

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

$y_i \in \{0, 1\}$ true label

$\hat{y}_i \in [0, 1]$ predicted

$$E_{w,b} = -\sum_n [y_n \ln \hat{y}(x_n) + (1 - y_n) \ln (1 - \hat{y}(x_n))] \quad (1)$$

$E_{w,b} \geq 0 \rightarrow$ has a min because both $\ln \leq 0$

$$\hat{y}(x_n) = \sigma(w^T x_n + b) = \frac{1}{1 + e^{-(w^T x_n + b)}} = \sigma(z)$$

$$\frac{\partial E}{\partial w} \rightarrow \frac{\partial E_n}{\partial w} = \frac{\partial E_n}{\partial \hat{y}(x_n)} \cdot \frac{\partial \hat{y}(x_n)}{\partial (z_n)} \cdot \frac{\partial z_n}{\partial w}$$

$$\frac{\partial E_n}{\partial \hat{y}(x_n)} = \frac{-y_n}{\hat{y}(x_n)} + \frac{1 - y_n}{1 - \hat{y}(x_n)} = \frac{-y_n + \hat{y}(x_n)}{\hat{y}(x_n)(1 - \hat{y}(x_n))}$$

$$\frac{\partial \hat{y}(x_n)}{\partial (z_n)} = \frac{\partial \sigma(z)}{\partial (z_n)} = \sigma(z_n)(1 - \sigma(z_n)) = \hat{y}(x_n)(1 - \hat{y}(x_n))$$

$$\frac{\partial z_n}{\partial w} = x_n$$

$$\rightarrow \frac{\partial E_n}{\partial w} = \frac{\hat{y}(x_n) - y_n}{\hat{y}(x_n)(1 - \hat{y}(x_n))} \times \hat{y}(x_n)(1 - \hat{y}(x_n)) \times x_n$$

$$\rightarrow \frac{\partial E_{w,b}}{\partial w} = \sum_n (\hat{y}(x_n) - y_n) x_n$$

$$\frac{\partial E}{\partial b} = \sum_n (\hat{y}(x_n) - y_n)$$

we can prove that the $E_{w,b}$ function has a minimum because it has a lower bound and the loss function is convex: with GD we can find the local (global) minimum.

Using gradient descent we can say:

$$w = w + \eta \sum_n (y_n - \hat{y}(x_n)) x_n$$

η is learning rate

$$b = b + \eta \sum_n (y_n - \hat{y}(x_n))$$

We compute $\hat{y}(x_n)$ for each sample using the current w and b .

then using GD we update w and b .

We repeat the process until the cost function $E_{w,b}$ converges to a minimum.

basically we wouldn't see any difference between updated w, b and previous ones, so we can say we have a minimum point at this situation.

to prove we have a minimum we use Hessian matrix \rightarrow

$$H = \begin{bmatrix} \frac{\partial^2 E}{\partial w^2} & \frac{\partial^2 E}{\partial w \partial b} \\ \frac{\partial^2 E}{\partial b \partial w} & \frac{\partial^2 E}{\partial b^2} \end{bmatrix}$$

$$|H| > 0$$

each term in H is non-negative $\rightarrow H$ is PSD $\rightarrow E$ is convex

$$\frac{\partial^2 E}{\partial w^2} = \sum_n \frac{\partial}{\partial w} ((\hat{y}(x_n) - y_n) x_n) = \sum_n \hat{y}(x_n) (1 - \hat{y}(x_n)) x_n^2 \geq 0$$

$$\frac{\partial^2 E}{\partial b^2} = \sum_n \hat{y}(x_n) (1 - \hat{y}(x_n)) \geq 0$$

$$\frac{\partial^2 E}{\partial w \partial b} = \sum_n \hat{y}(x_n) (1 - \hat{y}(x_n)) x_n$$

$$\nabla E(w, b) = 0 \rightarrow \text{minimum}$$

Covariate Shift:

(2) $\frac{1}{\sqrt{n}}$

Covariate shift refers to the change in the input distribution to a neural network layer during training. In neural network, as the input propagates through each layer, the distribution of the activations (input values to each layer) can shift due to weight updates. This shift can slow down training, as each layer has to constantly adjust to the new distribution of inputs, which affects the learning process.

