

Pbr:

Assume we are given a dataset of N observations

Input: $x^1, x^2, x^3 \dots x^N \{x^i\}$ $x^i \in \mathbb{R}^n \rightarrow \text{Vector}$

Output: $y^1, y^2, y^3 \dots y^N \{y^i\}$ $y^i \in \mathbb{R} \rightarrow \text{Scalar}$

straight $x \leftarrow x^i = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ ~~y^i~~ Each x^i is a vector
 \leftarrow cursor x

Assume that the data was generated by some function

$$f(x): \mathbb{R}^n \rightarrow \mathbb{R}$$

[Not known. Maybe no closed form. Eg Image, Sound clip]

We try to model f by an approximation \hat{f}

$$\hat{f}(x, w) = w_0 + \sum_{j=1}^{m-1} w_j \phi_j(x)$$

\downarrow Parameters $\underbrace{\hspace{2cm}}$ Basis functions

Linear in Parameters. Basis functions can be non-linear.

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix} \in \mathbb{R}^m \rightarrow \text{Parameter vector}$$

$w_0 \rightarrow$ Bias parameter

$$\phi_0(x) = 1$$

$$\Phi(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \\ \phi_m(x) \end{bmatrix} \rightarrow \text{Feature vector} \in \mathbb{R}^m$$

Aim: Find w s.t. $\hat{f}(x, w) = \Phi^T(x) \cdot w$
 Example: $\hat{f}(x, w)$ is close to f for all data points.

1. Linear Basis: $x \in \mathbb{R}^n$ $m = n+1$

$$\phi_j(x) = x_j \quad j \in [0 \dots n] \quad \phi_0(x) = 1$$

Let $x \in \mathbb{R}^2 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$\phi_0(x) = 1, \phi_1(x) = x_1, \phi_2(x) = x_2$

$\hat{f}(x, w) = w_0 + w_1 x_1 + w_2 x_2$

2. Polynomial Basis: $x \in \mathbb{R}$

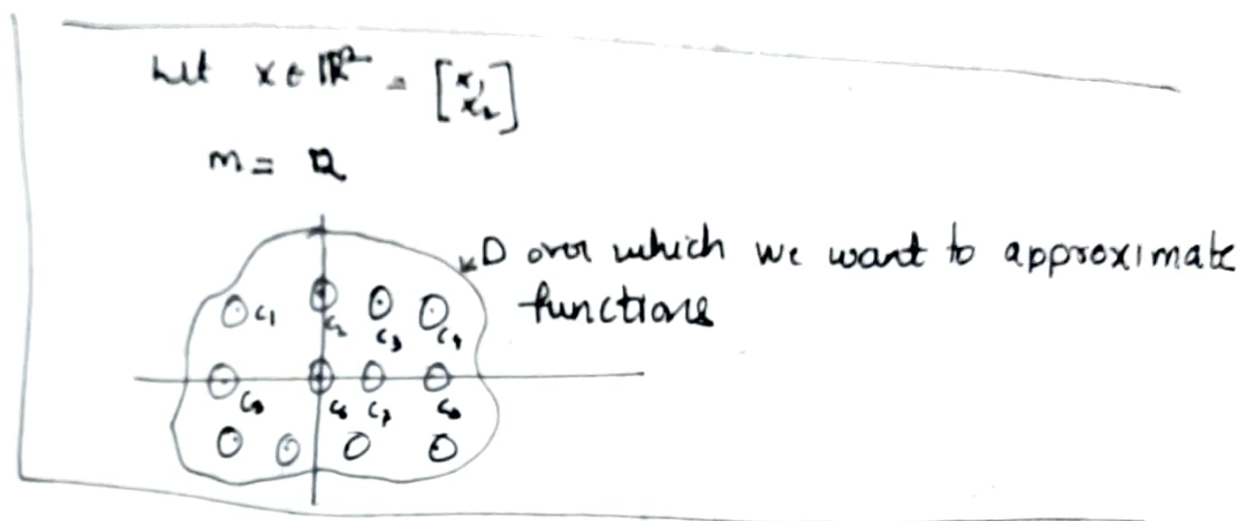
$$\phi_j(x) = x^j$$

Let $m = 2$

$\phi_0(x) = 1, \phi_1(x) = x, \phi_2(x) = x^2$

$\hat{f}(x, w) = w_0 + w_1 x + w_2 x^2$

3. Radial Basis functions: $x \in \mathbb{R}^n$, $\phi_j(x) = \exp\left(-\frac{\|x - c_j\|_2^2}{\sigma_j^2}\right)$ ↔
 Centres: $c_j \in \mathbb{R}^n$



Aim: Find w s.t. \hat{f} and f are close for each data point.

Sum of squares Error: (Sample by sample)

$$E_D(w) = \frac{1}{N} \sum_{i=1}^N \{ \overset{\text{Actual data}}{y_i} - \underset{\text{Data predicted by } \hat{f}}{\hat{y}_i} \}^2$$

$$= \frac{1}{N} \sum_{i=1}^N \{ y_i - \phi^T(x_i) \cdot w \}$$

Select w s.t. $E(w)$ minimized

Gradient descent

Calculate gradient of $E(w)$.

