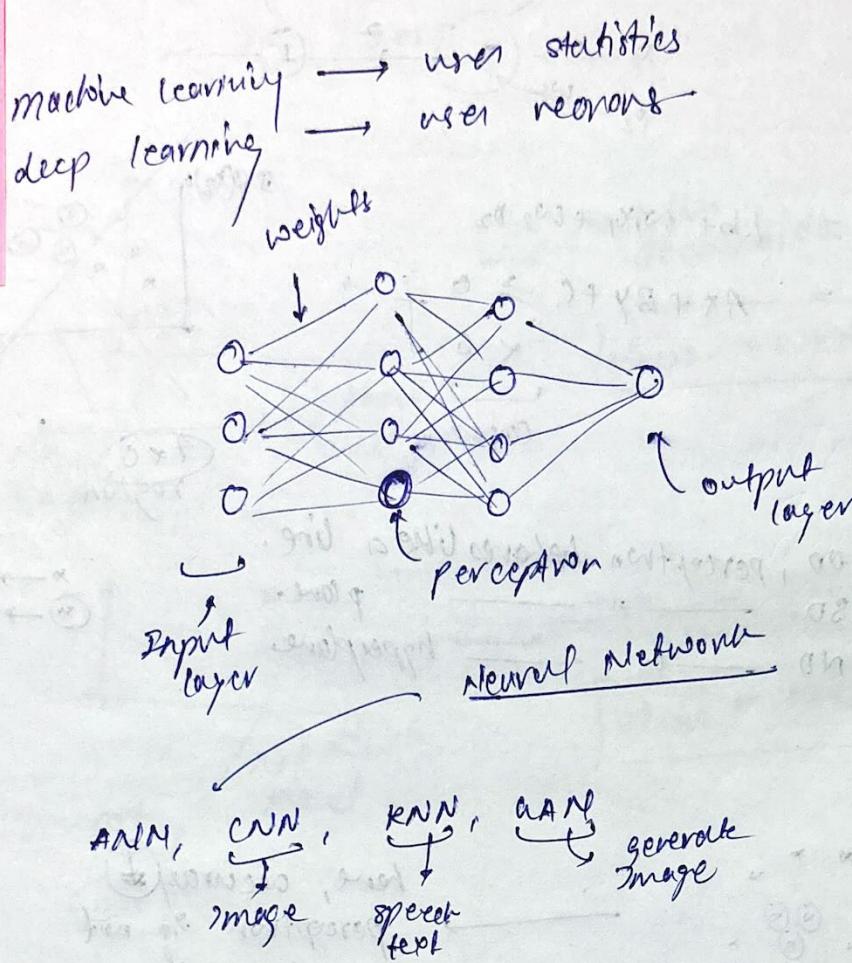
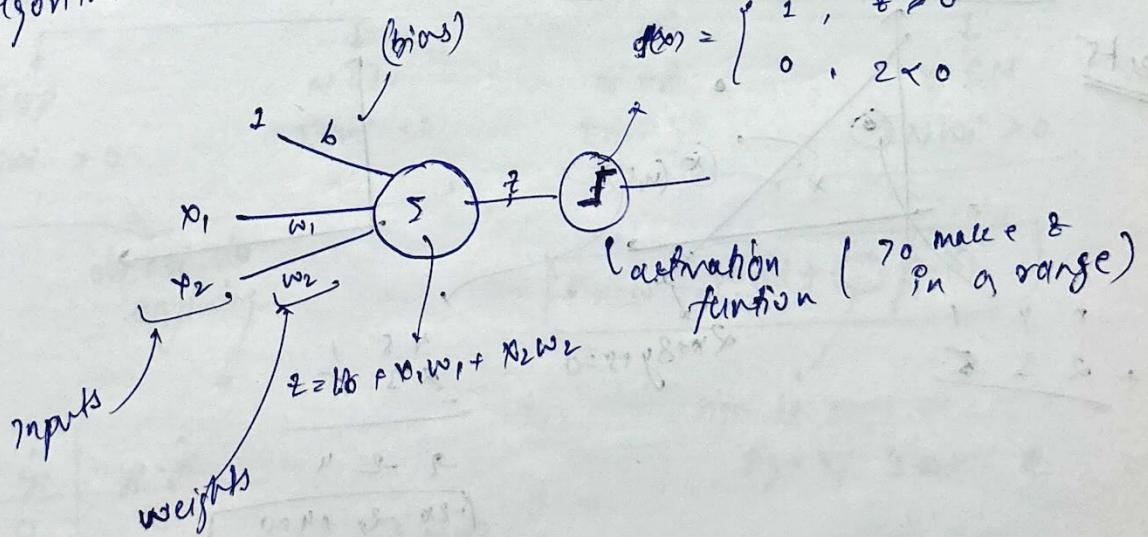


Deep learning

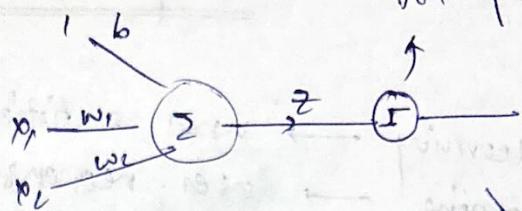


lect⁴

Perception
 Algorithm (mathematical model/function)



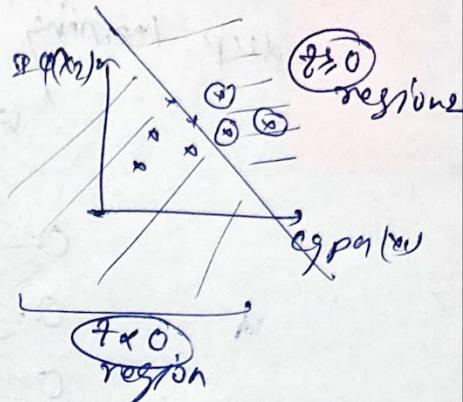
cgpa	sq	salary
x_1	x_2	



$$f = b + w_1 x_1 + w_2 x_2$$

$$= Ax + By + C \geq 0$$

$\underbrace{\quad}_{\text{cogen}}$



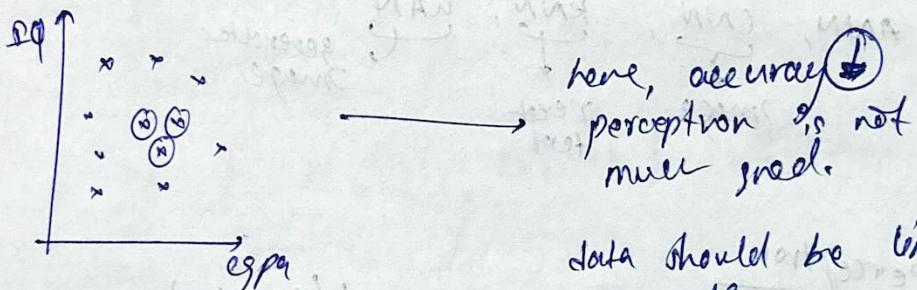
Note) ① In 2D, perceptron behaves like a line.

② 3D \longrightarrow plane.

③ ND \longrightarrow hyperplane

\rightarrow unplaced
 \circlearrowleft placed.

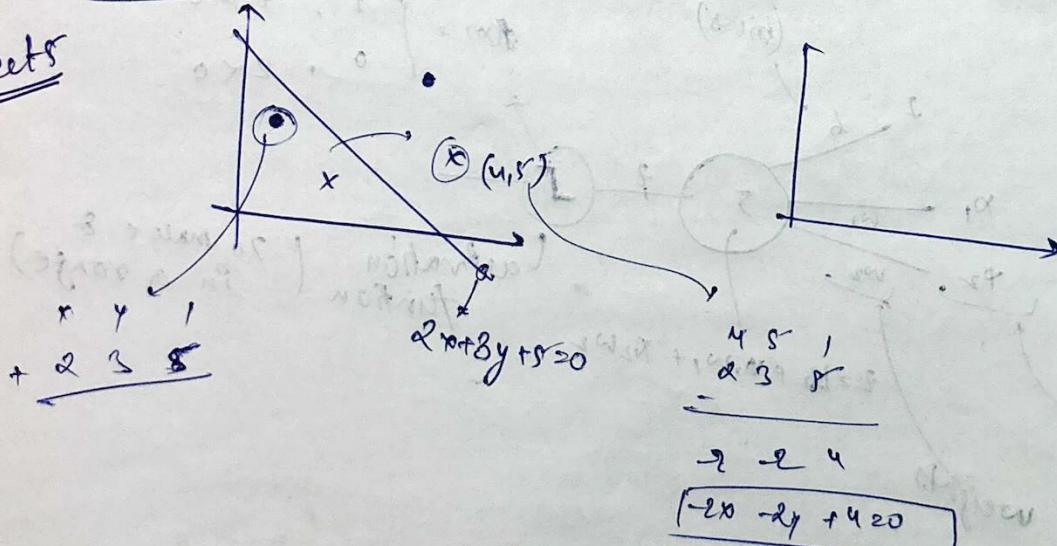
But



Data should be linearly
separable.

Solu with perceptron trick (method 1)

lets

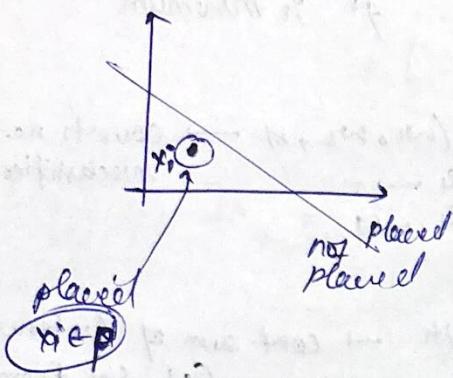


$$\text{Coefnew} = \text{Coefold} - \eta(\text{coord})$$

$$z = w_1 x_1 + w_2 x_2 + w_0 \cdot x_0$$

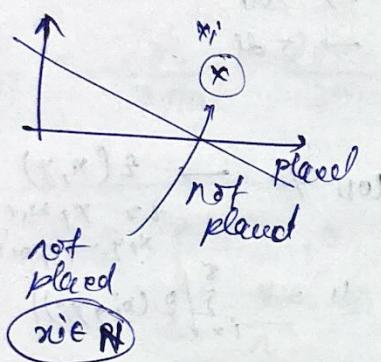
$\boxed{x_0 = 1}$

$$z = \sum_{i=0}^2 w_i x_i$$



move line downwards

$$w_{\text{new}} = w_{\text{old}} - \eta \cdot x_i$$



move line upwards

$$w_{\text{new}} = w_{\text{old}} + \eta \cdot x_i$$

Finally

total cases

$x_i \in P$

$$\sum w_i x_i > 0$$

$x_i \in N$

$$\sum w_i x_i < 0$$

$x_i \in P$

$$\sum w_i x_i < 0$$

$x_i \in N$

$$\sum w_i x_i > 0$$

do not do anything

$$w_{\text{new}} = w_{\text{old}} + \eta \cdot x_i$$

sign is same as sign of $\sum w_i x_i$

y_i	\hat{y}_i	$y_i - \hat{y}_i$
0	0	0
1	1	0
0	1	-1
1	0	1

for all cases

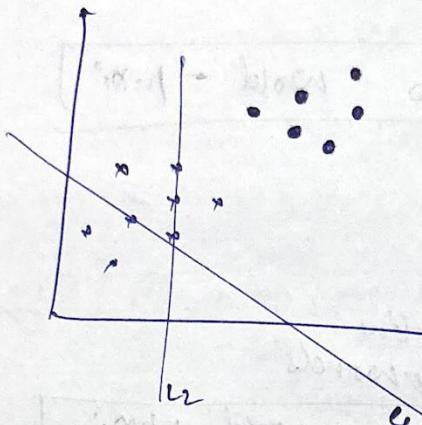
$$w_{\text{new}} = w_{\text{old}} + \eta(y_i - \hat{y}_i) \cdot x_i$$

