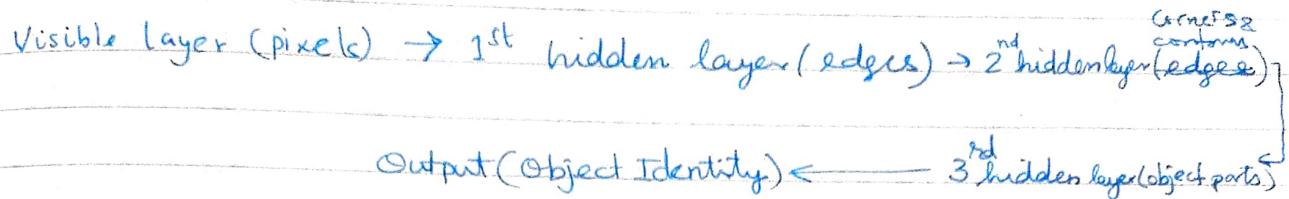


DL Defn: Computational Graph

NN evolved in the year 1952.

1952 - SGD

1958 - Perceptron

- Learnable weights

Block region of AI (N-L Problem - XOR)

1986 - Backpropagation

- Multi-Layer Perceptron

1995 - Deep convolutional NN

- Digit Recognition.

Issues at that time: Database availability and computationally efficient machines.

1. Big Data

- Larger datasets
- Easier collection & storage

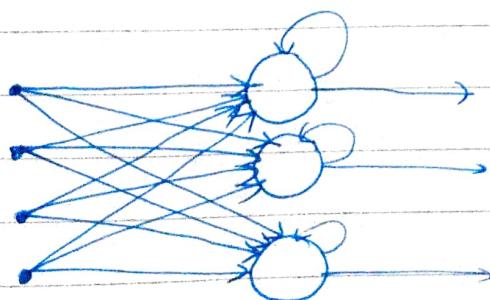
2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable

3. Software

- Improved Techniques
- New models
- Tool boxes

Recurrent Neural Network:



Eg:- LSTM (Long Short-Term Memory)

Gated Recurrent Unit (GRU)

Gradient Descent (GD)

→ Iterative, first-order, optimization algorithm

5. Function requirements:

→ Differentiable

$$w \leftarrow w - \eta (\nabla f(w)) \rightarrow \text{differentiable}$$

→ convex.

↳ line segment connecting two segments fn.'s points lies on or above the curve.

$$f(\lambda x_1 + (1-\lambda)x_2) =$$

mathematical way to f': second derivative > 0 .

$$f(x) = x^2 - x + 3$$

$$f'(x) = 2x - 1$$

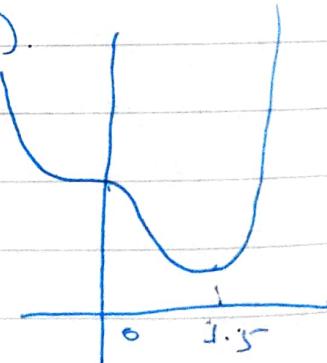
$$f''(x) = 2 > 0 \text{ (convex fn.)}$$

→ Saddle Points

→ Quasi-convex functions (minimax points).

$$20. \quad \underline{y} := f(x) = x^4 - 2x^3 + 2.$$

$$\begin{aligned} \frac{d}{dx}(f(x)) &= 4x^3 - 6x^2 \\ &= 4x^2(x - \underline{1.5}) \end{aligned}$$



$$25. \quad \frac{d^2}{dx^2}(f(x)) = 12x^2 - 12x = 12x(x - 1) \quad \underline{x=0}, \underline{f''(x)=0}.$$

$$\underline{x=1.5}, \underline{f''(x)=9}.$$

$x < 0$: convex

$0 < x < 1$: concave.

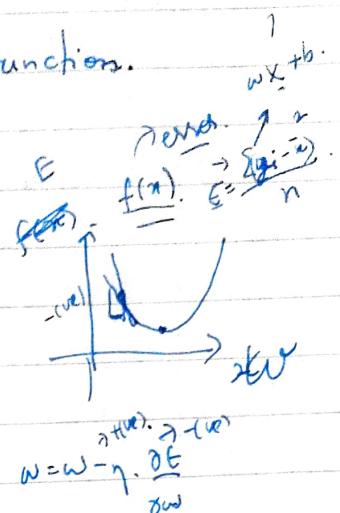
$x > 1$: convex.

Intuitively, gradient is the slope of the function at a particular direction.

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix} \text{ for } n\text{-dimensional } \cancel{\text{fun}} \text{ functions.}$$

$$w = w - \eta \cdot \frac{\partial E}{\partial w} \quad \rightarrow \text{opp. direction to max. change.}$$

$$\frac{\partial E}{\partial w} =$$



Algorithm Gradient Descent :

Basic / Vanilla Gradient Descent

13/02/2024

1) Initialize random weight, w

2) Go until converge

3) $w \leftarrow w - \eta \cdot \frac{\partial E}{\partial w}$

4) Until error is minimized.

