



15/12/2025

data collection

data preprocessing

Build the AI model

- Training set
- Test set
- validation set

70% - Training
{Shuffled and
fed to the AI
model}

30%

- Test (15%)
- validation (15%)

enable early stopping

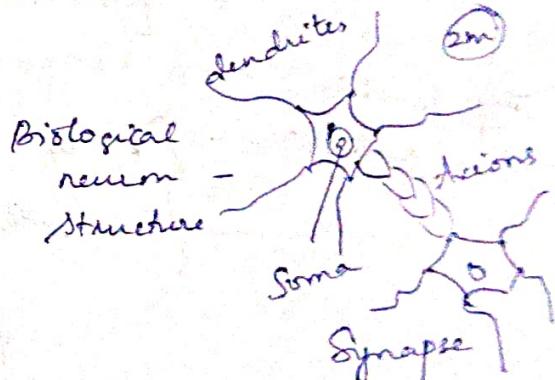
performance measure - Accuracy

- Confusion matrix

		P	N
P	TP	FN	
	FP	TN	

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- 3 types - Sensory neuron
of neuron
Motorway neuron
Interconnection neuron



biological neuron

Artificial Neuron

- call it as a neuron

- call it as a computational
node

- Inputs processed through
dendrites

- Inputs

- Soma - cell body

- processing centre

- Axion

- output

- Synapse

- connection or a junction that connect the output of one neuron to the input of the other neuron

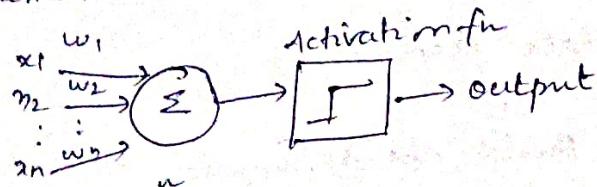
- less complex

- Complex

- less connected and perform
very less or complex tasks

McCulloch + Pitts neuron model

- mathematical model to denote artificial neuron



$$\text{Sum} = \sum_{i=1}^n x_i w_i$$

output = $f(\text{Sum})$ Threshold based activation function

$$f(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

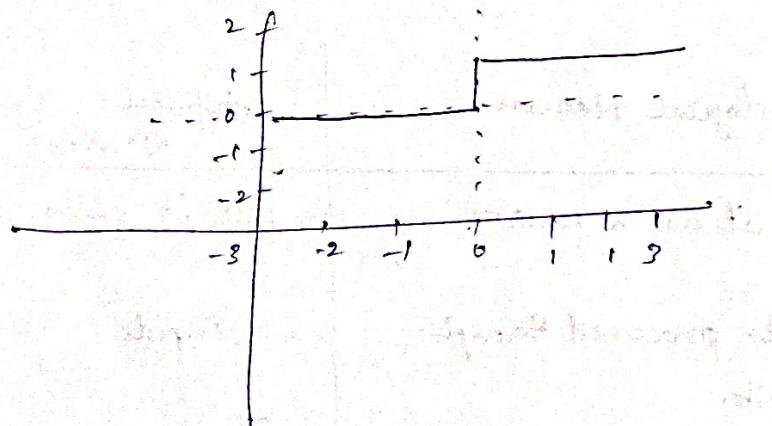
$$y_p = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$f(\text{Sum})$ or $f(\text{Sum})$

$$\begin{aligned} x^T w &= [x_1 \ x_2 \ x_3] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \\ &= [x_1 w_1 + x_2 w_2 + x_3 w_3]_{1 \times 1} \end{aligned}$$

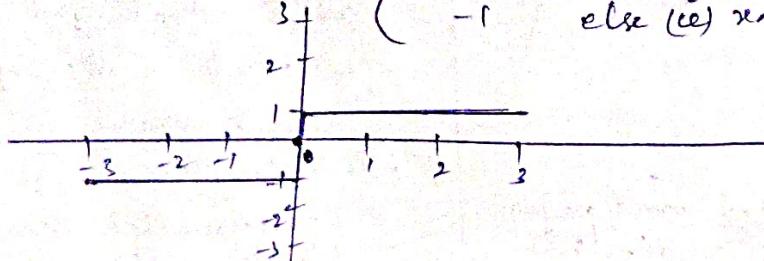
Step function:

$$\text{out}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$



Sigmoid activation fn: (-1, 0, 1)

$$\text{out}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{else (i.e.) } x < 0 \end{cases}$$



$$x = [1 \ 1 \ 1 \ 2]$$

$$w = [0.5 \ 0.5 \ 0.6]$$

$$\sum x_i w_i = 0.5 + 0.5 + 1.2 = 2.2 > 1$$

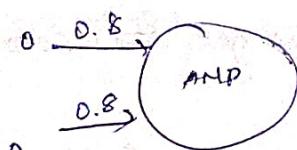
$$\text{out}(x) = 1$$

$x_1 \quad x_2 \quad y_{\text{AND}}$

0	0	0
0	1	0
1	0	0
1	1	1

$$T=0.5$$

$$T=1$$



$$\text{out}(0,0) = 0 \times 0.8 + 0 \times 0.8 = 0 < 1 = 0$$

$$\text{out}(0,1) = 0 \times 0.8 + 1 \times 0.8 = 0.8 < 1 = 0$$

$$\text{out}(1,0) = 1 \times 0.8 + 0 \times 0.8 = 0.8 < 1 = 0$$

$$\text{out}(1,1) = 1 \times 0.8 + 1 \times 0.8 = 1.6 > 1 = 1$$

16/12/2025

Perception model

- binary classifier

- used to solve the linearly separable problems

Steps - perception learning algo :

1. $\{(x_1, t_1), (x_2, t_2), (x_3, t_3) \dots (x_n, t_n)\}$

$$c_1 = +1 \quad c_2 = -1$$

2. Initialize the weights randomly $\in [-1, 1]$

3. Initialize the learning rate $\in [0, 1] \Rightarrow \eta$

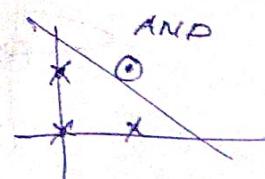
4. (x_i, t_i) for each input, do the forward propagation



$$\text{Sum} = \sum_{i=1}^n x_i w_i + \text{bias}$$

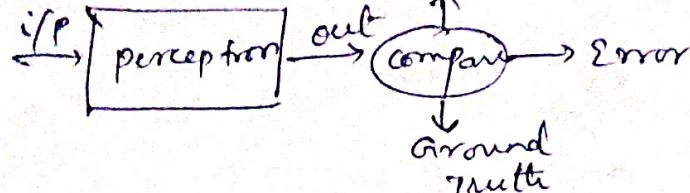
$$\text{bias } (x_0, w_0) \quad \downarrow \\ x_0 = 1$$

$$\text{Sum} = \sum_{i=0}^n x_i w_i$$



Then, the activation function must be applied
 $\text{out} = \text{Sign}(\text{net}) \xrightarrow{\text{signum}}$

$$(a) \text{out} = \begin{cases} 1 & \text{if net} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$



5) find the error: $\text{err} = t - \text{out}$

$$w_i = w_i + \Delta w_i$$

6) update: $\Delta w_i = \text{error} * \eta * x_i$

7) repeat the above steps till we get the required error. (st-6)

Q

$$\text{Err}(t) = t - \text{out} \quad (\text{Ground truth} - \text{predicted})$$

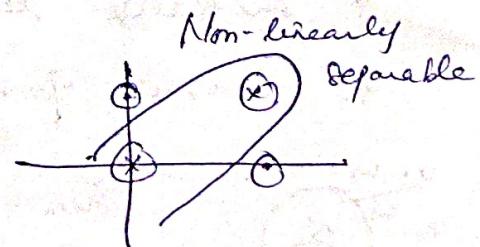
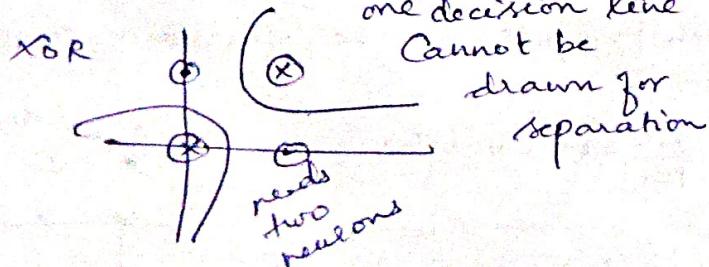
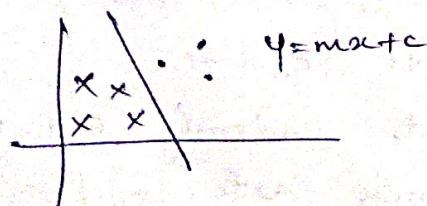
$t = 1 \quad \text{out} = 1 \quad \text{err} = 0 \quad \text{no update}$

