

Deep Learning Examples.

1. Self driving car
2. Digital Assistants : Alexa, Siri
3. Text to speech Synthesis
4. Language identificatⁿ
5. Optical character recognition
6. Translation.
7. Image Caption Generation - Surveillance cameras.
8. Video to textual description (Annotation to video).
9. Synthetic passing for NLP

Artificial Intelligence.

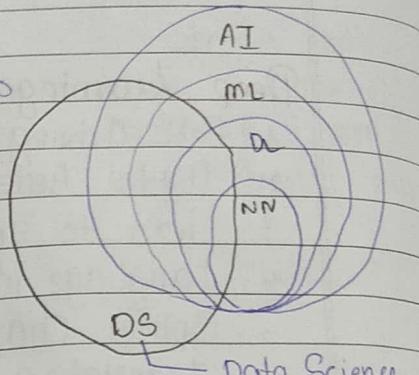
- It is just a program or machine's ability to learn & think for itself, in such a way that it is capable of developing skills as it practises & gains experience. Its purpose is to support & enhance human activities.
- Enables machine to think w/o human intervention & take decision.
- eg. are given above.

Machine Learning.

- Subset of AI.
- Provides statistical tools to explore & analyze data
- Method of data analysis that automates analytical model building
- Branch of AI based on the idea that systems can learn from data, identify patterns etc. & make decision.

Neural Networks.

- 1) makes use of neurons to transmit data as input & output values.
- 2) They are used to transfer data by network of connections.
- 3) 3 layers
 - Input layer
 - Hidden layer
 - O/P layer



Deep Learning.

- 1. Composed of several hidden layers while neural networks consists upto 3 layers.
- 2. Subset of ML.
- 3. This got created for like can machine to learn to think how actually human brain learn to think.
- 4. In DL, we create architecture called multi-neural network architecture. Idea of Deep neural network is to mimic human brain.

Data Science.

- 1. will work on above AI, ML, DL, using mathematical tools such as statistics, probability, Linear algebra, Differential calculus etc.

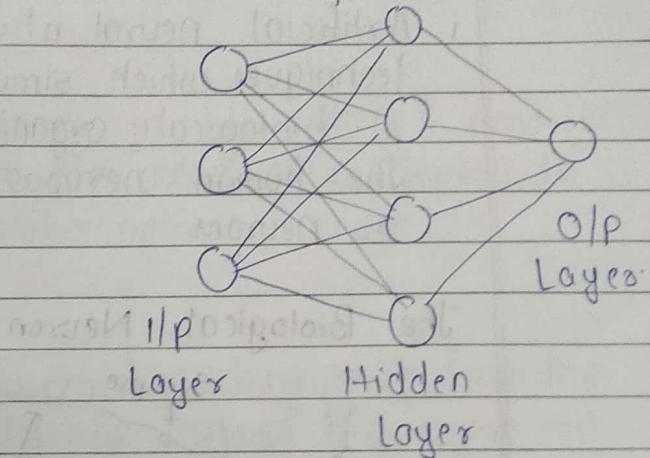
Some more examples of applied deep learning include the foll.

- 1. Automated image sharpening.
- 2. Automating image upscaling.
- 3. Wavenet: generating human speech that can immitiate anyone's voice.
- 4. Speech reconstruction from silent video.
- 5. Image autofill for missing regions.
- 6. Automated image captioning (process of generating textual description).

+ Turning hand-drawn doodles into stylized artwork.

- In Deep learning, there are 3 techniques.
 1. ANNs : Problems that involved with the numbers are solved with the help of ANN.

fig. Basic NN.



2. Convolution NN (CNN) : If the i/p is in images, then we use CNN. Advanced CNN known as Transfer learning.

3. Recurrent Neural N/w : If data is time series or kind of data, we use RNN.

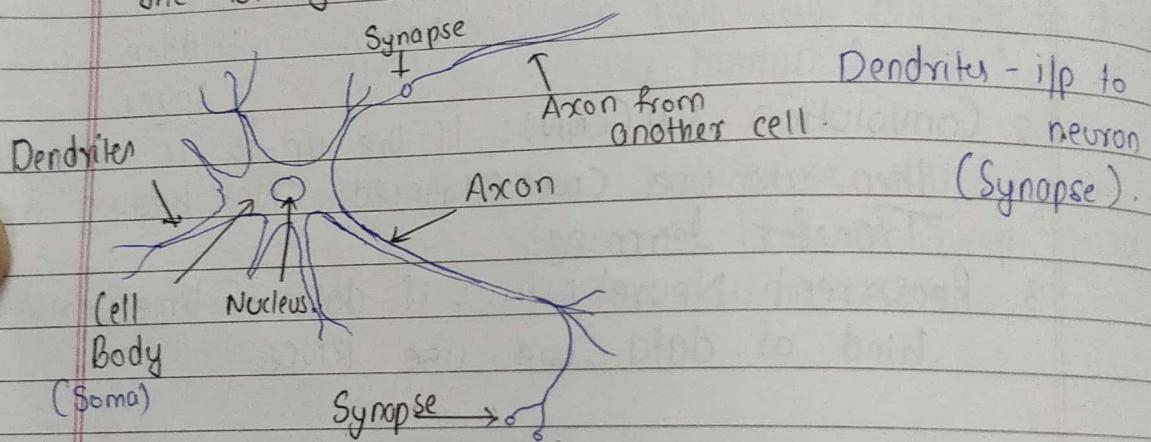
Tue.
07/01/20

INTRODUCTION TO NEURAL NETWORKS

* Introduction (Josh)

1. Artificial neural nw (ANNs) are mlc learning techniques which simulate the mech of learning in biological organisms.
The human nervous system contains cells, referred as neuron.

The Biological Neuron.



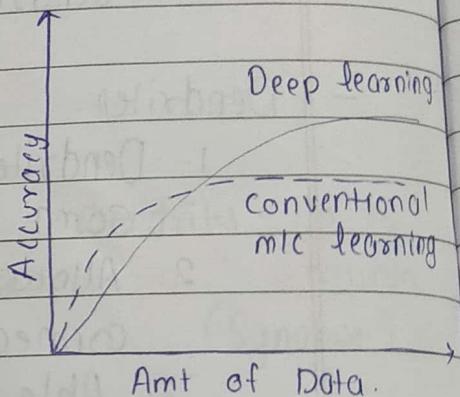
- The biological neuron is a nerve cell that provides the fundamental functional unit for the nervous systems of all animals.
- Neurons exist to communicate info & pass electrochemical impulses.
- The neuron is made up of a nerve cell consisting of soma that has many dendrites but only one axon.

- Axons are nerve cell fibers with special cellular extensions that come from the cell body. The axon can branch hundreds of times.
- Dendrites are thin structures that extend from main body.
- Synapses
 1. Synapses - connecting jn b/w axon & dendrite.
 2. Sends sig from the axon of a neuron to dendrite of another neuron.
- Dendrites.
 1. Dendrites have fibers branching out from the soma in a bushy nw around the nerve cell.
 2. Allows the cell to receive signals from connected neighbouring neurons.
 3. Able to perform multiplication by the dendrite weight value.
 4. ↑ or ↓ in ratio of neuro transmitters.
- Axons
 1. Single, long fiber extending from the main soma. (1cm in length & 100 times dia of soma)
 2. Branches 100s of times & connects to other dendrites.
- Information flow across the biological neuron.
 1. Synapses that ↑ the potential & those that decrease the potential.
 2. Neuron also have been shown to form new connections every time & even migrate.
 3. These combined mechanics

- Total connections in the human brain.
- 1. Biological NN(brain) are composed of roughly 86 billion neurons connected to many other neurons.
- 2. More than 500 trillion connectⁿ b/w neurons in the human brain.

* Humans Vs Computers : Stretching the limits of Artificial Intelligence. (Ch 9).

- for smaller data mlc learning is better coz of more choices, ease of model interpretation.



* Neural Networks (from Josh.)

1. NN nlws are computational model that shares some properties with animal brain in which many simple units are working in parallel with no centralized control unit.
2. The behaviour of NN is shaped by its nw arch.
3. Network's arch can be defined by.

No. of neurons.

No. of layers

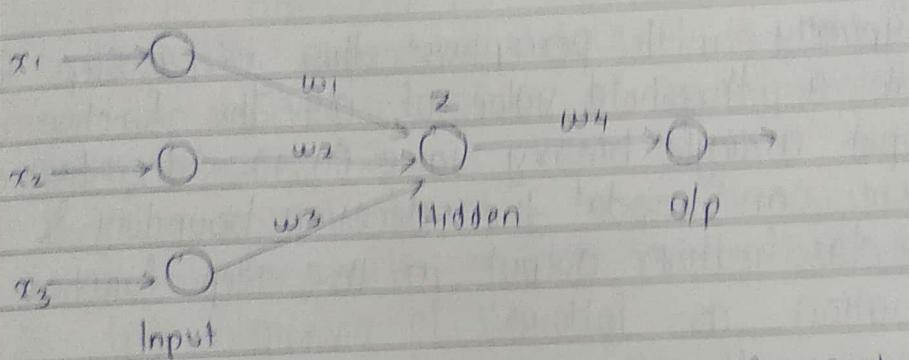
Types of connectⁿ b/w them.

Wednesday
08/10/20

- From biological to artificial:
- 1. Animal brain is fundamental components of the mind
 - 2. Research has shown ways of mapping functioning of human brain & track signals as they move through neurons

The Perceptron

The perceptron is a linear model used for binary classification. In the field of neural networks the perceptron is considered an artificial neuron.



$$y = x_1w_1 + x_2w_2 + x_3w_3 + \text{bias}$$

= $\text{Act}(y)$
↳ Activation fn.
= $z \cdot w_4$

- Types of Activation fn.

$$\text{Act}(y) \rightarrow \text{Sigmoid} = \frac{1}{1+e^{-y}}$$

↳ ReLU

The perceptron is a linear-model binary classifier with a simple input-output relationship as depicted.

It shows we're summing n no. of inputs times their associated weights & then sending this "net input" to a step function with a defined threshold.

