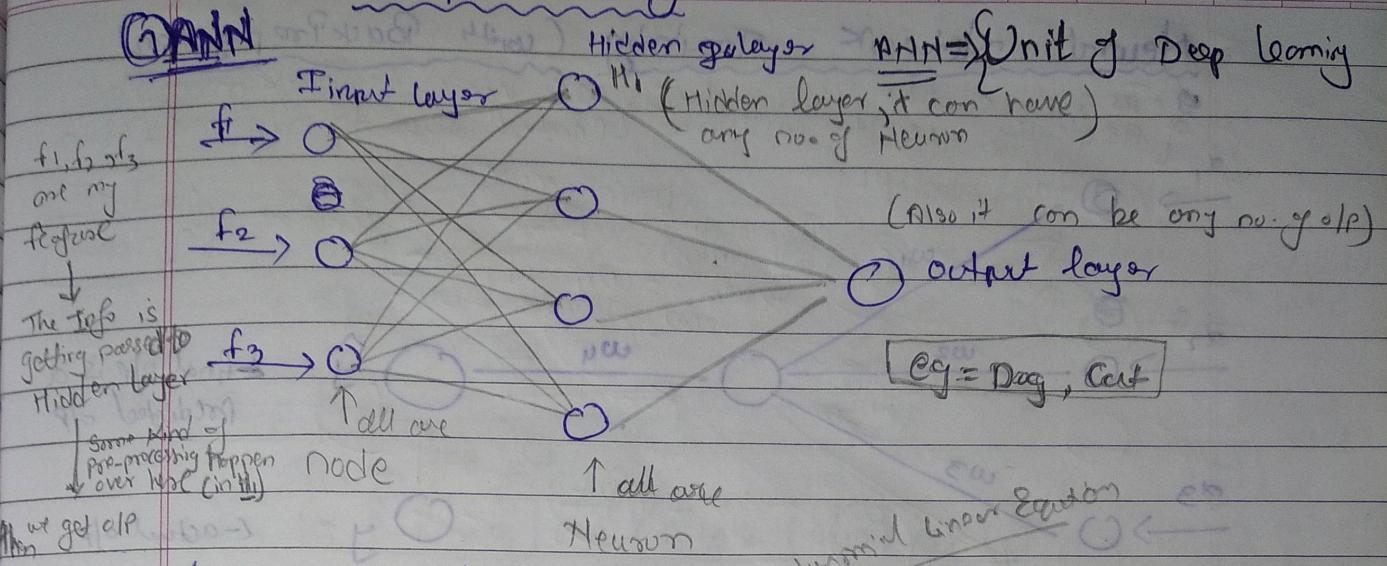


Deep Learning

① ANN



* How Neural Network Works

$$y = w_1x_1 + w_2x_2 + w_3x_3 + b$$

↑ I/P layer

$a_i \rightarrow$

w_1

x_2

x_3

$z = \text{Act}(y)$

If we are going to do summation (Σ) then passing to our sigmoid function (Step 1, Step 2)

O/P Layer

w_4

$z = z * w_4$

$= \text{Act}(z)$

Forward

Neuron Propagation

Step 1) weighted sum Step 2)

\sum (weight and I/P)

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

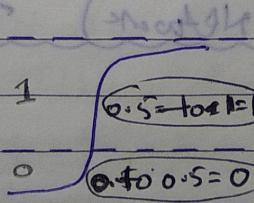
$\text{Act}\left(\sum_{i=1}^n w_i x_i\right)$

$\text{Act} \sum$ (weight and I/P)

* Activation Function

① Sigmoid Af

$$\frac{1}{1+e^{-y}}$$



Sigmoid function transform your y between 0 and 1

② Relu Af

$$\text{max}(y, 0)$$

$$y = \sum_{i=1}^n w_i x_i + b_i$$

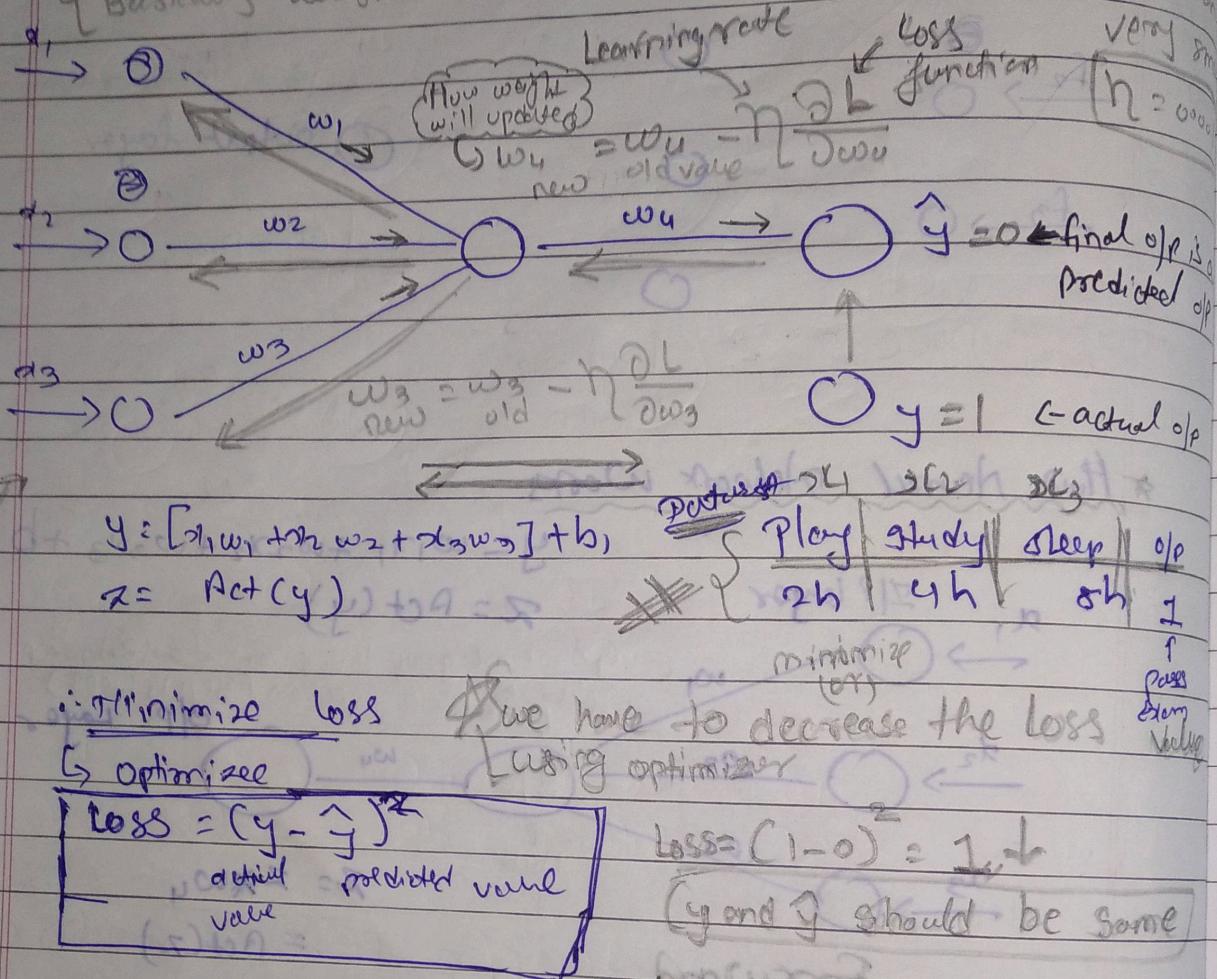
$$\text{Act}(y)$$

$$y = \sum_{i=1}^n w_i x_i + b_i$$

$$(x - v)^2 \text{max}(x - v, 0) = 0$$

$$(x - v)^2 \text{max}(x + v, 0) = +v$$

- * Neural Network Training (with Back Propagation)
 - To reduce the loss value, we have to back propogate
 - Basically weights will be updated due to back propagation



Note :- Sigmoid function transform the value between 0 to 1.

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \leftarrow \text{multiple record}$$

* Forward And Backward propagation

* MNN (Multi Layer Neural Network) Training

* Above NN Explanation

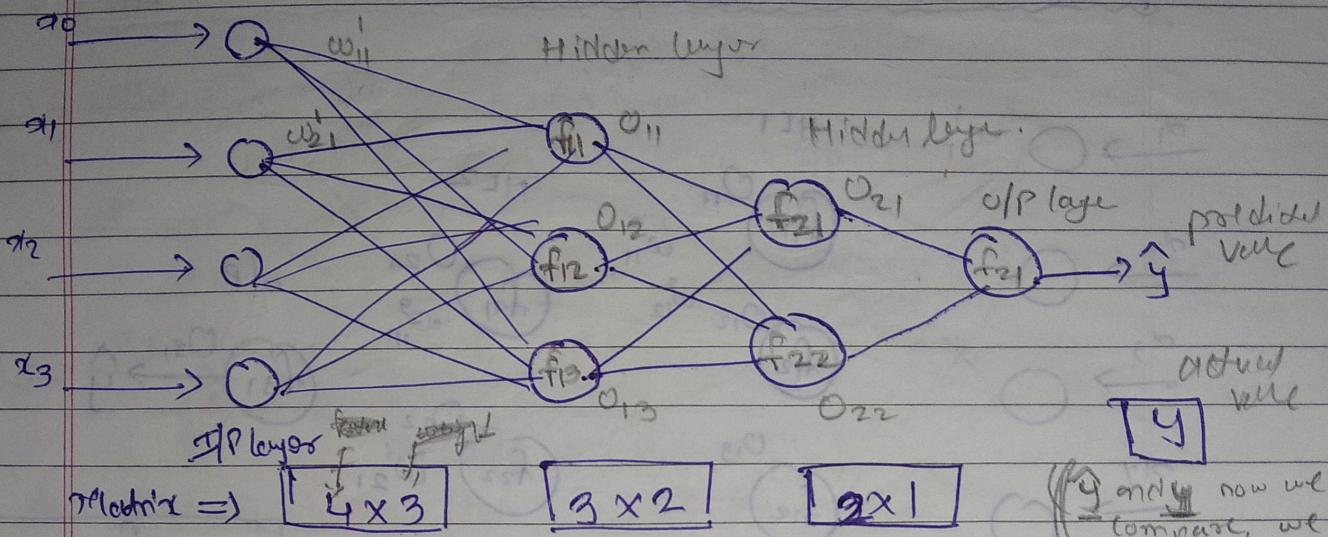
$\hat{y} \rightarrow$ predicted op

I already know that my training dataset the op for this kind of Data is $\rightarrow 1$. Suppose my \hat{y} is predicting as 0. we need to compare these whether \hat{y} and y are almost same. This difference can be found out using loss.

