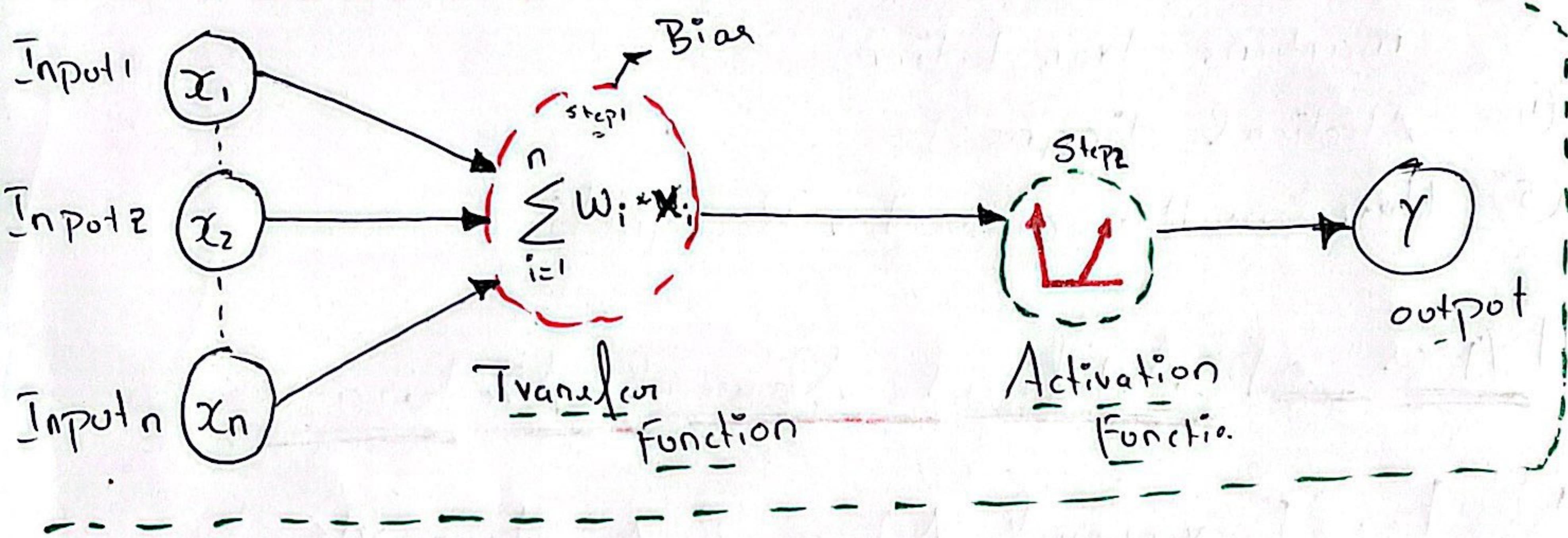


Artificial Neural Network

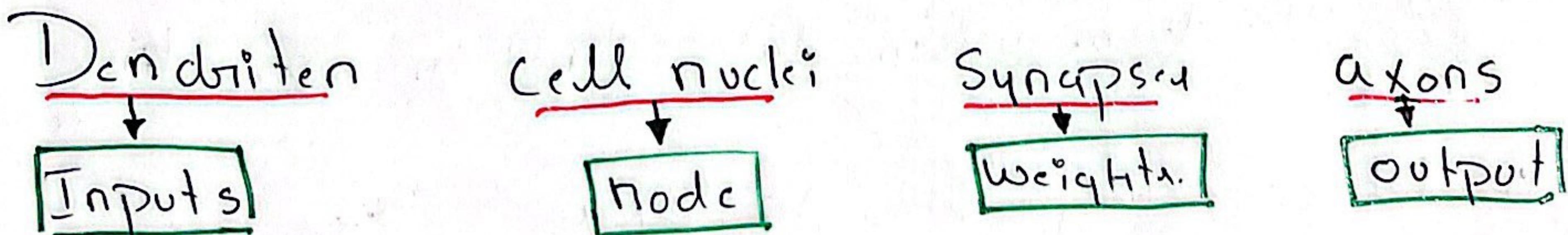
Classifier

Architecture



A computational network grounded on natural neural networks that construct the structure of the human brain is known as Artificial Neural Networks.

