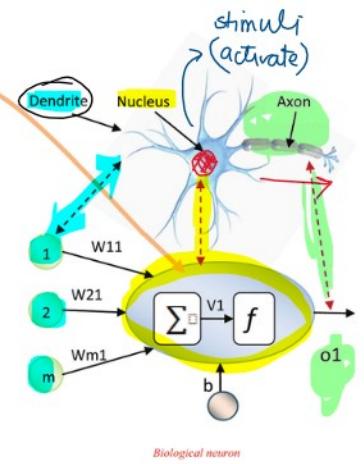
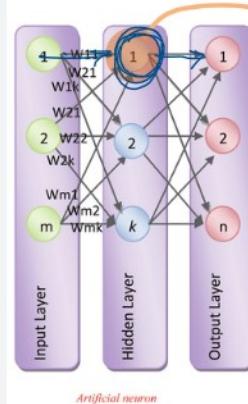
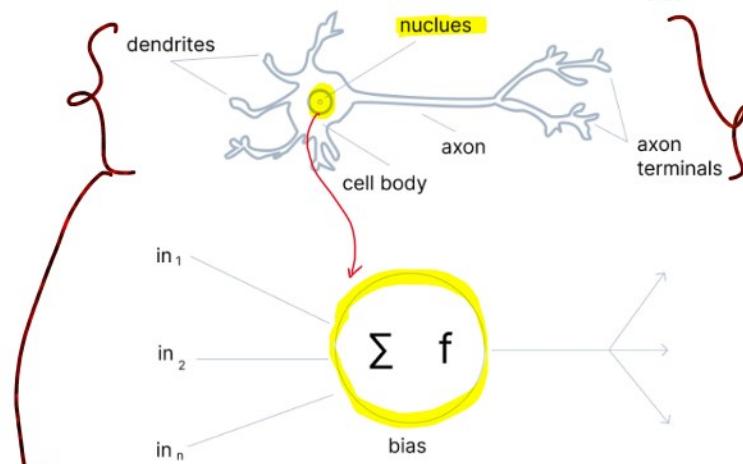


Activation Functions

08 September 2024 20:04

Activation Function helps the neural network to use **important information** while suppressing irrelevant data points

- Activation function decides whether a neuron should be fired or not
(activated)



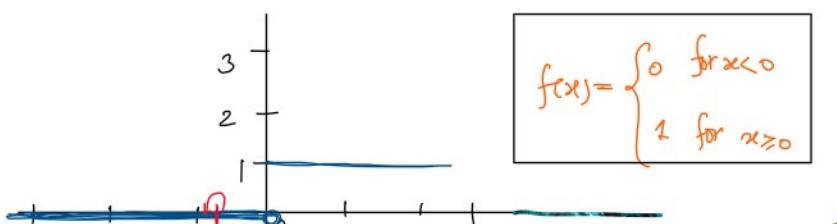
Depending on the nature and intensity of these input signals, the brain processes them to decide whether the neuron should be activated or not. } purpose / role of activation function.

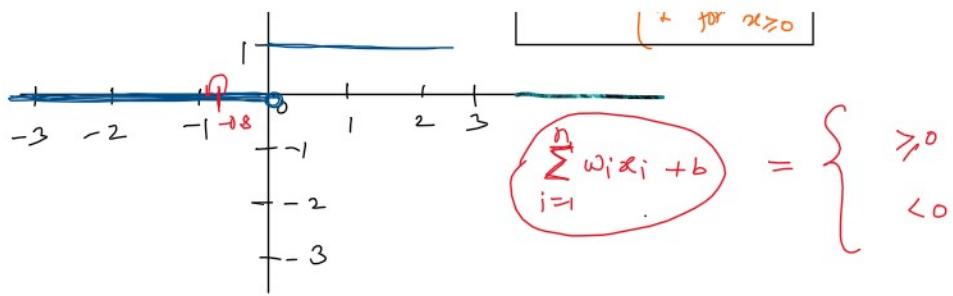
In deep learning, activation function serve this purpose.