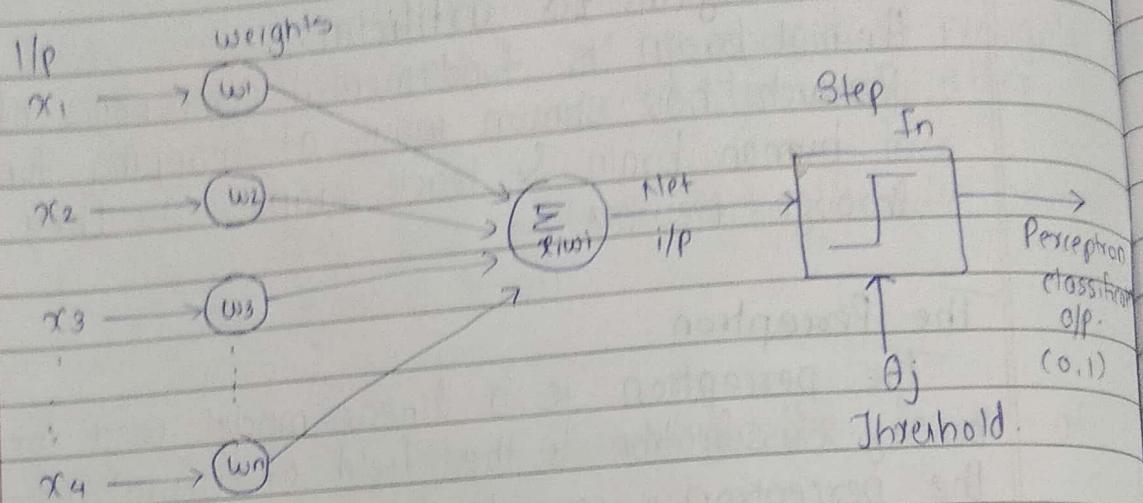


fig. Single layer perceptron.

- Typically with perceptrons, this is a step function with a threshold value of 0.5. This function will output a single binary value (0, 1), depending on input.
- We can model the decision boundary & the classification output in the step function equation as follows:

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

- To produce the net input Act fn, we take dot product of i/p & connect'g wts.

Function Parameter	Description
w	- Vectors of real valued weights on the connect'
$w \cdot x$	- Dot product ($\sum_i w_i x_i$)
n	- No. of i/p's to perception
b	- The bias term (input values does not affect bias).

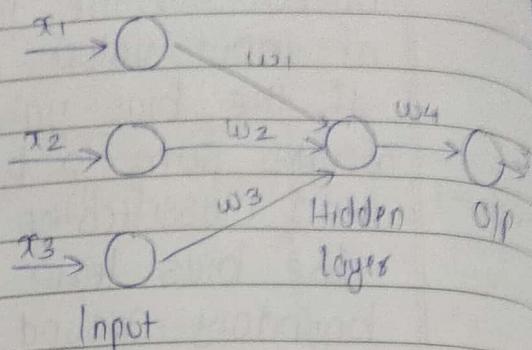
- The o/p of the step fn (activation fn) is the output for the perceptron & gives us a classification of input values.
- If the bias value is -ve , it forces the learned weights sum to be a much greater value to get a 1 classification o/p.
- The bias term in this capacity moves the decision boundary around for the model.
- Input values dont affect the bias, but bias is learned through perceptron learning algo.
- The algo will not terminate if the learning input is not linearly separable.
- A linearly separable dataset is one for which we can find the values of a hyperplane that will cleanly divide the two classes of the data set.
- The perceptron learning algorithm initializes the weight vector with small random values or 0.0s at the beginning of training.
- The perceptron learning algorithm takes each input record as we can see in the fig. & computes the output classification to check against the actual classification label.
- To produce the classification , the columns (features) are matched up to weights where n is the no. of dimensions in both our inputs & weights
- The first input value is the bias input , which is always 1.0 because we dont affect the bias input.
- The first weight is our bias term in this diag.
- The dot product of the input vector & weight vector gives us the i/p to our actn fn.

Sat.
11/01/20

Page 10

Types of Activation fn.

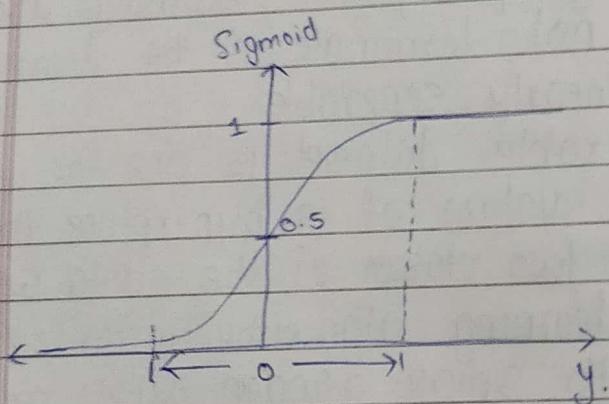
1. Linear
2. Non-Linear
3. Sigmoid (Logistic)
4. Tanh
5. ReLU



$$\text{Sigmoid fn.} = \frac{1}{1 + e^{-y}}$$

$$y = x_1w_1 + x_2w_2 + x_3w_3 + \text{bias}$$

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



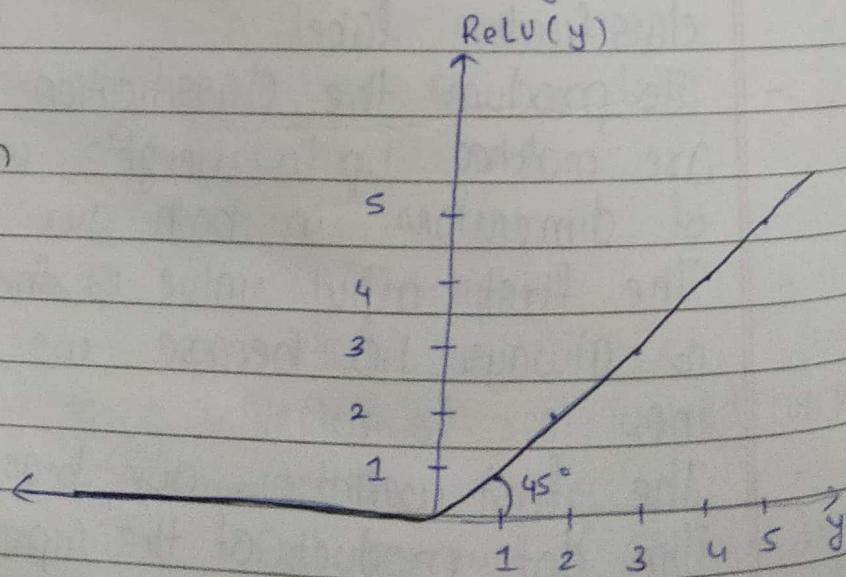
Mostly used in
o/p layer

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

$$y = -\text{ve}, \max(-\text{ve}, 0) = 0$$

$$y = +\text{ve}, \max(+\text{ve}, 0) = +\text{ve} = y$$

Mostly used in
hidden layer



* The Basics Architecture of Neural Networks (Chau Ag)

Single Computational Layer: The Perceptron.

1. In any type of application, one would have historical cases in which the class variable is observed
2. other (current) cases in which the class variable has not yet been observed but needs to be predicted.

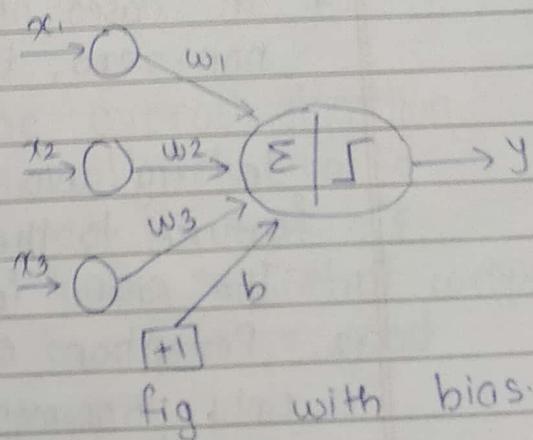
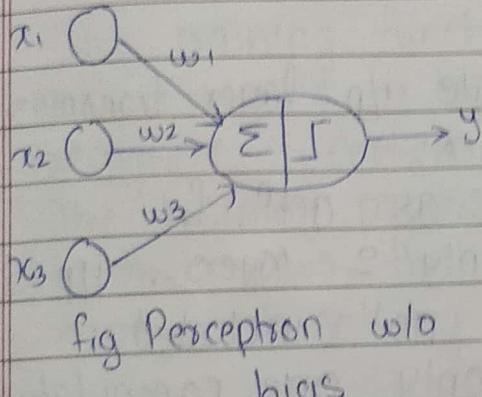


fig. basic Arch of Perceptron.

2. The input layer contains d nodes that transmit the d features $\bar{x} = [x_1 \dots x_d]$ with edges of weight $\bar{w} = [w_1 \dots w_d]$ to an o/p node.
3. The input layer does not perform any computation in its own.
4. The linear function $\bar{w} \cdot \bar{x} = \sum_{i=1}^d w_i \cdot x_i$ is computed at the output node.
5. Subsequently, the sign of this real value is used in order to predict the dependent variable of x . Therefore, the prediction \hat{y} is computed as follows.

$$\hat{y} = \text{Sign}\{\bar{w} \cdot \bar{x}\} = \text{Sign}\{\sum_{j=1}^d w_j \cdot x_j\}$$

6. The sign function maps a real value to either

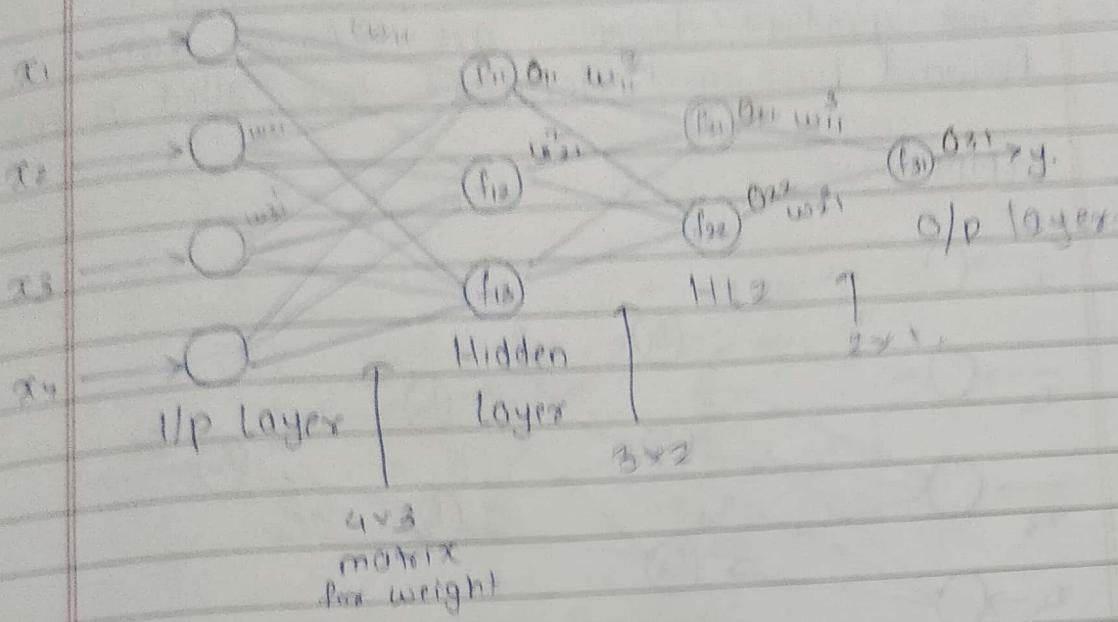
11. \hat{y} , which is appropriate for binary classification.
12. Note the cap mark on top of the variable \hat{y} to indicate that it is predicted value rather than an observed value.
13. The error of the prediction is therefore $E(\hat{y}) = \hat{y} - y$.
14. In cases where the error value $E(\hat{y})$ is non-zero, the weights in the neuron.

