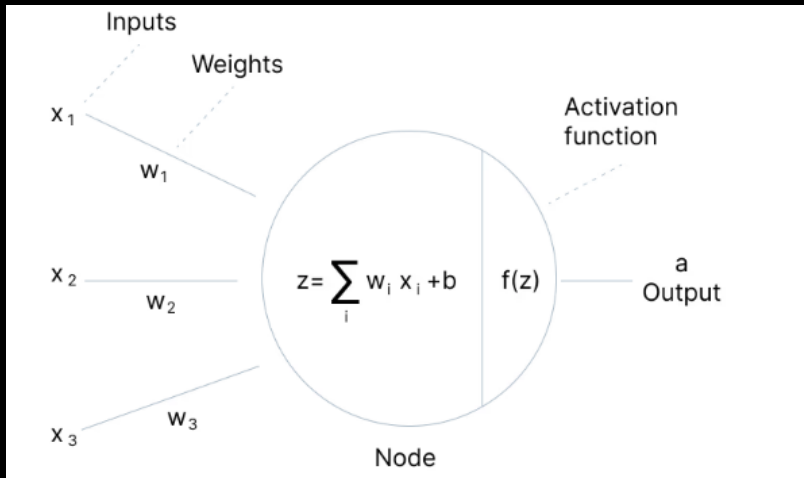


Last Class \longrightarrow Optimizers ADAM, RMSPROP

Today's Agenda:-

i) Activation Functions



(i) Step

(ii) Linear

(iii) Sigmoid

(iv) Relu

$$z = wx + b = -5$$

$$f(z) = (wx + b) \text{ step function}$$

$$= 0$$

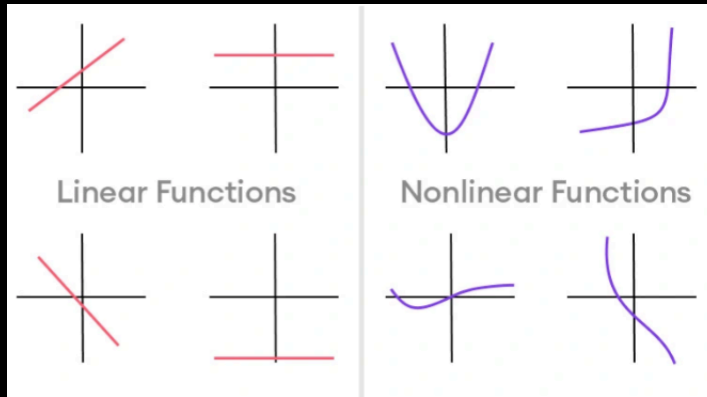
Mathematical Operations

$$f(z) = (wx + b) \text{ Linear function / Identity Function}$$

$$= (wx + b) = -5$$

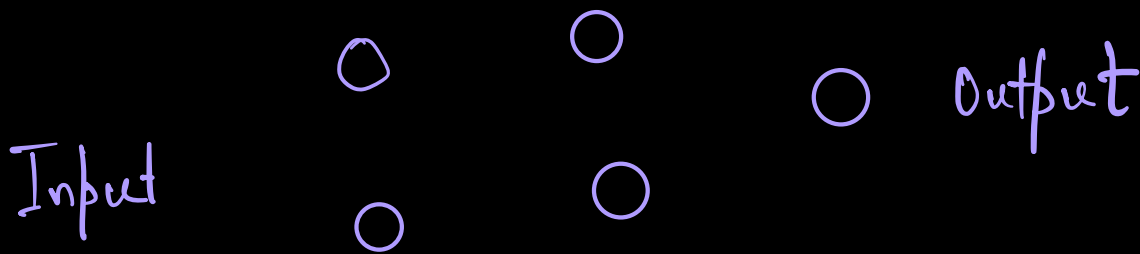
Why?

i) To non-linearity.



