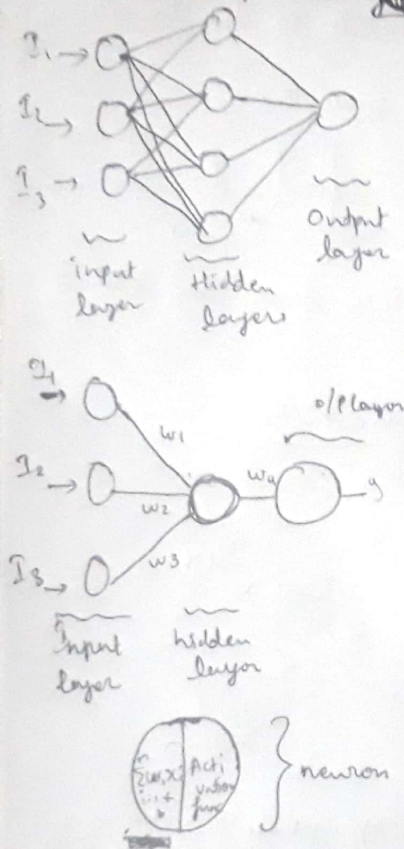
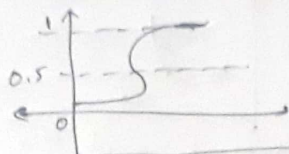


Deep learning



Activation function

① Sigmoid function



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

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

w = weights
 x = input
 b = bias

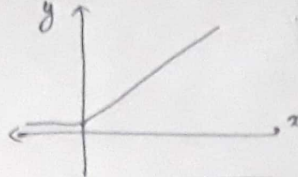
- if sig. func value ≤ 0.5
- if sig. func value > 0.5

0 \rightarrow not activated
1 \rightarrow activated.

$$① y = w_1 I_1 + w_2 I_2 + w_3 I_3 + \text{bias}$$

$$② z = \text{Act}(y)$$

② ReLU - Rectified Linear Unit



$$\text{Relu} = \max(z, 0)$$

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

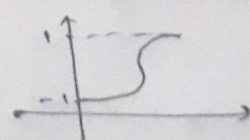
if Relu output is -ve

if Relu output is +ve

{ That corresponding value

• ReLU popular than Sigmoid

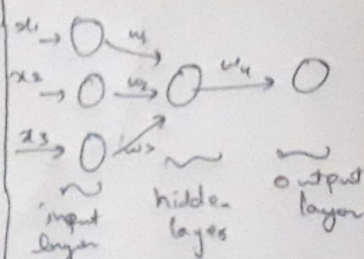
③ Hyperbolic tangent



range = -1 to 1

$$\tanh x = \frac{\sinh x}{\cosh x}$$

Back propagation



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

y = actual
 \hat{y} = predict

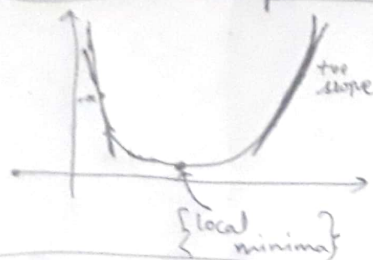
- We try to reduce the loss \rightarrow optimisers
- We back propagate and change weights acc to optimiser.

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

α = learning rate

- Once the weight is modified again we do forward propagation and calculate loss.
- No. of times we do this depends on no. of epochs.

↳ Gradient Descent - Optimiser



$$W_{new} = W_{old} - \alpha \frac{\partial L}{\partial W}$$

when slope is -ve i.e. $\left\{ \frac{\partial L}{\partial W} = -ve \right\}$

W_{new} = increases

- then it starts to converge to the local minima.
- How fast it does depends on α

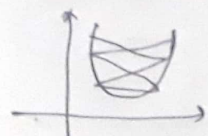
when slope is +ve i.e. $\left\{ \frac{\partial L}{\partial W} = +ve \right\}$

W_{new} = decreases

- then it starts to converge to local minima
- How fast it does depends on α

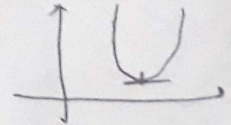
↳ α 'Learning Rate': needs to be lesser.

