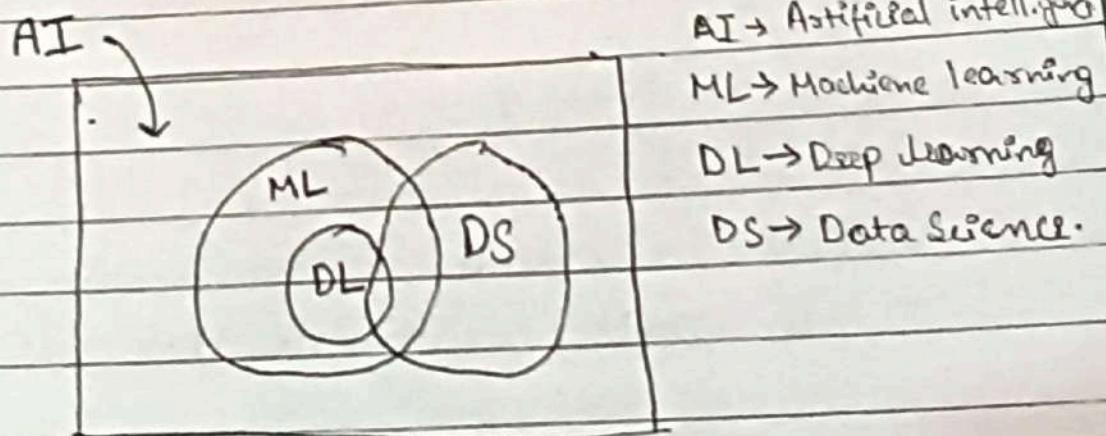


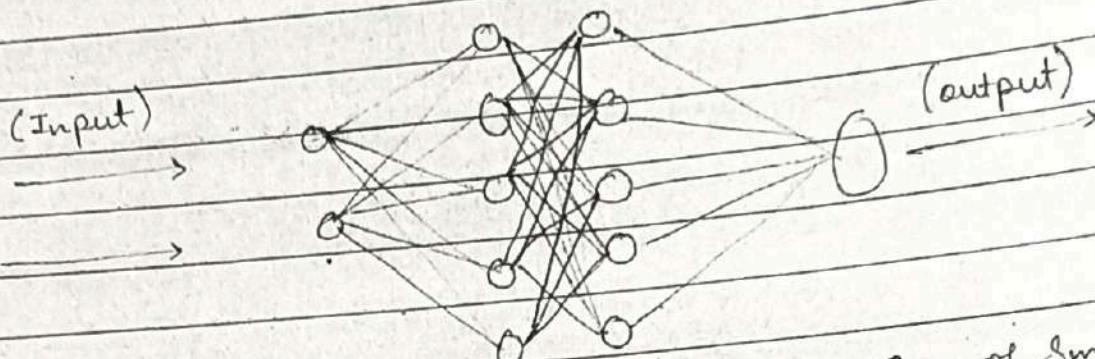
Introduction to Neural Network, Loss Function, Optimizers, Gradient Descent, SGD, Adagrad, RMSprop, Adam

Deep learning

- ANN → Artificial Neural Network (Tabular kind of dataset)
- CNN → Convolutional neural Network (Image, Videos)
- Image classification → CNN, Transfer Learning.
- Object detection → RCNN, FAST RCNN, FACTSR, RCNN, SSD
YOLO.
- Object Segmentation → Self driving car
- RNN → Recurrent Neural Network → Text, Time Saver, Data



A Simple Neuron
 Deep learning is a specialized subset of machine learning that uses layered neural networks to stimulate human decision-making.



A Neural Networks in AI is a collection of small computing units called neurons that take incoming data and learn to make decisions over time.

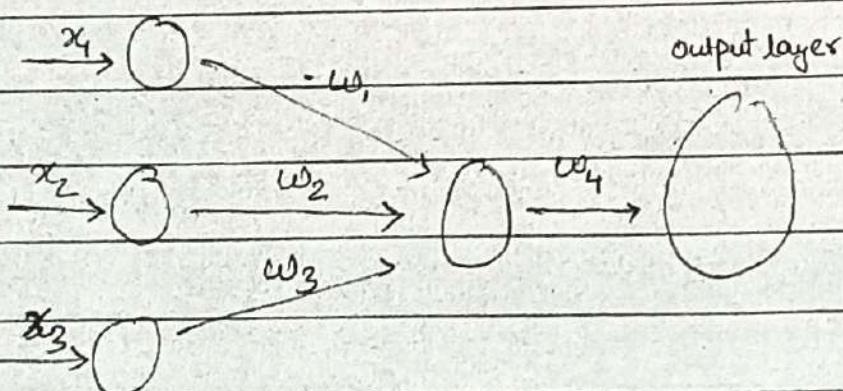
Perceptron

A perceptron is a neural network unit (an artificial neuron) that does certain computation to detect features or business intelligence in the input data.

$$\text{Step 1} \quad y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}$$

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

Input layer



Neuron

Step 1

Step 2

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

- The function output is not centered on 0, which will reduce the efficiency of weight update
- The Sigmoid function performs exponential operations, which is slower for computers.

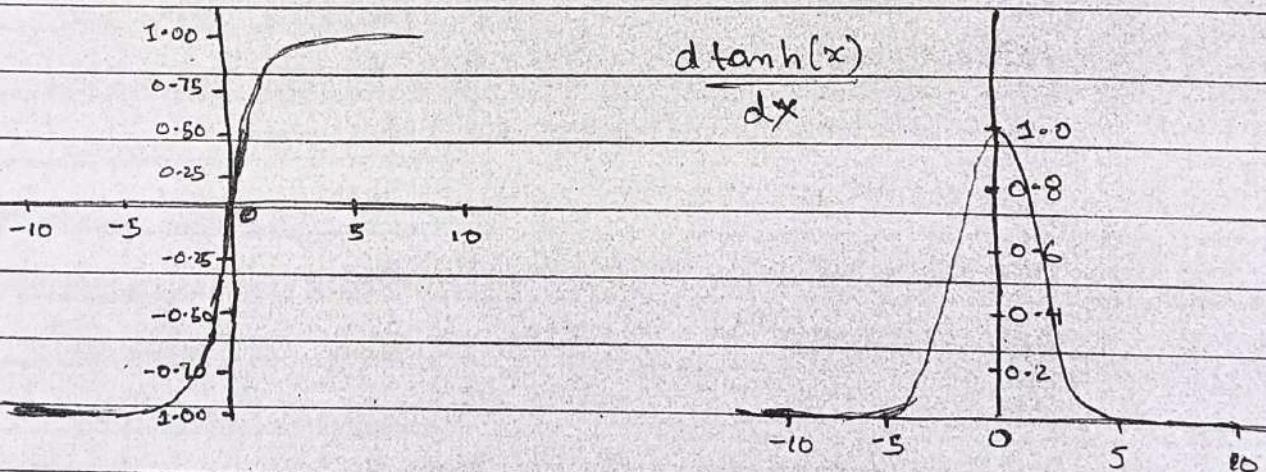
\Rightarrow Advantages of Sigmoid Function

- Smooth gradient, preventing "jumps" in output values.
- Output values bound b/w 0 and 1, normalizing the output of each neuron
- clear predictions, i.e very close to 1 or 0

\Rightarrow Sigmoid has three major disadvantages .

- Prone to gradient vanishing.
- Function output is not zero-centred
- Power operations are relatively time consuming .

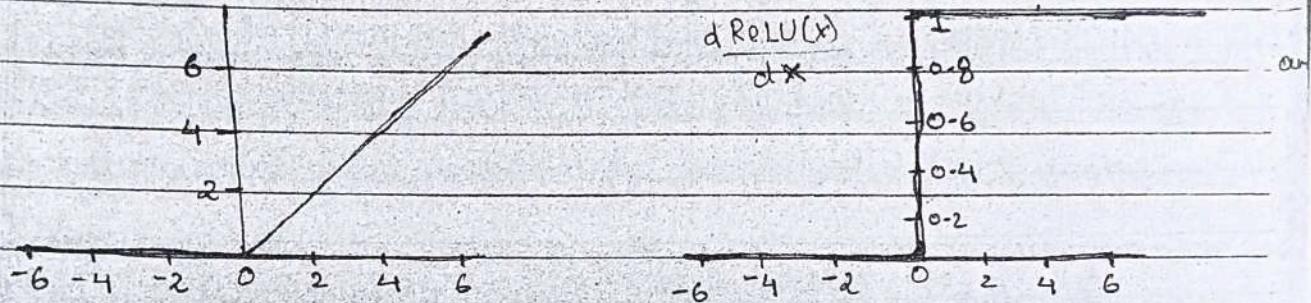
(2) Tanh function.



\rightarrow Tanh is a hyperbolic tangent function. The curves of tanh function and Sigmoid function are relatively similar.

- when the input is large or small, the output is almost smooth & the gradient is small, which is not conducive to weight update.
- The difference is the output interval. The output interval of tanh is 1
- The whole function is 0-centric, which is better than Sigmoid
- In general binary classification problems, the tanh function is used for the hidden layer and the Sigmoid function is used for the output layer. However, these are not static, and the specific activation function to be used must be analyzed according to the specific problem, or it depends on debugging.

③ ReLU function $[\max(0, x)]$

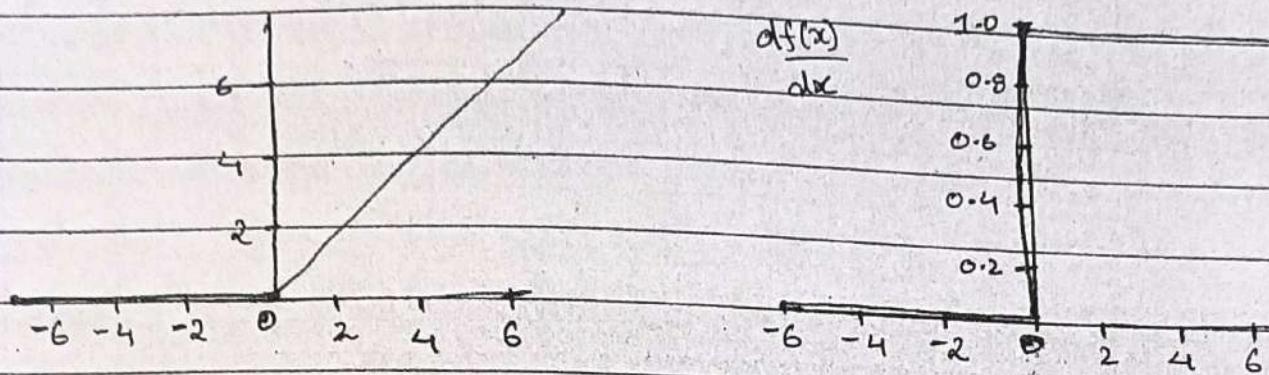


- whenever your $x > 0$, then the output will be that particular x . If $x < 0$, then the output will be 0.
- This is mostly used in hidden layers because it solves the vanishing gradient problem.
- Compared to the Sigmoid function and the tanh function, it has the following advantages:
 - when the input is positive, there is no gradient saturation problem.
 - The calculation speed is much faster.

→ Disadvantages.

- When the input is -ve, ReLU is completely inactive.
- The output of the ReLU function is either 0 or a positive number, which means that the ReLU function is not a 0-centric function.