Activations :- Positive gradient
Negative gradient

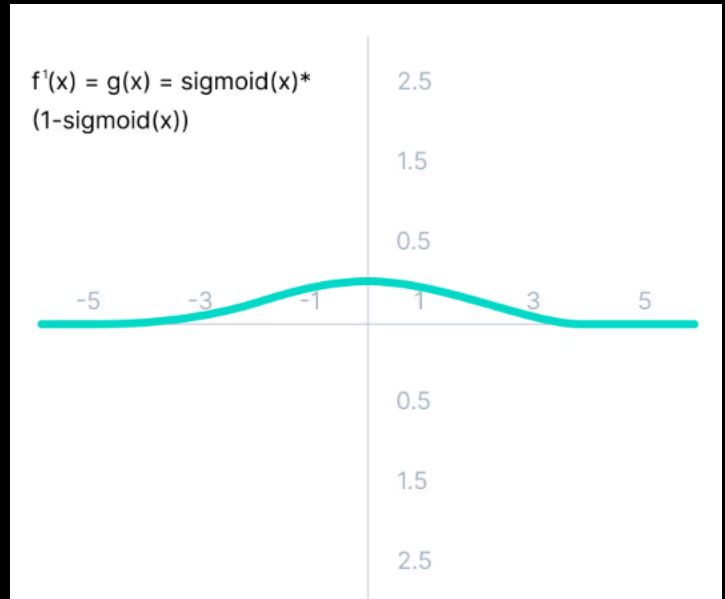
Neural Network:-



* Sigmoid *

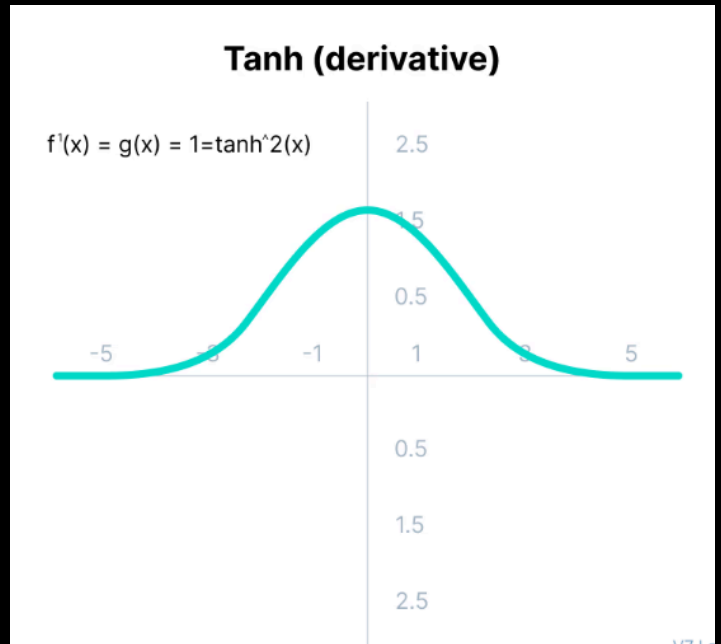
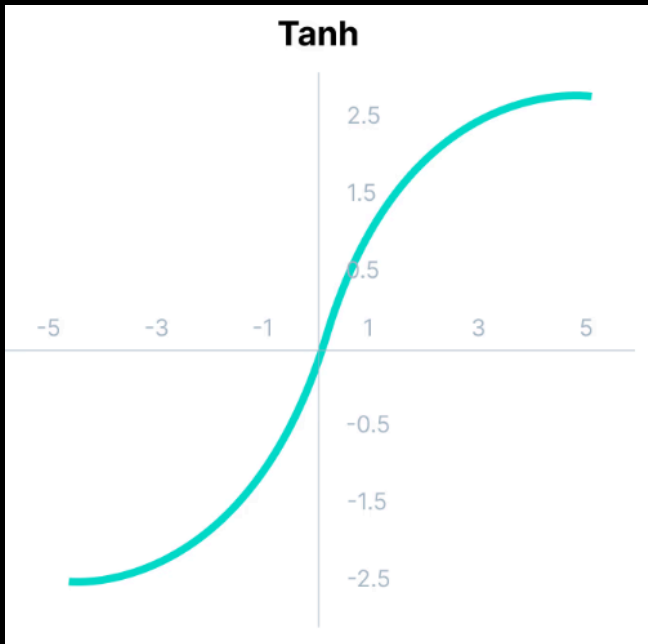
$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = \text{sigmoid}(x) \cdot (1 - \text{sigmoid}(x))$$