- Activation function is to transform summed weighted input from a node into an output value $(\sum_{i=1}^n w_i x_i)$ to be fed to the next hidden layer or output

① Binary Step Function

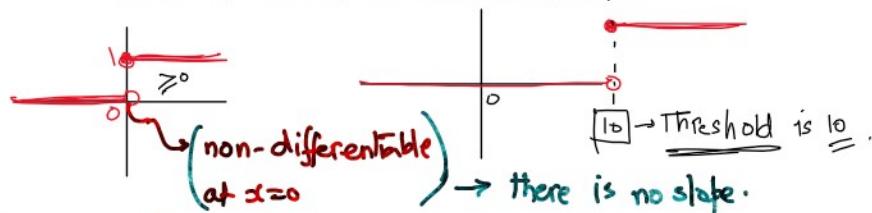
- it depends on a threshold value which decides whether a neuron should be activated or not.



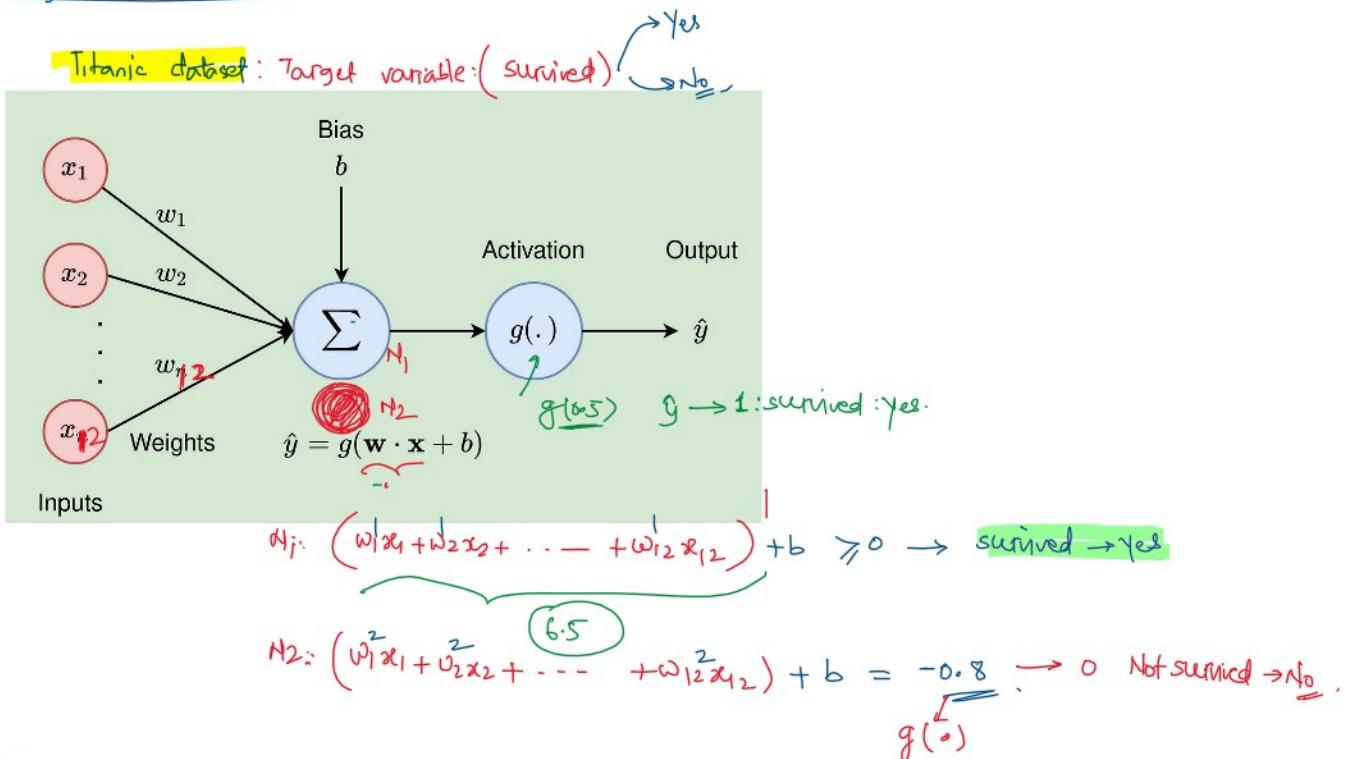


Pro tip interview

- also known as step function or **threshold function**



① Binary classification problem



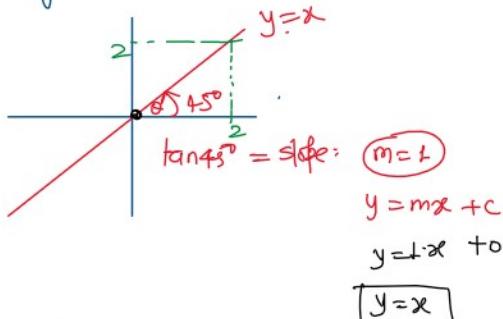
Limitations

- ① It cannot provide multi-class values - hence it can't be used for multi-class classification problems
- ② It doesn't have a derivative at the threshold point which complicates training neural networks using step as the activation function
 - gradient-based optimization techniques can't directly optimize networks with this activation function

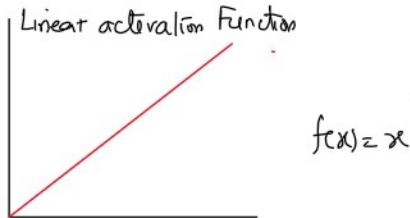
- gradient-based optimization techniques can't directly optimize networks with this activation function

