

Neural Network

* Neural Network :-

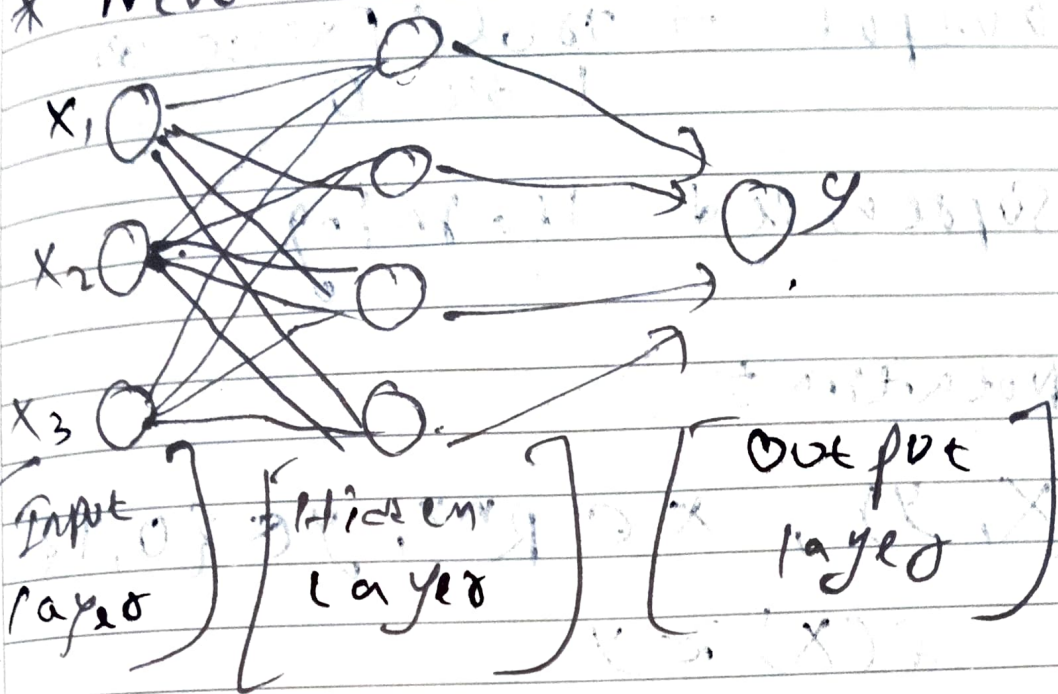


Fig : Standard NN

CNN \rightarrow Image recognition,
Object detection & NLP,
financial time series analysis

RNN \rightarrow Sequential data \rightarrow
NLP, time-series analysis

→ Low Binary Classification

Output → True / False or
1 or 0

Supervised Learning

Notation:

$(X, y), X \in \mathbb{R}^{n \times d}, y \in \{0, 1\}$

$f(X) \approx y$

Logistic Regression

$$z = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$= w^T x$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} \in \{0, 1\}$$

$$\hat{y} = \frac{1}{1 + e^{-(w^T x)}}$$

$x_1, x_2, x_3, \dots, x_n \rightarrow$ input
 $w_1, w_2, \dots, w_n \rightarrow$ weight

$w_0 \rightarrow$ bias (intercept)

* Loss function

* Gradient Descent
update rule

Given x , $\hat{y} = P(y=1 | x)$

$$x \in \mathbb{R}^{n \times 1}$$

parameters: $w \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$

$$* \hat{y} = \sigma(w^T x + b)$$

σ

Loss function:

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

for m no. of data:

Cost function

$$J(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$J(w, b)$$

Gradient Descent Rule:

Repeat Σ

$$w := w - \alpha \frac{dJ(w)}{dw}$$

}

$$w := w - \alpha dw$$

$\alpha \rightarrow$ learning rate

$$dw = X(\hat{y} - y) / m$$

$$w := w - \alpha \frac{dJ(w)}{dw}$$

$$\frac{dJ(w)}{d\hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$dz = \frac{dJ(w)}{d\hat{y}} = \hat{y} - y = db$$

$$\frac{dJ(w)}{dw} = (\hat{y} - y) \cdot X = dw$$

For multiple sample:

$$\frac{dJ(w)}{dw} = \frac{1}{m} X^T (\hat{y} - y)$$

$$\frac{dJ(w, b)}{db} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})$$

* Vectorized ~~the~~ dot product:

$$z = n p \cdot \text{dot}(W, x)$$

$$\left[\begin{array}{l} w \downarrow x = w \cdot x \\ \text{dot product} \\ W^T x \end{array} \right]$$

$$W := W - \alpha \frac{dW}{dw}$$

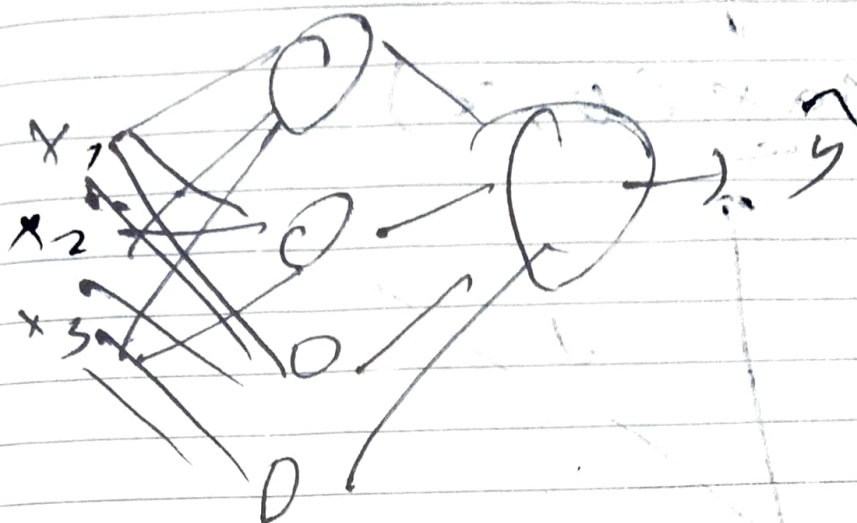
$$b := b - \alpha \frac{db}{db}$$

* Logistic Regression Cost function

$$\text{If } y = 1 : p(y|x) = \hat{y}$$

$$\text{If } y = 0 : p(y|x) = 1 - \hat{y}$$

$$p(y|x) = (\hat{y})^y (1 - \hat{y})^{1-y}$$



$$z^{[1]} = W^{[1]} x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

* Activation function :-

Sigmoid function is an activation function for previous cases.

(1) tanh function works better for hidden layer.

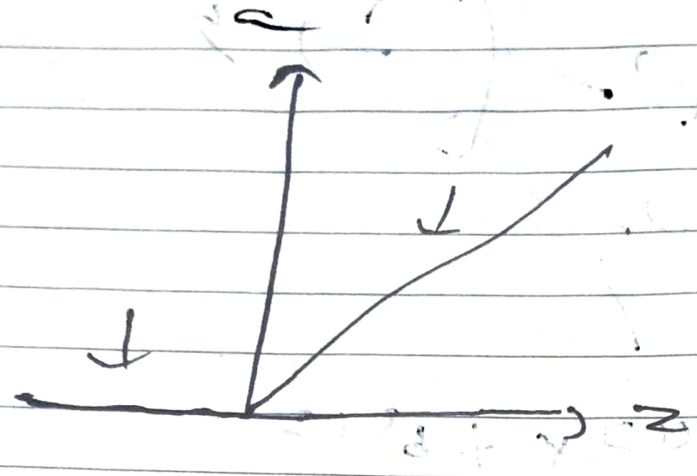
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

For output layer \rightarrow sigmoid is better.