$t = -1 \quad \text{out} = -1 \quad \text{err} = 0 \quad \text{no update}$

$t = 1 \quad \text{out} = -1 \quad \text{err} = 1 + 1 = 2 \quad \text{update (increase weight)}$
 $\Delta w_i = \text{error} * \eta * x_i$

$t = -1 \quad \text{out} = 1 \quad \text{err} = -1 + -1 = -2 \quad \text{update (decrease weight)}$

perceptron as a binary classifier



Problem :

$$\begin{array}{cccc}
 i_1 & i_2 & t \\
 \hline
 x_1 & 0 & 0 & -1 \\
 x_2 & 0 & 4 & +1 \\
 x_3 & 1 & 1 & +1 \\
 x_4 & 1 & -3 & -1
 \end{array}
 \quad w_0 = w_1 = w_2 = 1 ; \eta = 1$$

$$x_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad x_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix} \quad x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad x_4 = \begin{pmatrix} 1 \\ -3 \end{pmatrix}$$

Augment 1 to all inputs:

~~$$x_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad x_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix} \quad x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad x_4 = \begin{pmatrix} 1 \\ -3 \end{pmatrix}$$~~

~~$$x_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad x_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$$~~



~~$$x_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad x_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix} \quad x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad x_4 = \begin{pmatrix} 1 \\ -3 \end{pmatrix}$$~~

~~$$w_p = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad w = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$~~

$$s_1 = w^T x_p = (1 \ 1 \ 1) \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$= 1 \times 1 + 1 \times 0 + 1 \times 0$$

$$= 1$$

$$t = -1 \quad \text{and} \quad \text{out} = 1$$

$$\text{err} = -1 - 1 = -2$$

$$w_0 = w_0 + Aw_0 \quad ; \quad w_0 = (1 \ 1 \ 1) + 1 \times -2 \times (1 \ 0 \ 0)$$

$$w = w + Aw \quad = (1 \ 1 \ 1) + (-2 \ 0 \ 0)$$

$$= (-1 \ 1 \ 1)$$

$$x_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix} \quad w = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

$$s_2 = (-1+1) \begin{pmatrix} 1 \\ 0 \\ 4 \end{pmatrix}$$

$$= -1 + 0 + 4$$

$$= 3$$

$$t = +1, \text{ out} = \text{sign}(3) = 1$$

$$\text{err} = t - \text{out} = 1 - 1 = 0$$

so, no update in weight

$$x_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad w = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

$$s_3 = w^T x_3$$

$$= (-1+1) \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$= -1 + 1 + 1$$

$$= 1$$

$$\text{out} = \text{sign}(1) = 1$$

$$\text{err} = t - \text{out} = 1 - 1 = 0$$

so, no update in weight

$$x_4 = \begin{pmatrix} 1 \\ 1 \\ -3 \end{pmatrix} \quad w = \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}$$

$$s_4 = w^T x_4$$

$$= (-1+1-3) \begin{pmatrix} 1 \\ 1 \\ -3 \end{pmatrix}$$

$$= -1 + 1 - 3$$

$$= -3 + 1$$

$$= -3$$

$$\text{out} = \text{sign}(-3) = 0$$

$$\text{err} = t - \text{out} = -1 - 0 = -1$$

$$w = w + \eta w$$

$$w = \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}$$

$\eta \rightarrow 1$, then the epoch converges quickly

η should be as small as possible

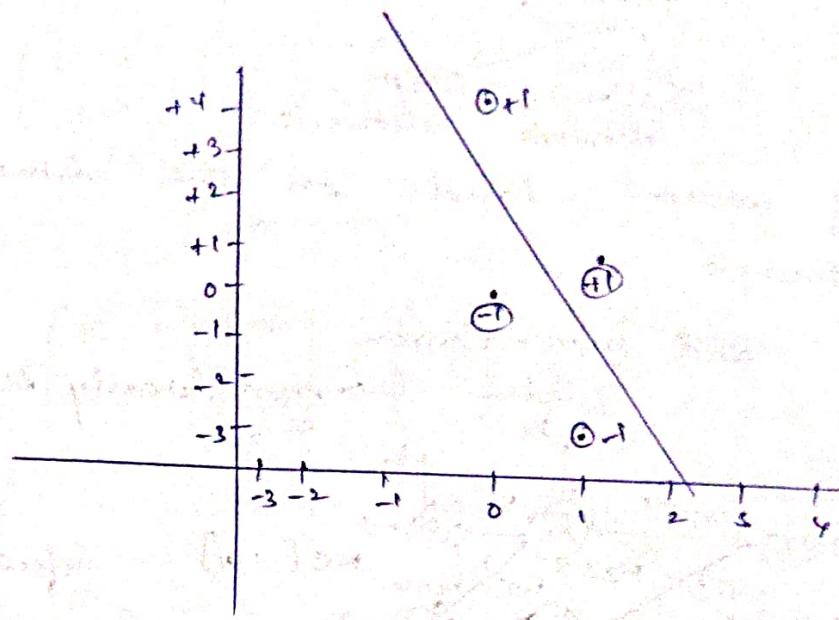
$$\text{out} = \text{sign}(-3) = -1$$

$$\text{err} = t - \text{out}$$

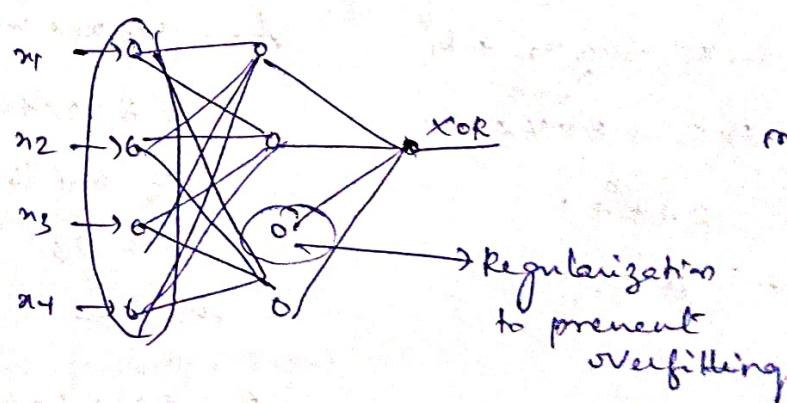
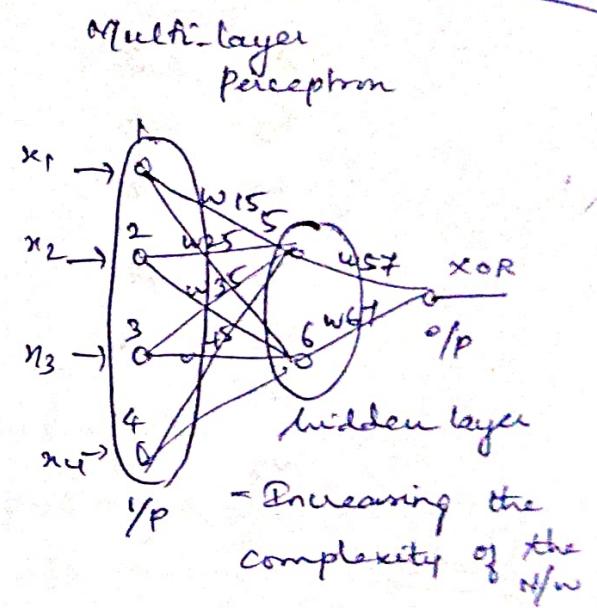
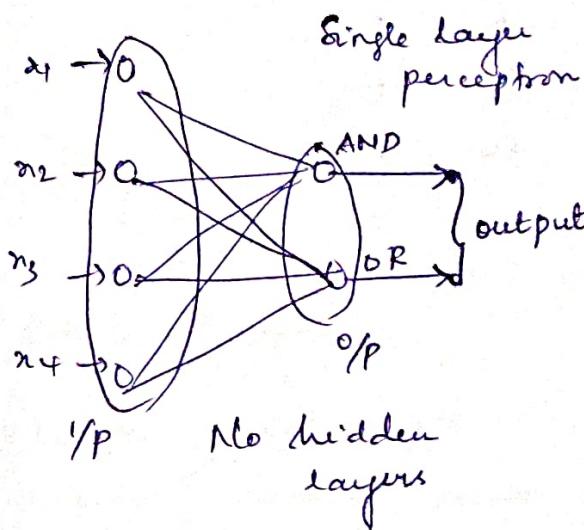
$$\text{err} = -1 - (-1)$$

