# Softmax
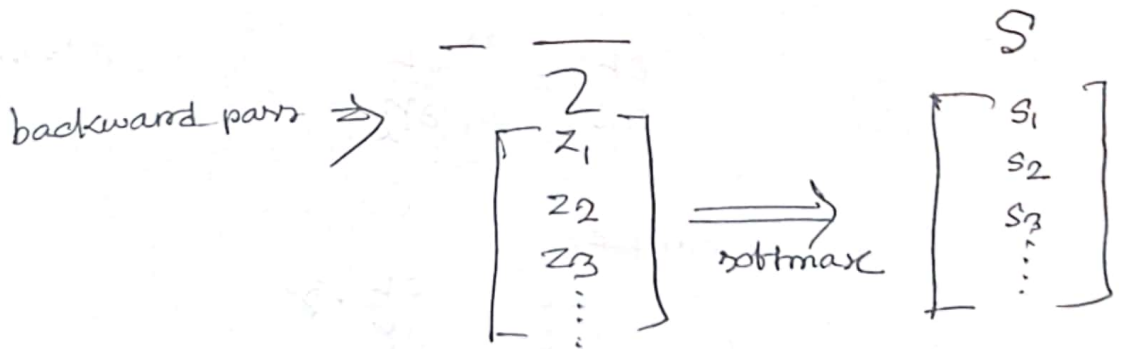
given outputs from prev. layer $= \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \end{bmatrix}$

for every $z_i$,

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum\limits_{3=1}^{n} e^{z_3}}$$

$$= \frac{e^{z_i - max(z)}}{\sum\limits_{3=1}^{n} e^{z_3 - max(z)}}$$

$\left[ \text{to stop overflow} \right]$

backward pass $\Rightarrow$ $\overset{Z}{\begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \end{bmatrix}} \xrightarrow[\text{softmax}]{} \overset{S}{\begin{bmatrix} s_1 \\ s_2 \\ s_3 \\ \vdots \end{bmatrix}}$

$$\frac{\partial s_i}{\partial z_3} = \begin{cases} s_i(1-s_i) & \text{when } i=3 \\ -s_i s_3 & \text{when } i \neq 3 \end{cases}$$

$J_n = \boxed{\text{diag}(S) - SS^T \dots}$

## Batch normalization

### forward

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$Y = \gamma \cdot \hat{x} + \beta.$$

### backward

$$\frac{\partial L}{\partial \gamma} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial \gamma}$$

$$= \left( \sum_{batch} \hat{x}_i \times \frac{\partial L}{\partial Y_i} \right)$$

$$= \sum_i \frac{\partial L}{\partial Y_i} \hat{x}_i$$

sum because gamama has effect on all the samples in the batch.

$$\frac{\partial L}{\partial \gamma} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i} \hat{x}_i$$

$$\frac{\partial L}{\partial \beta} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i} \qquad [m = \text{batch size}]$$

for a sample,

$$\frac{\partial L}{\partial \hat{x}} = \frac{\partial L}{\partial Y} \cdot \gamma.$$

$$\frac{\partial L}{\partial \sigma^2} = \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial \hat{X}} \times \frac{\partial \hat{X}}{\partial \sigma^2}$$

$$= \frac{\partial L}{\partial Y} \times \frac{\partial Y}{\partial \hat{X}} \times \left(x - \mu\right) \times -\frac{1}{2} \times \left(\sqrt{\sigma^2 + \epsilon}\right)^{-3/2}$$

for many samples $\Rightarrow$

$$\frac{\partial L}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i} \times \frac{\partial Y_i}{\partial \hat{X}_i} \times (x_i - \mu) \times -\frac{1}{2} \times \left(\sqrt{\sigma^2 + \epsilon}\right)^{-3/2}$$

$\not{\cancel{\times}}$ same way ( chain rule and partial derivative ),

we can find, $\frac{\partial L}{\partial X}$ and $\frac{\partial L}{\partial \mu}$ .

$\not{\cancel{\times}}$ sum over $x_i$, $\hat{X}$, or $Y$ is important

for $Y, \sigma^2, \beta, \mu$ because they are

same over all samples.