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial J(w)}{\partial w} \quad \boxed{\text{criterion}}$$

5) Return ' w '

Stochastic Gradient Descent (SGD)

1) Initialize w

2) Pick sample 'i' out of 'n' samples

3) Do

4) $w \leftarrow w - \eta \cdot \frac{\partial J_i(w)}{\partial w}$

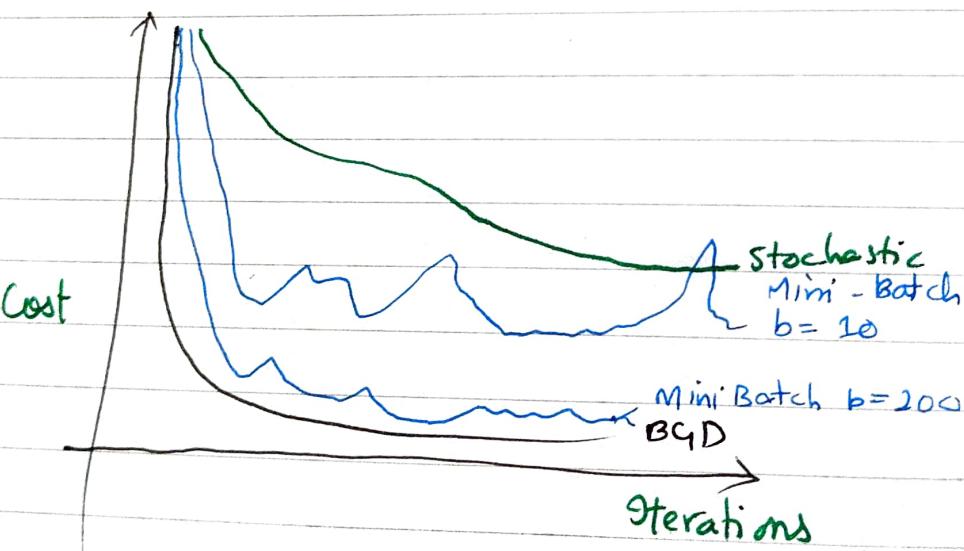
5) Until $J(w)$ is minimized.

6) Return ' w '

Mini Batch Gradient Descent

- 1) Initialize w'
 - 2) Pick some sample of batch 'B'
 - 3) ~~do~~
 - 4) $w \leftarrow w - \eta \cdot \sum_{i=1}^B \frac{\partial J_i(w)}{\partial w}$
 - 5) Until $J(w)$ is minimized
 - 6) Return w'
- 10) BGD : $\boxed{\frac{\partial J(w)}{\partial w}}$ → computationally intensive to compute
- SGD: Faster to compute gradient, but noisy.

Error minimization with GD :-



LOGISTIC REGRESSION :

$$y = mx + c.$$

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$x = x_1, x_2, \dots$$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$\theta_1 = m, \theta_2 = c$.

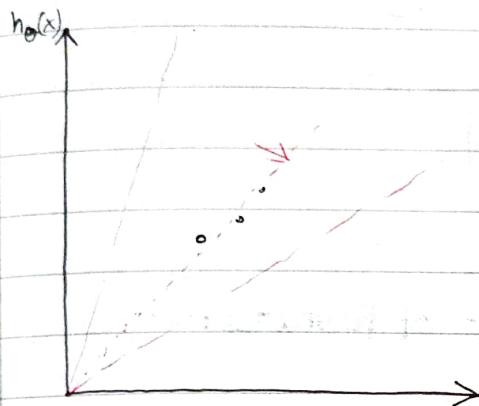
X = given F.V.

\hat{y} = Predicted O/P

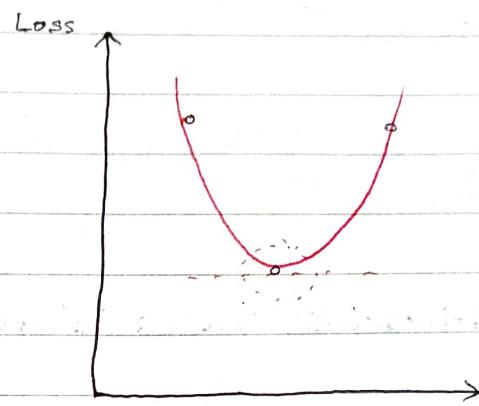
y = Actual O/P

n = no. of samples.

$h_{\theta}(x) = \theta_0 + \theta_1 x \Rightarrow$ univariate linear regression.



Hypothesis function θ



14/02/2024

Cost Function θ

$$y = mx + c$$

$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

x - Input

y - True label

$h_{\theta}(x)$ - Predicted value

$$\theta = ? \mid h_{\theta}(x) = y$$

$$\frac{\partial E}{\partial \theta} = \frac{\partial}{\partial \theta} E(\theta_0, \theta_1) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^i) - y^i)^2$$

$$= \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^i) - y^i)^2$$

$$h_{\theta}(x_i) = \theta_0 + \theta_1 x_i$$

$$\frac{\partial E}{\partial \theta_0} = \frac{2}{2} \sum_{i=1}^n (h_\theta(x^i) - y^i) \cdot \frac{\partial (h_\theta(x^i) - y^i)}{\partial \theta_0} \quad \text{True label}$$

$$\boxed{\frac{\partial E}{\partial \theta_0} = \sum_{i=1}^n (h_\theta(x^i) - y^i)}$$

$$\frac{\partial E}{\partial \theta_1} = \frac{2}{2} \sum_{i=1}^n (h_\theta(x^i) - y^i) \cdot \frac{\partial (\theta_0 + \theta_1 x^i - y^i)}{\partial \theta_1}$$

$$\boxed{\frac{\partial E}{\partial \theta_1} = \sum_{i=1}^n (h_\theta(x^i) - y^i), x^i}$$

Multi-variate case: 'd' number of features:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

Use vectorized approach: Compute $\vec{\theta}^T \vec{x}$.

