameters (K_p , K_i , and K_d) in specific tasks. We will conduct more in-depth researches for more general cases in the future. We will also investigate how to associate the PID optimizer with an adaptive learning rate for DNNs/RNNs optimization in future works.

REFERENCES

- [1] O. Russakovsky *et al.*, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, Dec. 2015.
- [2] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. 19th Int. Conf. Comput. Statist.*, Paris, France, Aug. 2010, pp. 177–186.
- [3] J. Zhang, "Gradient descent based optimization algorithms for deep learning models training," Mar. 2019, *arXiv:1903.03614*. [Online]. Available: <https://arxiv.org/abs/1903.03614>
- [4] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus, "Exploiting linear structure within convolutional networks for efficient evaluation," in *Proc. Adv. Neural Inf. Process. Syst.*, 2014, pp. 1269–1277.
- [5] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," May 2014, *arXiv:1405.3866*. [Online]. Available: <https://arxiv.org/abs/1405.3866>
- [6] X. Zhang, J. Zou, K. He, and J. Sun, "Accelerating very deep convolutional networks for classification and detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 1943–1955, Oct. 2016.
- [7] N. Vasilache, J. Johnson, M. Mathieu, S. Chintala, S. Piantino, and Y. LeCun, "Fast convolutional nets with fbfft: A GPU performance evaluation," Dec. 2014, *arXiv:1412.7580*. [Online]. Available: <https://arxiv.org/abs/1412.7580>
- [8] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," in *Proc. Int. Conf. Learn. Representations*, 2015.
- [9] D. J. Im, M. Tao, and K. Branson, "An empirical analysis of the optimization of deep network loss surfaces," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2017.
- [10] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [11] L. Bottou, "Online learning and stochastic approximations," in *Online Learning in Neural Networks*, D. Saad, Ed., Cambridge, U.K.: Cambridge Univ. Press, 1998, pp. 9–42. [Online]. Available: <http://dl.acm.org/citation.cfm?id=304710.304720>
- [12] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Proc. Int. Conf. Mach. Learn.*, 2013.
- [13] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *J. Mach. Learn. Res.*, vol. 12, pp. 2121–2159, Feb. 2011.
- [14] G. Hinton, N. Srivastava, and K. Swersky, "Neural networks for machine learning lecture 6a overview of mini-batch gradient descent," Univ. Toronto, Toronto, ON, Canada, 2012. [Online]. Available: https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- [15] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014.
- [16] K. Ogata, *Discrete-Time Control Systems*. Upper Saddle River, NJ, USA: Prentice-Hall, 1995, vol. 2.
- [17] L. Lessard, B. Recht, and A. Packard, "Analysis and design of optimization algorithms via integral quadratic constraints," *SIAM J. Optim.*, vol. 26, no. 1, pp. 57–95, Jan. 2016.
- [18] L. Wang, W. Cluett, and T. Barnes, "New frequency-domain design method for PID controllers," *IEE Proc. Control Theory Appl.*, vol. 142, no. 4, pp. 265–271, Jul. 1995.
- [19] K. Heong Ang, G. Chong, and Y. Li, "PID control system analysis, design, and technology," *IEEE Trans. Control Syst. Technol.*, vol. 13, no. 4, pp. 559–576, Jul. 2005.
- [20] A. L. Salih, M. Moghavvemi, H. A. F. Mohamed, and K. S. Gaeid, "Modelling and PID controller design for a quadrotor unmanned air vehicle," in *Proc. IEEE Int. Conf. Autom., Qual. Test., Robot. (AQTR)*, vol. 1, May 2010, pp. 1–5.
- [21] P. Rocco, "Stability of PID control for industrial robot arms," *IEEE Trans. Robot. Autom.*, vol. 12, no. 4, pp. 606–614, 1996.
- [22] P. Zhao, J. Chen, Y. Song, X. Tao, T. Xu, and T. Mei, "Design of a control system for an autonomous vehicle based on adaptive-PID," *Int. J. Adv. Robot. Syst.*, vol. 9, no. 2, p. 44, Aug. 2012. [Online]. Available: <https://doi.org/10.5772/51314>
- [23] P. S. Laplace, *Theorie Analytique Des Probabilites*. Roubaix, France: Courcier, 1820.
- [24] W. An, H. Wang, Q. Sun, J. Xu, Q. Dai, and L. Zhang, "A PID controller approach for stochastic optimization of deep networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018.
- [25] B. Polyak, "Some methods of speeding up the convergence of iteration methods," *USSR Comput. Math. Math. Phys.*, vol. 4, no. 5, pp. 1–17, Jan. 1964.
- [26] Y. Wei *et al.*, "HCP: A flexible CNN framework for multi-label image classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 9, pp. 1901–1907, Sep. 2016.
- [27] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 142–158, Jan. 2016.
- [28] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2017.
- [29] W. Ouyang *et al.*, "DeepID-Net: Object detection with deformable part based convolutional neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 7, pp. 1320–1334, Jul. 2017.
- [30] C. Farabet, C. Couprie, L. Najman, and Y. Lecun, "Learning hierarchical features for scene labeling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1915–1929, Aug. 2013.
- [31] I. Goodfellow *et al.*, "Generative adversarial nets," in *Proc. NIPS*, 2014.
- [32] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [33] M. Arjovsky, A. Shah, and Y. Bengio, "Unitary evolution recurrent neural networks," in *Proc. 33rd Int. Conf. Mach. Learn. (ICML)*. New York, NY, USA, Jun. 2016, pp. 1120–1128. [Online]. Available: <http://jmlr.org/proceedings/papers/v48/arjovsky16.html>
- [34] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015.
- [35] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning efficient convolutional networks through network slimming," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 2736–2744.
- [36] Y. He, P. Liu, Z. Wang, Z. Hu, and Y. Yang, "Filter pruning via geometric median for deep convolutional neural networks acceleration," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2019, pp. 4340–4349.
- [37] X. Dai, H. Yin, and N. K. Jha, "NeST: A neural network synthesis tool based on a grow-and-prune paradigm," *IEEE Trans. Comput.*, vol. 68, no. 10, pp. 1487–1497, Oct. 2019.
- [38] X. Du, Z. Li, and Y. Cao, "Cgap: Continuous growth and pruning for efficient deep learning," May 2019, *arXiv:1905.11533*. [Online]. Available: <https://arxiv.org/abs/1905.11533>
- [39] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 6869–6898, 2017.
- [40] R. Kidambi, P. Netrapalli, P. Jain, and S. Kakade, "On the insufficiency of existing momentum schemes for stochastic optimization," in *Proc. Inf. Theory Appl. Workshop (ITA)*, Feb. 2018.
- [41] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, "Pruning filters for efficient convnets," Aug. 2016, *arXiv:1608.08710*. [Online]. Available: <https://arxiv.org/abs/1608.08710>
- [42] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016.
- [43] G. Huang, Z. Liu, L. V. D. Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017.
- [44] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proc. IEEE Eur. Conf. Comput. Vis. (ECCV)*, Oct. 2016, pp. 630–645.
- [45] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *Proc. Brit. Mach. Vis. Conf.*, 2016.
- [46] Y. Nesterov, "A method of solving a convex programming problem with convergence rate $O(1/k^2)$," *Sov. Math. Doklady*, vol. 269, no. 3, pp. 543–547, 1983.
- [47] J. C. Maxwell, "On governors," *Proc. Roy. Soc. London*, vol. 16, pp. 270–283, Dec. 1867.
- [48] N. Minorsky, "Directional stability of automatically steered bodies," *J. Amer. Soc. Nav. Eng.*, vol. 34, no. 2, pp. 280–309, Aug. 2010.
- [49] N. Qian, "On the momentum term in gradient descent learning algorithms," *Neural Netw.*, vol. 12, no. 1, pp. 145–151, Jan. 1999.
- [50] M. R. Spiegel, *Advanced Mathematics*. New York, NY, USA: McGraw-Hill, 1991.

