

Overview

Fractional Adaptation of Activation Functions In Neural Networks



Authors:

Julio Zamora Esquivel, Jesus Adan Cruz Vargas, Paulo Lopez-Meyer, Hector Cordourier Maruri, Jose Rodrigo Camacho Perez, Omesh Tickoo

Hyper parameters of a NN model





Activation Function



Name 🗢	Plot ♦	Equation +	Derivative (with respect to x) •	Inverse square root unit (ISRU) ^[14]		$f(x)=rac{x}{\sqrt{1+lpha x^2}}$	$f'(x) = \left(rac{1}{\sqrt{1+lpha x^2}} ight)^3$
Identity	_/	f(x) = x	f'(x) = 1	Inverse square root linear unit (ISRLU) ^[14]		$f(x) = egin{cases} rac{x}{\sqrt{1+lpha x^2}} & ext{for } x < 0 \ x & ext{for } x \geq 0 \end{cases}$	$f'(x) = egin{cases} \left(rac{1}{\sqrt{1+ax^2}} ight)^3 & ext{for } x < 0 \ 1 & ext{for } x \geq 0 \end{cases}$
Binary step		$f(x) = egin{cases} 0 & ext{for } x < 0 \ 1 & ext{for } x \geq 0 \end{cases}$	$f'(x)=egin{cases} 0 & ext{for }x eq 0\ ? & ext{for }x=0 \end{cases}$	Square		$\int_{f(x)}^{1} \frac{1}{x - \frac{x^2}{4}} : 0 \le x \le 2.0$	$f'(x) = 1 \pm \frac{x}{2}$
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}} [1]$	f'(x) = f(x)(1 - f(x))	(SQNL) ^[11]		$ \begin{pmatrix} f(x) & - \\ x + \frac{x^2}{4} & : -2.0 \le x < 0 \\ -1 & : x < -2.0 \end{pmatrix} $	$f(x) = 1 + \frac{1}{2}$
TanH		$f(x)= anh(x)=rac{\left(e^{x}-e^{-x} ight)}{\left(e^{x}+e^{-x} ight)}$	$f'(x) = 1 - f(x)^2$	Rectified linear unit (ReLU) ^[15]		$f(x) = egin{cases} 0 & ext{for } x < 0 \ x & ext{for } x \geq 0 \end{cases}$	$f'(x)=egin{cases} 0 & ext{for } x<0\ 1 & ext{for } x\geq 0 \end{cases}$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	Bipolar rectified linear unit		$f(x_i) = egin{cases} ReLU(x_i) & ext{if } i ext{ mod } 2 = 0 \ -ReLU(-x_i) & ext{if } i ext{ mod } 2 eq 0 \end{cases}$	$f'(x_i) = egin{cases} ReLU'(x_i) & ext{if } i ext{ mod } 2=0 \ -ReLU'(-x_i) & ext{if } i ext{ mod } 2 eq 0 \end{cases}$
ArSinH		$f(x) = \sinh^{-1}(x) = \ln\left(x + \sqrt{x^2 + 1}\right)$	$f'(x) = \frac{1}{\sqrt{x^2 + 1}}$	(BRELO)			
ElliotSig ^{[9][10][11]} Softsign ^{[12][13]}		$f(x) = rac{x}{1+ x }$	$f'(x) = rac{1}{(1+ x)^2}$	Leaky rectified linear unit (Leaky ReLU) ^[17]		$f(x) = egin{cases} 0.01x & ext{for } x < 0 \ x & ext{for } x \geq 0 \end{cases}$	$f'(x) = egin{cases} 0.01 & ext{for } x < 0 \ 1 & ext{for } x \geq 0 \end{cases}$
Exponential linear unit (ELU) ^[20]		$f(lpha,x) = egin{cases} lpha(e^x-1) & ext{for } x \leq 0 \ x & ext{for } x > 0 \end{cases}$	$f'(lpha,x) = egin{cases} f(lpha,x)+lpha & ext{for } x\leq 0 \ 1 & ext{for } x>0 \end{cases}$	Parameteric rectified linear unit (PReLU) ^[18]		$f(lpha,x) = egin{cases} lpha x & ext{for } x < 0 \ x & ext{for } x \geq 0 \end{cases}$	$f'(lpha,x) = egin{cases} lpha & ext{for } x < 0 \ 1 & ext{for } x \geq 0 \end{cases}$
Scaled exponential linear unit (SELU) ^[21]		$\begin{split} f(\alpha,x) &= \lambda \begin{cases} \alpha(e^x-1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \\ \text{with } \lambda &= 1.0507 \text{ and } \alpha = 1.67326 \end{split}$	$f'(lpha,x)=\lambdaiggl\{egin{array}{cc} lpha(e^x) & ext{for } x<0\ 1 & ext{for } x\geq 0 \end{array} ight.$	Randomized leaky rectified linear unit (RReLU) ^[19]		$f(lpha,x) = egin{cases} lpha x & ext{for } x < 0_{[3]} \ x & ext{for } x \geq 0 \end{cases}$	$f'(lpha,x) = egin{cases} lpha & ext{for } x < 0 \ 1 & ext{for } x \geq 0 \end{cases}$
S-shaped rectified linear activation unit (SReLU) ^[22]		$f_{t_l,a_l,t_r,a_r}(x) = \begin{cases} t_l + a_l(x - t_l) & \text{for } x \leq t_l \\ x & \text{for } t_l < x < t_r \\ t_r + a_r(x - t_r) & \text{for } x \geq t_r \end{cases}$	$f_{t_l,a_l,t_r,a_r}'(x) = egin{cases} a_l & ext{for } x \leq t_l \ 1 & ext{for } t_l < x < t_r \ a_r & ext{for } x \geq t_r \end{cases}$	SoftExponential ^[28]		$f(lpha,x) = egin{cases} -rac{\ln(1-lpha(x+lpha))}{lpha} & ext{for } lpha < 0 \ x & ext{for } lpha = 0 \ rac{e^{lpha x}-1}{lpha} + lpha & ext{for } lpha > 0 \end{cases}$	$f'(lpha,x) = egin{cases} rac{1}{1-lpha(lpha+x)} & ext{for } lpha < 0 \ e^{lpha x} & ext{for } lpha \geq 0 \end{cases}$
Adaptive piecewise linear (APL) ^[23]		$f(x) = \max(0, x) + \sum_{s=1}^{S} a_i^s \max(0, -x + b_i^s)$	$f'(x) = H(x) - \sum_{s=1}^S a_i^s H(-x+b_i^s)^{[4]}$	Soft Clipping ^[29]		$f(lpha,x)=rac{1}{lpha}\lograc{1+e^{lpha x}}{1+e^{lpha(x-1)}}$	$f'(lpha,x)=rac{1}{2}\sinh\Bigl(rac{p}{2}\Bigr)\mathrm{sech}\Bigl(rac{px}{2}\Bigr)\mathrm{sech}\Bigl(rac{p}{2}(1-x)\Bigr)$
SoftPlus ^[24]	\square	$f(x) = \ln(1+e^x)$	$f'(x) = \frac{1}{1+e^{-x}}$	Sinusoid ^[30]	$\overline{\Lambda}$	$f(x) = \sin(x)$	$f'(x) = \cos(x)$
Bent identity		$f(x) = \frac{\sqrt{x^2 + 1} - 1}{2} + x$	$f'(x)=\frac{x}{2\sqrt{x^2+1}}+1$	Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0\\ \frac{\sin(x)}{2} & \text{for } x \neq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x = 0\\ \frac{\cos(x)}{x} - \frac{\sin(x)}{2} & \text{for } x \neq 0 \end{cases}$
Sigmoid Linear Unit (SiLU) ^[25] (AKA SIL ^[26] and Swish-1 ^[27])		$f(x) = x \cdot \sigma(x)^{[5]}$	$f'(x) = f(x) + \sigma(x)(1-f(x))^{[6]}$	Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$