$$\text{err} = 0$$

hence, no update



$$\text{OR: } \begin{matrix} -1 & 1 & 1 & 1 & 2=1 \\ t_1 & t_2 & t_3 & t_4 \end{matrix}$$



one/more hidden
layer.
↓
Shallow
network

More than
two
↓
Deep
Neural
Networks

22/12/2025

MLP

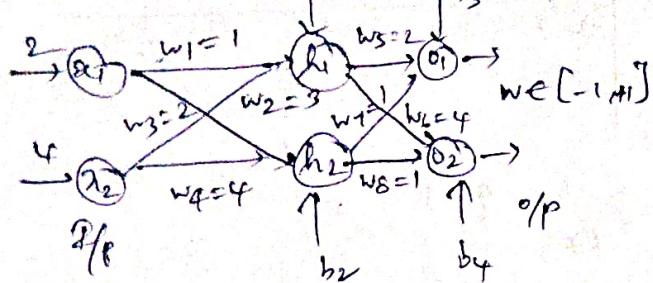
shallow
Network

Deep
network

No. of neurons - based on the dimensionality of the input sample

Convex and non-Convex

hidden \rightarrow non-linearly separable problem



feed forward
neural N/Ans
Acyclic graph

w_i	value
w_1	1
w_2	3
w_3	2
w_4	4
w_5	2
w_6	4
w_7	1
w_8	1

b_i	value
b_1	1
b_2	2
b_3	1
b_4	1

Solution:

$$\begin{aligned} h_1 &= w_1 x_1 + w_3 x_2 + b_1 \\ &= 1 \times 2 + 2 \times 4 + 1 \\ &= 2 + 8 + 1 \\ h_1 &= 11 \end{aligned}$$

$$\begin{aligned} h_2 &= w_2 x_1 + w_4 x_2 + b_2 \\ &= 3 \times 2 + 4 \times 4 + 1 \\ &= 6 + 16 + 1 \\ &= 23 \\ h_2 &= 23 \end{aligned}$$

$$\begin{aligned} o_1 &= h_1 w_5 + h_2 w_7 + b_3 \\ &= 11 \times 2 + 23 \times 1 + 1 \\ &= 22 + 23 + 1 \\ &= 46 \end{aligned}$$

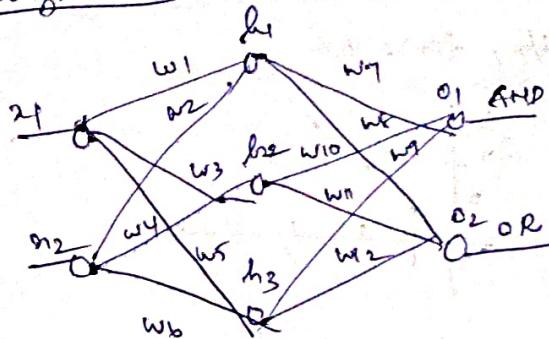
$$\begin{aligned} o_2 &= h_1 w_6 + h_2 w_8 + b_4 \\ &= 11 \times 4 + 23 \times 1 + 1 \\ &= 44 + 23 + 1 \\ &= 68 \end{aligned}$$

By matrix representation:

$$\begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2+1 \\ 6+16+1 \end{bmatrix} = \begin{bmatrix} 11 \\ 23 \end{bmatrix}$$

$$\begin{bmatrix} o_1 \\ o_2 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 4 & 1 \end{bmatrix} \begin{bmatrix} 11 \\ 23 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 22+23+1 \\ 44+23+1 \end{bmatrix} = \begin{bmatrix} 46 \\ 68 \end{bmatrix}$$

Assignment:



w_i value

w₁ 1

w₂ 2

w₃ 1 $b_i = 1$

w₄ 3 for i=1 to 3

w₅ 1 $b_4 = 2$

w₆ 4 $b_5 = 1$

w₇ 1

w₈ 2

w₉ 1

w₁₀ 1

w₁₁ 3

w₁₂ 3

Activation functions:

- ↳ To introduce the non-linearity
- ↳ To learn more complex problems & patterns
- ↳ Giving and - not giving the neurons

Linearly separable properties

↳ Homogeneity (Scaling)

$$f(cx) = k f(x)$$

② ↳ Additive property (Superposition)

$$f(a+b) = f(a) + f(b)$$

If the rps not proportional to the o/p's, then they are said to be non-linearly separable problems, doesn't satisfy the superposition properties.

Exponential - $a e^x$

Logarithmic - $\ln x$

Trigonometric - $\cos x$

Polynomial - $x^2 + bx + c$

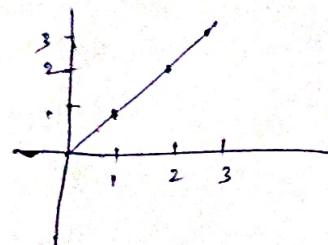
Activation function → linear

Non-linear

Linear activation function

① identity

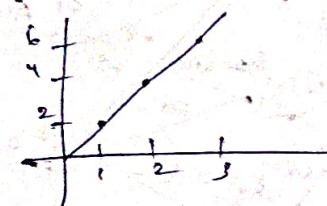
$$f(x), f(x) = x$$



$$f(x) = ax$$

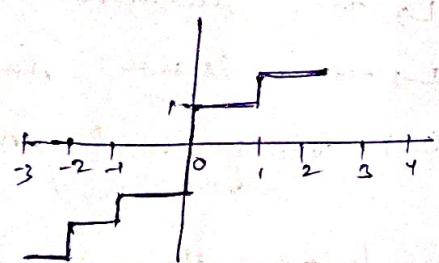
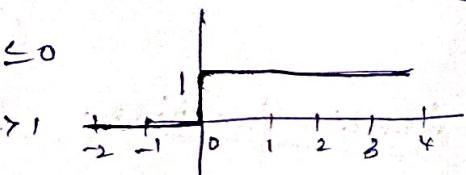
$$f'(x) = a$$

$$f(x), f(x) = 2x$$



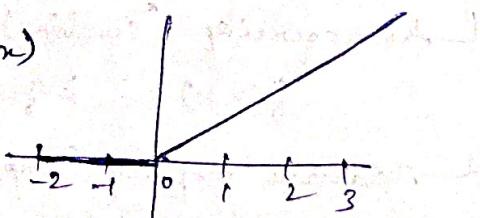
② Binary Step Activation Function (threshold activation)