④ Leaky ReLU function $f(x) = \max(0.01x, x)$

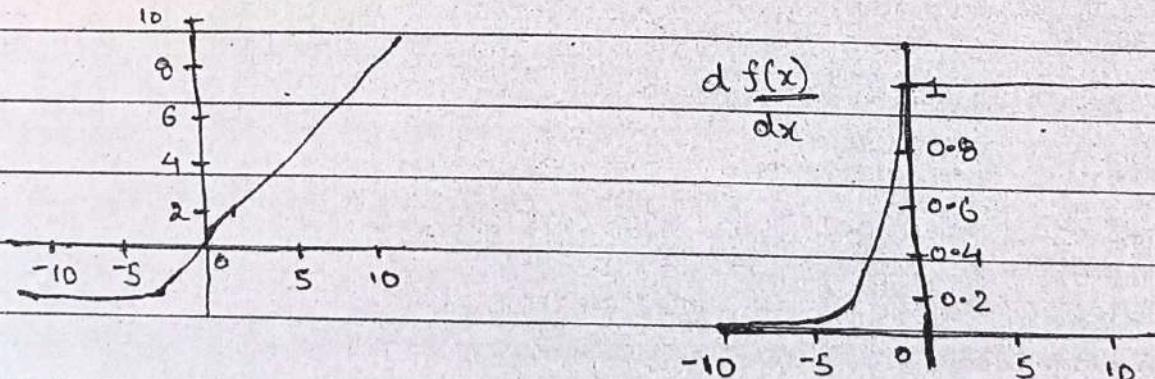


→ ~~ReLU~~ Leaky ReLU has all the advantages of ReLU, plus there will be no problems with Dead ReLU.

→ But in actual operation, it has not been fully proved that Leaky ReLU is always better than ReLU.

→ People proposed to set the first half of ReLU 0.01x instead of 0. Another intuitive idea is a parameter-based method, ~~ReLU~~ parametric ReLU : $f(x) = \max(\alpha x, x)$, where α can be learned from back propagation.

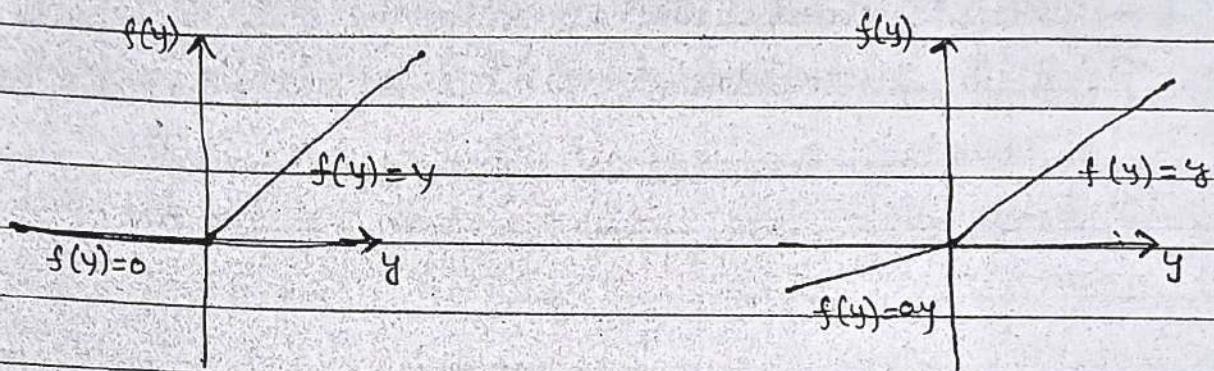
⑤ ELU (Exponential Linear Units) function



- No Dead ReLU issues
- The mean of the output is close to 0, zero-centered.

→ One small problem is that it is slightly more computationally intensive. Similar to Leaky ReLU, although theoretically better than ReLU, there is currently no good evidence in practice that ELU is always better than ReLU.

⑥ PReLU (Parametric ReLU)



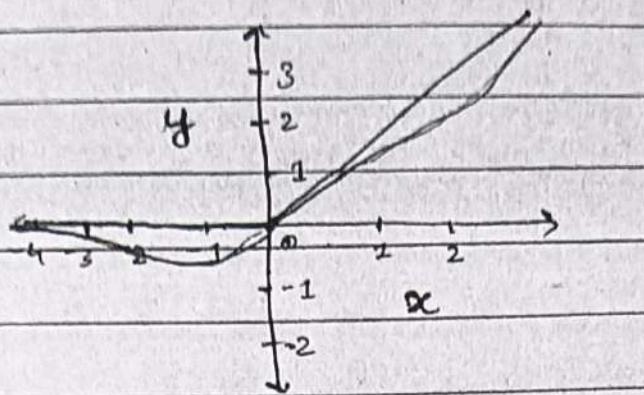
- PReLU is also an improved version of ReLU.
- In the negative region, PReLU has a small slope, which can also avoid the problem of ReLU death.

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

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

- if $(a_i = 0)$, f becomes ReLU
- if $(a_i > 0)$, f becomes leaky ReLU
- if (a_i) is a learnable parameter, f becomes PReLU.

⑦ Swish (A Self-Gated) Function



ReLU
Swish

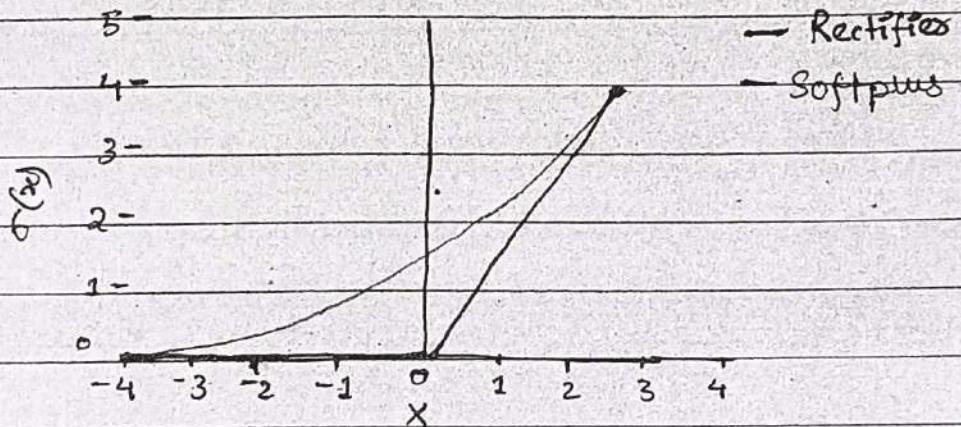
→ $y = x * \text{Sigmoid}(x)$ ← formula, and this is called Self-gated

→ This is used in LSTMs

→ we use this type of activation function when we have our neural network > 10 layers

→ It is zero centric and it also solves the dead-relu

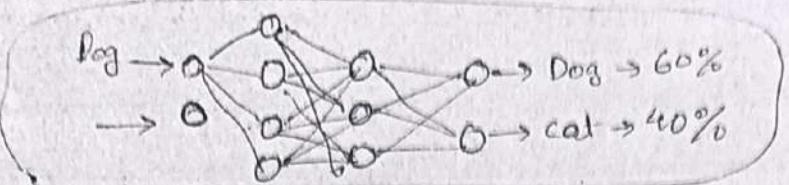
⑧ Softplus



— Rectifier
— Softplus

→ we know that derivative of zero is not possible in real world Scenario, most of the time, but we may get sometimes when we solving, so if you get zero in that specific way, using a ReLU will not serve the purpose because we don't know the derivative of 0. So we can use Softplus instead.

$$\rightarrow f(x) = \ln(1 + e^x)$$



⑦ Softmax → we use this when we have more than 2 outputs.

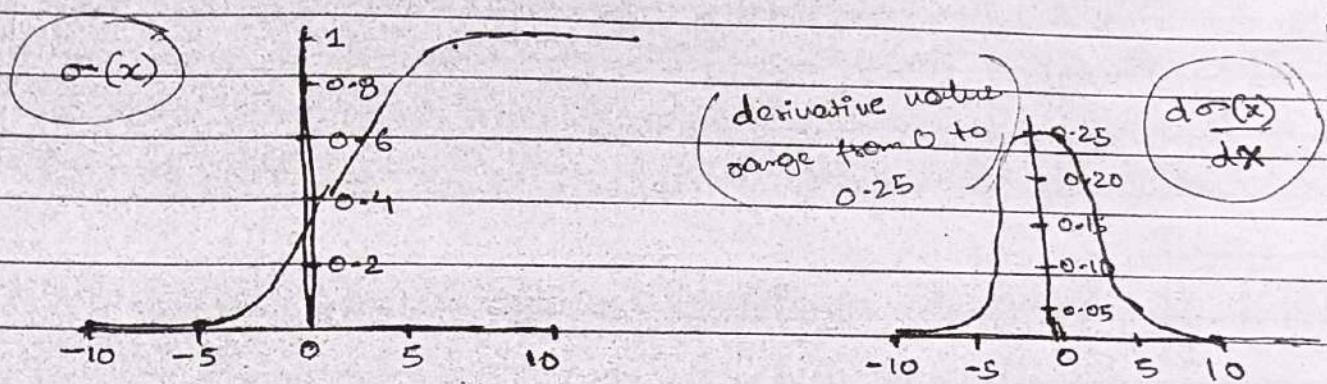
$$\rightarrow S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, j = 1, 2, \dots, K$$

Note Generally speaking, these activation functions have their own advantages and disadvantages. There is no statement that indicates which ones are not working, and which activation functions are good - All the good and bad must be obtained by experiments.

Name	Plot	Equation
Identity		$f(x) = x$
Binary Step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Logistic (a.k.a Sigmoid / Soft Step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$
Tanh		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$
Rectified Linear unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Leaky (ReLU)		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Softmax	for $i = 1, \dots, J$	$f_i(\vec{x}) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$
Maxout		$f(\vec{x}) = \max_i x_i$

* Commonly used activation functions

① Sigmoid function

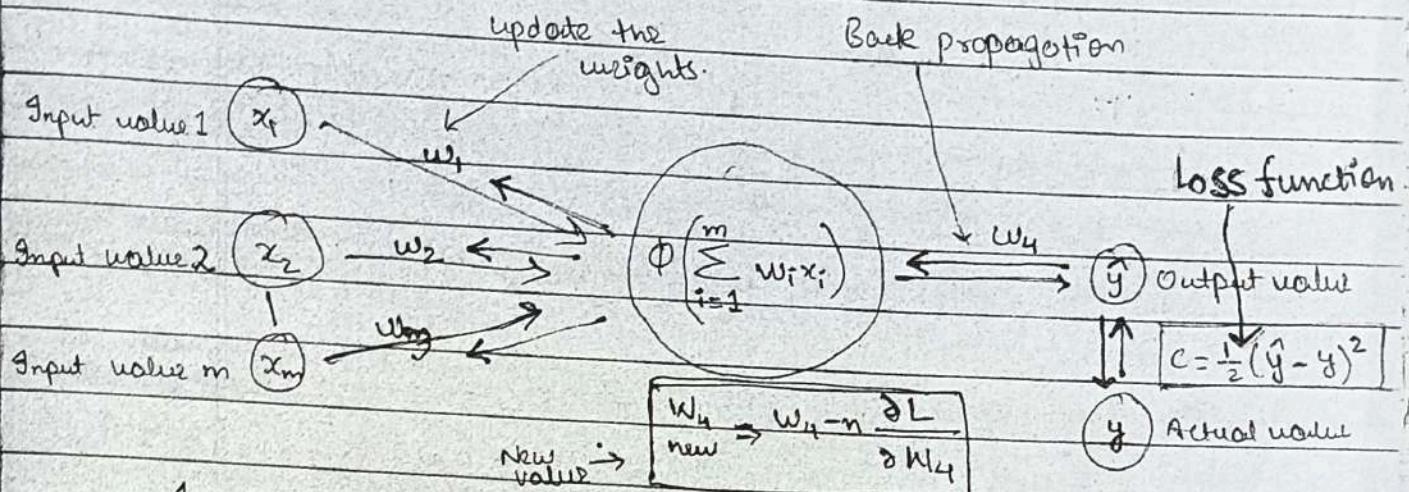


⇒ The Sigmoid function is the most frequently used activation function in the beginning of deep learning. It is a smoothing function that is easy to derive.

