

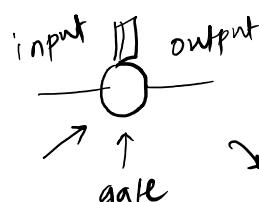
What are Activation Functions?

31 May 2022 14:49

In artificial neural networks, each neuron forms a weighted sum of its inputs and passes the resulting scalar value through a function referred to as an activation function or transfer function. If a neuron has n inputs then the output or activation of a neuron is

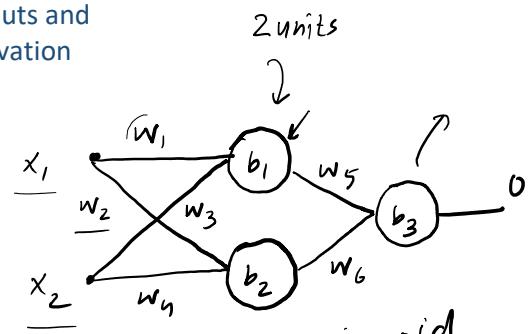
$$a = g(w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b)$$

This function g is referred to as the activation function.



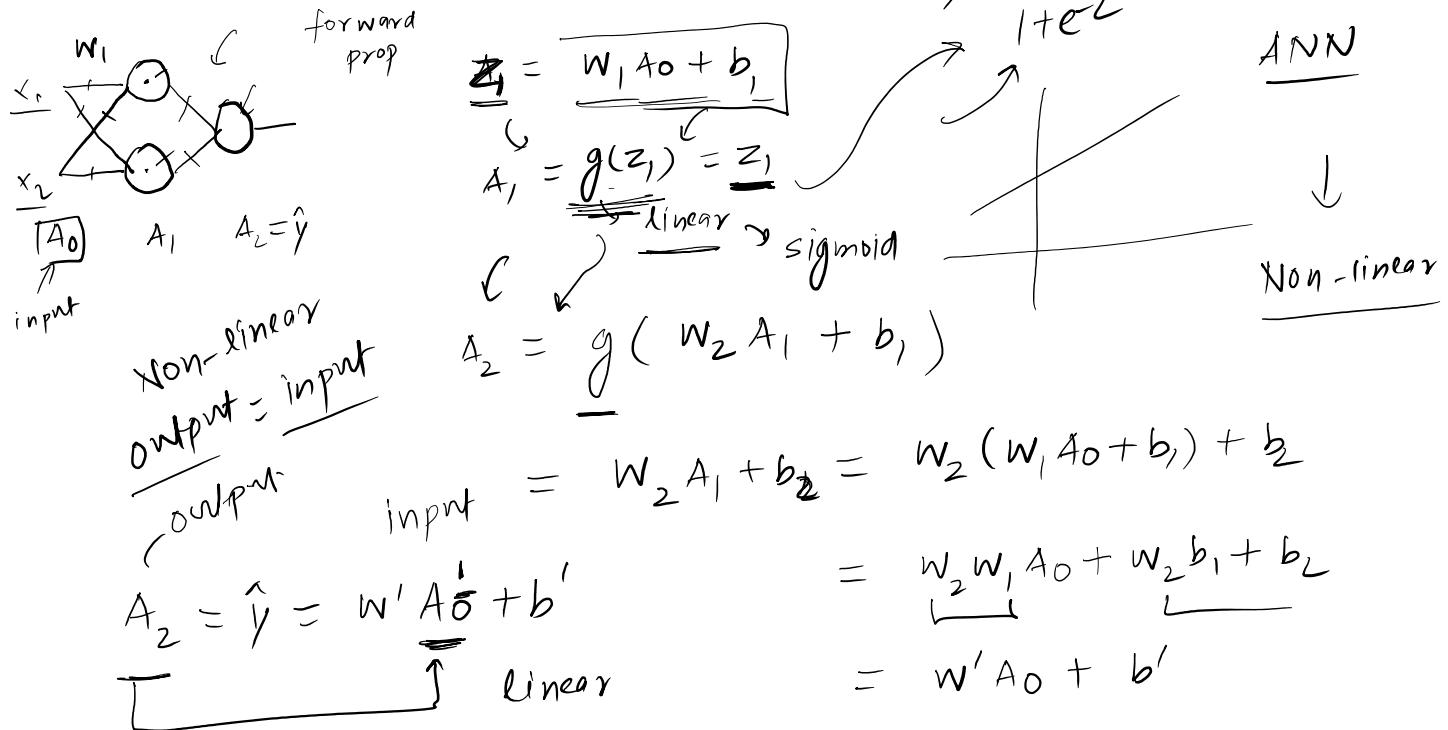
activated
(Y/N)
how much

$$\underbrace{g}_{\text{activation}}(w_1x_1 + w_2x_2 + b_1)$$



Why Activation Functions are needed?