Higher α



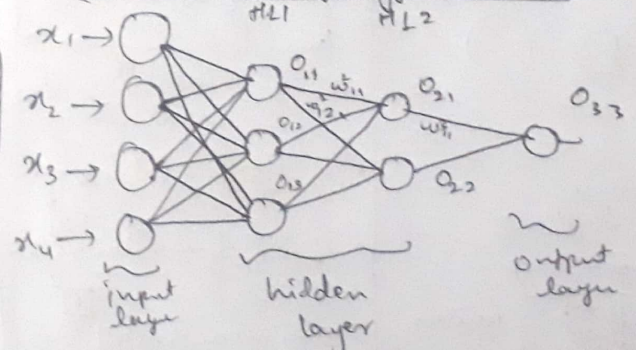
does not reach the local minima

Lower α



Takes time, but best is

↳ Chain Rule of Differentiation:



To optimise, we backpropagate & change the weights using optimiser to: $W_{new} = W_{old} - \alpha \frac{\partial L}{\partial W}$

Now how do we compute $\frac{\partial L}{\partial W}$ {using chain rule}

$$\frac{\partial L}{\partial W_{11}} = \frac{\partial L}{\partial O_{33}} \cdot \frac{\partial O_{33}}{\partial W_{11}}$$

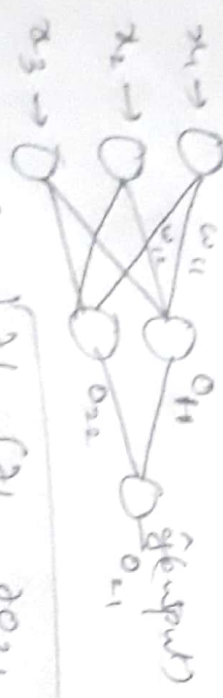
Chain rule
see where the weight is getting affected and change it

$$\text{eg: } \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{33}} \cdot \frac{\partial o_{33}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}^2}$$

But that neuron has also ~~connected~~ has another edge, so:

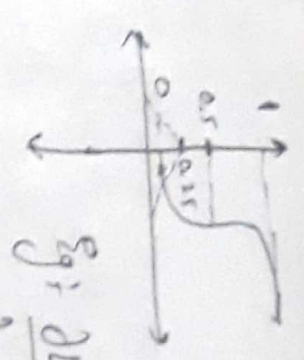
$$\frac{\partial L}{\partial w_{11}^2} = \left[\frac{\partial L}{\partial o_{33}} \cdot \frac{\partial o_{33}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}^2} \right] + \left[\frac{\partial L}{\partial o_{33}} \cdot \frac{\partial o_{33}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial w_{11}^2} \right]$$

→ Vanishing gradient problem in Sigmoid



$$\frac{\partial L}{\partial w_{11}} = \left[\frac{\partial L}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}} \right] + \left[\frac{\partial L}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{12}} \cdot \frac{\partial o_{12}}{\partial w_{11}} \right]$$

When the activation function is sigmoid the derivative of it ranges from $\boxed{0 \text{ to } 0.25}$



$$\frac{\partial \sigma(z)}{\partial z} \rightarrow \boxed{0 \text{ to } 0.25}$$

$$\text{eg: } \frac{\partial L}{\partial w_{11}} = \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}}$$

$$0.20 \times 0.02 = 0.004$$

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

$$w_{\text{new}} = 2.46$$

$$= 2.5 - (0.04)$$

As the layers increase, the term $\frac{\partial L}{\partial w}$ becomes smaller so that the effect of it is very very small becoming

$$w_{\text{new}} = w_{\text{old}}$$

$\frac{\partial L}{\partial w}$ is affecting very very small so

we call it as vanishing gradient.

So the vanishing gradient problem is due to sigmoid

Exploding gradient problem:-

- In exploding gradient problem, generally the weight is given higher value due to which $\eta \frac{\partial L}{\partial w_{ij}}$ has a bit higher value than the w_{new} & w_{old} will vary a lot.
- So that the gradient descent does not converge to local minima.

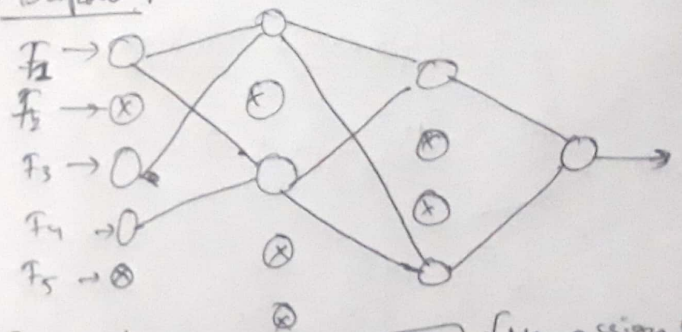
Drop out & Regularization:-

- When more layers are there in Neural Network, overfitting tends to occur. But we will never face underfitting.

• ~~Ways to solve Regularization~~ \rightarrow $P(L_1, L_2)$

• Way to solve Overfitting \rightarrow Regularization L_1, L_2
Dropout

Dropout:-



Dropout ratio (P): $0 \leq P \leq 1$ \rightarrow You assign P value by doing hyperparameter tuning.