② Linear Activation Function

Identity Function $y = x$

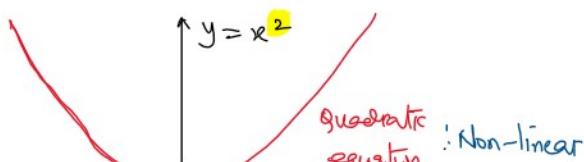


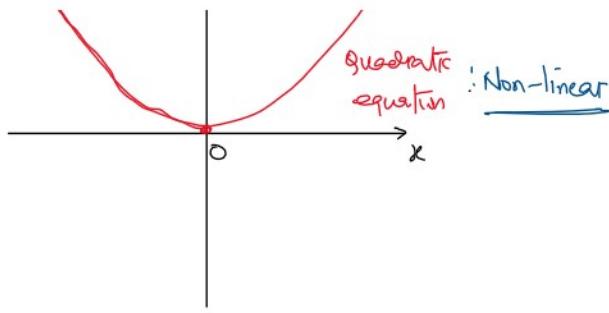
- it is also known as 'no activation' or 'identity function'.
- this activation function doesn't do anything to the weighted sum of the input.



Non-linear activation functions

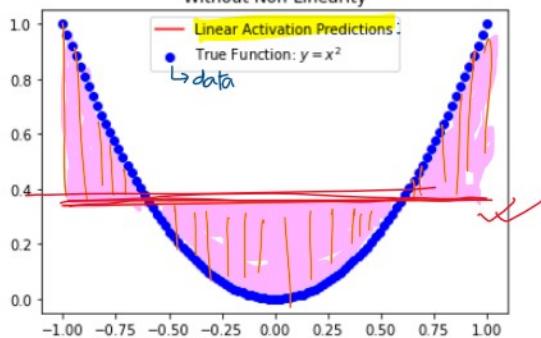
- Linear activation function is simply a linear regression model.
- it does not allow the model to create complex mappings b/w network's input and outputs
- Non-linear activation functions enable the model to learn and capture complex patterns.





scenario #① Modeling a Non-linear function without non-linear activation functions

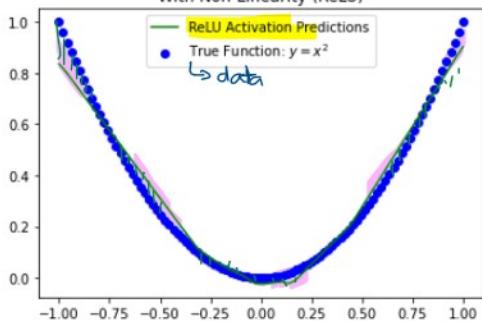
Without Non-Linearity



scenario #② Modeling a non-linear function with non-linear activation function



With Non-Linearity (ReLU)



- Non-linear functions allow backpropagation as the derivative function would be closely related to the input.

