# Activation Functions



step

sigmoid

$$\frac{1}{1+e^{-x}} \quad \frac{e^x}{e^x+1}$$

leaky ReLU

$$\max(0, x)$$

$$\max(0.1x, x)$$

tanh

$$\frac{2}{1+e^{-2x}} - 1$$

$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$

linear

ReLU

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$f(x) = \max(0, x)$$

$$Z$$

$$S$$

Softmax

$$0 \to 1$$

$$\Sigma \Rightarrow 1$$

$$z$$

$$y$$

$$y_i = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$w^{new} = w^{old} - \eta \nabla_w l(w)$$

$$= w^{old} - \frac{1}{N} \eta \sum \nabla_w l_i(w)$$

$$z_1 = f(w_{11}x_1 + w_{21}x_2 + w_{31}x_3 + b)$$

$$z_j = f\left(\sum_i w_{ij}x_i + b_j\right)$$

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

$$W^1 = \begin{pmatrix} w_{11} & w_{21} & w_{31} & \cdots & w_{d1} \\ w_{12} & w_{23} & w_{32} & & \\ & & \vdots & & \\ w_{13} & w_{23} & w_{33} & \cdots & w_{dj} \\ w_{1h'} & \cdots & \cdots & & w_{dh'} \end{pmatrix}$$

$W^1$  $h' \times d$ matrix

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

$$\sqrt{\sum_i w_{ij} x_i + b}$$

$$\hat{y} = f^n\left(\cdots f^3\left(W^3 f^2\left(W^2 f^1\left(W^1 x + b^1\right) + b^2\right) + b^3\right)\cdots\right)$$

$h^3 \times 1$

$h' \times d$   $d \times 1$   $h' \times 1$

$\underbrace{h' \times d}_{h' \times 1}$

$\underbrace{h^2 \times h'}_{h^2 \times 1}$   $h' \times 1$   $h^2 \times 1$

$$\text{Step} \left( \omega_1 x_1 + \omega_2 x_1 + \omega_3 x_3 + b \right)$$

Sigmoid

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

$$= \frac{e^x}{e^x + 1}$$

$$2\sigma(x) - 1$$

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

Sigmoid function
Tanh function

Sigmoid

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

# ReLU

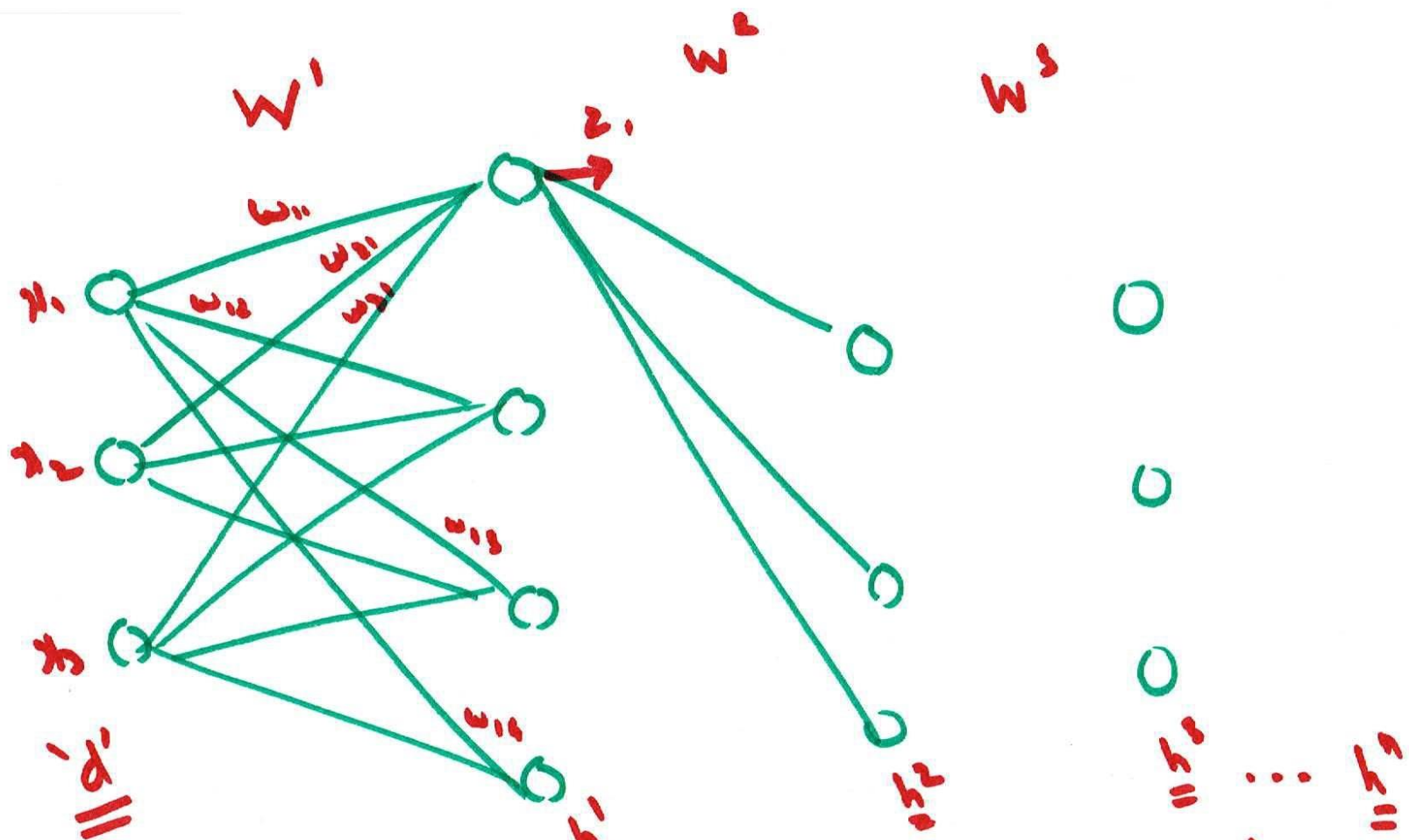$$= \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$= \max(0, x)$$