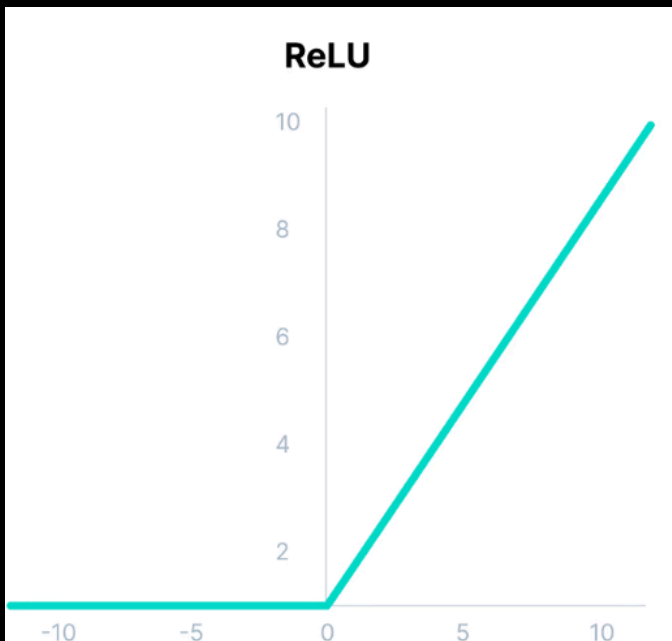
$W_{\text{new}} \approx W_{\text{old}}$, No learning
 No weight update
 Vanishing Gradient Problem $(-3, 3)$

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

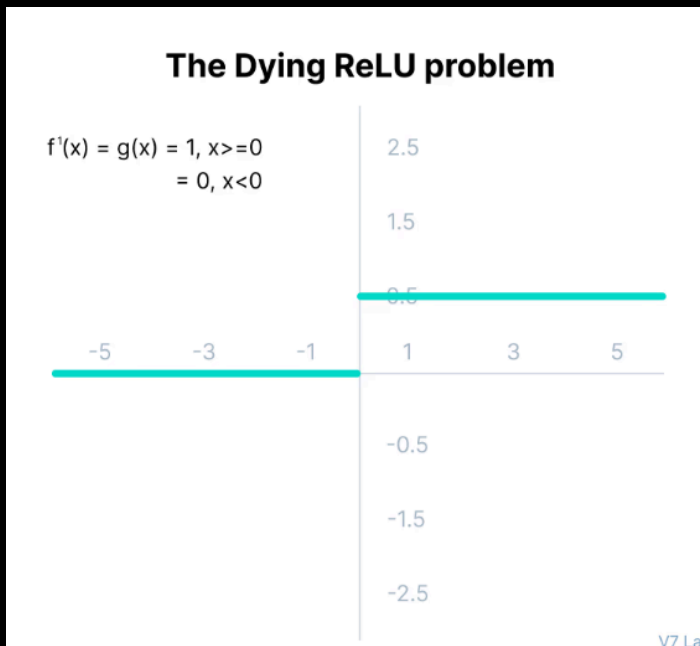
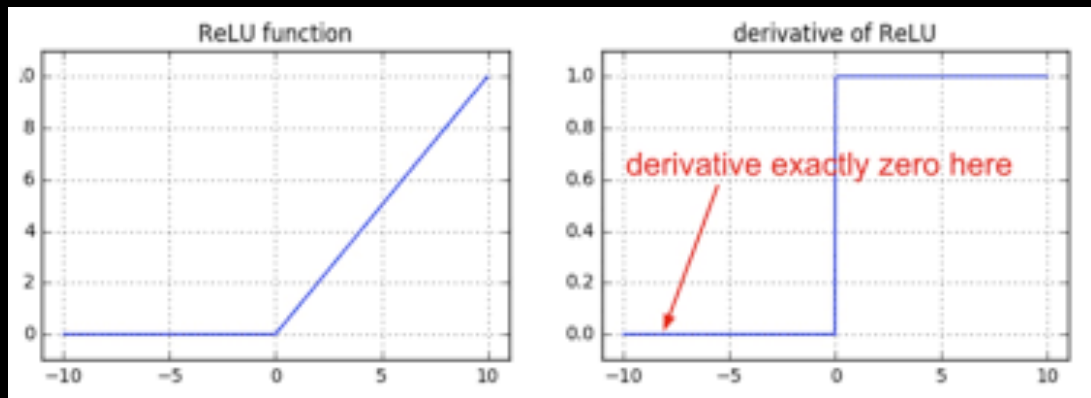


* Vanishing Gradient Problems *

ReLU \rightarrow Rectified Linear Units



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



(i) Negative gradients = zero

Backprop, some neurons are not updated.