Lect 6

solution with loss function (method 2)

$$\text{loss fn} = \underbrace{f(w_1, w_2, b)}_{\downarrow}$$

$w_1, w_2, b \rightarrow$ for which f_n is minimum



① $\int f(w_1, w_2, b) \rightarrow$ counts no. of misclassified points
 $u_1 \rightarrow 5$
 $u_2 \rightarrow 1$

② $\int \text{length} \rightarrow \text{cont sum of distance}$
 $u_1 \rightarrow \sum d_i$
 $u_2 \rightarrow \sum d_i$

③ $\int \text{loss fn} \rightarrow \ell(x, y)$
 $x = x_1 w_1 + x_2 w_2 + b$
 $x, y \rightarrow \text{point}$
 $u = \sum_{i=1}^n |\ell(x_i, y_i)|$

scikit-learn Impl

$$\boxed{\text{Loss} = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))}$$

$$f(x_i) = w_1 x_1 + w_2 x_2 + b$$

y_i	$\hat{f}(x_i)$	$\max(0, -y_i \hat{f}(x_i))$
1	1	0
0	0	0
1	0	0
0	1	1

$$= |f(x_i)|$$

$$= |f(x_i)|$$

only for handling sign.

solving with gradient descent

$$\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b}$$

$$\text{loss function} = \frac{1}{n} \sum \max(0, -y_i f(x_i)),$$

where $f(x_i) = w_1 x_1 + w_2 x_2 + b$

$$= \begin{cases} 0 & , -y_i f(x_i) \leq 0 \\ -y_i f(x_i) & , -y_i f(x_i) > 0 \end{cases}$$

$$(L) \text{ loss function} = \sum \begin{cases} 0 & , -y_i f(x_i) \leq 0 \\ \frac{-y_i f(x_i)}{n} & , -y_i f(x_i) > 0 \end{cases}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \cdot \frac{\partial f(x_i)}{\partial w_1}$$

$$\boxed{\frac{\partial L}{\partial w_1} = \sum \begin{cases} 0 & , -y_i f(x_i) \leq 0 \\ \frac{-x_i}{n} \cdot x_{i1} & , -y_i f(x_i) > 0 \end{cases}}$$

$$\boxed{\frac{\partial L}{\partial w_2} = \sum \begin{cases} 0 & , -y_i f(x_i) \leq 0 \\ \frac{-x_i}{n} \cdot x_{i2} & , -y_i f(x_i) > 0 \end{cases}}$$

$$\boxed{\frac{\partial L}{\partial b} = \sum \begin{cases} 0 & , -y_i f(x_i) \leq 0 \\ \frac{-1}{n} & , -y_i f(x_i) > 0 \end{cases}}$$

for i in epochs:

$$w_2 = w_2 - \eta \left(\frac{\partial L}{\partial w_2} \right)$$

$$w_1 = w_1 - \eta \left(\frac{\partial L}{\partial w_1} \right)$$

$$b = b - \eta \left(\frac{\partial L}{\partial b} \right)$$

Note if activation fn is sigmoid

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

$$\text{where, } z = w_1x_1 + w_2x_2 + b$$

and cost fn is binary cross entropy

$$L = -y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

output \rightarrow probability.

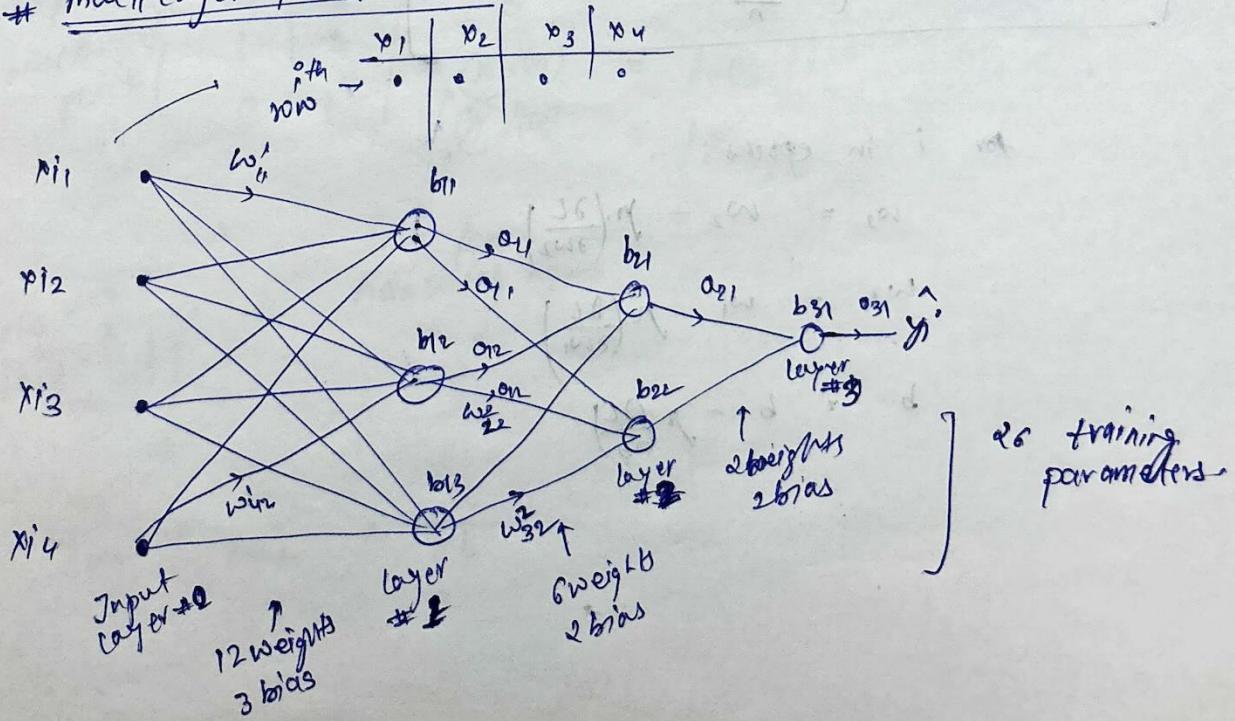
here, perceptrons logistic regression

Note

loss fn	activation fn	output
Hinge loss	step	perceptron ($0 \text{ or } 1$)
log-loss (binary cross entropy)	sigmoid	logistic reg ($0 \text{ or } 1 \rightarrow \text{Probability}$)
categorical cross entropy	softmax	($0 \text{ or } 1$)
MSE	linear	linear regression.

Recap

multilayer perceptron (MLP) Notation



$o_{ij}, b_{ij} \rightarrow$ $i =$ current layer
 $j =$ current node
 output bias

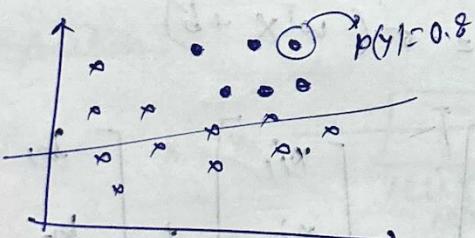
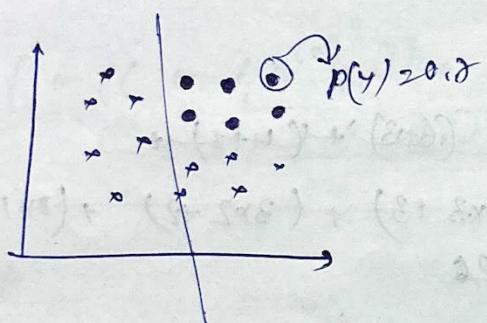
$b_{ij} \rightarrow$ current node
 current layer

~~Ref~~ \rightarrow best current node
~~j = current node~~
~~keeping layer.~~

$w_{kj}^i \rightarrow$ $i =$ current node
 $j =$ going node
 $k =$ going layer.

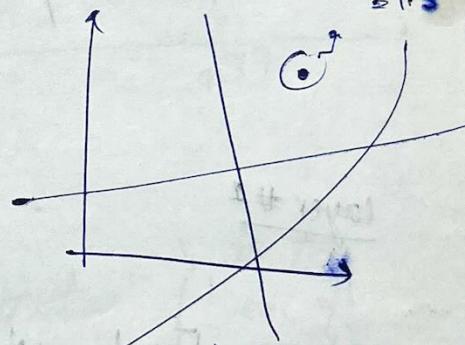
layer going
 $w_{kj}^i \rightarrow$ node going
 node coming from

Unit 9] MLP



can also we weight sum

$$p(y) = 0.2 + 0.8 = 1.0$$

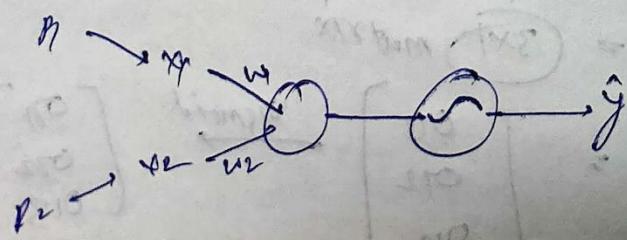


can also give bias

sigmoid

$$\frac{1}{1+e^{-1.5}} = 0.82$$

smoothing



Lect 10

Forward propagation

gpa	2D	20th m	10th m	placed

x_{11}

○

x_{12}

○

x_{13}

○

○

x_{14}

○

○

3

trainable parameters = $(4 \times 3) + (4 \times 2) + (3 \times 2 + 2) + (2 \times 1 + 1)$
 $= 26$

Layer #1

prediction = $\sigma(w^T x + b)$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}^T \cdot \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$(4 \times 3)^T \rightarrow (3 \times 4)$ $u \times 1$
 $(3 \times 1) + (3 \times 1)$

$= \textcircled{3 \times 1} \text{ mat } \times \text{vec}$

$$\Rightarrow \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \end{bmatrix} \xrightarrow{\text{sigmoid}} \begin{bmatrix} 0.7 \\ 0.72 \\ 0.73 \end{bmatrix}$$

layer #2

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}^T \cdot \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix} = \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \end{bmatrix}$$

$$\begin{bmatrix} w_{11}^3 & w_{12}^3 \\ w_{21}^3 & w_{22}^3 \end{bmatrix}^T \cdot \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \end{bmatrix} = \begin{bmatrix} o_{31} \end{bmatrix}$$

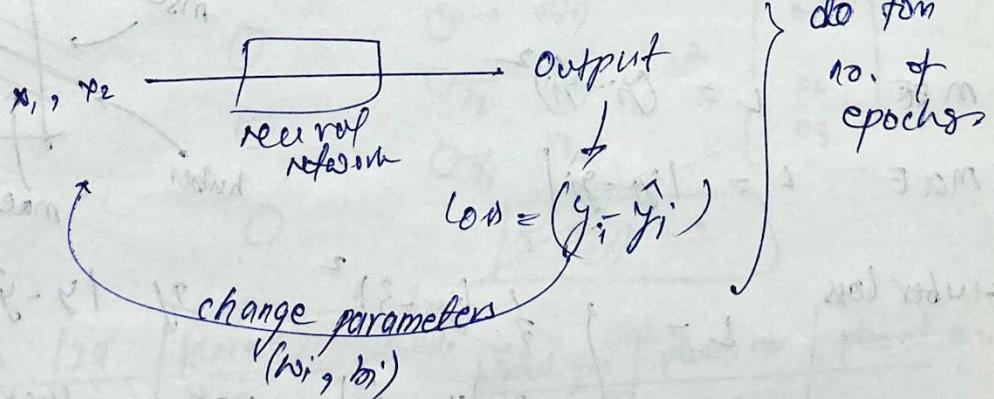
\dots

$$o\left(\begin{bmatrix} o_{21} \end{bmatrix}\right) = \begin{bmatrix} o_{31} \end{bmatrix}$$

$$\text{output} = o \left(o \left(o \left(o \left(a^{[0]}w^{[1]} + b^{[1]}\right) \right) \underbrace{w^{[2]} + b^{[2]}}_{a^{[2]}} \right) \right) \underbrace{w^{[3]} + b^{[3]}}_{a^{[3]}}$$

Layer 4/12/13/14

loss function



same as linear regg

loss fn in DL

Regression

- mse (L2 loss)
- mae
- huber loss

Classification

- binary crossentropy
- categorical
- hinge loss

GAN

- discriminator loss
- min max gain loss

object detection

- focal loss

Auto encoders

- KL divergence

embedding

- Triplet loss.