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x_0 \\ 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

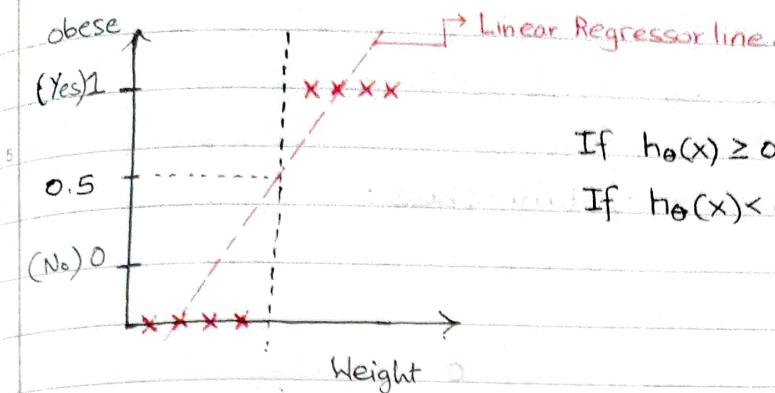
$$\vec{\theta}^T \vec{x} = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_d \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = h_\theta(x).$$

Gradient Descent Algorithm:

Repeat until convergence!

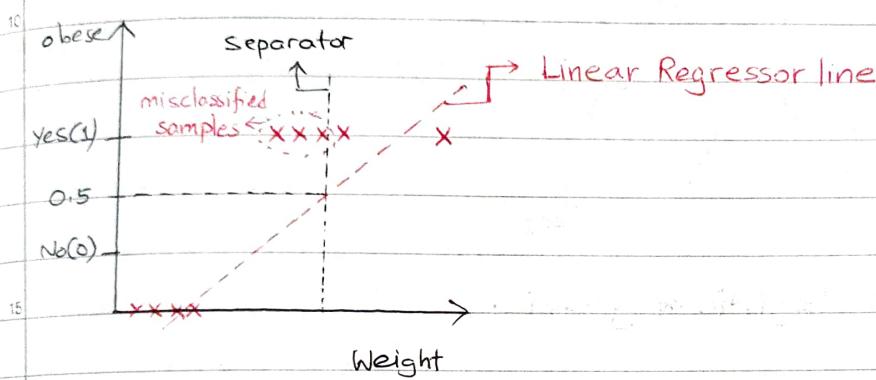
$$\left. \theta_i := \theta_i - \alpha \cdot \frac{\partial}{\partial \theta_i} (E(\theta_0, \theta_1, \dots, \theta_d)) \right\}, \alpha \leq i \leq d, i \in \mathbb{Z}^+$$

Issue with linear regression: outliers



If $h_0(x) \geq 0.5$, predict $y=1$.

If $h_0(x) < 0.5$, predict $y=0$.



Classification $y=0$ or 1 .

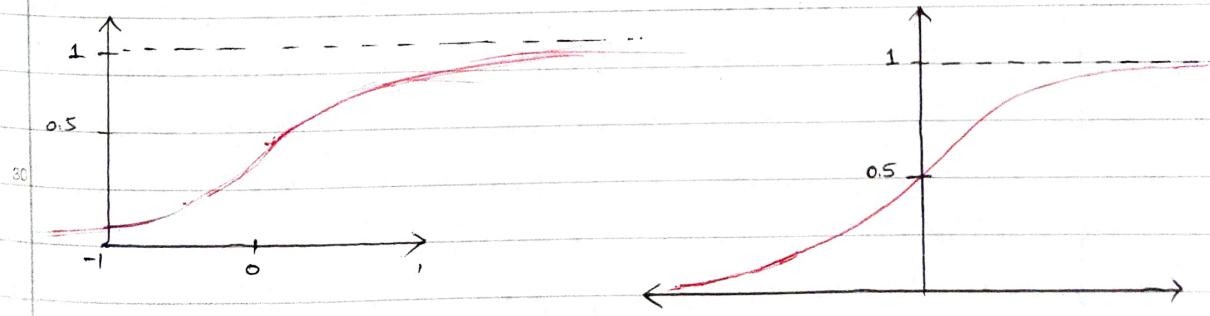
But here $h_0(x)$ can be >0 or <1 .

So, Linear Regression is not very suitable for classification.

Hypothesis of Logistic Regression: $0 \leq h_0(x) \leq 1$.

Previously, $h_0(x) = \phi^T x$.