Vectors: Regression Matrix

$$\Phi = \begin{bmatrix} \phi^T(x^1) \\ \phi^T(x^2) \\ \vdots \\ \phi^T(x^N) \end{bmatrix} \in \mathbb{R}^{N \times m} \quad Y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^N \end{bmatrix} \in \mathbb{R}^N$$

$$\hat{y} = \begin{bmatrix} \hat{y}^1 \\ \vdots \\ \hat{y}^N \end{bmatrix} = \begin{bmatrix} \phi^T(x^1) \cdot w \\ \phi^T(x^2) \cdot w \\ \vdots \\ \phi^T(x^N) \cdot w \end{bmatrix} = \begin{bmatrix} \phi^T(x^1) \\ \phi^T(x^2) \\ \vdots \\ \phi^T(x^N) \end{bmatrix} w = \Phi w \in \mathbb{R}^N$$

$$E(w) = \frac{1}{N} (y^i - \hat{y}^i)^2 = \frac{1}{N} \|y - \hat{y}\|_2^2$$

$$= \frac{1}{N} (y - \hat{y})^T (y - \hat{y}) \quad \text{where } \hat{y} = \Phi w$$

$$= \frac{1}{N} (y^T y - y^T \Phi w - w^T \Phi^T y + w^T \Phi^T \Phi w)$$

$$\mathbb{R}^m \ni \nabla E(w) = \frac{1}{N} (-\Phi^T y - \Phi^T y + 2 \Phi^T \Phi w)$$

since $w \in \mathbb{R}^m$

$$= \frac{2}{N} (-\Phi^T y + \Phi^T \Phi w)$$

Gradient is a function of w

1] Gradient descent:

Constant step size: η
 (learning rate)

Algo: Given $\{x^i\}$ $\{y^i\} \rightarrow$ Dataset

Compute Φ , and Y

$W \leftarrow \text{rand}(m) \rightarrow$ Initialize $W \in \mathbb{R}^m$ randomly

for $i=1$: max-steps

$$\nabla E = \frac{2}{N} (-\Phi^T Y + \Phi^T \Phi W)$$

$$W = W - \eta \nabla E$$

end

Problems:

1. Requires all the data to calculate gradient.
2. What if data is so large that it does not fit in memory

2] Stochastic gradient descent:

Idea: 1. Do not ~~take~~ use whole data

2. Use one sample at a time

3. Go through the data sequentially.

$$L(w) = \frac{1}{N} \sum (y^i - \phi_v(x^i)^T w)^2$$

$$\nabla L = \frac{2}{N} \left(-\cancel{\Phi^T} \cdot \cancel{Y} - \cancel{\Phi_v[i]^T} \cancel{\Phi[i]} w \right)$$

Algo: Given $\frac{2}{N} (y^i - \phi_v(x^i)^T w) \cdot \phi_v(x^i)$
 $\{x^i\}$ and $\{y^i\} \rightarrow$ Data set

~~Compute~~ $w = \text{rand}(m)$ number of
epoch
 for $i=1:\text{max_steps}$

One gradient step on the average. \downarrow epoch

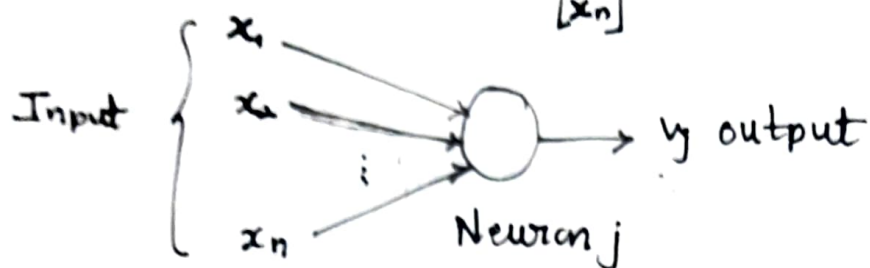
$\left\{ \begin{array}{l} \text{for } j=1:N \\ \nabla L = \frac{2}{N} (y^i - \phi_v(x^i)^T w) \cdot \phi_v(x^i) \\ \text{end} \end{array} \right.$
 end

3] Batch gradient descent: In between. $\lceil \bar{N} \rceil$ sized batches of data

$$L(w) = \frac{1}{N} \sum_{j=1}^{\bar{N}} (y^i - \phi_v(x^i)^T w)$$

Neural Networks

Neuron: Input: $x \in \mathbb{R}^n = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$ Output $y_j \in \mathbb{R}$



$$h_j = \sum_{i=1}^n w_{ji} x_i + b_j \rightarrow \text{bias term}$$

$$y_j = f(h_j) \quad f: \mathbb{R} \rightarrow \mathbb{R} \text{ Activation function}$$

$$\text{Let } W_j = \begin{bmatrix} w_{j1} \\ w_{j2} \\ \vdots \\ w_{jn} \end{bmatrix} \in \mathbb{R}^n$$