Loss v/s cost fn

$$(y_i - \hat{y}_i)^2 \rightarrow \text{for single data}$$

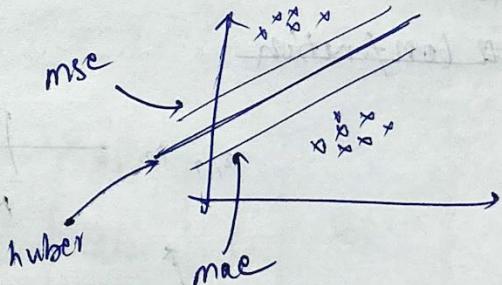
] loss function fn

$$\sum \frac{(y_i^0 - \hat{y}_i^0)^2}{n} \rightarrow \text{for nle data}$$

] cost fn

MSE $L = (y_i - \hat{y}_i)^2$

MAE $L = |y_i - \hat{y}_i|$



Huber loss $L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$

Binary cross Entropy (log loss)
 → classification (two classes $\rightarrow 0/1$)

$$L = -y \cdot \log(\hat{y}) - (1-y) \cdot \log(1-\hat{y})$$

→ output must be sigmoid activation.

hidden layer activation can be anything.

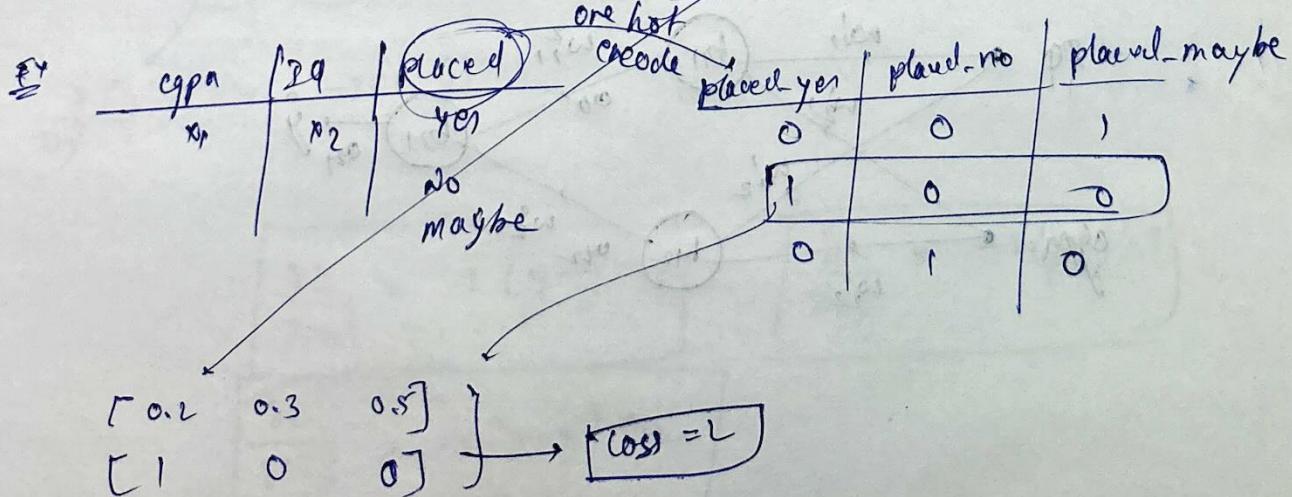
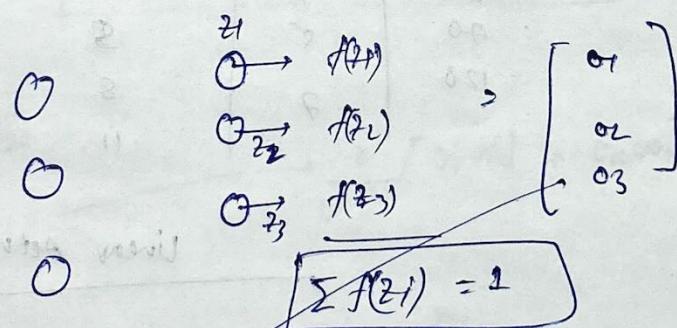
categorical cross entropy

→ classification (multi-class)

$$L = -\sum_{j=1}^K y_i \log(\hat{y}_i), \quad K = \text{no. of classes}$$

→ output must be softmax activation
 hidden layer Act can be anything

$$f(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$



cgpa	29	placed
7	70	Yes
8	80	No
6	60	maybe

categorical
cross entropy

spare categorical
cross entropy

→ placed
one hot encode

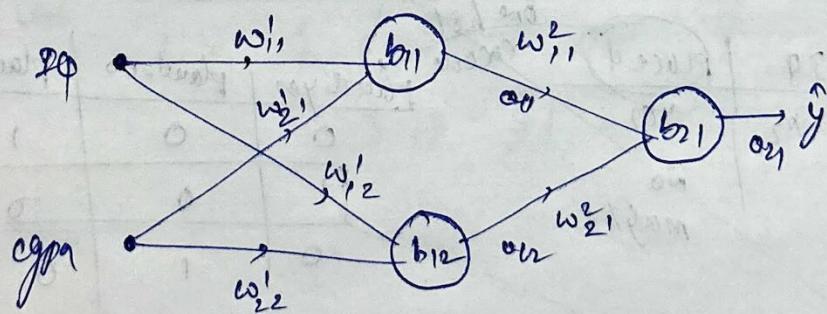
placed
↓
ordinal encode
(1,2,3)

Part 15

Back propagation

29	Cgpa	YPA
80	8	3
80	9	5
70	8	5
120	7	8
30		11
80		

Linear Actv fn.



steps

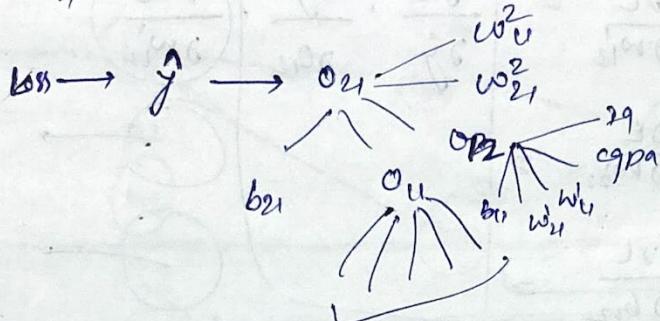
→ 1) Initialize w_{11}, b_1 , anything
 $\boxed{w \rightarrow 1, b \rightarrow 0}$

2) select a student

3) predict (\hat{y}_1) → forward prop

4) choose a cost fn. (mse)

$$O_{11} = w_{11}^2 O_U + w_{21}^2 O_{12} + b_{21}$$



5) update

weights and bias
gradient descent

$$\boxed{w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}}$$

$$\frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial y} = -2(y - \hat{y})$$

$$\frac{\partial}{\partial w_{11}^2} [O_{11} w_{11}^2 + O_{12} w_{21}^2 + b_{21}]$$

$$= O_{11}$$

$$\boxed{\frac{\partial L}{\partial w_{11}^2} = -2(y - \hat{y}) O_{11}}$$

$$\boxed{\frac{\partial L}{\partial b_{21}} = -2(y - \hat{y}) \cdot 1}$$

$$\boxed{\frac{\partial L}{\partial w_{21}^2} = -2(y - \hat{y}) O_{12}}$$

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial y} \right) \cdot \frac{\partial y}{\partial w_1}$$

$$\frac{\partial L}{\partial w_{12}} = \left| \begin{array}{c} \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_{12}} \\ \frac{\partial L}{\partial w_{22}} \\ \frac{\partial L}{\partial b_{12}} \end{array} \right|$$

$\frac{\partial L}{\partial w_{12}} \rightarrow \text{sp} = \pi_1$

$\frac{\partial L}{\partial w_{22}} \rightarrow \text{cpa} = \pi_2$

$\frac{\partial L}{\partial b_{12}} \rightarrow 1$

$$\frac{\partial y}{\partial \omega_{11}} = \frac{\partial}{\partial \omega_{11}} \left[\omega_0^2 \alpha_0 + \omega_1^2 \alpha_1 + b_{21} \right] = \omega_1^2$$

$$\frac{d\hat{y}}{dt} = \frac{d}{dx_0} \left[\frac{\dots}{\dots} \right] = \omega_1^2$$

now update the parameters
and repeat for n no. of epochs.

lect 16 → code implementation

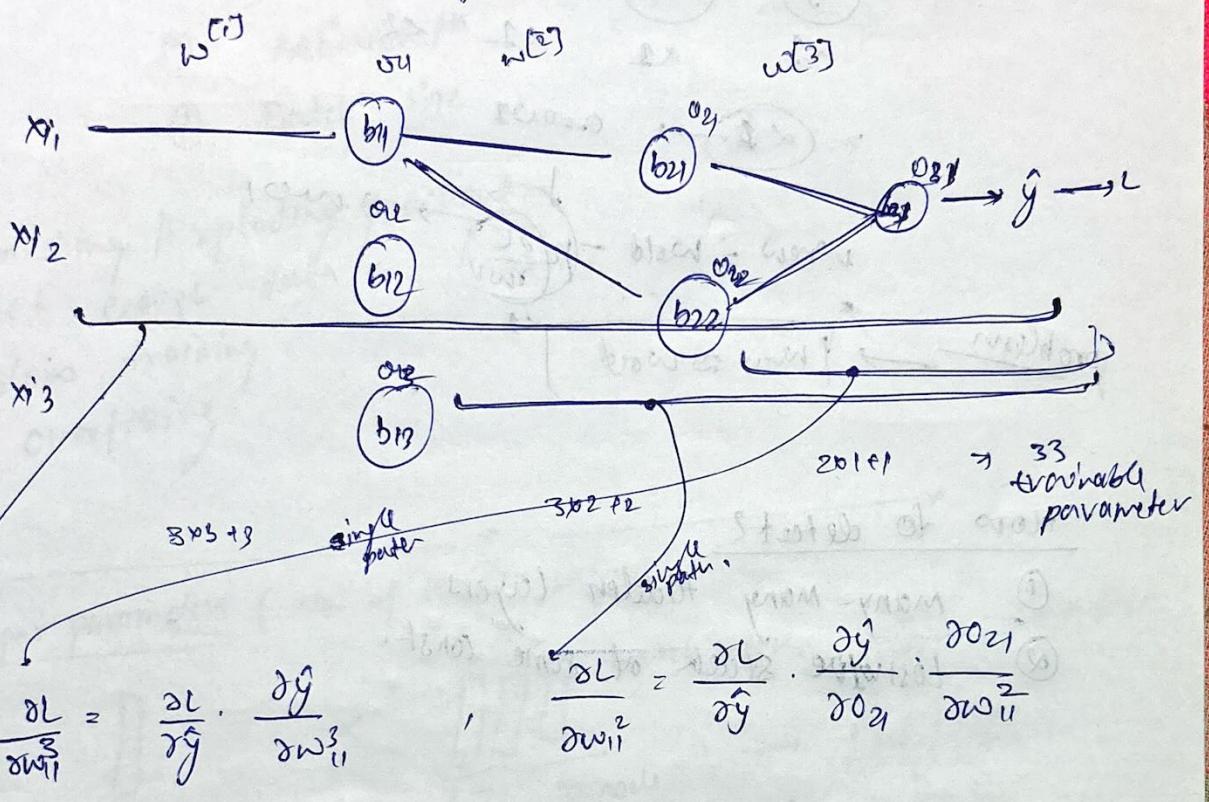
* gradient descent-

lect 17 → MLP Explained in detail

lect 18

MLP memorization

It means to store the computed result for some input (solving Fibonacci no with dp).



multiple paths

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \left[\left(\frac{\partial y}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{21}} \cdot \frac{\partial w_{21}}{\partial w_{11}} \right) + \left(\frac{\partial y}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial w_{22}} \cdot \frac{\partial w_{22}}{\partial w_{11}} \right) \right]$$

Vanishing Gradient problem

where found?

