

Q.1.) $f \rightarrow k=3^{\text{nd}}$ Fully connected MLP

g is a generic 1-Lipschitz activation

x is a $m \times 1$ column vector

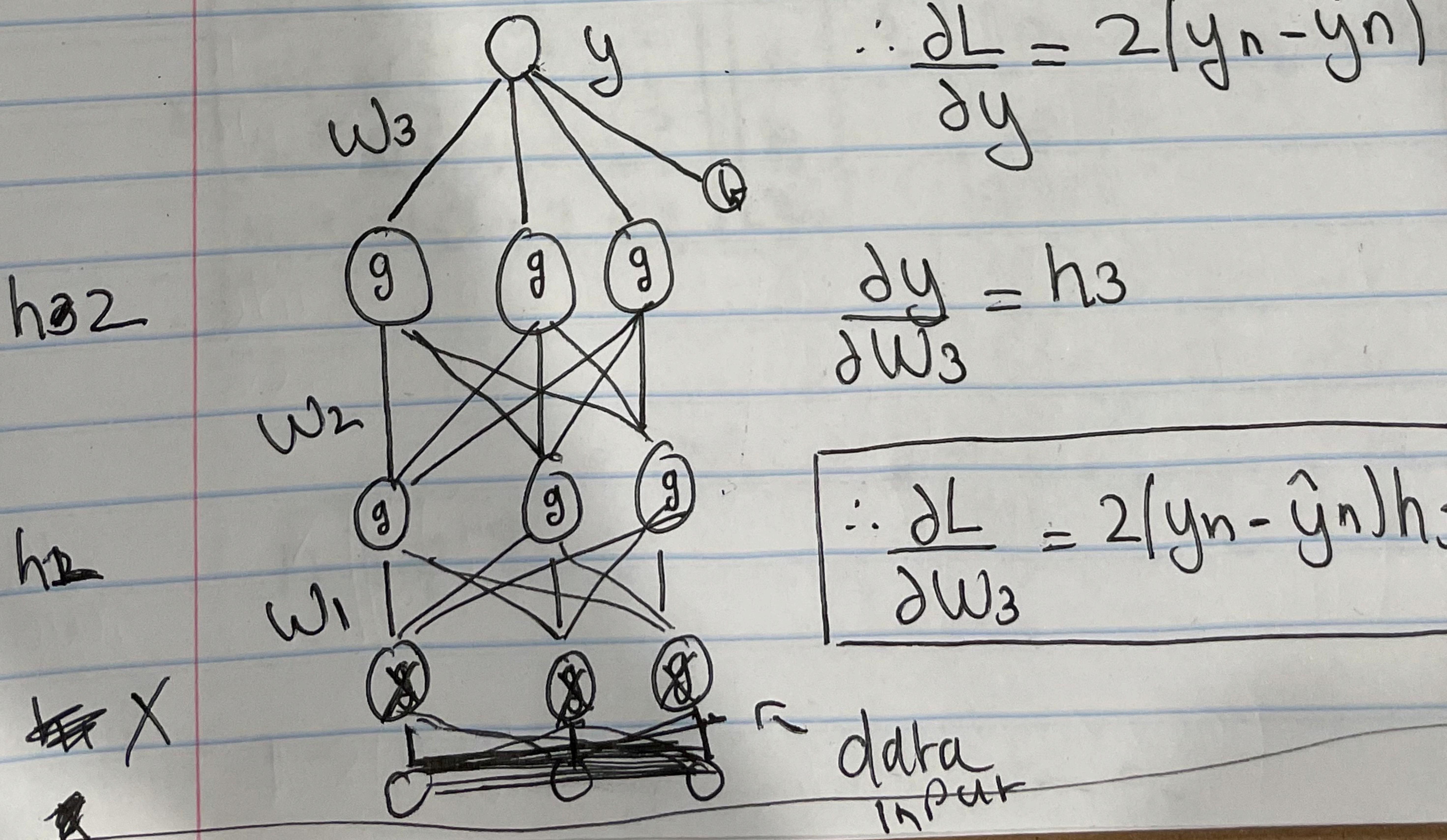
$$h_1 = g(W_1^T X + b_1)$$

$$h_2 = g(W_2^T X + b_2)$$

$$\hat{y} = W_3^T h_2 + b_3 \quad \therefore \hat{y} = W_3^T X + b_3$$

$$L(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

Find: $\frac{\partial L}{\partial W_3}, \frac{\partial L}{\partial W_2}, \frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial b_3}, \frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial b_1}$