$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_{old}} \text{ (Gradient Descent)}$$

{ Gradient \rightarrow Slope
Decent \rightarrow Downward }

PAGE No.	/ /
DATE	/ /

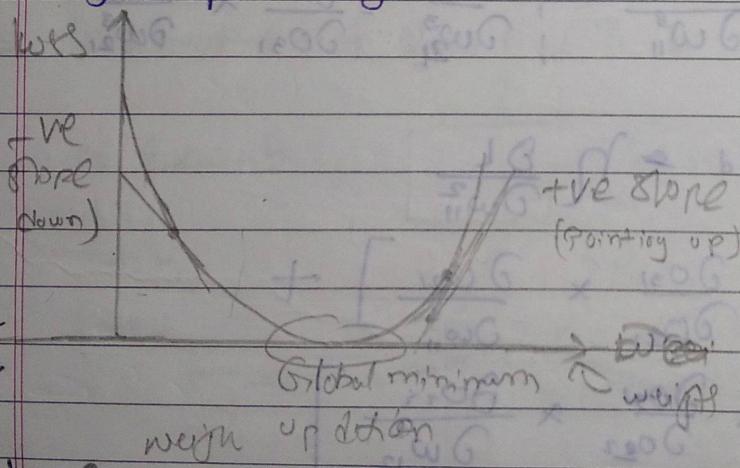


$$\text{Loss} = (y - \hat{y})^2$$

Derivatives means finding the slope.

The main aim is to reduce the loss value, for getting similar value y and \hat{y} .

↪ we use Optimizer (GD) $\Rightarrow w_{new} = w_{old} - \eta \frac{\partial L}{\partial w}$
weight update formula



The main aim of GD, is, we should be reaching to the Global minimum

Note { Unless and until we don't reach or get Global minima the update of weights will happen with the help of Back propagation }

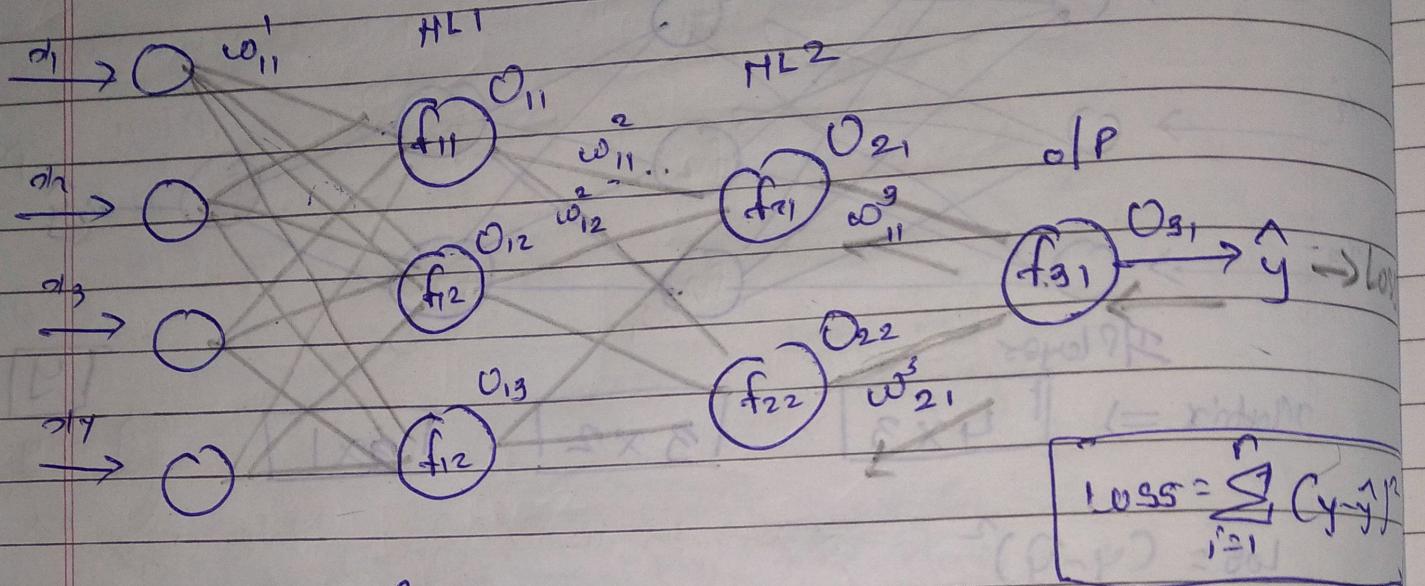
⇒ Once we get the op then work that we compute the Loss $L(y, \hat{y})$, the loss need to optimize using optimize

Notes $\frac{dy}{dx}$ what is this dy/dx ?

↪ we are going to find out Tangent or Slope over the single point

* Chain Rule In Back Propagation

$w_{11}^3 \rightarrow$ Impact the o/p of O_{31}



$$\text{LOSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$(i) \quad w_{11}^3 \text{ new} = w_{11}^3 \text{ old} - \eta \frac{\partial L}{\partial w_{11}^3}$$

update

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3} \quad ; \quad \frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial w_{21}^3}$$

$$(ii) \quad w_{11}^2 \text{ new} = w_{11}^2 \text{ old} - \eta \frac{\partial L}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{11}^2} = \left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}^2} \right] +$$

$$\left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial w_{11}^2} \right]$$

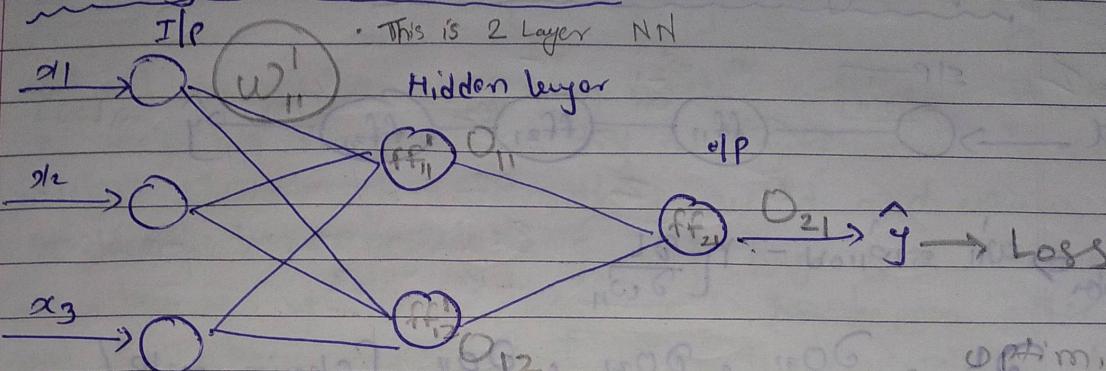
Note: Sigmoid transform the value between 0 to 1
 Note: The derivative of Sigmoid will be specifically between 0 to 0.25

PAGE NO.