15. In fig w/o bias, single i/p layer transmits the features to the o/p layer.
16. The sign in ~~several~~ as actn fn.
17. Perception contains only 2 layers, i/p layer only transmits hence not included.
18. Therefore it has only one computational layer & is called as single layer feed forward net.
19. In many settings, there is an invariant part of prediction which is referred to as bias.
- eg consider a setting in which the feature variables are mean centered, but the mean of the binary class prediction from $\{1, -1\}$ is not 0.
20. This will tend to occur in situations in which the binary-classification distribution is highly imbalanced.

Invariant definition: an entity, quantity etc that is unaffected by a particular transformation.

Use of additional input, which is called bias, of neuron improves properties of the neuron. It

- The bias value allows the actn fn to be shifted left or right, to better fit the data. It is always 1.
- A weight represent the strength of connectn b/w units. If the weight from node 1 to node 2 has greater magnitude, it means neuron 1 has greater influence over neuron 2.
- Weights in NN can, & will, become +ve or -ve depending on the training data.
-
- The training function is the overall algorithm that is used to train the neural network to recognize a certain input & map it to an o/p.
- The training process involves finding a set of weights in the nws that proves to be good or good enough of solving specific problem.
- A learning fn deals with individual wts & thresholds & decides how these would be manipulated.



$$w_{\text{new}} = w_{\text{old}} + \eta \cdot \frac{\partial L}{\partial w}$$

Learning rate.

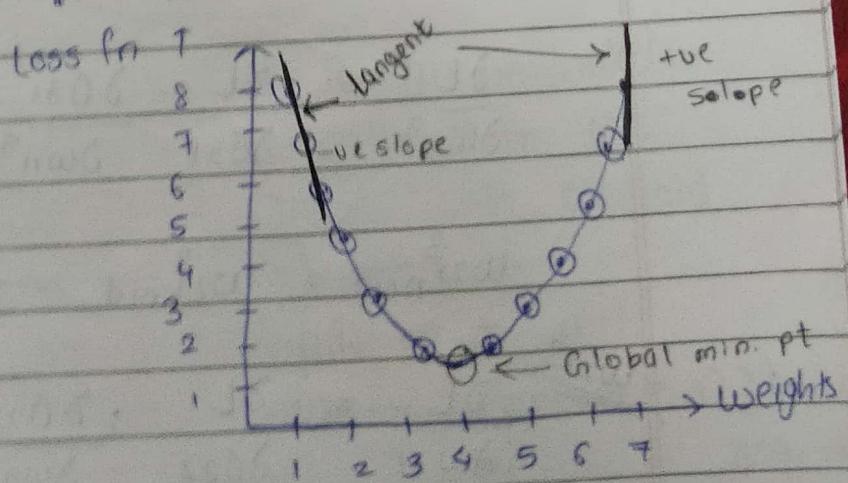
$\text{Loss fn} = (y - \hat{y})^2$
 predicted value.
 reduce optimiser.

$$\Rightarrow \frac{1}{m} (y - \hat{y})^2$$

where $m = \text{no. of o/p neurons}$.

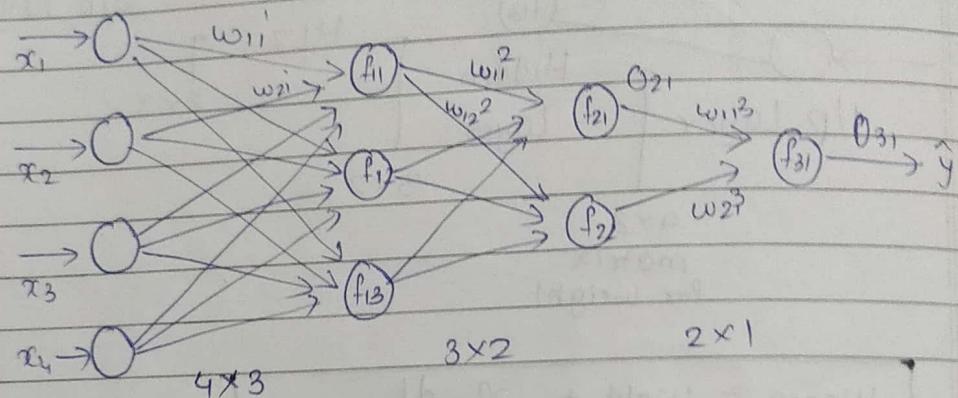
Gradient descent

update the weights such a way that \hat{y} to be closer to y . find that more efficient weights to be updated. to get the best predicted value \hat{y} .



- Optimizers will reduce the loss value.
- Some gradient descent are - Adam, SGD
- forward backward iterations are called as epoch

MNN



l/p	Hidden	Hidden	O/p
Layers	Layer 1	Layer 2	Layers

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

$$w_{11}^3 \text{ new} = w_{11}^3 \text{ old} - \eta \left[\frac{\partial L}{\partial w_{11}^3} \right]$$

$$w_{11}^3 = w_{11}^3 \cdot o_{21} + w_{21}^3 \cdot o_{22}$$

Chain rule on derivative in NN

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

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

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial w_{21}}$$

$$\underline{w_{ii}^2} \quad w_{ii}^2_{\text{new}} = w_{ii}^2_{\text{old}} - \eta \frac{\delta L}{\delta w_{ii}^2}$$

$$O_{21} = w_{ii}^2 \cdot O_{ii}$$

$$\frac{\delta L}{\delta w_{ii}^2} = \frac{\delta L}{\delta O_{31}} \times \frac{\delta O_{31}}{\delta O_{21}} \times \frac{\delta O_{21}}{\delta w_{ii}^2}$$

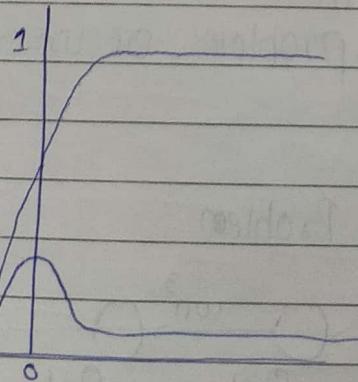
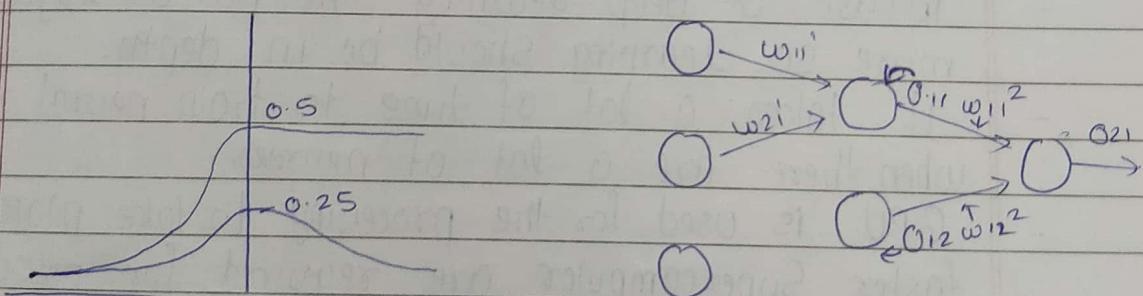
$$\underline{w_{i2}^2} \quad w_{i2}^2_{\text{new}} = w_{i2}^2_{\text{old}} - \eta \frac{\delta L}{\delta w_{i2}^2}$$

$$O_{22} = w_{i2}^2 \cdot O_{i2}$$

$$\frac{\delta L}{\delta w_{i2}^2} = \frac{\delta L}{\delta O_{31}} \cdot \frac{\delta O_{31}}{\delta O_{22}} \cdot \frac{\delta O_{22}}{\delta w_{i2}^2}$$

Vanishing Gradient Problem.

$$w_{ii}^{\text{new}} = w_{ii}^{\text{old}} - \eta \frac{\delta L}{\delta w_{ii}}$$



- It is an unwanted situation in NN
- Occurs with sigmoid activation fn
- Update of w_{ii}

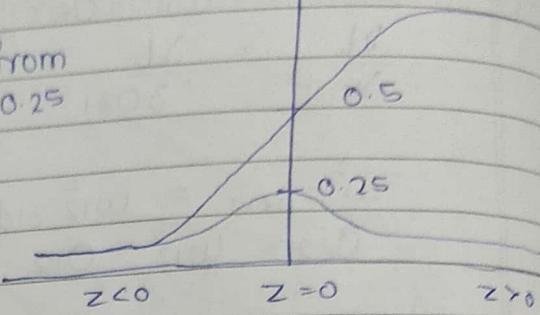
$$w_{ii}^{\text{new}} = w_{ii}^{\text{old}} - \eta \frac{\delta L}{\delta w_{ii}^{\text{old}}}$$

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

$$\frac{\partial L}{\partial w_{ii}} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial w_{ii}}$$

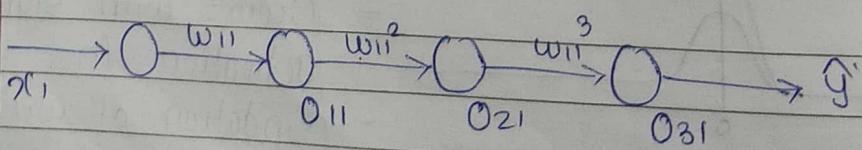
$$\text{sigmoid}(z) = \frac{1}{1+e^{-z}}$$

$\delta(\text{Sigmoid})$ varies from
 $\frac{\partial L}{\partial z}$ $0 - 0.25$



- since $\frac{\partial L}{\partial w_{ii}}$ is a very small value that is subtracted with w_{ii} .
- Using ReLU function the $\frac{\partial L}{\partial w_{old}}$ will give a higher value which gives the difference.
- In case of deep learning, the no. of layers are more, so learning should be in depth.
- It takes a lot of time to train neural netw when there are a lot of neurons.
- GPU is used for the processing to take place faster. Supercomputers are required for processing the large neural networks.
- Vanishing gradient problem occurs because of activation function.

Exploding Gradient Problem.



