

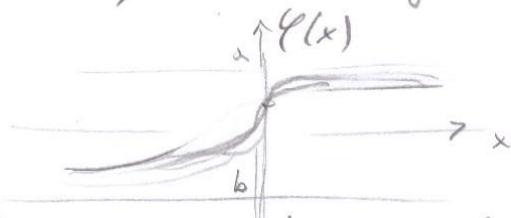
APPENDIX

ACTIVATION FUNCTIONS

In neural network architectures there are various non linear one dimensional functions. Their ~~work~~ role inside the networks, that use mostly linear connections, causes to maintain non linear overall characteristics that help to implement approximations to arbitrary multi-dimensional functions. Some of them are used for classifiers and some for gating operations. Their ~~use~~ ^{use} in the architecture is discussed in the various chapters of this book. Hence we limit ourselves to list many of them pointing to their gradient and main properties.

SIGMOIDAL ACTIVATION FUNCTIONS

Many activation functions used in neural network architectures, have a "sigmoidal" shape.



This behavior continues for $x \rightarrow -\infty$ and for $x \rightarrow \infty$ maintaining a non decreasing behavior $\forall x$.
Their threshold is usually set at $x=0$, but it can be shifted $y(x-t)$.
By controlling the scale of x , we can obtain
"quasi step functions" for large, or
"quasi linear" for x small.

More specifically any sigmoidal function^{AP2} can be scaled as

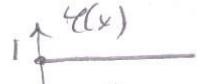
$$f(x) = \beta \sigma(\alpha x) \quad \text{where usually } \beta, \alpha > 0$$

- for $x \rightarrow \infty$ we have $f(x) \rightarrow \beta$ for $x > 0$ and $f(x) \rightarrow 0$ for $x < 0$
and for $x \rightarrow -\infty$ a linear behaviour around $x=0$.

Some peculiarities are in their derivatives and properties. That will be detailed in the following. Note that all cumulative distribution functions (cdf) have a sigmoidal behaviour. In fact most of them listed in the following are directly derived from typical random variable cdfs and pdfs.

STEP FUNCTION

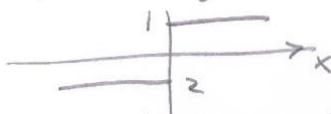
The step function is the first sigmoidal function to ever occur.



$$u(x) = u(x)$$

Also its shifted version SIGN FUNCTION

$$c(x) = 2u(x) - 1 = \text{sgn}(x)$$



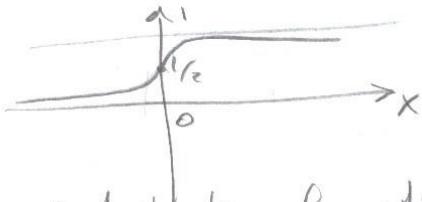
has the same behaviour with zero derivative everywhere except $\frac{\partial}{\partial x} x=0$ where they do not exist. In neural network architecture it is preferable to use smooth versions of these functions for more effective backpropagation of the error.

APPENDIX I

THE LOGISTIC FUNCTION

By far the most common sigmoid activation function is the logistic

$$\ell(x) = \ell(x) = \frac{1}{1+e^{-x}}$$



It is a valid CDF (cumulative distribution function)
For small values of x is almost linear $\approx x + \frac{1}{2}$
and has a number of interesting properties:

(1) THE DERIVATIVE:

$$\boxed{\ell'(x) = \ell(x)(1-\ell(x))}$$

Proof:

$$\ell'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{e^{-x} - 1 + 1}{(1+e^{-x})^2} = \frac{(1+e^{-x}) - 1}{(1+e^{-x})^2}$$

$$= \frac{1}{1+e^{-x}} - \frac{1}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}}\right) = \ell(x)(1-\ell(x))$$