Now, $h_0(x) = g(h_0(x))$ and $g(z) = \frac{1}{1+e^{-z}}$

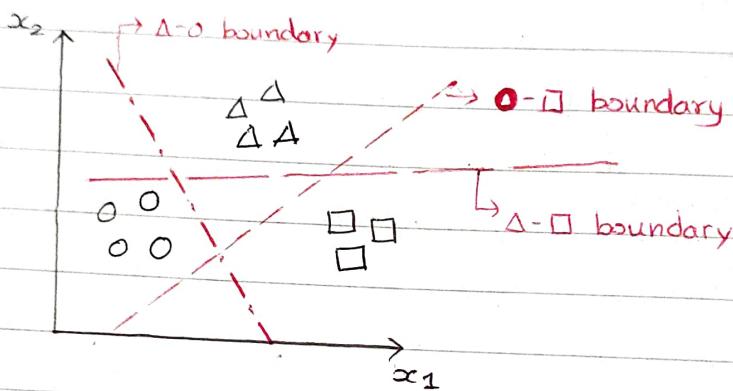


Trained weight : is also known as "Frozen weight".

For testing, compute $P(y=1|x, \theta)$

$$\text{Also, } P(y=0|x, \theta) = 1 - P(y=1|x, \theta)$$

Logistic Regression : Multi-class :



Linear vs Logistic Regression

Linear Regression

Logistic Regression.

1) Continuous variables

How many features (N) = 5
 Target (C) = $\{0, 1\}$ # = 2.

Confusion Matrix : where the model gets confused.

		Actual	
		T	F
Predicted	T	TP	FP
	F	FN	TN

$$\text{Accuracy} = \frac{\text{Total correctly classified samples}}{\text{Total samples}}$$

$$= \frac{(TP + TN)}{(TP + TN + FP + FN)}$$

Spam Detection: Precision is an important metric.

TN : 570

TP : 350

FP : 50

FN : 30

$$\text{Accuracy} = \frac{570 + 350}{920}$$

$$\text{Precision} = \frac{350}{350 + 50}$$

$$\text{Recall} = \frac{350}{350 + 30}$$

Multi-class :

Three types of averaging

macro:

compute for each class and take weighted average of them.

		Predicted		
		T	F	
		0	1	$\rightarrow EC_1$
Actual	T	(90)	(10)	$TP_0 = 90$
	F	(30)	(20)	$TP_1 = 20$
		(FP)		$FN_0 = 10$
		+ (P)	F	$FN_1 = 30$
<u>Class I:</u>		T		
(A)		F		

10. $P_0 = \frac{90}{100}, \cancel{\frac{P_0 = 20}{100}} = \underline{0.9}$

$P_1 = \frac{20}{30+20} = 0.4$

15. $P_{macro} = \frac{0.9 + 0.4}{2} = \underline{0.65}$

$P_{weighted} = P_0 \times \text{Support}_0 + P_1 \times \text{Support}_1$

$$= 0.9 \times \frac{100}{150} + 0.4 \times \frac{50}{150}$$

20. (i) ROC (ii) Cohen's kappa (iii) quad. weight kappa

$TP_A = 10, TP_B = 15, TP_C = 8, TP_D = 5$.

25. $FP_A = 3, FP_B = 4, FP_C = 2, FP_D = 1$.

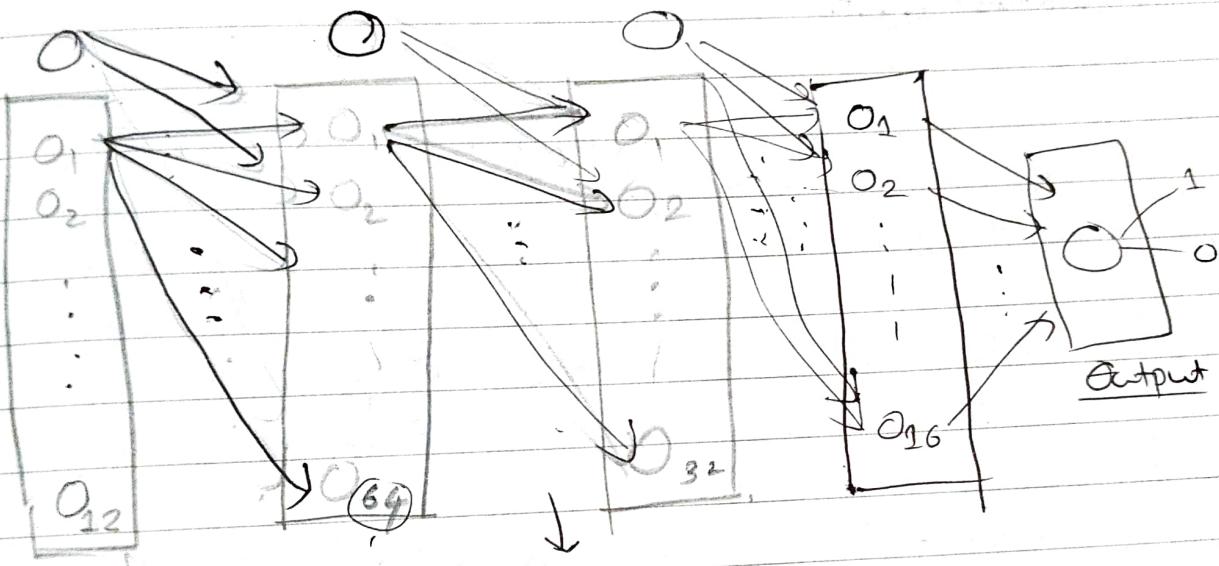
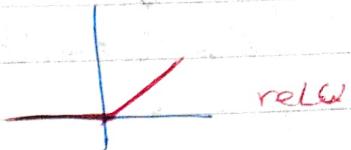
$FN_A = 3, FN_B = 3, FN_C = 2, FN_D = 2$.