$$f(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 1 \end{cases}$$



③ ReLU Activation Function

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



not fixed - dead - sparse

→ Make the models converge quickly

→ Better than sigmoid and tanh activation fun

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

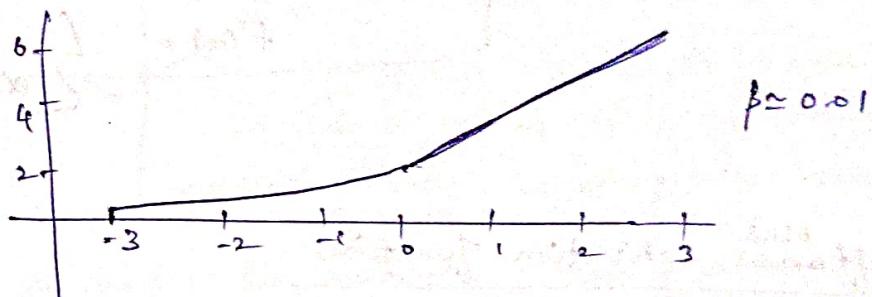
- vanishing gradient problem

$$\alpha = -4, f(\alpha) = 0$$

$$\alpha = 17, f(\alpha) = 17$$

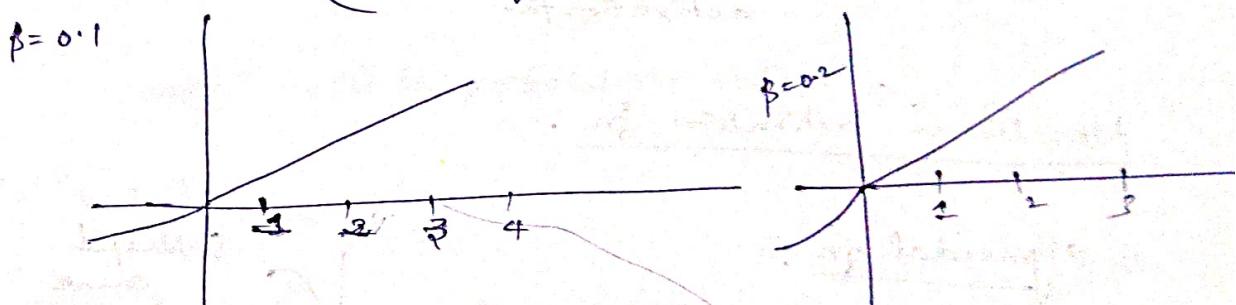
④ Leaky ReLU Activation fn

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



⑤ parametric ReLU Activation fn

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

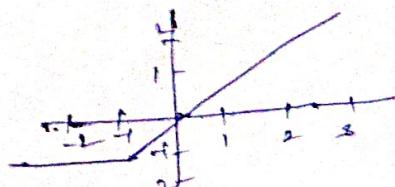


⑥ Randomized ReLU Activation fn

β value can chosen randomly

used for regularization

⑦ Shifted ReLU Activation fn

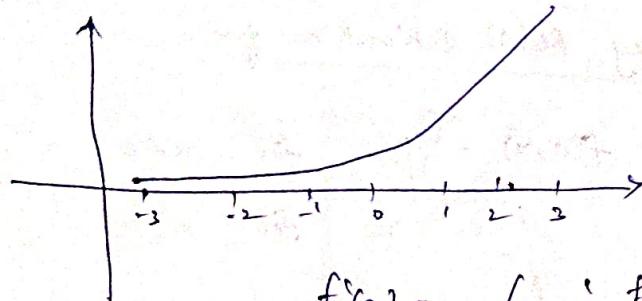


$$g(x) = \begin{cases} x & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$

① Exponential linear unit Activation fn (ELU) :

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

$f(x, \alpha) \in [0, \infty]$
log fn for non-zero value



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

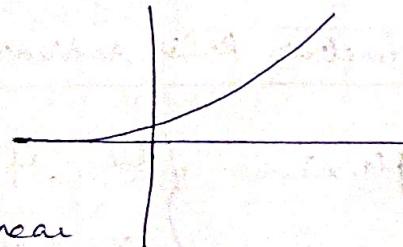
② Softmax activation function

$\in (0, \infty)$

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

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

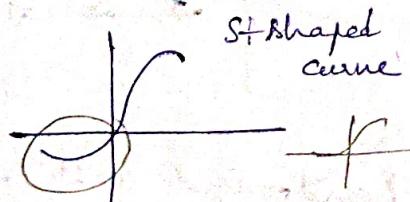
- used as a non-linear activation fn



Non-linear activation fn :

1) Sigmoidal fn :

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



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

$$x = 4$$

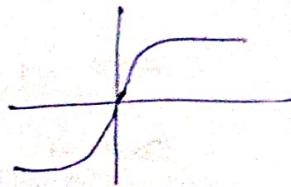
$$f(x) = \frac{1}{1 + e^{-4}} \approx 0.98$$

2) Tanh (Hyperbolic tangent activation fn)

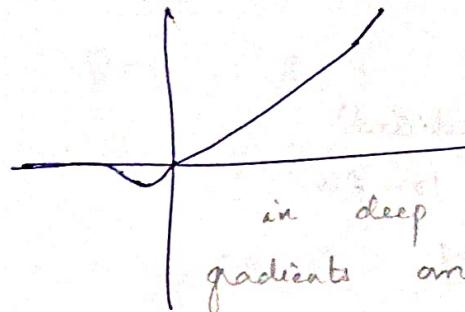
$$f(x) = \frac{e^x - e^{-x}}{1 + e^{-2x}}$$

$$= \frac{e^x - 1 - e^{-2x}}{4e^{-2x}}$$

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



3) Sigmoid activation fn:



$$f(x) = x + G(\beta x)$$

- used to improve learning
in deep neural N/w by providing smooth
gradients and avoiding dead neurons.

Softmax function for multi-class classification problem

4) Softmax fn.

Let inputs $\rightarrow x_i$ for $i = 1$ to n

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

output will be probabilistic value

$$x_i = [x_{i1}, x_{i2}, \dots]$$

$$f(x_1) = \frac{e^4}{e^4 + e^2 + e^1} \approx 0.8438$$

$$f(x_2) = \frac{e^2}{e^4 + e^2 + e^1} \approx 0.1142$$

$$f(x_3) = \frac{e^1}{e^4 + e^2 + e^1} \approx 0.042$$

Loss functions

Loss \rightarrow difference b/w Actual & predicted
 cost fn \rightarrow Average of all loss fns

$$\text{Loss} \rightarrow L(y, \hat{y})$$

$$\text{cost} \rightarrow \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i) \quad \dots \text{mean error}$$

1. Regression loss fn
2. classification loss fn

① Regression loss fn

$$\text{error} = y - \hat{y} \quad (\text{actual - predicted})$$

mean error = Average of loss fns

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Ex:

$$y = [5 \ 6 \ 8 \ 10 \ 12]$$

$$\hat{y} = [4 \ 7 \ 9 \ 10 \ 8]$$

y	\hat{y}	$y - \hat{y}$	$ y - \hat{y} $	$(y_i - \hat{y}_i)^2$
5	4	1	1	1
6	7	-1	1	1
8	9	-1	1	1
10	10	0	0	0
12	8	4	4	16
		—	—	19
		3	7	

$$ME = \frac{3}{5} = 0.6 \quad MAE = \frac{7}{5} = 1.4 \quad MSE = \frac{19}{5} = 3.8$$

MAE - sensitive to the outliers, not differentiable

MSE (Mean square error) $= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \dots$ differentiable

- easy to interpret

- not robust to outliers

$$RMSE = \sqrt{MSE}$$

SVM

- suitable for binary

classification

- Robust to noise

Huber loss:

$$\text{Huber}(y, f(x)) = \begin{cases} \frac{1}{2} (y - f(x))^2 & \text{for } |y - f(x)| \leq \delta \\ \delta(|y - f(x)| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

Avg Huber loss - Avg of huber loss over the sample

$$MAE \leq \text{Huber} \leq MSE$$

loss

Huber loss is a regression

that combines

- Mean Squared Error Loss

$$y=6 \quad \hat{y}=4 \quad \delta=4 \quad : \quad y - \hat{y} = 6 - 4 = 2 \quad \text{for small errors}$$

$$\text{Huber}(y, \hat{y}) \Rightarrow y - \hat{y} = 2 \leq 4$$

- Mean Absolute Error (MAE)

for large errors

$$\frac{1}{2} (6 - 4)^2 = \frac{1}{2} \times 2^2 = 2 \quad \text{It is robust to outliers}$$

Hinge loss $\in [-1, 1]$

Intuition:

$$\text{Hinge}(CD) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i \cdot \mathbf{w} \cdot \mathbf{x}_i)$$

\hookrightarrow mostly used for binary classification

Small errors \rightarrow behave like MSE (smooth learning)

Large errors \rightarrow behave like MAE (less sensitive to outliers)

$$y = [1, 1, -1, -1]$$

Goal: Be accurate but not

$$\hat{y} = [0.5, 1.5, 0.8, -1.4]$$

Goal: Be accurate but not disturbed by extreme values

y	\hat{y}	$y_h(y)$	$1 - y_h(y)$	0/p
1	0.5	0.5	0.5	0.5
1	1.5	1.5	-0.5	0
-1	0.8	-0.8	1.8	1.8
-1	-1.4	1.4	-0.4	0

$$\text{Hinge}(CD) = \frac{1}{4} \sum (0.5, 0, 1.8, 0)$$

$$= \frac{1.8 + 0.5}{4} = 0.575$$

Robust to outliers means that a method or loss function does not get affected by extreme or unusual data points

29/12/25

② Classification based Error functions

outcome \rightarrow categorical

Two types

- └ Binary cross entropy / log loss function
- └ categorical cross entropy

Binary cross entropy

$$L(y, \hat{y}) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$y=1, \hat{y}=0.3$$

$$L(1, 0.3) = -(1 \log 0.3 + (1-1) \log 0.3)$$

$$= -\log 0.3$$

$$= 1.7369 \quad 0.5229$$

Multivariate / categorical cross entropy function

$$L(y, \hat{y}) = -\sum_{i=1}^N y_i \log (\hat{y}_i)$$

$N \rightarrow$ total no. of possible outcomes

$y_i \rightarrow$ one hot vector of the actual value

$\hat{y}_i \rightarrow$ one hot vector of the predicted value

Consider three items red, green and blue

$$y(1 \ 0 \ 0) \quad \hat{y}(0 \ 1 \ 0) \quad b(0 \ 0 \ 1)$$

lets assume that possible o/p of the neural n/w is probability $\rightarrow p(0.5, 0.4, 0.1)$