(2) CONNECTIONS TO THE HYPERBOLIC TANGENT (discussed in the following)

$$\boxed{\ell(x) = \frac{1}{2} \left(1 + \tanh \frac{x}{2}\right)}$$

Proof

$$\frac{1}{2} \left(1 + \frac{e^{\frac{x}{2}} - e^{-\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}}\right) = \frac{1}{2} \frac{e^{\frac{x}{2}} + e^{-\frac{x}{2}} + e^{\frac{x}{2}} - e^{-\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}} = \frac{e^{\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}} = \frac{1}{1+e^{-x}}$$

(3) THE INVERSE

AFL4

Setting $p = l(x)$

$$\left[x = \ln \frac{p}{1-p} \right] \quad x \in]0,1[$$

(4) Given a set of real numbers $\alpha_1, \alpha_2, \dots, \alpha_N$

$$l(\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_N x_N) = \frac{p_1^{\alpha_1} p_2^{\alpha_2} \dots p_N^{\alpha_N}}{p_1^{\alpha_1} p_2^{\alpha_2} \dots p_N^{\alpha_N} + (1-p_1)^{\alpha_1} (1-p_2)^{\alpha_2} \dots (1-p_N)^{\alpha_N}}$$

where $p_i = l(x_i)$

Proof:

$$\text{If } p = l(x) \quad e^{-x} = \frac{1-p}{p} \quad \text{and} \quad e^{-\alpha x} = \left(\frac{1-p}{p} \right)^{\alpha}$$

Therefore

$$\frac{1}{1 + e^{-\alpha_1 x_1 - \alpha_2 x_2 - \dots - \alpha_N x_N}} = \frac{1}{1 + \left(\frac{1-p_1}{p_1} \right)^{\alpha_1} \left(\frac{1-p_2}{p_2} \right)^{\alpha_2} \dots \left(\frac{1-p_N}{p_N} \right)^{\alpha_N}}$$

which gives the result.

Interesting corollaries are the sum and the difference.

$$l(x_1 + x_2) = \frac{p_1 p_2}{p_1 p_2 + (1-p_1)(1-p_2)}$$

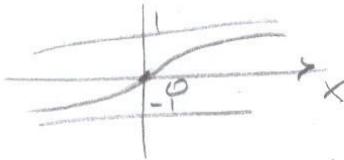
$$l(x_1 - x_2) = \frac{\frac{p_1}{p_2}}{\frac{p_1}{p_2} + \frac{1-p_1}{1-p_2}} = \frac{p_1(1-p_2)}{p_1(1-p_2) + p_2(1-p_1)}$$

THE HYPERBOLIC TANGENT

AF.5

Strictly related to the logistic function, is the hyperbolic function.

$$t(x) = \tanh x \triangleq \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



Its peculiarity is in being between -1 and 1.

For small values of x is linear: $t(x) \approx x$.

Used for large scale - $t(x)$ vs x .

We list some of the relevant properties:

① RELATION TO THE LOGISTIC

$\tanh x$ and $\ell(x)$ are really the same function scaled in the argument and shifted in value.

$$\tanh x = 2\ell(2x) - 1$$

Proof

$$2 \frac{1}{1+e^{-2x}} - 1 = \frac{2-1-e^{-2x}}{1+e^{-2x}} = \frac{1-e^{-2x}}{1+e^{-2x}} = \frac{e^{2x}(e^x - e^{-x})}{e^{2x}(e^x + e^{-x})} = \tanh x$$

② DERIVATIVE

$$\frac{d}{dx} \tanh x = (1 + \tanh x)(1 - \tanh x) = 1 - \tanh^2 x = 1 - t(x)^2$$

