
ADAPTIVELY CUSTOMIZING ACTIVATION FUNCTIONS FOR VARIOUS LAYERS

A PREPRINT

Haigen Hu^{1,2}, Aizhu Liu^{1,2}, Qiu Guan^{1,2}, Xiaoxin Li^{1,2}, Shengyong Chen^{1,3}, and Qianwei Zhou^{1,2,*}

¹College of Computer Science and Technology, Zhejiang University of Technology

²Key Laboratory of Visual Media Intelligent Processing Technology of Zhejiang Province

³School of Computer Science and Engineering, Tianjin University of Technology

ABSTRACT

To enhance the nonlinearity of neural networks and increase their mapping abilities between the inputs and response variables, activation functions play a crucial role to model more complex relationships and patterns in the data. In this work, a novel methodology is proposed to adaptively customize activation functions only by adding very few parameters to the traditional activation functions such as Sigmoid, Tanh, and ReLU. To verify the effectiveness of the proposed methodology, some theoretical and experimental analysis on accelerating the convergence and improving the performance is presented, and a series of experiments are conducted based on various network models (such as AlexNet, VGGNet, GoogLeNet, ResNet and DenseNet), and various datasets (such as CIFAR10, CIFAR100, miniImageNet, PASCAL VOC and COCO). To further verify the validity and suitability in various optimization strategies and usage scenarios, some comparison experiments are also implemented among different optimization strategies (such as SGD, Momentum, AdaGrad, AdaDelta and ADAM) and different recognition tasks like classification and detection. The results show that the proposed methodology is very simple but with significant performance in convergence speed, precision and generalization, and it can surpass other popular methods like ReLU and adaptive functions like Swish in almost all experiments in terms of overall performance. The code is publicly available at <https://github.com/HuHaigen/Adaptively-Customizing-Activation-Functions>. The package includes the proposed three adaptive activation functions for reproducibility purposes.

Keywords Adaptive activation function · Adaptable parameters · Various layers · Deep learning

1 Introduction

Activation functions play a key role during the training process of neural networks, and considerable attention has been paid to explore standard activation functions over the past years. Especially, with the remarkable development of Deep Neural Networks (DNN) in various computer vision applications, such as image classification [He et al., 2016, Krizhevsky et al., 2012, Tan et al., 2017], image segmentation [Chen et al., 2017], object detection [Girshick et al., 2014, Jiang et al., 2016, He et al., 2015], image enhancement [Lin et al., 2018, Tang et al., 2018], image retrieval [Yu et al., 2014, 2016] and tracking [Wu et al., 2016], Rectified Linear Unit (ReLU) [Nair and Hinton, 2010] has become extremely popular in the deep learning community in recent years. Owing to the significant improvements of ReLU in the deep neural networks, some extended versions are constantly springing up. For instance, Leaky ReLU (LReLU) [Maas et al., 2013] is proposed by replacing the negative part of the ReLU with a non-zero slope, while Exponential Linear Units (ELUs) [Clevert et al., 2015] can tend to converge cost to zero faster and produce more accurate results. All these extended versions can more or less achieve a certain effect in the respective fields.

However, there is hardly a generally accepted rule-of-thumb for the choice of activation functions owing to the fact that it solely depends on the problem at hand. Even the most popularly and commonly used activation function ReLU is not suitable for all datasets and network architectures. Therefore, adaptive activation functions have drawn more and more attention in recent years. For example, Maxout [Goodfellow et al., 2013] can approximate any convex

functions by selecting the maximum output value of multiple linear activation functions, but a large number of extra parameters are introduced, which causes large storage memory and high computation cost. In Parametric rectified linear unit (PReLU) [He et al., 2015], the slopes of negative part can be obtained by learning from data rather than the pre-defined fixed values, thus PReLU has theoretically all the advantages of ReLU and effectively avoids Dead ReLU. But in practice, it has not been fully confirmed that PReLU always surpasses ReLU. In 2017, an activation function with the property of "self-normalization" is proposed, named SELU [Klambauer et al., 2017], and it can avoid the problem of gradient vanishing and exploding, thereby leading to the feedforward neural network to obtain beyond state-of-the-art performance. However, the effectiveness of SELU in Convolutional neural networks (CNN) has not been confirmed. In the same year, swish [Ramachandran et al., 2017] with some complex characteristics, such as no upper and lower bound, smooth and non-monotonic, can perform better than ReLU on many deep models.

Although the existing adaptive activation functions are relatively more flexible than the traditional activation function owing to the adaptability, and have already achieved great improvements, they are limited to some specific application scenarios, and there are still many problems to be solved, such as low generalization capability and poor precision performance. For example, their performance often depends on some specific network models and data sets. In this work, a novel methodology is proposed to explore the optimal activation functions with more flexibility and adaptability only by adding few additional parameters to the traditional activation functions such as Sigmoid, Tanh and ReLU. The proposed methodology can avoid local minimums and accelerate convergence only by introducing very few parameters to the fixed activation functions, thereby increasing the precision, reducing the training cost and improving the generalization performance.

The primary contributions of our work are summarized as follows:

- A novel methodology is proposed to customize activation functions with more flexibility and adaptation for various layers only by introducing very few parameters to the traditional activation functions such as Sigmoid, Tanh, and ReLU.
- A theoretical analysis for accelerating the convergence and improving the performance is presented by taking an activation function of one layer as an example without loss of generality, and an experimental study is performed by comparing the weight increments between two successive epochs in different layers during the training process between the proposed ARELU and ReLU on CIFAR100 based on VGGNet.
- The proposed *ARELU* is a generalized form of the *ReLU-based* versions, while *ReLU* and *PReLU* are the special cases of the proposed *ARELU*.

The rest of the paper is organised as follows. Section 2 introduces the related work, and the proposed methodology is presented in Section 3. Section 4 presents the analysis for our methodology. Section 5 details the experimental results for comparison and validation. Section 6 concludes the paper.

2 Related work

Over the last few decades, many various activation functions have been proposed in the artificial neural network community. According to whether the parameter or shape of an activation function is learnable or variable during the training phase, activation functions can be divided in two categories: fixed activation functions and adaptive activation functions.

2.1 Fixed activation functions

Fixed activation functions indicate that the parameters or shapes can not be modified during the training phase (shown in Fig. 1), and the most common fixed activation functions can be fallen into three categories: Logistic (Sigmoid), Hyperbolic Tangent (Tanh) and Rectified Linear Activation (ReLU).