How BN solves covariate shift:

BN helps mitigate covariate shift by normalizing the input of each layer to have a mean of zero and a standard deviation of one (or some other desired values controlled by learnable params). By doing this, BN stabilizes the distribution of layer inputs, making the training faster and more stable. Essentially, BN ensures that each layer receives inputs with a consistent distribution, reducing the impact of covariate shift.

How BN helps generalization:

BN adds a small amount of noise to the activations due to the random sampling of batches. This noise acts as a regularizer, similar to dropout, which helps prevent overfitting and improves the network's generalization ability. Since each mini-batch is normalized independently, BN introduces variations in the input seen by the model during training. This added noise can encourage the model to learn more robust features that generalize better to unseen data.

$$\text{since } \hat{x}_i = x_i - \mu \quad \left(\mu = \frac{1}{n} \sum_{k=1}^n x_k \right) \rightarrow \frac{\partial \hat{x}_i}{\partial x_j} = \delta_{ij} - \frac{1}{n}$$

$$\text{for calculating } \frac{\partial L}{\partial x_j} \Rightarrow \frac{\partial L}{\partial x_j} = \sum_{i=1}^n \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_j} \quad \hookrightarrow \begin{cases} \text{if } i=j \rightarrow \delta_{ij}=1 \\ \text{if } i \neq j \rightarrow \delta_{ij}=0 \end{cases}$$

$$\leadsto \text{given } y_i = \gamma \hat{x}_i + \beta, \text{ we have } \frac{\partial y_i}{\partial \hat{x}_i} = \gamma. \text{ so } \frac{\partial L}{\partial x_j} = \sum_{i=1}^n \frac{\partial L}{\partial y_i} \cdot \gamma \cdot \left(\delta_{ij} - \frac{1}{n} \right)$$

$n=1$ $\mu = x_1$ $\hat{x}_1 = x_1 - x_1 = 0$ the derivative $\frac{\partial L}{\partial x_1}$ would simply depend on the behavior of y_1 , since there is no variation in a single sample case.

$n \rightarrow \infty$, the effect of subtracting the mean stabilizes across a large batch, so the normalization becomes more consistent.

$$\frac{\partial \hat{y}_k}{\partial z_i^{(2)}} \rightarrow \hat{y}_i = \frac{e^{z_i^{(2)}}}{\sum_{j=1}^K e^{z_j^{(2)}} \rightarrow \frac{\partial \hat{y}_k}{\partial z_i^{(2)}} = \begin{cases} i=k \rightarrow \hat{y}_k(1-\hat{y}_k) \\ i \neq k \rightarrow -\hat{y}_k \hat{y}_i \end{cases} \quad \text{from chain rule} \quad \text{③ سوال}$$

$$\frac{\partial L}{\partial z_i^{(2)}} \rightarrow \begin{matrix} y \text{ is one-hot} \\ y_k=1 \\ y_i=0 \quad k \neq i \end{matrix} \quad L = \sum_{i=1}^K y_i \log(\hat{y}_i) \xrightarrow{y_k=1} L = -\log(\hat{y}_k) \quad \text{ب}$$

$$\rightarrow \frac{\partial L}{\partial z_i^{(2)}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i^{(2)}}$$

$$\rightarrow \text{for } i=k \text{ (the correct class): } \frac{\partial L}{\partial z_k^{(2)}} = -\frac{1}{\hat{y}_k} \cdot \hat{y}_k(1-\hat{y}_k) = \hat{y}_k - 1$$

$$\text{for } i \neq k : \frac{\partial L}{\partial z_i^{(2)}} = -\frac{1}{\hat{y}_k} \cdot (-\hat{y}_k \hat{y}_i) = \hat{y}_i$$

$$\text{So we have: } \frac{\partial L}{\partial z_i^{(2)}} = \begin{cases} \hat{y}_k - 1 & \text{if } i=k \\ \hat{y}_i & \text{if } i \neq k \end{cases}$$

$$\frac{\partial L}{\partial z^{(2)}} = \hat{y} - y, \quad \frac{\partial L}{\partial a^{(1)}} = \frac{\partial L}{\partial z^{(2)}} \cdot (w^{(2)})^T = (\hat{y} - y) \cdot (w^{(2)})^T \quad \frac{\partial L}{\partial w^{(1)}}? \quad (✓)$$

dropout = binary mask $\rightarrow M = \begin{cases} 1 \\ 0 \end{cases}$ $\begin{matrix} 1-P = 0.8 \\ P = 0.2 \end{matrix} \rightarrow \frac{\partial L}{\partial \hat{a}^{(1)}} = (\hat{y} - y) \cdot (w^{(2)})^T \otimes M$
 element wise