DATE

7/1

Vanishing Gradient Problem



weight updation

$$\text{Bias UP dotted}$$

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

$$w'_{11, new} = w'_{11, old} - \eta \frac{\partial L}{\partial w'_{11, old}}$$

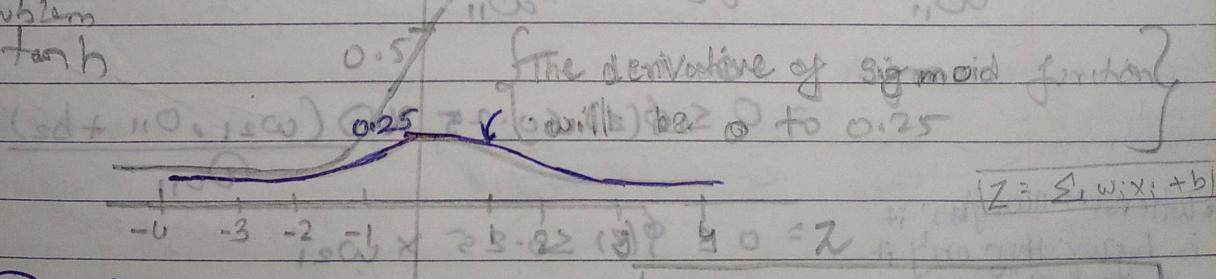
$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial o_{21}}{\partial w_{11}} \times \frac{\partial o_{11}}{\partial w_{11}} \quad (\text{chain rule})$$

Note: The derivative of Sigmoid will be value between 0 to 0.25

$$z = \sum w_i x_i + b$$

Some problem with tanh

$$\text{Sigmoid Function } \frac{1}{1+e^{-z}}$$



$$\frac{\partial \sigma(z)}{\partial z} > 0 \rightarrow 0.25$$

$$\text{range: } 0 \leq \sigma(z) \leq 0.25$$

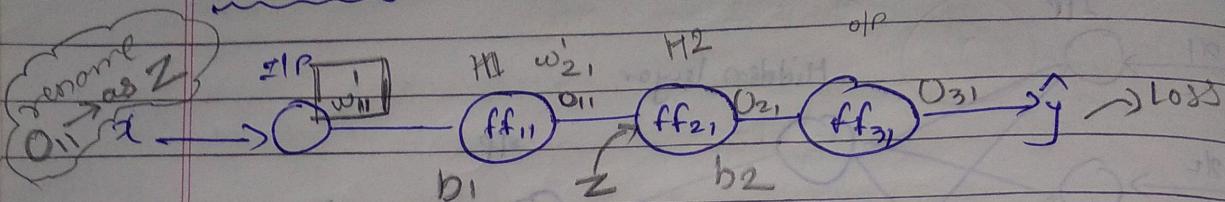
Before late 2000s, Researchers were not able to create a Deep Neural Network in ANN, the reason is that we have to use Sigmoid in each and every function in each and every neuron. Because of this Sigmoid we are facing problems called Vanishing Gradient Problem.

(Note: Activation function called ReLU was not invented)

Vanishing \Rightarrow Problem due to Activation function
 Exploding \Rightarrow Problem due to weight Initialization

PAGE NO. / /
DATE / /

Exploding Gradient Problem



weight up-gradation $\Rightarrow w_{11}' = w_{11} \cdot \eta \frac{\partial L}{\partial w_{11}}$ Exploding gradient problem takes here

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}} \quad [\text{chain rule}]$$

$$O_{21} = \phi(z) \quad \frac{1}{1+e^{-z}} \quad 0 \leq \phi(z) \leq 0.25$$

$$Z = w_{21} \cdot O_{11} + b_2$$

Note: Exploding Gradient Problem, it does not just happen because of ~~Sigmoid~~ Sigmoid function. The main reason this exploding gradient problem happens is because of weights. Let's take $\frac{\partial O_{21}}{\partial O_{11}}$ and let's compute this value (i.e. derivative value over here)

$$\frac{\partial O_{21}}{\partial O_{11}} = \frac{\partial \phi(z)}{\partial z} \times \frac{\partial z}{\partial w_{11}} \quad (\text{Assume Sigmoid Act function})$$

Consider my weight initialize here is 500.

$$\text{Then multiply } \Rightarrow 0.25 \times 500 = 0 \leq \phi(z) \leq 0.25 \times \frac{\partial (w_{21} \cdot O_{11} + b_2)}{\partial O_{11}}$$

• whenever weight if is higher than only it will perform "Exploding"

$$= 0 \leq \phi(z) \leq 0.25 \times w_{21}$$

Input If the weights are higher, what may happen you may converge to the global minimum point.

• If my weights are higher i usually get a higher value of derivative, but ~~the~~ derivative of sigmoid is between 0 to 0.25, but because of weights this derivative value is becoming larger and when it becomes larger because of the chain rule, as i compute a deep ANN

this particular derivative [i.e. $\frac{\partial L}{\partial w_1}$] with respect to w_1 ,

~~it~~ will become a very big number and when it becomes a very big number if i try to apply that in the weight updation formula, then word and w_{new} will be completely different will have huge gap between them and because of that, what will happen after each and every back propagation it will never reach the global minimum point. So that is why weight updation is important.

* Drop Out & Regularization [Refer Random Forest]

- whenever we have multilayered NN underfitting will never happen, because we will be having multiple layers in a multilayer NN.
- If we have one layer NN at that underfitting will usually happen. ~~always~~
- For multilayered NN we will never face underfitting but, overfitting will be a problem.

Note • There are two basic ways to solve and overfitting problem

(1) Regularization

(i) L1 (ii) L2

In RF we create
multiple DT

Random forest

DT \rightarrow overfitting

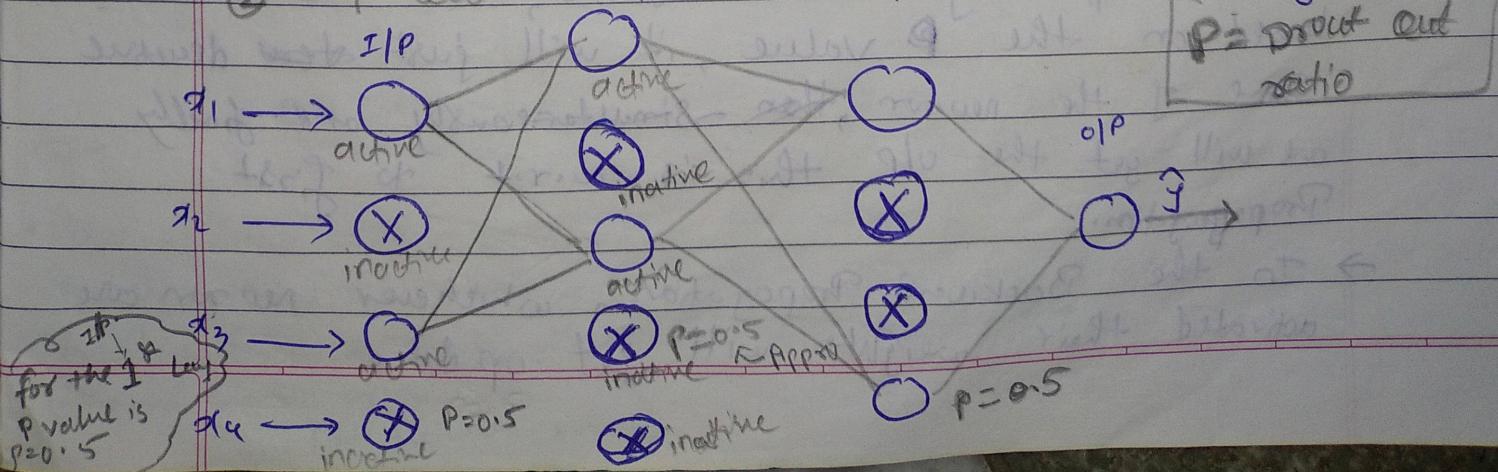
(2) Drop out

HL1

HL2

Subset of features.

P = drop out ratio



→ w.r.t Probability Features will be selected Randomly
P

PAGE NO.	/ /
DATE	/ /

- To implement drop out layer, we basically do is they select our dropout ratio
- Drop out ratio between $0 \leq P \leq 1$

The drop out ratio indicate that how we did for the random forest like a selected sample or subset of features similarly over here we'll be selecting subset of features from the i/p layer, Similarly we'll be selecting subsets of Activation function or the hidden neurons in the hidden layer, we will not be selecting everything we will just be selecting subset of features in each and every hidden layer along with the i/p layer features.

• Refer Diagram after Reading for this ↑

- when my forward and backward propagation will be going on when i select the $P=0.5$, it will ~~select~~ randomly ~~select~~ select some features & will deactivate them, when it deactivate them suppose in this case my first second and third node are deactivate, so my input will get pass all the activation function now in my second layer also in my first hidden layer you should see that again I've selected $P=0.5$
- Randomly it will select some of the Activation function some of the neuron over here, it will deactivate them all the processing will be same, after the best based on the P value, it will just ~~deactivate~~ deactivate some of the neuron, then simultaneously and finally we will get the o/p this is w.r.t to first Propagation.
- In the Backward Propagation, whichever neuron all activated their weights will get updated.

Note → for Test Data all our neurons get connected for test data, these will be no deactivated or activated neurons everything will be connected.

• $w \times p$ (wxp) will happen for each and every weights in each and every layer.

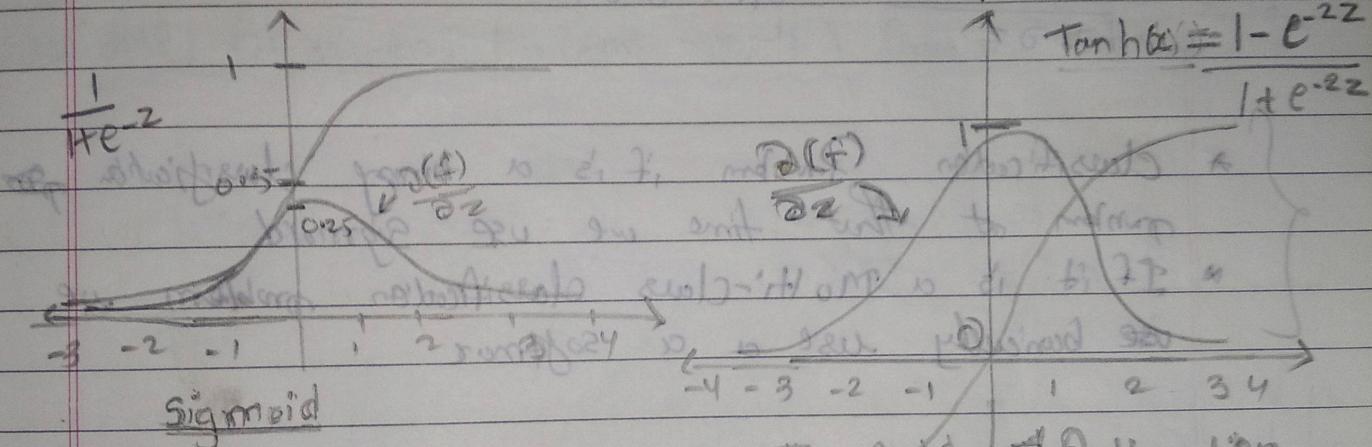
• for test Data just multiply wxp

• How we can select P-value?