- ① deep neural networks like 10-15 hidden layers
- ② $0.1 \times 0.1 \times 0.1 \times 0.1 = 0.001$
- ③ sigmoid / tanh activation f^n . generally

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial z_1} \right) * \left(\frac{\partial z_1}{\partial x_1} \right) * \frac{\partial z_2}{\partial w_1} * \frac{\partial z_3}{\partial w_1}$$

x_1 x_1 x_2 x_2

$$= \left(\frac{\partial L}{\partial z_1} \right) \rightarrow 0.0001$$

$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_1}$

x_2

problem \rightarrow $w_{new} \approx w_{old}$

How to detect?

- ① many-many hidden layers
- ② loss value stuck at some const.

How to solve?

① make any other neural network with less layers.

$$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \rightarrow \begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$$

- ② using RELU activation f^n .
- $$= \max(0, z)$$

$$= \begin{cases} 0, & z < 0 \\ z, & \text{otherwise} \end{cases}$$

③ Residual network
Building block

- ④ proper weight init $\xrightarrow{\text{glorot, fanin}}$

- ⑤ Batch normalization

lecture

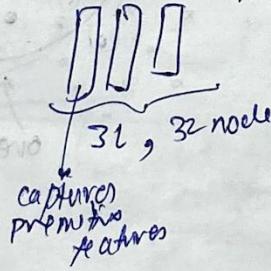
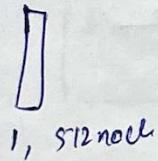
How to Improve a Artificial Neural Network (ANN)

① Fine tuning MN hyperparameters:

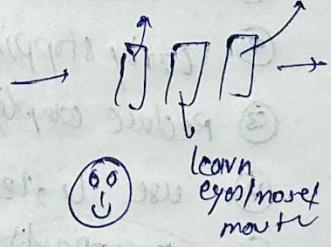
- ↳ ① hidden layers
- ② node per layer
- ③ learning rate
- ④ optimizer
- ⑤ activation fn
- ⑥ Batch size

- ⑦ vanishing / exploding gradient -
- ⑧ Not enough data
- ⑨ slow training
- ⑩ overfitting

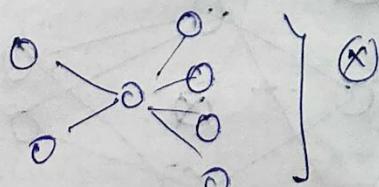
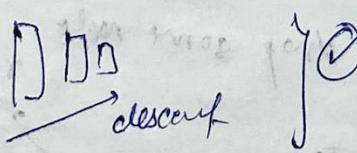
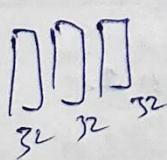
① Hyper parameters (No. of hidden layers)



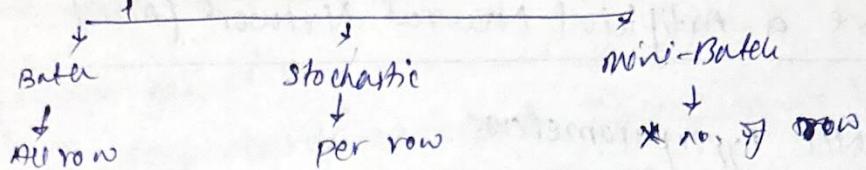
much better. learn lines



② Node in hidden layers



③ Batch size

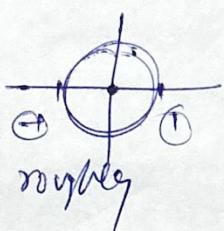


④ Epoch → use early stopping

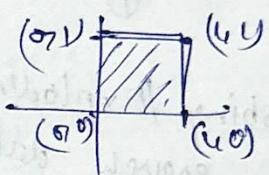
Lect 23

Do scaling

standard



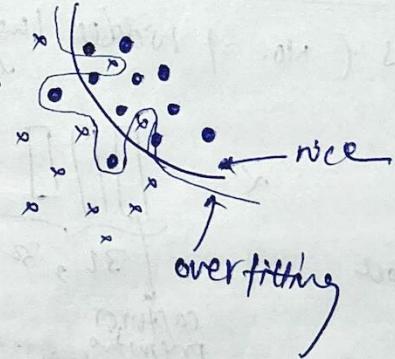
min max



Lect 24

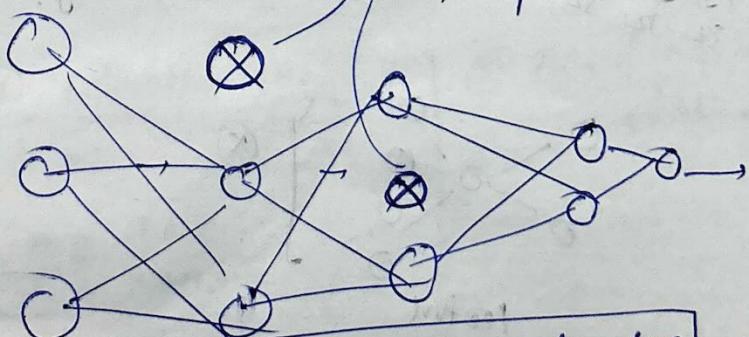
Dropout layers

- ① Add more data
- ② early stopping
- ③ Reduce complexity
- ④ use L1, L2
- ⑤ **Dropouts**



you drop some node

removed randomly for epoch 2.



⇒ Randomly remove some node for each epoch.

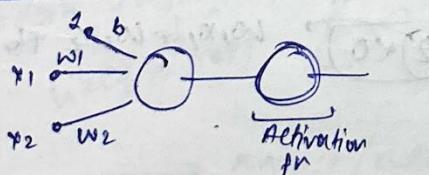
\rightarrow
 ANN $P(10 \rightarrow 50\%)$
 RNN $P(20 \rightarrow 30)$
 CNN $P(40 \rightarrow 50)$
 ↑
 prop of dropout

Lecture 2
 # Regularization
 more complex NN \propto more capable
 make curve fit training data.
 (over fitting).

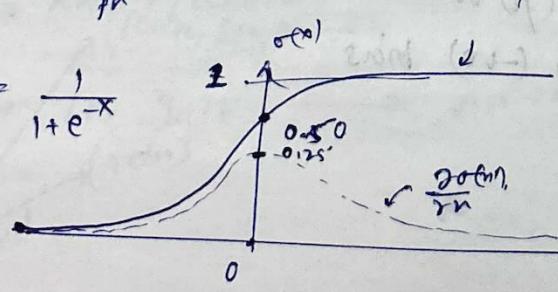
$$\begin{aligned}
 \text{Cost} &= \frac{1}{n} \sum L(\hat{y}_i, y_i) + \text{penalty term.} \\
 (\text{regularization}) &\equiv \frac{\partial}{\partial w} [w_1^2 + w_2^2 + \dots + w_n^2] \\
 (\text{regularization}) &\equiv \frac{\lambda}{n} \sum \|w_i\|^2
 \end{aligned}$$

$$\boxed{w_{new} = (1 - \mu \lambda) w_{old} - \mu \left(\frac{\partial L}{\partial w_{old}} \right) \text{ weight decay}}$$

Lecture 2
 # Activation function



① Sigmoid
 $\sigma(x) \in [0, 1]$

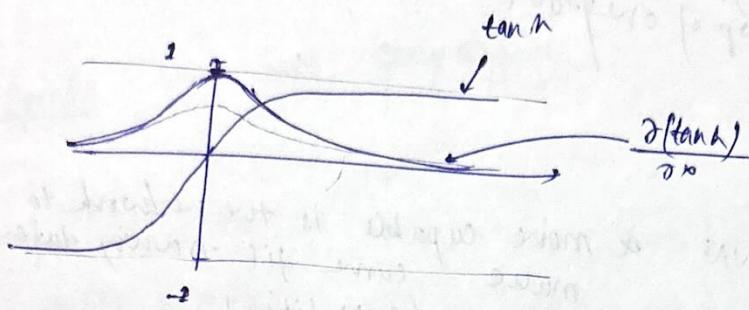


$\sigma(-x) = \frac{e^{-x}}{1 + e^{-x}}$
 $\sigma(n) \in [0, 1]$
 vanishing gradient problem occurs.

can be used as probability in output layer

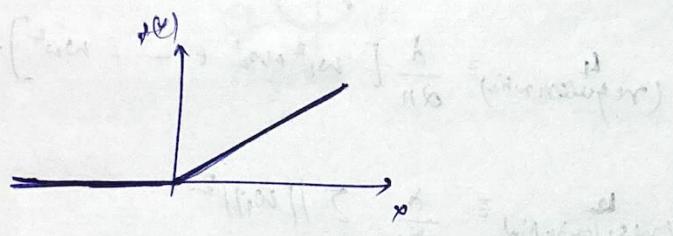
Tanh

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



ReLU

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



Dying ReLU problem

Dead neuron \Rightarrow Neur. to be output zero.

If one time it give zero, then always it will give zero.

Hence If 50% N are dead its a problem.

Dead neuron \rightarrow $\frac{\partial L}{\partial w} = 0 \rightarrow$ hence no update in parameters.

happens bcoz of

when $\sum x_0$, $w_0, x_1, + w_1 x_2 + b$

① high ($+ve$) rate

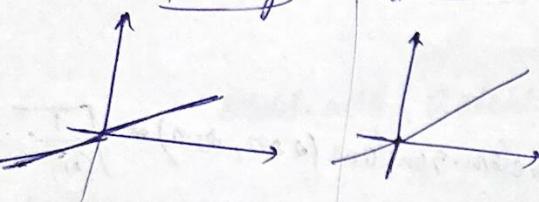
② high ($-ve$) bias

ReLU variants

linear

Leaky

parametric



$$f(x) = \max(0.01x, x)$$

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

α is variable

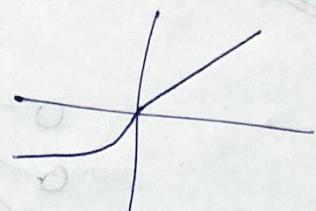
non-linear

ELU

SELU

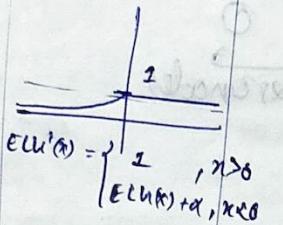


$$\text{ELU}(x) = \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$$

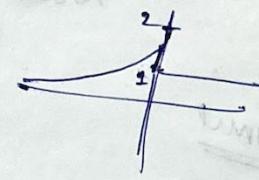


$$\text{SELU}(x) = \lambda \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$$

$$\begin{aligned} \lambda &= 1.673 \\ \alpha &= 1.0587 \end{aligned} \quad \left. \begin{array}{l} \text{fixed} \\ \text{parameters} \end{array} \right.$$



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



$$\text{SELU}'(\alpha) = \begin{cases} x, & x \geq 0 \\ \alpha e^x, & x < 0 \end{cases}$$

lect 09

weight initialization

① $w_i = 0$ (Don't do it)

Sigmoid AF

② Non-zero const (Don't do it)