Sigmoid Sigmoid function is a common S-like function or S-like growth curve, and is normally used to refer specifically to the logistic function. It can map any real value to the range [0,1], thereby being interpreted as a probability, defined as follows:

$$\text{Sigmoid}(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

It is differentiable, and the derivative is derived as follows:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (2)$$

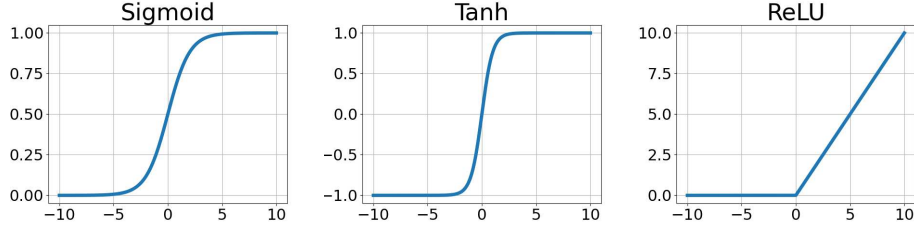


Figure 1: An illustration of fixed activation functions with a fixed shape.

Note that the gradient $\sigma'(x) \rightarrow 0$ as $\sigma(x) \rightarrow 0$ or $\sigma(x) \rightarrow 1$, meaning that, when the output of Sigmoid saturates for a large positive or negative inputs (i.e. the curve becomes parallel to x -axis shown in Fig. 1), the gradients are almost zero. Due to the zero gradient, the weights are no longer updated and the networks will not learn, thus the neuron dies, thereby causing the vanishing gradient problem. Besides, Sigmoid outputs are not zero-centered, and it can indirectly introduce undesirable zig-zagging dynamics in the gradient updates for the weights.

Tanh Tanh function, a hyperbolic tangent function, graphically looks very similar to Sigmoid. Actually, the Tanh is simply a scaled Sigmoid, such that its outputs range from -1 to 1, defined as follows:

$$\text{Tanh}(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(2x) - 1 \quad (3)$$

Like the Sigmoid, Tanh is also affected by the vanishing gradient problem. But unlike the Sigmoid, its output is zero-centered, the negative inputs will be mapped strongly negatives and the zero inputs will be mapped near zero. Therefore, the non-linearity of Tanh is always preferred to that of Sigmoid, and it has been widely used in deep learning & machine learning, especially in classification scenarios between two classes.

ReLU ReLU is a very simple and efficient activation function that has been widely used in almost all deep learning domains, especially in CNNs, defined as

$$\text{ReLU}(x) = \max(0, x) \quad (4)$$

Owing to the simpler mathematical operations, ReLU is far more computationally efficient than Tanh and Sigmoid. Besides, ReLU can solve parts of the saturation problem only in the positive region. Whereas for the negative inputs, the results contain one or more true zero values (called a sparse representation) to accelerate learning and simplify the model in representational learning, but the weights and biases are not updated owing to the zero gradient during the backpropagation process, thereby causing the dying ReLU problem.

2.2 Adaptive activation functions

Adaptive activation functions refer primarily to the functions that the parameters or shapes are trained and learned along with other parameters in neural networks (shown in Fig. 2), thereby adaptively varying with training data. In other words, the main idea of this kind of functions is to search a good function shape using knowledge given by the training data. For example, PReLU [He et al., 2015] replaces the fixed slope α of LReLU [Maas et al., 2013] with a trainable parameter α_i in the negative region. Whereas Swish [Ramachandran et al., 2017] is a recently proposed activation function with no upper bound, lower bound, smooth, and non-monotonic characteristic, and it can be loosely viewed as a bridging function between the linear function and the ReLU function. Other similar activation functions like FReLU [Qiu et al., 2018] and PELU [Trottier et al., 2017] have achieved performance improvements in some specific tasks.

Although the existing adaptive activation functions has shown to improve the network performances significantly, thanks to properties such as no saturation feature, flexibility and adaptivity, exploring the optimal and appropriate activation functions is still an open field of research, and there is still potential room for improvement in various scenarios, especially for complex datasets and different models.

3 Methodology

The training of neural networks is essentially a non-convex optimization problem, in which the optimal weight parameters can be searched and found by using the back-propagation algorithm, so that the functional subspace will be

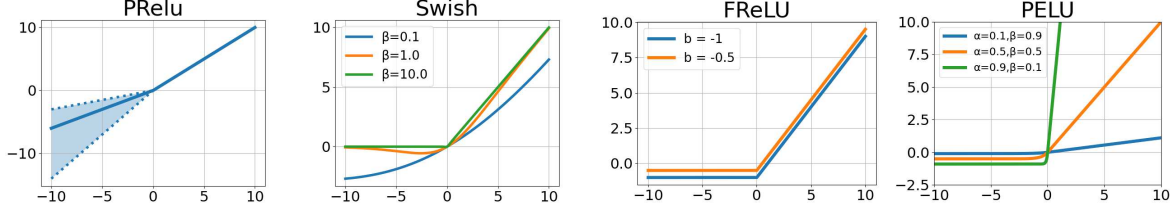


Figure 2: An illustration of various adaptive activation functions. The shapes of activation functions can be controlled and adjusted by some parameters, which are trained along with other parameters in the neural networks.

explored and determined by the activation function. Adaptive activation functions refer to the functions that adapt themselves to the network inputs, therefore they can learn hyper-parameters to adapt the parameters of the affine transformation to a given input, and thereby increase the flexibility and the representation ability of network models.

In this work, we attempt to construct a new parameter learning method for each layer only by introducing a few parameters to the fixed activation functions, and the general form in the i^{th} layer with activation functions can be defined as follows

$$f_A(a_i, b_i, c_i, d_i, z) = b_i f(a_i z + c_i) + d_i \quad (5)$$

where $f(\cdot)$ represents a traditional activation function (fixed activation function). a_i , b_i , c_i and d_i are four learnable parameters in the i^{th} layer, and they can adapt to the different tasks according to the complexity of input data so as to efficiently avoid falling into local minimums. z denotes the weighted sum of inputs, including the bias term, defined as

$$z = wx + b \quad (6)$$

w and b indicate weights and bias, respectively. x is an input vector.

In practice, the proposed adaptive activation function is very simple, and it is composed of two embedded linear equations, namely: internal linear equation

$$f_{in}(a_i, c_i, z) = a_i z + c_i \quad (7)$$

and external linear equation

$$f_{ex}(b_i, d_i) = b_i f_{in} + d_i \quad (8)$$

Therefore Equation (1) is rewritten as

$$f_A(a_i, b_i, c_i, d_i, z) = f_{ex}(b_i, d_i) = b_i f_{in}(a_i, c_i, z) + d_i \quad (9)$$

In the following sections, the effectiveness and advantages of the proposed methodology will be verified by taking some common fixed activation functions as baselines, such as *Sigmoid*, *Tanh* and *ReLU*, and the corresponding adaptive activation functions are named *ASigmoid*, *ATanh* and *ARELU*, respectively. According to Equation (5), these functions are respectively defined as

$$ASigmoid : f_{Asigmoid} = b_i Sigmoid(a_i z + c_i) + d_i \quad (10)$$

$$ATanh : f_{Atanh} = b_i tanh(a_i z + c_i) + d_i \quad (11)$$

$$ARELU : f_{Arelu} = maximum(a_i z + c_i, b_i z + d_i) \quad (12)$$

In *ASigmoid* and *ATanh*, a_i and c_i are respectively used to scale the inputs of Sigmoid and Tanh, while b_i and c_i scale the outputs, simultaneously.