Dead Neurons

$$z = wx + b$$

if $z < 0$ for all samples:-

Output = 0

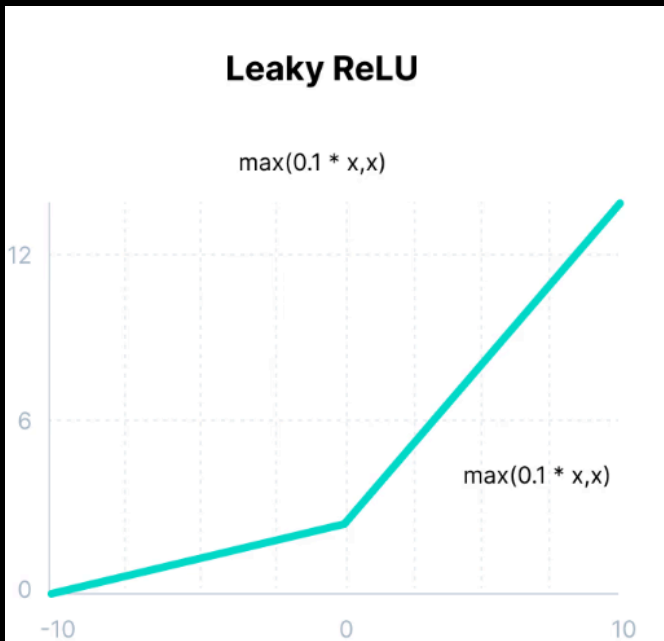
gradient = 0

learning X

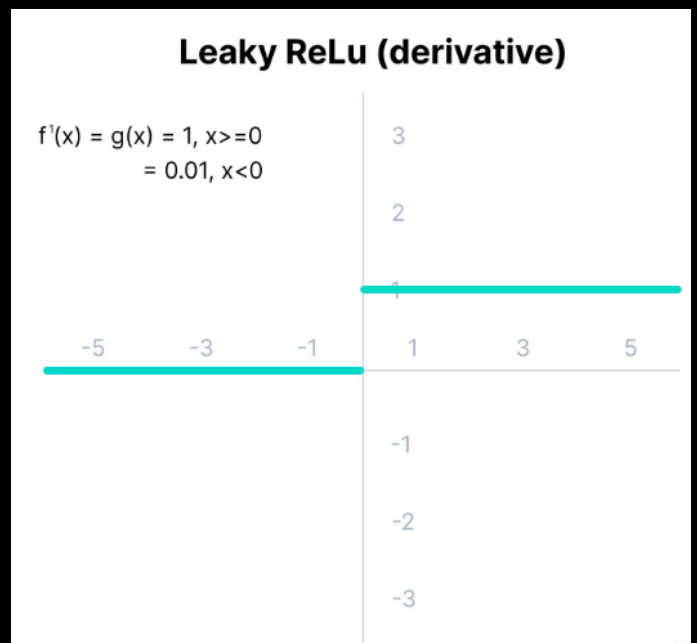
Why?

- 1) Two gradients for negative inputs
- 2) Bad weight initialization
- 3) Imbalanced Data

Leaky ReLU



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

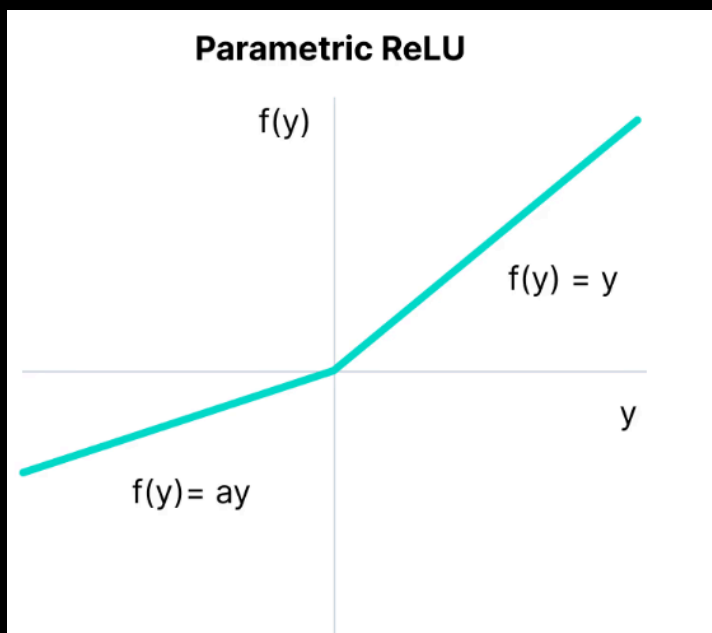


Gradient of our negative input will be very small value.