③ Init it with random values. \otimes
Random too small / Large varies.
(like 0.001) $\beta \sim \mathcal{N}(0, 1)$

lect 10

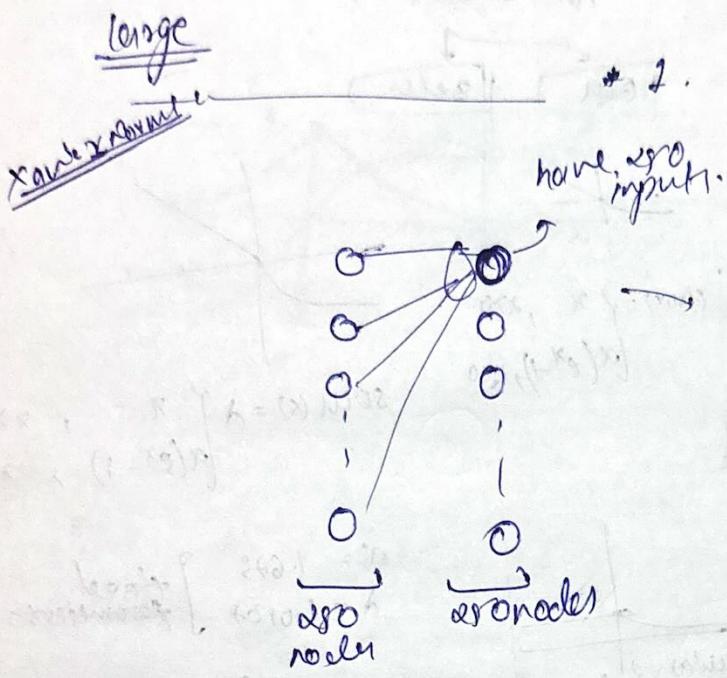
solⁿ

xavier / glorot \rightarrow normal, uniform
(tanh, sigmoid)

He init \rightarrow normal, uniform

Suggested to use

Small
np.random. random(280, 280) * 0.01



np. random. random(280, 280) * $\sqrt{\frac{1}{n}}$

be nonlin

weights = $C \rightarrow \sqrt{\frac{2}{fan\text{-in}}}$ no. of inputs.

random uniform

[-limit, limit] weights range

limit = $\sqrt{\frac{6}{fan\text{-in} + fan\text{-out}}}$

He Uniform

[-limit, limit]

limit = $\sqrt{\frac{6}{fan\text{-in}}}$

Lec 31

Batch Normalization

$$\begin{aligned} \text{mean} &= 0 \\ \sigma^2 &= 1 \end{aligned}$$

↓
you also normalize the
inputs of hidden layer
which were output of other
hidden layer.

`model.add(BatchNormalization())`

Lec 32

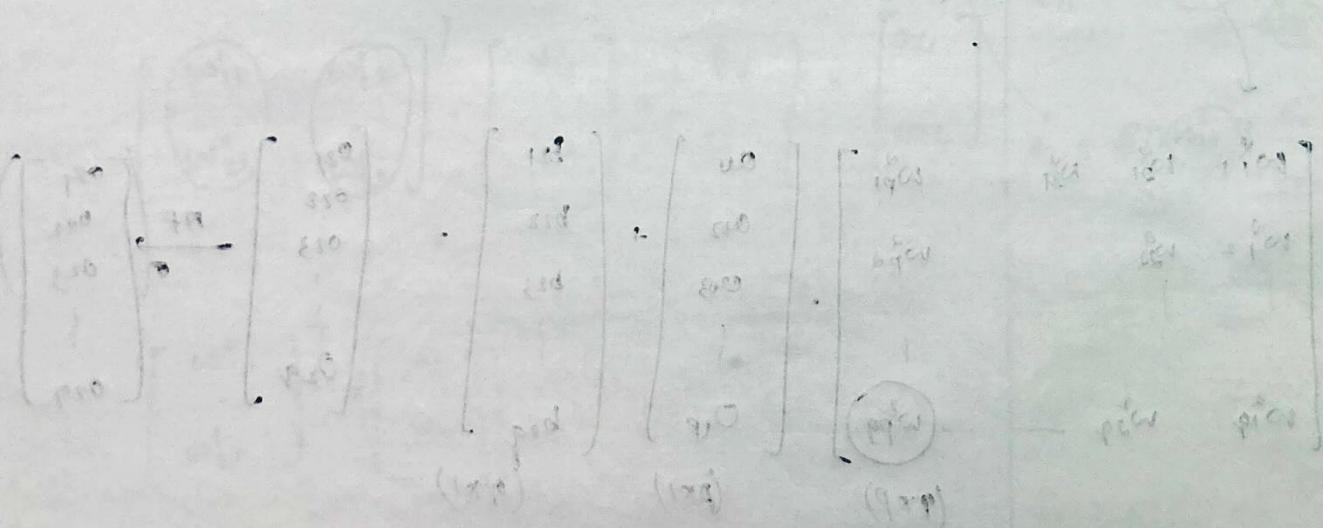
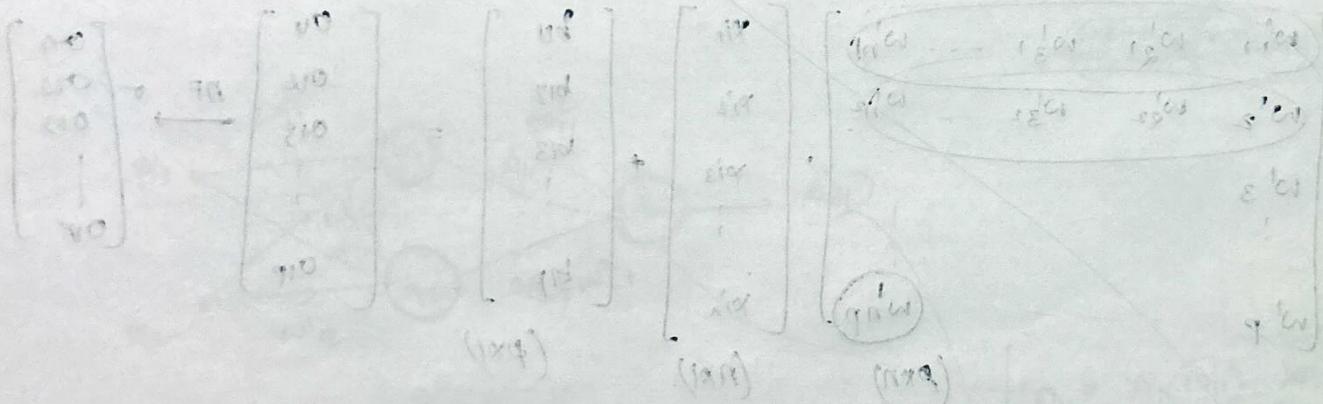
Optimizers

(To be continued)

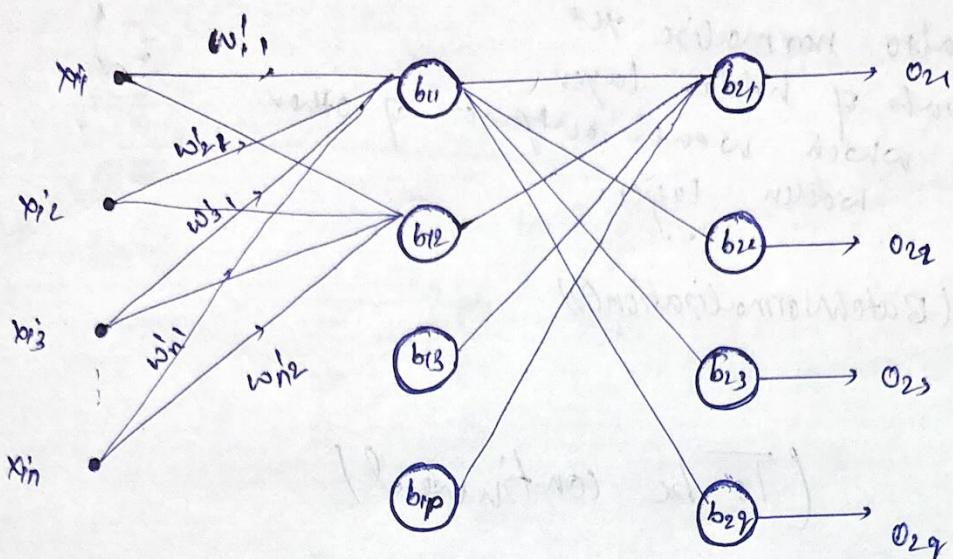
$$\begin{aligned} \text{AdamOptimizer} \\ \text{SGD} = \text{NNA} \\ \text{RMSProp} \end{aligned}$$

$$\begin{aligned} \text{Adagrad} = \text{NNA} \\ \text{Adadelta} \\ \text{Adam} \end{aligned}$$

$$\begin{aligned} \text{AdamOptimizer} \\ \text{Adagrad} \end{aligned}$$



Forward propagation (FAINT) project



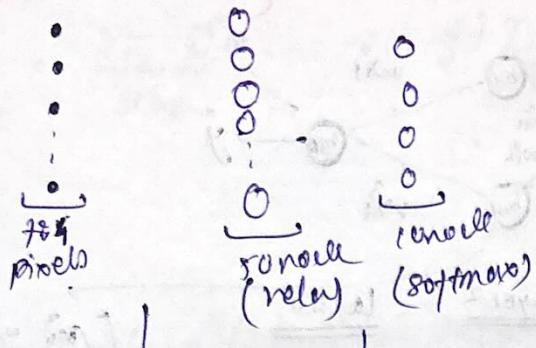
$n = \text{no. of nodes}$
Layer #0

$p = \text{no. of nodes}$
 $\text{AP}^n = \text{ReLU}$
layer #1

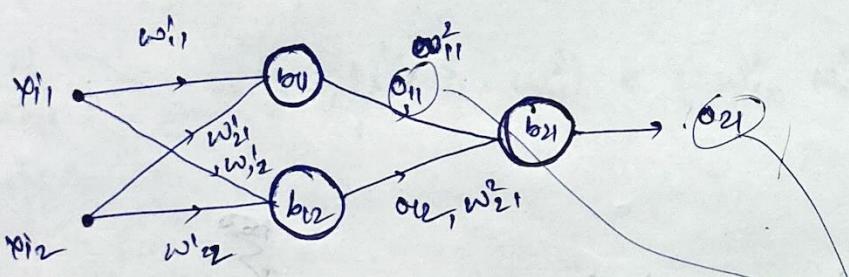
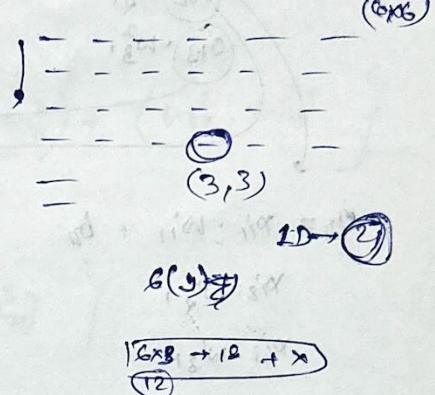
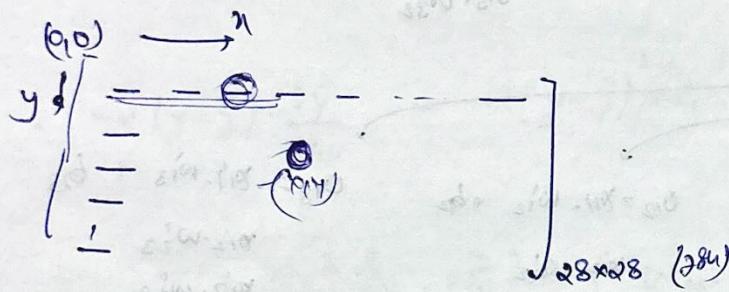
$q = \text{no. of nodes}$
 $\text{AP}^n = \text{softmax}$
layer #2

$$\begin{bmatrix}
 w_{11}^1 & w_{12}^1 & w_{13}^1 & \dots & w_{1p}^1 \\
 w_{12}^1 & w_{12}^2 & w_{13}^2 & \dots & w_{1p}^2 \\
 w_{13}^1 & w_{12}^3 & w_{13}^3 & \dots & w_{1p}^3 \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 w_{1p}^1 & w_{1p}^2 & w_{1p}^3 & \dots & w_{1p}^p
 \end{bmatrix} \cdot \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ x_{in} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ \vdots \\ b_{1p} \end{bmatrix} = \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \\ \vdots \\ o_{1p} \end{bmatrix} \xrightarrow{\text{AP}} \sigma \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \\ \vdots \\ o_{2q} \end{bmatrix}$$

$$\begin{bmatrix}
 w_{21}^1 & w_{22}^1 & w_{23}^1 & w_{2p}^1 \\
 w_{22}^1 & w_{22}^2 & w_{23}^2 & w_{2p}^2 \\
 \vdots & \vdots & \vdots & \vdots \\
 w_{2q}^1 & w_{2q}^2 & w_{2q}^3 & w_{2q}^p
 \end{bmatrix} \cdot \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \\ \vdots \\ o_{1p} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \\ b_{23} \\ \vdots \\ b_{2q} \end{bmatrix} = \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \\ \vdots \\ o_{2q} \end{bmatrix} \xrightarrow{\text{AP}} \sigma \begin{bmatrix} o_{21} \\ o_{22} \\ o_{23} \\ \vdots \\ o_{2q} \end{bmatrix}$$



$$(28 \times 28) \text{ f}(28 \times 28) \rightarrow (28 \times 1) + (50 \times 1)$$