31 May 2022 14:50



Ideal Activation Function

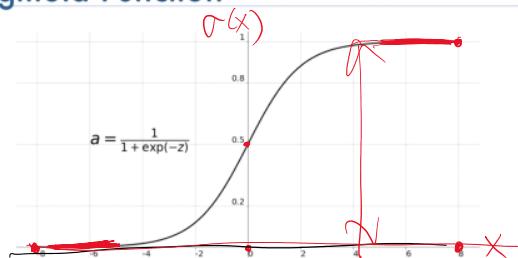
31 May 2022 14:50

- ✓ $y = f(x)$ linear non-linear $\sigma(z) = \frac{1}{1+e^{-z}}$
- 1) Non-linear ✓ Universal Approx Theorem
 - 2) Differentiable gradient des \rightarrow derivatifs \rightarrow ReLU (0)
 - 3) Computationally inexpensive derivative \rightarrow simple easy fast training slow
 - 4) Zero centered zero cent normalized mean = 0 normalized data input later layers faster \rightarrow tanh
 - 5) Non-saturating Sigmoid $\rightarrow [0, 1]$ $f(x) = \max(0, x)$
 tanh $\rightarrow (-1, 1)$ ReLU
- update x $[0, 1] \rightarrow$ backprop $\xrightarrow{\text{saturating AF}}$ vanishing gradient problem
 $0.0001 \leftarrow 0.0001 \leftarrow 0.01$