Significantly, when $a_i = 1$ and $b_i = c_i = d_i = 0$, the negative part of the *ARELU* is replaced with a zero slope, while the slope of the positive part is fixed. In this case, *ARELU* is actually degenerated to a standard *ReLU*, given as

$$ReLU : f_{ReLU} = maximum(0, z) \quad (13)$$

Furthermore, when $b_i = 1$ and $c_i = d_i = 0$, the slope (i.e., the parameter a_i) of the negative part is adjustable, which means that the parameter can learn from data rather than be obtained by the pre-defined. Under these conditions, *ARELU* is evolved to *PReLU* when , given as

$$PReLU : f_{Prelu} = maximum(az, z) \quad (14)$$

Therefore, *ARELU* is a generalized form of the *ReLU-based* versions, while *ReLU* and *PReLU* are the special cases of the proposed *ARELU*.

Above, it can be clearly seen that our method only adds four parameters for each layer. For the entire network model, $4i$ parameters should be added. This parameter amount and calculation amount is negligible compared with the entire network model.

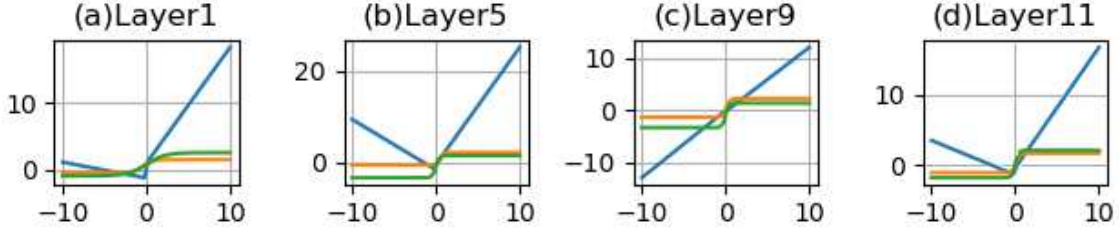


Figure 3: The shapes of ASigmoid (green), ATanh (orange) and AReLU (blue) at different layers during the training process on CIFAR100 based on VGG.

4 Analysis

A convex loss function $L(\cdot)$ for the linear weighted combination of each activation function applied to an input is defined to find the optimal weights by adopting suitable optimization strategies based on the back-propagation algorithm. Thus, the training process of the network model is essentially an iterative optimization process for the weight parameters by minimizing the loss function $L(\cdot)$ in the functional subspace.

4.1 Theoretical Analysis

In order to facilitate analysis, we just take an activation function of one layer as an example without loss of generality. Suppose that a neural network with traditional activation functions is given as

$$y = f(z) \quad (15)$$

For the update process of weights, the partial derivative chain is defined as follows.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w} \quad (16)$$

$$\frac{\partial y}{\partial z} = f'(z) \quad (17)$$

Meanwhile, the weight is updated as follows.

$$\tilde{w} = w - \eta \frac{\partial L}{\partial w} \quad (18)$$

where η is learning rate. Equation (16) and (17) are substituted into (18), we can obtain the weights update equation for a common activation function as follows.

$$\tilde{w} = w - \eta f'(z) \frac{\partial L}{\partial y} \frac{\partial z}{\partial w} \quad (19)$$

Considering that the proposed adaptive activation functions consist of two linear equations, for simplicity, we just consider the internal linear function and omit its intercept term in a certain layer, given as

$$\tilde{y} = f(az) \quad (20)$$

where \tilde{y} is the output of an adaptive activation function. The hyperparameter a represents a generalization form for scale inputs in any layers, and it can be inferred to fine-tune the learning rate so as to speed up the update of weights, and the corresponding derivation is given as follows.

For the output \tilde{y} , the partial derivative is given as

$$\frac{\partial \tilde{y}}{\partial z} = a f'(z) \quad (21)$$

With Equations (16), (18) and (21), the update process of the weight is given as

$$\tilde{w}_1 = w - \eta a f'(z) \frac{\partial L}{\partial \tilde{y}} \frac{\partial z}{\partial w} \quad (22)$$

By comparison between Equations (19) and (22), the learning rate of the adaptive activation function can be written as:

$$\tilde{\eta} = \eta a \quad (23)$$

From Equation (23), we can adaptively adjust the learning rate $\tilde{\eta}$ by using the hyperparameter a . Simultaneously, the optimization of the hyperparameter a in neural networks is similar to the hyperparameter w , then the update process of a is achieved by using the chain rule.

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial a} \quad (24)$$

$$\frac{\partial \tilde{y}}{\partial a} = z f'(z) \quad (25)$$

With Equations (24) and (25)

$$\tilde{a} = a - \eta z f'(z) \frac{\partial L}{\partial \tilde{y}} \quad (26)$$

With Equations (6) and (26)

$$\tilde{a} = a - \eta w x f'(z) \frac{\partial L}{\partial \tilde{y}} \quad (27)$$

Therefore, the adaptive activation function is to achieve rapid convergence by the adapting learning rate, and this method is achieved by adjusting the weight w and the parameter a mutually to speed up learning in neural networks and lead to higher classification precision.

Besides, the internal linear equation has also its respective intercept, which contributes to tuning the parameters from another vertical direction during training process, thereby avoiding involving local extremum.

Similarly, the external linear equation has the same effects for accelerating the convergence and improving the performance. More importantly, the internal equation is embedded within the external, such case will enable the optimization toward a global optimum solution more efficiently from all directions.

4.2 Experimental Analysis

Owing to the fact that each layer has its own independent activation function, the optimal hyperparameters of each layers can obtained by learning from the respective complexity of input data, and their values will vary with the input data characteristics. Therefore, the obtained weights will be optimal, and the corresponding activation functions are different with different layers. Fig. 3 shows the visualization of different layers during the training process on CIFAR100 based on VGG, and different shapes of ASigmoid, ATanh and AReLU at different layers indicate that these functions can learn the optimal hyperparameters from the inputs of the respective layers, which would lead to the enhancement and improvement of the fitting capability and the accuracy of the networks.