$$\frac{\partial \hat{a}_i^{(1)}}{\partial z_i^{(1)}} = \begin{cases} 1 & z_i^{(1)} > 0 \\ 0.1 & z_i^{(1)} < 0 \end{cases}$$

$$\mathbf{I} \rightarrow \mathbf{I}(z^{(1)} > 0) z^{(1)} + 0.1 \mathbf{I}(z^{(1)} \leq 0) z^{(1)}$$

we write down leaky relu

A calculation: $\frac{\partial L}{\partial z^{(1)}} = \frac{\partial L}{\partial \hat{a}^{(1)}} \otimes \text{leaky relu}'(z^{(1)})$ chain rule

$$\rightarrow \frac{\partial L}{\partial z^{(1)}} = ((\hat{y} - y) \cdot (w^{(2)})^T) \otimes M \otimes \text{leaky relu}'(z^{(1)})$$

$\left. \begin{matrix} 1 & z_i^{(1)} > 0 \\ 0.1 & z_i^{(1)} \leq 0 \end{matrix} \right\} \leftarrow$

calculate gradient with respect to $w^{(1)}$:

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial z^{(1)}} \cdot x^T \quad \rightarrow \frac{\partial L}{\partial w^{(1)}} = ((\hat{y} - y) \cdot (w^{(2)})^T) \otimes M \otimes \text{leaky relu}'(z^{(1)}) \cdot x^T$$

سوال 4

$J_y = \begin{bmatrix} \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} & \frac{\partial y}{\partial z} \end{bmatrix}$ → each component represents the rate of change of y with respect to variables u, v , and z is essentially the gradient vector:

$$J_y = \begin{bmatrix} \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} & \frac{\partial y}{\partial z} \end{bmatrix}$$

each component represents the rate of change of y with respect to one of the variables. the jacobian captures this first-order information.

The Hessian matrix H_y is the square matrix of second order partial derivatives of y .

$$H_y = \begin{bmatrix} \frac{\partial^2 y}{\partial u^2} & \frac{\partial^2 y}{\partial u \partial v} & \frac{\partial^2 y}{\partial u \partial z} \\ \frac{\partial^2 y}{\partial v \partial u} & \frac{\partial^2 y}{\partial v^2} & \frac{\partial^2 y}{\partial v \partial z} \\ \frac{\partial^2 y}{\partial z \partial u} & \frac{\partial^2 y}{\partial z \partial v} & \frac{\partial^2 y}{\partial z^2} \end{bmatrix}$$

each element H_{ij} represents the rate of change of the partial derivative with respect to one variable when another variable changes, which gives us information about the concavity or convexity of y in that direction.

the jacobian vector (or gradient) J_y represents the direction and rate of the changes in y with respect to the variables u, v and z . if we are looking at how the function $\psi(u, v, z) = y(u, v, z)$ changes locally. the jacobian provides the first-order approximation.

The Hessian matrix H_y provides the second-order information, capturing how the rate of changes itself varies with each pair of variables. Allows us to understand the curvature.

$$\frac{\partial J_1}{\partial w_i} = -\delta_i x_i \left(y_d - \sum_{k=1}^n \delta_k w_k x_k \right)$$

$$E\left[\frac{\partial J_1}{\partial w_i}\right] = E\left[-\delta_i x_i \left(y_d - \sum_{k=1}^n \delta_k w_k x_k \right)\right] = -x_i \left(E[\delta_i] y_d - \sum_{k=1}^n w_k x_k E[\delta_i \delta_k] \right)$$

$$E[\delta_i \delta_k] = \sigma^2 + 1 \quad (i=k) \quad \text{if } i \neq k \rightarrow E[\delta_i \delta_k] = E[\delta_i] E[\delta_k] = 1$$

$$E[\delta^2] = \text{var}(\delta_i) + E[\delta_i]^2$$

therefore $\rightarrow E\left[\frac{\partial J_1}{\partial w_i}\right] = -x_i \left(y_d - w_i x_i (\sigma^2 + 1) - \sum_{\substack{k=1 \\ k \neq i}}^n w_k x_k \right)$

without regularization $J_1 = 0.5 \left(y_d - \sum_{k=1}^n w_k x_k \right)^2$

$\delta_k \sim \text{Normal}(1, \sigma^2)$

Gaussian DO applied $J_{1, \text{regul}} = 0.5 E_{\delta} \left[\left(y_d - \sum_{k=1}^n \delta_k w_k x_k \right)^2 \right]$