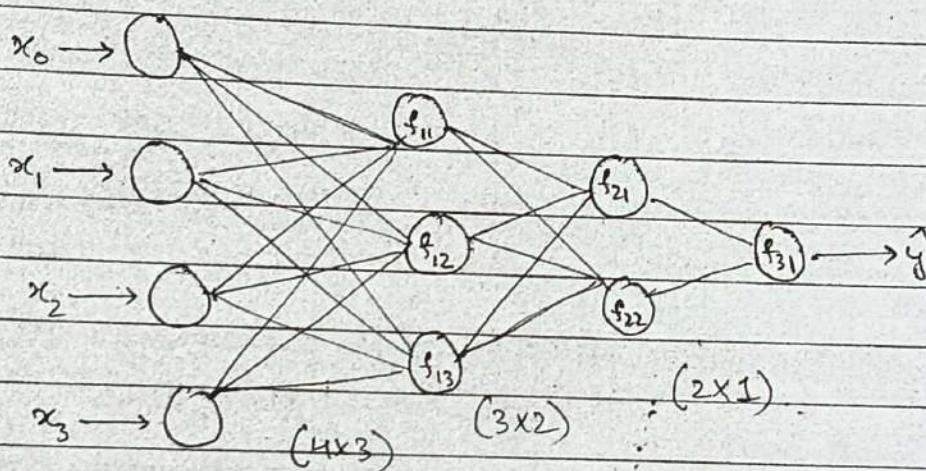
- ⇒ The function itself have certain defects:
- when the input is slightly away from the co-ordinate origin, the gradient of the function becomes very small, almost zero. In the process of neural network backpropagation, we all use the chain rule of differential to calculate the differential of each weight (w). when the backpropagation passes through the Sigmoid function, the differential on this chain is very small. Moreover, it may pass through many Sigmoid functions, which will eventually cause the weight (w) to have little effect on the loss function. which is not conducive to the optimization of the weight. This problem is called gradient Saturation or gradient dispersion.

In case of Back propagation, our \hat{y} and y values are different. So we take Back propagation. here we update the w_4 value first the w_1, w_2 & w_3 . Once we get the new value, again the front propagation will go ahead, this process will keep on happening until the Loss value is completely reduced i.e. $\hat{y} = y$.

★ Back propagation



★ Training Multilayer Neural Network by Gradient Descent



- for comparing \hat{y} and y we use loss function ($L = (y - \hat{y})^2$)
- here we are trying to see the difference b/w y and \hat{y} and our main aim is to reduce that loss value.
- then only we'll be getting similar results for y & \hat{y}
- If we want to reduce the loss function, we basically use an optimizer such as gradient descent.
- Gradient descent work is to update ~~all~~ all the weights in such a way they should be able to get a \hat{y} similar to y .
- This updating the weights to each layer is known as

$$\text{for -ve slope} \rightarrow W_{\text{new}} = W_{\text{old}} + \eta \frac{\partial L}{\partial w}$$

formula for updating the weights

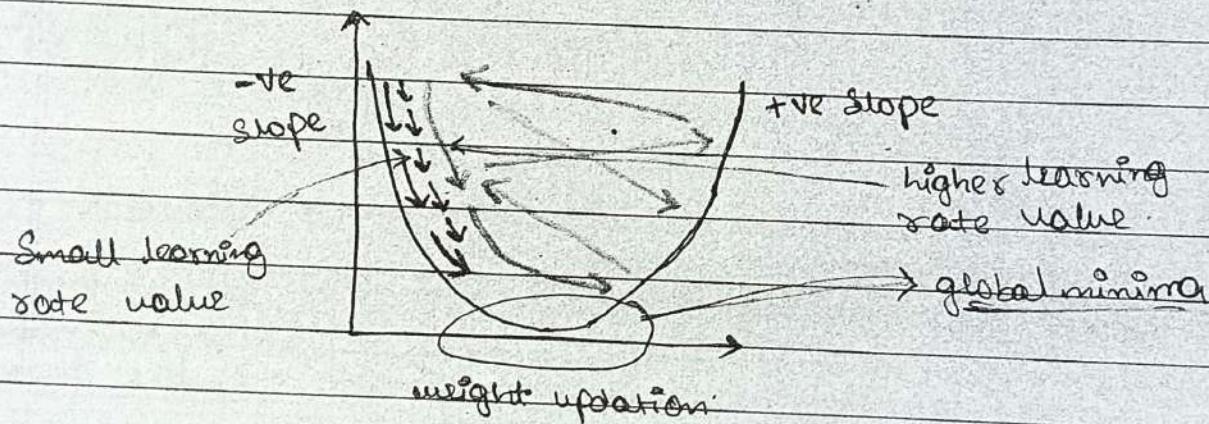
$$\text{for +ve slope} \rightarrow W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

$\rightarrow W_{\text{new}} = \text{new weight}$
 $\rightarrow W_{\text{old}} = \text{old weight}$
 $\rightarrow \eta = \text{learning rate}$

Q) why do we need to find derivative? $\rightarrow \frac{\partial L}{\partial w} \rightarrow \text{derivative of the loss}$

Ans) derivative is nothing but finding the slope $\rightarrow \frac{\partial w}{\partial w} \rightarrow \text{derivative of that specific weight}$

- * The main aim of Gradient Descent is to reach in "global minima". Unless and until we don't reach to * global minima, the updation of the weights will happen with the help of back propagation.
- * If the loss value is decreasing, that will also indicate that we are reaching towards global minima.



- If we take a small learning rate value, then it will take more time to reach global minima point.
- If we take a high learning rate value, then it will take less time to reach, but most of the time it will start jumping ~~here & there~~ b/w -ve slope and +ve slope. and it may never reach the global minima point.

→ So that's why we should select a "learning value"?
 for such a way that it should be a value that
 is suitable with respect to gradient descent.

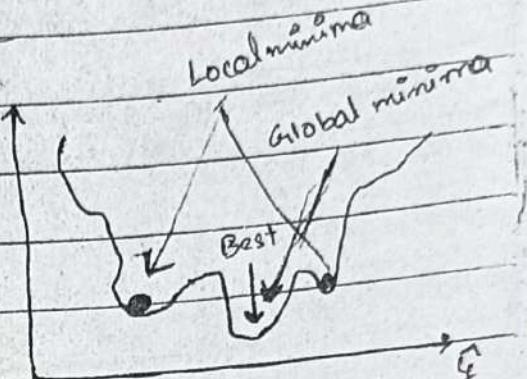
★ Local Minima and Global Minima

→ The point at which a function takes the minimum value is called global minima. However, when the goal is to minimize the function and solved using optimization algorithms such as gradient descent, it may so happen that function may appear to have a minimum value at different points.

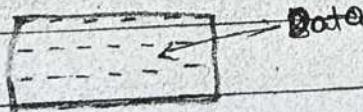
→ Those several points which appear to be minima but are not the point where the function actually takes the minimum value are called local minima.

→ Machine learning algorithms such as gradient descent algorithms may get stuck in local minima during the training of the models.

→ In the order to avoid this we use Stochastic gradient descent.



Batch - GD



→ May end up in local minima

→ Slow

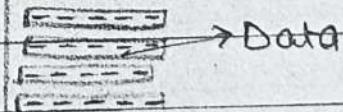
→ more deterministic

minibatch - GD



→ If stuck in local minima, some noisy steps can lead the user out of them.

Stochastic - GD

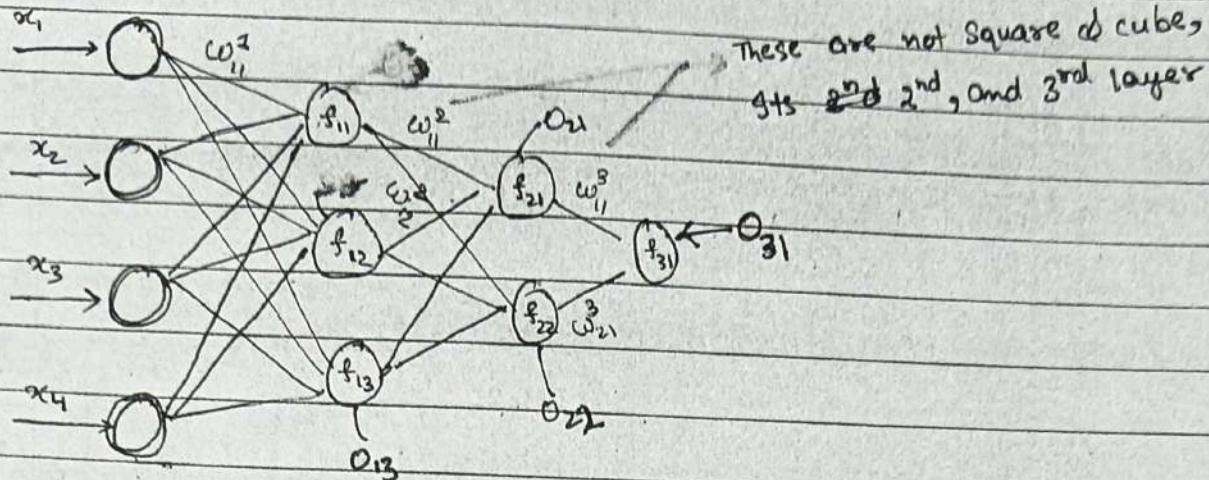


→ Local minima problem can be resolved

→ Fast

→ Less deterministic

★ Chain Rule in Back Propagation



* If we want to update w_{11}^3 (1st layer from back side)

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

$$\rightarrow W_{11}^3 \text{ new} = W_{11}^3 \text{ old} - \eta \frac{\partial L}{\partial w_{11}^3}$$

we can also write as

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}$$

where $\frac{\partial O_{31}}{\partial w_{11}^3}$ is the derivative of the output of the previous layer

* If we want to update w_{21}^3

$$W_{21}^3 \text{ new} \rightarrow W_{21}^3 \text{ old} - \eta \frac{\partial L}{\partial w_{21}^3}$$

can be written as

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{21}^3}$$

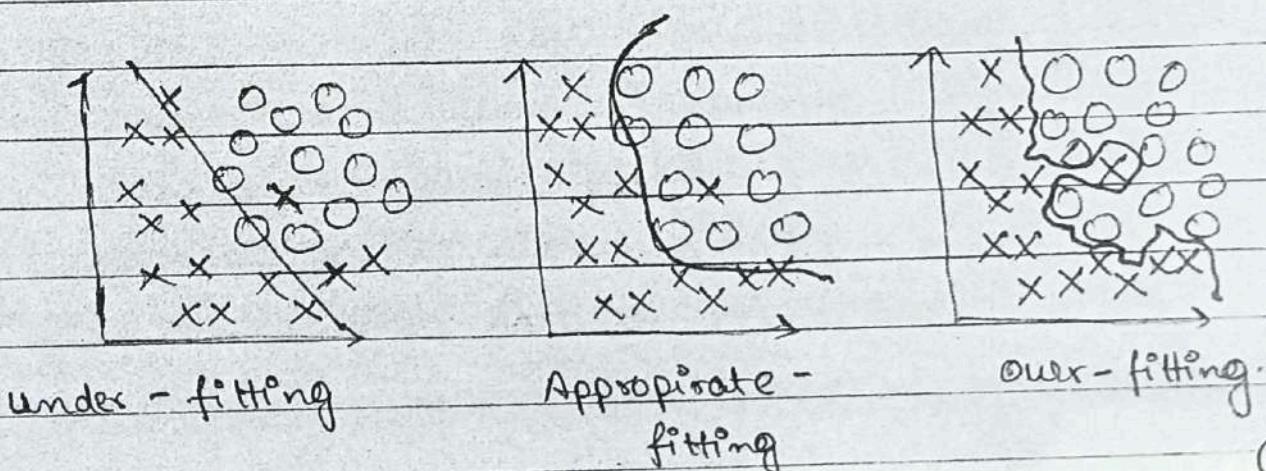
* If we want to find the derivative of w_{11}^2

$$\rightarrow w_{11}^2 \text{ new} \rightarrow w_{11}^2 \text{ old} - \eta \left(\frac{\partial L}{\partial w_{11}^2} \right)$$