From Biological Aspect



* ANN, denrites from biological neural networks represent input, cell. and ANN is a non-linear Statistical model. ①

That demonstrate a complex relationship between input and output in order to uncover a new pattern.

Application

- ① image recognition
- ② speech recognition.
- ③ machine translation
- ④ Medical diagnosis.
- ⑤ Handwritten character Recognition

Types of Artificial Neural Networks.

① Feed Forward Neural Network:

- Information flow in FFNN is single direction.
- No feedback loops.
- These neural networks are commonly employed in supervised learning. for classification.
 - Image recognition.
- We use them when data is not in consecutive order.
- FF are comparable to CNN

② Feed back Neural Network.

- Feedback loops are heart of ANN
- used when data is sequential or time-dependent.
- FB are comparable to RNN.

Pre Processing Steps for ANN

for (historical text, numerical data)

Step 1: Importing the data.

Step 2: Understanding the distribution of the data.

Step 3: Data cleaning (No null value).

Step 4: Encoding (categorical to numerical)

Step 5: Feature Scaling (Normalization)

Step 6: Train - Test - Split.

Pre Processing for Image Classification

Step 1: Image Resizing

The images need to be resized to a standard size to ensure that the model is not affected by the size of the image.

Step 2: Image Normalization

The pixel value of the image need to be normalized to ensure that they are in the same range and do not affect the model.

Step 3: Image Augmentation

To increase the size of the training dataset, images can be augmented by applying transformation such as rotation, flipping and scaling.

Step 4: Image Filtering:

To remove any noise or unwanted features from the images, they can be filtered using techniques such as median filtering & Gaussian filtering.

Step 5 Image Segmentation:

To isolate the objects of interest in the image, they can be segmented using techniques such as thresholding, edge detection, and morphological operation.

Step 6: Image Histogram Equalization:

To adjust the brightness and contrast of the images, they can be equalized using techniques such as Histogram Equalization.

Step 8: Image Convolution:

The image needs to be encoded into numerical vectors converted to a format that is compatible with the model. Such as grayscale, ${}_{\text{our RGB}}^{\text{(H)}}$

Step 8 : Image Encoding

The labels for the images need to be encoded into numerical values to be used by the model.

Step 9 : Image Preprocessing Pipeline

The Preprocessing steps need to be combined into a Pipeline to ensure that they are performed consistently and efficiently.

Optimizers used by ANN algorithm.

Problems and How did they overcome it?

① Gradient Descent

Problem: It's too expensive to calculate the gradient if the size of the data is huge.

- can't determine how far to travel along the gradient for nonconvex function.

↓
To
Overcome

② Stochastic Gradient Descent.

How did they overcome: instead of taking the whole dataset for each iteration we randomly select the batches of data. we take few samples from the dataset

Problem: our path had noise.

- higher no of iterations to reach local minimum
- computation time increases

↓
To
Overcome

③ Mini Batch Gradient Descent

How did it overcome problem?

- Instead of taking single Sample at a time
Batch of Samples were taken to reduce the time.

↓ Problems: as it has noise in Path.
to overcome.

④ Stochastic Gradient Descent with Momentum

SGD with Momentum overcomes this problem by:

- Hence momentum does help in faster convergence of the loss function.

We use Exponential weighted average.
Which used in Time Series.

↓ Problem: Not completely solve Noise problem
But reduces efficiency. if the momentum is increased possibility of reaching GM will be less.

As of now let's keep this optimizers in Best

Because, it reduces the noise, time, computation cost.

SGD with Momentum was all about weights and bias.

①

AdaGrad (Adaptive Gradient Descent)

This optimizer is requiring learning rate.

In Previous Opt we know η (learning rate is constant).

But AdaGrad η will be Adaptive.

Problem: It will decrease fast so it hard to find global minima. And because of Squared gradient denominator. L R becomes extremely small. So model unable to acquire more knowledge.

To overcome this

②

Ada delta and RMS Prop.

~~Both~~ RMS Prop Adaptive to Improve.

AdaGrad: It uses Exponential Moving average which similar to momentum, decrease to speedup.