$$\therefore h_j = w_j^T x + b_j$$

$$y_j = f(h_j) = f(w_j^T x + b_j)$$

Examples of activation functions

1. Sigmoid (Logistic): $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

2. Hyperbolic tangent: (tanh): $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

$$\tanh'(x) = 1 - \tanh^2(x)$$

3. ReLU: $f(x) = x \quad \text{if } x \geq 0$
 $= 0 \quad \text{if } x < 0$

$$f'(x) = 1 \quad \text{if } x > 0$$

$$= 0 \quad \text{if } x < 0$$

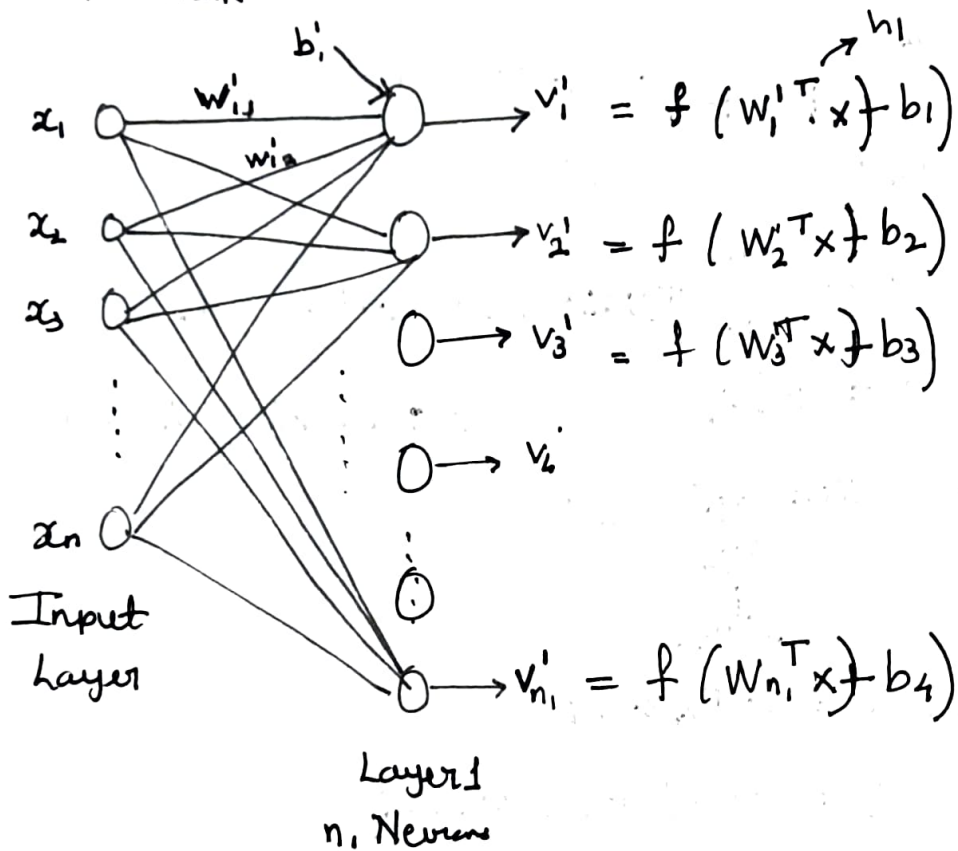
$$f'(x) = 0 \quad \text{at } x = 0$$

↓
Convention.

Neural Network: 1. Network of neurons.

2. Stack neurons on top of one another to form neural network

Input: $x \in \mathbb{R}^n$



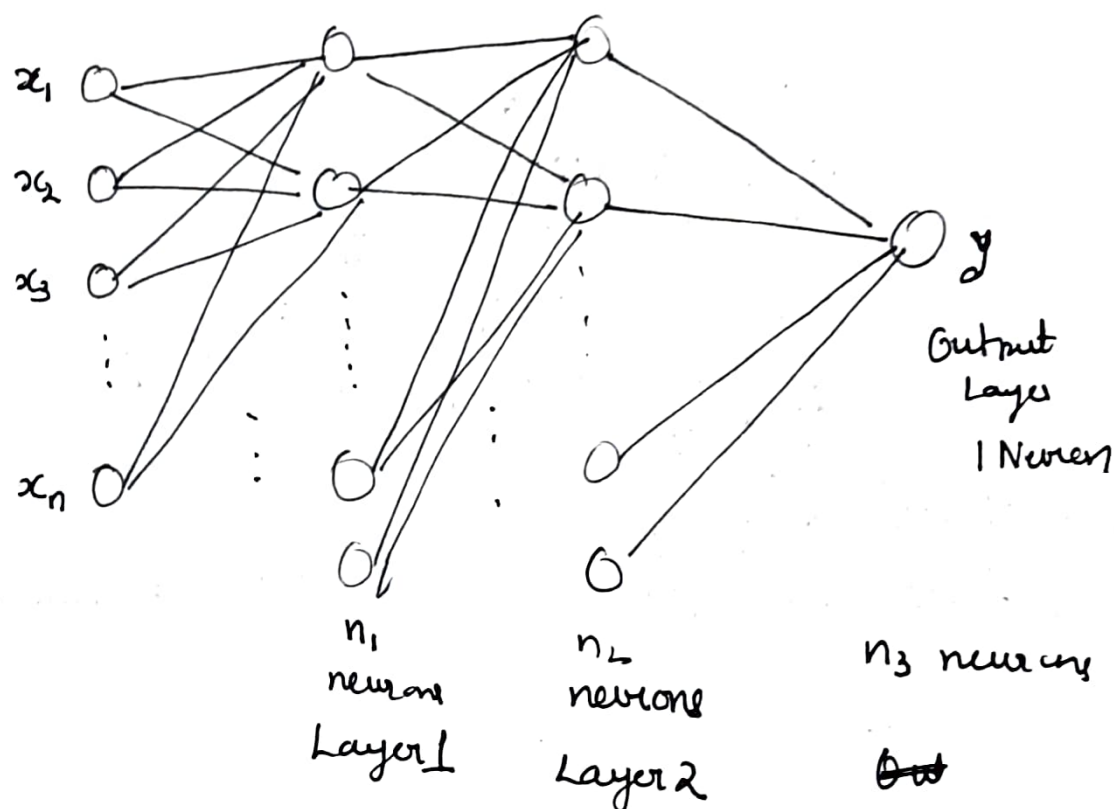
$$h_{0,1} = \begin{bmatrix} h_1' \\ h_2' \\ \vdots \\ h_{n_1}' \end{bmatrix} = \begin{bmatrix} W_1'^T x + b_1' \\ W_2'^T x + b_2' \\ \vdots \\ W_{n_1}'^T x + b_{n_1}' \end{bmatrix} = \begin{bmatrix} W_1'^T \\ \vdots \\ W_{n_1}'^T \end{bmatrix} x + \begin{bmatrix} b_1' \\ b_2' \\ \vdots \\ b_{n_1}' \end{bmatrix}$$

$$h_1 = W_1 x + b_1$$

$$V_1 = \begin{bmatrix} f(h_1') \\ f(h_2') \\ \vdots \\ f(h_{n_1}') \end{bmatrix} = f \left(\begin{bmatrix} h_1' \\ \vdots \\ h_{n_1}' \end{bmatrix} \right) = f(h_1) \rightarrow \text{elementwise application.}$$

$$V_1 = f(W_1 x + b_1)$$

Add more layers



$$v_1 = f(w_1 x + b_1) \quad v_2 = f(w_2 v_1 + b_2) \quad v_3 = y = f(w_3 v_2 + b_3)$$