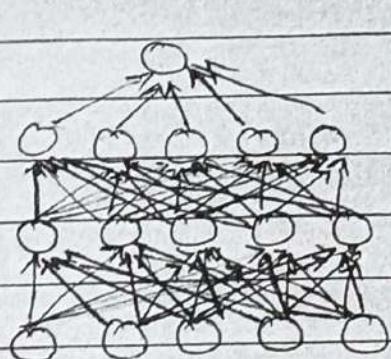
$$\rightarrow \frac{\partial L}{\partial w_{11}^2} = \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial w_{11}^2} \right] + \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial w_{12}^2} \right]$$

* Vanishing Gradient Problem

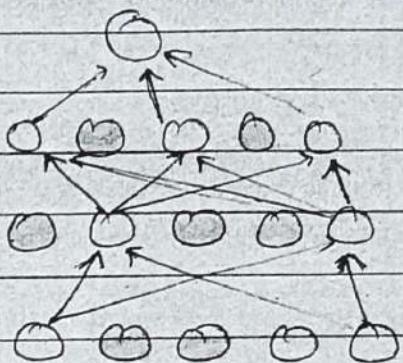
- When the elements of the gradient become exponentially small so that the update of the parameters with the gradient becomes almost insignificant.
- As the number of layers increases, the derivative value decreases and at one situation you will see that the value of 'old weight' and 'new weight' are matching. (Predicted output (y) and given output (g) are different. the grad be same)
- So this is the problem we face in Sigmoid & tanh
- This is the reason why Sigmoid is bad and it shouldn't be used in each and every layer.
- This is where ~~ReLU~~ ReLU comes into existence.
- Sigmoid function ranges b/w 0 to 1 & Derivative of Sigmoid function ranges b/w 0 to 0.25
- ↑ Dropout Regularization



→ 'Dropout' in machine learning refers to the process of randomly ignoring certain nodes in a layer during training.



a) Standard Neural Network



b) After applying dropout

→ Dropout is used as a regularization technique — it prevents overfitting by ensuring that no units are codependent

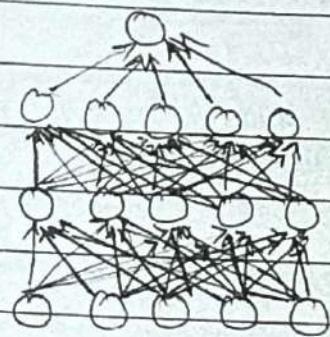
Q) why will dropout help with overfitting?

- It can't rely on one input as it might be dropped out at random
- Neurons will not learn redundant details of inputs.

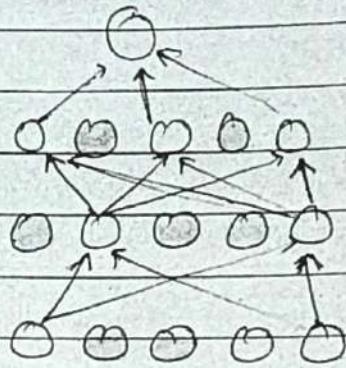
★ Dead neuron

- A neuron that cannot be activated, i.e., contributes nothing to feature extraction or classification. In practise, when you see a feature map is almost 0 no matter what kind of input you feed, there is very likely a dead neuron.
- The problem with ReLU is that Sigmoid function always converts it into the range 0 to 1 which creates dead neuron.

- Activation function → Activation functions are mathematical equations that determine the output of a neural network.
- 'Dropout' in machine learning refers to the process of randomly ignoring certain nodes in a layer during training.



a) Standard Neural Network



b) After applying dropout

- Dropout is used as a regularization technique — it prevents overfitting by ensuring that no units are codependent

Q) why will dropout help with overfitting?

- It can't rely on one input as it might be dropped out at random
- Neurons will not learn redundant details of inputs.

★ Dead neuron

- A neuron that cannot be activated, i.e., contributes nothing to feature extraction or classification. In practice, when you see a feature map is almost 0 no matter what kind of input you feed, there is very likely a dead neuron.
- The problem with ReLU is that Sigmoid function always converts it into the range 0 to 1 which creates dead neuron.

→ That's why we use Leaky ReLU instead. In order to solve the Dead ReLU Problem, people proposed to set the final half of ReLU (0-0.01) instead of (0).

GD, SGD, Minibatch SGD \rightarrow These all optimizers were helping us to reduce the loss function with respect to the weights

★ weight initialization

- \rightarrow weights should be small
- \rightarrow weights should not be same
- \rightarrow weights should have good variance

→ There are those weight initialization techniques.
(we don't have to remember this, just understand)

① uniform Distribution

$$\rightarrow w_{ij} \sim \text{Uniform } [a, b] \quad \begin{matrix} -1 & , & 1 \\ \sqrt{\text{fan_in}} & , & \sqrt{\text{fan_in}} \end{matrix} \quad \begin{matrix} \text{fan_in} & \text{fan_out} \end{matrix}$$

Hidden Layers

② Xavier / Glorot \rightarrow This works good with Sigmoid AD

\rightarrow Xavier Normal

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

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

\rightarrow Xavier Uniform

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

④ He init

$$w_{ij} \sim U\left[-\frac{\sqrt{6}}{\sqrt{\text{fan_in}}}, \frac{\sqrt{6}}{\sqrt{\text{fan_in}}}\right] \quad w_{ij} \sim N(0, \sigma)$$

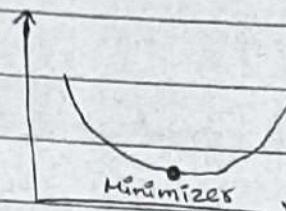
$$\sigma = \sqrt{\frac{2}{\text{fan_in}}}$$

\Rightarrow He uniform

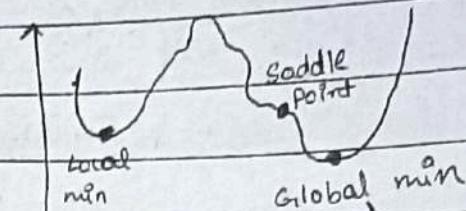
\Rightarrow He normal

If all the points present
in the same curve

\star Difference b/w convex & Non-convex function



Convex function

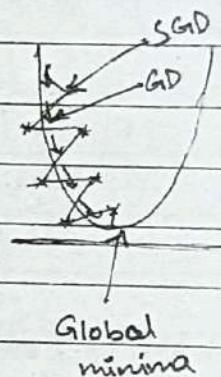
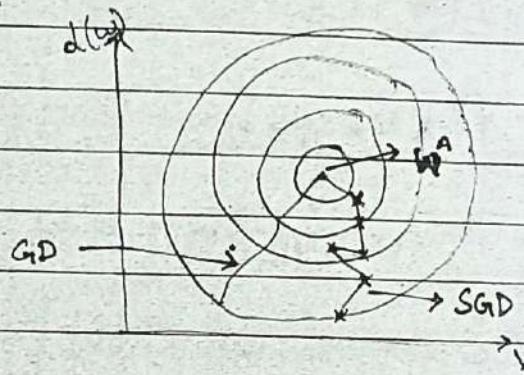


Non-convex function

\rightarrow we usually see this is ML Problems \rightarrow we find this in DL problems

SGD with momentum

This shows how your weights are getting updated to best
minima weight



\Rightarrow we can see that we ~~have~~ always have a
noisy data in case of SGD, thus for the convergence
it will take more time. So everytime we use SGD,
we should try to remove that noise from the data. In
the order to remove that noise, there is an concept
called 'Exponential moving average'

$$(b_1, b_2, \dots, b_n)$$

\Rightarrow Suppose ~~we~~ have data points at different time intervals
(t_1, t_2, \dots, t_n)

$t_1 \rightarrow b_1$
$t_2 \rightarrow b_2$
\vdots
\vdots

So, the Exponential moving average formula

$$\Rightarrow V = \gamma \times V_i + b_j$$

V = variable

γ = gamma

Need not have to learn,
just remember

t_1	t_2	\dots	t_n
b_1	b_2	\dots	b_n

finding the Exponential moving average $\Rightarrow \gamma = 0.5$

$$\rightarrow V_{t_1} = b_1 \quad (1)$$

$$\rightarrow V_{t_2} = \gamma \times V_1 + b_2 \quad (2)$$

$$= 0.5 \times b_1 + b_2 \quad [\text{from 1}]$$

$$\rightarrow V_{t_3} = \gamma + V_2 + b_3$$

$$= \gamma (\gamma \times V_1 + b_2) + b_3 \quad [\text{from 2}]$$

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

$$= 0.25 b_1 + 0.5 b_2 + 1 \times b_3$$

This is how moving average is calculated by
 $b_1, b_2 \text{ & } b_3$

\Rightarrow Now how do we apply the same formula in our G-D

$$\rightarrow w_{\text{new}} = w_{\text{old}} - n \times \frac{\delta L}{\delta w_{\text{old}}} \quad \begin{array}{l} \text{we write this instead} \\ \text{in this case} \end{array}$$

This is momentum

$$\rightarrow w_{\text{new}} = w_{\text{old}} - \left[\gamma V_{t-1} + n \frac{\delta L}{\delta w_{\text{old}}} \right] \quad (\because \gamma = \text{gamma})$$

$$\rightarrow \therefore V_{t-1} = 1 \times \left[\frac{\delta L}{\delta w_{\text{old}}} \right]_t + \gamma \left[\frac{\delta L}{\delta w_{\text{old}}} \right]_{t-1} + \gamma^2 \left[\frac{\delta L}{\delta w_{\text{old}}} \right]_{t-2}$$

\therefore this is how we should be able to reduce the noise

Implementation details

Initially

$$V_{d\omega} = 0 \quad V_{db} = 0$$

← weight

← bias

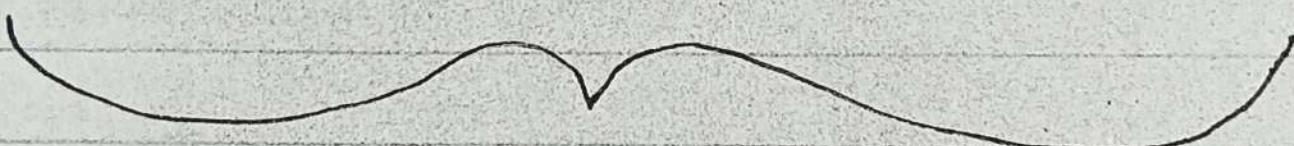
On iteration t inside epoch, compute
 $d\omega, db$ on the current mini batch.

$$\rightarrow V_{d\omega t} = \beta V_{d\omega t-1} + (1-\beta) \frac{\partial L}{\partial \omega_{t-1}}$$

$$\rightarrow V_{dbt} = \beta V_{dbt-1} + (1-\beta) \frac{\partial h}{\partial b_{t-1}}$$

$$\Rightarrow \boxed{\begin{aligned} \omega &= \omega - \eta V_{d\omega} \\ b &= b - \eta V_{db} \end{aligned}}$$

← by this we
will achieve
smoothening of
convergence.