30. $TN_A = 32, TN_B = 25, TN_C = 35, TN_D = 40$.

Hyperparameters - values are predetermined prior to start of learning process.

Eg:-

- 1) Learning Rate \rightarrow Comp. complexity \uparrow
- 2) Batch size \rightarrow Small \rightarrow Time complexity \uparrow . for execution.
- 3) epochs
- 4) Activation fn.
- 5) No. of layers / No. of neurons.



Trainable Parameters

weight / biases
 $= 12 \times 64 + 1 \times 64$
 $= 832$

$64 \times 32 + 1 \times 32$
 $= 2080$

528 17

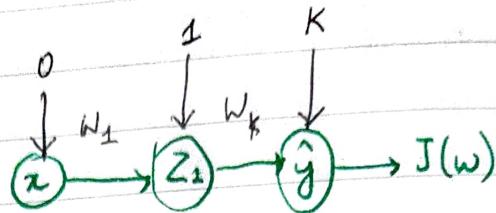
Vanishing and Exploding gradient Problem:-

activation function : non-linear, differentiable

5 deep learning

↓
Deep Neural Network

↓
100 - 150 layers



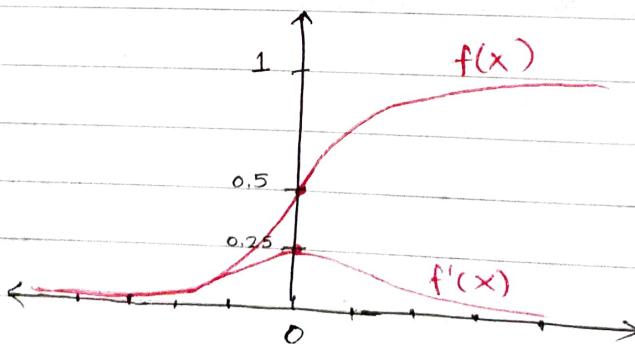
10

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial E}{\partial w}$$

$$15 \quad \frac{\partial J(w)}{\partial w_1} = \frac{\partial J(w)}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

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

20



25

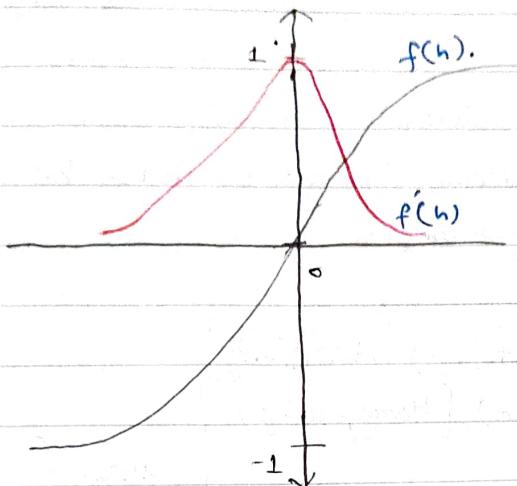
Maximum value of derivative of sigmoidal function is

$$30 \quad \text{say } \frac{\partial L}{\partial w} = 0.015 \quad (0.2 \times 0.15 \times 0.05)$$

$$w_{\text{new}} = w_{\text{old}} - (1) \cdot 0.015 = 2.4985$$

→ The gradient of activation function is small, that's why the issue occurs.

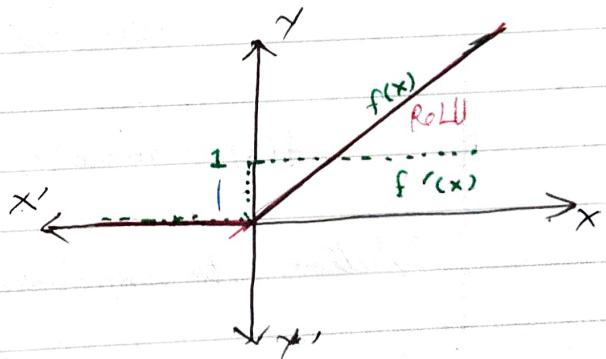
→ tanh:



Again, it can result in vanishing gradient problem.

→ ReLU (Rectified Linear Unit), Leaky ReLU, Parametric ReLU.
(saturates in one direction).