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

$$\frac{\delta L}{\delta w_{11}} = \frac{\delta L}{\delta o_{31}} \cdot \frac{\delta o_{31}}{\delta o_{21}} \cdot \frac{\delta o_{21}}{\delta o_{11}} \cdot \frac{\delta o_{11}}{\delta w_{11}}$$

$$o_{21} = \text{Act}(z)$$

$$z = w_{11}^2 \times o_{11} + b_2$$

↑
result of
step 1.

biases to
then 2nd neuron.

$$\frac{\delta L}{\delta w_{11}} = \frac{\delta L}{\delta o_{21}} \times \frac{\delta o_{21}}{\delta o_{11}} \times \frac{\delta o_{11}}{\delta w_{11}}$$

$$\frac{\delta o_{21}}{\delta o_{11}} = \frac{\partial [\text{Act}(z)]}{\partial o_{11}} = \frac{\partial [\text{Act}(z)]}{\partial z} \times \frac{\partial z}{\partial o_{11}} \quad (\text{chain rule})$$

(0 to 0.25)

Substituting for z ,

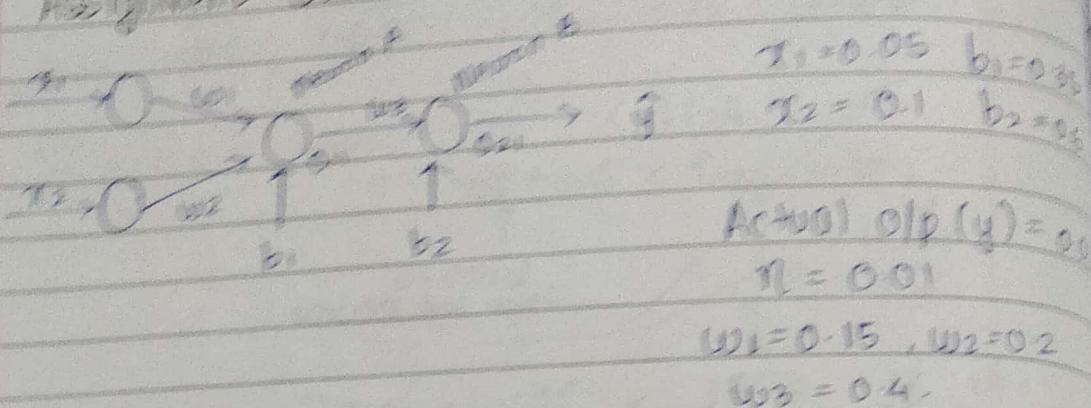
$$(0 \text{ to } 0.25) \times \frac{\partial [w_{11}^2 \times o_{11} + b_2]}{\partial o_{11}}$$

$$(0 \text{ to } 0.25) \times [w_{11}^2 + 0]$$

$$(0 \text{ to } 0.25) \times w_{11}^2$$

This will usually be a large value when compared to vanishing gradient problem. So there will be a variation in w_{11} new & w_{11} old. So it will never reach a global minima region.

Assignment 0234-20



Analyse the given NN for one epoch with
1) Forward Propagation 2) Backward Propagation

Neuron A

$$z_1 = w_1 x_1 + w_2 x_2 + b_1 = 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35 \\ = 0.3775$$

$$o_{11} = \text{Act}(z_1) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-0.3775}} = \frac{1}{1.6856} = 0.5933$$

Neuron B

$$z_2 = w_3 o_{11} + b_2 = 0.4 \times 0.5933 + 0.6 = 0.8373.$$

$$o_{21} = \text{Act}(z_2) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-0.8373}} = \frac{1}{2.8627} = 0.411 = \hat{y}$$

$$\text{Error} = (y - \hat{y}) = (0.01 - 0.411) = -0.401$$

$$\text{Loss} = (y - \hat{y})^2 = (-0.401)^2 = 0.1608.$$

To update the weight w_3 ,

$$w_{3\text{new}} = w_{3\text{old}} - n \frac{\partial L}{\partial w_3} \quad \dots \dots \dots (1)$$

Now,

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial w_3} \quad \dots \dots \dots (2)$$

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

$$\begin{aligned} \frac{\partial L}{\partial o_{21}} &= \frac{\partial [(y - o_{21})^2]}{\partial o_{21}} = 2(y - o_{21}) \times (-1) \\ &= -2(y - o_{21}) \\ &= -2(0.01 - 0.411) = 0.8 \end{aligned}$$

$$\frac{\partial o_{21}}{\partial w_3} = \frac{\partial o_{21}}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_3}$$

$$\text{we know, } o_{21} = \text{Act}(z_2) = \frac{1}{1 + e^{-z_2}}$$

$$\begin{aligned} \frac{\partial o_{21}}{\partial z_2} &= \frac{\partial (\text{Act}(z_2))}{\partial z_2} = \text{Act}(z_2) \times (1 - \text{Act}(z_2)) \\ &= 0.411 \times (1 - 0.411) = 0.242 \end{aligned}$$

$$\text{Now, } z_2 = w_3 o_{11} + b_2$$

$$\frac{\partial z_2}{\partial w_3} = o_{11} + 0 = 0.5933$$

In eq (2)

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial o_{21}} \times \left\{ \frac{\partial o_{21}}{\partial z_2} \times \frac{\partial z_2}{\partial w_3} \right\} = 0.802 \times 0.242 \times 0.5933 = 0.1153$$

In eqn (1)

$$w_3 \text{ new} = 0.4 - 0.01 \times 0.1153 = 0.3988$$

Now to update the weight w_1 ,

$$w_1 \text{ new} = w_{1 \text{ old}} - \eta \frac{\partial L}{\partial w_1} \quad \dots \dots (3)$$

Now,

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial z_2}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

\downarrow
 $\frac{\partial}{\partial w_1} [w_1 x_1 + w_2 x_2] = x_1 + 0 + 0 = 0$

\downarrow
 $\frac{\partial}{\partial z_1} (\text{Act}(z_1)) = \text{Act}(z_1) \times (1 - \text{Act}(z_1)) = 0$

\downarrow
 $\frac{\partial}{\partial o_{11}} [w_3 o_{11} + b_2] = w_3 + 0 = 0.4$

$0.802 \quad 0.242$

$$\text{Now, } = 0.802 \times 0.242 \times 0.4 \times 0.2413 \times 0.05 = 0.94$$

$$w_1 \text{ new} = 0.15 - 0.01 \times 0.93 \times 10^{-3} = \underline{\underline{0.14999}}$$

next class - 04/02/20

Drop out layers in Multi Neural Network.

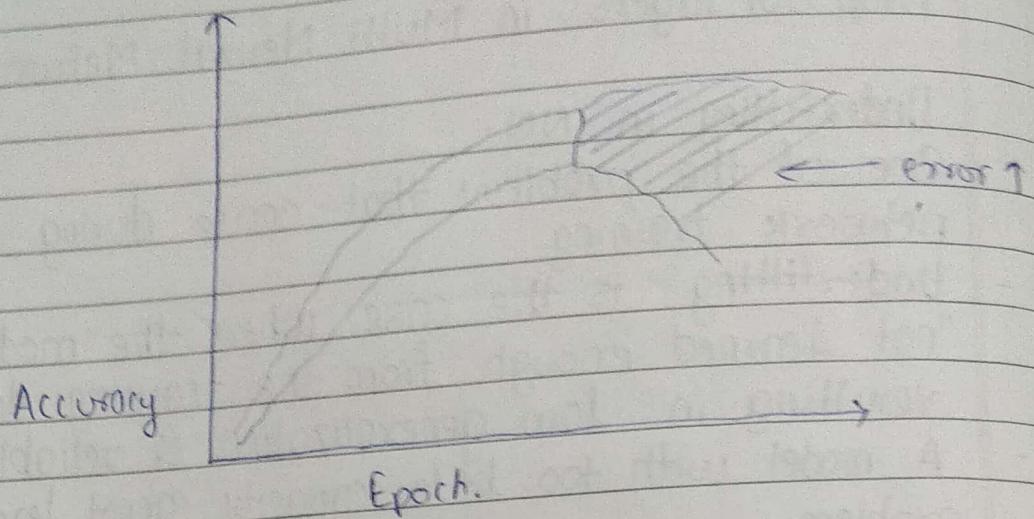
Underfitting in NN.

- One of the problems that occur during neural network training.
- Underfitting is the case where the model has "not learned enough" from the training data, resulting in low generalization & reliable predictn
- A model with too little capacity cannot learn the problem.
- Too little learning will make the model to perform poorly on the training dataset & on new data.
- An underfit model has high bias & low variance

Overfitting in NN

- One of the problems that occur during neural network training is called overfitting.
- The error on the training set is driven to a very small value, but when new data is presented to the network the error is large.
- The network has memorized the training examples, but it has not learned to generalize to new situations.
- A model with too much capacity can learn it too well & overfit the training dataset.
- With too much learning, the model will perform well on the training dataset & poorly on new data.
- An overfit model has low bias & high variance. The model learns the training data too well & performance varies widely with new unseen.

D



Good fit Model.

A model that suitably learns the training dataset & generalizes well to the old outdoor dataset

How to handle underfitting problem?

Underfitting can easily be addressed by increasing the capacity of the nlu.

How do you handle overfitting problem?

(i) Reduce the nlu's capacity by removing layers or reducing the number of neuron elements in the hidden layers.

(ii) Use Dropout layers which will randomly remove certain features by setting them to zero.

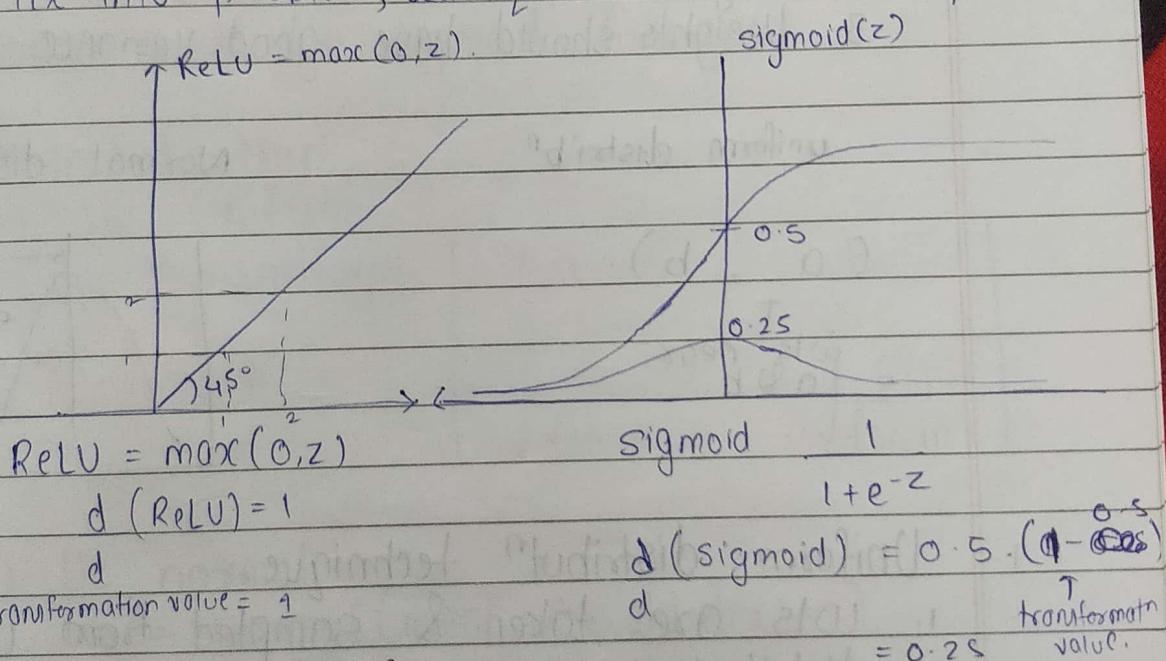
- Dropout layers are an easy & effective way to prevent overfitting in NN.