- Most of non-linear activation functions allow gradient-based optimization → which means that they are differentiable
 - ↓
 - compute the gradients of the loss w.r.t each weight and use that information to update the weights during training.

and use that information to update the weights during training.

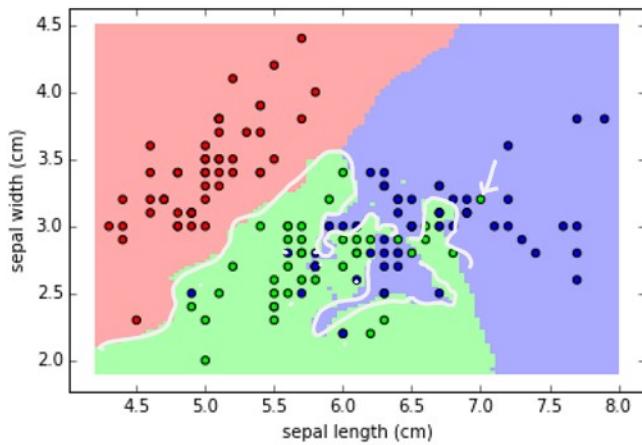
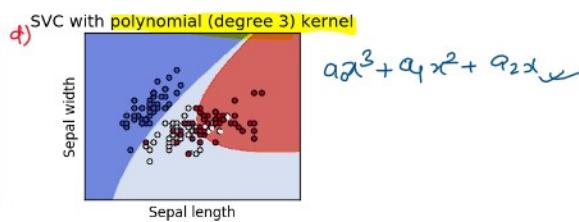
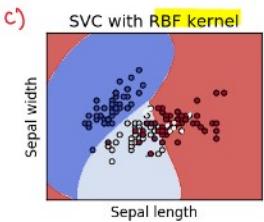
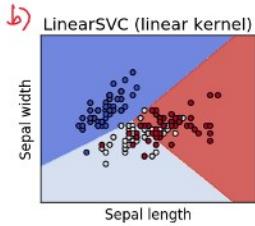
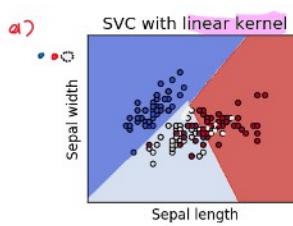
$$\text{loss function} \quad \text{error} = \underline{\text{actual}} - \underline{\text{predicted}}$$

→ helps the model to get the best possible fit for data distribution (curve)