Moreover, compared with the traditional adaptive activation functions, two embedded linear equations with intercepts can accelerate the weights adjustment. For further verification, the change curves of the weight increments Δw between two successive epochs in various layers are visualized along with the training process (shown in Fig. 4). The results clearly show the amplitudes of the increments Δw by using the proposed AReLU are much larger than those of the traditional ReLU in the early training stages, then the increments of the two methods converge, which means the proposed methods can provide faster weight updates than the traditional methods. Consequently the proposed methods can improve greatly the convergence speed and reduce the computational burden. Meanwhile, the large amplitudes of the increments can also help to avoid falling into a local optimum when training artificial neural networks with gradient-based learning methods and backpropagation.

5 Experiments

In this section, a series of experiments are implemented to verify and evaluate the effectiveness of the proposed methodology based on the three baseline activation functions such as Sigmoid, Tanh and ReLU. Considering the fact that ReLU is the most common activation function used in neural networks, there exist many derivatives of ReLU, and some typical derivatives like LReLU and PReLU are selected for comparison to highlight the effectiveness of the parameterization method in the activation function. Whereas swish, as an outstanding activation function, is used to demonstrate the state-of-the-art performance of the proposed adaptive activation function. Firstly, many comparison experiments between the proposed functions and its corresponding baseline functions have been conducted by

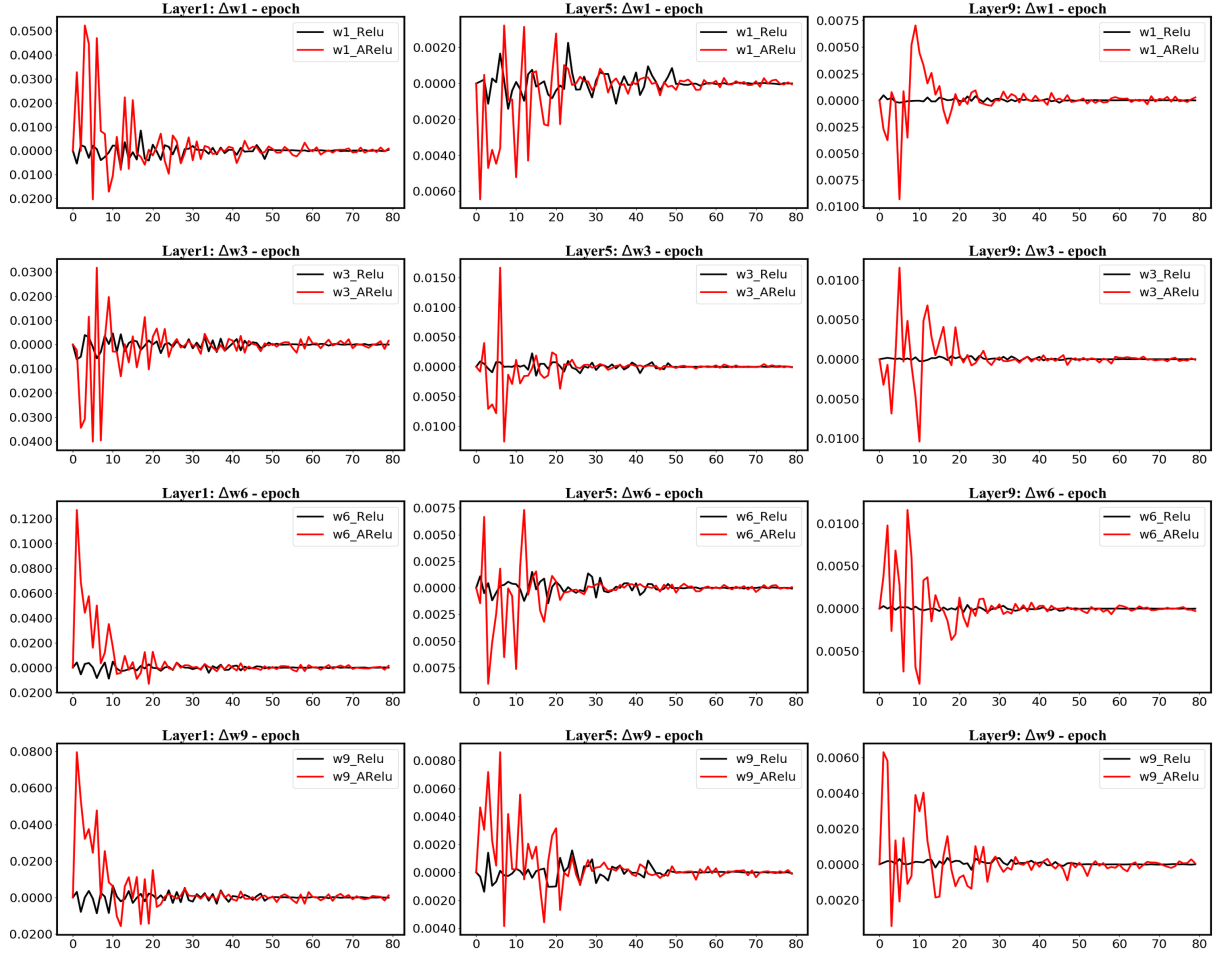


Figure 4: The change curves of the weight increments Δw between two successive epochs in different layers during the training process by using ARELU and ReLU on CIFAR100 based on VGGNet. Column 1,2 and 3 represent the 1th, 5th and 9th layer, respectively. Rows 1-4 illustrate four various weights by using ARELU and ReLU.

using Stochastic gradient descent (SGD) [Cramer, 1946] on the datasets of CIFAR10 and CIFAR100 based on different network models, such as AlexNet [Krizhevsky et al., 2012], VGG [Simonyan and Zisserman, 2014], GoogleNet [Szegedy et al., 2015], ResNet [He et al., 2016] and DenseNet [Huang et al., 2017]. Then, some experiments are implemented to further verify the validity and suitability in various optimization strategies, such as SGD, Momentum [Qian, 1999], AdaGrad [Duchi et al., 2011], AdaDelta [Zeiler, 2012] and ADAM [Kingma and Ba, 2014]. Finally, a series of comparison experiments are conducted to further verify the effectiveness, suitability and generalization ability on other more complicated datasets like miniImageNet [Vinyals et al., 2016], PASCAL VOC [Everingham et al., 2012] and COCO [Lin et al., 2014].

5.1 Experimental setup

We test the proposed adaptive activation functions on CIFAR10 and CIFAR100 based on AlexNet, VGGNet, GoogleNet, ResNet and DenseNet. The detailed experimental setup is illustrated in Fig. 5.

Note that Dense block consists of $\{\overset{bn}{\rightarrow} conv(1, 1, 2k) \xrightarrow{drop0.8, bn} conv(3, 3, 2k) \xrightarrow{drop0.2}\} \times 6, \times 12, \times 48, \times 32$, respectively, and the transition layer is shown in Fig. 3, which corresponds to the sequence $conv(1, 1) - avg_pool(2, 2)$. Moreover, the growth rate is $k=24$ for all.