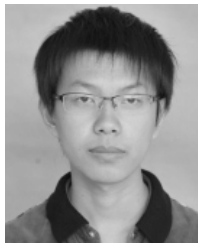
- [51] K. De Jong, "An analysis of the behavior of a class of genetic adaptive systems," Ph.D. dissertation, Dept. Comput. Commun. Sci., Univ. Michigan, Ann Arbor, MI, USA, 1975.
- [52] J. G. Ziegler and N. B. Nichols, "Optimum settings for automatic controllers," *Trans. ASME*, vol. 64, no. 11, pp. 759–765, 1942.
- [53] G. E. Robert and H. Kaufman, *Table Laplace Transforms*. Philadelphia, PA, USA: Saunders, 1966.
- [54] H. K. Khalil, *Nonlinear Systems*. Upper Saddle River, NJ, USA: Prentice-Hall, 1996.
- [55] F. F. Li, "Tiny imagenet visual recognition challenge," Sandford Artif. Intell. Lab., Stanford Univ., Stanford, CA, USA, 2015. [Online]. Available: <https://tiny-imagenet.herokuapp.com/>
- [56] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [57] A. Krizhevsky, "Learning multiple layers of features from tiny images," M.S. thesis, Univ. Toronto, Toronto, ON, Canada, 2009.
- [58] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017.
- [59] T. Coijmans, N. Ballas, C. Laurent, Ç. Gülçehre, and A. Courville, "Recurrent batch normalization," 2016, *arXiv:1603.09025*. [Online]. Available: <https://arxiv.org/abs/1603.09025>
- [60] C. Laurent, G. Pereyra, P. Brakel, Y. Zhang, and Y. Bengio, "Batch normalized recurrent neural networks," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Mar. 2016, pp. 2657–2661, doi: [10.1109/ICASSP.2016.7472159](https://doi.org/10.1109/ICASSP.2016.7472159).
- [61] Y. Gal and Z. Ghahramani, "A theoretically grounded application of dropout in recurrent neural networks," in *Proc. Adv. Neural Inf. Process. Syst. 29th, Annu. Conf. Neural Inf. Process. Syst.* Barcelona, Spain, Dec. 2016, pp. 1019–1027.