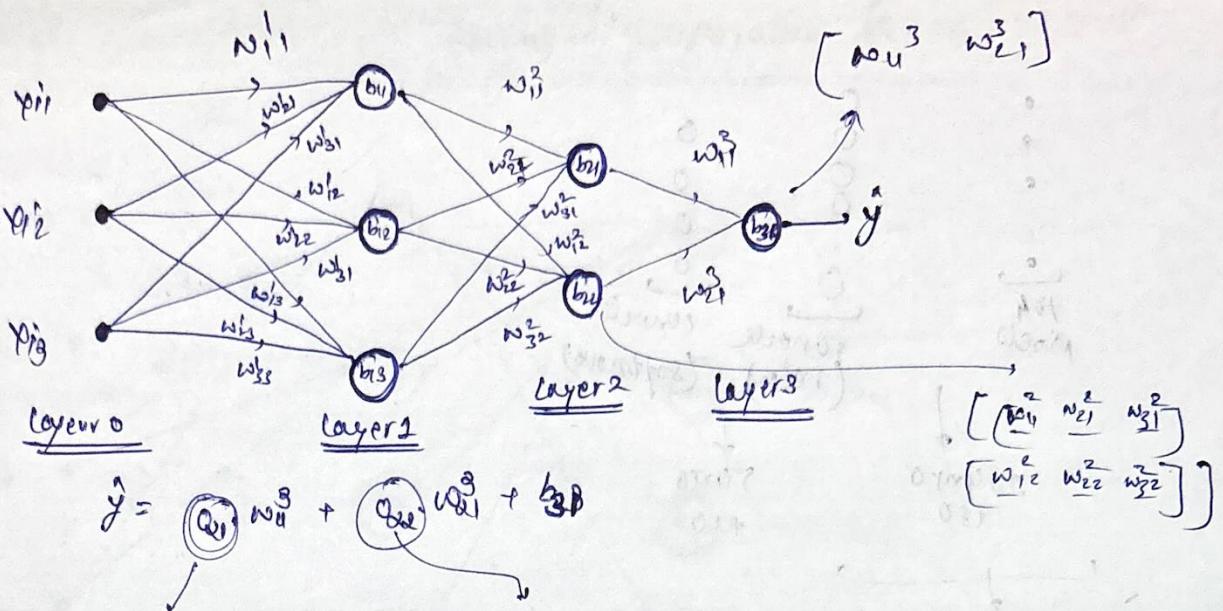


$$\begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix}^T \cdot \begin{bmatrix} x_{1,1} \\ x_{1,2} \end{bmatrix} + \begin{bmatrix} b_{1,1} \\ b_{1,2} \end{bmatrix} = \begin{bmatrix} o_{1,1} \\ o_{1,2} \end{bmatrix}$$

$$o_{1,1} = x_{1,1} \cdot w_{1,1} + x_{1,2} \cdot w_{1,2} + b_{1,1}$$

$$o_{1,2} = x_{1,1} \cdot w_{2,1} + x_{1,2} \cdot w_{2,2} + b_{1,2}$$

$$\begin{bmatrix} w_{1,1} \\ w_{2,1} \end{bmatrix}^T \cdot \begin{bmatrix} o_{1,1} \\ o_{1,2} \end{bmatrix} + [b_{2,1}] = [o_{2,1}]$$



$$o_{j1} = \sum_{i=1}^3 x_i w_{ji} + b_{ji}$$

$$o_{j1} = \sum_{i=1}^3 x_i w_{ji} + b_{ji}$$

$$o_{j1} = \sum_{i=1}^3 x_i w_{ji} + b_{ji}$$

$$x_1 w_{11} + b_{11}$$

$$x_2 w_{21} + b_{21}$$

$$x_3 w_{31} + b_{31}$$

$$o_{j2} = \sum_{i=1}^3 x_i w_{ji} + b_{ji}$$

$$x_1 w_{12} + b_{12}$$

$$x_2 w_{22} + b_{22}$$

$$x_3 w_{32} + b_{32}$$

$$o_{j3} = \sum_{i=1}^3 x_i w_{ji} + b_{ji}$$

$$x_1 w_{13} + b_{13}$$

$$x_2 w_{23} + b_{23}$$

$$x_3 w_{33} + b_{33}$$

$$L = (y - \hat{y})^2$$

Layer 3

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial (y - \hat{y})^2}{\partial w_{11}^3} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial w_{11}^3} = -2(y - \hat{y}) \cdot o_{j1}$$

$$\frac{\partial L}{\partial w_{21}^3} = -2(y - \hat{y}) \cdot o_{j2}$$

Layer 2 (o_{j1}, o_{j2})

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial (y - \hat{y})^2}{\partial w_{11}^2} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial w_{11}^2} = -2(y - \hat{y}) \cdot o_{j1}$$

$$\frac{\partial L}{\partial w_{21}^2} = -2(y - \hat{y}) \cdot o_{j2}$$

Layer 1

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial (y - \hat{y})^2}{\partial w_{11}^1} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial w_{11}^1} = -2(y - \hat{y}) \cdot \left(w_{11}^3 \cdot w_{11}^2 \cdot x_{11} + w_{21}^3 \cdot w_{21}^2 \cdot x_{11} \right)$$

$$\frac{\partial L}{\partial w_{21}^1} = -2(y - \hat{y}) \cdot \left(w_{11}^3 \cdot w_{11}^2 \cdot x_{12} + w_{21}^3 \cdot w_{21}^2 \cdot x_{12} \right)$$

layer 3

$$\frac{\partial L}{\partial b_{31}} = \frac{\partial (y - \hat{y})^2}{\partial b_{31}} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial b_{31}} = -2(y - \hat{y}) \cdot 1$$

layer 2

$$\frac{\partial L}{\partial b_{21}} = -2(y - \hat{y}) \left(\frac{\partial \hat{y}}{\partial b_{21}} \right) = -2(y - \hat{y}) \cdot w_{11}^3$$

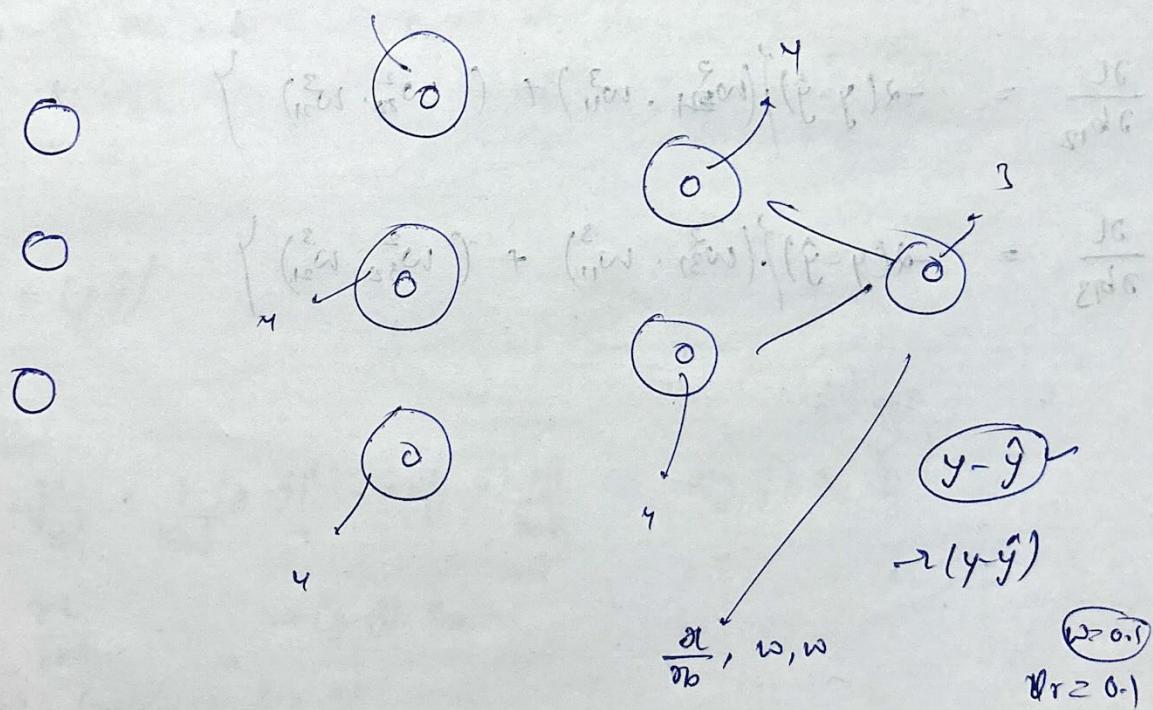
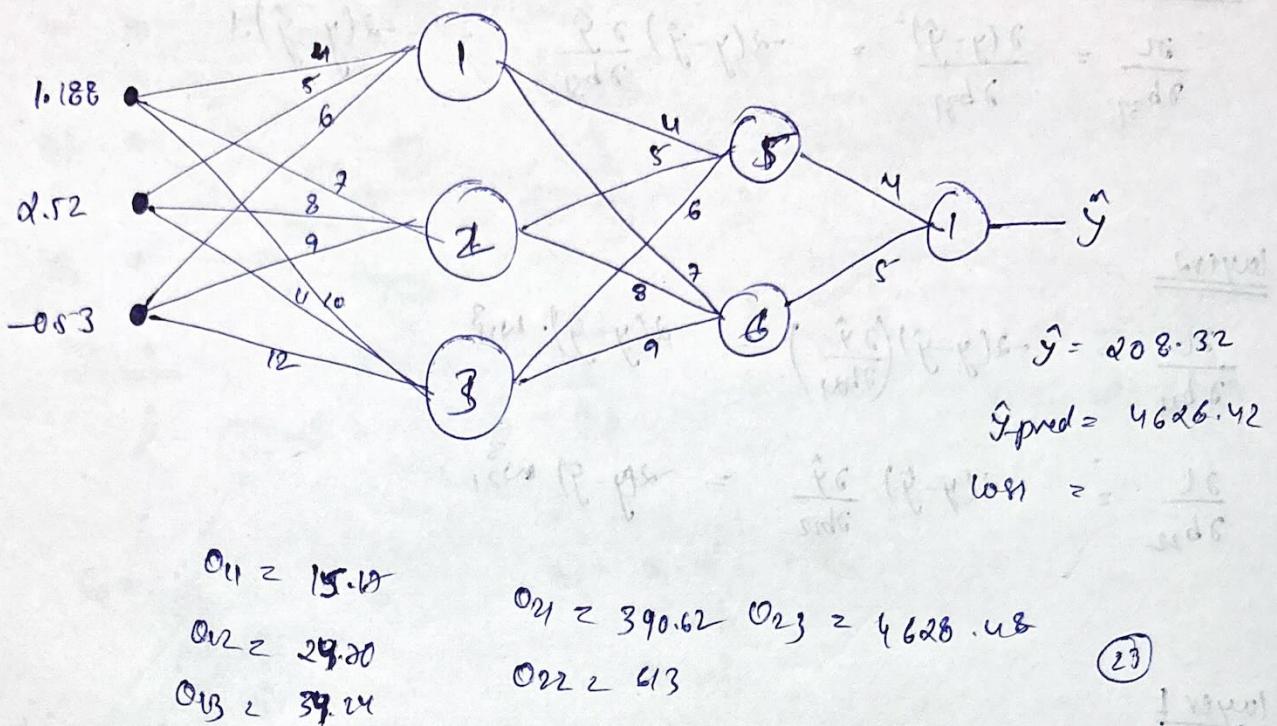
$$\frac{\partial L}{\partial b_{22}} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial b_{22}} = -2(y - \hat{y}) w_{21}^3$$

