

Question 01

- (a) Assuming this is a Regression Problem and, that we use only one datum at a time during training.

$$\text{Loss} = \frac{1}{2} (y - t)^2 \leftarrow \text{Squared Error Loss.}$$

- (b) Forward Pass:-

1	8	0.5	6
3	1.1	-1	-8
-0.5	-6	-3	-1.1
-1	-8	-0.5	-6

Maxpooling →

8	6
-0.5	-0.5

(i)

Maxpooling

8	6
-0.5	-0.5

: values of the matrix U .

(ii)

↓ $\text{LReLU}(\cdot)$ activation function. = $\begin{cases} x & \text{if } x \geq 0 \\ 0.1x & \text{otherwise.} \end{cases}$

8	6
-0.05	-0.05

: values of the matrix V

(iii)

network output $y =$

(iii) Network output $y = \sum_{i=1}^4 v_i w_i$

$$W = \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 10 & 4 \\ \hline \end{array}$$

$$y = (8 \times 1) + (6 \times 2) + (-0.05) \times (10) + (-0.05 \times 4)$$

$$= 8 + 12 + (-0.5) + (-0.2)$$

$$\boxed{y = 19.3}$$

(c) Backward Pass:-

Assumption $\Rightarrow \frac{\partial L}{\partial y} = 2.0$ [where: $L = \frac{1}{2} (y - t)^2$]

(i) Consider forward pass data flow

$$X \rightarrow \text{mexpol}() \rightarrow u \rightarrow \text{LReLU}() \rightarrow v \rightarrow \underbrace{\text{Dense}}_{(v \cdot w)} \rightarrow y \Rightarrow \text{loss}$$

* Using Chain Rule

$$\frac{\partial L}{\partial v} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial v}$$

$$= \frac{\partial L}{\partial y} \times \frac{\partial (y = w \cdot v)}{\partial v}$$

$$\frac{\partial L}{\partial v} = \frac{\partial L}{\partial y} \times w$$

∴ By Applying elementwise

$$\frac{\partial L}{\partial v_1} = \frac{\partial L}{\partial y} \times w_1$$

$$= 2 \times 1$$

$$\underline{\underline{\frac{\partial L}{\partial v_1} = 2}}$$

(ii) $\frac{\partial L}{\partial v} =$

$\frac{\partial L}{\partial v_1} = 2 \times 1$	$\frac{\partial L}{\partial v_2} = 2 \times 2$
$\frac{\partial L}{\partial v_3} = 2 \times 10$	$\frac{\partial L}{\partial v_4} = 2 \times 4$

 $=$

2	4
20	8

(iii) Using the chain Rule recursively:

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times \frac{\partial v}{\partial u}$$

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times \frac{\partial [\text{ReLU}(u)]}{\partial u} \quad \text{Since } v = \text{ReLU}(u)$$

$$u \geq 0 \quad \swarrow \quad \searrow \quad u < 0$$

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times \frac{\partial (v = u)}{\partial u}$$

↓

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times 1$$

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times \frac{\partial (v = 0.1u)}{\partial u}$$

↓

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial v} \times 0.1$$

$$\therefore \frac{\partial L}{\partial u} = \begin{cases} \frac{\partial L}{\partial v} & ; \text{ for } u \geq 0 \\ 0.1 \frac{\partial L}{\partial v} & ; \text{ otherwise} \end{cases}$$

$$\therefore \frac{\partial L}{\partial u} =$$

$\frac{\partial L}{\partial v}$	$\frac{\partial L}{\partial v}$
$0.1 \frac{\partial L}{\partial v}$	$0.1 \frac{\partial L}{\partial v}$

 \Rightarrow

2	4
2	0.8

$$(iv) \quad \frac{\partial L}{\partial x} = \frac{\partial L}{\partial u} \times \frac{\partial u}{\partial x}$$

$$u = \max(x)$$

$$\frac{\partial L}{\partial x} = \begin{cases} \frac{\partial L}{\partial u} \times 1 \\ \frac{\partial L}{\partial u} \times 0 \end{cases}$$

where x is maximum (maximum)

otherwise.

$$\frac{\partial L}{\partial x} =$$

0	2	0	4
0	0	0	0
2	0	0	0
0	0	0.8	0

Gradient becomes zero as there is no affect from those weights.