AlexNet		VggNet		GoogleNet		ResNet		DenseNet	
input	32 x 32 x3	input	32 x 32 x3	input	32 x 32 x3	input	32 x 32 x3	input	32 x 32 x3
conv	1 x 11 x 96	conv	3 x 3 x 64	conv	7 x 7 x 64	conv	7 x 7 x 64	conv	7 x 7 x 2k
LRN		conv	3 x 3 x 64	maxpool	3 x 3	bn		BN	
maxpool	3 x 3	bn		conv	3 x 3 x 192	maxpool	3 x 3	conv	1 x 1 x 2k
conv	5 x 5 x 256	conv	3 x 3 x 128	conv	3 x 3 x 192	conv	1 x 1 x 64	dropout	0.8
LRN		conv	3 x 3 x 128	maxpool	3 x 3	conv	3 x 3 x 64 x 3	BN	
maxpool	3 x 3	bn		inception_3a		conv	1 x 1 x 256	conv	3 x 3 x k
conv	3 x 3 x 384	conv	3 x 3 x 256	inception_3b		conv	1 x 1 x 128	dropout	0.8
conv	5 x 5 x 256	conv	3 x 3 x 256	maxpool	3 x 3	conv	3 x 3 x 128 x 4	Transition Layer	
fc	2048	conv	3 x 3 x 256	inception_4a		conv	1 x 1 x 512	BN	
dropout	0.9	bn		inception_4b		conv	1 x 1 x 256	conv	1 x 1 x 2k
fc	10/100	maxpool	2 x 2	inception_4c		conv	3 x 3 x 256 x 6	dropout	0.8
output	10/100	conv	3 x 3 x 512	inception_4d		conv	x 1 x 1024	BN	
		conv	3 x 3 x 512	inception_4e		conv	1 x 1 x 512	conv	3 x 3 x k
		conv	3 x 3 x 512	maxpool	3 x 3	conv	3 x 3 x 512 x 3	dropout	0.8
		bn		inception_5a		conv	x 1 x 2048	Transition Layer	
		maxpool	2 x 2	inception_5b		avg_pool	7 x 7	BN	
		conv	3 x 3 x 512	avg_pool	8 x 8	fc	10/100	conv	1 x 1 x 2k
		conv	3 x 3 x 512	dropout	0.9	output	10/100	dropout	0.8
		conv	3 x 3 x 512	fc	10/100			BN	
		bn		output	10/100			conv	3 x 3 x k
		maxpool	2 x 2					dropout	0.8
		fc	8192					Transition Layer	
		dropout	0.9					BN	
		fc	7168					conv	1 x 1 x 2k
		dropout	0.9					dropout	0.8
		fc	10/100					BN	
		output	10/100					conv	3 x 3 x k
								dropout	0.8
								BN	
								GAP	7 x 7
								fc	10/100
								output	10/100

Figure 5: The network architectures of AlexNet, VGGNet, GoogleNet, ResNet and DenseNet on CIFAR10 and CIFAR100.

5.2 CIFAR10

In these experiments, the proposed adaptive activation functions (10)~(12) are applied on CIFAR10 dataset based on the models (shown in Fig. 5), respectively. All trainings are implemented for no less than 80 epochs with a 64-batch size and without data augmentation by using SGD with fixed learning rate schedule of 0.001, 0.0001 and 0.00001 with the training process, respectively.

Fig. 6 shows the convergence curves (top row) and the area enclosed by the convergence curves (bottom row) during the training process. It is obvious that the smaller the area, the faster the convergence speed, and it is also clearly shown from the results that the proposed activation functions can surpass the corresponding baseline functions and LReLU for different network models in terms of convergence speed. Thereinto, AReLU can obtain the fastest convergence speed among these activation functions, especially on DenseNet, ResNet and VGG16, it can converge much faster than other activation functions. Tables 1 and Table 2¹ show the quantified results in precision, and the proposed methodology has an overall advantage. Table 1 illustrates the comparison results between the proposed methods and their respectively corresponding baselines. From the results, the proposed methodology can effectively applied to the classic fixed activation functions, and can surpass the corresponding baseline functions on most network models. For instance, ASigmoid can surpass the corresponding baseline function Sigmoid on all models, while ATanh can also obtain better precision than the corresponding Tanh on the models of VGGNet, ResNet and DenseNet, and the obtained precision of ATanh in AlexNet and GoogleNet are only slightly lower than those of the corresponding Tanh. Table 2 shows the comparison results between AReLU and other adaptive activation functions. AReLU, except PReLU on the AlexNet and VGGNet, can overall obtain higher precision than other adaptive functions on various models. Note that some traditional adaptive functions like PELU and FReLU are not suitable for some network models owing to the lack of convergence during the training process. While the proposed methodology can apply to various deep learning models and has better generalization performance than other traditional methodologies.

¹Indicates that the activation function does not converge in this mode.

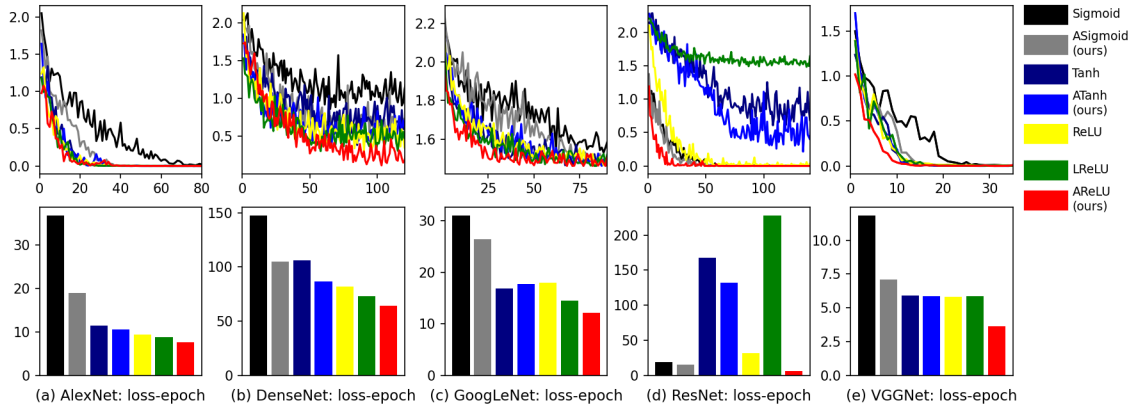


Figure 6: The top row represents the loss-epoch convergence curves for various activation functions on CIFAR10 by using SGD optimizer on various models during the training process, such as AlexNet, DenseNet, VGGNet, ResNet and DenseNet. Whereas the bottom row illustrates the area enclosed by the corresponding convergence curves in top row, and the smaller area means the faster convergence speed.