$$\frac{d}{dx} (a-x)^2 = 2(x-a)$$

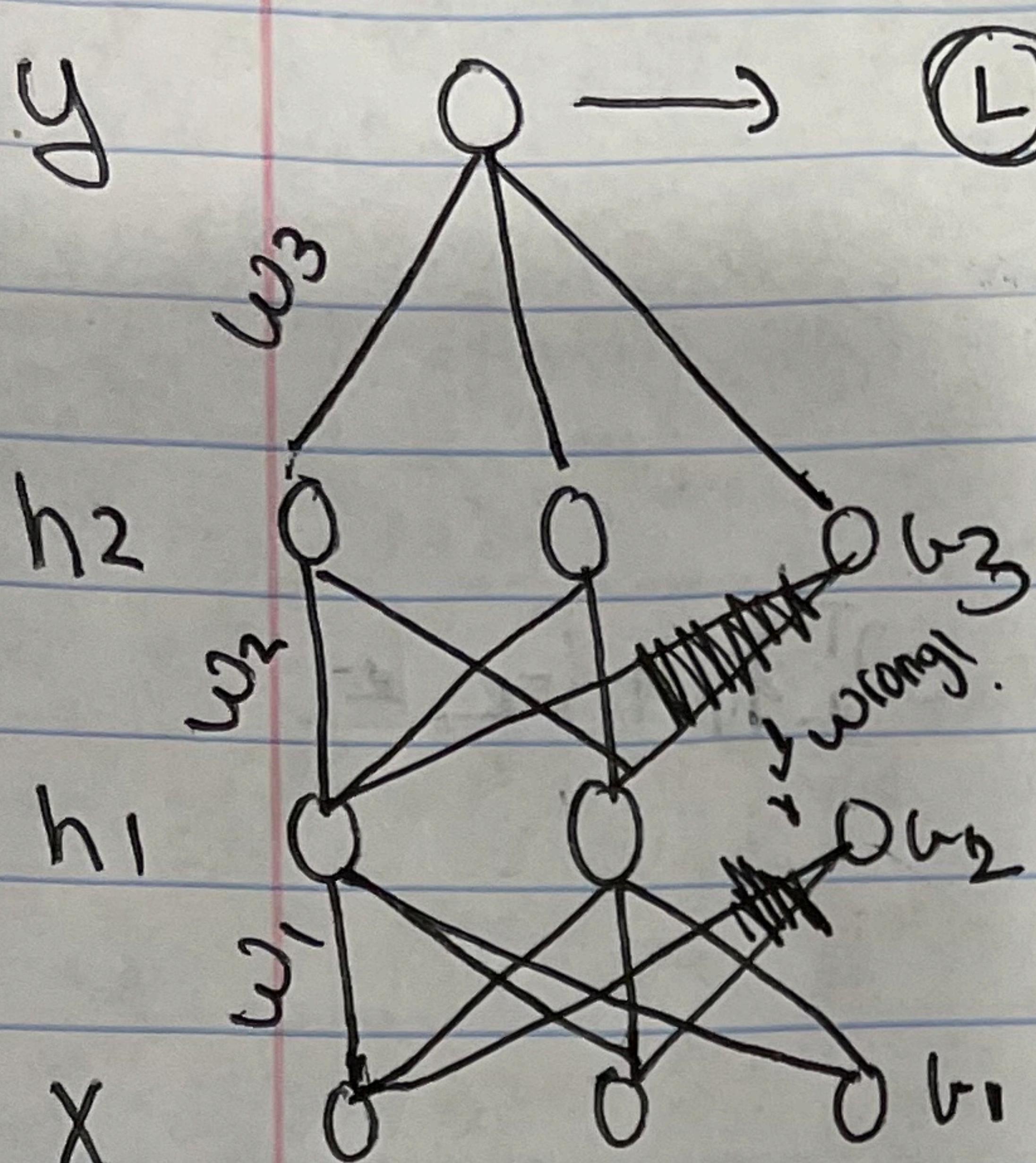
$$x^2 + x^2 - 2ax$$

$$2x - 2a = 2(x-a)$$

let

$$z_i = w_i^T h_{i-1} + b_i$$

y



$$y = w_3^T h_2 + b_3$$

$$h_2 = g(w_1^T h_1 + b_2)$$

$$h_1 = g(w_1^T x + b_1)$$

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \rightarrow \text{MSE}$$

$$\frac{\partial L}{\partial \hat{y}} = \cancel{2(\hat{y} - y)} \quad 2(\hat{y} - y_i)$$

$$\frac{\partial y}{\partial w_3} = h_2 \quad \therefore \frac{\partial L}{\partial w_3} = 2(\hat{y} - y_i) h_2$$

$$\frac{\partial y}{\partial b_3} = 1 \quad \therefore \frac{\partial L}{\partial b_3} = 2(\hat{y} - y)$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial z_2} \quad \therefore \frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