$$\max(0.01x, x)$$

linear $= x$

$$\text{Step} = \mathcal{J}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$\mathcal{J}(w_1 x_1 + w_2 x_2 + b)$$

x   y

Output nodes

hidden layer

Classifier

S: sigmoid, tanh
Softmax

Sigmoid ✓

tanh ✓

Reg → linear

ReLU ✓

~~linear~~

$\hat{a} + \hat{b}(a + bx)$

# Softmax

$$z_1$$
$$z_2$$
$$z_3$$
$$z_4$$

$$\frac{e^{z_i}}{\sum_j e^{z_j}}$$

$$\rightarrow \quad \frac{e^{z_1}}{\sum_j e^{z_j}}$$

$$\rightarrow \quad \cdot \quad \frac{e^{z_2}}{\sum_j e^{z_j}}$$

$$\rightarrow \quad \cdot$$

$$\rightarrow \quad \cdot \quad \frac{e^{z_4}}{\sum_j e^{z_j}}$$

0.1

0.2

0.4

0.3

___

1

$x$     $\textcircled{y}$     $\hat{y}$

$\begin{array}{c} \fbox{-10} \\ \fbox{+12} \\ \fbox{-5} \\ 6 \end{array}$ $\longleftrightarrow$ $\begin{array}{c} 3 \\ 30 \\ -16 \\ 0 \end{array}$ $\Big\}$

$\begin{array}{c} -10 \\ +12 \\ -5 \\ 6 \end{array}$ $\begin{array}{c} 3 \\ 30 \\ -10 \\ 0 \end{array}$

$\overset{x}{\underset{=}{x}} \longrightarrow \Big\{ \text{[network]} \Big\} \longrightarrow \textcircled{\hat{y}}$



Loss Function     $L(y, \hat{y})$

Re

$$\boxed{\frac{1}{N} \sum_i \left( y_i - \hat{y}_i \right)^2}$$

$L_2$ loss
MSE
SSE

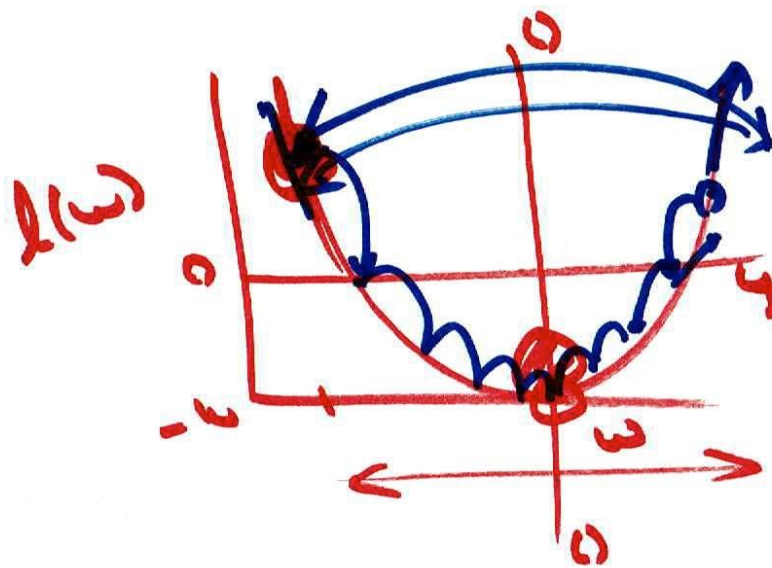$$L(y, \hat{y}) = l(w)$$

Classific $L(y, \hat{y}) = -\left(y_i \log(\hat{y}_i) + (1-y_i)(\log(1-\hat{y}_i))\right)$

Cross entropy loss

How min $\dfrac{L(y, \hat{y})}{ta}$ by ~~Every~~ Changing

$\underline{W^1, W^2, \ldots W^n}$

my

$w_{11}, w_{12} \ldots$



$l(w)$

$\boxed{y = x^2 - 10} = \boxed{-10}$

$\dfrac{dy}{dx}^2 \boxed{2x = 0}$

$x = 0$

$\boxed{\dfrac{dL}{dw}} = \boxed{\phantom{xxxx}} = 0$

$\dfrac{dl}{dw}$

$w^{new} = 3 - \ell \dfrac{\nabla_w l}{\cancel{x}}$

↑
learning rate

$$GD$$

$$w^{new} = w^{old} - \eta \nabla_w l(w)$$

$$= w^{old} - \eta \frac{1}{N} \sum \nabla_w l_i(w)$$

$$\boxed{SGD}$$

$$w^{new} = w^{old} - \eta \nabla_w l_i(w)$$

$$\boxed{w^{new} = w^{old} - \frac{1}{M} \eta \sum \nabla_w l_i(w)}$$

over a min batch

Loss

$$L = \frac{1}{N} \sum \underset{\text{Function of } \omega \ (\ell(\omega))}{\text{\#}} \left( \underline{y_i} = f^n \left( \dots f \left( w^2 f^1 (w^1 x + b^1) + b^2 \right) \right)_i \right)^2$$

Chain Rule

$$f(g(h(\dots x \dots)))$$

$$\frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dx}$$

Back Propagation