Proof

$$\frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

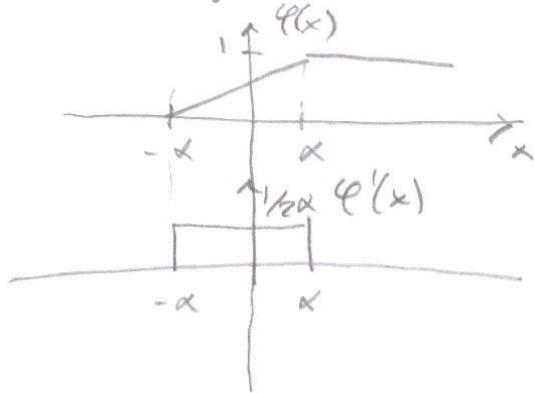
$$= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}}\right)^2 = 1 - \tanh^2 x$$

$$= (1 + \tanh x)(1 - \tanh x) = 1 - \tanh^2 x; \quad t'(x) = 1 - t(x)^2$$

PIECEWISE LINEAR

AF6

A sigmoidal function can be built using linear segments as



$$\Phi(x) = \begin{cases} 0 & x < -\alpha \\ \frac{x}{2\alpha} & -\alpha < x < \alpha \\ 1 & x > \alpha \end{cases}$$

$$\Phi'(x) = \begin{cases} \frac{1}{2\alpha} & -\alpha < x < \alpha \\ 0 & \text{else} \end{cases}$$

$\Phi(x)$ is clearly the cdf of a uniform random variable. $\Phi(x)$ is continuous but discontinuous in its derivative.

When $\alpha = 1$, it is called HARD SIGMOID and can be written as

$$\Phi(x) = \max(0, \min(1, \frac{x+1}{2}))$$

Easy to verify (and to code)

Note the the piecewise linear $\Phi(x)$ can be also written scaling the hard sigmoid

$$\Phi(x) = \max(0, \min(1, \frac{\frac{x}{\alpha} + 1}{2}))$$

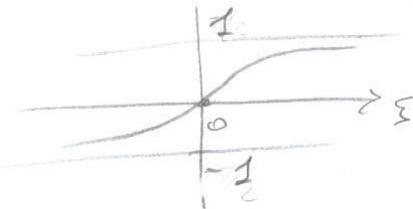
Hyperbolic Functions

AF7

TRIGONOMETRIC ARCTANGENT

The trigonometric arctangent is naturally a hyperbolic function.

$$\ell(x) = \frac{2}{\pi} \operatorname{arctan} x$$



with the well-known derivative

$$\ell'(x) = \frac{2}{\pi} \frac{1}{1+x^2}$$

SOFT SIGN

the soft sign is similar to the hyperbolic tangent, but it converges to its asymptotic values polynomially rather than exponentially

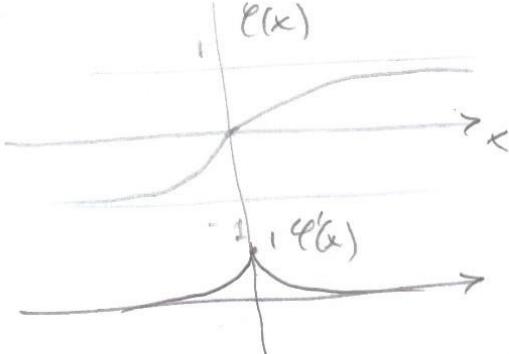
$$\ell(x) = \frac{x}{|x| + 1}$$

Its derivative is

$$\ell'(x) = \left(\frac{1}{|x| + 1} \right)^2$$

Proof

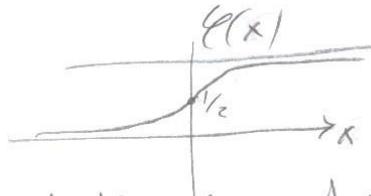
$$\begin{aligned} \ell'(x) &= \frac{|x| + 1 - x \cdot \frac{\partial |x|}{\partial x}}{(|x| + 1)^2} = \frac{|x| + 1 - x \cdot \frac{x}{|x|}}{(|x| + 1)^2} \\ &= \frac{|x| + 1 - |x|}{(|x| + 1)^2} = \frac{1}{(|x| + 1)^2} \end{aligned}$$



$$\begin{aligned} \text{since } \frac{\partial |x|}{\partial x} &= \frac{x}{|x|} \\ \text{and } \frac{x^2}{|x|} &= \frac{|x|^2}{|x|} = |x| \end{aligned}$$