Non-linear decision boundaries

Note: sepal width and sepal length

Algorithm: Support vector classification

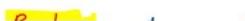


digit recognition

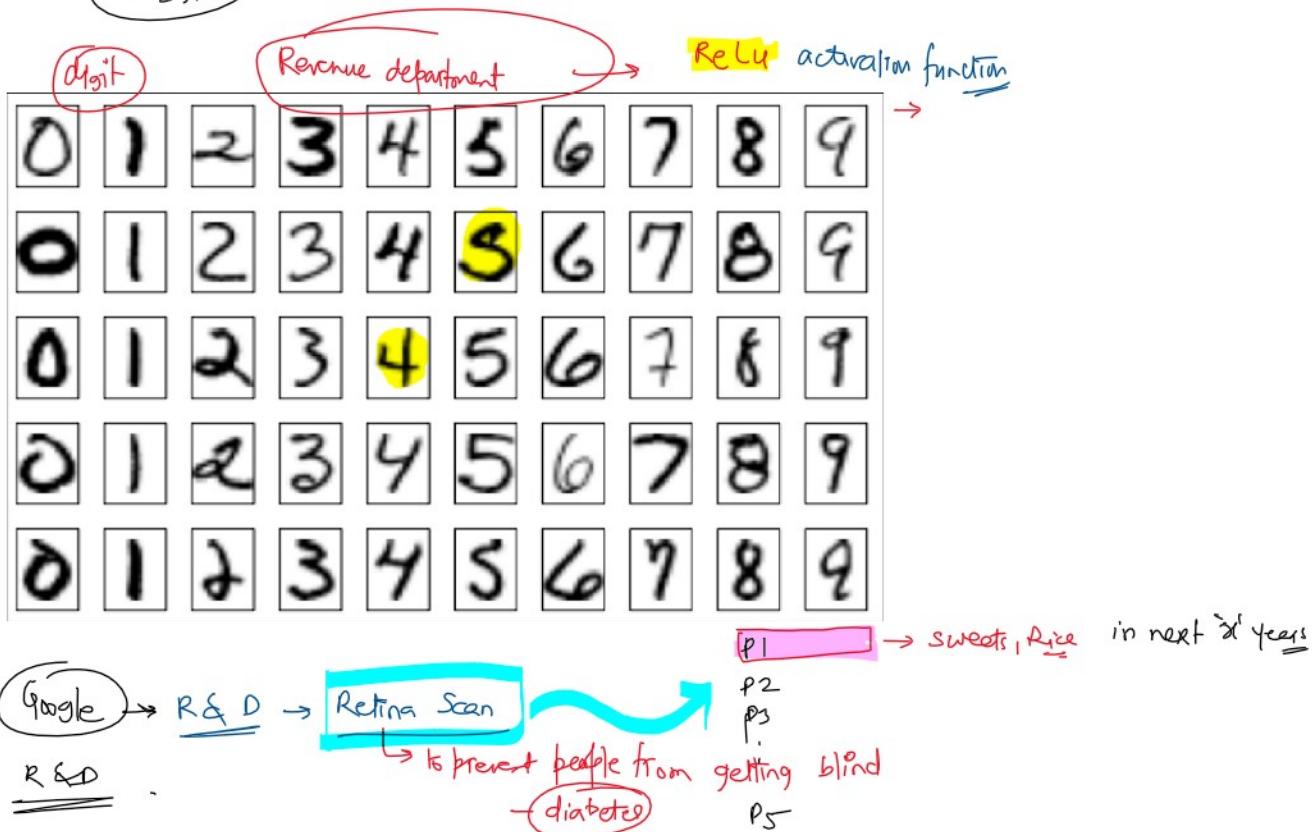


1 2 3 4 5

6-25%



6-255



Nonlinearity activation function helps in back propagation

$$\# \quad y = 2x + 3 \quad | \quad y = x$$

$$\frac{dy}{dx} = 2 \quad | \quad \frac{dy}{dx} = 1$$

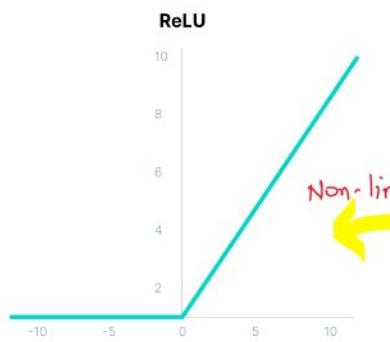
gradient or derivative is constant

- If all the layers are linear activations, the derivative of the output w.r.t any input or weight will be constant which basically means that the model can't effectively learn relationships
- However if we introduce non-linearities using non-linear activation functions such as ReLU, each layer will have its own gradient (unique) and that eventually helps each layer to learn useful patterns which contribute to minimizing