Haoqian Wang (Member, IEEE) received the B.S. and M.E. degrees from Heilongjiang University, Harbin, China, in 1999 and 2002, respectively, and the Ph.D. degree from the Harbin Institute of Technology, Harbin, in 2005.

He was a Post-Doctoral Fellow with Tsinghua University, Beijing, China, from 2005 to 2007. He has been a Faculty Member with the Graduate School at Shenzhen, Tsinghua University, Shenzhen, China, since 2008, where he has also been an Associate Professor since 2011 and the Director of Shenzhen

Institute of Future Media Technology. His current research interests include generative adversarial networks, video communication, and signal processing.



Yi Luo received the B.E. degree in Xidian University, Xi'an, China, in 2019. He is currently pursuing the M.E. degree with the Department of Automation, Tsinghua University, Beijing, China.

He is currently a Research Assistance with the Graduate School at Shenzhen, Tsinghua University, Shenzhen, China. His research interest includes optimization in deep learning.



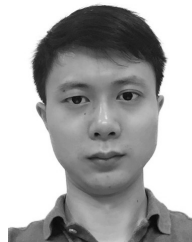
Wangpeng An received the B.E. degree from the Kunming University of Science and Technology, Kunming, China, in 2012, and the M.E. degree from the Department of Automation, Tsinghua University, Beijing, China, under the supervision of Prof. Q. Dai.

His research interests include face attribute recognition, generative adversarial networks, and deep learning optimization.



Qingyun Sun received the B.S. degree from the School of Mathematical Sciences, Peking University, Beijing, China, in 2014. He is currently pursuing the Ph.D. degree with the Department of Mathematics, Stanford University, Stanford, CA, USA.

His research interests include the mathematical foundation for artificial intelligence, data science, machine learning, algorithmic game theory, multi-agent decision making, optimization, and high-dimensional statistics.



Jun Xu received the B.Sc. degree in pure mathematics and the M.Sc. degree in information and probability from the School of Mathematics Science, Nankai University, Tianjin, China, in 2011 and 2014, respectively, and the Ph.D. degree from the Department of Computing, The Hong Kong Polytechnic University, in 2018, under the supervision of Prof. D. Zhang and Prof. L. Zhang.

He was a Research Scientist with the Inception Institute of Artificial Intelligence, Abu Dhabi, United Arab Emirates. He is currently an Assistant Professor with the College of Computer Science, Nankai University.



Lei Zhang (Fellow, IEEE) received the B.Sc. degree from the Shenyang Institute of Aeronautical Engineering, Shenyang, China, in 1995, and the M.Sc. and Ph.D. degrees in control theory and engineering from Northwestern Polytechnical University, Xi'an, China, in 1998 and 2001, respectively.

From 2001 to 2002, he was a Research Associate with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong. From 2003 to 2006, he was a Post-Doctoral Fellow with the Department of Electrical and Computer Engineering, McMaster University, Hamilton, ON, Canada. In 2006, he joined the Department of Computing, The Hong Kong Polytechnic University, as an Assistant Professor, where he has been the Chair Professor since 2017. His research interests include computer vision, image and video analysis, pattern recognition, and biometrics. He has authored or coauthored more than 200 articles in those areas.

Prof. Zhang is a Senior Associate Editor of the *IEEE TRANSACTIONS ON IMAGE PROCESSING* and an Associate Editor of the *SIAM Journal of Imaging Sciences and Image and Vision Computing*. He is a Clarivate Analytics Highly Cited Researcher from 2015 to 2018. His publications have been cited more than 38 000 times in the literature. More information can be found in his homepage: <http://www4.comp.polyu.edu.hk/cslzhang/>