- A dropout layer randomly drops some of the connections b/w layers. This helps to prevent overfitting. We reduce the weights of the model to small during learning.

These techniques not only reduce overfitting but they can also lead to faster optimization of the model & better overall performance.

- In deep NN, there are many weights & bias parameters. When we have huge weights & bias parameters, the NN tends to overfit the data set problem on a particular use scales.
- In deep NN, underfit will never happen. Usually underfitting happens with one layer NN.

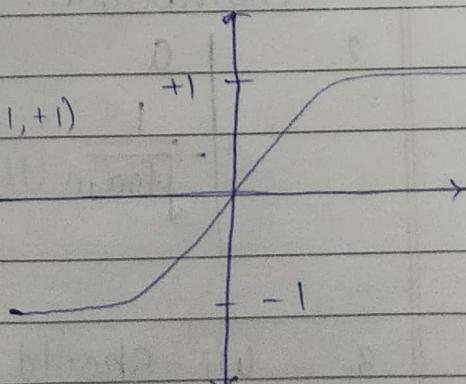
To fix this problem, technique used called dropout.



tanh Activation fn.

$$\tanh(z) = \frac{1-e^{-2z}}{1+e^{-2z}} = (-1, +1)$$

$$d(\tanh) \Rightarrow 0 \text{ to } 1.$$

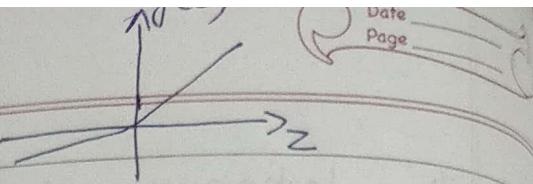


$$-1 \vee (1 - (-1))$$

$$1 - e^{-2x - \infty}$$

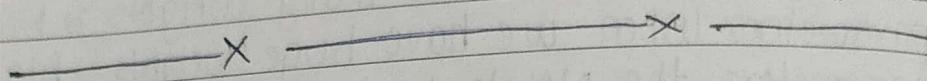
$$\frac{1 - e^{-\infty}}{1 + e^{-\infty}} = 0$$

Tue.
11/02/20



Leaky Relu

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise.} \end{cases}$$



Weight Initialization Techniques in NN

Key Points while initialising the weights.

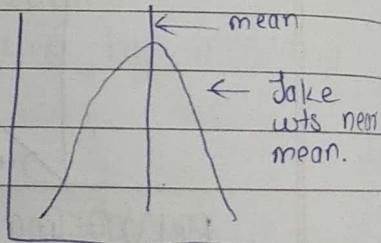
1. Weights should be small.
2. Weights should not be same.
3. Weights should have good variance.

uniform distribⁿ

$$(a, b)$$

↑
wts near
a & b

Normal distribⁿ

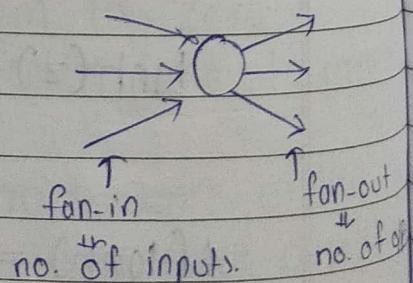


1. Uniform distribⁿ technique.

1. Wts are taken & sampled from uniform distribution.

2.

$$\left| \frac{1}{\sqrt{\text{fan-in}}} + \frac{1}{\sqrt{\text{fan-in}}} \right|$$



3. Wts should be in the above range.

4. Sigmoid actn fn works well with this distribⁿ.

$$GD \Rightarrow \text{Error} = (y - g)$$

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

$$\text{Cost} = \sum_{i=1}^{100} (y_i - g)^2$$

classmate

Date _____

Page _____

2. Xavier/Glorot distribution.

(i) Xavier/Glorot uniform distribution.

- w_i will be in the range from.

$$a = -\frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}} \quad b = \frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}$$

(ii) Xavier/Glorot normal distribution.

- $w_i = \text{Normal}(\text{mean}=0, \text{std dev}=\sigma)$

$$\sigma = \sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$$

- Works well for sigmoid function.

3. He initialization.

(i) He uniform

$$w_i = \text{Uniform.}$$

$$a = -\frac{\sqrt{6}}{\sqrt{\text{fan-in}}} \quad b = \frac{\sqrt{6}}{\sqrt{\text{fan-in}}}$$

(ii) He normal

- $w_i = \text{normal}(\text{mean}=0, \text{std dev}=\sigma)$

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

- Works well with ReLU actⁿ fn.

Optimizers

Gradient

Descent

Considers all the
data sets at once

Stochastic

Gradient

Descent

Wed.
12/02/20

classmate

Date _____
Page _____

Stochastic Gradient Descent.

$$\begin{aligned} \text{Error} &\rightarrow (y - \hat{y})^2 \\ \text{Loss} &\rightarrow (y - \hat{y})^2 \\ \text{Cost} &\rightarrow (y - \hat{y})^2 \end{aligned}$$

Instead of considering all data like GD, it only considers one data at a time.

Note:

$$\begin{aligned} \text{speed (GD)} &> \text{speed (SGD)} \\ \text{accuracy (GD)} &< \text{accuracy (SGD)} \end{aligned}$$

- To increase the speed, consider group of data at a time (i.e. subset of data). This maintains the accuracy (slightly less than b4) but ~~loss~~ much faster than before.

This is called as Minibatch SGD.

$$\text{Error} = (y - \hat{y})$$

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

$$\text{Cost} = \sum_{i=1}^k (y - \hat{y})^2 \quad \text{where } k < n$$

- Minibatch SGD is faster than GD.

It is a bad optimizer but converges to global minima faster than G.D.

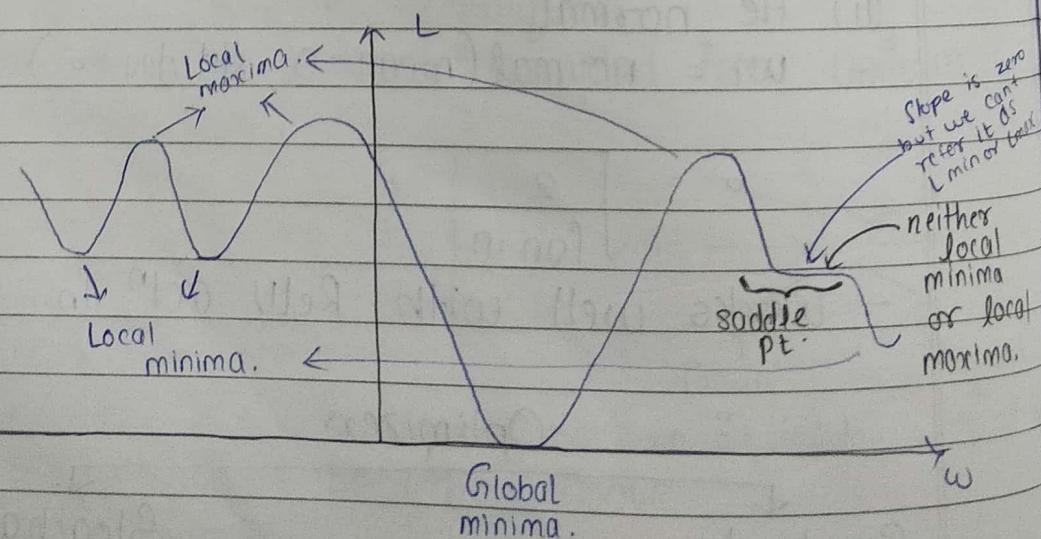


Fig.1. Gradient Descent curve.

Here noise is getting introduced in data ~~classmate~~
 This is an unwanted data structure removed
 in noisy data structure

- In GD the pts will be far away & on joining a smooth curve isn't obtained.
- In SGD there will be a large no. of pts \therefore more accuracy & on joining we get a smoother curve.

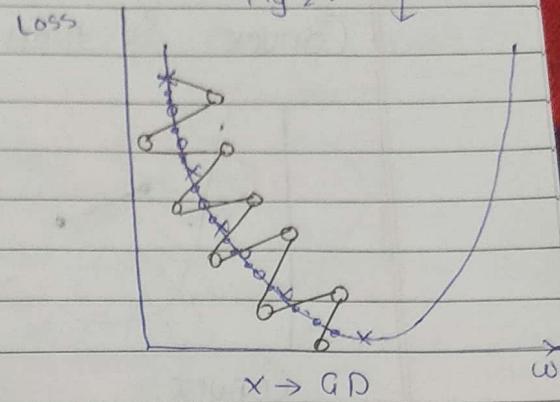


Fig. 2.

- In case of GD, memory & computational power reqd is more when compared to MSGD
- In that sense MSGD is faster than GD
 \therefore This is when the total data is very huge.

- In fig. 1. the loss $\uparrow \downarrow \uparrow \downarrow \dots$. At a certain point it reaches global minima (where $y - \hat{y} \approx 0$)
 \therefore (Here the slope will be zero).
 The remaining ones are local minima.
- The global minima is a point where the gradient descent curve completely converges.
 \therefore There is only one in each graph.
- The local minima are the points in which only in that locality, that is the best converging point.
- There will be advanced techniques with optimizers where these problems can be solved.
- The peak on the top are called local minima.
- The saddle point occurs occurs bw local minima & local maxima. The slope here is zero. It isn't global minima either.

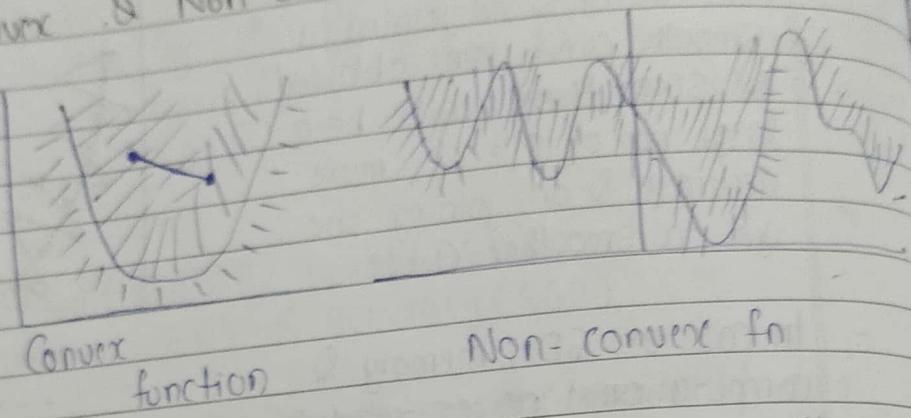
$$C = \frac{1}{2} (y_1 - \hat{y})^2 + \frac{1}{2} (y_2 - \hat{y})^2$$

classmate

Date _____

Page _____

Convex & Non-convex function.