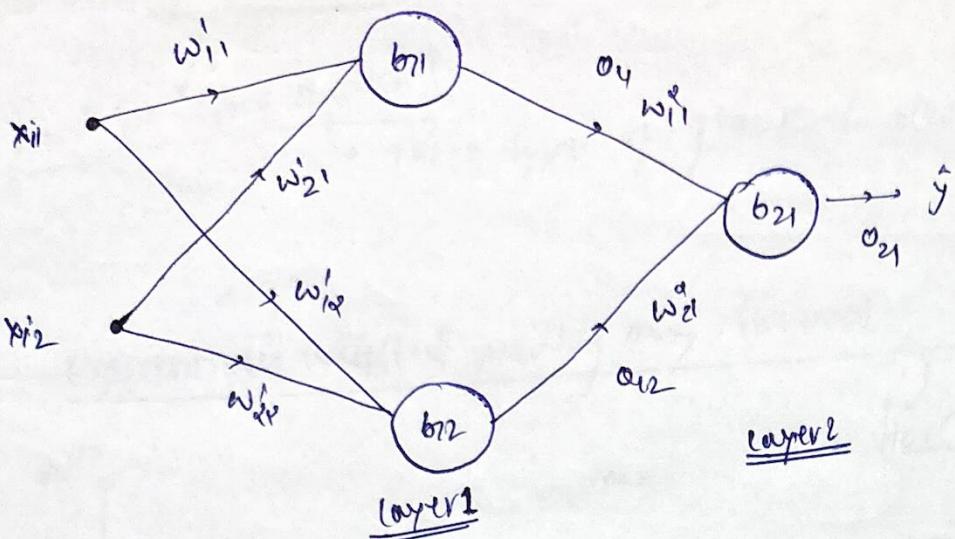
layer 1

$$\frac{\partial L}{\partial b_{11}} = -2(y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial b_{11}} = -2(y - \hat{y}) \left[(w_{11}^2 \cdot w_{11}^3) + (w_{12}^2 \cdot w_{21}^3) \right]$$

$$\frac{\partial L}{\partial b_{12}} = -2(y - \hat{y}) \left[(w_{12}^2 \cdot w_{11}^3) + (w_{22}^2 \cdot w_{21}^3) \right]$$

$$\frac{\partial L}{\partial b_{13}} = -2(y - \hat{y}) \left[(w_{13}^2 \cdot w_{11}^3) + (w_{32}^2 \cdot w_{21}^3) \right]$$





$$\begin{aligned}
 o_{11} &= x_{11} \cdot w_{11} + b_{11} \\
 o_{12} &= x_{11} \cdot w_{12} + x_{12} \cdot w_{21} + b_{12} \\
 o_{21} &= x_{11} \cdot w_{11} + x_{12} \cdot w_{21} + b_{21}
 \end{aligned}$$

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ 0.3 & 0.4 \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} = \begin{bmatrix} o_{11} \\ o_{12} \end{bmatrix} \rightarrow \sigma \left(\begin{bmatrix} o_{11} \\ o_{12} \end{bmatrix} \right) = \begin{bmatrix} 0.82 \\ 1.083 \end{bmatrix}$$

$$\begin{bmatrix} w_{11}^2 & w_{21}^2 \\ 0.5 & 0.6 \end{bmatrix} \begin{bmatrix} o_{11} \\ o_{12} \end{bmatrix} + \begin{bmatrix} b_{21} \\ 0.07 \end{bmatrix} = \begin{bmatrix} o_{21} \\ 0.07 \end{bmatrix} \rightarrow \sigma(o_{21}) = 1.083$$

$$L = \frac{(y - \hat{y})^2}{2} = \frac{(y - \hat{y})}{2(y - \hat{y})} = \frac{(10 - 1.083)^2}{2(10 - 1.083)} = \frac{8.918}{18.834} = -12.834$$

Layer 2

$$\frac{\partial L}{\partial w_{11}^2} = -\alpha (y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial w_{11}^2} = -\alpha (y - \hat{y}) \cdot o_{11} = -10.16$$

$$\frac{\partial L}{\partial w_{21}^2} = -\alpha (y - \hat{y}) \cdot o_{12} = -21.04$$

$$\frac{\partial L}{\partial b_{21}} = -\alpha (y - \hat{y}) \cdot 1 = -12.834$$

Layer 1

$$-8.918 = \frac{\partial L}{\partial w_{11}} = -\alpha (y - \hat{y}) \cdot x_{11} \cdot (w_{11}^2) = -\alpha (y - \hat{y}) \cdot x_{11} \cdot w_{11}^2$$

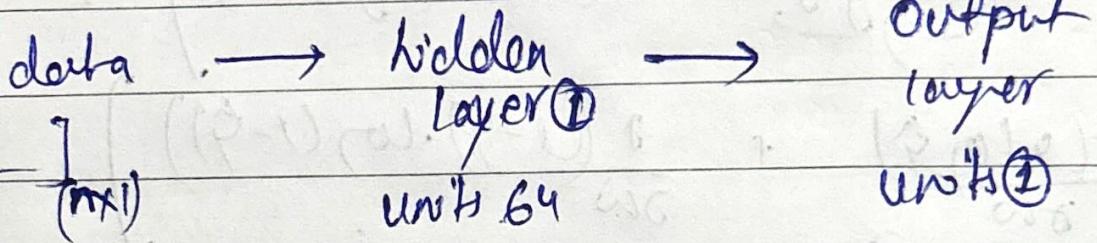
$$-12.834 = \frac{\partial L}{\partial w_{21}} = -\alpha (y - \hat{y}) \cdot x_{12} \cdot w_{21}^2$$

$$-8.918 = \frac{\partial L}{\partial b_{11}} = -\alpha (y - \hat{y}) \cdot w_{11}^2$$

$$\begin{aligned}
 \frac{\partial L}{\partial w_{11}} &= -\alpha (y - \hat{y}) \cdot x_{11} \cdot w_{11}^2 = -10.2 \\
 \frac{\partial L}{\partial w_{21}} &= -\alpha (y - \hat{y}) \cdot x_{12} \cdot w_{21}^2 = -21.4 \\
 \frac{\partial L}{\partial b_{11}} &= -\alpha (y - \hat{y}) \cdot w_{11}^2 = -10.2
 \end{aligned}$$

~~Notes~~

Both sigmoid



$$z = WX + B$$

$$(64 \times 1) (64 \times n) (n \times 1) (64 \times 1)$$

$$z = WX + B$$

$$(1 \times 64) (64 \times 1) (1 \times 1)$$

Layer 1

$$\begin{aligned} & \text{Input } X \text{ (n x 1)} \xrightarrow{\text{Matrix } W_1 \text{ (64 x n)}} \text{Intermediate } z_1 \text{ (64 x 1)} \\ & z_1 = \begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots & w_{1n} \end{bmatrix} \begin{bmatrix} - \\ - \\ - \\ \vdots \\ - \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ \vdots \\ b_{164} \end{bmatrix} \\ & \text{Output } A_1 \text{ (64 x 1)} \end{aligned}$$

$$\begin{aligned} & \text{Input } A_1 \text{ (64 x 1)} \xrightarrow{\text{Matrix } W_2 \text{ (2 x 64)}} \text{Intermediate } z_2 \text{ (2 x 1)} \\ & z_2 = \begin{bmatrix} w_{21} & w_{22} & \dots & w_{264} \end{bmatrix} \begin{bmatrix} 0_{11} \\ 0_{12} \\ 0_{13} \\ \vdots \\ 0_{64} \end{bmatrix} + \begin{bmatrix} b_{21} \end{bmatrix} \xrightarrow{\text{Sigmoid}} \begin{bmatrix} 0_{21} \end{bmatrix} \\ & \text{Output } A_2 \text{ (2 x 1)} \end{aligned}$$

R

$$\hat{z}_1 = w_1 \cdot x_1 + b_1$$

$$A_1 = \sigma(z_1)$$

$$\hat{z}_2 = w_2 \cdot A_1 + b_2$$

$$A_2 = \sigma(z_2)$$

No.

Date: / /

$$\frac{\partial L}{\partial w_i} = \frac{1}{m} \cdot x (y - \hat{y})^T$$

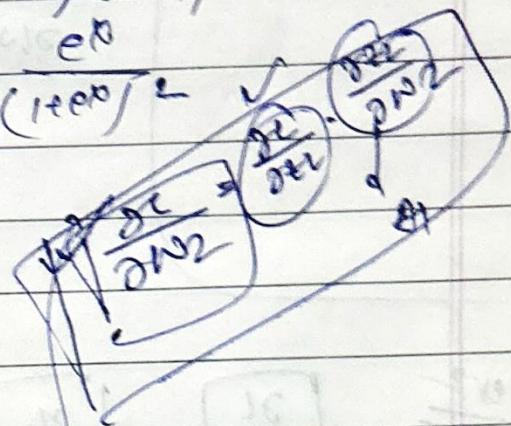
$$\frac{\partial L}{\partial b} = \frac{1}{m} \sum (\hat{y}_i - y_i)$$

$$(\text{obj}) = - (y \cdot \log \hat{y} + (1-y) \cdot \log(1-\hat{y}))$$

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

$$\hat{z}_1 = w_1 \cdot x + b_1$$

$$A_1 = \sigma(\hat{z}_1) =$$



$$\hat{z}_1 = w_1 \cdot x + b_1$$

$$\begin{cases} \hat{z}_1 = w_1 \cdot x + b_1 \\ a_1 = \text{ReLU}(\hat{z}_1) \end{cases} \quad \frac{\partial z_2}{\partial z_1}$$

$$g = \begin{cases} \hat{z}_2 = w_2 \cdot A_1 + b_2 \\ A_2 = \text{sigmoid}(\hat{z}_2) \end{cases}$$

$$L = -y \cdot \log \hat{y} - (1-y) \cdot \log(1-\hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}}$$

$$\begin{cases} \frac{\partial \hat{y}}{\partial z_2} = \sigma(z_2)(1-\sigma(z_2)) \\ = g(1-g) \end{cases}$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2}$$

$$\frac{\partial L}{\partial z_2} = -y(1-g) + g(1-y)$$

$$\begin{aligned} z_1 &= w_1 \cdot x + b_1 \Rightarrow a_1 = \text{relu}(z_1) \\ z_2 &= w_2 \cdot a_1 + b_2 \Rightarrow a_2 = \text{sigmoid}(z_2) = y \end{aligned}$$

$$\frac{\partial \hat{y}}{\partial w_1} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{222}{142} = \frac{w_2}{w_1}$$

$$\frac{\partial \sigma(z_2)}{\partial w_1} = \sigma(z_1)(1 - \sigma(z_2)) \cdot \frac{\partial z_2}{\partial w_1}$$

$$w_2 \cdot \left(\begin{array}{c} \cancel{w_1} \\ w_1 \end{array} \right) \rightarrow \left\{ \begin{array}{l} 1, z_1 > 0 \\ 0, z_1 \leq 0 \end{array} \right.$$