$g = 1$ Lipschitz

$$\frac{\partial z_2}{\partial w_2} = h_1; \quad \frac{\partial h}{\partial z_2} = \frac{\partial g(z_2)}{\partial (z_2)}; \quad \frac{\partial y}{\partial h_2} = w_3 \quad \frac{\partial L}{\partial y} = 2(\hat{y} - y)$$

$$\therefore \frac{\partial L}{\partial w_2} = h_1 \left(\frac{\partial g(z_2)}{\partial (z_2)} \cdot w_3 \cdot 2(\hat{y} - y) \right)^T$$

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

$$y = \mathbf{w}_3^T h_2 + b_3 \quad h_2 = g(z_2) \quad z_2 = \mathbf{w}_2^T h_1 + b_2$$

$$h_1 = g(z_1) \quad z_1 = \mathbf{w}_1^T \mathbf{x} + b_1$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial b_2}$$

$$\frac{\partial z_2}{\partial b_2} = 1 \quad \text{since } z_2 = \mathbf{w}_2^T h_1 + b_2$$

~~$$\frac{\partial z_2}{\partial b_2} = \frac{\partial z_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial b_2}$$~~

$$\therefore \frac{\partial L}{\partial b_2} = \frac{\partial (g(z_2))}{\partial z_2} \cdot \mathbf{w}_3 \cdot 2(\hat{y} - y)$$

Gradients w.r.t \mathbf{w}_1, b_1

$$\frac{\partial L}{\partial \mathbf{w}_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial b_2} \frac{\partial b_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial \mathbf{w}_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial b_2} \frac{\partial b_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$\frac{\partial h_1}{\partial z_1} = \frac{\partial g(z_1)}{\partial z_1} \quad \frac{\partial z_1}{\partial \mathbf{w}_1} = \mathbf{x} \quad \frac{\partial z_1}{\partial b_1} = 1$$

$$\therefore \frac{\partial L}{\partial \mathbf{w}_1} = \mathbf{x} \left(2(\hat{y} - y) \cdot \mathbf{w}_3 \cdot \frac{\partial (g(z_2))}{\partial z_2} \cdot \mathbf{w}_2 \cdot \frac{\partial g(z_1)}{\partial z_1} \right)^T$$

$$\frac{\partial L}{\partial b_1} = 1 \cdot 2(\hat{y} - y) \cdot \mathbf{w}_3 \cdot \frac{\partial g(z_2)}{\partial z_2} \cdot \mathbf{w}_2 \cdot \frac{\partial g(z_1)}{\partial z_1}$$

$$z_n = w^T x$$

let ~~$z_n = w^T x$~~ $\hat{y}_n = \sigma(z_n)$

$$Q.3.) L(D; w) = \sum -y_n \log(w^T x) - (1-y_n) \log(1-\sigma(w^T x))$$

$$\therefore \frac{\partial L}{\partial w} = \frac{\partial L}{\partial z_n} \frac{\partial z_n}{\partial w} \quad \frac{\partial z_n}{\partial w} = x_n$$

$$\begin{aligned} \frac{\partial L}{\partial z_n} &= \left(-\frac{y_n}{\hat{y}_n} + \frac{1-y_n}{1-\hat{y}_n} \right) \hat{y}_n (1-\hat{y}_n) \quad (\text{chain rule}) \\ &= (\hat{y}_n - y_n) \end{aligned}$$

~~$\frac{\partial z_n}{\partial w}$~~

$$① \therefore \frac{\partial L}{\partial w} = \sum_{n=1}^N \frac{\partial L}{\partial z_n} \frac{\partial z_n}{\partial w} = (\hat{y}_n - y_n) x_n$$

② Hessian = second derivatives

$$H = \frac{\partial^2 L}{\partial w \partial w^T}$$

$$\therefore H_{ij} = \sum_{n=1}^N \frac{\partial^2 L}{\partial w_i \partial w_j} = \sum_{n=1}^N \hat{y}_n (1-\hat{y}_n) x_i x_{nj}$$

\therefore The hessian = $x^T R x$ for $R =$
the diagonal matrix with entries
 $R_{nn} = \hat{y}_n (1-\hat{y}_n)$

the diagonal approximation is just

$$H_{ii} = \sum_{n=1}^N \hat{y}_n(1-\hat{y}_n)X_{ni}^2 \text{ and for } H_{ij} \text{ where } i=j \neq 0$$