- Area is divided into 2.
- Mark 2 pts & join them.
- The line remain in some region.
- Most Deep Neural Net will have non-convex fn
- Even if it crosses, it will be just once
- Here the pts will cross more than 1 region, multiple times.

Exponentially moving average.

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

$$V_1 = b_1$$

$$V_2 = \gamma V_1 + b_2 \quad 0 < \gamma \leq 1$$

$$\text{let } \gamma = 0.5$$

$$\therefore V_2 = 0.5V_1 + b_2$$

$$V_3 = \gamma V_2 + b_3$$

$$= \gamma(0.5V_1 + b_2) + b_3$$

$$V_3 = 0.25V_1 + 0.5b_2 + b_3$$

$$V_3 = 0.25b_1 + 0.5b_2 + b_3$$

This is included in NN to come out of noise data.

momentum.

Now, $w_{\text{new}} = w_{\text{old}} - (\gamma v_{t-1} + \eta \frac{\delta L}{\delta w})$

$$v_{t-1} = \frac{\delta L}{\delta w} \Big|_t + \gamma \frac{\delta L}{\delta w} \Big|_{t-1} + \gamma^2 \frac{\delta L}{\delta w} \Big|_{t-2} + \dots$$

SGD with momentum.

Smoothens the curve by removing noise from data points.

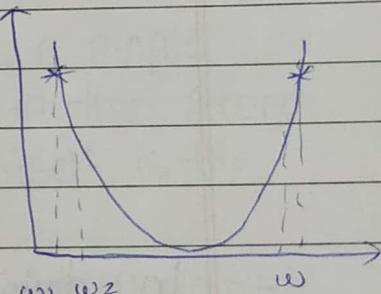
Adaptive gradient descent optimizers (Adagrad)

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

Dense \rightarrow non-zero

Sparse \rightarrow zero

x_1	x_2	x_3	x_n
0	1	1	0	1
0	0	0	1	0
1	1	0	1	0



$$w_t = w_{t-1} - \eta_t \frac{\delta L}{\delta w_{t-1}}$$

$$\eta_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

small +ve no.
(so that if $\alpha_t = 0$
 η_t shd not become ∞)

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

Exponentially weighted avg :

$$w_{\text{avg}} = \gamma w_{\text{avg}}^{(t-1)} + (1-\gamma) \frac{\delta L}{\delta w_{\text{avg}}^{(t-1)}}$$

$$\gamma = 0 \text{ to } 1$$

Hyperparameter:
Parameter that is set before training the NN.

e.g. [1/P
weights]

tuning
hyperparameter.

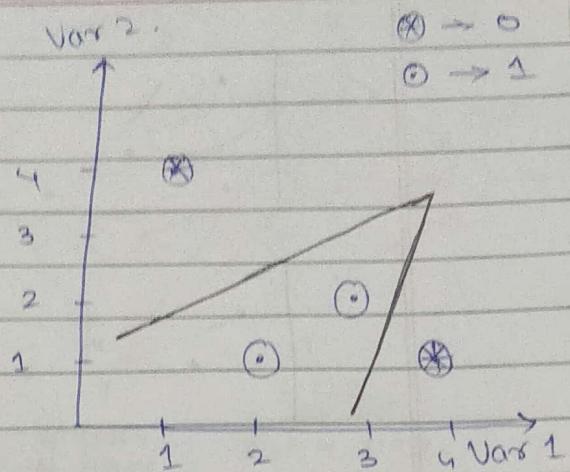
- These are used to take care of learning learning role.
- if there are γ are not used then we may have a problem of vanishing gradient.

e.g.	Variable 1	Variable 2	Class label
	1	4	0
	4	1	0
	3	2	1
	2	1	1

- In designing ANN archi,
 - Start by selecting the no. of neurons in the input & output layers.
 - 2 variables as inputs & hence 2 input nodes.
 - A binary classification problem, 1 output neuron
- For more than 2 classes, we'd need to create an output neuron for each class.
- To deducing the best number of hidden layers & neurons is to visualize the samples in a 2D graph.

"What is the minimum no. of connected lines to separate the classes?"
- No. of lines represents the no. of hidden neurons across all hidden layers.

- In this example, there is overlap b/w the two classes. Thus we cannot use just a single line to split the data. Minimum no. of lines is 2.

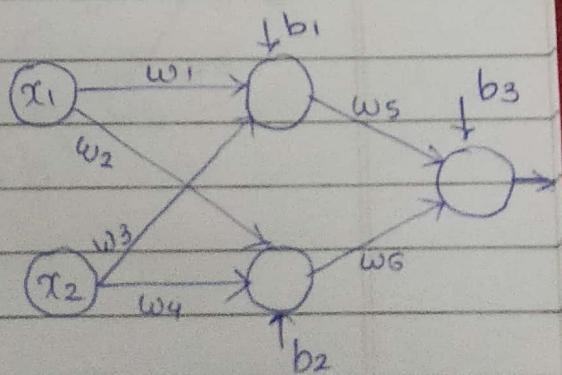


- No. of required lines is 2, NN will have 2 hidden neurons in the first hidden layer, where each neuron produces a line.

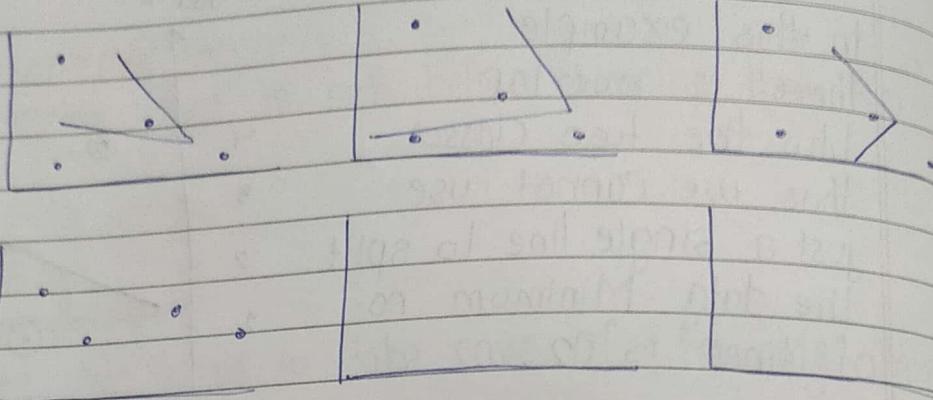
"How many hidden layers are required?"

- The preceding hidden layer connects these lines.
- The no. of connections defines the number of hidden neurons in the next hidden layer.
- Because there are just 2 lines, with a single connection (ie just single hidden neuron in the second hidden layer). The second hidden layer itself is the output layer.
- To conclude, 1 hidden layer with 2 hidden neurons.

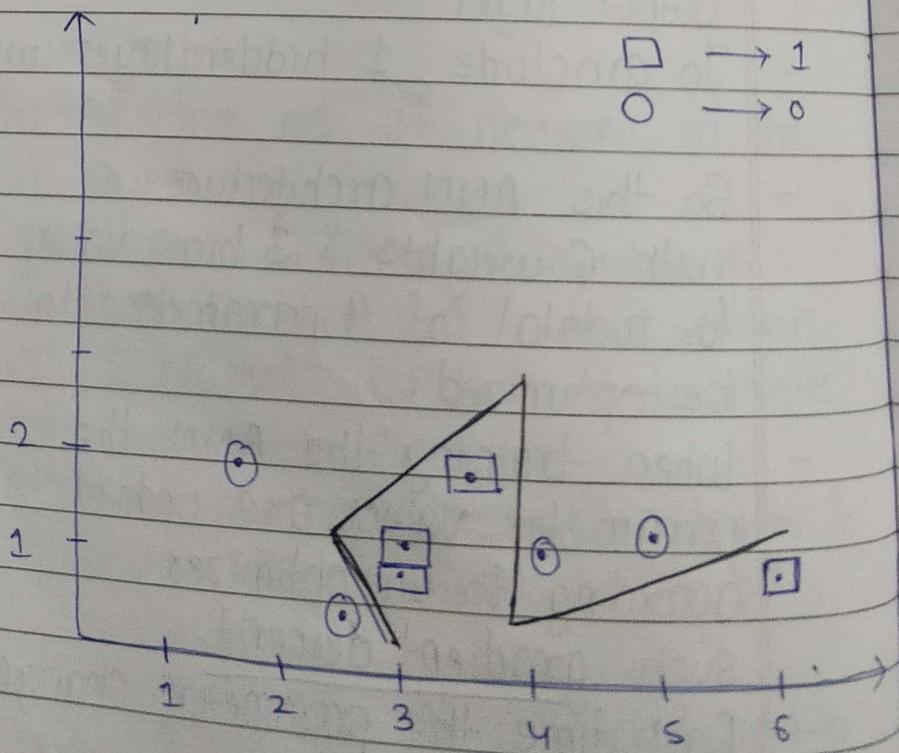
- So this ANN architecture is with 6 weights & 3 bias values for a total of 9 parameters to be optimized.



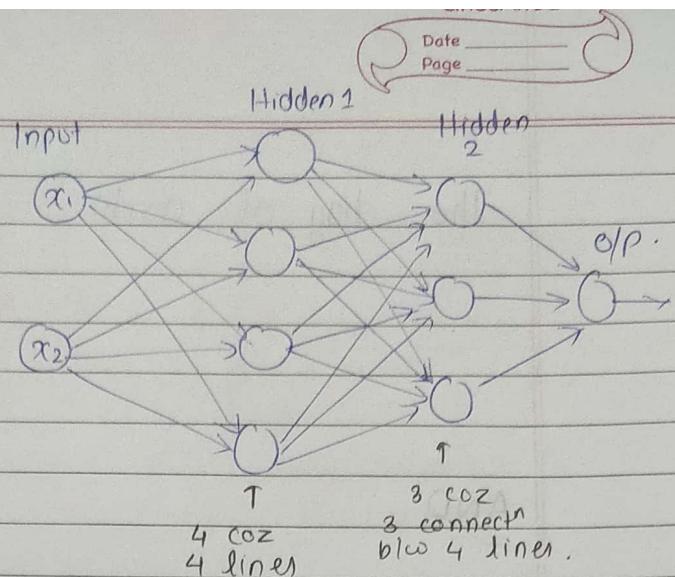
- When training the ANN, the parameter values are updated according to an optimizer, such gradient descent.
- Each time the parameters change, the lines are rearranged.