(d.)

(i)

14	-4
-6	3

U

ReLU →

14	$\alpha_0(4)$
$\alpha_0(-6)$	3

=

14	$-4\alpha_0$
$-6\alpha_0$	3

V

Answer

(ii)

$\frac{\partial L}{\partial V}$

2.0	-2.5
1.5	3.0

$$V_i = \begin{cases} u_i & \text{for } u_i \geq 0 \\ \alpha_0 u_i & \text{otherwise} \end{cases}$$

* Using chain rules

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial V} \times \frac{\partial V}{\partial u} = \begin{cases} \frac{\partial L}{\partial V} & \text{if } u \geq 0 \\ \alpha_0 \frac{\partial L}{\partial V} & \text{otherwise} \end{cases}$$

* Using chain rule again

$$\frac{\partial L}{\partial \alpha} = \frac{\partial L}{\partial V} \times \frac{\partial V}{\partial \alpha}$$

$$\frac{\partial L}{\partial \alpha} = \frac{\partial L}{\partial V} \times \frac{\partial V}{\partial \alpha}$$

if $u_i \geq 0$ otherwise

$$\frac{\partial L}{\partial \alpha} = \frac{\partial L}{\partial V} \times \frac{\partial u}{\partial \alpha}$$

$$\frac{\partial L}{\partial \alpha} = \frac{\partial L}{\partial V} \times \frac{\partial (\alpha u)}{\partial \alpha}$$

$$\therefore \frac{\partial L}{\partial \alpha} = \begin{cases} 0 & \text{for } (u_i \geq 0) \\ \frac{\partial L}{\partial V} \times u_i & \text{otherwise} \end{cases}$$

$$\Rightarrow \frac{\partial L}{\partial \alpha} =$$

0	$(-2.5 \times (-4))$
$1.5 \times (-6)$	0

Answer

→

$$\frac{\partial L}{\partial \alpha} \bigg|_{\alpha_0} =$$

0	10
-9	0

Question 02

(a) Definition of softmax = $p_i^o = \frac{\exp(Z_i)}{\sum_{j=1}^n \exp(Z_j)}$

(i) consider \Rightarrow $\text{softmax}(\bar{Z}_i) = \frac{\exp(Z_i - Z_{\max})}{\sum_{j=1}^n \exp(Z_j - Z_{\max})}$; $Z_{\max} = \max\{Z_i\}$
 $\forall i \in \{1, 2, \dots, n\}$

$$= \frac{\exp(Z_i) / \exp(Z_{\max})}{\sum_{j=1}^n \exp(Z_j) / \exp(Z_{\max})}$$

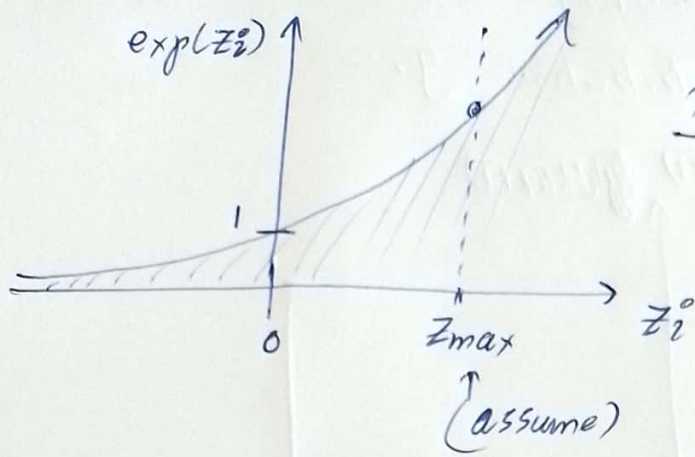
$$= \frac{\left\{ \frac{1}{\exp(Z_{\max})} \right\} \times \exp(Z_i)}{\left\{ \frac{1}{\exp(Z_{\max})} \right\} \times \sum_{j=1}^n \exp(Z_j)}$$