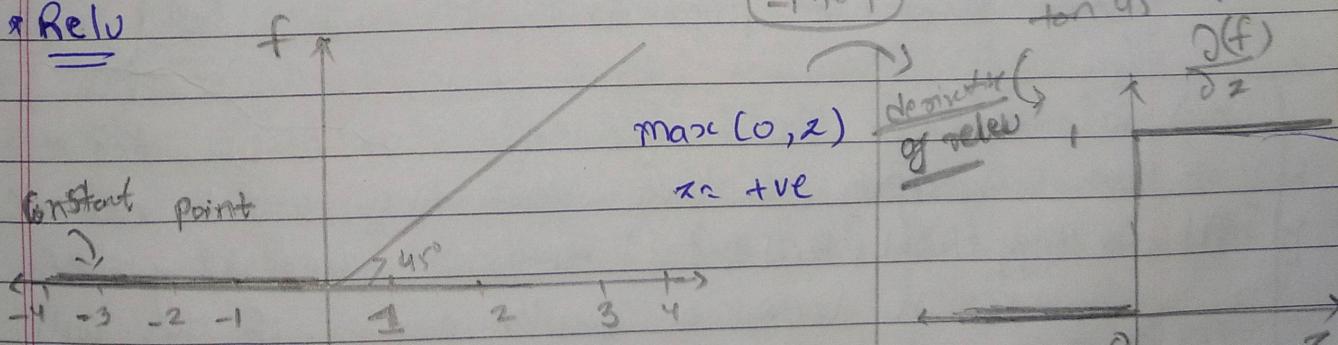
• Whenever a Deep NN is doing an overfitting, p-value should be little bit higher

* Rectified Linear Unit (ReLU)

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



* Relu



* @ to infinity

• Relu used in Hidden layer, Because, it solves the Problem of Vanishing Gradient Problem.
(Derivative $\rightarrow 1$ or 0 not between 1 or 0)

R. Derivative

$$\text{Here we have only 2 values } \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$\frac{\partial x}{\partial z} = 1$$

Leaky Relu

PAGE No. / / /
DATE / / /

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

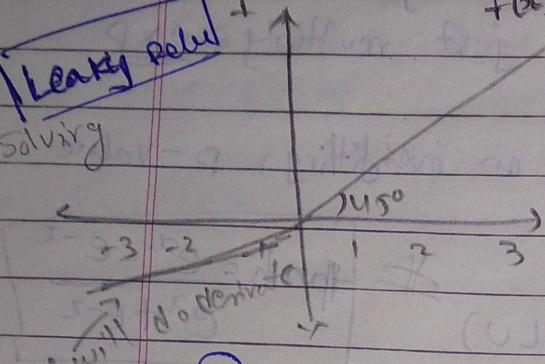
$$= (0 \times 1 \times 1) \text{ derivation value}$$

↳ Dead Neuron (Dead Relu) (or) Dead Activation function

$$f(z) = \max(0, 0.01z, z)$$

Leaky Relu

Solving



$$\left\{ \begin{array}{l} z \geq 0 \\ z < 0 \end{array} \right.$$

$$0.01(z) \quad z < 0$$

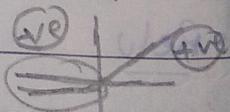
Small constant value

$$\frac{\partial (0.01) \times z}{\partial z} = 0.01 \quad (\Rightarrow 0.01)$$

↳ Here op will not be 0, like Relu

- * Classification Problem if it is a binary classification problem at that time we use Sigmoid
- * If it is a Multi-class classification problem we basically use a softmax

Note:- Sometimes Leaky Relu brings up sometime Vanishing Gradient problem, i.e. when $\approx w_{old}$



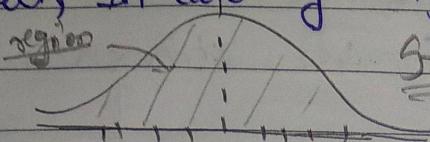
⇒ Reason - Suppose most of weights are (-ve), maximum no. weights are (-ve), that time wrt derivative ($\frac{\partial L}{\partial w}$) will be multiplying $(0.01, 0.01, 0.01)$ everywhere, because of these it will become a smaller number, if this will become a smaller no. ($\frac{\partial L}{\partial w}$) ← this will also become a smaller no.

⇒ When we updating along with (η) learning rate our $w_{new} \approx w_{old}$

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

Activation Functions - ELU, PReLU, Softmax, Swish, and Softplus

- The most specific thing about Normal Distribution Data or Gaussian Distribution Data, the Data is zero centred
- But, In case of Sigmoid it is not centred



↳ Data is zero centred

- tanh zero centred

$$\mu = 0$$

$$\sigma = 1$$

most of the data will be either (+ve) or (-ve), but it will be always be in this specific region which is centred towards 0

ELU (Exponential Linear Unit) function.