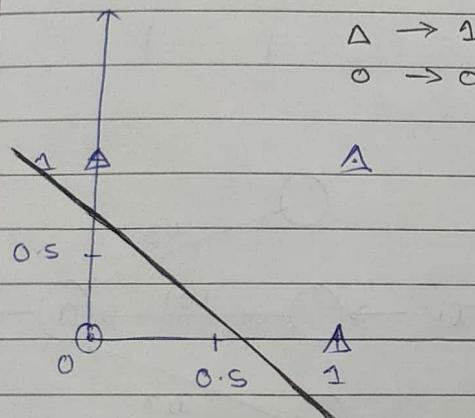
eg. 2.	Var 1	Var 2	Class Label.
5		1.5	0
4		1.25	0
1.5		1.125	0
2.5		1.0625	0
6		1.03125	1
3.5		1.015625	1
3		1.015625	1
3		1.015625	1



\therefore The archi



OR gate implementation in NN



OR gate

x_1	x_2	O/P
0	0	0
0	1	1
1	0	1
1	1	1

The o/p will be

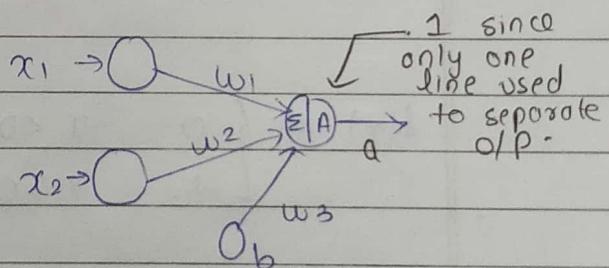
$$O/P = \text{ReLU}$$

$$(A) = \text{Activation}(w_1 x_1 + w_2 x_2 + w_3 b)$$

$$\text{Let } b = 1$$

$$w_3 = -0.5$$

$$w_1 = w_2 = 1.$$

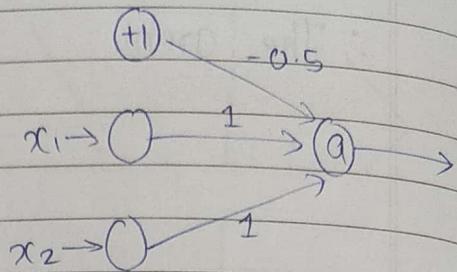


$$\therefore a = \text{Act}(x_1 + x_2 - 0.5)$$

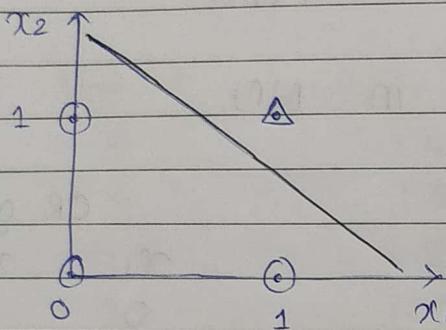
$$x_1 + x_2 - 0.5$$

x_1	x_2	$x_1 \text{ OR } x_2$	$x_1 + x_2 - 0.5$	a
0	0	0	-0.5	0
0	1	1	0.5	1
1	0	1	0.5	1
1	1	1	1.5	1

The diag of archi

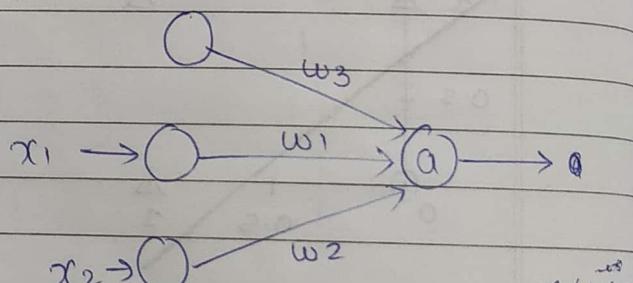


AND

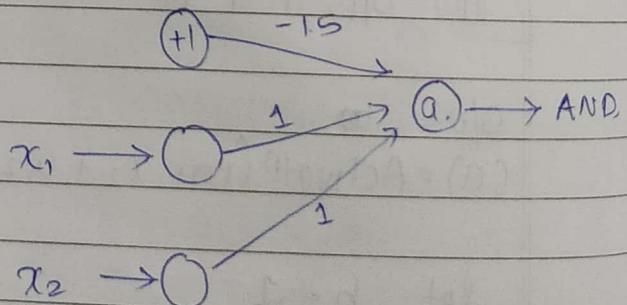


x_1	x_2	$x_1 \text{ AND } x_2$
0	0	0
0	1	0
1	0	0
1	1	1

$a = \text{Act}$



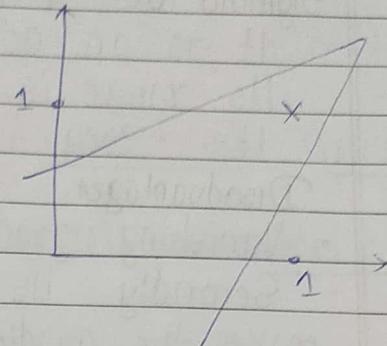
$$f(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$



X	Y	$X \text{ AND } Y$	$-1.5 + x_1 + x_2$	$a = f(w_1x_1 + w_2x_2 + w_3b)$
0	0	0	-1.5	0
0	1	0	-0.5	0
1	0	0	-0.5	0
1	1	1	0.5	1

XNOR

X_1	X_2	$X_1 \text{XNOR } X_2$
0	0	1
0	1	0
1	0	0
1	1	1

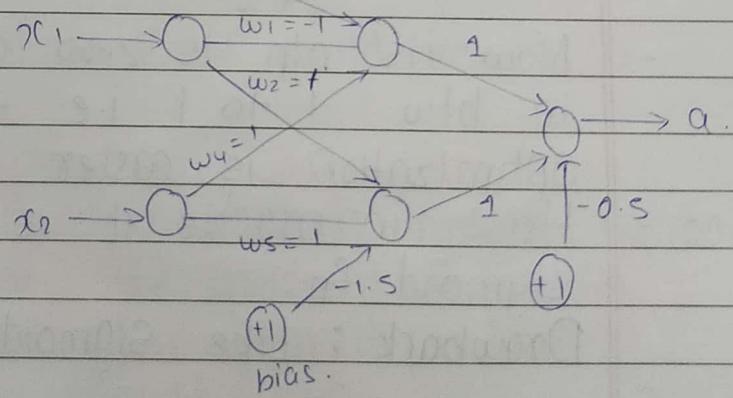


$$XNOR = \overline{X_1 \bar{X}_2 + \bar{X}_1 X_2}$$

$$\begin{aligned} X_1 \text{XNOR } X_2 &= \text{NOT}(X_1 \text{XOR } X_2) = \text{NOT}[(A+B) \cdot (A+B)] \\ &= (\bar{A}+\bar{B})' + (\bar{A}'+\bar{B}')' = (\underline{\bar{A}' \cdot \bar{B}'}) + (\underline{\bar{A} \cdot \bar{B}}) \end{aligned}$$

Inputs $\rightarrow \bar{A} \cdot \bar{B}$ \rightarrow OR \rightarrow XNOR.

$$-0.5 + x_1 + x_2$$



\rightarrow AND

$$+w_2x_2 + w_3b$$

— Activation fn decides whether a neuron should be activated or not by calculating weighted sum & further adding bias with it.

Sigmoid actⁿ fn:

- It is an actⁿ fn of form, $f(x) = \frac{1}{1+e^{-x}}$.
- Its range is b/w 0 & 1.
- It's easy to understand & apply.

Disadvantages.

- Vanishing gradient problem.
- Secondly its o/p is not zero centered So it makes the gradient updates go too far in different directⁿ. $0 < o/p < 1$ & it makes optimization harder.
- Sigmoids saturate & kill gradients.
- Sigmoids have slow convergence.
- Computationally expensive.

Hyperbolic Tangent fn - Tanh.

Its mathematical formula is

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

- Now its o/p is zero centered because its range is b/w -1 to 1 i.e. $-1 < o/p < 1$. Hence optimization is easier in this method.
- Hence in practice it is always preferred over sigmoid fn.

Drawback : Like sigmoid fn only.

ReLU - Rectified Linear Units.

- It has become very popular in the past couple of years.
- It was recently proved that it had 6 times improvement in convergence from Tanh fn.
- It's just $R(x) = \max(0, x)$ i.e if $x < 0, R(x) = 0$ & if $x > 0, R(x) = x$.

- Hence as seeing the mathematical form of this fn, we can see that it is very simple & efficient.
- It avoids & rectifies vanishing gradient problem.
- Almost all deep learning models use ReLU nowadays.
- Drawback:-
 - Its limitation is that it should only be used within hidden layers of a neural nw model.
 - Hence for o/p layers we should use a Softmax fn for a classification problem to compute the probabilities for the classes.
 - Another problem with ReLU is that some gradients can be fragile during training & can die, i.e. it can cause a weight update which will make it never activate on any data point again.
- Simply saying that ReLU could result in dead neurons.

- The Dying ReLU Problem: when inputs approach zero, or are -ve, the gradient of the fn becomes zero, the nw cannot perform backpropagation & cannot learn.
- To fix this problem another modification was introduced called Leaky ReLU to fix the problem of dying Neurons. It introduces a small slope to keep the updates alive.

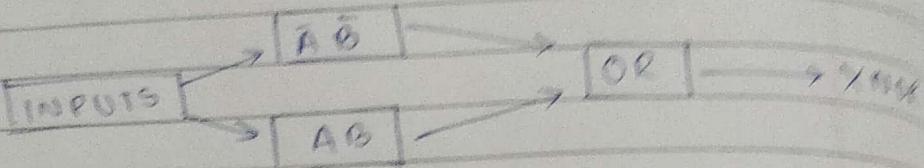
Drawback.

In Leaky ReLU results are not consistent & does not provide consistent predictions for -ve i/p values.