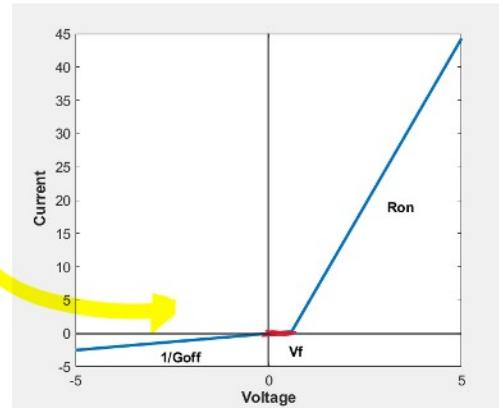
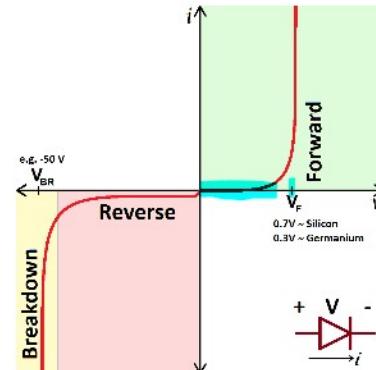
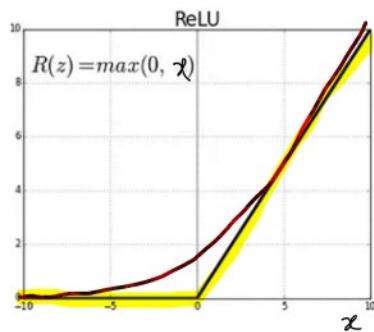
useful patterns which contribute to minimizing the loss function.

Common Non-linear activation functions

#① ReLU (Rectified Linear Unit)



- Although the above graph gives an impression of a linear function however ReLU is a non-linear activation function and has a derivative too.



ReLU is also known piece wise linear function
 It generates the output (Non-linear)
 as the input directly if input is positive,
 otherwise, it will output zero

$$\text{input} > 0, \text{ output} = \text{input}$$

$$\text{input} \leq 0, \text{ output} = 0$$

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

→ Neurons will be deactivated if

output of transformation is less than or equal to '0'

Ex:

Input = -5	Output = 0
Input = 0	Output = 0
Input = 5	Output = 5

* It's simple and computationally efficient, leading to faster training and convergence.

Since only a certain no. of neurons are activated, the ReLU function is far more computationally efficient

with compared to Sigmoid and Tanh functions

In general;

$$f(x) = \max(x, 0) \quad \text{threshold is '0' (without)}$$

$$\downarrow \quad T$$

$$f(x) = \begin{cases} x & \text{if } x > T \\ 0 & \text{if } x \leq T \end{cases}$$

$$\text{here } [T=0]$$

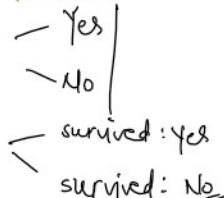
In case, $T = 10$ [with threshold $T=10$]

$$f(x) = \begin{cases} x & \text{if } x > 10 \\ 0 & \text{if } x \leq 10 \end{cases}$$

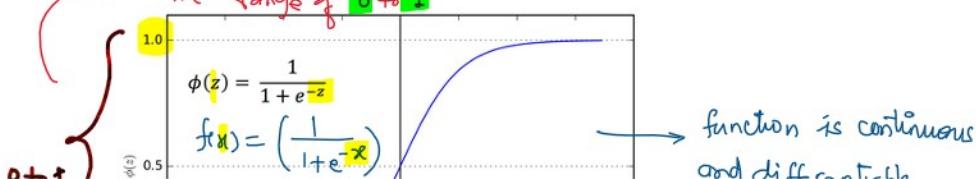
* ② Sigmoid Activation Function

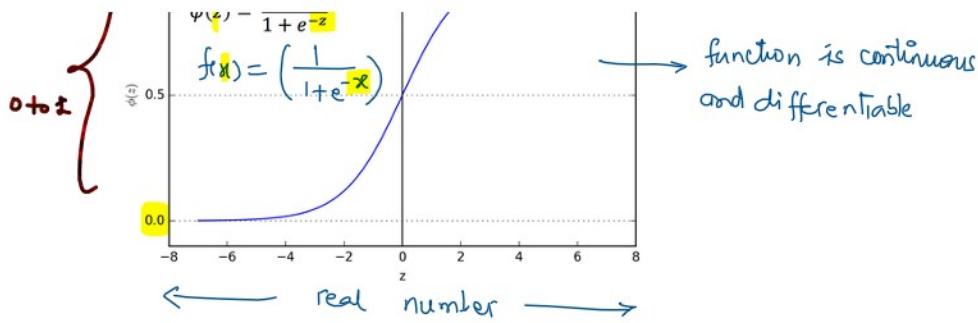
one of the most common activation functions, used for binary classification functions.