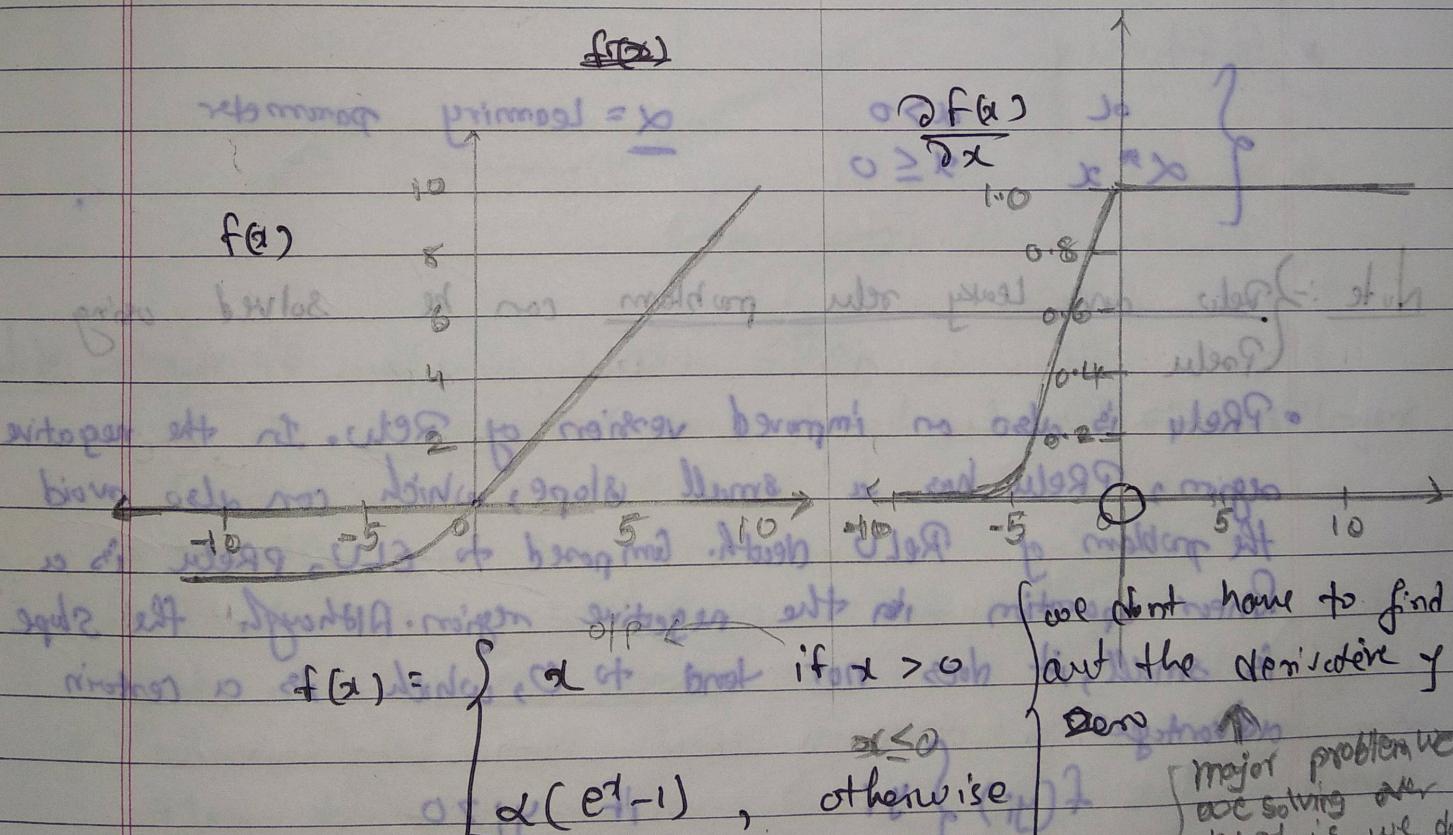
$$f(x)$$

$$\text{primoz} = x$$

$$0 \geq f(x)$$

$$0 \geq \frac{\partial f(x)}{\partial x}$$

$$x \in \mathbb{R}$$



$$\alpha(e^x - 1) \geq 0 \quad \forall x$$

$\alpha \leftarrow$ hyperparameter

learning parameter

it can be 0.01, 0.1, 0.3

problem have

$$\frac{\partial \alpha(e^x - 1)}{\partial x}$$

\propto determining w.r.t

Computationally Expensive
Some with tanh

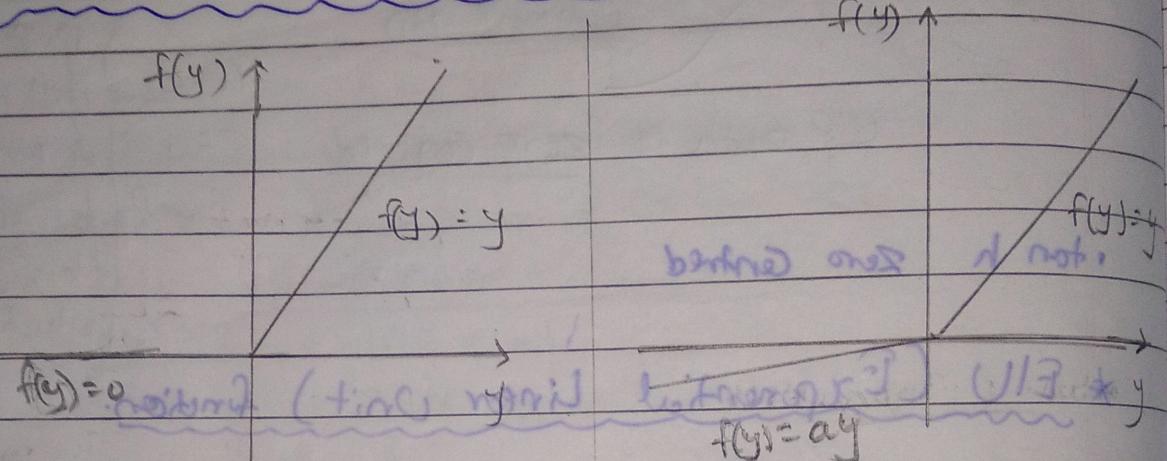
major problem we are solving over here is we don't have to find derivative of α

ELU is also proposed to solve the problem of RELU.
ELU has all the advantages of RELU

- No Dead RELU issues

The mean of the output is closer to 0, zero-centred.

PReLU (Parametric ReLU)



$$\left\{ \begin{array}{l} \text{or } \alpha > 0 \\ \alpha \neq 0 \end{array} \right.$$

α = Learning parameter

Note :- ReLU and Leaky ReLU problem can be solved using PReLU

PReLU is also an improved version of ReLU. In the negative region, PReLU has a small slope, which can also avoid the problem of ReLU death. Compared to ELU, PReLU is a linear operation in the negative region. Although the slope is small, it does not tend to 0, which is a certain advantage.

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

if $y_i \leq 0$ ($\alpha > 0$)
retaining negative \rightarrow

The parameter α is generally a number between 0 and 1, and is generally relatively small.

Such as in few cases. When $\alpha = 0.01$ we call PReLU as Leaky ReLU, it is regarded as a special case PReLU.

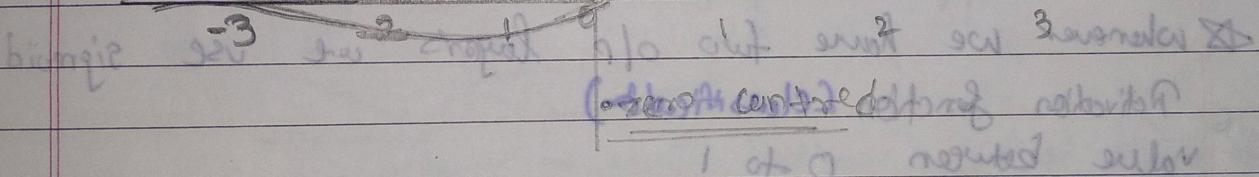
- Above, y_i is any input on the i^{th} channel and a_i is the negative slope which is a learnable parameter.

$\left\{ \begin{array}{l} \rightarrow \text{if } a_i = 0, f \text{ becomes ReLU} \\ \rightarrow \text{if } a_i > 0, f \text{ becomes Leaky ReLU} \\ \rightarrow \text{if } a_i \text{ is a learnable parameter, } f \text{ becomes PReLU} \end{array} \right.$

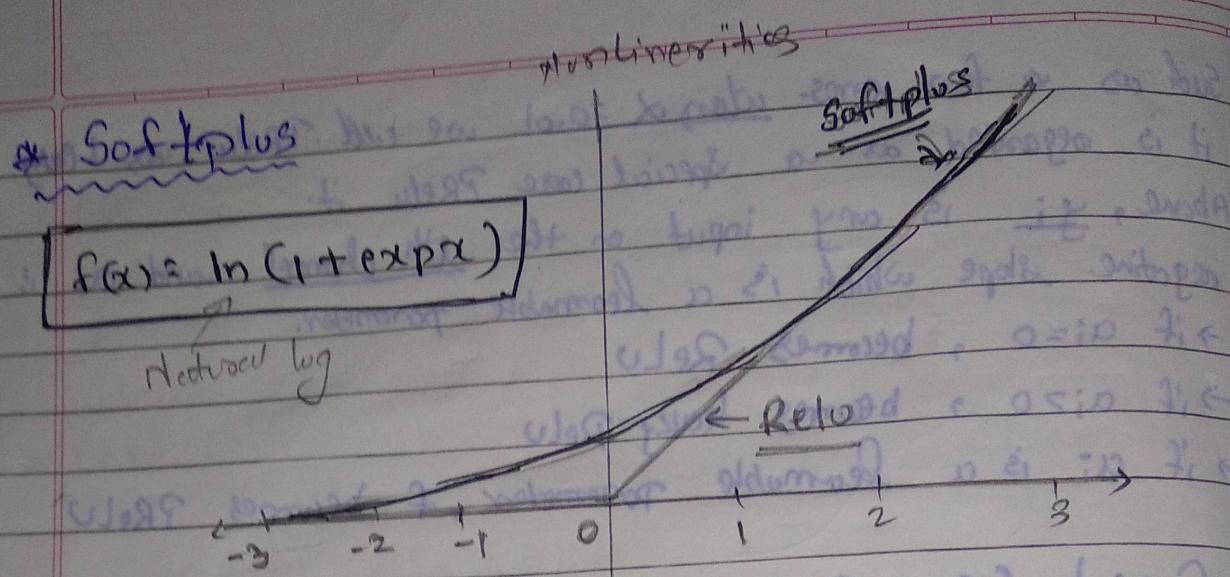
* Swish (A self-Gated) Function

$$\text{formula} = y_i = x_i * \text{Sigmoid}(x)$$

- Solving Dead ReLU problem



- Note → This is used in LSTM mostly
- We use Swish only when Neural Network > 4 layer (min greater 4 layer)
 - Less than 4 layer not work
 - The advantage of Self-Gating is that it only requires a simple input, while normal gating requires multiple scalar inputs.
 - Whenever we have Deep Neural Network, means min more than 4 layer we should use Swish function.



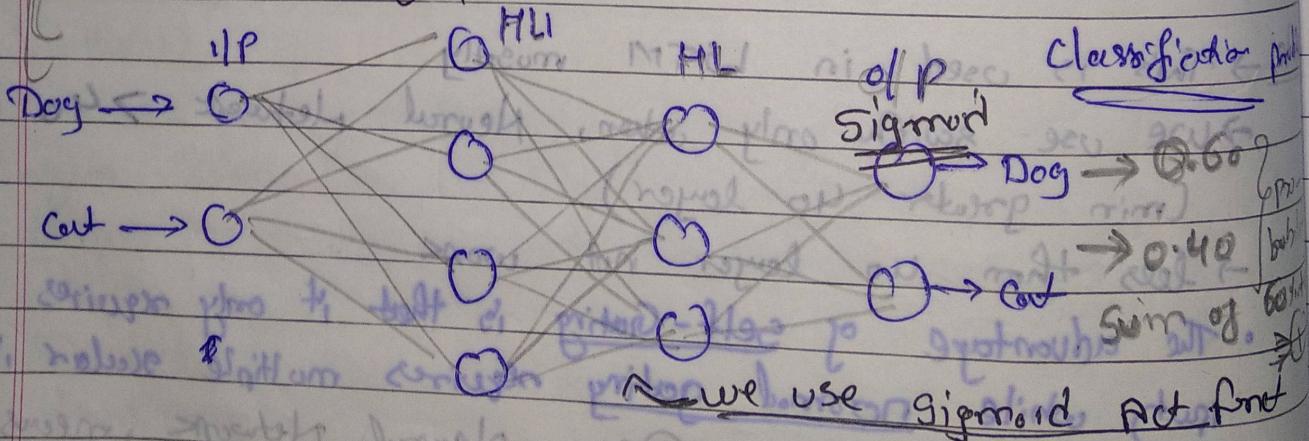
→ The softplus function is similar to the ReLU function, But it is relatively smooth. It is unilaterial suppression like RELU.

→ It has a wide acceptance range ($0, +\infty$)
→ Variant of ReLU only

* Softmax

Whenever we have two o/p layers we use sigmoid Activation function (Classification)

Value between 0 to 1



Hint:- Whenever we have more than 2 o/p layers we use softmax Activation function.