GRADIENT DESCENT WITH MOMENTUM

Implementation details

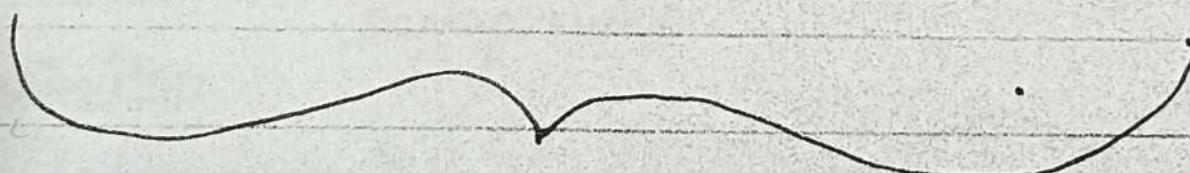
on iteration t in minibatch

compute $\frac{\partial h}{\partial \omega}$, $\frac{\partial L}{\partial \omega}$ on current minibatch

$$\rightarrow S_{d\omega} = \beta S_{d\omega} + (1-\beta) \left(\frac{\partial L}{\partial \omega} \right)^2$$

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

$$\boxed{w_t = w_{t-1} - \eta' \frac{\partial h}{\partial w_{t-1}}}$$
$$\boxed{b_t = b_{t-1} - \eta' \frac{\partial h}{\partial b_{t-1}}}$$



Algorithm for weight update

★ AdaGrad optimizer → Adaptive Gradient Descent

→ It is an algorithm for gradient-based optimization
 → In AdaGrad optimizer, there is no momentum concept
 So, it is much simpler compared to SGD with momentum.

→ The idea behind this is to use different learning rates for each parameter based on iteration.

→ The reason behind the need for different learning rates is that the learning rate for sparse features parameters needs to be higher compared to dense features parameters because the frequency of occurrence of sparse features is lower.

$\begin{bmatrix} \text{Sparse} \rightarrow (0, 1) \\ (\text{most of the values are } 0) \end{bmatrix}, \begin{bmatrix} \text{dense} \rightarrow \text{most of the values are non-zero} \end{bmatrix}$

So the main concept behind AdaGrad optimizer is that, 'can we use different learning rates with respect to number of feature, different weights and different iteration.'

$$\Rightarrow w_t = w_{t-1} - \eta \times \frac{\partial h}{\partial w_{t-1}}$$

$$\Rightarrow w_t = w_{t-1} - \left(\eta \frac{1}{t} \right) \times \frac{\partial h}{\partial w_{t-1}}$$

Initial learning rate

$$\Rightarrow \eta_t^1 = \frac{\eta}{\sqrt{\alpha_1 + \epsilon}}$$

$\because \epsilon \rightarrow \text{Small +ve number}$

$$\therefore \alpha_t \uparrow \Leftrightarrow \eta_t^1 \downarrow$$

$$\Rightarrow \alpha_t = \sum_i \left(\frac{\partial h}{\partial w_i} \right)^2$$

★ Advantages of using AdaGrad

- It eliminates the need to manually tune the learning rate
- Convergence is faster and more reliable than SGD when the scaling of the weights is unequal
- It is not very sensitive to the size of the master step

★ Disadvantages of using AdaGrad

- Its optimization method can result in aggressive monotonically decreasing learning rates.
- It is incredibly slow, this is because the sum of gradient squared only grows and never shrinks. (when we square the α_t value).
- $\alpha_t \uparrow \longleftrightarrow \gamma_t \downarrow$
- AdaDelta, RMSProp (for Root Mean Square Propagation) fixes this issue by adding a decay factor.

2) Adadelta or RMSProp

$$\text{Adagrad} \rightarrow w_{(t)}^{\text{new}} = w_{(t-1)}^{\text{old}} - \gamma_t \frac{\delta h}{\delta w_{\text{old}}}$$

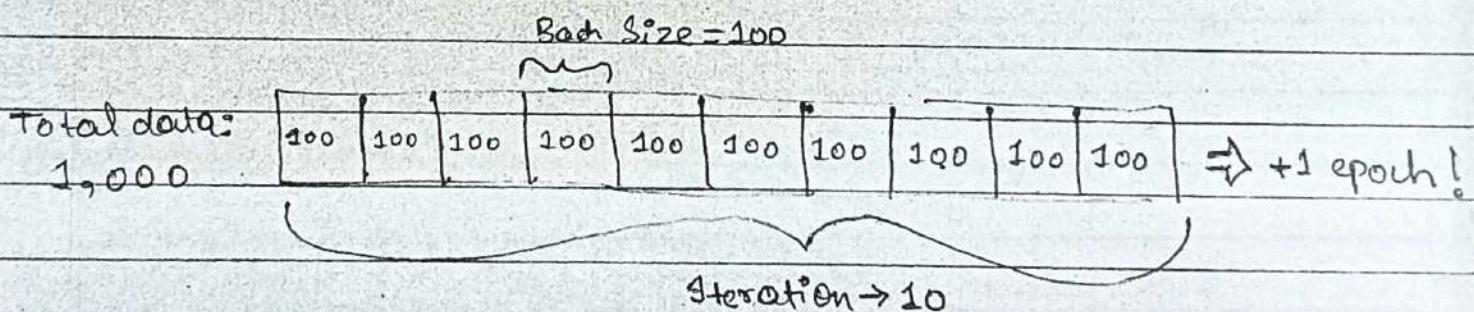
$$\gamma_t = \frac{\eta}{\sqrt{w_{\text{avg}} + \epsilon}}$$

$$w_{\text{avg}} = \gamma \times w_{\text{avg}}^{(t-1)} + (1-\gamma) \times \left(\frac{\delta L}{\delta w_t} \right)^2$$

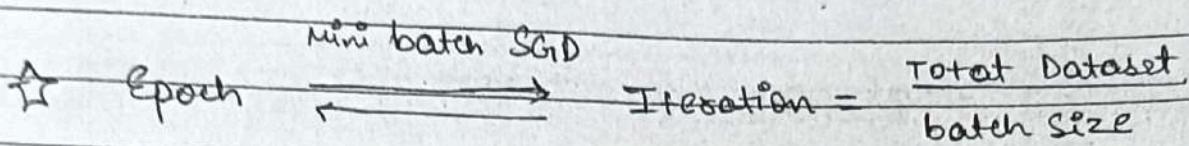
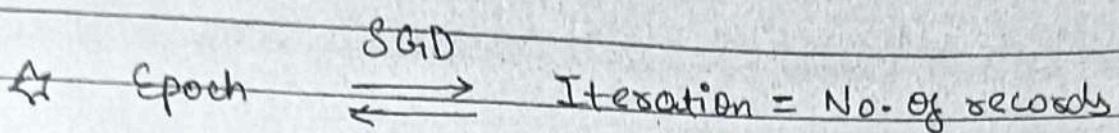
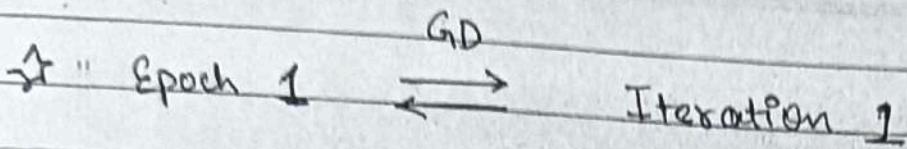
- It restricts α_t to become a very very big number.
- It should become small, but not very very small.
- Since it prevents α_t to become very very big numbers, so your learning rate will not be very very small and convergence will not take that much time.

★ Difference b/w Epoch and Iteration

- An epoch is defined as the number of times an algorithm visits the data set. In other words, epoch is one backward and one forward pass for all the training.
- Iteration is defined as the number of times a batch of data has passed through the algorithm. In other words, it is the number of passes, one pass consists of one forward and one backward pass.
- For example :-
- * consider a set of 1000 images having a batch size of 10 one each. iteration would process 10 images, it will take 100 such iteration to complete 1 epoch.



- Total Dataset/No. of records \rightarrow 10K
- Batch size \rightarrow 1K



★ Adam optimizer \rightarrow Adaptive Moment Estimation

- \Rightarrow Adam combines the best properties of the AdaGrad and RMSProp algorithms to provide an optimization algorithm that can handle sparse gradients on noisy problems.
- \Rightarrow with the help of momentum, you will get Smoothing, and with the help of RMS PROP, you will be able to change your learning rate in an efficient manner. Such that α_t value will also not go high.

ALGORITHM IMPLEMENTATION

$$V_{dw} = 0 \quad V_{db} = 0 \quad S_{dw} = 0 \quad S_{db} = 0$$

on iteration t

\Rightarrow compute $\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$ using current mini batch

$$\left. \begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1 - \beta_1) \frac{\partial L}{\partial w} \\ V_{db} &= \beta_1 V_{db} + (1 - \beta_1) \frac{\partial L}{\partial b} \end{aligned} \right\} \begin{array}{l} \text{Momentum} \\ \text{Smoothing} \\ \text{factors} \end{array}$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta) \left(\frac{\partial L}{\partial w} \right)^2$$

RMS problem

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

$$w_b = w_{t-1} - \frac{\eta \times V_{dw}}{\sqrt{S_{dw} + \epsilon}}$$

→ initial learning rate

Adam optimized

$$b_t = b_{t-1} - \frac{\eta \times V_{db}}{\sqrt{S_{db} + \epsilon}}$$

→ Smoothening factor

→ RMS properties

★ Bias Correction

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

$$V_{db}^{\text{correction}} = \frac{V_{db}}{(1-\beta_1^t)}$$

Momentum → Smoothening factor

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

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

RMS problem

$$\Rightarrow w_b = w_{t-1} - \frac{\eta \times V_{dw}^{\text{correction}}}{\sqrt{S_{dw}^{\text{correction}} + \epsilon}}$$

$$b_t = b_{t-1} - \frac{\eta \times V_{db}^{\text{correction}}}{\sqrt{S_{dw}^{\text{correction}} + \epsilon}}$$

* Multi class - Cross Entropy Loss

$$L(x_i, y_i) = - \sum_{j=1}^c y_{ij} * \log(\hat{y}_{ij})$$

↳ categorical Cross Entropy

→ we use this when we have more than 2 outputs.

★ Training ANN

Step 1 → Randomly initialise the weights to small numbers close to 0 (but not 0).

Step 2 → Input the first observation of your dataset in the input layer, each feature in one input node.

Step 3 → Forward-propagation : from left to right, the neurons are activated in a way that the impact of each neuron's activation is limited by the weights. Propagate the activation until getting the predicted result y .

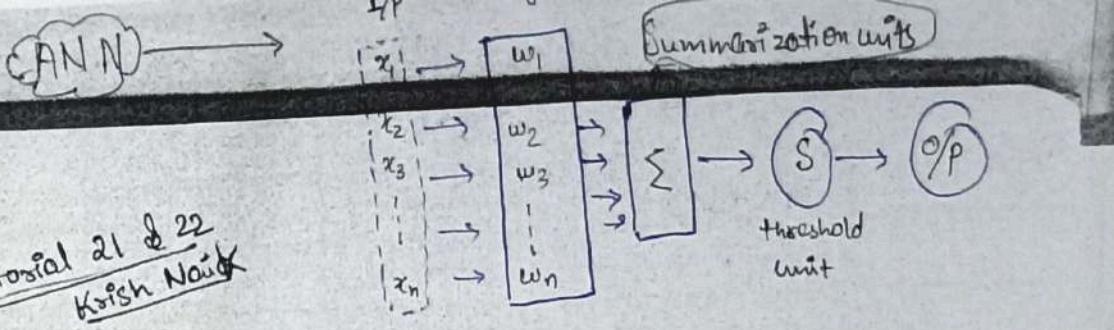
Step 4 → Compare the predicted result to the actual result. Measure the generated error.

Step 5 → Back-propagation : from right to left, the error is back-propagated. Update the weights according to how much they are responsible for the error. The learning rate decides by how much we update the weights.