Table 1: The classification precision of various fixed activation functions for different models on CIFAR10.

Methods	AlexNet	VGGNet	GoogleNet	ResNet	DenseNet
Sigmoid	0.8141±1.52e-3	0.8680±1.17e-3	0.8027±1.13e-3	0.8079±1.25e-3	0.7611±3.86e-3
ASigmoid(ours)	0.8469±1.68e-3	0.8819±2.28e-4	0.8517±4.07e-4	0.8087±1.36e-3	0.9076±8.16e-4
Tanh	0.8043±9.68e-4	0.8841±7.72e-4	0.8433±6.06e-4	0.5784±9.75e-4	0.8966±1.99e-5
ATanh(ours)	0.7900±5.31e-4	0.8904±3.05e-4	0.8342±7.99e-4	0.6014±8.11e-4	0.9163±1.87e-4
ReLU	0.8351±3.95e-4	0.8800±1.01e-4	0.8733±4.36e-5	0.6641±2.59e-3	0.9242±1.10e-4
LReLU[Maas et al., 2013]	0.8255±1.66e-3	0.9253±4.68e-5	0.8742±9.31e-4	0.6698±1.27e-3	0.9289±2.37e-5
AReLU(ours)	0.8331±2.48e-4	0.8328±1.26e-3	0.8773±1.83e-3	0.9230±5.51e-4	0.9538±5.58e-4

Table 2: The classification precision of various adaptive adaptive activation functions for different models on CIFAR10. The top results are highlighted in black bold and the second-best results in blue.

Methods	AlexNet	VGGNet	GoogleNet	ResNet	DenseNet
PReLU[He et al., 2015]	0.8558±4.61e-4	0.9344±5.89e-5	0.8551±1.79e-4	0.6522±3.20e-4	0.9231±1.40e-4
Swish[Ramachandran et al., 2017]	0.7557±2.20e-3	0.9171±7.98e-4	0.8710±3.25e-3	0.6446±3.39e-4	0.9276±1.87e-4
PELU[Trottier et al., 2017]	— ¹	—	0.8128±4.45e-4	0.7891±3.91e-4	—
FReLU[Qiu et al., 2018]	—	0.8558±2.44e-4	0.8694±5.21e-3	0.8726±3.71e-4	0.8214±4.11e-4
AReLU(ours)	0.8331±2.48e-4	0.8328±1.26e-3	0.8773±1.83e-3	0.9230±5.51e-4	0.9538±5.58e-4

5.3 CIFAR100

To further verify the validity and applicability, CIFAR-100 is selected to training based on several typical network models, such as AlexNet, DenseNet (shown in Fig. 5) and VGG-v. The Dataset is trained for 150 epochs with a 250-batch size and a fixed learning rate of 0.001, 0.0001 and 0.00001 with the training process, respectively. Note that VGG-v is an extended version of the VGGNet.

Table 3 shows the classification comparison results between the proposed methods and their corresponding baselines on AlexNet, VGG-v and DenseNet, respectively. Except AReLU in AlexNet, the proposed methods can surpass the respective corresponding baseline functions. Table 4² illustrates the comparison results between AReLU and other

²Indicates that the activation function does not converge in this method.

Table 3: The classification precision of various fixed activation functions for different models on CIFAR100.

Methods	AlexNet	VGG-v	DenseNet
Sigmoid	0.4662±2.49e-3	0.5545±5.89e-4	0.2561±1.23e-3
ASigmoid(ours)	0.5312±8.66e-4	0.6578±7.56e-4	0.6103±3.15e-3
Tanh	0.5058±6.06e-4	0.6007±6.90e-4	0.5960±2.53e-3
ATanh(ours)	0.5236±1.98e-4	0.6166±2.41e-4	0.6734±3.98e-4
ReLU	0.5701±1.18e-4	0.6972±4.99e-4	0.5616±1.09e-3
LReLU[Maas et al., 2013]	0.5500±1.09e-3	0.6991±5.53e-4	0.6914±5.12e-4
AReLU(ours)	0.5647±1.71e-4	0.7005±1.04e-3	0.7081±1.64e-3

Table 4: The classification precision of various adaptive activation functions for different models on CIFAR100.

Methods	AlexNet	VGG-v	DenseNet
PReLU[He et al., 2015]	0.5325±5.01e-4	0.6838±3.84e-5	0.5781±8.87e-3
Swish[Ramachandran et al., 2017]	0.5519±1.16e-3	0.6729±8.32e-5	0.7079±3.01e-3
PELU[Trottier et al., 2017]	— ¹	—	—
FReLU[Qiu et al., 2018]	—	0.1534±2.17e-4	0.1325±2.16e-4
AReLU(ours)	0.5647±1.71e-4	0.7005±1.04e-3	0.7081±1.64e-3

adaptive functions. From the results, AReLU can achieve the best precision performance on the three network models. Similarly, PELU and FReLU cannot converge to desired loss values.

5.4 Validity and practicability in various optimization strategies

Gradient descent algorithms are often used as a black-box optimizer in neural networks, and different optimization strategies have great influence on the performance of activation functions in practice. Therefore, to further verify the validity and practicability in various optimization strategies, the best method AReLU is selected as activation function based on GoogleNet and ResNet, and a series of comparison experiments are achieved on CIFAR10 among various optimizers, such as SGD, Momentum, AdaGrad, AdaDelta and ADAM.

Fig. 7 shows the convergence curves by using various activation functions with various optimization strategies on GoogleNet and ResNet, respectively. The results show that AReLU converges faster than all other activation functions on the GoogleNet model. Whereas on the ResNet model, it is also obvious for AReLU to have an overall convergence advantage, especially compared with AdaGrad and AdaDelta. These results indicate that the proposed AReLU can accelerate convergence, thereby reducing the training cost.

Table 5 further reveals that the proposed AReLU can achieve better overall performance than other activation functions based on different optimization strategies and different network models. Except ReLU with a Momentum optimizer on GoogleNet and Swish with an ADAM optimizer on ResNet, the proposed AReLU surpasses other activation functions with various optimization strategies on both the network models, and the obtained precision performance is far better than other methods. While AReLU is only slightly worse than ReLU with a Momentum optimizer on GoogleNet and Swish with an ADAM optimizer on ResNet, respectively. Significantly, AReLU with SGD can generally achieve the best precision performance among these activation functions with various optimization strategies on both models. Fig. 8 illustrates the convergence of the proposed AReLU with various optimizers, and the results show that SGD has faster convergence speed than all other optimizers on both models, especially on ResNet.

The above results shows the proposed adaptive activation functions have faster convergence speed and higher precision than traditional activation functions. And it suggests the proposed methodology can avoid local minimums and accelerate convergence, thereby increasing the precision, reducing the training cost and improving the generalization performance.