Sigmoid Activation Function

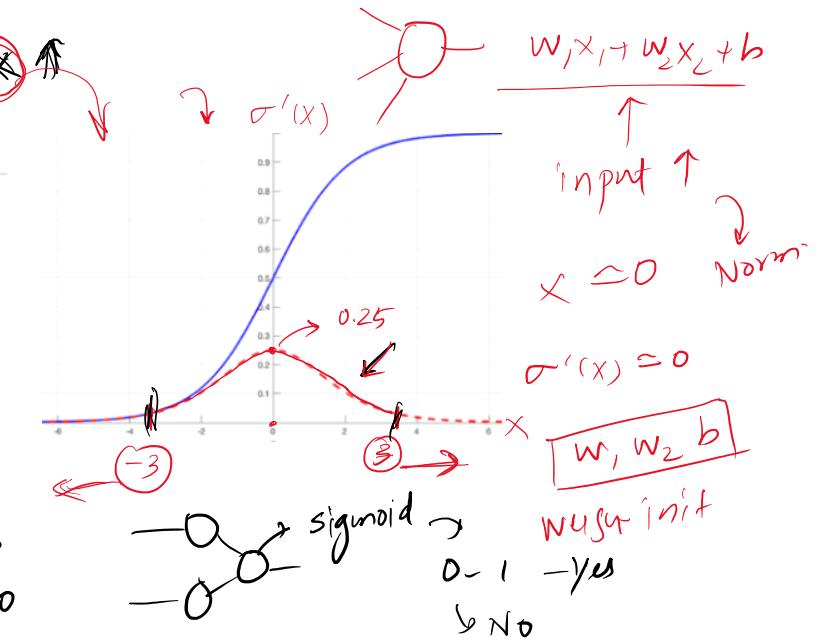
31 May 2022 14:50

Sigmoid Function



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

Yes
No

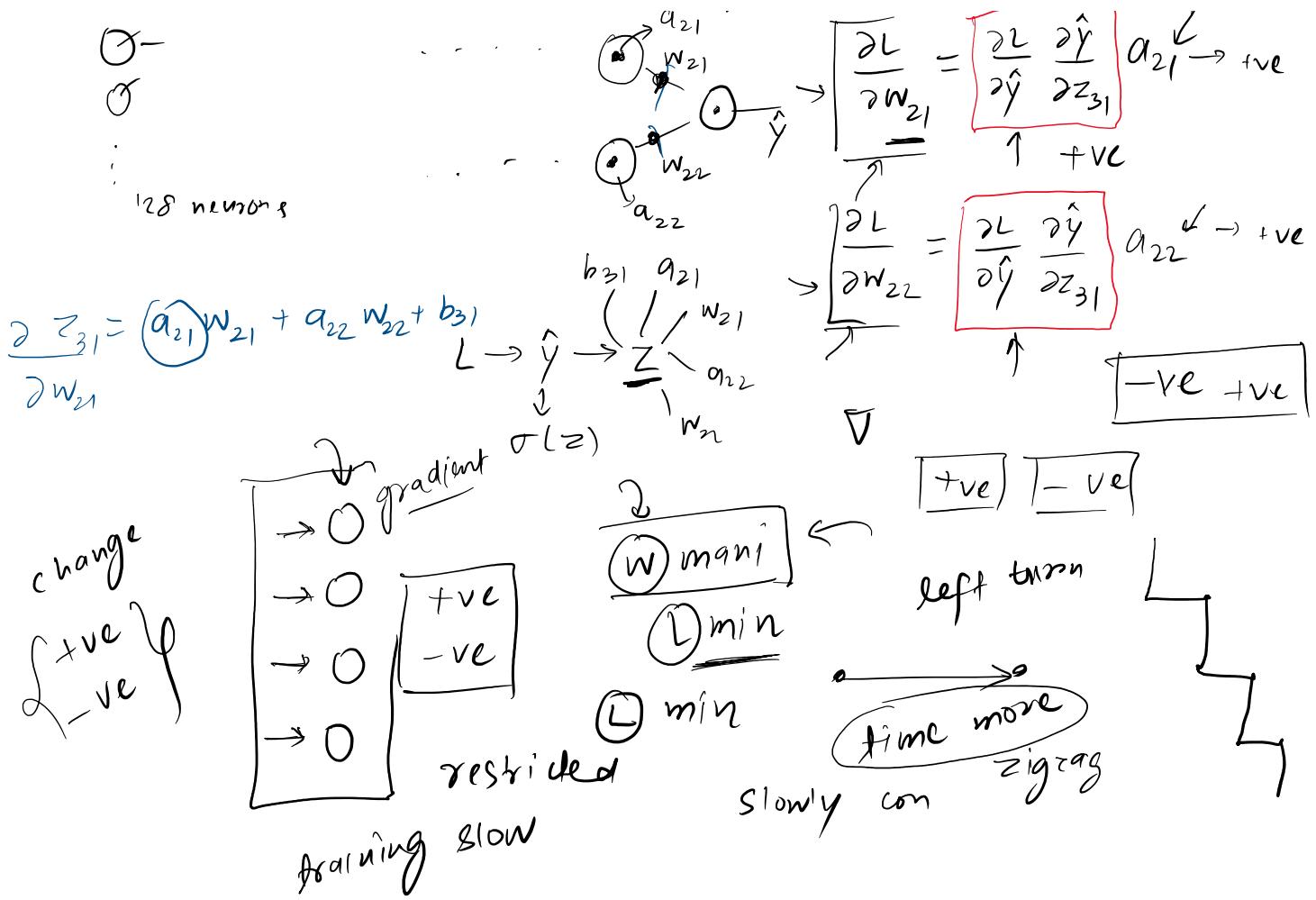


Advantages

- 1) $[0, 1] \rightarrow$ probability \rightarrow output layer \rightarrow Binary classification
- 2) Non-linear \rightarrow Non-linear data \rightarrow good option
- 3) Differentiable \rightarrow Backprop \rightarrow $\left[\frac{\partial L}{\partial w} \right]$

Disadvantages

- 1) Saturating function $\rightarrow [-\infty, \infty] \rightarrow [0, 1]$
 - ↳ Vanishing gradient problem
 - ↳ Backprop \rightarrow update $w_n = w_o - \eta \frac{\partial L}{\partial w}$
 - ↳ No update will take place
 - ↳ training X
- 2) Non Zero centered
 - ↳ Normalized \downarrow
 - ↳ mean = 0
 - ↳ a_{21} w_{21} \downarrow
 - ↳ $\frac{\partial L}{\partial w_{21}} = \left[\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_2} \right] a_{21} \downarrow$
 - ↳ +ve -ve convergence problem $[0 \rightarrow 1]$



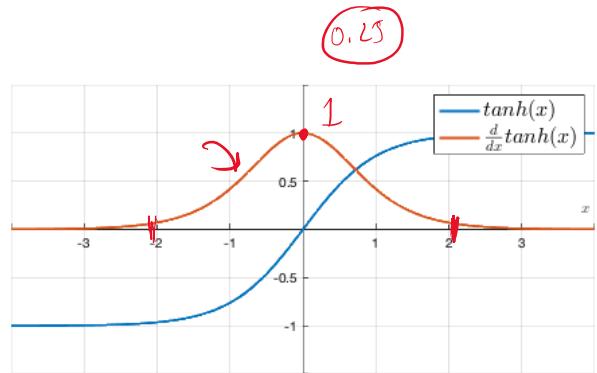
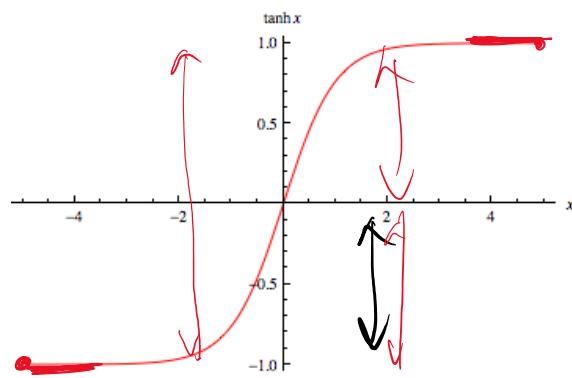
3) $\sigma(x) = \frac{1}{1+e^{-x}}$ \rightarrow computation
~~exp~~ time

hidden layer
 sigmoid
 output \rightarrow b cpm

Tanh Activation Function

31 May 2022 14:50

(0,1) (-1,1)