Find the cross entropy for red

$$\begin{aligned} CE(p) &= -[1 * \log(0.5) + 0 * \log(0.4) + 0 * \log(0.1)] \\ &= -(-0.501) \\ &= 0.501 \end{aligned}$$

$$CE(g) = -(0 \cdot \log 0.5 + 1 \cdot \log 0.4 + 0 \cdot \log 0.1)$$

$$= 0.3979$$

$$CE(b) = -(0 \cdot \log 0.5 + 0 \cdot \log 0.4 + 1 \cdot \log 0.1)$$

$$= 1$$

Role of computational graph and chain rule

Computational graph

- Acyclic graph that shows the graphical representation of the mathematical model
- To show how the derivatives are calculated during the backpropagation process

Two types of nodes



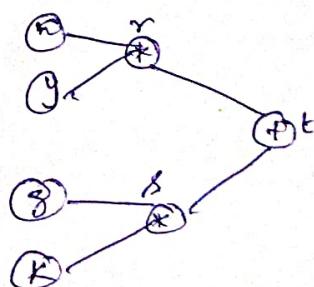
(ips, bias,
weights)

Consider a simple example:

$$f = xy + zk$$

$$x = xy \quad s = zk$$

$$f = s + t$$



chain rule:

- used to find derivative of the computational functions
- functions may be univariate or multivariate
- composite: combination of both

$$g(x, y, m, n, s)$$

$$u = g(f(x))$$

Chaining

↳ sequence of combination of the functions is called the chaining process

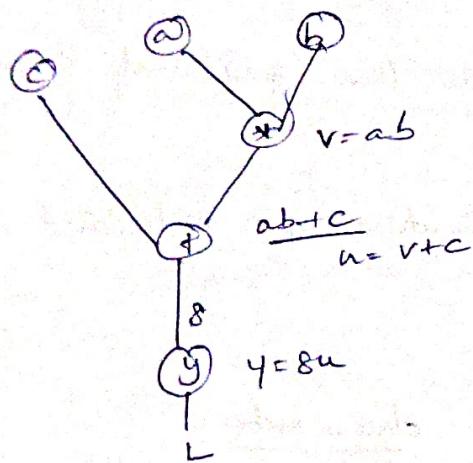
Consider derivative of $h(u) = u$
then, derivative $\{y(h(u))\}$?

$$\frac{\partial L}{\partial u} = \frac{\partial h}{\partial u} \times \frac{\partial u}{\partial u}$$

The partial derivative of the Composite function is product of partial derivative of a w.r.t x and derivative of a with respect to x .

Normal W/Fs.

$$f \circ (P^2(P^1(F(u))))$$



$$f(a, b, c) = (ab+c)8$$

$$1) V = ab$$

$$2) U = v + c$$

$$3) Y = 8u$$

$$\frac{\partial y}{\partial u} = 8 \quad \frac{\partial u}{\partial v} = 1 \quad \frac{\partial v}{\partial c} = 1$$

$$\frac{\partial y}{\partial v} = \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial v} = 8 \times 1 = 8$$

$$\frac{\partial y}{\partial c} = \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial c} = 8 \times 1 \times 1 = 8$$

$$\frac{\partial v}{\partial a} = b \quad \frac{\partial v}{\partial b} = a$$

$$\begin{aligned} \frac{\partial y}{\partial a} &= \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial a} \\ &= 8 \times 1 \times b = 8b \end{aligned}$$

$$\begin{aligned} \frac{\partial y}{\partial b} &= \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial b} \\ &= 8 \times 1 \times a = 8a \end{aligned}$$

Backpropagation Algorithm :

Two stages

{ Forward stage
Backward stage

- Algo:
1. Read the input
 2. Assign the input to the input layer

3. Compute all the neurons of hidden layer

$$\text{net}_j = w_0 + \sum_{i=1}^n w_{ij} o_i$$

$$o_j = \frac{1}{1+e^{-\text{net}_j}}$$

4. Compute the difference between the predicted o/p and the target o/p as an error

$$\text{o/p layer: } \delta_j = (t_j - o_j) o_j (1 - o_j)$$

$$\text{hidden layer: } \delta_k = o_j (1 - o_j) \sum \delta_k w_{jk}$$

$$\begin{aligned} \Delta w_{ij} &= \eta \delta_j o_i \\ w_{ij} &= w_{ij} + \Delta w_{ij} \end{aligned} \quad \left. \begin{array}{l} \text{update the weights} \\ \text{model parameters} \end{array} \right.$$

$$o_j' = o_j + \Delta o_j \rightarrow \text{update the bias}$$

optimization Algorithms

Problems
↓
two types

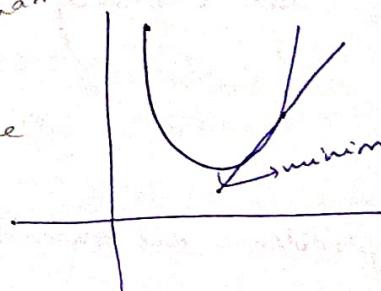
Constrained

unconstrained

optimization algo's search for critical points
because - Minimum loss occurs at a critical point
parameters - update naturally (loss as gradient)
- They indicate possible solution
model parameters
hyperparameters
parameters not learned
from the data, but
associated with the
algo's

$$\min f(x, \theta) \quad \xrightarrow{\text{model}} \text{model parameters}$$

$x \in X$ search space



$$f(x) = x^2 + 2x + 1$$

$$f'(x) = 2x + 2$$

$$2x + 2 = 0$$

$$2x = -2$$

$$x = -1$$

$$f''(x) = 2 > 0$$

\therefore the \rightarrow local minimum

$$f(-1) = (-1)^2 + 2(-1) + 1 = 0$$

Substitute the critical points in
the 2nd order derivatives

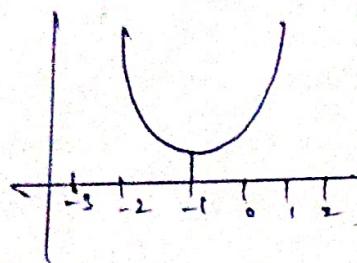
1) \therefore the \rightarrow critical point is

local minimum

2) \therefore the \rightarrow local maximum

3) \therefore \rightarrow inconclusive

Gradient descent algo



- an optimization algo used to minimize a loss (cost) function by iteratively updating model parameters in a direction that reduces the error, the most.

Multivariate function - $f(x_1, x_2, \dots, x_n)$

$$\nabla f(x) = \left(\frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_1}, \frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_2}, \dots, \frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_n} \right)$$

↙
first order derivative

↗ one shot vector

Gradient descent algorithm

- 1) stochastic gradient
- 2) Adaptive gradient
- 3) Adam optimizer

Optimization phase

A problem where we must choose the best solution from a set of possible solutions (Search space) by maximizing or minimizing an objective function, subject to certain constraints

features of optimizer

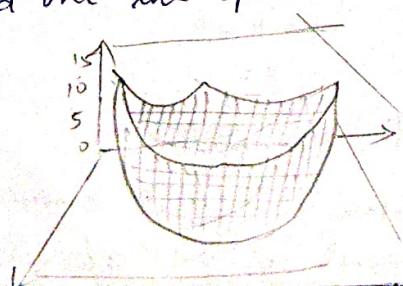
- ↳ minimizing the loss function
- ↳ efficiency
- ↳ handling of non-convex functions
- ↳ To provide a complex loss landscape
- ↳ Adaptation to the data
- ↳ To improve the convergence of the model

Decision Variables - The values we are allowed to change

(e.g) weights w , bias, learning rate

- Search space:
- Search space - Set of all possible values of the decision variables contains all the possible candidate solutions.
- ↳ optimization problems are framed as search space
 - ↳ All the candidate outputs are framed as search space
 - ↳ optimization problems will find out the optimal or best soln from the search space

Candidate take one possible
take from the
answer from the
search space



Two types of optimization problems

Convex

(have a clear boundary b/w the classes)

(eg) Linear regression,
SVM (with linear kernel)

hyper-parameter tuning

↳ tune the hyperparameters so that the model's performance can be improved

↓
values set before training, not learned from the data
having rate

(eg) learning rate

↳ it tells how much data that model uses before updating weights

Batch size - No. of samples used for weight update

Random search:

- ↳ Sample the input randomly from the distribution
- ↳ check for all combination of inputs, evaluate the model and check the results
- ↳ pick the best combination of inputs that generalizes the model well
- ↳ less memory required, can be done parallelly, thus computationally intensive.