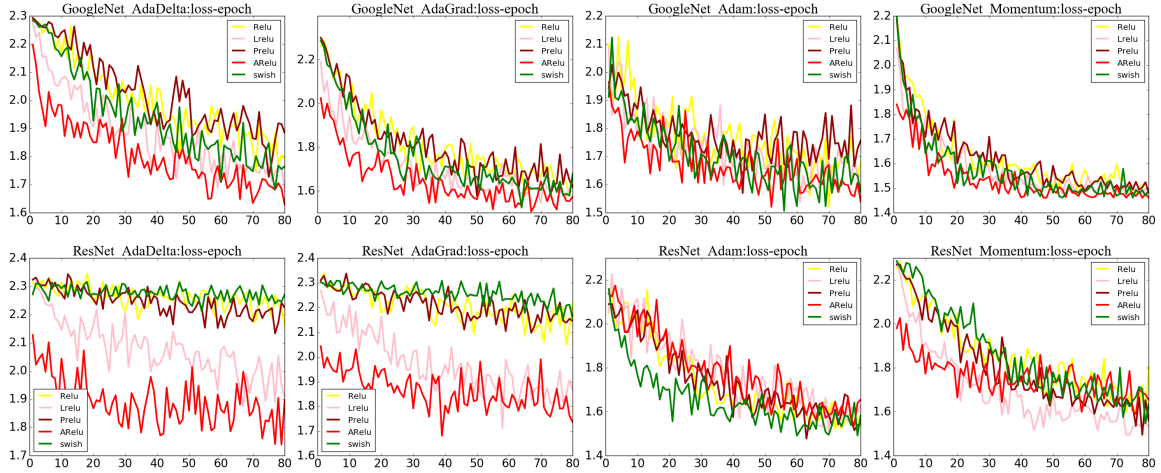


Figure 7: The convergence curves by using various optimization strategies on GoogleNet (top row) and ResNet (bottom row).

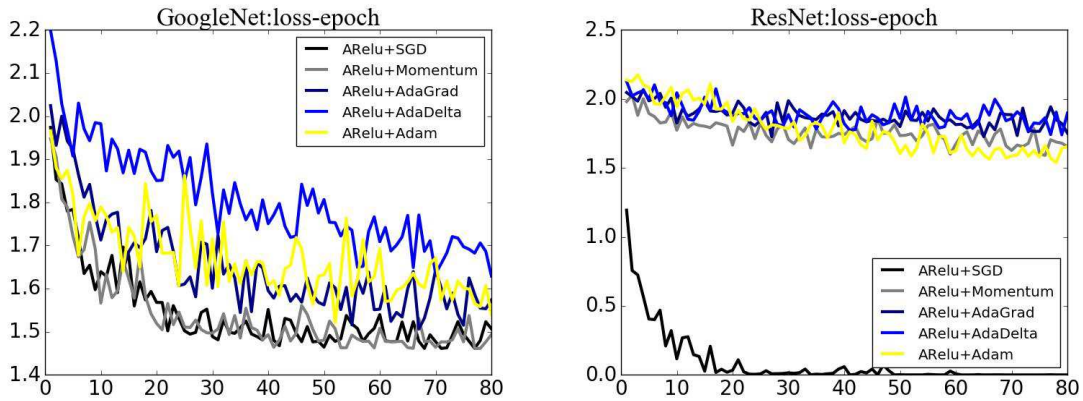


Figure 8: The convergence curves of AReLU by using various optimization strategies on GoogleNet and ResNet.

Table 5: Classification precision comparisons between various activation functions by using different optimization strategies and models on CIFAR10. The top results are highlighted in black bold and the second-best results in blue.

Models	Methods	SGD	Momentum	AdaGrad	AdaDelta	ADAM
GoogleNet	ReLU	0.8733±4.36e-5	0.8951±7.30e-5	0.6065±1.53e-4	0.5686±3.31e-4	0.8797±1.08e-3
	LReLU[Maas et al., 2013]	0.8742±9.31e-4	0.8498±4.20e-4	0.6623±3.22e-3	0.6521±1.92e-3	0.8654±8.78e-5
	PReLU[He et al., 2015]	0.8551±1.79e-4	0.8138±2.14e-4	0.6278±8.44e-6	0.5367±1.45e-3	0.8177±2.34e-6
	Swish[Ramachandran et al., 2017]	0.8710±3.25e-3	0.8772±1.39e-3	0.6747±1.58e-3	0.6139±2.66e-4	0.8743±1.92e-3
	AReLU(ours)	0.8773±1.83e-3	0.8651±8.44e-5	0.7408±6.38e-3	0.7221±2.57e-3	0.8910±3.37e-4
ResNet	ReLU	0.6641±2.59e-3	0.6274±1.99e-3	0.4253±7.38e-4	0.3553±1.57e-3	0.8519±6.20e-4
	LReLU[Maas et al., 2013]	0.6698±1.27e-3	0.6635±1.40e-3	0.5162±1.61e-4	0.4046±3.87e-3	0.8065±1.07e-3
	PReLU[He et al., 2015]	0.6522±3.20e-4	0.6325±2.11e-4	0.3493±2.07e-3	0.2846±4.28e-3	0.7889±8.08e-4
	Swish[Ramachandran et al., 2017]	0.6446±3.39e-4	0.6099±7.02e-5	0.4032±3.41e-3	0.3112±3.31e-3	0.8909±6.84e-5
	AReLU(ours)	0.9230±5.51e-4	0.7862±1.38e-3	0.7341±2.92e-3	0.7155±5.56e-4	0.8805±2.31e-4

Table 6: The classification precision of various activation functions on miniImageNet and PASCAL VOC based on ResNet50. The top results are highlighted in black bold and the second-best results in blue.

Methods	miniImageNet	PASCAL VOC
ReLU	0.7562±4.60e-3	0.5933±5.70e-3
PReLU[He et al., 2015]	0.7813±1.20e-3	0.6066±2.00e-4
Swish[Ramachandran et al., 2017]	0.7134±2.50e-3	0.5154±2.30e-3
PELU[Trottier et al., 2017]	— ¹	—
FReLU[Qiu et al., 2018]	—	—
AReLU(ours)	0.7751±1.40e-3	0.6092±3.02e-3

Table 7: The detection precision of various activation functions for different methods on various datasets. The top results are highlighted in black bold and the second-best results in blue.