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

$$f'(x) = (1 - \tanh^2(x))$$

Advantages

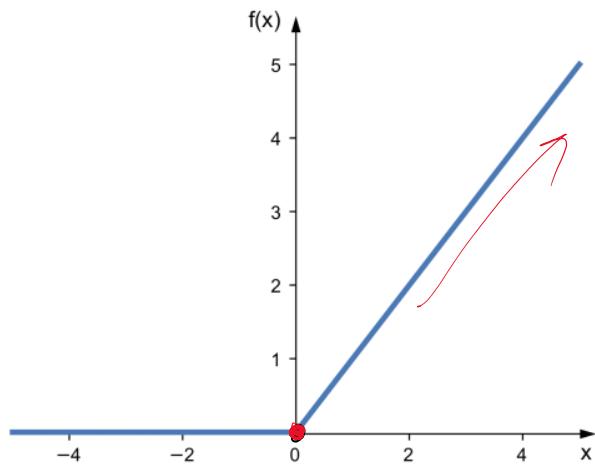
- 1) Non-linear
- 2) Differentiable
- 3) zero centered $\begin{cases} \rightarrow +ve \\ \rightarrow -ve \end{cases}$ by training faster Sigmoid

Disadvantage

- saturating function → Vanishing gradient prob
- computationally exp slow

Relu Activation Function

01 June 2022 16:43



activation