$$= \frac{\exp(Z_i)}{\sum_{j=1}^n \exp(Z_j)} \equiv \text{softmax}(Z_i)$$

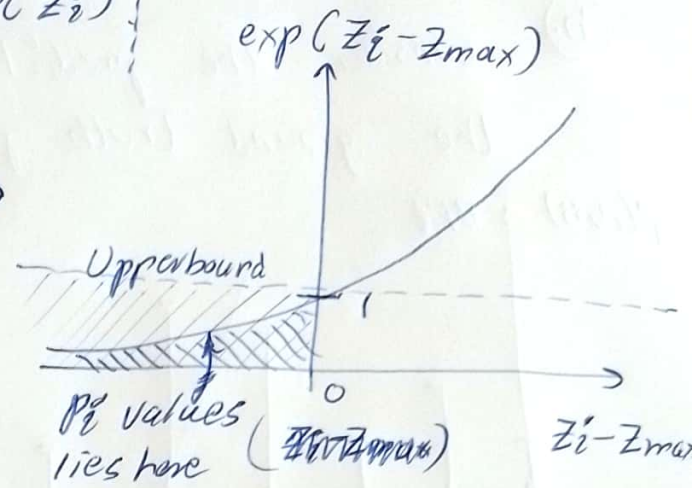
\Downarrow

$$\therefore \text{softmax}(Z) = \text{softmax}(\bar{Z}) //$$

(7)

(ii) Consider the plot of the $\exp(z_i)$ 

Transform



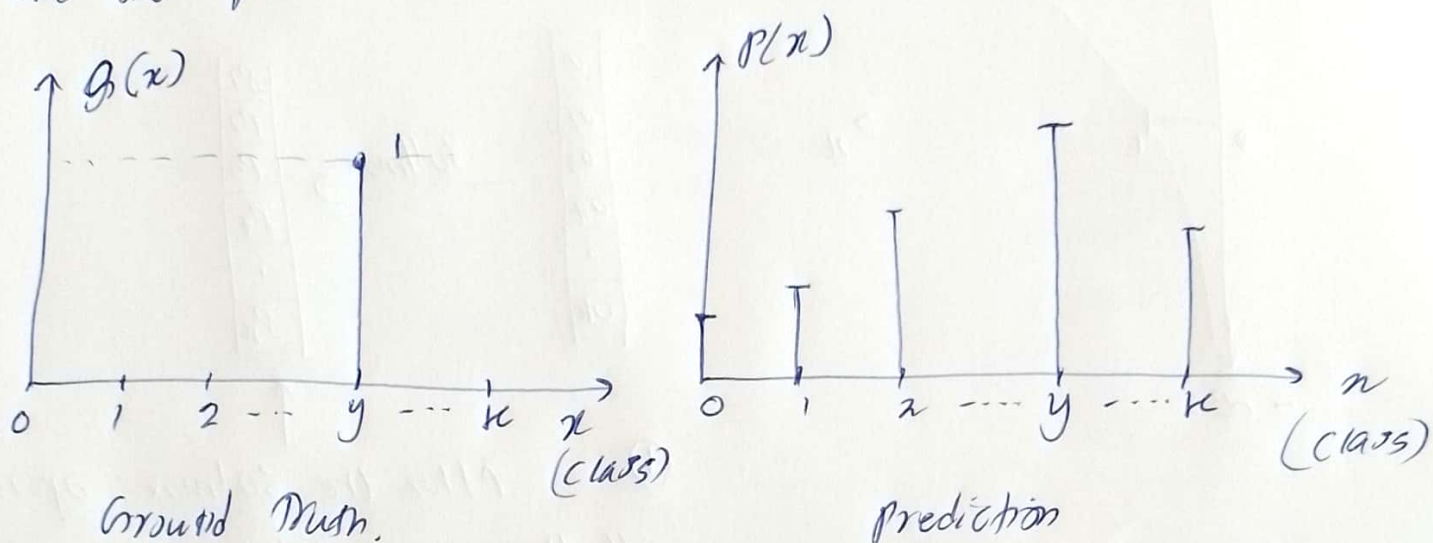
* Subtracting the max of the logit vector is done, to "improve the numerical stability" when training deep learning models.