$f = \text{Identity}$

$$y = w_3 v_2 + b_3$$

NN is a function:

$$y = w_3 f(w_2 v_1 + b_1)$$

Parameters:

$$y = w_3 f(w_2 f(w_1 x + b_1))$$

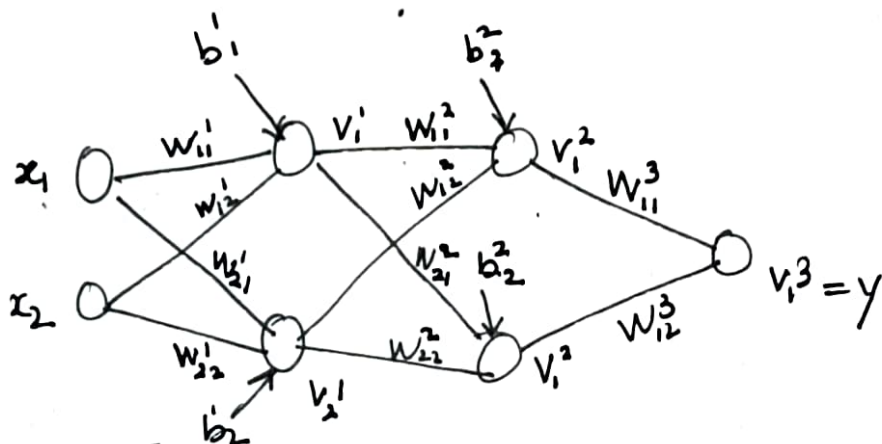
$$P = [w_1, b_1, w_2, b_2, w_3, b_3]$$

$y = \tilde{f}(x, P) \rightarrow$ Compare to linearly parametrized networks.

\rightarrow Parameters occur nonlinearly in output

Makes computing gradients a challenge.

eg] A simple NN with $x \in \mathbb{R}^2$, $n_1=2$, $n_2=2$, $n_3=1$



$$\text{Step 1]} \quad h_1 = \begin{bmatrix} h_1^1 \\ h_2^1 \end{bmatrix} = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^1 \\ b_2^1 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} h_1^1 \\ h_2^1 \end{bmatrix}} \right\} \begin{matrix} 6 \text{ params} = n \times n_1 + n_1 \\ v_1 = \begin{bmatrix} v_1^1 \\ v_2^1 \end{bmatrix} \end{matrix}$$

$$\text{Step 2]} \quad h_2 = \begin{bmatrix} h_1^2 \\ h_2^2 \end{bmatrix} = \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \end{bmatrix} \begin{bmatrix} v_1^1 \\ v_2^1 \end{bmatrix} + \begin{bmatrix} b_1^2 \\ b_2^2 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} h_1^2 \\ h_2^2 \end{bmatrix}} \right\} \begin{matrix} 6 \text{ params} \\ n_1 \times n_2 + n_2 \end{matrix}$$

$$\text{Step 3]} \quad y = v_3^3 = \begin{bmatrix} w_{11}^3 & w_{12}^3 \end{bmatrix} \begin{bmatrix} v_1^2 \\ v_2^2 \end{bmatrix} + \begin{bmatrix} b_1^3 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} v_1^2 \\ v_2^2 \end{bmatrix}} \right\} \begin{matrix} 6 \text{ params} \\ n_2 \times n_3 + n_3 \end{matrix}$$

$$y = NN(x)$$

$$E(P) = \frac{1}{N} \sum_{i=1}^N \left\{ NN(x^i) - y^i \right\}^2$$

$$L^1 = \frac{1}{N} \sum$$

$$W_1 = W_1 + \eta \nabla E_i(P) W_1 \rightarrow \text{How to compute these gradients.}$$

$$W_2 = W_2 + \eta \nabla E_i(P) W_2$$

⋮

Backpropagation: Computing derivatives Algorithmically

①

$W = L(z)$ be a function of z

$z = f(x, y)$ be a function of x and y

$$L(z) = L(f(x, y))$$

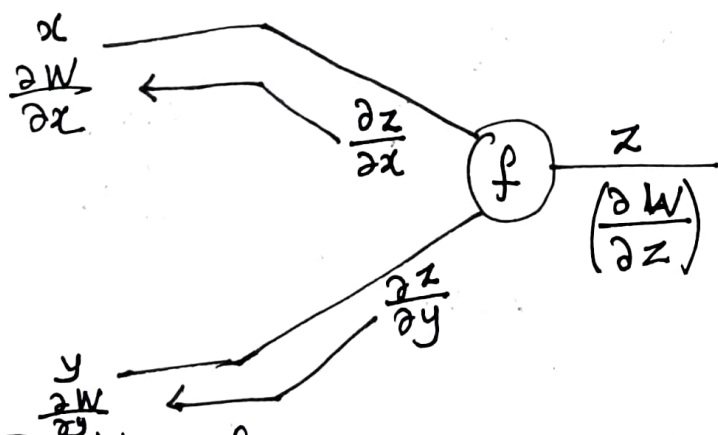
L is indirectly a function of x and y

$$\therefore \frac{\partial W}{\partial x} = \frac{\partial W}{\partial z} \cdot \frac{\partial z}{\partial x}$$

$$\frac{\partial W}{\partial y} = \frac{\partial W}{\partial z} \cdot \frac{\partial z}{\partial y}$$

} Chain Rule.