$$f(x) = \underline{\max(0, x)} \quad f(x) = x$$

Advantage

1) Non-linear

2) Not saturated in the +ve region

3) Computationally inexpensive

4) Converge → faster → sigmoid
↳ fanh

$$\boxed{f(x) = \max(0, x+1)}$$

$$- \max(0, x-1)$$

Disadvantage

$$x < 0 \rightarrow 0$$

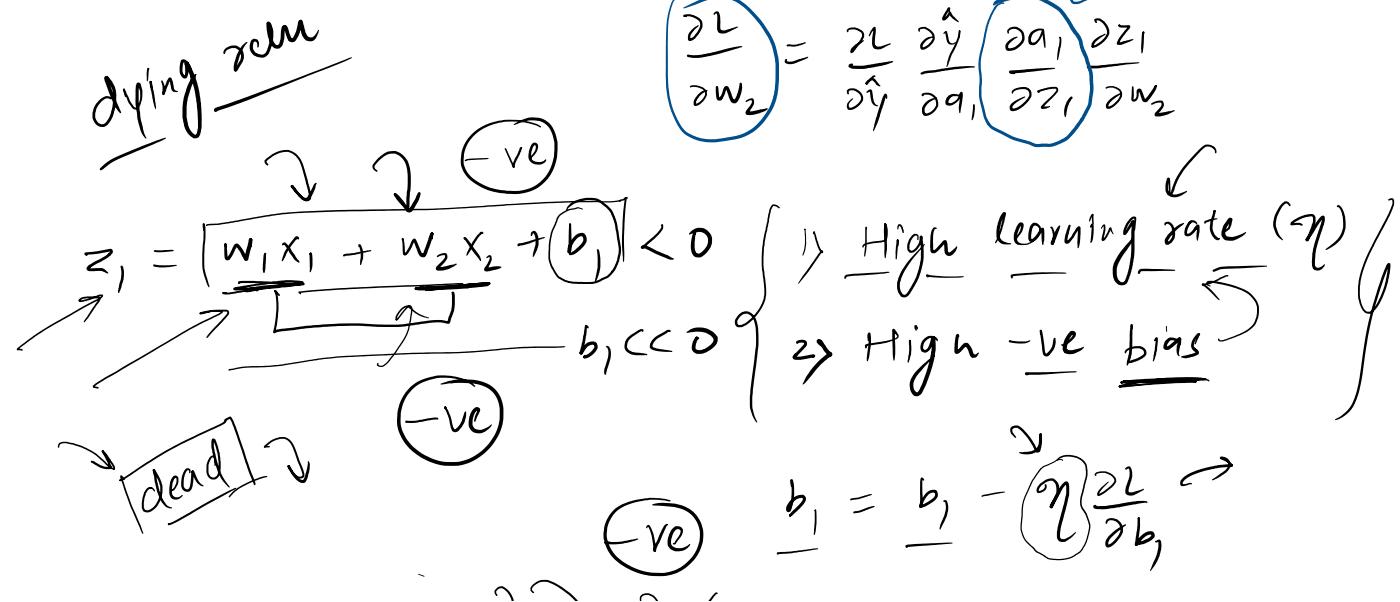
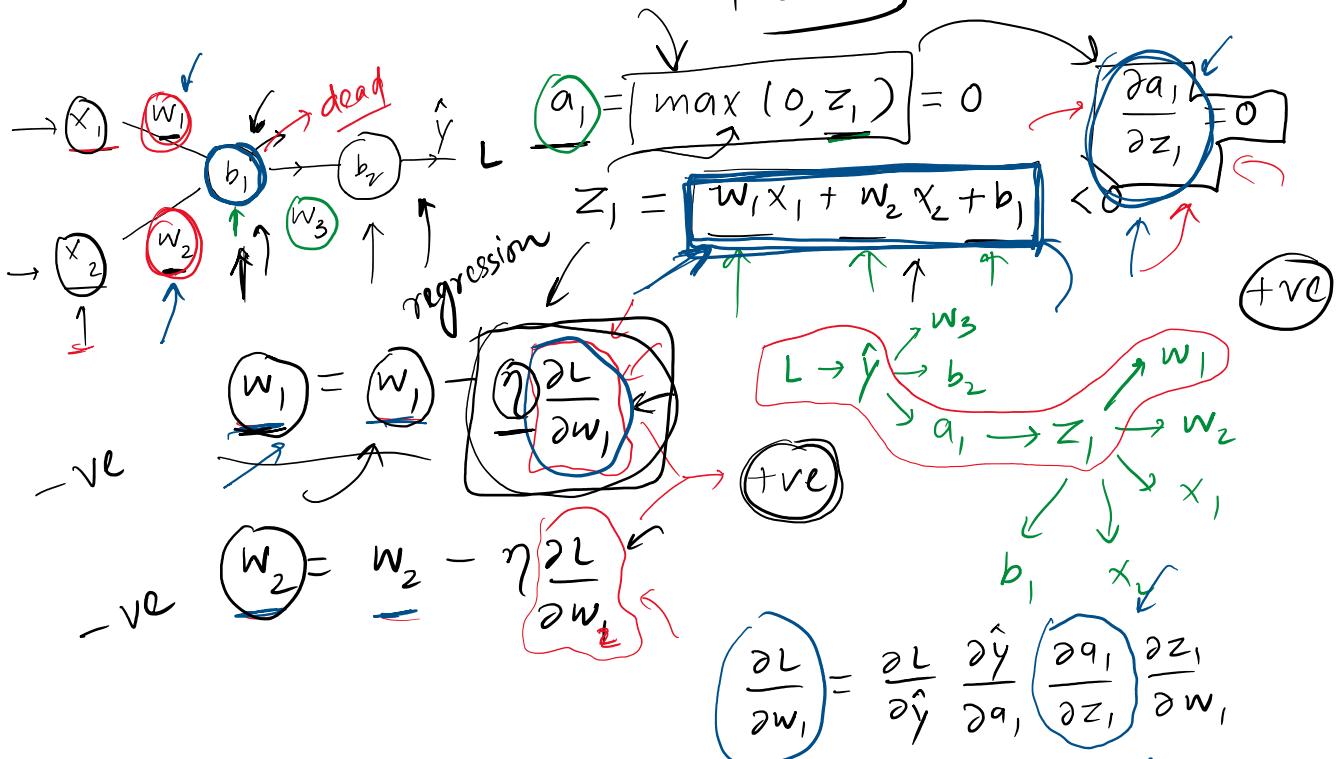
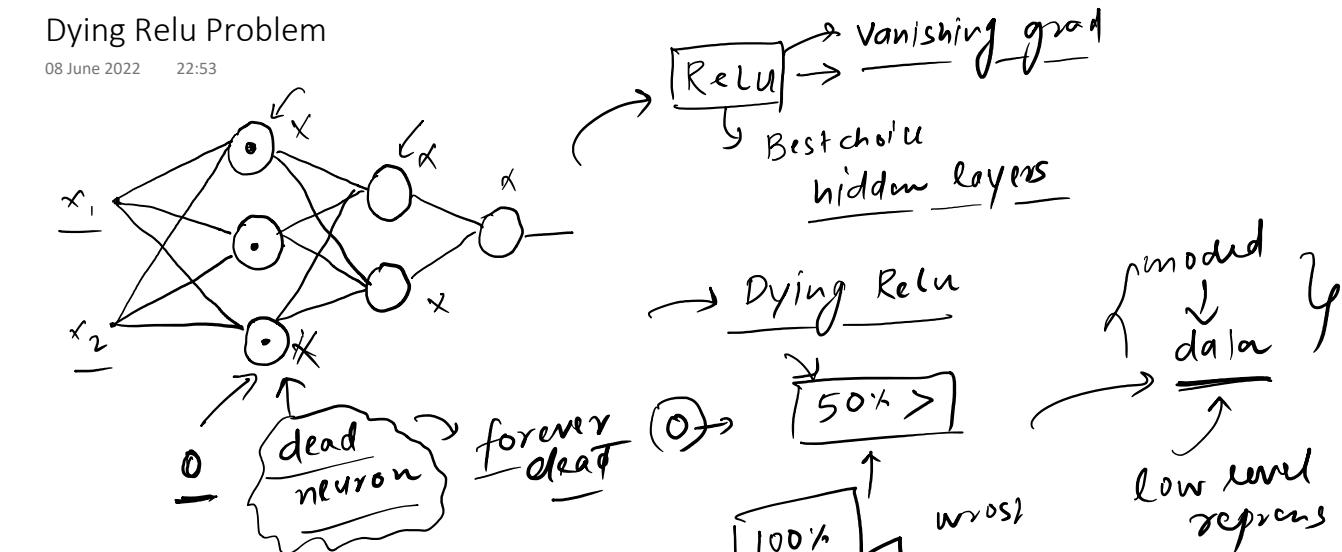
$$\rightarrow \text{Differentiability} \quad x \geq 0 \rightarrow 1$$