Ada delta: focus on learning rate component.
delta is difference b/w current and new weight

Adadelta completely removes the use of learning rate parameter completely. By replacing it with D.

③

Adam Optimizer. (Best in Market).

It is the combination of Momentum and RMSProp. Both adaptive learning rate and smoothing is also done together. It can handle spars gradient in noisy Problem.

Activation Function's Problem

and How to Overcome it.

① Sigmoid:

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

used for the output of Binary Probabilistic Function. Gives smooth & slow converging in gradient and clear prediction with $\frac{1}{1+e^{-x}}$.

Problem:

= Prone to Vanishing Gradient problem.

Derivative will not exceed 0.25 at this point it will be constant.

- Not zero centric (always giving positive values).
- computationally expensive (exponentials involved)

Note: Should used this function for Output layers in Binary classification.

② Tanh:

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

- It is 0 centric value (-1 to 1).

Problem : ① Prone to Vanishing Gradient function
② computation cost.

③ ReLU (Rectified Linear Unit).

$$\boxed{\max(0, x)}$$

It solves VGD By -ve value will be 0 and +ve value will be any positive integer.

Problem: -ve input will be inactive neuron.
 $\underset{=}^{(0)}$

④ Leaky Relu: $f(x) = \max(0.01x, x)$.

-ve value will be small value.

It solves dead neuron problem.

⑤ ELU (Exponential Linear Unit).

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^{-x} - 1), & \text{otherwise.} \end{cases}$$

- overcomes dead relu problem; smoother convergence for any negative axis value.

① It Behave like Step function.

Problem: more computational intensive.

⑥ PReLU (Parametric Relu): $f(y_i) = \begin{cases} y_i & \text{if } y_i > 0 \\ a_i y_i & \text{if } y_i \leq 0 \end{cases}$

• It has learning parameter function which fine-tuned the activation based on learning rate. (no constant value).

- ⑦ Soft Max:
$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, 1, 2, \dots, K$$
- used when output layer is multiclass problem.
max probability, one with highest probability value. it will be consider.
- ⑧ Swish (A soft Gradiel function)