~~layer 2~~

$$\left[\frac{\partial L}{\partial w_2} \right], \left[\frac{\partial L}{\partial b_2} \right] \hat{y}(1-\hat{y})$$

$$\frac{\partial L}{\partial w_2} = \left(\frac{\partial L}{\partial y} \right) \cdot \left(\frac{\partial y}{\partial w_2} \right) + \left(\frac{\partial L}{\partial w_2} \right)$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} \rightarrow 1.$$

Layer 2

$$\left| \frac{\partial}{\partial w_1} \right|, \left| \frac{\partial}{\partial b_1} \right|$$

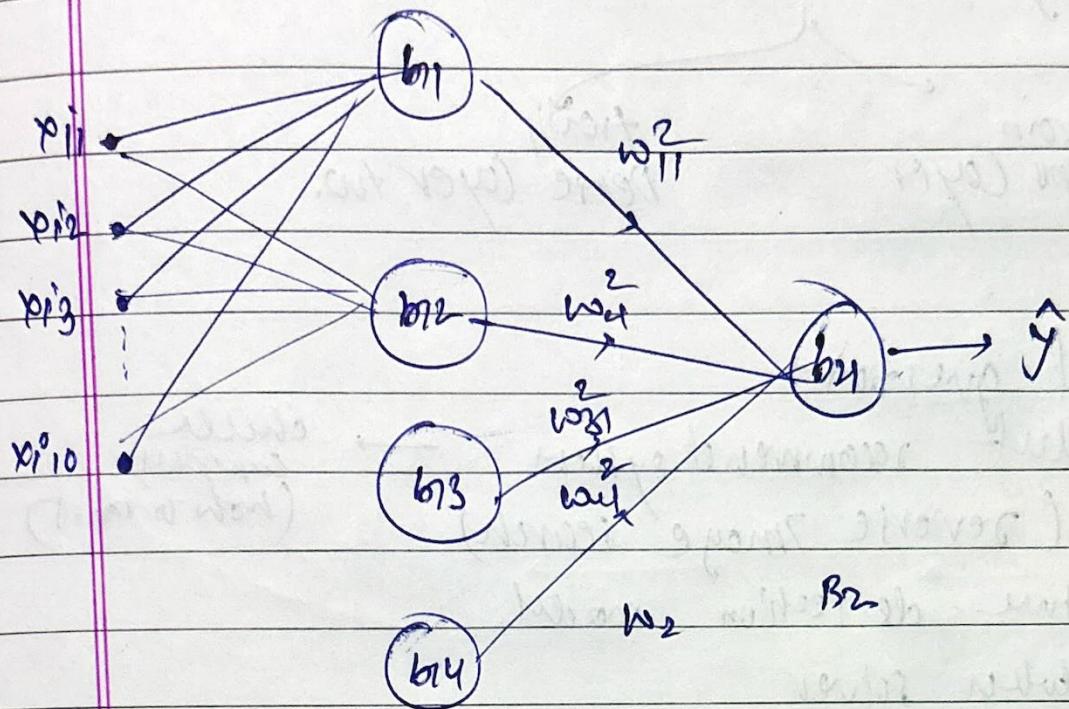
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_1}.$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z_2} \cdot \frac{\partial z_2}{\partial b}$$

$$\frac{\partial z_2}{\partial w_1} = \frac{w_2}{w_1} \left(\frac{\partial A_1}{\partial w_1} \right) \rightarrow (u \times 1)$$

$(1 \times 4) T \rightarrow (u \times 1)$

No. _____
Date: / /



$$W_1 \begin{bmatrix} x_{i,1} & x_{i,2} & x_{i,3} & x_{i,10} \end{bmatrix} \cdot \begin{bmatrix} w_{i1} \\ w_{i2} \\ w_{i3} \\ w_{i4} \end{bmatrix} + B_1 \begin{bmatrix} b_{i1} \\ b_{i2} \\ b_{i3} \\ b_{i4} \end{bmatrix} = A_1 \begin{bmatrix} u \times 1 \end{bmatrix}$$

$$W_2 \cdot A_1 + B_2 = A_2$$

$$(1 \times 4) (u \times 1) \quad (1 \times 1) \quad (1 \times 1)$$

$$\frac{\partial z_2}{\partial w_1} = W_2 \frac{\partial A_1}{\partial w_1} \times \frac{\partial A_2}{\partial w_1} \rightarrow x_{i1}$$

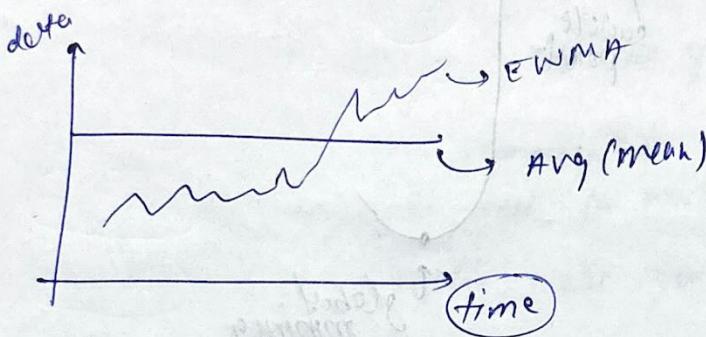
$$\begin{cases} 1, z_1 > 0 \\ 0, z_1 \leq 0 \end{cases}$$

Lecture 32 # optimizers → methods which are used to find weights and biases.

what we have studied? → three types of gradient descent.

Lecture 33

Exponentially weighted moving avg (EWMA)



for time series data

weights
Data_{t2} > Data_{t1}
where t₂ > t₁

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

data at time (t)

EWMA at time(t)

EWMA out (t+1)

$v_0 = 0$ or $v_0 = \theta_0$
more accurate result

Date	θ (temp)
1	28
2	13
3	17
4	31
5	43

EWMA is avg of last this much time

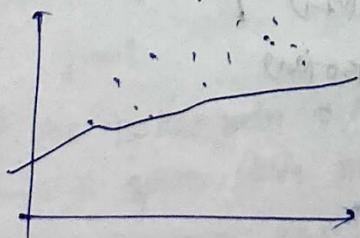
$$v_1 = (0.9) \cdot v_0 + (0.1) \cdot 13$$

$$v_1 = 1.3$$

$$v_2 = (0.9)(1.3) + (0.1)(14)$$

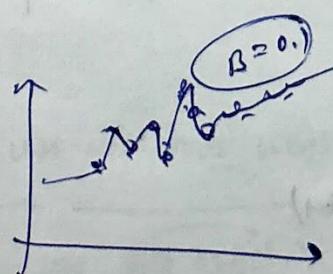
v_3, v_4, v_5

$\beta = 0.98$



$\beta \uparrow$
 $\beta \downarrow$

more weightage to prev points
more weightage to newer points



$\beta = 0.1$

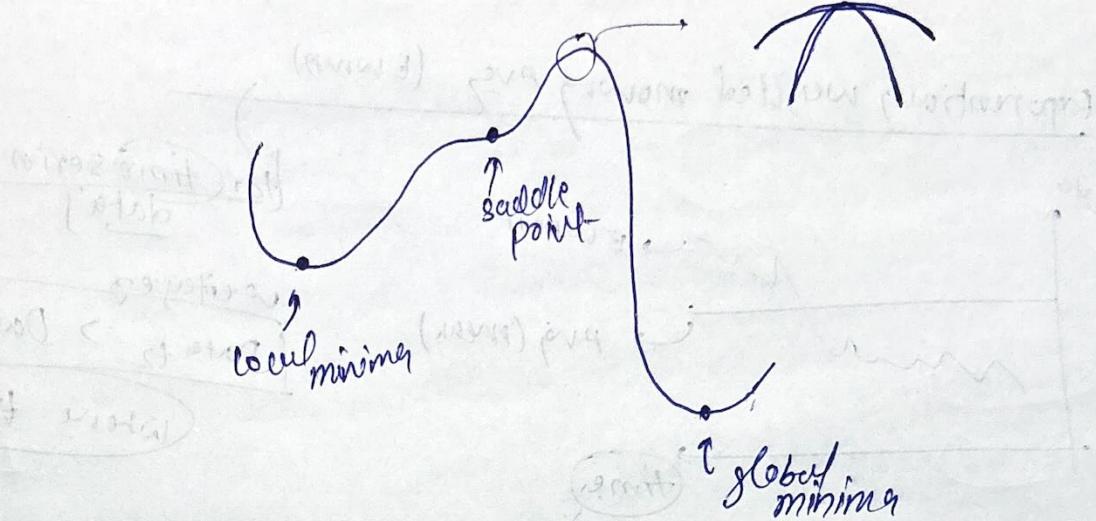
In parabola

clf. [column name]. cwm (alpha = 0.1). means

Part 34]

sto. gradient descent
with momentum

curvature.



problems

- ① local minima
- ② saddle point
- ③ High curvature

make training slow.

momentum technique.
is good in all three cases.

Ball

speed \uparrow

If 4 people tell you, you are moving right

then you have your speed

else if said right and said left, you move right with slow speed.

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

$v \rightarrow$ velocity

$$v_{t+1} = v_t - \eta \cdot \nabla v_t$$

$$v_t = \beta \cdot (v_{t-1}) + \eta \cdot (\Delta w_t)$$

BS (or)

Implementation

[lect 25]

Nesterov Accelerated Gradient (NAG)

$$w_{new} = w_{old} - \left(\beta \frac{\partial J}{\partial w} + \beta V_t \right)$$

\downarrow
clip'n Acceleration

Here

$$w_{new} = \underbrace{w_{old} - \beta V_t}_{\text{At end of } \frac{\partial J}{\partial w}}$$

Here your magnitude
is less for moving.

→ But NAG can make you settle in local minima. however.

In Verbal

e.g. keras optimizers. SGD

learning_rate = 0.01, momentum = 0, nesterov = 0, name = "SGD", ~~kwarg~~

SGD

0
False

SGD with momentum

momentum = 0.01
nesterov = False

Nesterov

= 0.01
= True

[lect 26] # AdaGrad

Adaptive Gradient → learning rate keep on changing.

Apply when most of the
feature values are zero. (sparse)

Why?

It's column mai Bahut zado o honge (sparse) uss clip'n mai weight/mas term update honge
and —————— sparse delta nahi hoga —————— Bahut update honge
that's why NAG, SGD are slow. use AdaGrad

Lect 37] # Root mean square prop. (RMS prop).
↓
Improvement over Adam grad.

First you move very quickly towards non-sparse data,
then towards sparse data.

Bcoz of this what happens is, if there were
some movement left in non-sparse dirn it
cannot converge as the lr may become too low.
to solve this RMS prop.

* Bent optimizing technique. (No disadvantage)

Lect 38] # Adam optimizer

→ on low

Lect 39] # kerov Tuner (Hyperparameter)

import kerovtuner as kt

def build_model(hp):

model = Sequential()

model.add(Dense(32, return_ →))

model.compile(optimizer=hp.Choice('optimizer', ['adam', 'sgd']))

return model.

tuner = kt.RandomSearch(build_model, objective='val_accuracy',
max_trials=5)

tuner.search(xtrain, ytrain, epochs=5, validation_data=(xtest, ytest))