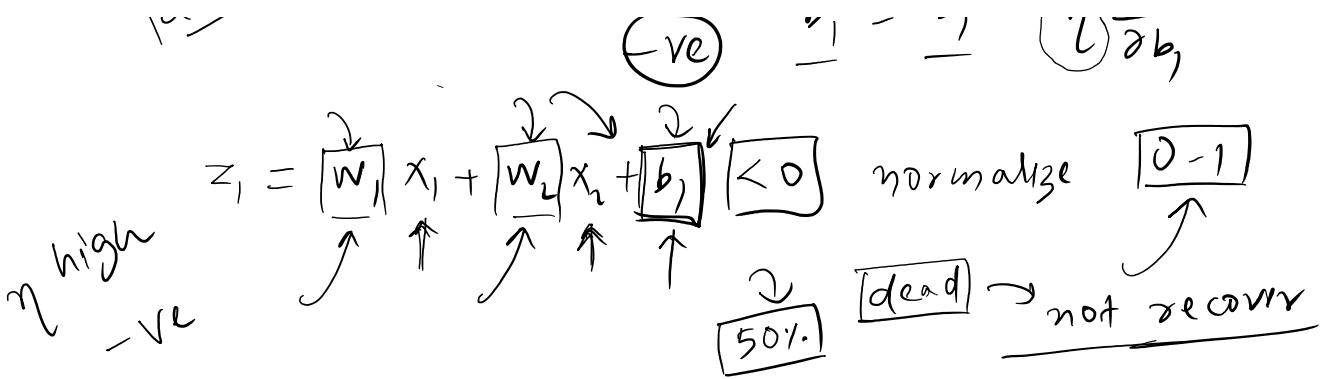
→ Non-zero centered → sigmoid → Batch Normalization
normalize →