Computational graph of f :



Final
 $L \rightarrow$ Output function

1. Functions form nodes of computational graph.
2. x and y are input variables of node f .
3. z is the output variable of node f .
4. At each variable we also store derivative of output function w.r.t that variable.

5. At each node compute local derivatives ~~to~~
 (i.e. Derivative of output variable w.r.t input variable)

and place it on the corresponding input branch

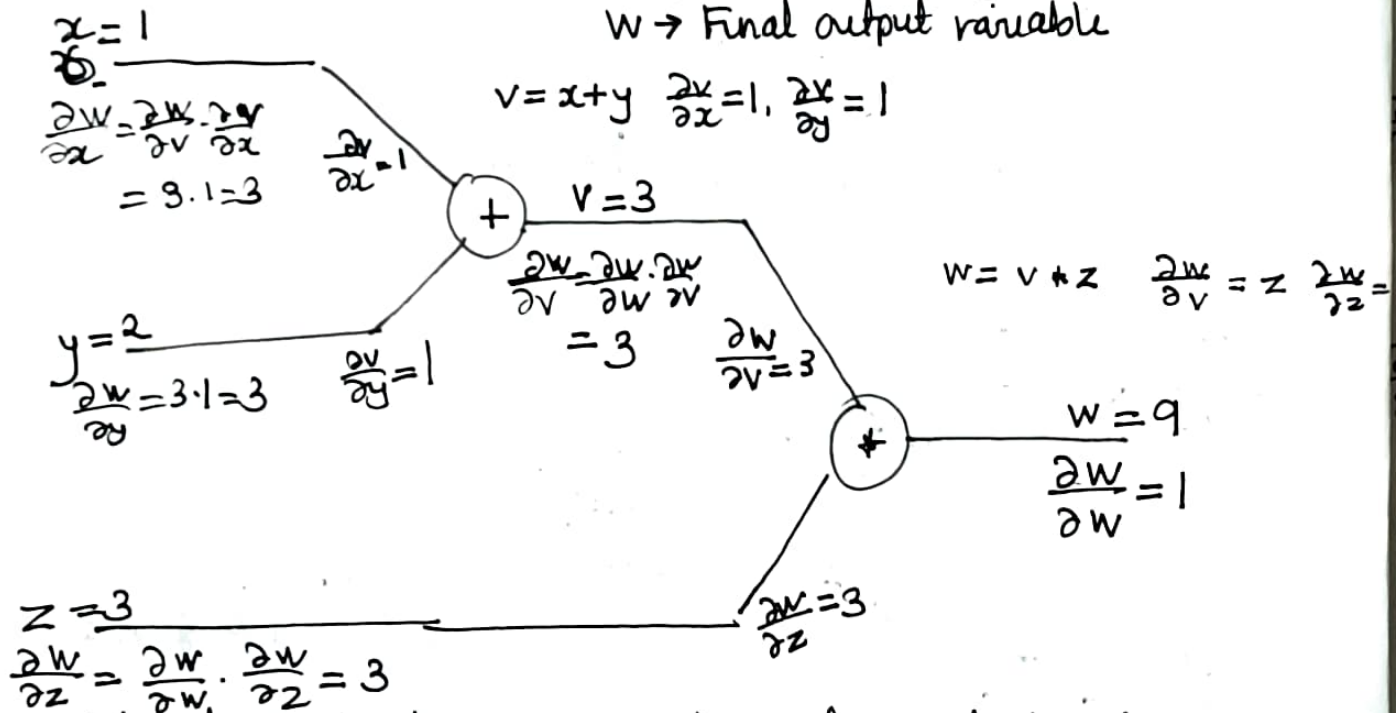
6. Follow the path from output variable to the input variable to get the corresponding input.

$$\frac{\partial W}{\partial x} = \frac{\partial W}{\partial z} \cdot \frac{\partial z}{\partial x}$$

$$\frac{\partial W}{\partial y} = \frac{\partial W}{\partial z} \cdot \frac{\partial z}{\partial y}$$

Ex. $W = (x+y) \cdot z$ We want $\frac{\partial W}{\partial x}, \frac{\partial W}{\partial y}, \frac{\partial W}{\partial z}$