GAUSSIAN CDF

$$\varphi(x) \triangleq \int_{-\infty}^x N(\xi; 0, 1) d\xi$$



AF8

The advantage here is that the derivative is the gaussian function

$$\varphi'(x) = N(x; 0, 1)$$

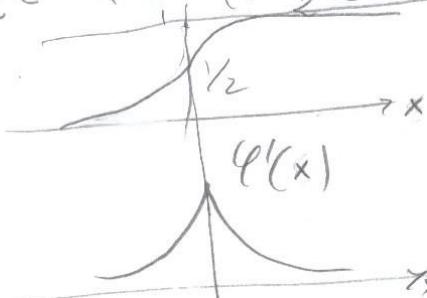


The disadvantage is in the fact that $\varphi(x)$ cannot be computed in closed form, but needs a numerical algorithm. Most computer languages have $\varphi(x)$ available using their function library. Perhaps in the form of an ~~error function~~ $\text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$,

$$\varphi(x) = \frac{1}{2} \left(1 + \text{erf} \left(\frac{x}{\sqrt{2}} \right) \right) \quad (\text{see my book see random variables})$$

LAPLACIAN CDF

$$\varphi(x) = \frac{1}{2} e^{X^2} u(-x) + \left(1 - \frac{1}{2} e^{-|x|} \right) u(x)$$



$$\varphi'(x) = \frac{1}{2} e^{-|x|}$$

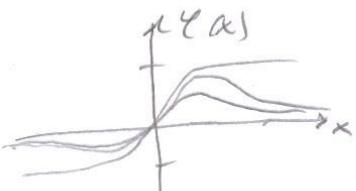
This sigmoidal function has an exponential growth in the transition region.

AFG

SOBOL'eva MODIFIED HYPERBOLIC TANGENT (SMHT)

A generalization of the hyperbolic tangent is the Soboleva modified function

$$q(x) = \frac{e^{ax} - e^{-bx}}{e^{cx} + e^{-dx}}$$



By controlling the values of a, b, c, d we can obtain various behaviours that go from the tan-like

$$a = b = c = d = 1$$

to other non-saturated behaviours.

We leave the many details here for brevity, suggesting to the interested reader to look at the specific literature.

NON SIGMOIDAL ACTIVATIONS

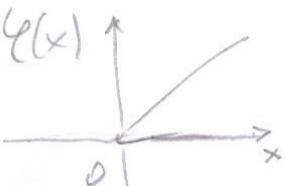
Since the original simplest neural network, the perceptron, has at the output a sigmoidal function, sigmoids have been used extensively in many architectures. However sigmoid functions may saturate and during learning may cause large plateaux in the cost functions.

Saturation in the first layers of a neural network may limit the expressiveness of some architectures as the higher layers may have to work on saturated values.

The choice of the most appropriate activation function is application-specific and we report here some other non-sigmoidal choices that have become popular in the latest architectures.

THE RECTIFYING NON LINEARITY (ReLU - Rectifying Logic Unit)

By far this is the most common choice



$$y(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} = \max(0, x)$$

It is piecewise linear and its derivative is the step function $\delta(x)$



$$\delta(x)$$

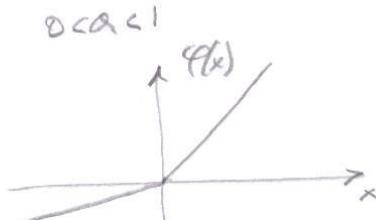
This function essentially blocks all the negative values and lets the positive one to be copied at the output.