* learning of the NN will be slow,
time consuming *

$$W_{old} \approx W_{new}$$

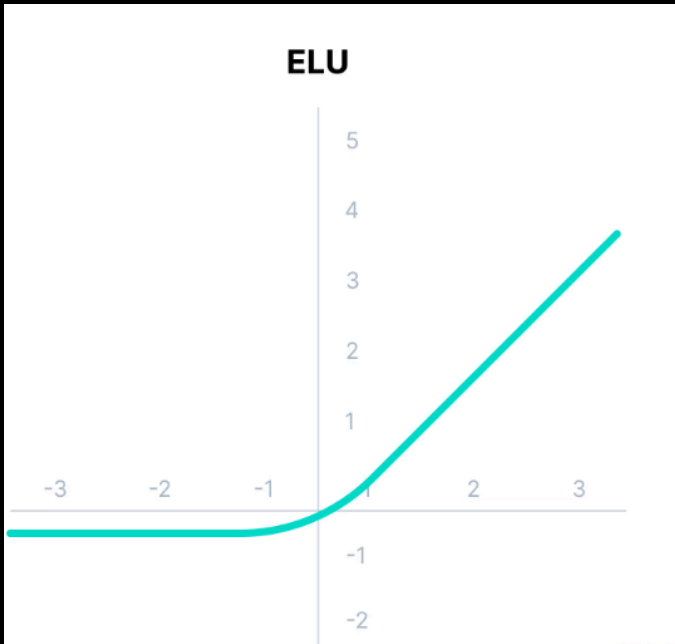
PRELU (Parametric ReLU)



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

a = slope parameter for
our negative values.

Exponential Linear Units (ELU)

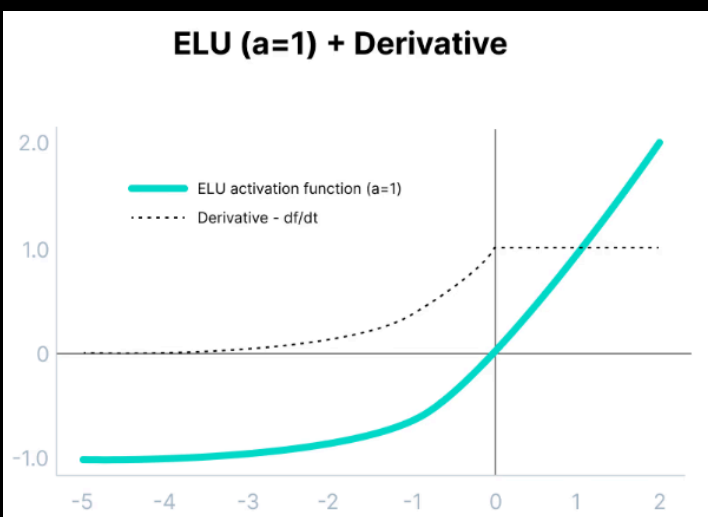


modify the slope of the negative gradients.

log curve

$$f(x) = \begin{cases} x & \text{for } x \geq 0 \\ \alpha(e^x - 1) & \text{for } x < 0 \end{cases}$$