$W \rightarrow$ Final output variable



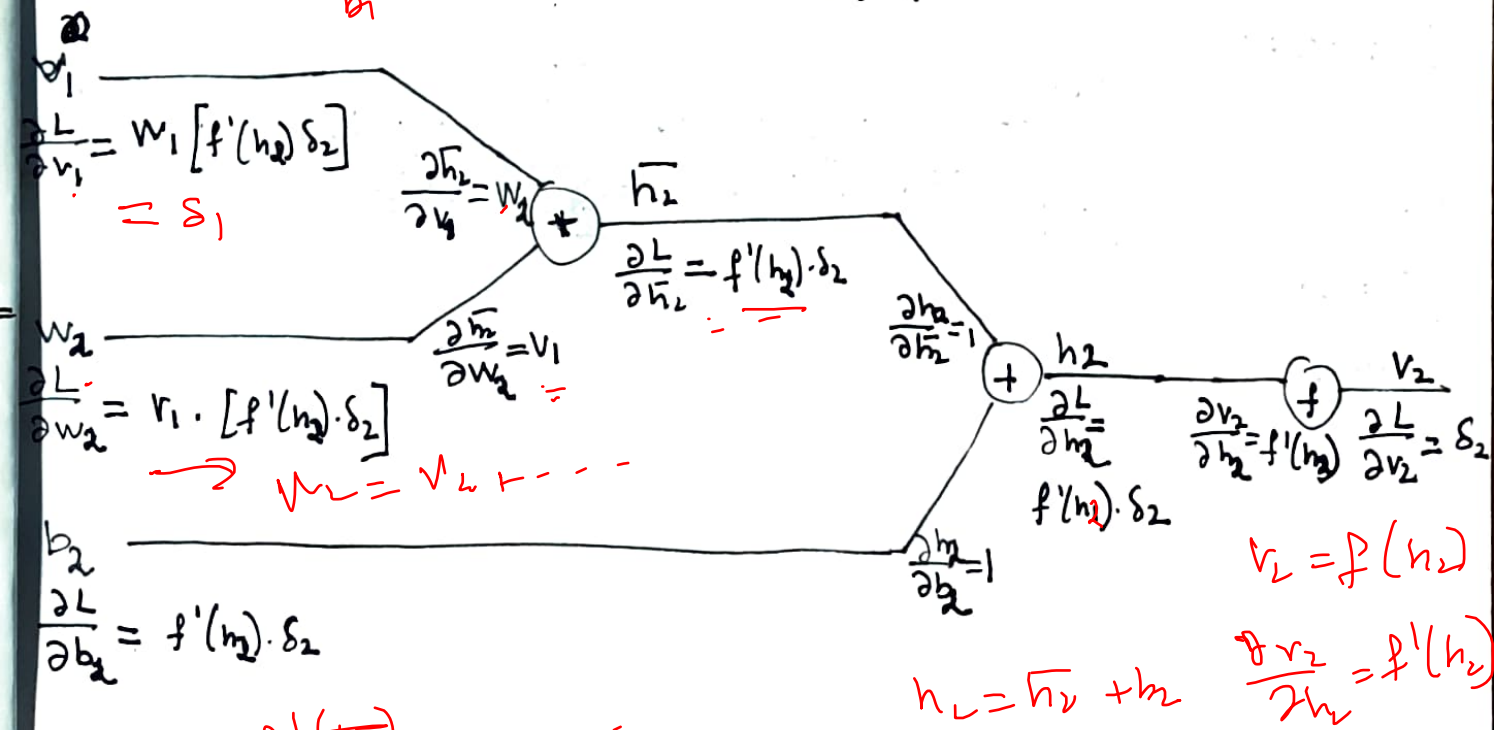
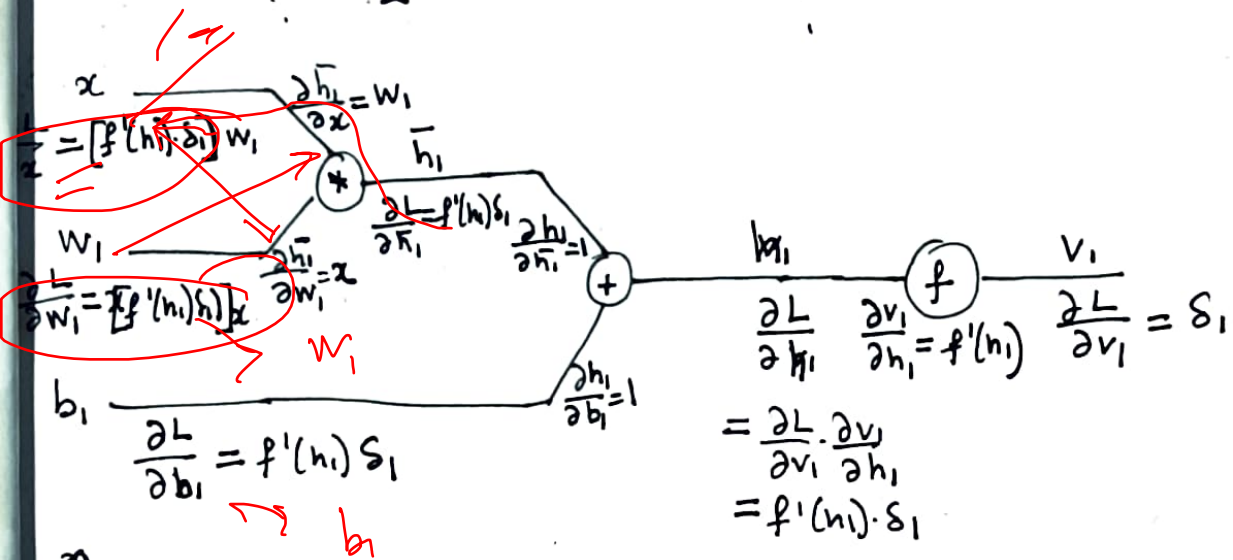
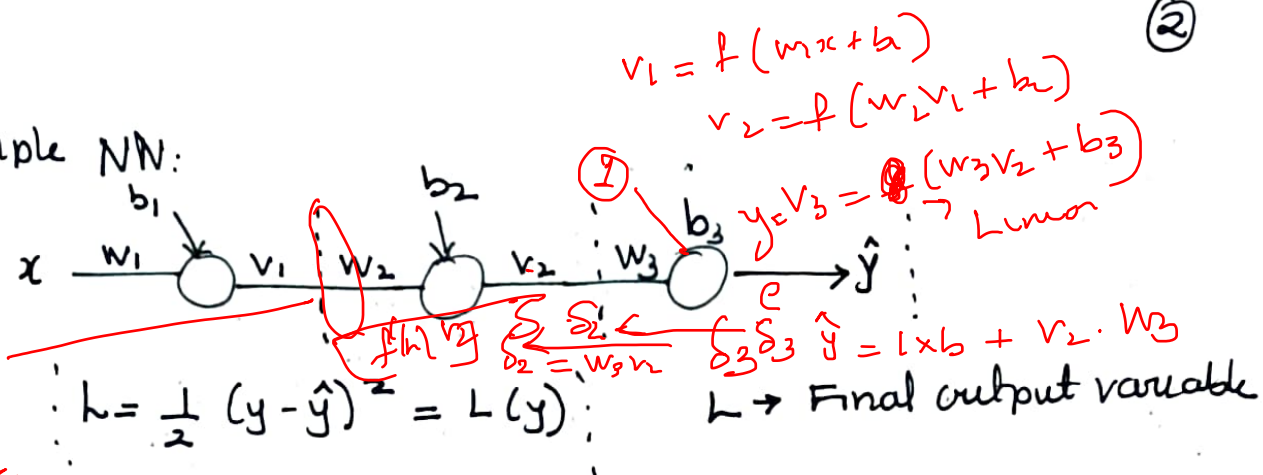
Step 1: Evaluate comp. graph in forward direction