Step 6 → Repeat steps 1 to 5 and update the weights after each observation (Reinforcement learning) or

Repeat steps 1 to 5 but update the weights only after a batch of observations (Batch learning)

Step 7 → when the whole training set passed through the ANN, that makes an epoch. Redo more epochs.



★ Convolution

→ Convolution is a mathematical operation on two functions (f and g) that produces a third function (h) that expresses how the shape of one is modified by the other. The term convolution refers to both the resulting function and to the process of computing it.

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

*

+1	0	-1
+2	0	-2
+2	0	-1

=

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

(4x4)

↓
Suppose → white block

↓
=

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

↓
Transforming this image using min-max Scores

★ Padding in CNN

- Padding basically extends the area of an image in which a convolutional neural network processes.
- In order to work the Kernel/filter with processing in the image, padding is added to the outer frame of the image to allow for more space for the filter to cover in the image.
- Adding padding to an image processed by an CNN allows for a more accurate analysis of image.

before padding

$n-f+1$

$n=6-3+1$

If we want to get my output dimension as 6
--

$n-f+1=6$

$n=6+3+1$

after padding

$P=1, n+2P-f+1$

$6+2\times 1-3+1$

use its image recognition and tasks that involve the processing of pixel data.

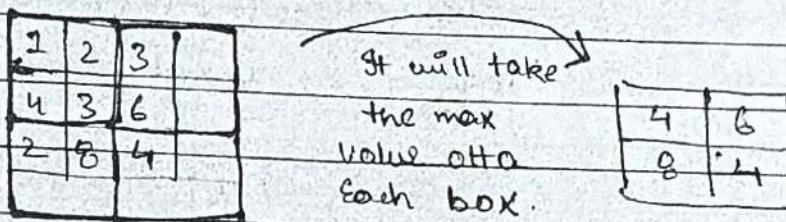
ANN

CNN → Convolutional neural networks

- uses fully connected layers → uses Partially connected layers
Size of connection depends on size of the filter
- can be used only for small images → can be used for any image
- considers entire image for learning the patterns → considers only the part of the image, which is in Receptive Field of the Filter..
- Number of parameters will be very high → Number of Parameters will be very less.
- Doesn't learn the patterns in Spatial Hierarchy i.e., higher layers of ANN don't combine the lower layers outputs → Learns Spatial Hierarchy of Patterns i.e., higher layers of CNN are formed by combining lower layers. This helps in identifying the patterns more effectively than ANN
- Success rate is high → Success rate is more than ANN
- Scalability is lower → Scalability is higher
- neurons of fully connected → neurons are partially connected
- has 2 basic layers → input & output → It has input, output, convolution and Pooling layers
- It has 1D to 3D → It has 3D Neuron Volume
- Not translation invariant: Once a regular DNN has learned to recognize a pattern in one location, it can recognize it only in that location. So we have to train the image again (can't reuse the learning (weights) can be reused previous weights) if it is Rotated or Shifted.
- Translation invariant: Once the CNN has learned to recognize a pattern in one location, it can recognize it in any location. So even if the image is Rotated or Shifted.

★ Max Pooling

- Max Pooling is a Pooling operation that calculates the maximum value for patches of a feature map and uses it to create a downsampled (pooled) feature map. It is usually used after a convolutional layer.
- It is used to reduce the dimensions of the feature maps. Thus it reduces the Number of parameters to learn and the amount of computation performed in the network.
- This makes the model more robust to variations in the position of the features in the input image.



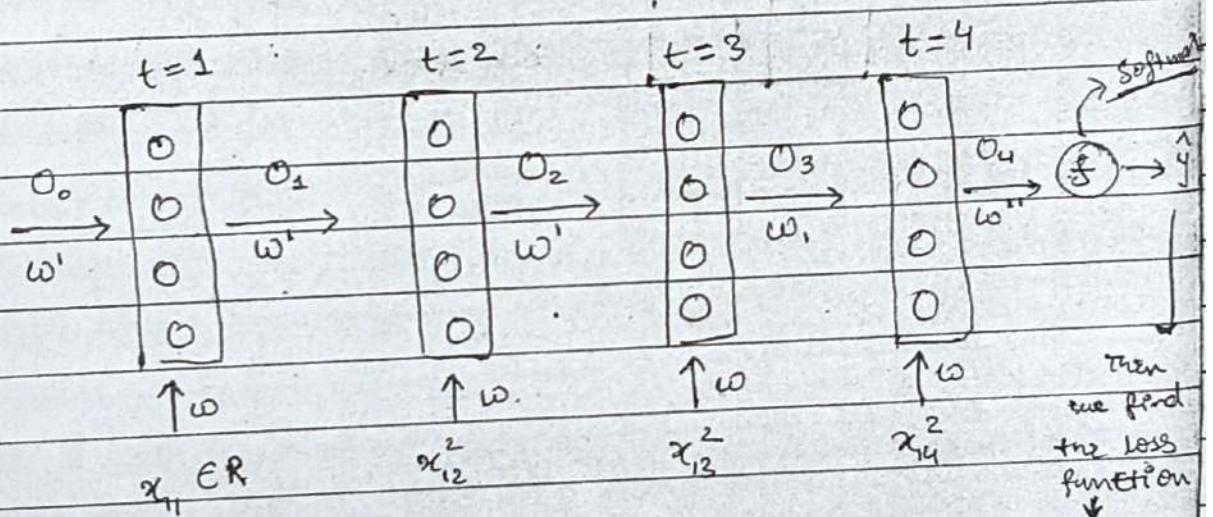
★ Data Augmentation

- Data augmentation in data analysis are techniques used to increase the amount of data by adding slightly modified copies of already existing data or newly created synthetic data from existing data. It acts as a regularizer and helps reduce overfitting when training a machine learning model.
- It basically flip, horizontal Shifts, vertical Shifts or zoom the data in the order to give more variance.
- Because even if we provide flipped image to CNN, it should be able to predict what image it is.

use of convolution layers which basically use filters to help recognise important feature in an image. There are so many features which can help distinguish a particular image. All these features are being learned internally. When we use a deep CNN, weights can be millions.

RNN Deep Learning

- Recurrent neural networks (RNN) are the state of the art algorithm for sequential data and are used by Apple's Siri and Google's voice search. It is one of the algorithms behind the scenes of the amazing achievements seen in deep learning over the past few years.
- Recurrent Neural Networks (RNN) are a type of Neural Network where the output from the previous step is fed as input to the current step. RNN's are mainly used for, Sequence Classification - Sentiment classification of video classification. Sequence Labelling - Part of, Speech tagging & named entity recognition.



$$\rightarrow o_1 = f(x_{11}w + o_0 w')$$

$$\rightarrow o_2 = f(x_{12}w + o_1 w')$$

$$\rightarrow o_3 = f(x_{13}w + o_2 w')$$

$$\rightarrow o_4 = f(x_{14}w + o_3 w')$$

$$\text{Loss} \rightarrow (\hat{y}_1 - y_1)$$

Our main aim is to reduce the loss function

- ⇒ forward propagation has happened from left to right. just with respect to one statement. Each and every words can be represented into vectors of some (d) dimensions.

Transfer Learning

whenever we train CNN on a bunch of Images, convolution network make use of convolution layers which basically use filters to help recognise important feature in an image. There are so many feature which can help distinguish a particular image. All these features are being learned internally. When we use a deep CNN, weights can be million when million of parameters need to be learn, it's gonna take a lot of time. So transfer learning is a research problem in machine learning (ML) that focuses on storing knowledge gained while solving one problem and applying it to a different but related problem.

for eg → knowledge gained while learning to recognize cars could be applied when trying to recognize trucks.

→ weight updation is important so that we can reach global minima in the particular gradient descent. When million of parameters need to be learned, it's gonna take a lot of time. So transfer learning is a technique.

and that depends on the type of vocabulary, No. words that are present in that vocabulary.

- The output (4) depends on x_{14} and output (3)
- Output (3) depends on x_{13} and output (2)
- Output (2) depends on x_{12} and output (1)

⇒ So by this way, the sequence information is always maintained, it can't be discarded. Which was the problem in bow, TF-IDF.

→ weight only gets assigned during backward propagation, not forward.

★ Backward propagation in RNN

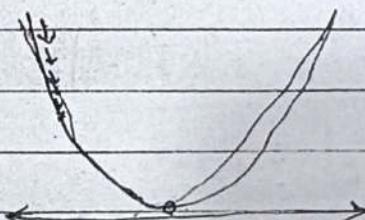
→ When the back propagation is getting applied, we just finding the derivative and updating the weights from backward.

→ As soon as we find derivative, which is actually a slope, we will be able to find out the weights, we will update the weights and this will go on for many iterations.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_4} \cdot \frac{\partial o_4}{\partial w} \quad \text{for } 1^{\text{st}} \text{ we from back.}$$

★ Problems in Simple RNN

RNN suffer from the problem of vanishing gradients. The gradients carry information used in the RNN parameter update and when the gradient becomes smaller and smaller, the parameter updates become insignificant which means no real learning is done.



NLP

Complete Road Map to prepare NLP

→ PyTorch

→ Keras

→ TensorFlow

(Best

Transformers

→ Spacy

→ Natural Language
Analyses with NLTK

(Bidirectional LSTM RNN, Encoders/Decoders, Attention Models

(Text Preprocessing Level 3 - Word Embeddings, Word2vec

Understanding Recurrent Neural Networks, LSTM, GRU

Get the Understanding of Artificial Neural Network

Solve Machine Learning Usecases

Text preprocessing - Gensim, Word2Vec, Avg. word2vec

Text Pre-processing Level-2, Bag of words, TFIDF, Unigram, Bi, n-Gram

Text preprocessing Level I - Tokenization, Lemmatization,
Stopword, stemming, Part of Speech (POS) Tagging

Natural language processing helps computers communicate with humans in their own language and scales other language-related tasks. For example, NLP makes it possible for computers to read text, hear speech, interpret it, measure sentiment and determine which parts are important.

⇒ Tokenization

- The first thing you need to do in any NLP project is text preprocessing.
- Tokenization can separate sentences, words, characters, or subwords. When we split the text into sentences, we call it sentence tokenization. For words, we call it words tokenization.

⇒ Stemming and Lemmatization

- Stemming and Lemmatization both generate the foundation sort of the inflected words and therefore the only difference is that stem may not be an actual word whereas, lemma is an actual language word.
- Lemmatization takes more time.

word	Poterstemmer Stemming	wordNetLemmatizer Lemmatization
information	inform	information
informative	inform	informative
computers	comput	computer
feet	feet	foot

- Problem with respect to Stemming is that it produces intermediate representation of the word may not have any meaning.
- Stop words → words that actually does not put much value in the paragraph.

★ Bag of words

A bag of words is a representation of text that describes the occurrence of words within a document. It involves two things: A vocabulary of known words. A measure of the presence of known words.

Eg →

Sentence 1 = He is a good boy

Sentence 2 = She is a good girl

Sentence 3 = Boy & girl are good

11 Vectors

(BOW)

	good	boy	girl	O/P
Sent 1	1	1	0	-
Sent 2	1	0	1	-
Sent 3	1	1	1	-

→ Disadvantages

- It gives the same preference to all the words.
- If the new sentences contain new words, then our vocabulary size would increase and thereby, the length of the vectors would increase too. Additionally, the vectors would also contain many 0s, thereby resulting in a sparse matrix (which is what we would like to avoid).
- It is messy and techniques like machine learning algorithms prefer well defined fixed-length inputs & outputs.

TF - IDF

→ TF (Term frequency) $\Rightarrow \frac{\text{No. of repetition of words in Sentence}}{\text{No. of words in Sentence}}$

IDF (Inverse Document Frequency) $\Rightarrow \log \left(\frac{\text{No. of Sentences}}{\text{No. of Sentences containing words}} \right)$

So we have to multiply TF & IDF to quantify words in a set of documents.