Groups of activation function







Group2

International Conference on Pattern Recognition

ICPR2020



Groups of Activation function





Grouping the activation function



Fractional calculus



$$y' = \frac{\mathrm{d}y}{\mathrm{d}x}, y'' = \frac{\mathrm{d}^2 y}{\mathrm{d}x}, y''' = \frac{\mathrm{d}^3 y}{\mathrm{d}x}$$

Is it possible to generate the 1.5 Derivative?



Fractional derivative



Gamma Function



Plot of Gamma



Fractional Derivative



Generalized ReLU



$$f(x) = \begin{cases} x & x > 0\\ 0 & x \le 0 \end{cases}$$
$$g(x) = \begin{cases} D^a x & x > 0\\ 0 & x \le 0 \end{cases}$$
$$g(x) = \frac{x^{1-a}}{\Gamma(2-a)}$$

0

0



Generalized Sigmoid

Softplus $g(x) = ln(1 + e^x)$ $g(x) = D^a ln(1 + e^x)$ $g(x) = \lim_{h \to 0} \frac{1}{h^a} \sum_{n=0}^{\infty} (-1)^n \frac{\Gamma(a+1)\ln(1+e^{(x-nh)})}{\Gamma(n+1)\Gamma(1-n+a)}$





Generalized Hyperbolic Tangent





Generalized swish



Adaptive activation function



The neurons of the NN were modified to adjust not only their weights, but also their activation function. The parameter **a** represents the derivative order of the activation function, producing a generalization, it means an activation function that can morph from softplus to sigmoid to bell-like shape, for the best result.





Learning ReLU

$$g(x) = \frac{x^{1-a}}{\Gamma(2-a)}$$

$$\frac{\partial}{\partial a}g(x) = -\left[\frac{\Gamma'(2-a)}{\Gamma(2-a)^2} + \frac{\ln(x)}{\Gamma(2-a)}\right]x^{1-a}$$

Using Digamma function

$$\frac{\partial}{\partial a}g(x) = -\left[\psi(2-a) + \ln(x)\right]g(x)$$

Adaptive activation function



The results using this technology in comparison with academic state-of-the-art neural nets shows performance improvements for image recognition datasets, with no increase in parameter size.

- The optimal solution is not a combination of the known activation functions: It is a fractional value
- Using AAF in 18 layers gives 95.08% accuracy : 18(AAF)Layers = 1000 Layers in CIFAR10
- Two Percent Improvement with similar number of parameters

CIFAR10 Results

Table 1. Comparing ResNet18* (with adaptive activation function) Vs reported ResNet topologies

Neural Network	Depth	#Parameters	Accuracy%
ResNet18	18	11M	93.02
ResNet50	50	25.6M	93.62
ResNet100	100	44.5M	93.75
ResNet18*	18	11M	95.08





Activation function Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ $\ln(1+e^x)$ $D^a \ln(e^{-bx} + e^x)$ Hyperbolic Tangent $y = \tanh(x)$ $\ln(e^{-x} + e^x)$





International Conference on Pattern Recognition ICPR2020

18

#