Grid search algo:

↳ explore all possible combinations

↳ Take the hyperparameters with discrete set of values

↳ find all combinations of input

A(1,2) B(3,4)
 ↑ Parameters

↳ (1,2) (3,4)

(1,4) (2,2)

local minimum - the point where the function value is smaller than the nearby values

Global minimum - the point where the val is smallest over the entire search space decision boundary

(multiple minimum - need to find the global minimum)

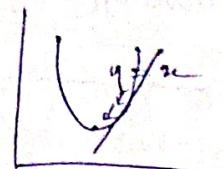
→ 1st order optimization algo

→ 2nd order optimization algo

(eg) Neural networks, k-means clustering

Tuning the model parameters

1. Gradient descent algorithm



↳ Meaning of gradient - slant/slope of the surface

↳ minimize the loss function

$$f(x_1, x_2) \quad \text{direction } \left(\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2} \right)$$

$$\text{magnitude } \sqrt{\left(\frac{\partial y}{\partial x_1}\right)^2 + \left(\frac{\partial y}{\partial x_2}\right)^2}$$

$$f(x_1, x_2) = x_1 + x_2$$

$$\text{gradient vector} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right) \rightarrow \text{direction of the gradient.} \\ = (1, 1)$$

$$\text{magnitude} = \sqrt{1+1} = \sqrt{2}$$

$$\text{dir} = \frac{\nabla f(x_{cur})}{\|\nabla f(x_{cur})\|} = \left(\frac{x_1}{\sqrt{2}}, \frac{x_2}{\sqrt{2}} \right)$$

represented with norm values

Steps

$$1. w_0$$

$$2. \Delta w_0 = w_0 - \frac{\partial f}{\partial w_0} \times \eta$$

$$3. w_{t+1} = w_t + \eta \frac{\partial f}{\partial w_t}$$

$$(eg) \quad x^2 + 2x + 1$$

$$x_0 = 1 \quad \eta = 0.1$$

$$f(x) = x^2 + 2x + 1$$

$$f'(x) = 2x$$

$$f''(x) = 2$$

$$x_{t+1} = x_t - \eta \frac{\partial f}{\partial x_t}$$

$$x_0 = x_0 - 0.1 \times \frac{\partial f}{\partial x_0}$$

$$x_1 = 1 - 0.1 \times (2(1) + 2)$$

$$x_1 = 1 - 0.1 \times (4)$$

$$x_1 = 1 - 0.4$$

$$x_1 = 0.6$$

$$x_2 = x_1 - \eta \frac{\partial f}{\partial x_1}$$

$$= 0.6 - 0.1(2(0.6) + 2)$$

$$= 0.6 - 0.1(3.2)$$

$$= 0.6 - 0.32$$

$$= 0.28$$

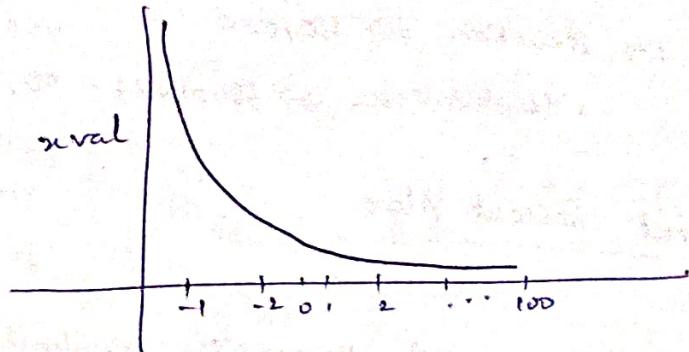
$$n_3 = x_2 - \eta \frac{\partial F}{\partial n_2} \quad f(n) = 2(0.28) + 2$$

$$= 0.28 - (0.1) \times 2.56$$

$$= 0.024$$

$$= 0.56 + 2 = 2.56$$

Loss functions keep on decreasing each iteration



- Batch Gradient descent
- 1) Take all training samples
 - 2) compute predictions for all samples
 - 3) calculate total cost for all samples
 - 4) compute gradient w.r.t. entire dataset
 - 5) update parameters & repeat until convergence
- model parameters are updated using the entire training dataset at once in each iteration

2. Batch gradient descent

algo

- Initialize the model parameters randomly and the learning rate
- Repeat until the iterations converges or number of iterations is reached
 - ↳ For entire dataset, compute the cost function

$$\nabla J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla J(\theta, x_i, y_i)$$

where m is the number of samples in the training set

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta)$$

- (2) - advantage and disadvantage of gradient descent algo

Advantages ↳ converge to a global minima is smoother compared to stochastic gradient descent

↳ stable convergence

↳ Convergence is guaranteed for convex functions not for non-convex functions

↳ suitable for small datasets

Disadvantages

- ↳ computationally intensive for large datasets
- ↳ memory requirement is higher.

Q) Consider 100 data points with five features

$$\text{computations} = 5 \times 100 = 500$$

$$100 \text{ iteration} \Rightarrow 50,000$$

$$1 \text{ iteration - 1 sec} \Rightarrow 50,000 \times 1 = 50,000 \text{ sec.}$$

Q) Stochastic gradient descent Algo :

- model updates its parameters after seeing one training sample at a time

Algo

- Initialize the model parameters randomly
- Initialize the learning rate randomly
- repeat until iterations converge or no. of iterations is reached

↳ shuffle the dataset

↳ Training sample randomly, compute gradient of the ^{loss} cost of the current point w.r.t model parameter

working - Exit model params

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$

- pick one training sample

- predict and compute loss

Advantages

↳ frequent updates

↳ memory requirement is less

↳ Adaptive for online learning

- update params
repeat for all samples
(one epoch)

- continue for multiple epochs

Disadvantages

↳ computationally very intensive

↳ data comes one by one

↳ model updates immediately after each sample

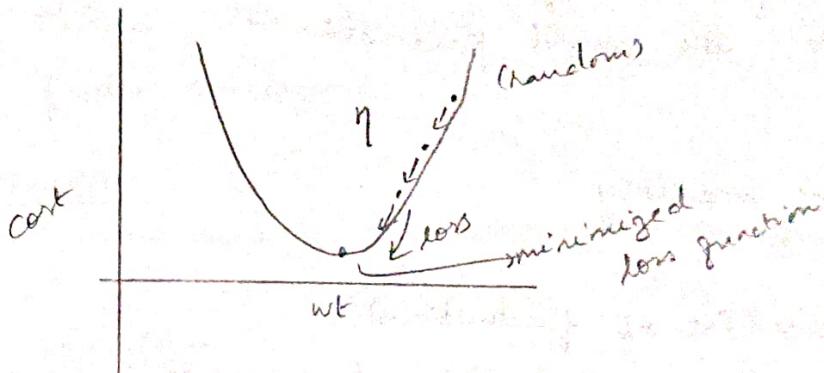
↳ No need to store all part data

05/02/2026

Steps in Gradient descent :-

1. Initialize all the model parameters randomly
2. Find out gradient of the loss fn with respect to each of the parameters
3. update ^{model} the parameters in order to minimize the loss function

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$



Gradient descent with Armijo Goldstein condition

→ To reduce the objective function and loss function using line search algorithm

$$f(x^{t+1} + \eta \nabla f(x^t)) - f(x^t) \geq c \eta \nabla f(x^t)^T$$

$c \in \{0, 1\}$

Gradient descent with full Armijo Goldstein condition

uses 2nd order gradient

$$f(x^{t+1} + \eta \nabla f(x^t)) - f(x^t) \geq c \cdot \eta (\nabla f(x^t))^T + \frac{1}{2} \eta^2 \| \nabla f(x^t) \|^2$$

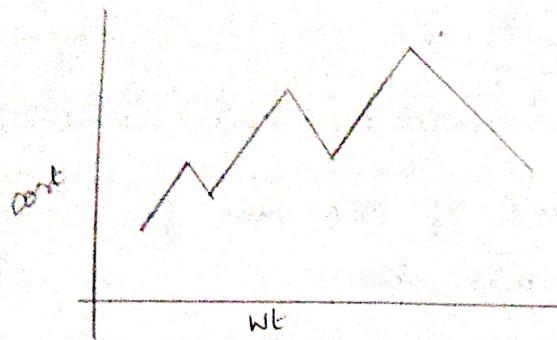
$$+ H(x) \cdot \nabla f(x^t)$$

↳ Hessian matrix

Variants of gradient descent Algo's

1. Stochastic gradient descent algo

- update model parameters after each and every sample



not a smooth curve,
lot of fluctuations

steps:

1. Select one training (e.g) randomly
2. compute gradient of the loss function
3. update the model parameters

Advantages

- less memory

disadvantages

- noise (lot of fluctuations)
- requires more iterations to converge

Mini Batch Gradient descent Algo

steps

- 1) Select the mini-batch with fixed size
- 2) Compute the mean gradient of the mini batch
- 3) update the model parameters
- 4) Repeat all the above steps for all remaining mini-batches.