- Eg → Sentence 1 → good boy
- Sentence 2 → good girl
- Sentence 3 → boy girl good

↓

TF		
	good	boy
good	1/2	1/2
boy	0	1/2
girl	0	1/3

X
(3x3)

IDF	
words	IDF
good	$\log \left(\frac{3}{3} \right) = 0$
boy	$\log \left(\frac{3}{2} \right)$
girl	$\log \left(\frac{3}{2} \right)$

	f_1 good	f_2 boy	f_3 girl
Sent. 1	0	$\frac{1}{2} \times \log \left(\frac{3}{2} \right)$	0
Sent. 2	0	0	$\frac{1}{2} \times \log \left(\frac{3}{2} \right)$
Sent. 3	0	$\frac{1}{3} \times \log \left(\frac{3}{2} \right)$	$\frac{1}{3} \times \log \left(\frac{3}{2} \right)$

⇒ when we were using BOW, we just had values like 0 and 1. So, TF-IDF gives semantic meaning to the data.

★ word 2 Vec

Bog of words and TF-IDF Problems

→ Both BOW and TF-IDF approach semantic information is not stored. TF-IDF gives importance to uncommon words.

→ There is definitely chance of over fitting.

Solution — Word 2 Vec

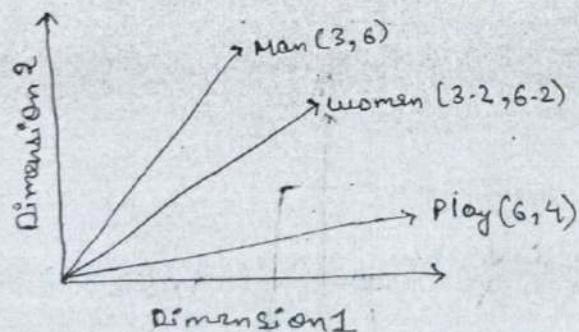
→ In this specific model, each word is basically represented as a vector of 32 or more dimension instead of a single number.

→ Here the semantic information and relation b/w different words is also preserved.

→ Another great advantage of word 2 Vec approach is that the size of the embedding vector is very small.

Stopwords → words that actually does not put much value in the paragraph.

Visual Representation - word 2 Vec



$$\Rightarrow \text{King} - \text{Man} + \text{woman} = \text{Queen}$$

Steps to create word2vec

- Tokenization of the sentences
- Create Histograms
- Take most frequent words
- Create a matrix with all the unique words. It also represent the occurrence relation between the words

∴ we use PorterStemmer() for Stemming purpose

Gated Recurrent Unit (GRU)

LSTM

GRU

→ 3 Gates: Input, output, forget

→ 2 Gates: reset, update

→ More accurate on longer sequence less efficient

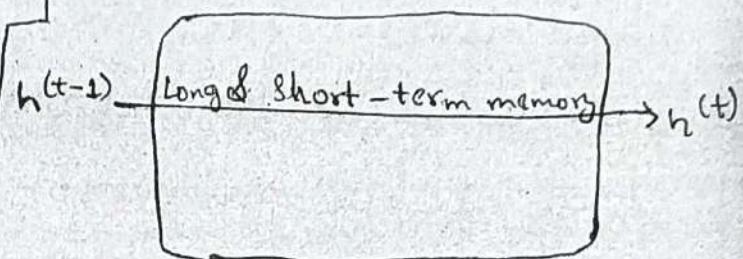
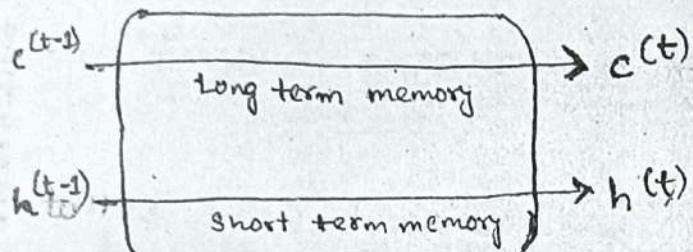
→ More efficient computation wise. Getting more popular.

→ Invented: 1995 - 1997

→ Invented: 2014

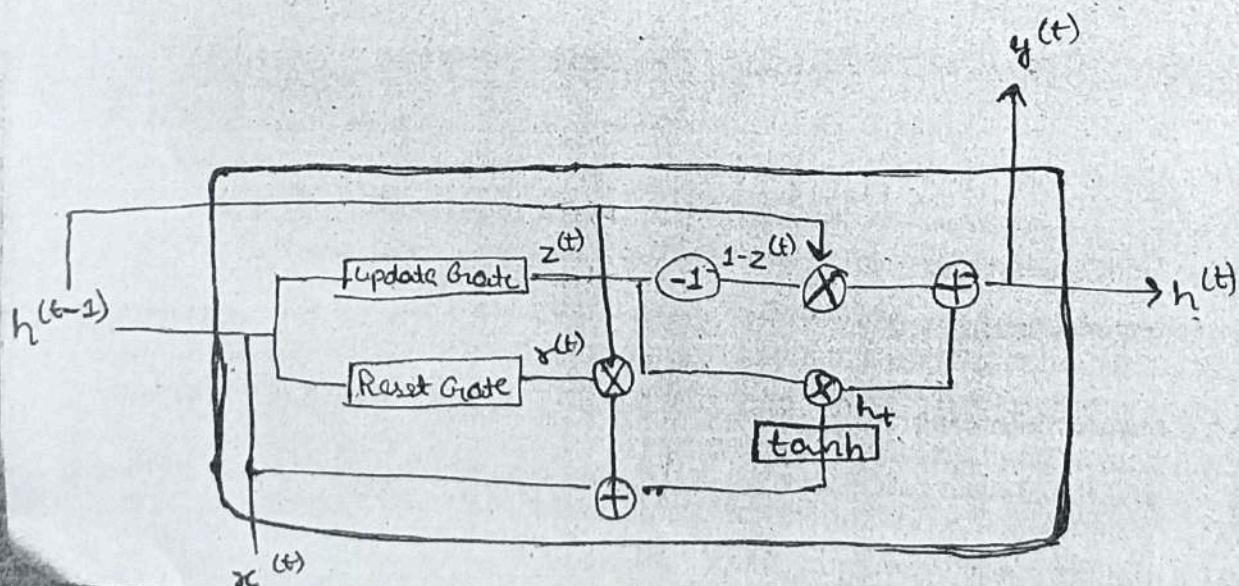
LSTM

GRU



→ GRU is a lightweight version of LSTM where it combines long and short term memory into its hidden state.

→ LSTM has two cell state and hidden state here it has only hidden state which can combine both long & short term memory



$$\rightarrow z_t = \sigma_g (w_z x_t + U_z h_{t-1} + b_z)$$

$$\rightarrow r_t = \sigma_g (w_r x_t + U_r h_{t-1} + b_r)$$

$$\rightarrow \hat{h}_t = \phi_h (w_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

$$\rightarrow h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t$$

Variables :-

• x_t → input vector

• z_t → update gate vector

• h_t → output vector

• r_t → reset gate vector

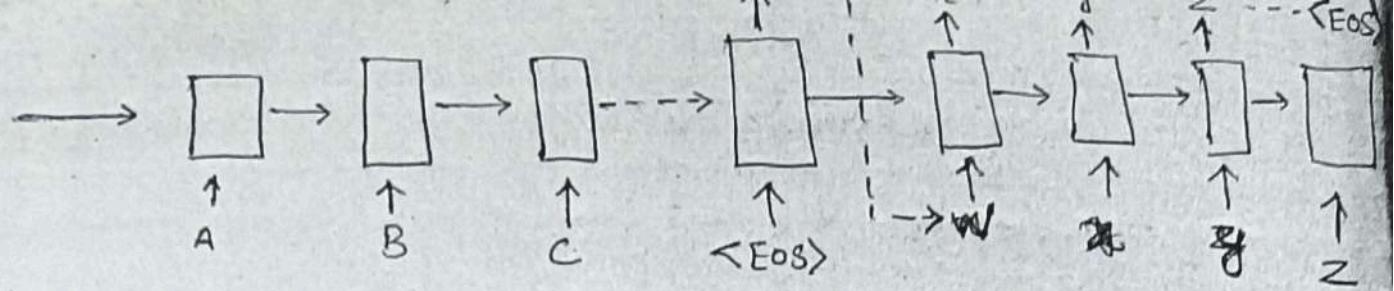
• \hat{h}_t → candidate activation
vector

• w, U and b : parameter
matrices & vectors.

The crux of this whole GRU unit is, it combines short-term memory and long term memory into one state and its also little bit lightweight than LSTM.

Seq2Seq Learning with Neural Networks

- In Seq2Seq, my model takes a sequence of information, and it output some sequence of information.
- Suppose we want to convert English to French, or converting image to text. (Image captioning).
- for Encoders, we don't have the output at every timestamp.
- When we get an 'context vector' out of 'Encoder', this vector is then passed to the 'Decoder'. For decoders, we do have the output at every timestamp.
- The input sequence and the output sequence need not



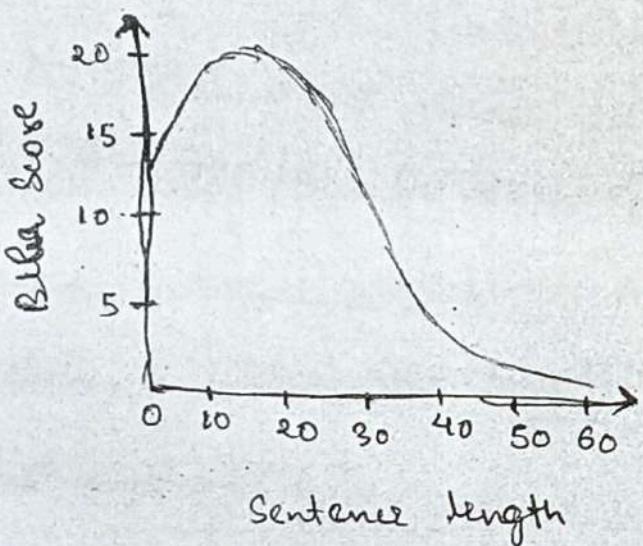
(Encoders)

(Decoders)

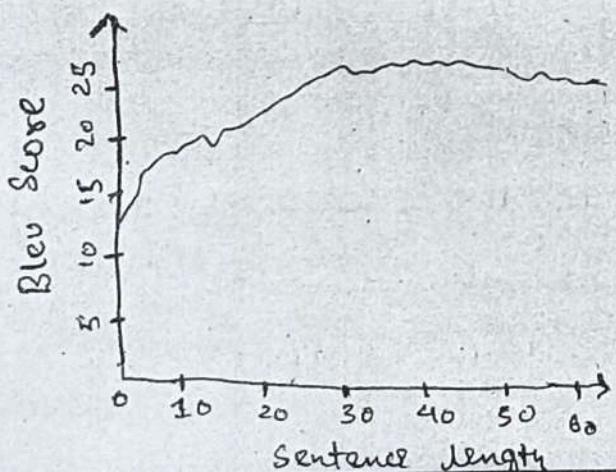
→ EOS = End of stream

→ One disadvantage is that, whenever you have long sentence, the accuracy is not that good. → disadvantage of Encoders & Decoders

- BLEU → To measure the machine translation effectiveness, we will evaluate the closeness of the machine translation to human reference translation using a metric known as BLEU - Bilingual Evaluation Understudy.