Note: Sigmoid and softmax Activation function are kept for last phase

PAGE No. _____
DATE 01/02/2013 / 01/04

softmax $\alpha_j = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}$
 $[0.6, 0.2, 0.1, 0.1]$
 $\sum = 1$

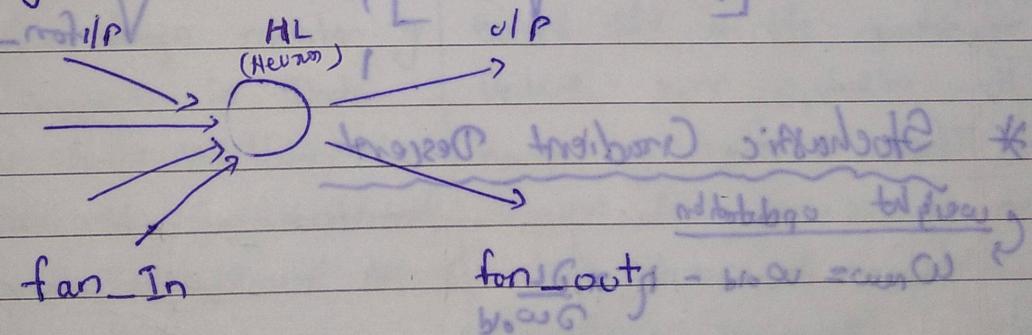
Image $\rightarrow O$

$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, j = 1, 2, \dots, K$$

$$\Rightarrow x_j = w^* i/p + b$$

* Weight Initialization Techniques in NN

- Key points => 1) weights should be small
- 2) weights should not be same
- 3) weights should have good variance.



(1) Uniform Distribution Technique

\rightarrow (work nicely with sigmoid fnct)

weights are basically selected from uniform distribution

$$w_{ij} \sim \text{Uniform}(-\frac{1}{\sqrt{\text{fan_in}}}, \frac{1}{\sqrt{\text{fan_in}}})$$

• weights were selected from uniform distribution with some raw and bias values. $w_{ij} \sim \text{Uniform}(-1, 1)$

② Xavier / Glorot Techniques (work nicely with Sigmoid / ReLU)

(Weights are taken from H1D with some value like 0.01)

① Xavier Normal

Normal Distribution $\mu = 0$, $\sigma^2 = \frac{1}{n}$

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fan_in} + \text{fan_out}}}$$

$$(d + g) / i \cdot w = 10$$

③ He init (work nicely with ReLU)

① He Uniform

weights are taken from uniform dist with some a and b values

$$w_{ij} \approx U \left[-\frac{\sqrt{6}}{\text{fan_in}}, \frac{\sqrt{6}}{\text{fan_in}} \right]$$

② He Normal

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fan_in}}}$$

* Stochastic Gradient Descent

weight update

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

* Suppose we have 1000 data points

$$\frac{\partial L}{\partial w_{\text{old}}}$$

→ 1 data point

→ Gradient Descent

At every epoch taking only 1 data point at a time

where k is always less than N (total datapoint)

$\hookrightarrow k$ data point → Mini Batch SGD

② Xavier or Uniform

Uniform Dists

$$w_{ij} \sim U(a, b)$$

Weights are taken from uniform distribution

$$a = -\frac{\sqrt{6}}{\sqrt{\text{fan_in} + \text{fan_out}}}, b = \frac{\sqrt{6}}{\sqrt{\text{fan_in} + \text{fan_out}}}$$

ReLU activation function

weights are selected from H1D with $u=0$ and $u>0$

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fan_in}}}$$

ReLU activation function

slope of not zero

slope of zero

Global min

Mini Batch SGD

$$\text{Loss} = \sum_{i=1}^k (y_i - \hat{y}_i)^2$$

In many CNN models like VGG, CNN we basically used this technique of mini batch SGD.

Gradient Descent

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

In GD we require Higher Computational memory

SGD

In linear Reg we used in ML

$$\text{Loss} = (y_i - \hat{y}_i)^2$$



GD

→ In Gradient Descent, if there will be a huge number of records suppose in Data

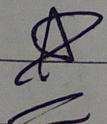
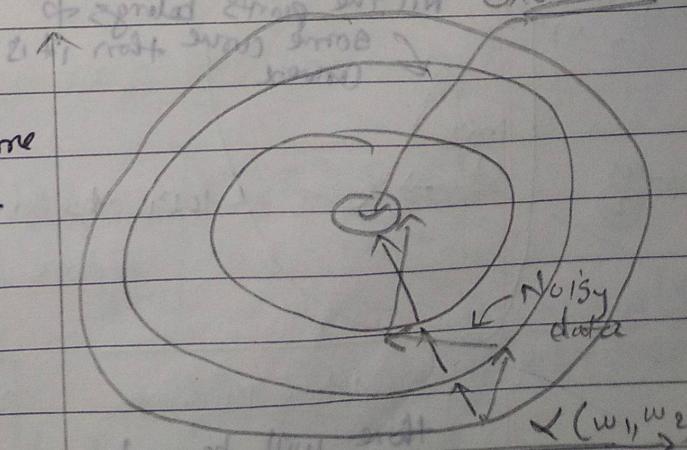
Let you have 10 million records so it becomes computationally powerful computation will require a lot of resources

• GD will take less time to reach Global Min point

• Mini Batch SGD will take some time to reach Global min point

Mini Batch SGD

→ In case of Mini Batch SGD Based on the particular K value your updation will be happening and it will also require less number of resources when compared to the Gradient Descent. Suppose this is the point from which to reach Global min.



$$\left[\frac{\partial L}{\partial w_{i,0}} \right] \underset{\text{Sample}}{\approx} \left[\frac{\partial L}{\partial w_{i,k}} \right]_{\text{Population}}$$

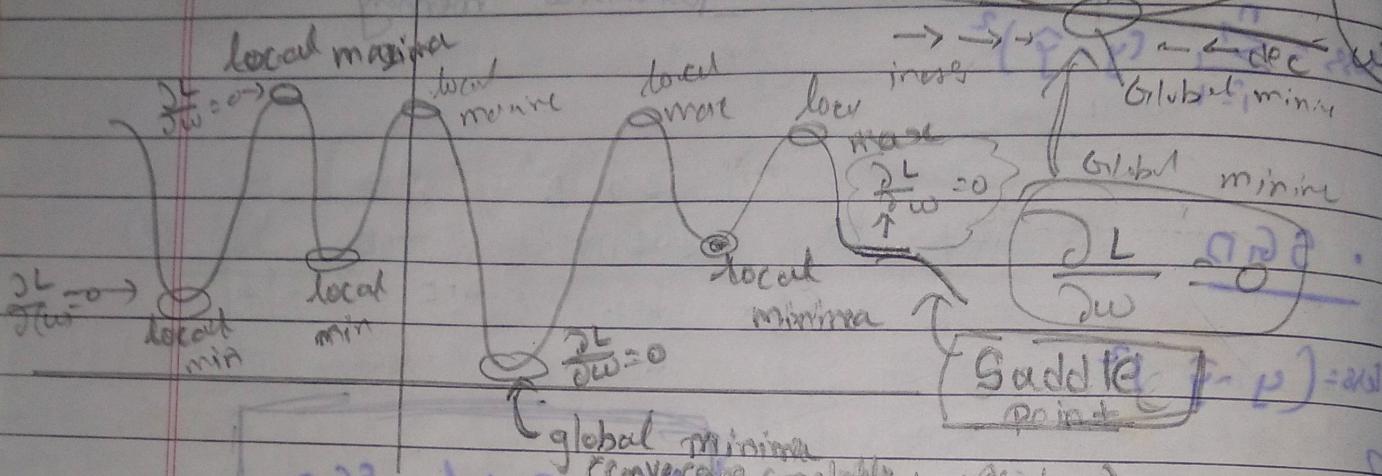
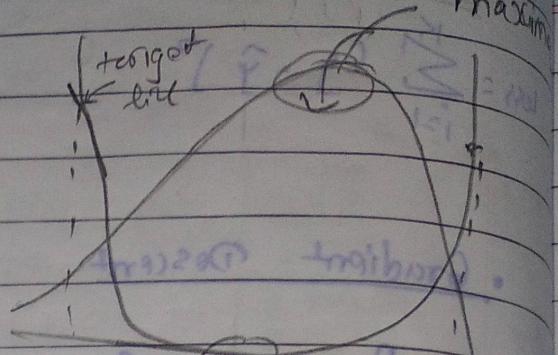
zig-zag \Rightarrow mini Batch

straight line \Rightarrow (GD)

* Global Minimum ∇ Local Minimum and global maximum

optimizer \Rightarrow GD, SGD, mini Batch SGD

$$L(w) = \sum_{i=1}^k (y_i - \hat{y}_i)^2$$



Saddle point $\frac{\partial f}{\partial w} = 0$

(converging completely towards down)