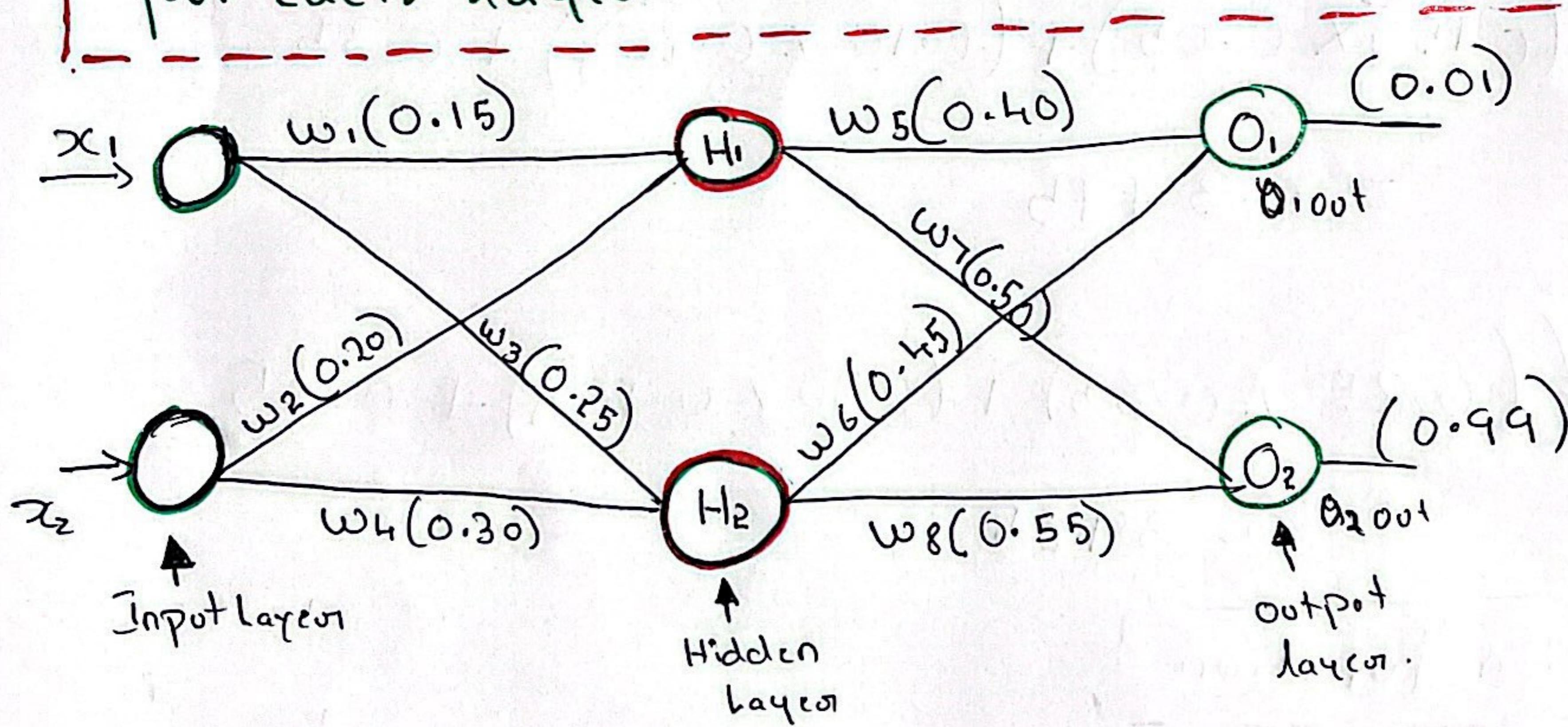
$$y = x^* \text{Sigmoid}(x).$$
- used when there are more than 10 layers,
with LSTM (Long Short Term Memory) Network.
- Problem: computationally expensive.
- ⑨ Maxout:
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$
- It is learnable Activation function.
 Parameter is, over comes w.r.t Generalized ReLU & Leaky ReLU). output will be different activation values.
- ⑩ Softplus:
$$f(b) = \ln(1 + \exp x)$$
- Similar to ReLU & rectified Smooth, can handle SGD.
- Problem: computationally expensive.

ANN

Math Calculation

Math with Algorithm Steps

Step 1: Define the architecture of the network including number of layers, the number of neurons in each layer, Activation function for each layer.



Input: $x_1 = 0.05$, $x_2 = 0.10$ Actual output

$T_1 = 0.01$, $T_2 = 0.99$

Weights taken randomly \rightarrow Bias

$$w_1 = 0.15$$

$$w_2 = 0.20$$

$$w_3 = 0.25$$

$$w_4 = 0.30$$

$$w_5 = 0.40$$

$$w_6 = 0.45$$

$$w_7 = 0.50$$

$$w_8 = 0.55$$

$$b_1 = 0.35$$

$$b_2 = 0.60$$

$$h = 0.6$$

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

$$\text{Activation function: } \frac{1}{1 + e^{-x}}$$

Step 2: Feed forward Propagate the input data through the network, calculating the Activation of each neuron.

Summation function = $\sum_{i=0}^i w_i x_i + b$

$$H_1 = (w_1 x_1 + w_2 x_2) + b_1 \quad H_1 (\text{In})$$

$$H_2 = (w_3 x_1 + w_4 x_2) + b_2 \quad H_2 (\text{In})$$

$$\begin{aligned} H_1 &= ((0.15 \times 0.05) + (0.10 \times 0.20)) + 0.35 \\ &= 0.3775 \end{aligned}$$

$$\begin{aligned} H_2 &= ((0.25 \times 0.05) + (0.30 \times 0.10)) + 0.35 \\ &= 0.3925 \end{aligned}$$

Activation function,

$$\begin{aligned} H_1 \text{out}_{\sigma(H_1)} &= \frac{1}{1 + e^{-H_1}} \\ &= \frac{1}{1 + e^{-0.3775}} \\ &= 0.5932 \end{aligned}$$

$$\begin{aligned} H_2 \text{out}_{\sigma(H_2)} &= \frac{1}{1 + e^{-H_2}} \\ &= \frac{1}{1 + e^{-0.3925}} \\ &= 0.59688 \end{aligned}$$