AFII

PARAMETRIC RECTIFIED LINEAR UNIT

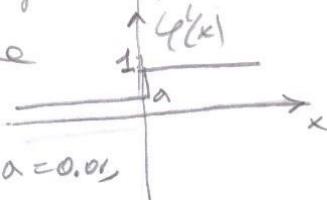
(PReLU)

$$\varphi(x) = \begin{cases} ax & x \leq 0 \\ x & x > 0 \end{cases}$$



This function is a slight modification of the ReLU, allowing a "leakage" in the negative part.

The derivative is a positive stepwise constant function.

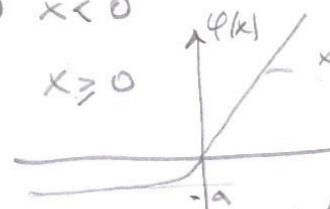


When a is small, typically $a=0.01$,

This is called LEAKY RECTIFIED LINEAR UNIT (Leaky ReLU)

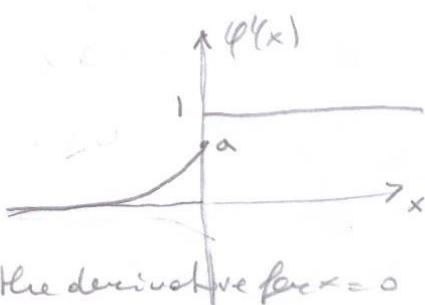
EXPONENTIAL LINEAR UNIT (ELU)

$$\varphi(x) = \begin{cases} a(e^x - 1) & x < 0 \\ x & x \geq 0 \end{cases}$$



This unit is a continuous modification of the ReLU in the negative part where it tends exponentially to $-a$. The derivative is

$$\varphi'(x) = \begin{cases} ae^x = \varphi(x) + a & x < 0 \\ 1 & x \geq 0 \end{cases}$$



There is a discontinuity in the derivative for $x=0$ unless $a=1$.

This function can be scaled to get the

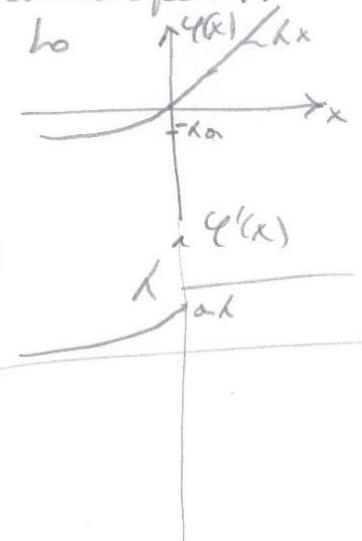
$$\varphi(x) = \begin{cases} \lambda a(e^x - 1) & x < 0 \\ \lambda x & x \geq 0 \end{cases}$$

SCALED EXPONENTIAL LINEAR UNIT (SELU)

AF12

To control the slope of the linear part,
the derivative is similarly to
above

$$Q(x) = \begin{cases} \lambda e^x - Q(x) + \lambda a & x > 0 \\ \alpha & x \leq 0 \end{cases}$$



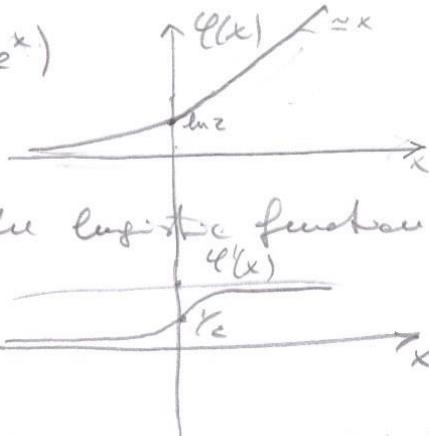
SOFTPLUS

AF13

Another smooth version of the ReLU is the softplus activation

$$\ell(x) = \ln(1 + e^x)$$

Note that for large x the function is $\approx x$.