Methods	Faster-RCNN (PASCAL VOC)			YOLOv2 (PASCAL VOC)			FCOS (COCO)		
	AP50	AP75	mAP	AP50	AP75	mAP	AP50	AP75	mAP
ReLU	0.8501	0.6266	0.7661	0.5924	0.1459	0.3970	0.4974	0.3426	0.3232
PReLU[He et al., 2015]	— ¹	—	—	0.1777	0.0176	0.0931	0.4943	0.3399	0.3214
Swish[Ramachandran et al., 2017]	0.7111	0.4981	0.6658	0.6250	0.2432	0.4621	0.4923	0.3418	0.3220
PELU[Trottier et al., 2017]	—	—	—	—	—	—	—	—	—
FReLU[Qiu et al., 2018]	0.8486	0.6273	0.7669	0.2604	0.0627	0.1670	0.5030	0.3473	0.3274
AReLU(ours)	0.8530	0.6271	0.7675	0.6384	0.2459	0.4687	0.5009	0.3490	0.3279

5.5 More complicated datasets

Two more complicated datasets miniImageNet[Vinyals et al., 2016] and PASCAL VOC[Everingham et al., 2012] are used to further test the validity of the proposed methodology based on ResNet50. Table 6³ shows the performance comparisons of classification precision between AReLU and other adaptive functions. The results indicate that AReLU and PReLU can obtain the best classification precision in PASCAL VOC and miniImageNet, respectively. Overall, their performances are nearly equal in the two datasets. Note that PELU and FReLU cannot converge, meaning that the two datasets are very complicated and challenging.

The above all experiments are conducted for various models, datasets and methods based on classification tasks, and it means that the proposed adaptive activation functions can overall obtain the best classification performance. To test the validity and practicability in other deep learning tasks, PASCAL VOC is used for the object detection tasks by respectively adopting Faster RCNN[Ren et al., 2015] and YOLOv2[Redmon and Farhadi, 2016] based on the proposed AReLU. Besides, a more complicated detection dataset COCO[Lin et al., 2014] is selected to further verify the effectiveness by employing FCOS[Tian et al., 2021]. Table 7⁴ shows performance comparisons of detection precision among various adaptive functions. From the results, the proposed AReLU can nearly achieve the best detection performance including AP50, AP75 and mAP among various adaptive functions based on various methods and datasets. It means that our method has a validity and practicability.

From the results of a series of comparison experiments, the proposed methods can achieve better performance in various scenarios, such as datasets, network models, optimization methods and deep learning tasks. The most significant reason is that our methodology has an internal and external bilinear structure, and the internal function is embedded within the external, such case will enable the optimization toward a global optimum solution more efficiently from all directions, thereby accelerating the convergence and improving the performance. More importantly, the proposed methodology only adds a small number of parameters (i.e., four parameters for each layer), and the number of parameters is negligible compared with millions of parameters in entire network model, which means the amount of network computing and the risk of over-fitting is only increased inconsiderably.

³Indicates that the activation function does not converge in this method.

⁴Activation function does not converge in this model.

6 Conclusions

In this work, a novel methodology is proposed to adaptively customize activation functions for various layers, and it will contribute to avoiding local minimums and accelerating convergence, thereby increasing the precision, reducing the training cost and improving the generalization performance. In this methodology, a small number of parameters are introduced to the traditional activation functions such as Sigmoid, Tanh and ReLU, and some theoretical and experimental analysis for accelerating the convergence and improving the performance is presented. To verify the effectiveness of the proposed methodology, a series of experiments are implemented on CIFAR10, CIFAR100, mini-ImageNet, PASCAL VOC and COCO by employing various network models such as VGGNet, GoogleNet, ResNet, DenseNet, various optimization strategies such as SGD, Momentum, AdaGrad, AdaDelta and ADAM, and various task like classification and detection tasks. The results show that the proposed methodology is very simple but with significant performance in convergence speed, precision and generalization, and it can surpass other popular methods like ReLU and swish in almost all experiments in terms of overall performance.

7 Acknowledgments

The authors would like to express their appreciation to the referees for their helpful comments and suggestions. This work was supported in part by Zhejiang Provincial Natural Science Foundation of China (Grant No. LGF20H180002 and GF22F037921), and in part by National Natural Science Foundation of China (Grant No. 61802347, 61801428 and 61972354), and the National Key Research and Development Program of China (Grant No. 2018YFB1305202).

References

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- Yunlan Tan, Pengjie Tang, Yimin Zhou, Wenlang Luo, Yongping Kang, and Guangyao Li. Photograph aesthetical evaluation and classification with deep convolutional neural networks. *Neurocomputing*, 228:165–175, 2017.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- Xiaoheng Jiang, Yanwei Pang, Xuelong Li, and Jing Pan. Speed up deep neural network based pedestrian detection by sharing features across multi-scale models. *Neurocomputing*, 185:163–170, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- Guimin Lin, Qingxiang Wu, Lida Qiu, and Xixian Huang. Image super-resolution using a dilated convolutional neural network. *Neurocomputing*, 275:1219–1230, 2018.
- Zhimin Tang, Linkai Luo, Hong Peng, and Shaohui Li. A joint residual network with paired relus activation for image super-resolution. *Neurocomputing*, 273:37–46, 2018.
- Jun Yu, Yong Rui, and Dacheng Tao. Click prediction for web image reranking using multimodal sparse coding. *IEEE Transactions on Image Processing*, 23(5):2019–2032, 2014.
- Jun Yu, Xiaokang Yang, Fei Gao, and Dacheng Tao. Deep multimodal distance metric learning using click constraints for image ranking. *IEEE transactions on cybernetics*, 47(12):4014–4024, 2016.
- Guoxing Wu, Wenjie Lu, Guangwei Gao, Chunxia Zhao, and Jiayin Liu. Regional deep learning model for visual tracking. *Neurocomputing*, 175:310–323, 2016.
- V. Nair and G.E. Hinton. Rectified linear units improve restricted boltzmann machines. In *in: ICML*, pages 807–814, 2010.
- Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3, 2013.

- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. *arXiv preprint arXiv:1302.4389*, 2013.
- Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In *Advances in neural information processing systems*, pages 971–980, 2017.
- Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.
- Suo Qiu, Xiangmin Xu, and Bolun Cai. Frelu: flexible rectified linear units for improving convolutional neural networks. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 1223–1228. IEEE, 2018.
- Ludovic Trottier, Philippe Giguere, and Brahim Chaib-Draa. Parametric exponential linear unit for deep convolutional neural networks. In *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 207–214. IEEE, 2017.
- Harold Cramer. *Mathematical methods of statistics*, princeton univ. Press, Princeton, NJ, 1946.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *In Advances in Neural Information Processing Systems Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge 2012 (voc2012) results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>, 2012.
- Tsung-Yi Lin, Michael Maire, and Serge J. et al. Belongie. Microsoft coco: Common objects in context. *CoRR*, abs/1405.0312, 2014.
- Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster r-cnn:towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015 (NIPS2015), December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99, 2015.
- Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. *arXiv preprint arXiv:1612.08242*, 2016.
- Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos:a simple and strong anchor-free object detector. 2021.