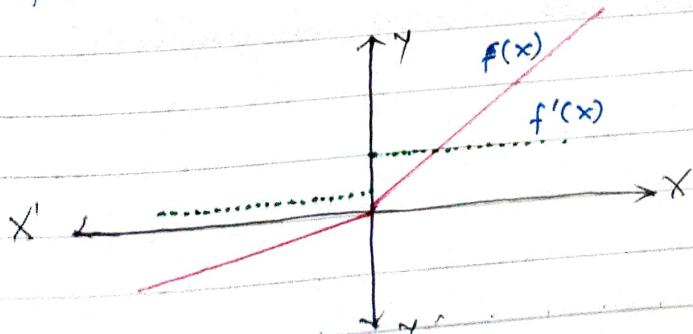
$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$



If $f(x) \leq 0$, then, $w_{old} = w_{new}$.

This condition is called as dead neuron.

Modification done to solve dead neuron problem:- Leaky ReLU.
(not vanishing gradient)



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

Date : 19/02/2024

$$O_1 = W_k \cdot Z_1$$

$$\hat{y} = f(W_k \cdot Z_1)$$

sigmoid

→ keeps oscillating, doesn't converge.

→ the weight vector changes too much because of large value of weight vector.

$$f(x+y) = f(x) + f(y) \quad [\text{Additivity}]$$

$$f(ax) = af(x). \quad [\text{Homogeneity}]$$

$$f(ax+by) = af(x) + bf(y).$$

Linear function

20/02/2024

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

⇒ step fn. ⇒ not continuous/differentiable

(Also gradient = 0 since value is constant in the two intervals.)

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

Solution to exploding gradient :

gradient clipping :

If $\|g\| > \text{threshold}$,

$$g \leftarrow \frac{\text{threshold} \times g}{\|g\|}$$

[same direction but lesser magnitude].

In leaky ReLU, $\alpha = 0.01$,

If $\alpha \neq 0.01$, then it is parameterized ReLU.

Heaviside step fn. $f(n)$

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

$f'(x) = \delta(x)$, where $\delta(x)$ is the dirac delta function.

ReLU function

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

$$f'(x) = \begin{cases} 0, & x < 0 \\ \text{undefined}, & x=0 \\ 1, & x > 0 \end{cases}$$

For practical purposes, $f'(x)=0$ or 1 at $x=0$ or
subgradients are used (value in between 0 and 1).

Activation Functions

The purpose of activation fn. is to add non-linearity of the neuron.

Characteristics of activation functions:

- Continuous and differentiable
- Non-decreasing function
- Not always constant

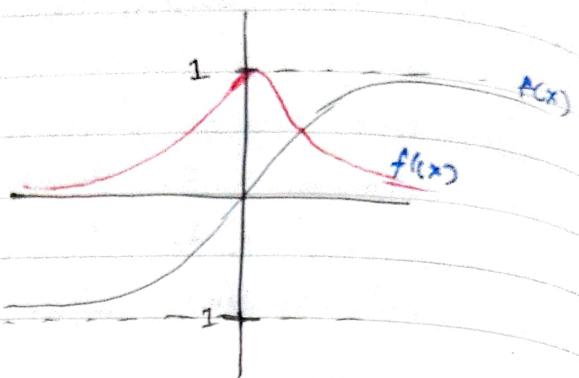
Sigmoid - used in o/p layer for binary classification.

→ Hyperbolic tangent

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

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

Vanishing gradient.



→ Inverse tangent

$$f(x) = \tan^{-1}(x)$$

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

→ ReLU

Most popular activation function since 2017.

It can result in dead neurons.

→ Leaky ReLU

Fixes the problem of dead neuron
Used only in hidden layers of NN.

$$f(x) = \begin{cases} 0.01x, & x < 0 \\ x, & x \geq 0 \end{cases}$$

→ Parameterized ReLU

$$f(x) = \begin{cases} dx, & x < 0 \\ x, & x \geq 0 \end{cases}, \quad d \neq 0.01$$

Maxout function

$$f(\vec{x}) = \max_i x_i$$

$$\frac{\partial f}{\partial x_j} = \begin{cases} 1 & \text{for } j = \arg \max_i x_i \\ 0 & \text{for } j \neq \arg \max_i x_i \end{cases}$$

$[0.5, 0.7, 0.1] \rightarrow [0, 1]$ → One Hot vector.

Exponential Linear Unit (ELU)

Function

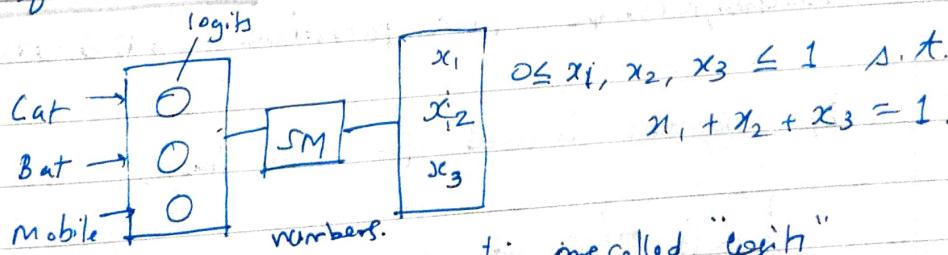
$$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \rightarrow \text{O/P is zero centered} \\ x & \text{for } x \geq 0 \quad (\text{Mean} = 0). \end{cases}$$

$$f'(\alpha, x) = \begin{cases} f(\alpha, x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Identity function - for O/P layer

Use: Regression.