Calculate O_1 (in and out)

$$O_1(\text{in}) = h_1(\text{out}) \times w_5 + h_2(\text{out}) \times w_6 + b_2$$
$$= 0.593 \times 0.4 + 0.596 \times 0.45 + 0.6$$
$$= \underline{\underline{1.105}}$$

$$O_1(\text{out}) = \frac{1}{1 + e^{-O_1(\text{in})}} = \boxed{\underline{\underline{0.7513}}}$$

Calculate O_2 (in and out)

$$O_2(\text{in}) = h_1(\text{out}) \times w_7 + h_2(\text{out}) \times w_8 + b_2$$
$$= 0.5932 + 0.50 + 0.5968 + 0.55 + 0.6$$

$$O_2(\text{out}) = \underline{\underline{1.22484}}$$

$$= \frac{1}{1 + e^{-O_2(\text{in})}} = \boxed{\underline{\underline{0.7729}}}$$

Step 3: Calculate the error: compare the predicted output to the actual output and calculate the error using a loss function such as Mean Squared Error.

Calculate logs (MSE)

$$\begin{aligned}
 E_{\text{total}} &= \sum \frac{1}{2} (\text{target} - o(p))^2 \\
 &= E_{o_1} + E_{o_2} \\
 &= \frac{1}{2} (0.01 - 0.7513)^2 + \frac{1}{2} (0.99 - 0.7729)^2 \\
 &= \frac{1}{2} (-0.7413)^2 + \frac{1}{2} (0.2171)^2 \\
 &= 0.274 + 0.0235 = 0.29837 \quad (\text{approximate})
 \end{aligned}$$

Step 4

Now to minimize the loss we use.

Backward Propagation Method.

Calculating Back Propagation.

w_5, w_6, w_7 and w_8 . (Output \rightarrow Hidden)

First lets adjust w_5 . \rightarrow learning rate.

$$w_5 = w_5 - h_i \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial o_{\text{out} O_1}} \times \frac{\partial o_{\text{out} O_1}}{\partial \text{net } O_1} \times \frac{\partial \text{net } O_1}{\partial w_5}$$

$$\frac{\partial \text{total}}{\partial w_{01}} = \text{Out}_{01} - \text{target}_{01}$$

$$0.751365 - 0.01$$

$$= 0.7413565$$

$$\frac{\partial \text{out}_{01}}{\partial \text{ncto}_1} = \text{out}_{01} (1 - \text{out}_{01})$$

$$= 0.751365 (1 - 0.751365)$$

$$= 0.186815602$$

$$\frac{\partial \text{ncto}_1}{\partial w_5} = \text{out}_{h1} = 0.59326992$$

Now calculating together

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_5} &= 0.7413565 \times 0.186815602 \times \\ &\quad 0.59326992 \\ &= \underline{0.8216704} \end{aligned}$$

$$w_{5_{\text{new}}} = w_5 - \eta \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$= 0.4 - 0.6 \times 0.8216704$$

$$= \underline{0.350699776}$$

$$(ii) w_6 = w_6 - \eta \left\{ \frac{\partial E_{\text{total}}}{\partial w_6} \right\}$$

$$\frac{\partial E_{\text{total}}}{\partial w_6} = \frac{\partial E_{\text{total}}}{\partial o_{\text{out}0}} * \frac{\partial o_{\text{out}0}}{\partial n_{\text{net}0}} * \frac{\partial n_{\text{net}0}}{\partial w_6}$$

$$w_6 = 0.408666186.$$

$$(iii) w_7 = w_7 - \eta \frac{\partial E_{\text{total}}}{\partial w_7}$$

$$\frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial o_{\text{out}0}} * \frac{\partial o_{\text{out}0}}{\partial n_{\text{net}0}} * \frac{\partial n_{\text{net}0}}{\partial w_7}$$

$$w_7 = 0.511301270$$

Step 5: Update weights of each connection in the networks. This optimization which is calculated in over steps. (w_1, w_2, w_3, w_4) (Hidden \rightarrow input layer)

First lets adjust w_1 .