Tue
03/03/20

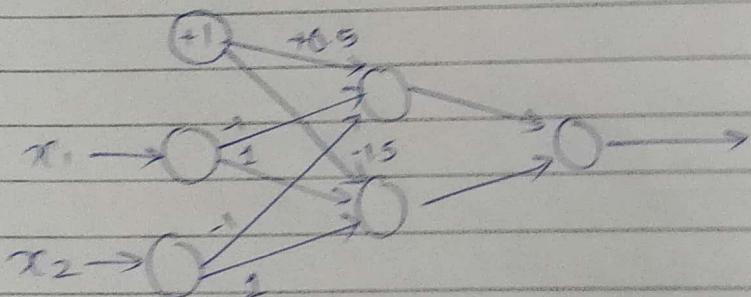
XNOR Gate

$$\begin{aligned}x_1 \text{ XNOR } x_2 &= \text{NOT}(x_1 \text{ XOR } x_2) \\&= \text{NOT}[(A+B) \cdot (A'B')] \\&= (A+B)' + (A'+B') \\&= A'B' + A'B\end{aligned}$$



	AND	XNOR
1	0	1
0 OR	0	0
0	0	0
0	1	1

This can
be 0 or 1
but if 1
it will
be XNOR
... its 0



Conclusion

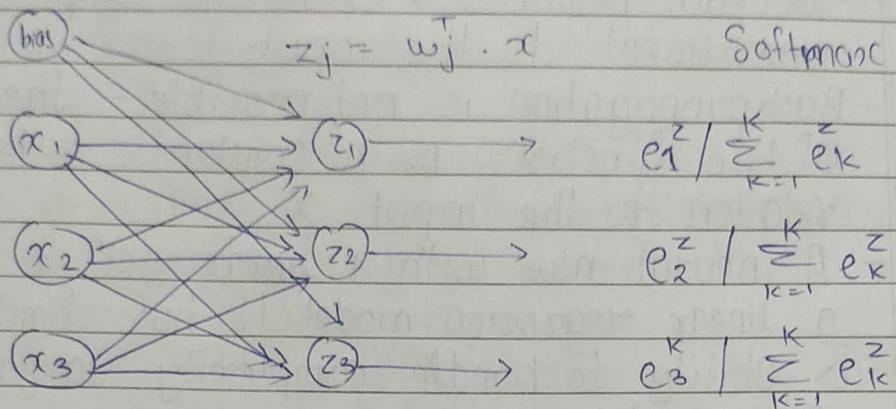
Sigmoid & tanh should not be used due to the vanishing gradient problem which causes a lot of problems degrades the accuracy & performance of deep neural net model.

- Instead we use ReLU which should only be applied to hidden layers
- And if your model suffers from dead neurons during training, we should use leaky ReLU fn.

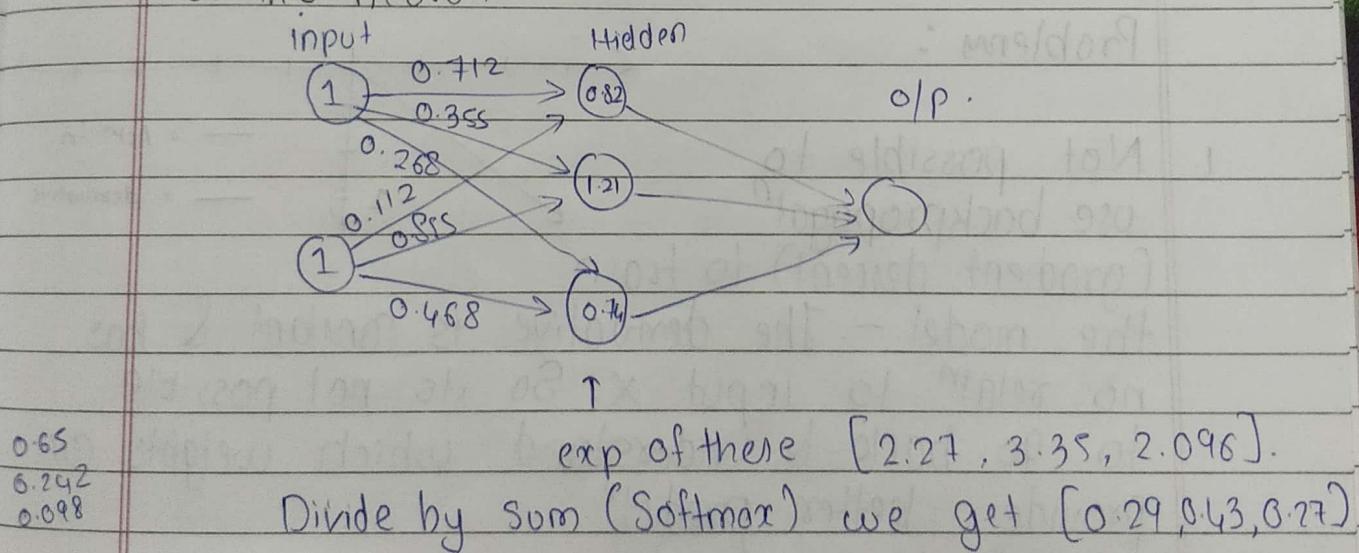
Other activation functions.

1. Softmax : $\frac{e^{a_i}}{\sum_k e^{a_k}}$

Multiclass Classification with NN & Softmax fn



What would be the Softmax activated value of this neuron?



e.g.

LOGITS
score.

$$y = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \Rightarrow \text{Softmax} \Rightarrow \begin{bmatrix} 0.65 \\ 0.242 \\ 0.098 \end{bmatrix} \approx [0.7, 0.2, 0.1]$$

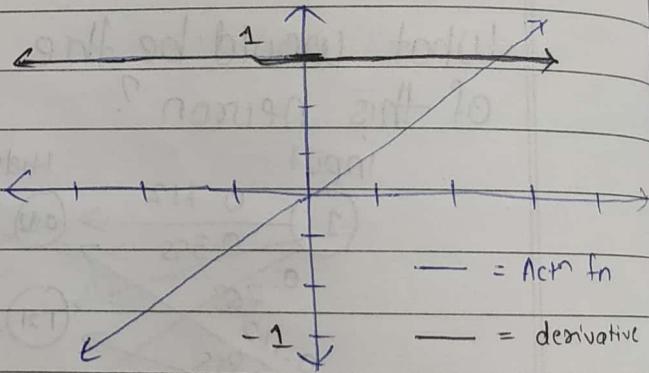
2. Linear Activatⁿ fn.

- Takes the form : $Act^n(x) = cx$.
- Tanh, Sigmoid are not linear

Linear or Identity Actⁿ fn.

- It takes the input, multiplied by the weights for each neuron, & creates an output signal proportional to the input
- Backpropagation is not possible - The derivative of the function is a constant, & has no relation to the input, x .
- A neural nw with a linear actⁿ fn is simply a linear regression model. It has limited power & ability to handle complexity varying parameters of input data.

Problems :



1. Not possible to use backpropagatⁿ

(gradient descent) to train

the model -. The derivative is constant & has no relatiⁿ to input x . So its not possible to go back & understand which weights can provide better predictⁿ.

2. All layers of neural nw collapse into one.

With Linear Actⁿ fn, no matter how many layers in the NN, the last layer will be a linear fn of the first layer. (because linear combinatⁿ of linear fn is still linear fn). Thus NN into just one layer.

∴ Not much used.

3. Non-Linear Actⁿ fn.

- Today's NN models use non-linear Actⁿ fn.
- They allow the model to create complex mappings b/w the n/w's inputs & outputs, which are essential for learning & modelling complex data such as images, video, audio & datasets which are non-linear or have high dimensionality.
- Non-linear fn address the problems of a linear Actⁿ fn.
- They allow backpropagation because they have a derivative fn which is related to the inputs.
- They allow "stacking" of multiple layers of neurons to create a deep neural n/w.
- Multiple hidden layers of neurons are needed to learn complex data sets with high level of accuracy.

Chp 2: Shallow Neural Networks.

Batch Normalization.

One of the important issues in NN is that the training of the network takes a long time for effectively deep networks.

By backpropagation, the NN learn how much error it did & correct themselves, i.e. correct their "weights" & "biases".

By this, NN learn the problem to produce outputs for given inputs.

Back propagation involves computing gradients for each layer & propagating it backward, hence the name.

But during backpropagation of errors to the weights & biases, we'll face a undesired property of Internal Covariate Shift. This makes the net too long to train.

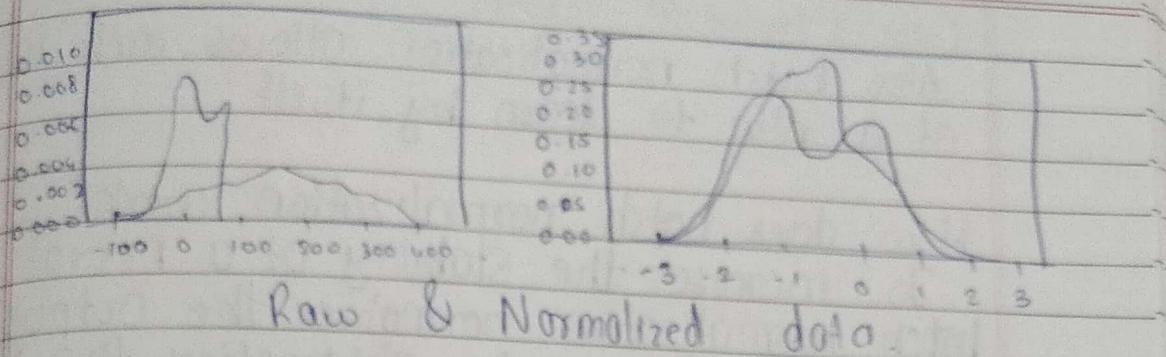
Internal Covariate Shift

During training, each layer is trying to correct itself for error made up during the forward propagation. But every single layer acts separately, trying to correct itself for the error made up.

More specifically, due to changes in weights of previous layers, the distribution of input values for current layer changes, forcing it to learn from new "input distribution".

Normalization.

In a dataset, all the features may not be in some range. It takes a lot of time to train for these kind of datasets.



Raw & Normalized data

We normalize the input layer by adjusting & scaling the features.

e.g. When we have features from 0 to 1 & some from 1 to 1000, we should normalize them to speed up learning.

Same thing also for the values in the hidden layers, that are changing all the time. & get 10 times more improvement in training speed.

Batch normalization reduces the amount by what the hidden unit values shift around.

To know shift around, let's have a deep network on cat detection. We train our data on only black cats' images.

If we now try to apply this network to data with colored cats, it is obvious, we're not going to do well.

The training set & the prediction set are both cats' images but they differ a little bit.

In other words, if an algorithm learned some X to Y mapping.

If the distribution of X changes, then we might need to ~~train~~ retrain the learning algo by trying to align the distribution of X with the distribution of Y .

- Also batch normalization allows each layer of a n/w to learn by itself.

How does batch normalization work?

To increase the stability of a neural n/w, batch normalization normalizes the output of a previous actⁿ layer by subtracting the batch mean & dividing by the batch standard deviation.

After this scale of actⁿ ops, the weights in the next layer may be optimal. SGD under this normalization if it's a way for it to maximize the loss fn.

If so, batch normalization adds two trainable parameters to each layer, so the normalized op is multiplied by a "standard deviation" & add a "mean".

In other words, batch normalization lets SGD do the normalization by changing only these two weights for each actⁿ, instead of losing stability of the n/w by changing all the weights.