use-case: Fraud detection



Probability of anything exists between the range of 0 to 1





Note: softmax activation function is a more generalized activation function which is used for multiclass activation function

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

differentiate
↓ gradient: w.r.t x'

$$f'(x) = \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right]$$

$$f'(x) = -1 \left(1+e^{-x} \right)^{-2} \cdot \frac{d}{dx} (1+e^{-x})$$

$$f'(x) = \frac{-1}{(1+e^{-x})^2} * (0 - e^{-x})$$

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

gradient

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

sigmoid

Derivative of sigmoid is:

$$g(x) = f'(x) = \text{sigmoid}(x) * [1 - \text{sigmoid}(x)]$$

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

vanishing gradient problem

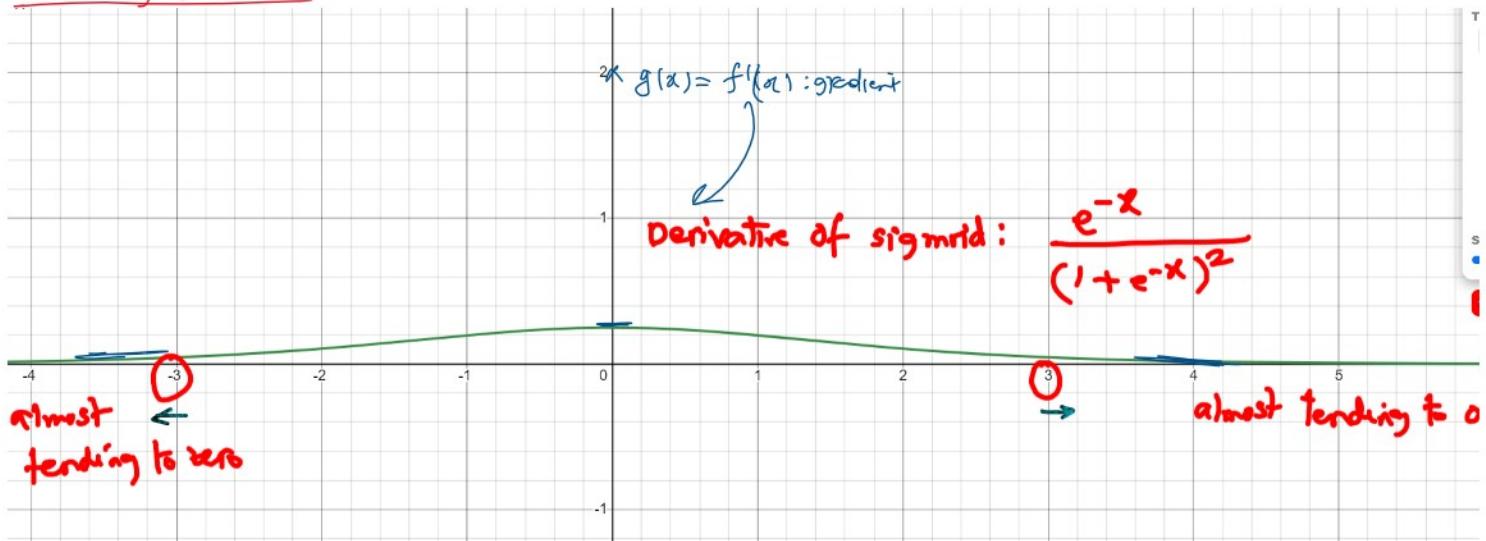
$$\begin{aligned} \frac{d}{dx} (x^n) &\quad (n=1) & \frac{d}{dx} (x^n) &= n x^{n-1} \cdot \frac{d}{dx} \\ &= -1 (x) & &= n x^{n-1} \end{aligned}$$