- During Training depending on the P value few neurons gets dropped out (inactive) only the active neurons does the operation after 1 process, again randomly other neurons gets activated and performs training.

- So that only few neurons participate in training.

So ~~the~~ the neurons are trained during training times.

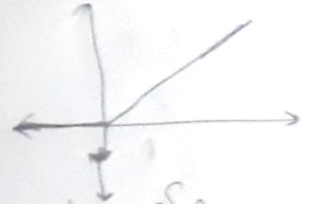
Test time: During the test time, all the neurons gets activated and connected.

Note for all the weights will be multiplied by the dropout ratio in each & every layer.

Weights = $W \times P$

↳ ReLU & Leaky ReLU

ReLU:

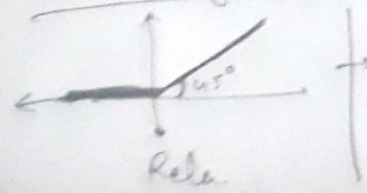


$$ReLU = \max(0, z)$$

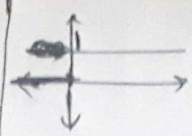
$$z = \sum_{i=1}^n W_i x_i + b_i$$

when $\begin{cases} z < 0 \rightarrow 0 \\ z > 0 \rightarrow z \end{cases}$

Derivative of ReLU:



Derivative \rightarrow slope
 $m = \tan \theta$
 $\theta = 45^\circ \rightarrow m = 1$
 $\theta = 0^\circ \rightarrow m = 0$



derivative $\begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$

Derivative of ReLU

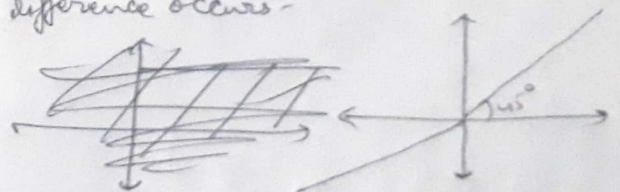
When derivative becomes 0 in $z < 0$ then

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

then $\frac{\partial L}{\partial W} = 0$ then $\eta \frac{\partial L}{\partial W} = 0$ therefore

dead neuron $(W_{new} = W_{old})$ so there is no change.

To overcome it we do Leaky ReLU i.e. when $z < 0$ we multiply z with a small quantity (α). so therefore a small difference occurs.



Therefore

$\begin{cases} z & z > 0 \\ 0.01(z) & z < 0 \end{cases}$

Weight Initialization

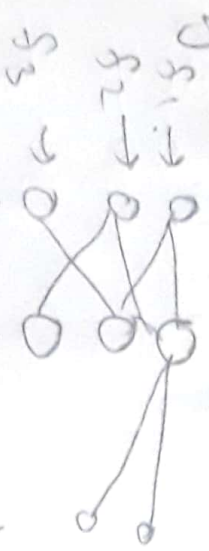
Key points

- ① Weights should be small
- ② Weights should not be ~~same~~ same
- ③ Weights should have good variance

Techniques used: ①: Uniform Distribution

$$w_{ij} \sim \text{Uniform} \left[-\frac{1}{\sqrt{\text{fan}_{in}}}, \frac{1}{\sqrt{\text{fan}_{in}}} \right]$$

The weights gets ~~evenly~~ assigned uniformly from the range $[a, b]$



$$\left[-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]$$

②: Xavier/Glorot Distribution

works well with sigmoid

Xavier Normal:

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

$$\sigma = \sqrt{\frac{2}{(\text{fan}_{in} + \text{fan}_{out})}}$$

③: He init + CRelu



the uniform

$$w_{ij} \sim \text{Uniform} \left[-\frac{6}{\sqrt{\text{fan}_{in}}}, \frac{6}{\sqrt{\text{fan}_{in}}} \right]$$

(a, b)

Xavier uniform:

$$w_{ij} \sim \text{Uniform} \left[-\frac{\sqrt{6}}{\sqrt{\text{fan}_{in}}}, \frac{\sqrt{6}}{\sqrt{\text{fan}_{in}}} \right]$$

the normal

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

$$\sigma = \sqrt{\frac{2}{\text{fan}_{in}}}$$

↳ Stochastic Gradient Descent

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

Now, if:

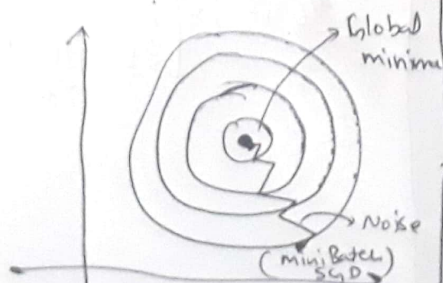
• $\frac{\partial L}{\partial W_{\text{old}}} \rightarrow$ When we consider all n points \rightarrow gradient descent

• $\frac{\partial L}{\partial W_{\text{old}}} \rightarrow$ 1 data point \rightarrow ~~SGD~~ SGD

• $\frac{\partial L}{\partial W_{\text{old}}} \rightarrow$ k data points \rightarrow mini Batch
 $k < n$ ~~SGD~~

$$\left[\frac{\partial L}{\partial W_{\text{old}}} \right] \approx \left[\frac{\partial L}{\partial W_{\text{old}}} \right]_{\text{GD}}$$

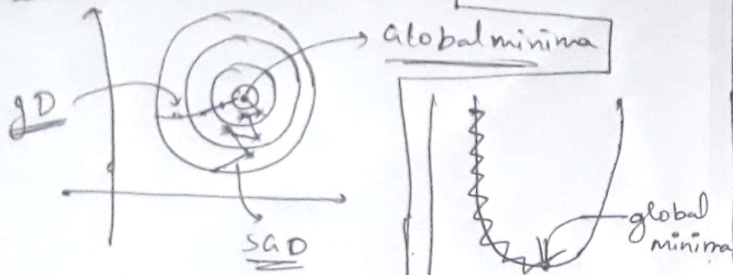
mini Batch SGD



There is a noise produced in ^{mini batch} Stochastic GD
 to remove it we use SGD with momentum

↳ Global Maxima & Local Maxima

SGD with momentum



- SGD from above it takes more time to converge that is due to the noise so we need to remove it.
- Exponential moving average is used to remove the noise

time	t_1	t_2	t_3	t_4	...	t_n
data point	b_1	b_2	b_3	b_4	...	b_n

$$V_{t1} = b_1 \quad \{0 \leq \gamma \leq 1\}$$

$$V_{t2} = (\gamma)(V_{t1}) + b_2$$

$$V_{t2} = (\gamma)(b_1) + b_2$$

$$V_{t3} = (\gamma)(V_{t2}) + b_3$$

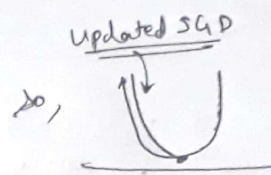
$$= (\gamma)(\gamma b_1 + b_2) + b_3$$

$$= \gamma^2 b_1 + \gamma b_2 + b_3$$

$$V_{t3} = \gamma^2 b_1 + \gamma b_2 + b_3$$

When $\gamma = 0.5$

$$V_{t3} = 0.25 b_1 + 0.5 b_2 + b_3$$



more importance is given to b_3

So formula becomes:-

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

$$W_{\text{old}} = \left[\gamma V_{t-1} + \eta \frac{\partial L}{\partial W_{\text{old}}} \right]$$

$$V_{t-1} = 1 \times \left[\frac{\partial L}{\partial W_{\text{old}}} \right]_t + \gamma \left[\frac{\partial L}{\partial W_{\text{old}}} \right]_{t-1} + \gamma^2 \left[\frac{\partial L}{\partial W_{\text{old}}} \right]_{t-2} + \dots$$

↳ Adagrad Optimizer:

In GD, SGD, Minibatch SGD we used to have the η (learning rate) fixed to each and every neuron.

• But in Adagrad optimizer, we have η varying learning rate at each ~~iteration~~ iterations.

$$w_t = w_{t-1} - \eta'_t \left(\frac{\partial L}{\partial w_{t-1}} \right)$$

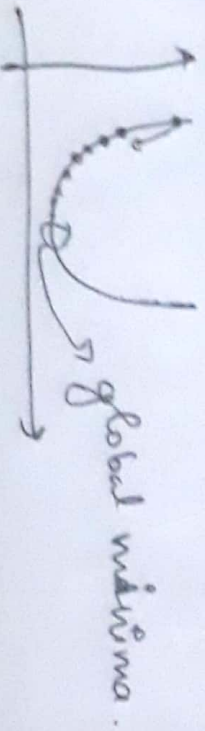
$t \rightarrow$ current, $t-1 \rightarrow$ previous.

$$\eta'_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \rightarrow \text{initial LR}$$

\rightarrow small the number

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

\rightarrow starts with high value then η'_t decreases
So w_t also ~~varies~~ varies.



Disadvantage with Adagrad:

~~at~~ The α_t parameter gets increased a lot after some time.

↳ Adadelta & RMSprop