* If we don't perform this transformation, the terms in the softmax function become very large leading to numerical overflow.

* This can cause issues when training the model, as the gradient of loss function w.r.t the logit values may become very large or NaN (Not a Number).

⑤

(i) Consider the probability distributions of the Ground Truth and the predicted values.



* Cross-entropy distance = $H(g, p) = - \sum_{x \in X} g(x) \cdot \log(P(x))$.

* However, due to the nature of the ground truth:

① $g(x) = 0 \quad \forall x \in \{1, 2, 3, \dots, k\} \setminus \{y\}$

② $g(y) = 1$: for the correct class

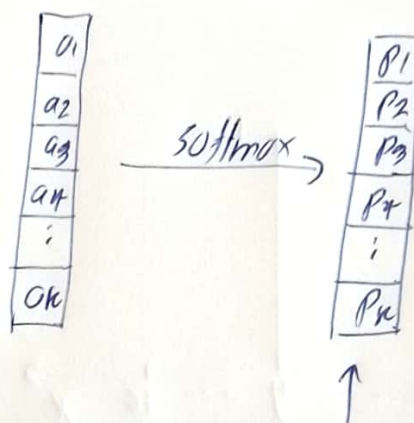
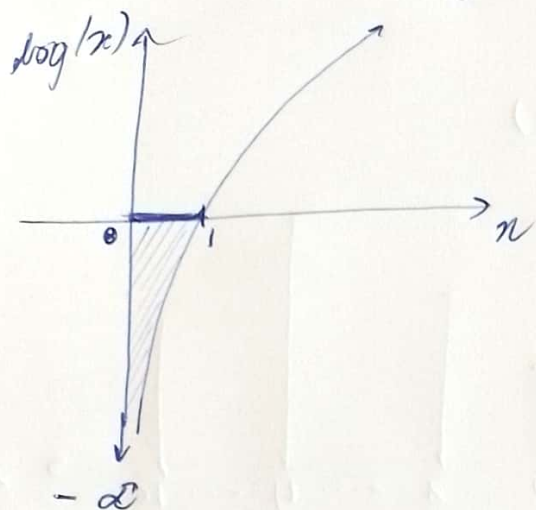
∴ Above summation can be further simplified as follows.

$$H(g, p) = \underbrace{[-g(1) \log(P_1)]}_{(=0)} + \underbrace{[-g(2) \log(P_2)]}_{(=0)} + \dots + \underbrace{[-g(y) \log(P_y)]}_{(=1)} + \dots + \underbrace{[-g(k) \log(P_k)]}_{(=0)}$$

$$\begin{aligned} \therefore H(g, p) &= L(g, p) = - \underbrace{g(y) \cdot \log(P_y)} \\ &= - 1 \times \log(P_y) \end{aligned}$$

∴ $L(g, p) = -\log(P_y)$

(ii) consider the $\log(x)$ function:



① After the softmax operation all the probability values comes to the range of $[0, 1]$.

② The corresponding (\log) values for the range $[0, 1]$ lies between in the range $(-\infty, 0]$.

③ consider the loss function ($\because -\log(p_y) \in [0, \infty)$)

$$-\log(x) = -\log(p_y)$$

if $(p_y = 1)$ which is the final goal, $(-\log(1) = 0)$.

Then the loss is zero.



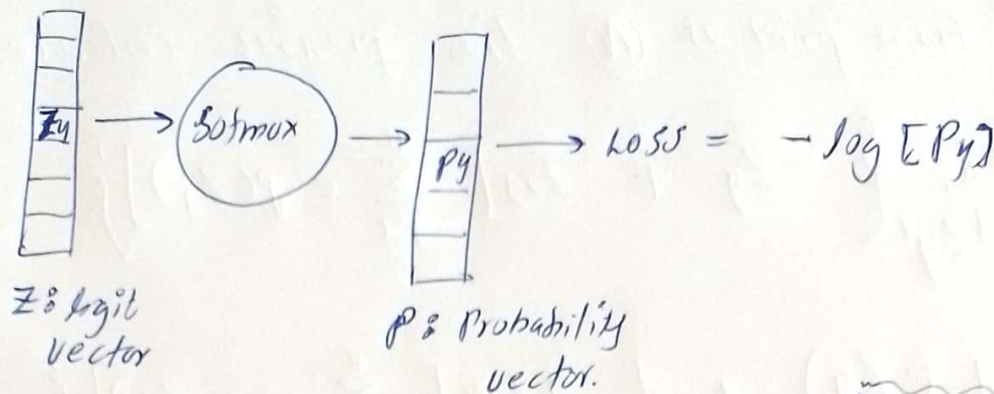
\therefore We try to minimize the $[-\log(p_y)]$.