Computationally expensive



Classification Problem

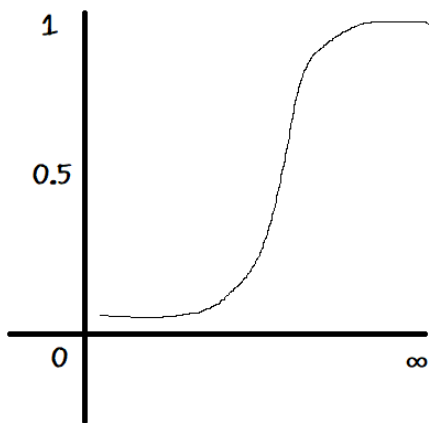
Binary Classification

(i) Sigmoid

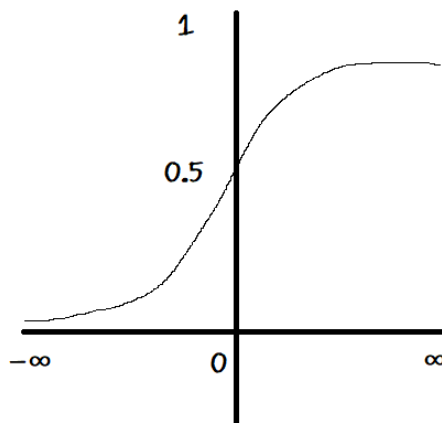
Multiclass

(ii) Softmax

Softmax



Sigmoid



Softmax

$$f(x) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

3 classes

a, b, c

$$= \frac{e^a}{e^a + e^b + e^c}$$

Softmax = combination of multiple sigmoids

example k class $a = \{0.8, 0.7, 0.6, 0.1\}$

Relative Probability

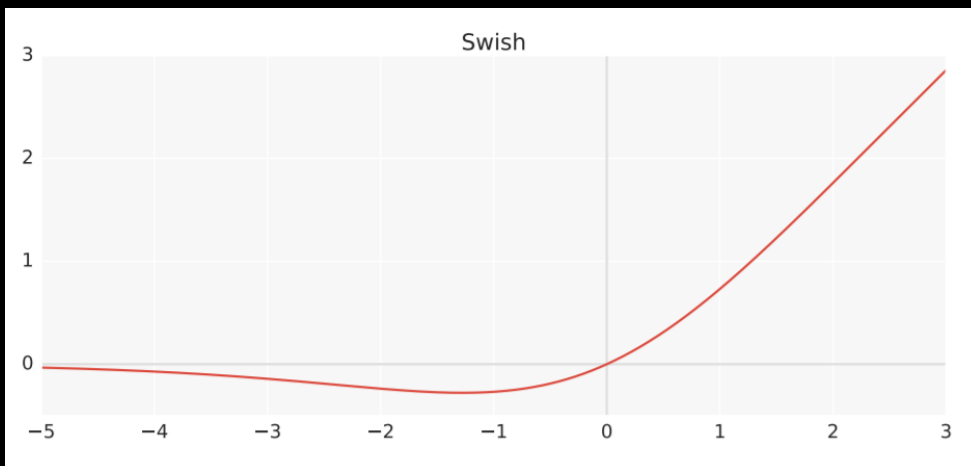
Sigmoid OR Softmax

→ output layer

Swish Activation

Self Gated Activation

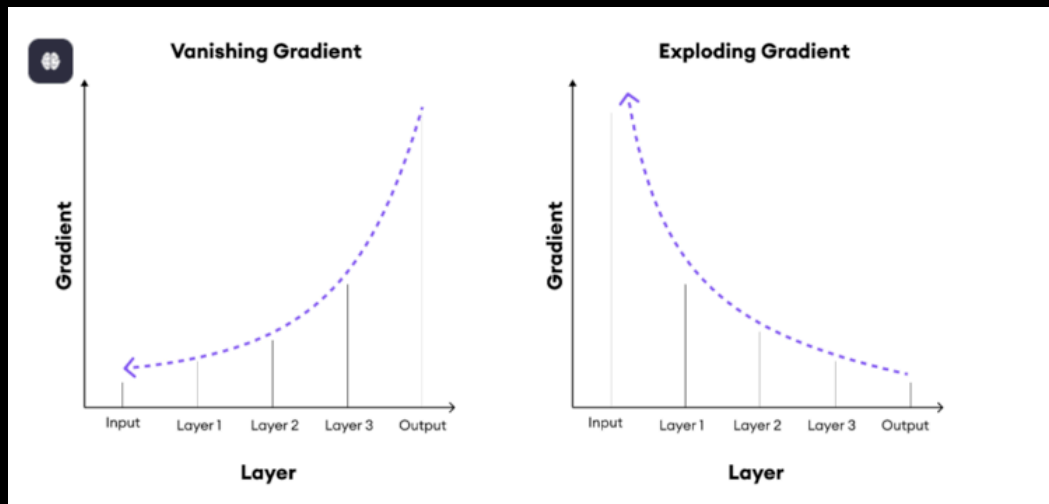
Google
↙
Human brain



$$J(x) = x \cdot \text{Sigmoid}(x)$$

Vanishing Gradient Problem

Exploding Gradient Problem



$$W_{\text{new}} \approx W_{\text{old}}$$

$$W_{\text{new}} > W_{\text{old}}$$

gradient will be very big, reaching towards infinity

$$(i) \quad W_{\text{new}} > W_{\text{old}}$$

$$66 + 0.058$$

$$\text{gradient} = 660000804$$

Generalization of Neural Networks

→ Dropout