Softmax function



$$f_i(x) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \quad \text{for } i = 1, \dots, n.$$

(Σ+) → $\frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$, form classes.

$$f'(x) = 1 - \sum_{i=1}^n f_i(x)$$

Regression Model (0-1) (S.I.) Classification Model (S.I.)
 ↳ Prediction (\hat{y}) ↳ classify/assign discrete value for the sample

→ Loss function.

$$E = (y - \hat{y})$$

$$(L_1) \text{MSE} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

$$(L_2) \text{MAE} = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

→ Huber Loss function

→ Loss function

Cross Entropy Loss

KL-divergence: $\mathbb{E}[\log(\mathbb{P}_{\text{true}})]$

Bhattacharya (dissimilarity)

Cross Entropy Loss:

$$- \left(\sum_{i=1}^c y_i \log(\hat{y}_i) \right)$$

CEEL (Cross Entropy Loss)

BEL (Binary Cross Entropy Loss)

Hinge loss function

→ Binary classification,

$$- \left(\sum_{i=1}^c y_i \log(\hat{y}_i) \right)$$

$$- \left(\sum_{i=1}^2 y_i \log(\hat{y}_i) \right) = - (y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2))$$

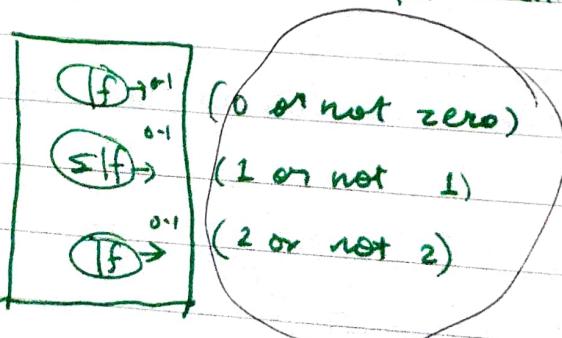
$$= - (y_1 \log(\hat{y}_1) + (1-y_1) \log(1-\hat{y}_1))$$

$$\boxed{\text{BCE} = - (y_1 \log(\hat{y}_1) + (1-y_1) \log(1-\hat{y}_1))}$$

Robust outlier - Less Penalize - MAE

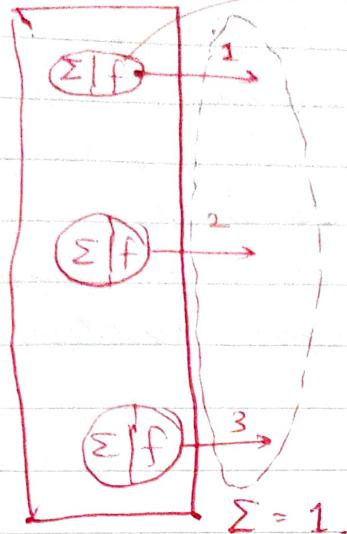
Misclassified Sample - More Penalize - MSE - More quick weight update.

If sigmoid is used for Multi-class,



$\sum \neq 1$ always.

If softmax is used \rightarrow s.m. (softmax).

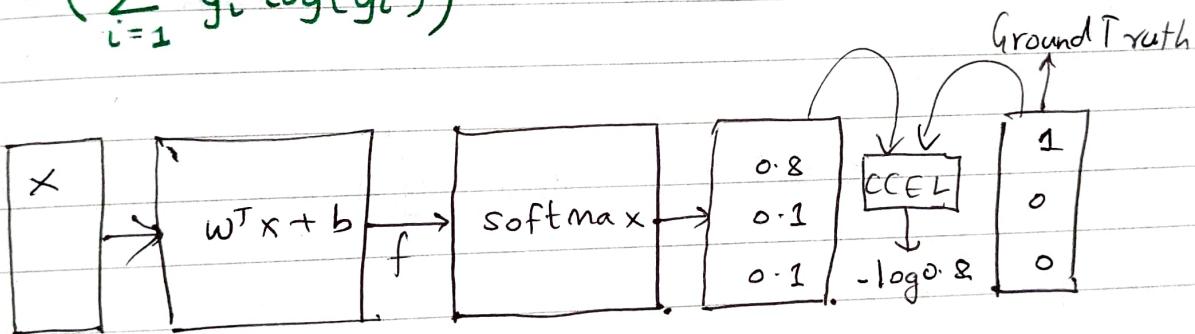


$$\frac{e^{f_i}}{\sum_{j=1}^C e^{f_j}} \quad (0-1).$$

$$\sum_{i=1}^C p_i = 1.$$

CCEL (Categorical Cross Entropy Loss).

$$- \left(\sum_{i=1}^C y_i \log(\hat{y}_i) \right)$$



Why cross entropy is used most commonly as compared to MSE?