Step 2: ~~Compute local derivative at all nodes~~ Start backprop
 a) Start from end of C.G.

b) ~~Move backwards~~ From end compute derivative at each variable and local derivatives at each node

②

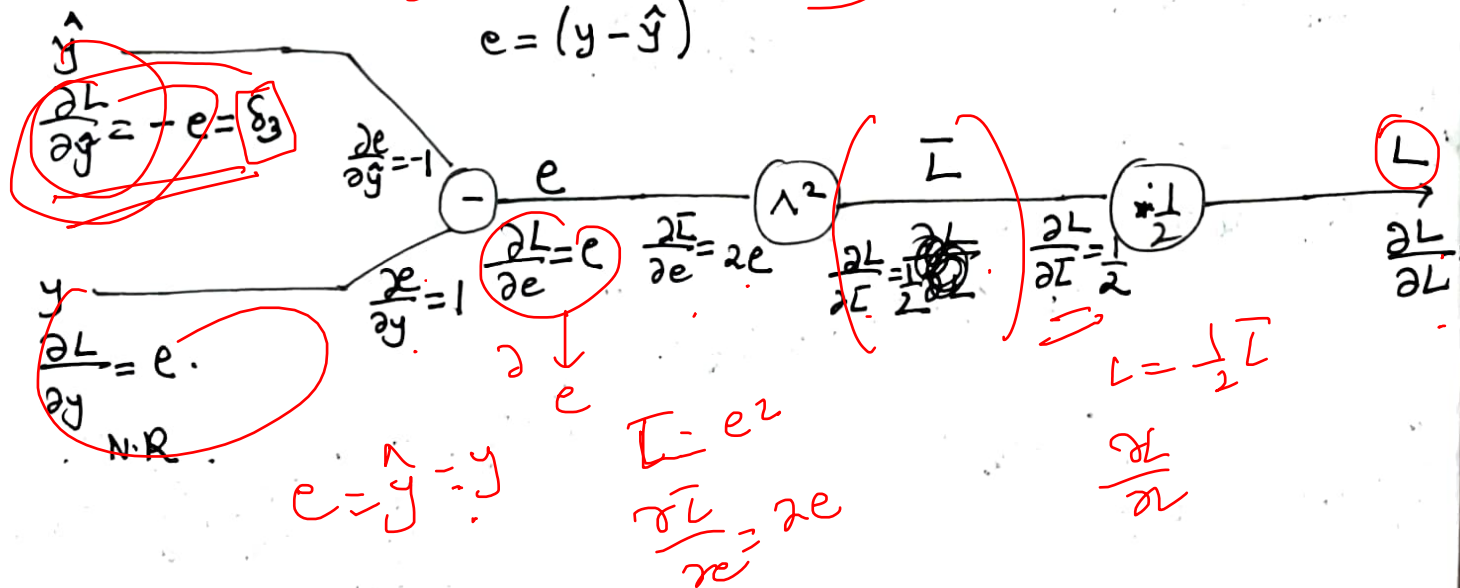
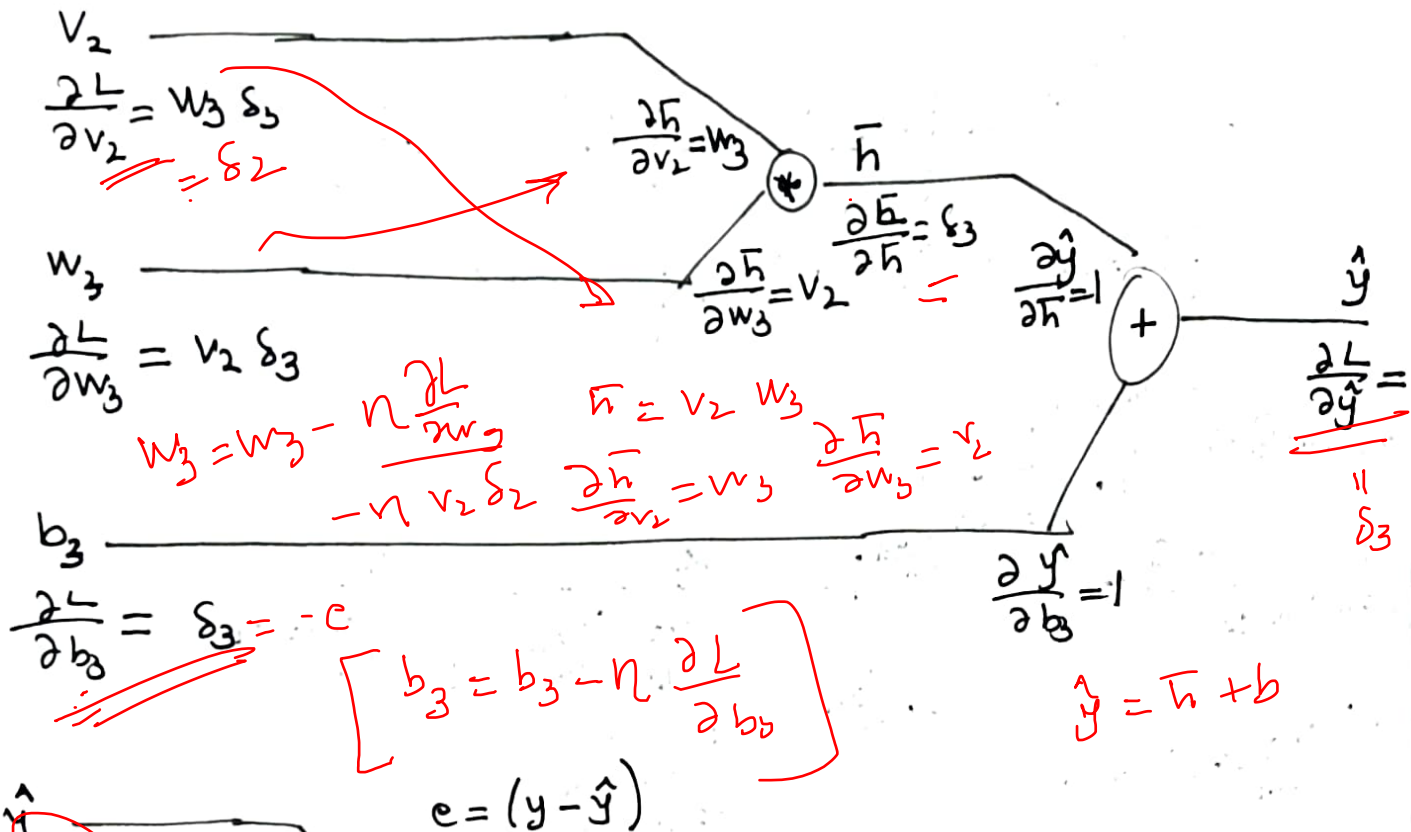
2] Simple NN:



~~$b_2 = f'(h_2)$~~
 $b_2 = b_2 - \eta f'(h_2) \delta_2$

$h_2 = v_2 + b_2$
 $\frac{\partial v_2}{\partial h_2} = f'(h_2)$

$$\hat{y} = \frac{1}{n} w_3 v_2 + b_3$$



go: Compute Gradient

1. Forward pass

$$v_1 = f(w_1 x + b_1)$$

$$v_2 = f(w_2 v_1 + b_2)$$

$$v_3 = w_3 v_2 + b_3 = \hat{y}$$

(Matrix Eq)

$$v_1 = f(w_1 x + b_1)$$

$$v_2 = f(w_2 v_1 + b_2)$$

$$v_3 = w_3 v_2 + b_3 = \hat{y}$$

2. $\delta_3 = -(y_i - \hat{y}_i) = -e$

$\delta_3 = -(y_i - \hat{y}_i) = -e$

$$\frac{\partial L}{\partial w_3} = [\delta_3] v_2$$

$$\frac{\partial L}{\partial b_3} = \delta_3$$

$n_3 = 1$ $\leftarrow f(x, w)$

no of rows

$$\mathbb{R}^{n_3 \times n_2}: \frac{\partial L}{\partial w_3} = \delta_3 \cdot v_2^T$$

$1 \times 1 \times n_2 = 1 \times n_2$
 $n_3 \times n_2 = n_3 \times n_2$
 $w_3 = w_3 - \eta \frac{\partial L}{\partial w_3}$

$$\mathbb{R}^{n_3}: \frac{\partial L}{\partial b_3} = \delta_3$$

$$\frac{\partial L}{\partial v_2} = \delta_2 = w_3 \delta_3$$

$$\mathbb{R}^{n_2}: \frac{\partial L}{\partial v_2} = w_3^T \cdot \delta_3 = \delta_2$$

$$\frac{\partial L}{\partial w_2} = [f_2'(h_2) \cdot \delta_2] \cdot v_1$$

$h_2 = w_2 v_1 + b_2$
 $\delta_2 = \delta_2$
 v_1 present

$$\mathbb{R}^{n_2 \times n_1}: \frac{\partial L}{\partial w_2} = [f_2'(h_2) \cdot \delta_2] \cdot v_1^T$$

$$\frac{\partial L}{\partial b_2} = [f_2'(h_2) \cdot \delta_2]$$

$$\mathbb{R}^{n_2}: \frac{\partial L}{\partial b_2} = [f_2'(h_2) \cdot \delta_2]$$

$$\frac{\partial L}{\partial v_1} = \delta_1 = w_2 [\delta_2 \cdot f_2'(h_2)]$$

$$\mathbb{R}^{n_1}: \frac{\partial L}{\partial v_1} = w_2^T [\delta_2 \cdot f_2'(h_2)]$$

$1 \times 1 \times n_2$

$$\frac{\partial L}{\partial w_1} = [f_1'(h_1) \cdot \delta_1] \cdot x$$

$$\mathbb{R}^{n_1 \times n}: \frac{\partial L}{\partial w_1} = [\delta_1 f_1'(h_1)] \cdot x^T$$

$$\frac{\partial L}{\partial b_1} = [\delta_1 f_1'(h_1)]$$

$$\mathbb{R}^{n_1}: \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial b_1} [\delta_1 f_1'(h_1)]$$