<https://www.khanacademy.org/math/ap-calculus-ab/ab-differentiation-1-new>
<https://www.khanacademy.org/math/ap-calculus-ab/ab-differentiation-1-new/ab-2-6a/v/derivative-properties-and-polynomial-derivatives>
<https://www.khanacademy.org/math/ap-calculus-ab/ab-differentiation-2-new/ab-3-1a/v/chain-rule-introduction>

$$\begin{aligned} \frac{d}{dx} (e^x) &= e^x & \frac{d}{dx} e^{-x} &= e^{-x} \cdot (-1) = -e^{-x} \\ \frac{d}{dx} (e^{-x}) &= -e^{-x} & & \\ \frac{d}{dx} (\text{constant}) &= 0 & & \end{aligned}$$

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

Vanishing gradient



As we can see from the above figure, the gradient values are only significant for range -3 to 3 and the graph gets much flatter outside this range.
<https://www.desmos.com/calculator>



As the gradient values approach zero, the network ceases to learn and suffers from the vanishing gradient problem.

For reference

c) Update rule:

* For $j=0$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} [J(\theta_0, \theta_1)]$$

→ is almost zero

Input → by model!

$\theta_0 \stackrel{\Delta}{=} \theta_0_{\text{new}}$ ✓

$\theta_0 \stackrel{\Delta}{=} \theta_0_{\text{old}}$

$\equiv \times 100 \quad \} \text{ time taking}$

#③ Softmax activation function

- it is a type of activation function that transforms a vector of real-valued numbers into a probability distribution

For ex. (For reference): drawing analogy.

$$\begin{bmatrix} 5 \\ 3 \\ 2 \\ 5 \end{bmatrix} \rightarrow \text{sum of real values} = 15$$

$\frac{5}{15} = \frac{1}{3} = 33.33\%$

$\frac{3}{15} = \frac{1}{5} = 20\%$

$\frac{2}{15} = \frac{2}{15} = 13.33\%$

$\frac{5}{15} = \frac{1}{3} = 33.33\%$

$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$

Total ✓ $\frac{z_i}{\sum_{j=1}^k z_j}$?

Example:

Raw prediction scores: $K = 12.26$

1	2.33	$P(\text{class 1})$	$\frac{e^{2.33}}{e^{2.33} + e^{-1.46} + e^{0.56}} = \frac{10.278}{12.26} = 0.8382$	83.82%
2	-1.46	$P(\text{class 2})$	$\frac{e^{-1.46}}{K} = \frac{0.232}{12.26} = 0.0189$	1.89%
3	0.56	$P(\text{class 3})$	$\frac{e^{0.56}}{K} = \frac{1.75}{12.26} = 0.1427$	14.27%
Raw scores				100%

Key characteristics

1. Normalization: Softmax activation function normalizes the input values into a probability distribution ensuring that the sum of all output values is 1.

2. Exponentiation: Leveraging exponentiation of the inputs, the softmax function in DL amplifies the difference b/w the inputs, making the largest value more

the softmax function ... - ... - ... - ...
b/w the inputs, making the largest value more
pronounced in the output probabilities

For a vector $z = [z_1, z_2, \dots, z_n]$, the softmax function for the i -th element is defined as:

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

Where:

- z_i is the i -th element of the input vector z ,
- e is the exponential function,
- n is the number of elements in the vector.

Application

① Neural Networks: In the output layer of a neural net for multi-class classification

For example: In Image classification

whether the given image is a cat, dog or car

② Logistic Regression (Multinomial) softmax is an extension
of the sigmoid function to perform multi-class classification

Key advantages of softmax function

- ① Multi-class classification
- ② Probability distribution - raw scores to probabilities
logits
- ③ Differentiable
↳ softmax is smooth, differentiable function,
making it suitable for use in backpropagation
- ④ Relative Importance: it emphasizes the relative differences
between the input values.
↓
one input significantly higher than others → assign a probability

n | . . . *cently he has 1 n then > , n) t₅*
w a n cou - | t₁ then

⑤ stable output