Optimization techniques

↳ Gradient Descent \Rightarrow

Optimization Tech./Algo. to max./min.
the criterion function $J(w)$

Loss function

$$\rightarrow L/E = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{MSE}$$

$$\rightarrow L/E = \sum_{i=1}^n |y_i - \hat{y}_i| \quad \text{MAE}$$

$$\rightarrow L/E = \sum_{i=1}^n y_i \log(\hat{y}_i)$$

$$\frac{\partial L}{\partial w}$$

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial L}{\partial w}$$

\rightarrow Batch G.D. (Computationally complex for large I/P size)

\rightarrow Stochastic G.D. (Longer Time to converge)  oscillating towards solution

\rightarrow MiniBatch G.D.

BGD - 1 epoch = 1 iteration

SGD - 1 epoch = 10k iterations

MBGD - 1000 - Mini Batch Size

1 epoch - $\frac{10k}{1000} = 10$ iterations

Assume

I/P size - 10k

GD optimization with momentum:-

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial L}{\partial w}$$

$$w_t \leftarrow w_{t-1} - \eta \cdot \frac{\beta \frac{\partial L}{\partial w}}{1 + \beta \frac{\partial L}{\partial w_{t-1}}} \quad \text{exponential weighted average.}$$

$$\begin{matrix} t_1 & t_2 & \dots & t_n \\ a_1 & a_2 & \dots & a_n \end{matrix}$$

$$v_{t1} = a_1$$

$$v_{t2} = \beta \cdot v_{t1} + (1-\beta) a_2. \quad \text{usually, } \beta = 0.95.$$

More weight to previous than current if $\beta = 0.95$.

$$v_{t3} = \beta \cdot v_{t2} + (1-\beta) a_3$$

$$= \beta \cdot (\beta \cdot v_{t1} + (1-\beta) a_2) + (1-\beta) a_3$$

$$v_{t3} = \beta^2 v_{t1} + \beta(1-\beta) a_2 + (1-\beta) a_3$$

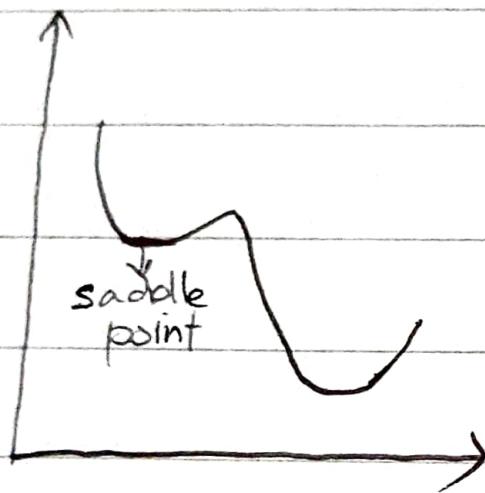
$$v_{t3} = \beta^2 a_1 + \beta(1-\beta) a_2 + (1-\beta) a_3$$

$$w_t \leftarrow w_{t-1} - \eta \cdot v_{dw_t}$$

$$v_{dw_t} \leftarrow \beta(v_{dw_{t-1}}) + (1-\beta) \cdot \frac{\partial L}{\partial v_{dw_t}}$$

$$v_{dw_t} \leftarrow \beta(v_{dw_{t-1}}) + (1-\beta) \left(\frac{\partial L}{\partial v_{dw_t}} \right)$$

15



20

- ① Change in gradient component
- ② Change in learning rate
- ③ Change in both GC and LR.

1.1

Gradient Descent with Momentum:

25

$$w_t = w_{t-1} - \eta \cdot \frac{\partial L}{\partial w_{t-1}}$$

30

$$w_t = w_{t-1} - \eta \cdot m_t$$

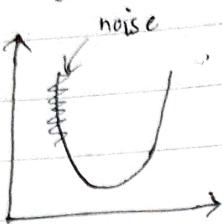
$$m_t = \underbrace{\beta \cdot m_{t-1}}_{0.95} + \underbrace{(1-\beta) \cdot \frac{\partial L}{\partial w_{t-1}}}_{0.05}$$

It helps to avoid getting stuck.

$$w_t \leftarrow w_{t-1} - \eta \cdot \nabla_{w_t}$$

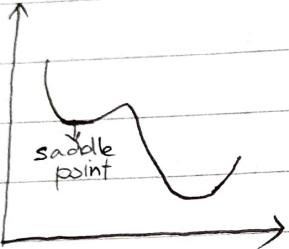
$$b \leftarrow b_{t-1} - \eta \cdot \nabla_{b_t}$$

Optimization Technique: $w; J(w)$



$$w_t = w_{t-1} - \eta \cdot \frac{\partial L}{\partial w}$$

Saddle points (when multiple local minima are present)



① Change in gradient component

② Change in learning rate

③ Change in both G.C and LR.

1.1 Gradient Descent with Momentum:

$$w_t = w_{t-1} - \eta \cdot \frac{\partial L}{\partial w_{t-1}}$$

$$w_t = w_{t-1} - \eta \cdot m_t$$

$$m_t = \beta \cdot m_{t-1} + (1-\beta) \cdot \frac{\partial L}{\partial w_{t-1}}$$

It helps to avoid getting stuck in saddle points

This computation occurs at time $t_0, t_1, t_2, \dots, t_n$.

How it comes?

1 epoch \Rightarrow more than one iteration

1.2 Newton Aculer by (NACL)

Not commonly used.