$$= 0.5 \underbrace{\left(y_d - \sum_{k=1}^n w_k x_k \right)^2}_{\text{non regularized}} + 0.5 \sigma^2 \underbrace{\sum_{k=1}^n w_k^2 x_k^2}_{\text{regularization term}}$$

non-regularized target

$$J_1 = 0.5 \left(y_d - \sum_{k=1}^n w_k x_k \right)^2$$

regularized target (Gaussian dropout):

$$J_{1, \text{regul}} = 0.5 \left(y_d - \sum_{k=1}^n w_k x_k \right)^2 + 0.5 \sigma^2 \sum_{k=1}^n w_k^2 x_k^2$$

$0.5 \sigma^2 \sum_{k=1}^n w_k^2 x_k^2$: penalty function for weights \rightarrow decrease overfitting
 \rightarrow smoothness during training because of random noise

Newton : $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$

$f(x) \approx f(\hat{x}) + f'(\hat{x})(x - \hat{x}) + \frac{f''(\hat{x})}{2}(x - \hat{x})^2$: 6 D's
taylor $\rightarrow 0$

error in k 'th iteration $e_k = x_k - \hat{x}$ $\longrightarrow f(x) \approx f'(\hat{x})e_k + \frac{f''(\hat{x})}{2}e_k^2$

calculation $\rightarrow x_{k+1} = x_k - e_k - \frac{f''(\hat{x})}{2f'(\hat{x})}e_k^2$ $\left| \frac{f''(\hat{x})}{2f'(\hat{x})} \right|$

$e_{k+1} = x_{k+1} - \hat{x} = -e_k - \frac{f''(\hat{x})}{2f'(\hat{x})}e_k^2 = -\frac{f''(\hat{x})}{2f'(\hat{x})}e_k^2 \Rightarrow |e_{k+1}| \approx C|e_k|^2$

e_k decreases quadratically \rightarrow newton method converges to x^* , which is the optimal point of function.

a) $\hat{y}_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \rightarrow L(z, y) = - \sum_{k=1}^K y_k \log \left(\frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \right) = - \sum_{k=1}^K y_k (z_k - \log \sum_{j=1}^K e^{z_j})$ السؤال 7

$\frac{\partial}{\partial z_i} \left(\sum_{k=1}^K y_k \log \left(\sum_{j=1}^K e^{z_j} \right) \right) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \sum_{k=1}^K y_k = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} = \hat{y}_i$

$y = \text{one-hot}$

$$\rightarrow \frac{\partial L}{\partial z_i} = -y_i + \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} = -y_i + \hat{y}_i \rightarrow \nabla_z L = \hat{y} - y$$

b) $H_{ij} = \frac{\partial^2 L}{\partial z_i \partial z_j} = \frac{\partial}{\partial z_j} (\hat{y}_i - y_i) = \frac{\partial}{\partial z_j} \left(\frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \right) = \frac{e^{z_i} \cdot \delta_{ij} \sum_{k=1}^K e^{z_k} - e^{z_i} \cdot e^{z_j}}{\left(\sum_{k=1}^K e^{z_k} \right)^2}$

where δ_{ij} is 1 if $i=j$ else is 0 (as previously said). so $H_{ij} = \hat{y}_i (\delta_{ij} - \hat{y}_j)$

$$H = \text{diagonal}(\hat{y}) - \hat{y} \hat{y}^T$$

2) $x \in \mathbb{R}^K \rightarrow x^T H x = x^T D(\hat{y}) x - x^T \hat{y} \hat{y}^T x = \sum_{i=1}^K \hat{y}_i x_i^2 - \left(\sum_{i=1}^K \hat{y}_i x_i \right)^2$

So we just need to say $(x^T \hat{y})^2 \leq \left(\sum_{i=1}^K \hat{y}_i \right) \left(\sum_{i=1}^K x_i^2 \hat{y}_i \right) = \sum_{i=1}^K \hat{y}_i x_i^2$ non-negative $\leftarrow (x^T \hat{y})^2$

$\rightarrow x^T H x \geq 0$ for all $x \in \mathbb{R}^K \rightarrow H$ is PSD الذي نحتاجه

c) As we concluded H is PSD and therefore its function is convex with respect to z and has an optimal point.