if $0 < p_y < 1$

$[-\log(p_y)]$ becomes a large positive value, making the loss very large.

\therefore By making this loss minimum we can increase the $Pr(\text{correct class})$.

© (i)



$$I. \frac{dL}{dz_y} = \frac{d}{dz_y} \left\{ -\log [p_y] \right\}$$

$$= \frac{d}{dz_y} \left\{ -\log \left[\frac{\exp(z_y)}{\sum_{j=1}^K \exp(z_j)} \right] \right\}$$

Recall:

$$\frac{d f(x)}{dx} = \frac{1}{f(x)} \times f'(x)$$

$$= \frac{d}{dz_y} \left\{ -\log [\exp(z_y)] + \log \left[\sum_{j=1}^K \exp(z_j) \right] \right\}$$

$$= \frac{d}{dz_y} \left\{ -z_y + \log \left[\sum_{j=1}^K \exp(z_j) \right] \right\}$$

$$= \frac{d}{dz_y} (-z_y) + \frac{d}{dz_y} \log \left[\sum_{j=1}^K \exp(z_j) \right] \rightarrow \textcircled{A}$$

$$= -1 + \frac{1}{\sum_{j=1}^K \exp(z_j)} \times \frac{d}{dz_y} \left[\sum_{j=1}^K \exp(z_j) \right]$$

$$= -1 + \frac{1}{\sum_{j=1}^K \exp(z_j)} \times \frac{d}{dz_y} \left[\exp(z_1) + \exp(z_2) + \dots + \exp(z_y) + \dots + \exp(z_K) \right]$$

$$\therefore \frac{dL}{dz_y} = -1 + \frac{\exp(z_y)}{\sum_{j=1}^K \exp(z_j)} \left(\frac{d}{dz_y} (\exp(z_i)) = 0 \text{ if } i \neq y \right)$$

$$\boxed{\frac{dL}{dz_y} = -1 + p_y} \leftarrow \text{Answer}$$

(ii) Starting from point of (A) from previous calculations: (11)

$$\frac{dL}{dz_{y'}} = \frac{d}{dz_{y'}} \left\{ -z_y + \log \left[\sum_{j=1}^k \exp(z_j) \right] \right\}$$

$$= \underbrace{\frac{d}{dz_{y'}} (-z_y)}_{=0} + \frac{d}{dz_{y'}} \left\{ \log \left[\sum_{j=1}^k \exp(z_j) \right] \right\}$$

$$= \frac{1}{\sum_{j=1}^k \exp(z_j)} \times \frac{d}{dz_{y'}} \left\{ \sum_{j=1}^k \exp(z_j) \right\}$$

$$= \frac{1}{\sum_{j=1}^k \exp(z_j)} \times \frac{d}{dz_{y'}} \left\{ \exp(z_1) + \exp(z_2) + \dots \right. \\ \left. + \exp(z_{y'}) + \dots \exp(z_y) + \dots + \exp(z_n) \right\}$$

$$\left(\begin{array}{l} \frac{d}{dz_{y'}} \exp(z_j) = 0 \\ \forall j \neq y' \end{array} \right)$$

$$\therefore \frac{dL}{dz_{y'}} = \frac{1}{\sum_{j=1}^k \exp(z_j)} \times \exp(z_{y'}) \quad \leftarrow (\text{softmax for } z_{y'})$$

\Downarrow

$$\boxed{\therefore \frac{dL}{dz_{y'}} = p_{y'}}$$

$$\forall y' \in \{1, 2, \dots, k\} \setminus \{y\}$$