↳ b/w stochastic and batch gradient descent

↳ requires medium amount of memory and less time to converge compared to SGD (stochastic Gradient Descent)

problem - it may stuck with local minima

Stochastic gradient descent with momentum

↳ helps to accelerate the convergence rate by smoothing out the noisy gradients of Stochastic gradient descent, thus reducing the fluctuations and improve the speed of

convergence compared to SGD Algo.

$$v_{t+1} = \beta v_t + (1-\beta) \nabla J(\theta_t) \quad \beta \in [0, 1]$$

where $v_t \rightarrow$ momentum at time t

$\beta \rightarrow$ momentum coefficient

$$\theta_{t+1} = \theta_t - \eta v_{t+1}$$

Advantages

- mitigate the oscillations
- Reduce the Variance
- faster convergence

disadvantages

- we have to tune momentum co-efficient β properly

Advance Optimizers:

1. AdaGrad (Adaptive Gradient descent algo)

- Adapts the learning rate for each and every model parameters based on the historical gradient info
- effective for the sparse (less-frequently used) features

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{g_t + \epsilon}} \nabla J(\theta_t)$$

where, $g_t \rightarrow$ sum of the squared gradients

$\epsilon \rightarrow$ small constant to avoid divide by zero error

Advantage

- adapt the learning rate for each of the parameter to improve the training efficiency

problem

- learning rate decays very quickly - so it leads to slow convergence of the model

2. RMS prop

- Improvement over the AdaGrad with decay factor to prevent the learning rate decreasing too rapidly.

$$E(g^2)_t = \gamma \cdot E(g^2)_{t-1} + (1-\gamma)(\nabla J(\theta_t))^2$$

\downarrow
Exponentially weighted Moving Average (EWMA)
of the square gradient
 $\gamma \in [0, 1]$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E(g^2)_t + \epsilon}} \nabla J(\theta_t)$$

Advantage

- prevent the excessive decay of learning rate

disadvantage

- Computationally intensive

3. Adam optimizer

- Adaptive gradient with momentum
- combines RMSprop with Momentum
- uses both first order and second order momentum of gradient to adapt learning rate for each parameter

Steps:

- update the first order momentum

$$m_t = \beta_1 \cdot m_{t-1} + (1-\beta_1) \cdot \nabla J(\theta_t)$$

- update the second order momentum

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla J(\theta_t))^2$$

- bias correction

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t} \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

- update the model parameters

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} + \hat{m}_t$$

where, β_1 and β_2 are the decay factors for the 1st order and 2nd order momentum respectively

Advantage - Faster convergence

disadvantage - Heavy computation, requires lot of memory

Comparison of optimizers:

optimizer	Advantages	disadvantages
SGD SGD	Simple and easy of implement	slow convergence
Mini Batch SGD SGD	Faster than SGD	Computationally expensive, stuck in local minima
SGD with momentum	faster convergence, reduced noise	Tune the β parameter properly
AdaGrad	Adaptive learning rate for all parameters	learning rate decays rapidly, slow convergence
RMS prop	Prevents fast decay of learning rate	Computationally intensive
Adam	fast convergence due to combining momentum and RMSprop	Memory intensive, computationally expensive.

Tuning the model params:

Steps in Gradient descent algo

- ↳ Init model params with random values
- ↳ Find the derivative of the univariate function $\frac{\partial f}{\partial x}$

$$f(x_1, x_2) = f(x_1, x_2) \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)$$

$$\text{magnitude (norm)} |\nabla f(x)| = \sqrt{\left(\frac{\partial f}{\partial x_1}\right)^2 + \left(\frac{\partial f}{\partial x_2}\right)^2}$$

(e.g) $f(x_1, x_2) = x_1 + x_2$

$$\nabla f(x_1, x_2) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right) = (1, 1)$$

$$|\nabla f(x)| = \sqrt{1^2 + 1^2} = \sqrt{2}$$

$$u = \frac{\nabla f(x_0, x_1)}{|\nabla f(x_0, x_1)|} = \frac{(1, 1)}{\sqrt{2}} = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)$$

direction of the gradient

1) w_0

$$2) \Delta w_0 = w_0 - \frac{\partial f}{\partial w_0}$$

$$3) w_{t+1} = w_t - \eta \Delta w_0$$

$$\textcircled{82} \quad f(x) = x^2 + 2x + 1$$

$$\eta = 0.1 \quad x_0 = 1$$

$$f'(x) = 2x + 2$$

$$x_1 = x_0 - \eta f'(x_0)$$

$$= 1 - (0.1)4$$

$$= 1 - 0.4$$

$$= 0.6$$

$$x_{t+1} = x_t - \eta f'(x_t)$$

$$f(x_0) = 2(1) + 2 = 4$$

$$f'(0.6) = 2(0.6) + 2$$

$$= 1.2 + 2$$

$$= 3.2$$

$$f'(0.28) = 2(0.28) + 2$$

$$= 0.56 + 2$$

$$= 2.56$$

$$x_2 = x_1 - \eta f'(x_1)$$

$$= 0.6 - 0.1 \times 3.2$$

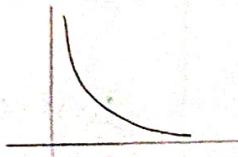
$$= 0.28$$

$$x_3 = x_2 - \eta f'(x_2)$$

$$= 0.28 - 0.1 \times 2.56$$

$$= 0.28 - 0.256$$

$$x_3 = 0.24$$



Batch gradient descent

1. θ

2. η , no. of iterations

3. Repeat until the model converges or the number of iterations reached

a) For entire dataset, compute cost J (θ)

$$\nabla J(\theta) = \frac{1}{m} \cdot \sum_{i=1}^m \nabla J(\theta, x_i, y_i)$$

b) update the model parameter

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$

c) return the optimal parameters

Advantage

- converge to global minima in a smoother way compared to stochastic gradient descent
- stable convergence
- suitable for small dataset

Disadvantage

- computationally intensive if dataset is large
- memory requirement will also be large if the dataset is large

Stochastic gradient descent algo

1. θ

2. η , no. of iterations

3. Repeat until the convergence or no. of iterations is reached

a) shuffle dataset

b) randomly choose a sample

c) $\nabla J(\theta, x, y)$

d) $\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$

Advantages

- Frequent updates support large datasets (as only one sample is randomly chosen)
- requires less memory
- Adaptive for online learning
- overcome the local minima by randomness

Disadvantage

- Due to randomness, noise is introduced, leads to noisy gradient leading to oscillations and fluctuations, leads to fail to reach the local minima.

Mini batch gradient descent:

1) θ , n_f , no. of iterations

2) no. of Batches (b), K Samples

3) Repeat until convergence or no. of iterations is reached

$$\nabla J(\theta) = \frac{1}{K} \sum_{i=1}^K \nabla J(\theta, x_i, y_i)$$

$\theta_{t+1} = \theta_t - n \nabla J(\theta)$

for each batch b

4) return the fine-tuned model parameters θ , end.

Advantage

- lies b/w SGD and batch GD

- steady and stable convergence

Disadvantage

- noise due to randomness

- no guarantee of convergence

Consider 10,000 data points

1) SGD = 10,000

2) BAD = 1

3) mini Batch(5) = $\frac{10000}{5} = 2000$

Concepts of Momentum

↳ Accelerate the model convergence by collecting the gradient in a variable Velocity

↳ Velocity as the smootherst Version of gradient

↳ Momentum is added with model parameter update

accelerate convergence speed, dampening oscillations, avoid the saddle points, efficient navigation of plateau, improves the generalization which leads to avoiding the overfitting

Two types

Traditional momentum

Nestronne momentum

Traditional : \rightarrow momentum coeff $\in [0, 1]$

$$v_t = \beta \cdot v_{t-1} + (1-\beta) \nabla J(\theta_{t-1})$$

(or)

$$v_{t+1} = \beta \cdot v_t + (1-\beta) \nabla J(\theta_t)$$

$$\theta_{t+1} = \theta_t - \eta v_{t+1}$$

Nestronne :

$$\text{future } v_{t+1} = \beta v_t + (1-\beta) \nabla J(\theta_t)$$

$$\text{past } v_{t-1} = \beta \cdot v_{t-2} + (1-\beta) \nabla J(\theta_t)$$

$$\theta_{t+1} = \theta_t - \eta (v_{t+1} + v_{t-1})$$

adjusting the model parameters based on the

past and future

Adaptive Gradient.