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

\rightarrow weight update formula

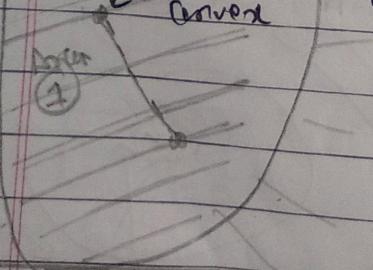
① Convex function

Linear Reg, Logistic Reg etc

→ usually occurs in non-linear ml

curve happens if of the DL techniques

All the points belongs to some curve then it is convex



Global Minima

there will be only

1 Global Minima

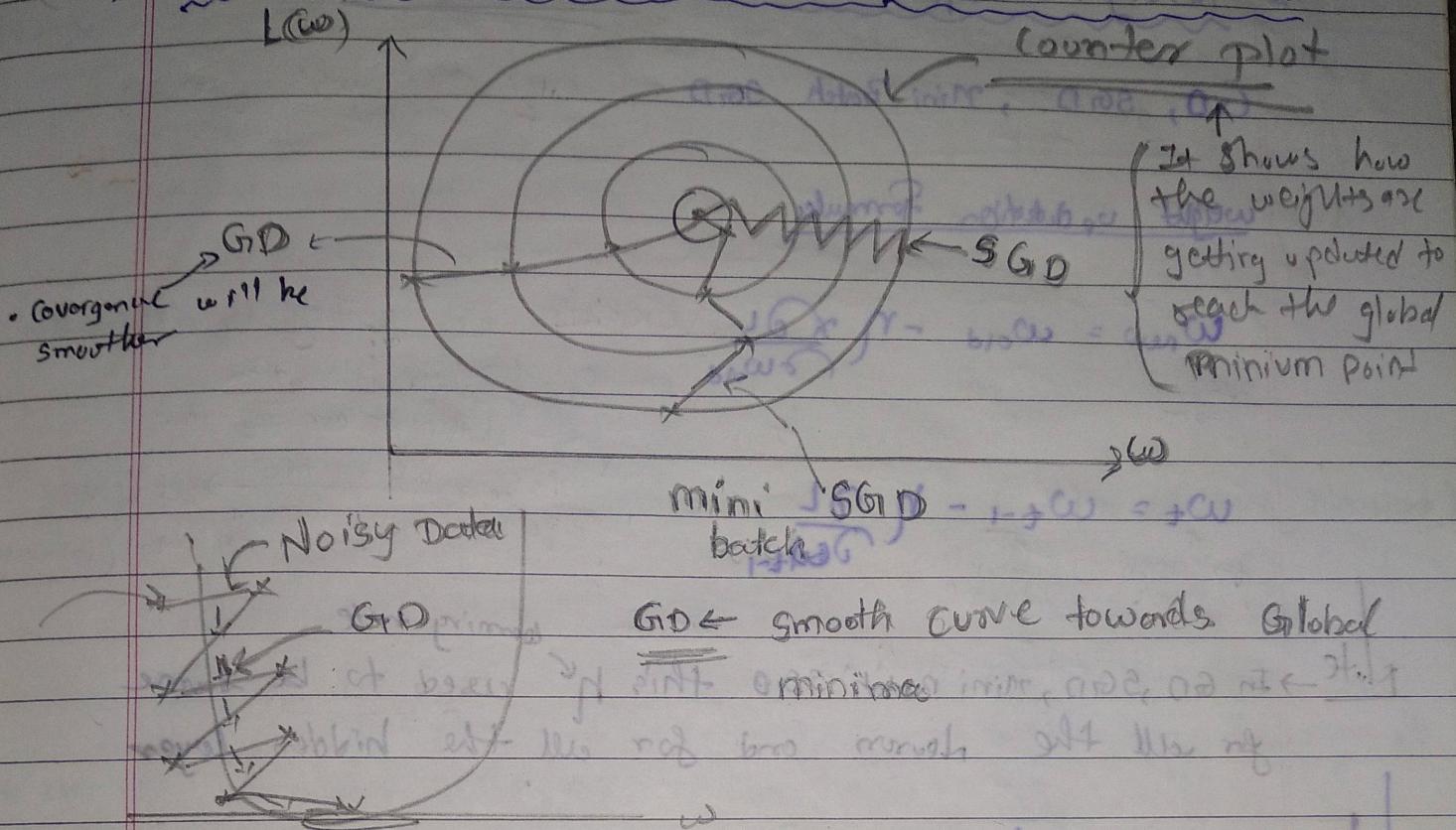
Convex function \Rightarrow always occur in ML like LR, LR

points are present in both region

W.C.W region

\rightarrow due to overfitted minima so many weights parameter so many neuron we are using each and every neuron needs to converge

Stochastic Gradient Descent with Momentum



Noisy Data will take more time for convergence i.e. it

takes more time to reach Global minimum.

→ How noisy Data can be removed with the help of

Exponential ^{weighted} moving average

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

Gamma

$$\gamma = 0.9$$

appropriate value

Value decreases the noise as we are going towards global minimum.

Momentum \leftarrow global minimum point

$$w_t = \left[\gamma v_{t-1} + \eta \frac{\partial L}{\partial w_{\text{old}}} \right]_{(t-1)}$$

$$0 \leq \gamma \leq 1$$

$$w_{\text{old}} = v^*$$

$$v_{t-1} = \gamma \left[\frac{\partial L}{\partial w_{\text{old}}} \right]_t + \gamma \left[\frac{\partial L}{\partial v_{t-1}} \right]_{t-1} + \gamma^2 \left[\frac{\partial L}{\partial w_{\text{old}}} \right]_{t-2}$$

forward direction this equation $= +v$

Exponentially increase

Adagrad Optimizers Adaptive Gradient Optimizer

weight update formulae

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \left\{ \begin{array}{l} \text{weights are getting updated} \\ \text{on epoch, or iteration} \end{array} \right.$$

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

$\uparrow \text{current time}$ $\downarrow \text{previous time}$

↓ then we can say that

$\begin{cases} w_{\text{new}} & \xrightarrow{\text{we can take}} w_t \text{ (current)} \\ w_{\text{old}} & \xrightarrow{\text{}} w_{t-1} \text{ (prev)} \end{cases}$

learning rate

Note → In GD, SGD, mini Batch SGD this η used to be same for all the neuron and for all the hidden layer

Note → The idea behind Adagrad optimizers is that, It should be different Learning Rate

Idea → we can use Different Learning Rate for each and every neuron and each and every hidden layer based on different - different iteration, weights etc.

→ Two types of features

① → Dense (Most of the features will be Non-zero)

② → Sparse (Most of the values are basically 0)

e.g.: Bow [feature value in 0 and 1 or 0.1 (Decimal)]

* Adagrad Equation

$$w_t = w_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}}$$

• t specifies w.r.t. iteration (like previous iteration, current iteration previous to previous iteration)

$\Delta_t \uparrow n_t$ $n_t \leftarrow \text{dash}$ n_t \downarrow constant learning rate
initialize $n = 0.01$

we notice this according to the Adagrad case

$$\sqrt{\Delta_t + \epsilon}$$

Patra

Q why do we required ' ϵ ' (Episoden) value

→ If the Δ_t value is "0" the whole value will be 0 then n_t will be 0 then, $w_t = w_{t-1} - \frac{n_t}{\Delta_t} \times \frac{\partial L}{\partial w_t}$ will be 0

$$w_t = w_{t-1} - \frac{n_t}{\Delta_t} \times \frac{\partial L}{\partial w_t}$$

$$w_t = w_{t-1} \quad (\text{Equal})$$

Notes to prevent this problem we use ϵ (small +ve number)

$\epsilon \rightarrow$ small +ve numbers

$$\Delta_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$$

here $\Delta_t \rightarrow$ will become high no because we squaring it

write as Δ_t here we are doing summation for all the iteration

→ Learning rate n_t decrease, As the iteration goes on

→ Weights Decreasing slowly as we go on

use :- we can see that n values changing w.r.t the factors, w.r.t different layers, w.r.t iteration.

Initial $\rightarrow w$

Initial it jump big step then goes on decreasing

Based on the no. of iteration

w goes on decreasing

finally converge global minimum point

Disadvantage (1) As the no. of iteration increase ↑
 At some time Δ_L will be very high ↑↑↑
 $\Delta_L = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$

↑
 To fit
 Training set
 We use
 Momentum to
 overcome this

* Adadelta and RMSprop Optimizer

$$\underbrace{\text{Adadelta}}_{\substack{\text{learning rate} \\ \text{should not become} \\ \text{very very small}}} \quad w_{\text{new}}^{(t)} = w_{\text{old}}^{(t-1)} - \eta' \frac{\partial L}{\partial w_{\text{old}}^{(t)}} \quad \begin{array}{l} \text{initial rate} \\ \text{work time} \\ \hookrightarrow \text{because it will change with each} \\ \text{iteration} \end{array}$$

$\eta' = \eta \rightarrow \text{constant}$

$\eta' t \quad \text{where, } \eta' t \approx \text{small +ve}$

$\text{when } \eta' \Delta_L = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$

- we taking ϵ because when $\Delta_L = 0$, then we trying to add ϵ value. so that anything suppose if we don't have the Δ_L will become 0. ~~it's possible~~ ~~it's not possible~~
- So anything Divide by 0 will give you a very big error to handle this we will use ϵ (small +ve no.)
- $\Delta_L \uparrow \eta' \uparrow$

As our NN becomes much deeper in Adadelta we use see it's that the learning rate became smaller, hence the convergence become

- To handle the Disadvantage we use
- Adadelta & RMSprop

Optimizers → update the w in the Backpropagation

- (1) GD → we take whole dataset \rightarrow (3) Mini Batch \rightarrow Batch Size
- (2) SGD → 1 Data at a time

To prevent this

$$\eta_t = \frac{0.95}{\text{Gammal}}$$

$$\eta_t = \frac{n}{\sqrt{\text{Wavg}_t + \epsilon}}$$

η_t will decrease very very slowly.