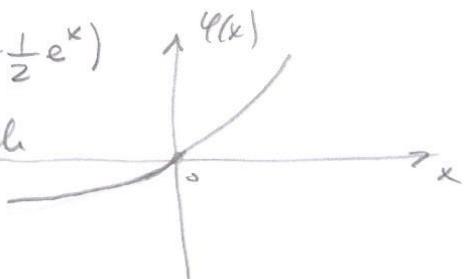
Its derivative is the logistic function

$$\ell'(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$

Sometimes this function is modified to be the SHIFTED SOFTPLUS

$$\ell(x) = \ln\left(\frac{1}{2} + \frac{1}{2}e^x\right)$$

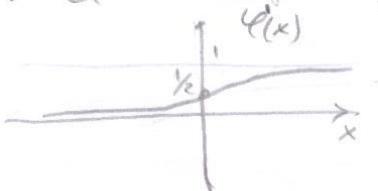
It is made to go through $x=0$



Obviously

$$\ell(x) = \ln\frac{1}{2}(1 + e^x) = \ln(1 + e^x) - \ln 2$$

and the derivative is still the logistic function.



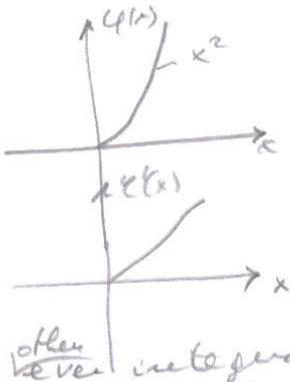
SQUARED RELU

AF14

Where the growth in the positive part needs to be longer than linear, we could use a larger order power, for example 2.

$$\varphi(x) = \begin{cases} x^2 & x \geq 0 \\ 0 & \text{else} \end{cases}$$

$$\varphi'(x) = \begin{cases} 2x & x > 0 \\ 0 & \text{else} \end{cases}$$



Clearly the exponent can be any ~~never~~ ^{other} integer. Some authors have found benefits in transformer architectures using squared ReLU.

SIGMOID LINEAR UNIT (SiLU) or SWISH-1

Any sigmoidal function can be easily modified to give a linear behaviour for large values of x and tend to zero with negative values of x .

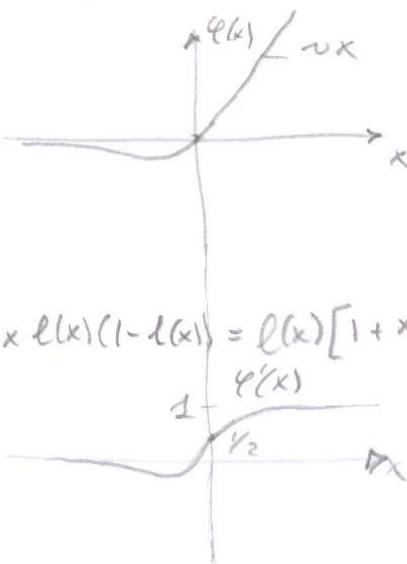
$$\varphi(x) = \frac{x}{1 + e^{-x}}$$

Lime

$$\varphi(x) = x \ell(x)$$

↑ logistic

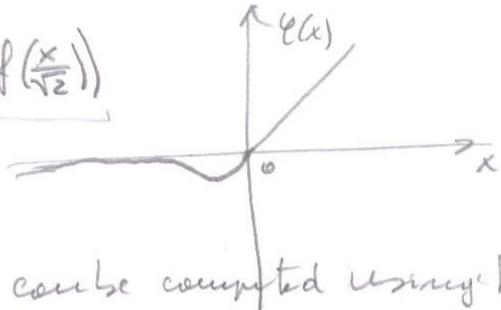
$$\varphi(x) = \ell(x) + x \ell(x) = \ell(x) + x \ell(x)(1 - \ell(x)) = \ell(x)[1 + x(1 - \ell(x))]$$