$$w_{1,\text{new}} = w_1 - \eta \left\{ \frac{\partial E_{\text{total}}}{\partial w_1} \right\}$$

$$\frac{\partial E_{\text{total}}}{\partial w_i} = \left[\frac{\partial E_{\text{total}}}{\partial \text{out}_{H_i}} \right] \times \frac{\partial \text{out}(H_i)}{\partial \text{net}_{H_i}} \times \frac{\partial \text{net}_{H_i}}{\partial w_i}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}(H_i)} = \frac{\partial E_{O_1}}{\partial \text{out}_{H_i}} + \frac{\partial E_{O_2}}{\partial \text{out}_{H_i}}$$

$$\frac{\partial E_{O_1}}{\partial \text{net}_{O_1}} \times \frac{\partial \text{net}_{O_1}}{\partial \text{out}_{H_i}}$$

$$\frac{\partial E_{O_2}}{\partial \text{net}_{O_2}} \times \frac{\partial \text{net}_{O_2}}{\partial \text{out}_{H_i}}$$

$$\frac{\partial E_{O_1}}{\partial \text{out}_{O_1}} \times \frac{\partial \text{out}_{O_1}}{\partial \text{net}_{O_1}}$$

$$\frac{\partial E_{O_2}}{\partial \text{out}_{O_2}} \times \left[\frac{\partial \text{out}_{O_2}}{\partial \text{net}_{O_2}} \right]$$

$$\frac{\partial E_{O_2}}{\partial \text{out}_{O_2}} = (\text{out}_{O_2} - \text{target } O_2)$$

$$= 0.772928465 - 0.99$$

$$= -0.217071535$$

$$\frac{\partial \text{out}_{O_2}}{\partial \text{net}_{O_2}} = \text{out}_{O_2} (1 - \text{out}_{O_2})$$

$$= 0.7729 (1 - 0.7729)$$

$$= 0.175510052.$$

$$\frac{\partial E_{01}}{\partial \text{out}_{01}} = (\text{out}_{01} - \text{target}_{01})$$

$$= (0.7513 - 0.01)$$

$$= \underline{0.7413}$$

$$\frac{\partial \text{out}_{01}}{\partial n_{\text{cto1}}} = \text{out}_{01}(1 - \text{out}_{01})$$

$$= (0.7513(1 - 0.7513))$$

$$= \underline{0.18684831}$$

$$\frac{\partial n_{\text{cto1}}}{\partial \text{out}_{h1}} = \text{on } O_1 \text{ from } h_1 \Rightarrow w_5$$

$$= 0.4.$$

$$\frac{\partial n_{\text{cto2}}}{\partial \text{out}_{h1}} = \text{on } O_2 \text{ from } h_1 \Rightarrow w_7$$

$$= 0.50$$

$$0.13849856 \times 0.4 + (-0.0380982 \times 0.50)$$

$$= 0.055399427 + (-0.019049119)$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}(h)} = \underline{0.036350306}$$

$$\frac{\partial \text{out}_h^i}{\partial w_i} = \text{out}_h^i(1 - \text{out}_h^i) \\ = 0.241300709$$

$$\frac{\partial \text{net}_h^i}{\partial w_i} = \frac{\partial}{\partial w_i} (\underbrace{w_1 x_1 + w_2 x_2 + b_1}_0) \\ = 0.05$$

$$\frac{\partial E_{\text{total}}}{\partial w_i} = 0.036350306 \times 0.241300709 \times 0.05 \\ = 0.00438568$$

$$w_{i,\text{new}} = w_i - \eta \frac{\partial E_{\text{total}}}{\partial w_i} \\ = 0.15 - 0.6 \times 0.00438568 \\ = \underline{\underline{0.1497368592}}$$

Step 6: Repeat Step 2-5 : continuously. feedforward, error calculation, Back Propagation, and weight update steps until the networks converge or until a pre-defined stopping condition is met.

Bias update also works similarly.
see weight updation.

formula

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

If we repeat the step 2-5 we will get
a. $w_2, w_3, b_{\text{bias}}$ values.

$$w_2 = 0.19956143$$

$$w_3 = 0.24975114$$

$$w_b = 0.29950229$$