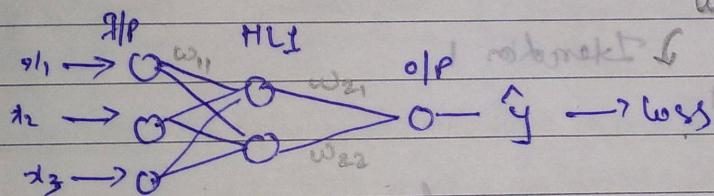
$\text{Wavg}_t = \text{Wavg}_{t-1} + (1-\gamma)(\frac{1}{n}w_t)$

weighted average at time t

Note:- This prevent η_t to become very very higher

* All Optimizers in One Video - SGD with Momentum, Nesterov, Polyak, RMSprop, Adam optimizers.

① Gradient Descent Optimizer



we calculate the error

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

→ we reduce the loss function with the help of optimizer

→ weight update formula

$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_{old}} \quad (\text{Gradient Descent})$$

① Epoch \rightarrow Complete Dataset is taken in GD

② Iteration

② SGD (convergence will be very slow)

Parallel \rightarrow Taking 1 record at a time.

Note:- Suppose I have 10k records,

for every Epoch, I have to do 10k Iterations of forward and backward

for, 1 Epoch \rightarrow 10k iteration

for, 20 Epoch \rightarrow 20 \times 10k iteration

How many no. of time you are training the Model

Epoch

④ Mini Batch SGD

→ suppose we have 10K records

→ for Epoch we will take Batch size

- Taking Batch size = 1000

PAGE No. _____
DATE _____

Convergence will be in zig-zag

Noise

Noise we are getting
Because we are creating a Batch Data

Dataset

20 Epoch

1 Epoch → 10 Iterations

10K
1000

10,000
1000

Batch size

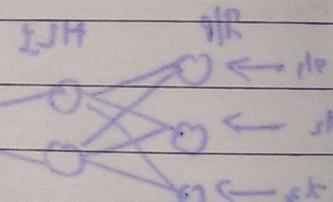
10 forward and Backward propagation

Epoch And Iteration

Suppose 10K Records

② SGD

① Epoch 1 $\xrightarrow{\text{forward}}$ Iteration 1 $\xleftarrow{\text{backward}}$



Epoch 1 $\xrightarrow{\text{Iteration 1}} \xrightarrow{\text{Iteration 2}} \dots$ based on no. of records

③

Mini batch SGD $\xrightarrow{\text{Batch size}}$ (Batch size)

Epoch \Leftrightarrow Iteration \Rightarrow 10 Iterations = $\frac{10,000}{1000}$ = 10 Iterations

- for every Epoch, for 10K records there will be 10 iterations

To restart if 100% of send is dropped from the last

Implementation details

Initially $V_{dw} = 0$ $V_{db} = 0$

On iteration t inside Epoch

compute $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial b}$ on the current minibatch

Moving Exponential average

$$V_{dw_t} = \beta V_{dw_{t-1}} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

$$V_{db_t} = \beta V_{db_{t-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

$$w_t = w - h^x V_{dw} \quad \rightarrow \text{smoothness of convergence}$$

$$b_t = b - h^x V_{db}$$

* Adam Optimizer (Adaptive Moment Estimation)

Adam combine i) Momentum (for smoothing) ii) RMS prop (Learning Rate) changing Learning Rate is efficient manner so that L value does not go High

Implementation detailed

Initially $V_{dw} = 0$ $V_{db} = 0$ $S_{dw} = 0$ $S_{db} = 0$

On iteration t ,

i) Compute $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial b}$ using current Mini-batch

ii) $V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \frac{\partial L}{\partial w}$ } w.r.t momentum

$V_{db} = \beta_1 V_{db} + (1-\beta_1) \frac{\partial L}{\partial b}$ }

$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \left(\frac{\partial L}{\partial w} \right)^2$ } w.r.t RMS prop

$S_{db} = \beta_2 S_{db} + (1-\beta_2) \left(\frac{\partial L}{\partial b} \right)^2$ }

$$w_t = w_{t-1} - \frac{h * \text{Value}_{\text{Exp moving avg}}^{\text{correction}}}{\sqrt{s_{dw} + \epsilon}}$$

initial h

$$b_t = b_{t-1} - \frac{h * v_{db}^{\text{correction}}}{\sqrt{s_{db}^{\text{correction}} + \epsilon}}$$

Bias correction

$$v_{dw}^{\text{correction}} = \frac{\text{Value}}{(1-\beta_1^t)}$$

Bias correction w.r.t Momentum

$$v_{db}^{\text{correction}} = \frac{v_{db}}{(1-\beta_2^t)}$$

$$s_{dw}^{\text{correction}} = \frac{s_{dw}}{(1-\beta_2^t)}$$

$$s_{db}^{\text{correction}} = \frac{s_{db}}{(1-\beta_2^t)}$$

These 2 forms
is the Updation

Adam Optimizer

$$\sigma^2 = \sigma^2 + \alpha - \omega^2 \quad \sigma = \omega V \quad \sigma = \omega b V$$

$$\text{Actual initial bias } \frac{16}{w_6} + \frac{16}{w_5} \text{ (avg)}$$

$$\left. \begin{array}{l} \text{initial bias } \frac{16(9-1) + \omega b V \cdot q}{w_6} = \omega b V \\ \frac{16(9-1) + \omega b V \cdot q}{w_5} = \omega b V \end{array} \right\}$$

* Loss function (Different types)

- Loss function = It is basically only for ~~only~~ meant for forward and backward propagation

$$\text{Loss func} \Rightarrow \frac{1}{2} (y - \hat{y})^2$$

- Cost function = we are passing Batch size of records

$$J = \frac{1}{t} \sum_{i=1}^t (y_i - \hat{y}_i)^2 \quad \dots t = \text{batch size}$$

(i) Squared Error loss

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

(MSE) Mean Squared error loss

$$\left(\sum_{i=1}^t (y_i - \hat{y}_i)^2 \right) / t = \text{MSE}$$

• Properties of squared error loss

- ① It is in the form of Quadratic Equation $ax^2 + bx + c$

- ② Plot the quadratic equation, we will get a Gradient Descent with only Global minimum
- ③ we don't get any local minimum

- ④ The Ridge loss penalizes the model for having large error by squaring them

- Disadvantage MSE loss function

→ it is not robust to outliers

① Loss function in Regression

Advantages & Disadvantages

(i) Squared Error loss

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

Convex
to Cost function

$$J = \frac{1}{t} \sum_{i=1}^t (y_i - \hat{y}_i)^2$$

.. t = Batch size

② Absolute Error Loss

local minima

$$\rightarrow L = |y - \hat{y}| \quad \text{absolute value}$$

$$\rightarrow J = \sum_{i=1}^n |y_i - \hat{y}_i| \quad \text{modulus}$$

Adv (i) L is more robust to outliers as compared to MSE

Dis (i) Computation is very difficult

③ Huber Loss

Combine both (i) MSE (ii) MAE most hardware friendly

$$\text{Loss} = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}, \quad \text{Quadratic Equation}$$

$$\text{out} \rightarrow \begin{cases} \text{Quadratic Equation} & \text{if } |y - \hat{y}| \leq \delta \\ \text{Linear Equation} & \text{otherwise} \end{cases}$$

Hyper parameter

④ Classification

④ Cross entropy

used in Logistic Regression specifically

$$\log loss \Rightarrow \text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

for computing \hat{y}
we use sigmoid

$$\begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \text{sigmoid}(z) = \frac{1}{1+e^{-z}}$$

new row between (i)

⑤ Multiclass Cross Entropy Loss

$$L(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \log(\hat{y}_{ij}) \quad \left. \begin{array}{l} \text{(Categorical)} \\ \text{Cross Entropy} \end{array} \right\}$$

i = row number

↳ Here y_i is one hot encoded target vector

$$\text{o/p} \rightarrow \text{Multiclass} \quad \left. \begin{array}{l} \text{y}_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}] \end{array} \right\}$$

Ground, Bad, Neutral → One hot Encoding.

Suppose i have f factors

	f_1	f_2	f_3	$0 \mid P$	Bits {internal}
2	3	4	5	Ground → [1 0 0]	it will be
4	5	6	7	Bad → [0 1 0]	in Bits
7	8	9	1	Neutral → [0 0 1]	

	f_1	f_2	f_3	y_{11}	y_{12}	y_{13}
2	3	4	5	[1 0 0]	0	0
4	5	6	7	0	1	0
7	8	9	1	0	0	1

→ 3 bits
vector

$$y_{ij} = \begin{cases} 1 & \text{if the } i\text{-th element is in class } j \\ 0 & \text{otherwise} \end{cases}$$

$\hat{y}_{ij} \Rightarrow$ softmax activation function { Multiclass classifier }

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

$$Z = W^T x + b$$

• Negatives of Loss Function (Classification most common)

→ Probabilistic losses

(i) Binary Crossentropy class [informally they will be doing one hot-encoding]

↳ Specifically used for Binary classes

$$g = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

(ii) Categorical Crossentropy class [and every class we have 0/p as 3 categories]

$$g_p \hat{y} = \begin{bmatrix} 0.3 & 0.2 & 0.5 \end{bmatrix}$$

{ which will be high for class 2 }
{ will be taking that particular class }

(iii) Sparse Categorical Crossentropy class

Here will be get Index as a output

$$\hat{y} = \begin{bmatrix} 0 & 1 & 2 \\ 0.3 & 0.2 & 0.5 \end{bmatrix}$$

\hat{y} will be maximum probability index [2]

* Environments

$$d + r^T w = 5$$