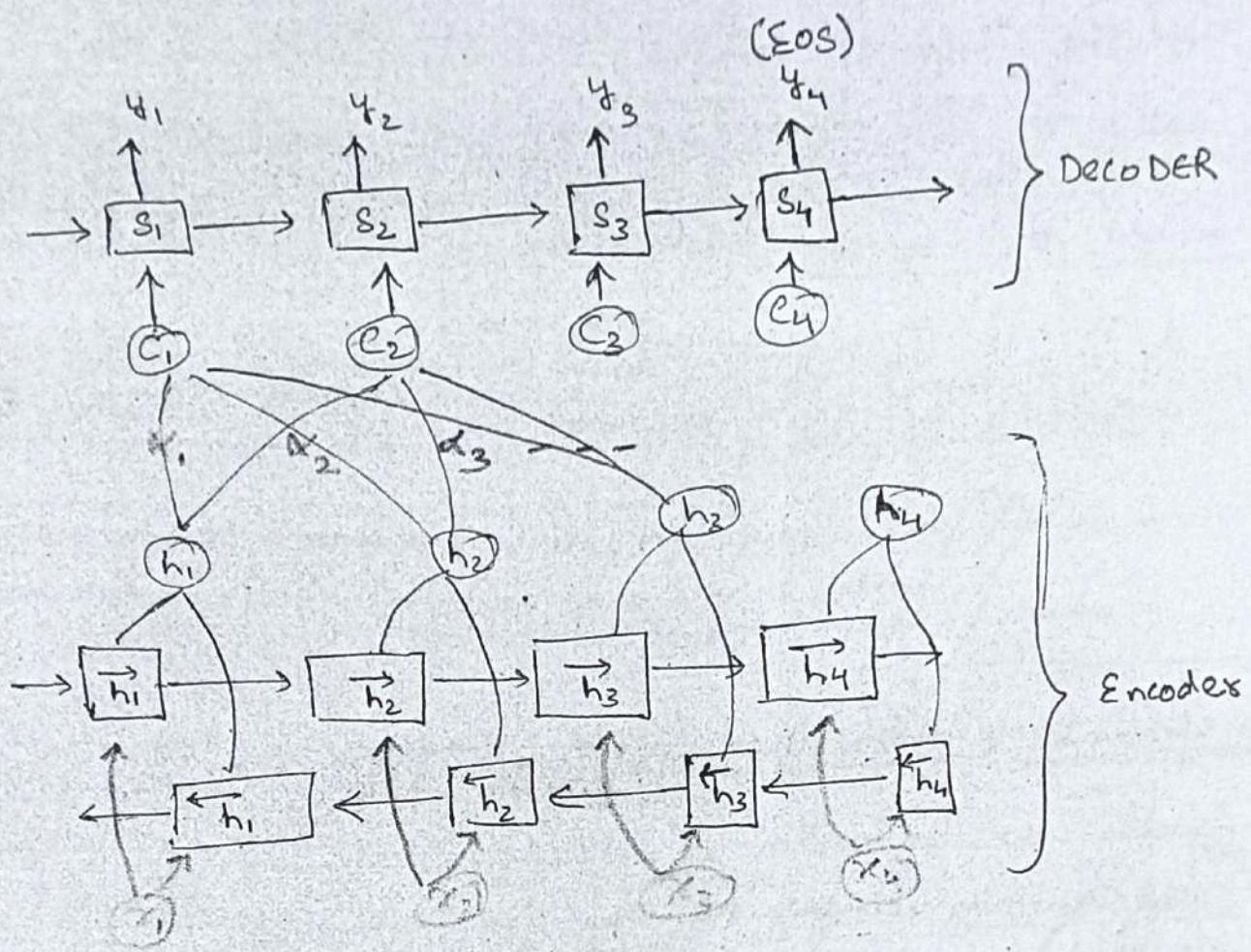


→ It shows that, as the Sentence length increasing the bleu Score is coming down.



→ In the order to solve the problems faced by seq2seq we use Bidirectional-LSTM RNN. as we can see in this graph, the accuracy has improved a lot.

★ BiLSTM RNN (Bidirectional LSTM RNN)



- x = Input
- h = Output of various hidden layers
- c = context vector
- α = (alpha) Sigmoid activation function
- y = output
- s = state of the decoder

→ Conditions with respect to α_i

- ① The sum of all the alphas within a time span should be equal to 1.
- ② we can calculate the context vector by multiplying the output of bidirectional LSTM (h) with the alpha values. like how in deep learning, we multiply inputs with the weights, we get some kind of output. In the similar way.

Feedforward Neural Network →

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

- The reason researchers may have applied this feed forward neural network, is basically for the back propagation. because we need to update these values of $\alpha_1, \alpha_2, \dots, \alpha_n$. unless and until we don't update the weights, you won't get good accuracy.
- By applying Softmax function, we can make sure that the value of (α_{ij}) is equal to 1.
- The weight (α_{ij}) of each annotation (h_j) is computed by:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

where $e_{ij} = a(\text{feed-forward neural network}, s_{i-1}, h_j)$

Transformers : Attention is all you need

- It adopts the mechanism of self-attention.
- It differentially weight the significance of each part of the input data.
- Like RNN, transformers are designed to handle sequential input data, such as natural language, for tasks such as translation and text summarization.
- However, unlike RNNs, transformers do not necessarily process the data in order. ~~rather~~, the attention mechanism provides context for any position in the input sequence.

Advantages of transformers

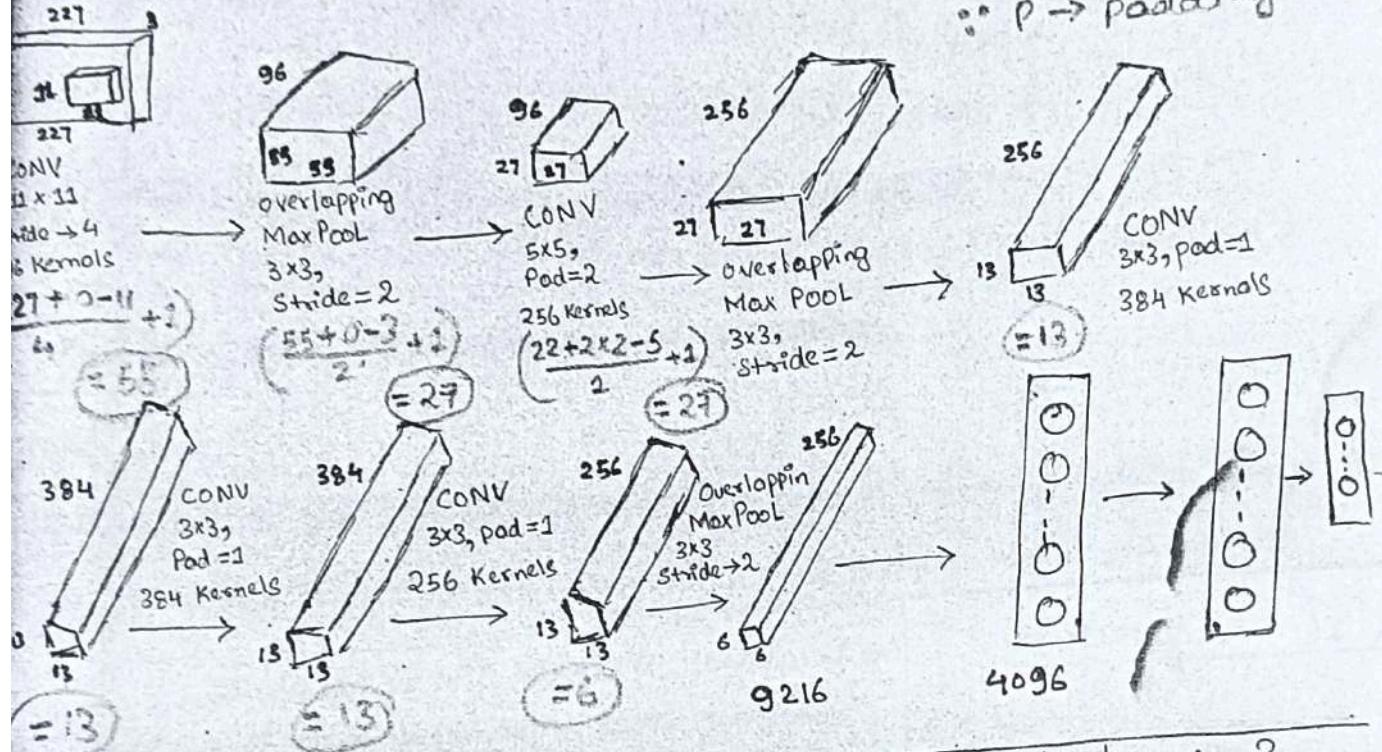
- ① They hold the potential to understand the relationship between sequential elements that are far from each other.
- ② They are way more accurate.
- ③ They pay equal attention to all the elements in the sequence.

- ④ Transformers boast of faster processing. In other words, they can process more data in lesser time.
- ⑤ They could work with virtually any kind of sequential data.
- ⑥ They can predict what could happen next.
- ⑦ Transformers serve to be helpful in anomaly detection.

ALEXNET ARCHITECTURE

$$\text{formula} \rightarrow \left(\frac{n+2p-f}{s} \right)$$

$\therefore p \rightarrow \text{padding}$



- ⑧ Alexnet architecture consists of 5 convolutional layers, 3 max-pooling layers, 2 normalization layers, 2 fully connected layers, and 1 Softmax layer.
- ⑨ Each convolutional layer consists of convolutional filters & a nonlinear activation function ReLU.
- ⑩ The pooling layers are used to perform max pooling.
- ⑪ Input size is fixed due to the presence of fully connected layers.
- ⑫ The input size is mentioned at most of the places as $224 \times 224 \times 3$ but due to some padding which happens it works out to be $227 \times 227 \times 3$. It has 60 million parameters.

→ AlexNet is a leading architecture for any object-detection task and may have huge applications in the computer vision sector of artificial intelligence problems. In the future, AlexNet may be adopted more than CNNs for image tasks.

Important things to learn in Reinforcement

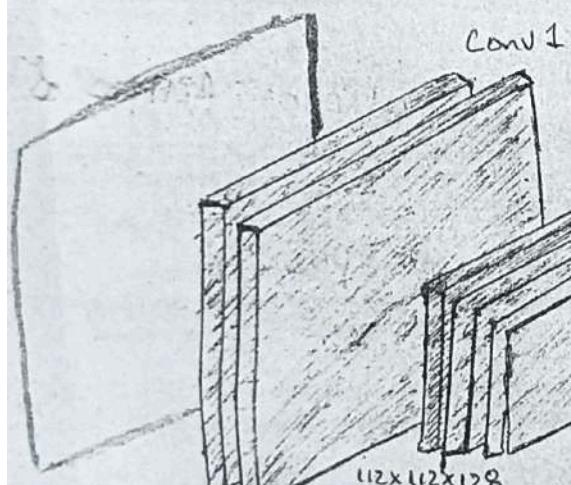
→ Deep Q-Networks ::

→ Deep Deterministic Policy Gradient

→ ~~Poison~~ ABC Algo

⇒ Q-Learning

VGG 16 Architecture



→ In Alexnet, we have lot of variations, something that you don't find in VGG16
→ That's why people prefer VGG16.
→ all convolution layers :-

- filter size = 3×3
- stride = 1
- padding = same

- convolution + ReLU
- max pooling
- fully connected + ReLU

Hummingbird

- Hummingbird is a library for compiling trained traditional ML models into tensor computations.
- Hummingbird allows users to seamlessly leverage neural network frameworks (such as PyTorch) to accelerate traditional ML models.
- In general, Hummingbird syntax is very intuitive and minimal. To run your traditional ML model on DNN frameworks, you only need to import hummingbird-ml and add to ("dnn_framework") to your code. ~~Below is an example:~~

TensorDash

TensorDash is an application that lets you remotely monitor your deep learning model's metrics and notifies you when your model training is completed or crashed.