Rectified Linear Unit

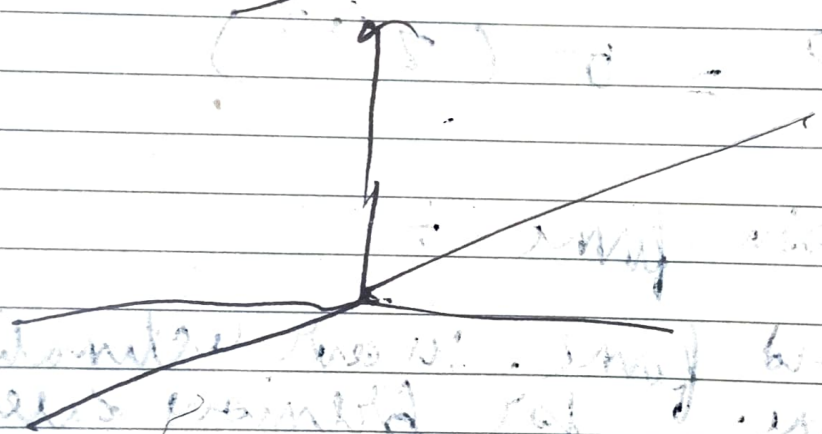
(ii) ReLU func.

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



mostly used

(iii) Leaky ReLU



$$a = \max(0.1z, z)$$

Leaky ReLU is a variation of the ReLU function. It is defined as:

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.1x & \text{if } x \leq 0 \end{cases}$$

$$f(5) = 5$$
$$f(-5) = -0.5$$

It is used in many neural networks to solve the problem of 'dying ReLU'.

* Derivative of activation func.

(i) sigmoid: $g(z) = \frac{1}{1 + e^{-z}}$

$$\frac{dg(z)}{dz} = g(z)(1 - g(z))$$

(ii) Tanh: $g(z) = \tanh(z)$

$$\frac{dg(z)}{dz} = 1 - (g(z))^2$$

(iii) ReLU: $g(z) = \max(0, z)$

$$g'(x) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

Derivative of ReLU: $g'(x) = 1$ if $x > 0$, 0 if $x < 0$

$$g'(x) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

* Gradient descent

$$\theta := \theta - \alpha \frac{dL}{d\theta}$$

(-) call

$\theta \rightarrow$ parameter (weight, bias)

$\alpha \rightarrow$ Learning rate

* Forward Propagation

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

* Back propagation

$$dz^{[2]} = a^{[2]} - y$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} (a^{[1]})^T$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

element wise product

$$dW^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$$

Dimension of weights + bias:

$m \rightarrow$ input

$n \rightarrow$ neuron

weight $\rightarrow (n, m)$

bias $\rightarrow (n,)$

Random initialization:

$w = np.random.randn((n, m))$
* 0.01

$b = np.zeros((n,))$

DNN

$L \rightarrow$ No. of layers

$n^{(l)} \rightarrow$ # Neurons in layer l
units

$a^{(l)} \rightarrow$ Activation in layer l

$$a^{(l)} = g^{(l)}(z^{(l)})$$

$w^{(l)}$
 $b^{(l)}$ \rightarrow weight for $z^{(l)}$

$X \rightarrow a^{(0)} \rightarrow$ input

* Forward Propagation in DNN:

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

* Back Propagation in DNN -

For layer l

Input $a^{[l]}$

Output $a^{[l-1]}$, $\Delta w^{[l]}$, $\Delta b^{[l]}$

$$\Delta z^{[l]} = \Delta a^{[l]} * g^{[l]'}(z^{[l]})$$

$$\Delta w^{[l]} = \Delta z^{[l]} \cdot a^{[l-1]}$$

$$\Delta b^{[l]} = \Delta z^{[l]}$$

$$\Delta a^{[l-1]} = w^{[l]T} \cdot \Delta z^{[l]}$$

$$\Delta z^{[l]} = w^{[l+1]T} \cdot \Delta z^{[l+1]} * g^{[l+1]'}(z^{[l+1]})$$

$$\Delta w^{[l]} = \frac{1}{m} \Delta z^{[l]} \cdot A^{[l-1]T}$$

$$\Delta b^{[l]} = \frac{1}{m} \text{np.sum}(\Delta z^{[l]}, \text{axis}=1, \text{keepdims=True})$$

$$\Delta A^{[l-1]} = w^{[l]T} \Delta z^{[l]}$$