↳ Dying ReLU problem

Dying Relu Problem

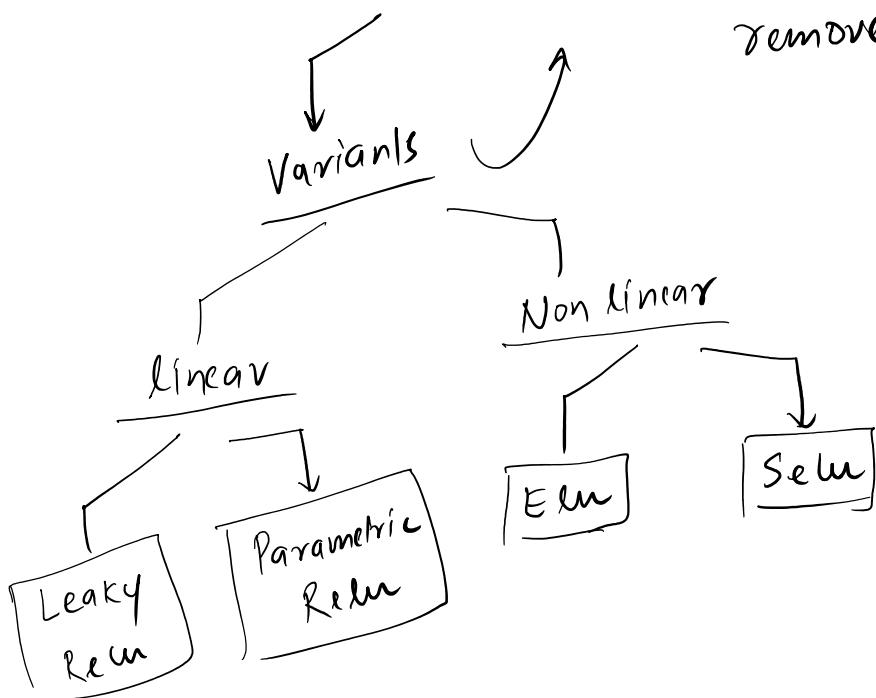
08 June 2022 22:53





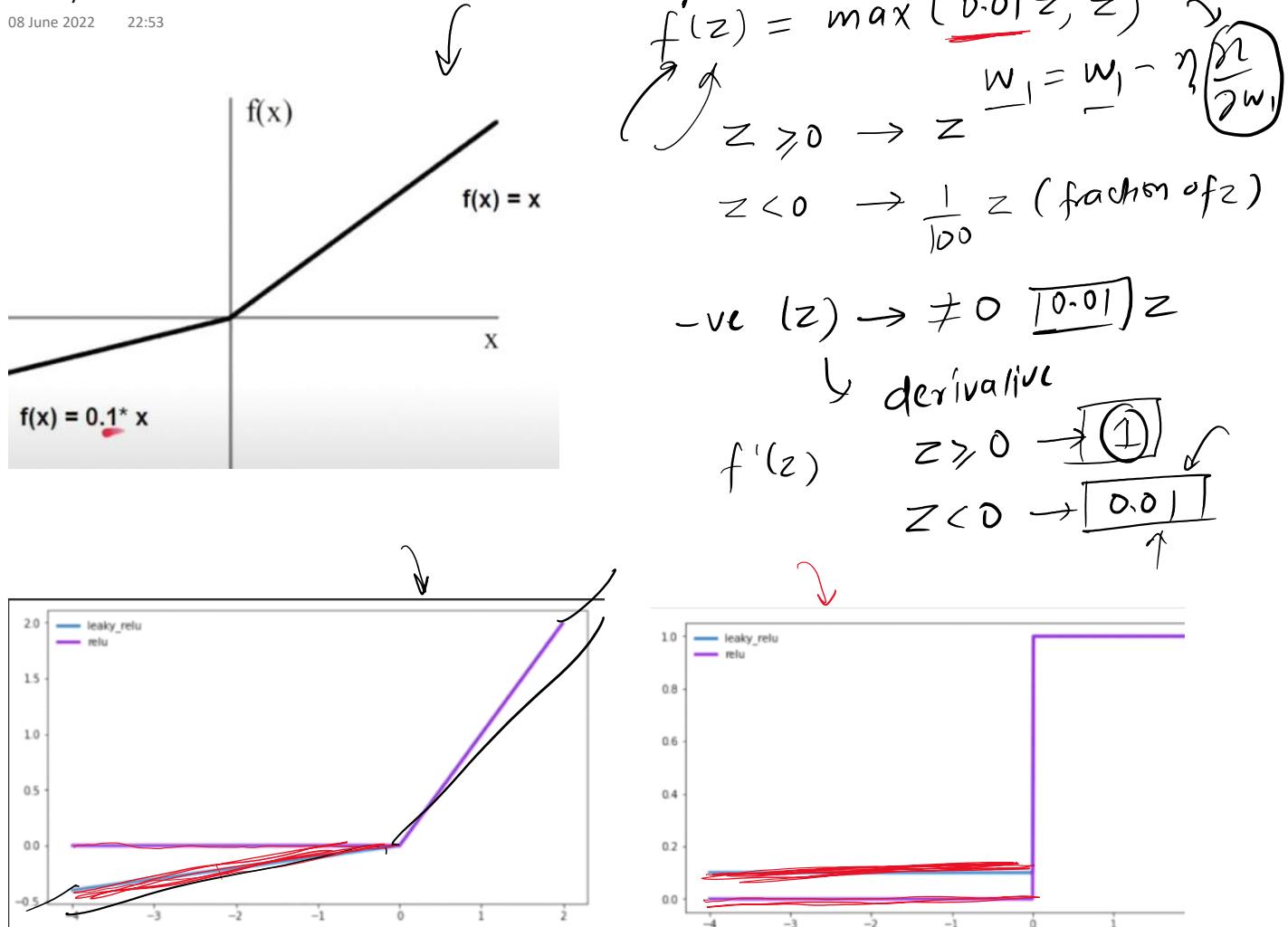
Solutions

- Set low learning rate
- bias \rightarrow +ve value \rightarrow $[0.01]$
- Don't use $\text{relu} \rightarrow$ variants (+ve keep)
remove



Leaky Relu

08 June 2022 22:53



Advantages

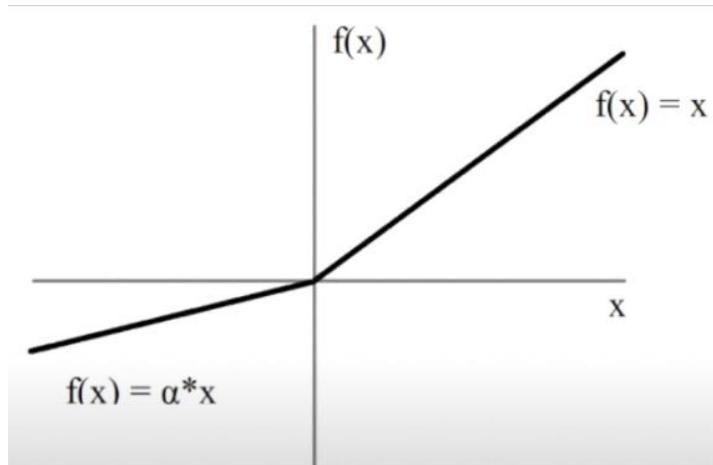
- Non-saturated \rightarrow Unbounded
- Easily computed
- No dying relu problem
- Close to 0 centered