GAUSSIAN ERROR LINEAR UNIT (GELU) AF15

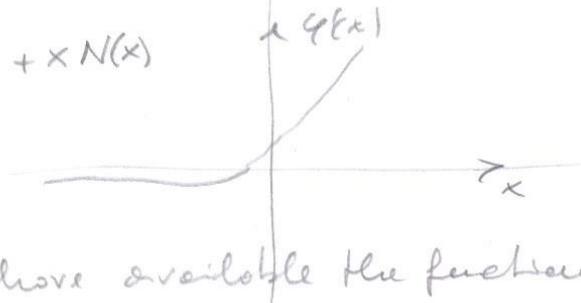
Also the gaussian cdf sigmoid can be modified to be

$$\ell(x) = x \frac{1}{2} \left(1 + \operatorname{erf} \left(\frac{x}{\sqrt{2}} \right) \right)$$



Also the derivative can be computed using the results on the gaussian cdf sigmoid

$$\ell'(x) = x \ell(x) = \ell(x) + x N(x)$$

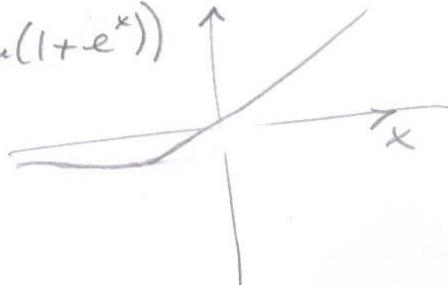


clearly one has to have available the functions $\ell(x)$ and $N(x)$.

AF16

LOG-SOFTPLUS ERROR activation FUNCTION (SERF)

$$Q(x) = x \operatorname{erf}(\ln(1+e^x))$$

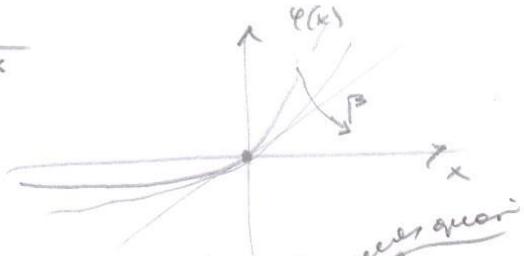


SWISH

AFIT

The swish function is a generalization of the Sigmoid linear Unit. (SLU) because it has a parameter to control its smoothness.

$$\ell(x) = \frac{x}{1 + e^{-\beta x}}$$



For large values of x the function becomes quasi-linear, but β can control the overall behavior. For large β the function is almost linear.

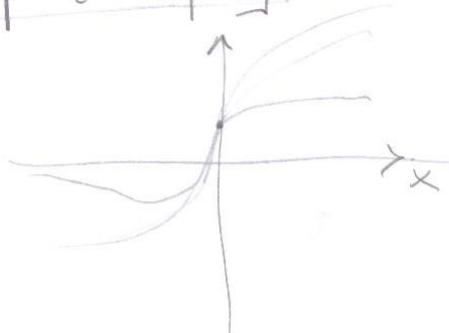
The parameter β is often learned allowing the unit to adjust itself to a linear/nonlinear behavior. Note that changing β is not the same as scaling x because the coefficient is only at the denominator.

Since $\frac{1}{1 + e^{-\beta x}} = \ell(\beta x)$ and $\ell'(\beta x) = \beta \ell(\beta x)(1 - \ell(\beta x))$

$$\ell(x) = x \ell(\beta x)$$

$$\ell'(x) = \ell(\beta x) + x \ell'(\beta x) = \ell(\beta x) + x \beta \ell(\beta x)(1 - \ell(\beta x))$$

$$\boxed{\ell'(x) = \ell(\beta x) [1 + \beta x (1 - \ell(\beta x))]}$$



AF\8

More activation functions have been proposed for various architectures. Various authors have reported advantages of using a function rather than another one. The results are unfortunately very specific to the application and to the algorithms utilized. Most of the time the best choice is determined at the end of a trial-and-error phase.

There are other parametric activation functions proposed that learn their parameters. They are based on responsibilities of basis functions or properly-chosen parameterizations. In these cases the idea is that the most appropriate shape of non-linearity has to emerge from learning.