- Considering Square of Gradients in Normalized Variable for all model parameters

Step:

1) Init - θ, η, n , no. of iterations

2) Repeat until convergence or no. of iterations is reached

$$a) s = \sigma g \cdot g$$

$$b) \nabla \theta = \eta \times \frac{1}{\sqrt{s+c}} \cdot g \quad \left| \begin{array}{l} \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{s+c}} \nabla J(\theta_{t+1}) \\ s = \sigma g \cdot g \end{array} \right.$$

$$c) \theta = \theta - \Delta \theta$$

→ return the updated θ

Advantage - adapts l-r for each parameter

disadvantage - l-r decays quickly

Speciality of Adagrad - personalized l-r for each parameter

- adaptation is based on gradient history

- effective for sparse and high dimensional data

ob/ob/ob

RMS prop:

- ↳ improvement over the AdaGrad
- ↳ using decaying average or Moving average of partial derivative

RMS prop

- * uses decaying avg or moving avg of the partial derivative

- * Exponential weighted to moving that gives higher weightage to the ^{recent} ~~moving~~ updates rather than the past history

$$* y_0 = x_0$$

$$y_n = \rho x_n + (1-\rho) y_{n-1}$$

$\rho \rightarrow$ decay factor

when $\rho=1$, $y_n = x_n$

$\rho=0$, $y_n = y_{n-1}$

$$\rho \in [0, 1]$$

$$x_0 \quad x_1 \quad x_2 \quad x_3 \quad x_4$$

$$5 \quad 10 \quad 12 \quad 14 \quad 16$$

$$\rho = 0.2$$

$$y_0 = x_0 = 5$$

$$y_1 = \rho x_1 + (1-\rho) y_0$$

$$= 0.2 \times 10 + (1-0.2) \times 5$$

$$= 0.2 \times 10 + 0.8 \times 5$$

$$= 2 + 4 = 6$$

AdaGrad

- * Sum of square gradients

- * Square of partial derivatives grow monotonically, hence leads to slow rate of convergence

$$y_2 = 0.2 \times 12 + (1-0.2) \times 6 = 7.2$$

$$y_3 = 0.2 \times 14 + (1-0.2) \times 7.2 = 8.56$$

algo

1) $\theta, \eta, \text{ no. of iterations, } \epsilon$

2) Repeat until the model converges or the no. of iterations is reached

a) SWMA of the gradient

$$\hat{g} = \rho \hat{g} + (1-\rho) g \cdot g$$

b) update the parameters

$$\Delta \theta = \eta \times \frac{1}{\sqrt{\text{ote}}} \cdot g$$

$$E(g^2)^t = \gamma \cdot E(g^2)^{t-1} + (1-\gamma) (\nabla J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \eta \sqrt{\frac{1}{E(g^2)^t + \epsilon}} \cdot t \cdot \nabla J(\theta_t)$$

c) Apply update

$$\theta = \theta - \Delta \theta$$

Specificity

↳ Adaptive learning for each parameter

↳ Reduction in oscillation towards minimum so convergence is stable

↳ Computationally very effective

Adam optimizer (Adaptive Moment Estimation):

↳ This combines both gradient descent with momentum (RMS prop)

↳ Adaptive lr for each parameter because of RMS prop

↳ Bias correction to remove the noise

$$m_t = \beta_1 m_t + (1-\beta_1) \nabla J(\theta_t)$$

$$v_t = \beta_2 v_t + (1-\beta_2) (\nabla J(\theta_t))^2$$

2) Bias correction:

$$\hat{m}_t = \frac{m_t}{1-\beta_1 t}$$

Variance correction:

$$\hat{v}_t = \frac{v_t}{1-\beta_2 t}$$

3) Model parameter update

$$\theta_{t+1} = \theta_t - \frac{\eta \cdot m}{\sqrt{v + \epsilon}}$$

Algo:

1) Init: θ , η and no. of iterations, ϵ , p_1, p_2

2) Repeat until convergence or the no. of iterations is reached

$$a) m = p_1 \cdot m + (1-p_1) \cdot g$$

$$b) v = p_2 \cdot v + (1-p_2) \cdot g^T \cdot g$$

$$c) \hat{m} = \frac{m}{1-p_1 \epsilon} \quad \hat{v} = \frac{v}{1-p_2 \epsilon}$$

d) update

$$\Delta \theta = \eta \cdot \frac{1}{\sqrt{\hat{v} + \epsilon}} \cdot \hat{m}$$

$$e) \theta = \theta - \Delta \theta$$

Second order optimizer

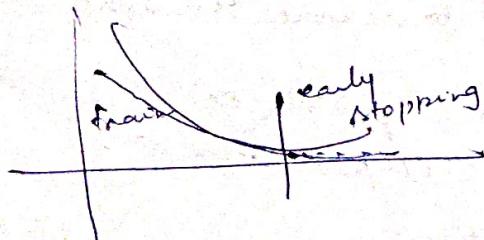
Newton's method

Quasi
Newton's
method

Regularization

↳ done in order to prevent overfitting (works well in training, but poor performance in test sample)

regularization and Generalization

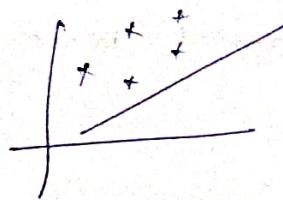


Error - diff b/w actual & predicted

Bias - Avg of all errors (cost)

Variance - 2nd order derivative \rightarrow inconsistency and erratic predictions because of fluctuations
Generalization error \rightarrow bias + variance

underfitting:



Condition

1) high bias
and high
variance

model

underfitting
useless



Remarks

2) high bias
and low
variance

underfitting

3) low bias
and high
variance

overfitting

4) low bias
and low
variance

Ideal

Difference b/w overfitting and underfitting :

underfitting

* TOO simplistic to
capture real time
scenarios and failed
to capture the
underlying patterns
in data

overfitting

* Very Complex, fits the
data as well as
the noise and fluctuation
present in the data

* characterized by high
bias and low variance

* characterized by low bias
and high variance

* performance is poor
during training and testing

* performance good in
training, poor on testing

* poor generalization

* poor generalization

M_1, M_2 - models

↓
Single complex

↓

we can choose complex model M_2

provided $\text{performance}(M_2) > \text{performance}(M_1)$

$$M_1 = 2x^2 + 3$$

$$M_2 = 6x^4 + 7x^2 + 2x^2 + 3$$

$$L_1 = |x_1| + |x_2| + |x_3| + \dots + |x_n|$$

norm:

$$M_1 = 2$$

$$M_2 = 6 + 7 + 2 = 15$$

$$\text{L}_2 \text{ norm: } \|x\|^2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

$$M_1 = \sqrt{4} = 2$$

$$M_2 = \sqrt{6^2 + 7^2 + 2^2}$$

$$L_\infty = |x_k|$$

$$= \max(|x_1|, |x_2|, \dots, |x_n|)$$

$$L_p = \max_i \left(\sum_{i=p}^n |x_i|^p \right)^{1/p}$$

Regularization ways

↳ Neural N/w Initialization

↳ Neural N/w Configuration

↳ Neural N/w Regularization

Weight Initialization

↳ zero weight initialization

↳ Random Initialization $\in (-0.5, 0.5)$

↳ Xavier (or) Glorot Initialization

$$w \approx u \left[\frac{-\sqrt{6}}{2 \cdot f_{\text{ain,in}} + f_{\text{out}}} , \frac{+\sqrt{6}}{2 \cdot f_{\text{ain,in}} + f_{\text{out}}} \right]$$

$$N(0,1) \sim \frac{\sqrt{1}}{f_{\text{ain,in}}}$$

$$\sigma = \frac{2}{\sqrt{f_{\text{ain,in}} + f_{\text{out}}}}$$

4) He uniform distribution

$$w \approx u \left[-\sqrt{\frac{6}{f_{\text{ain,in}}}} + \sqrt{\frac{6}{f_{\text{out}}}} \right]$$

5) LeCun Initialization

- uses gaussian distribution

$$\sigma = \sqrt{\frac{3}{f_{\text{ain,in}} + f_{\text{out}}}}$$

width - no. of neurons in each layer (maximal no. considered)

depth - no. of hidden layers in the model

$$N_i \leq N_h \leq N_o$$

$$N_h \leq 2N_o$$

$$N_h = \frac{N_o}{N_i} (\alpha N_i + N_o)$$