Disadv

-ve/+ve

Parametric Relu

08 June 2022 22:53



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{otherwise} \\ 0.01x & \end{cases}$$

$\alpha \rightarrow$ trainable parameter

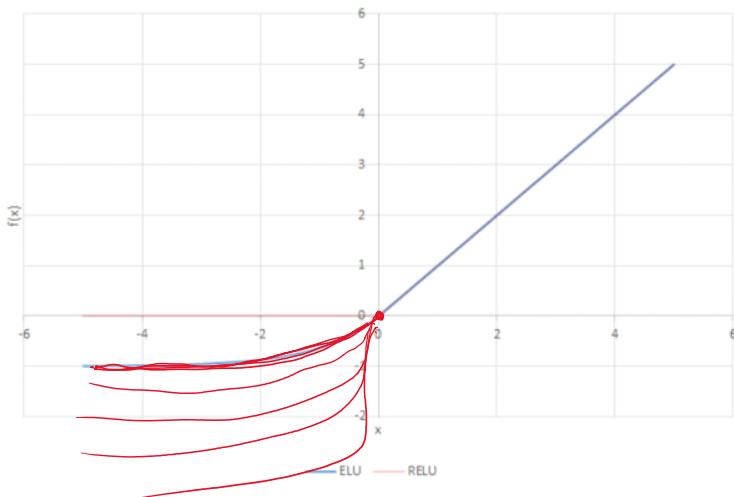
flexibility

depend on data

Elu - Exponential Linear Unit

09 June 2022 00:29

performs better than ReLU

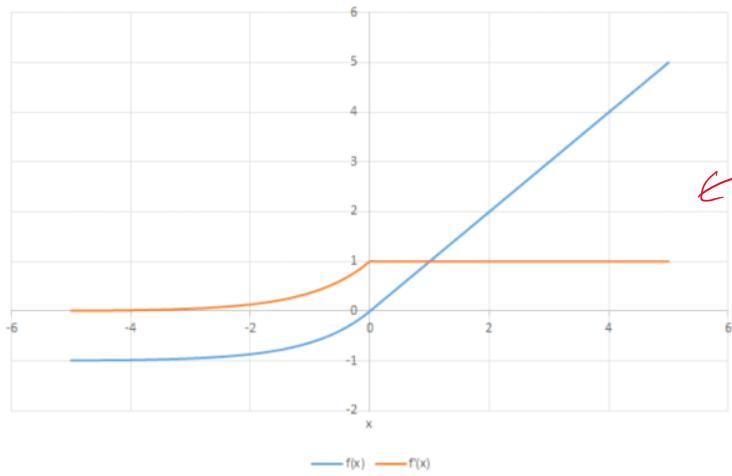


$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

0.1 to 0.3

$$\text{constant}$$

$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \text{ELU}(x) + \alpha & \text{if } x \leq 0 \end{cases}$$



Advantages

- Close to zero centered
- Convergence faster
- Better generalized
- Dying ReLU x
- Always continuous as well

Disadv
→ computation expensive (ex)

Selu - Scaled Exponential Linear Unit

09 June 2022 00:29

Recent \rightarrow new \rightarrow 90+ pages



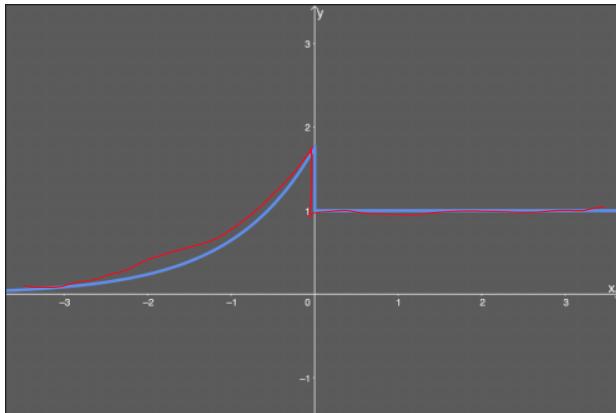
$$\text{SELU}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$

$$a \approx 1.6732632423543772848170429916717$$

$$\lambda \approx 1.0507009873554804934193349852946$$

\rightarrow fixed ✓
train bnx

$$\text{SELU}'(x) = \lambda \begin{cases} 1 & \text{if } x > 0 \\ \alpha e^x & \text{if } x \leq 0 \end{cases}$$



Self-normalizing \rightarrow activation
 \downarrow
normalized

$m = 0 \quad \sigma = 1$

converge faster