

$$L(y, \hat{y}) = \sum_{i=1}^{10} y_i \log(\hat{y}_i)$$

1) softmax layer

$$\frac{\partial L}{\partial \hat{y}_i} = -y_i / \hat{y}_i \quad i \in [1, 10]$$

2) output of Dense layer

$$\frac{\partial L}{\partial z_j} = \hat{y}_j - y_j \quad j \in [1, 10]$$

3) Dense layer

$$\frac{\partial L}{\partial w_{jk}} = \frac{\partial L}{\partial z_j} \times x_k \quad x_k = k^{\text{th}} \text{ input of dense layer}$$

$$\frac{\partial L}{\partial b_j} = \frac{\partial L}{\partial z_j}$$

4) Maxpool 2

$$\frac{\partial L}{\partial \tilde{x}_i} = \frac{\partial L}{\partial z_j} \times 1 \quad \text{if } x_i = \text{max in pooling window}$$

$$= 0 \quad \text{otherwise}$$

5) Conv 2 sigmoid 2

where, $\sigma'(n) = \sigma(n)(1 - \sigma(n))$

~~$\delta L / \delta n$~~

$$\delta L / \delta y_{jk}^2 = \sum_i \sum_{p,q} \delta L / \delta \hat{n}_{ipq}^2 \times w_{ipq}^2 \times \sigma'(y_{jk}^2)$$

where \hat{n}_{ipq}^2 is output of maxpool 1 and corresponds to ~~max~~ index
max value at position (i, p, q) of the input of maxpool 1
& $\sigma'(y_{jk}^2)$ is evaluated at output of Conv 2 sigmoid 2 layer,

~~6) Lossy bias~~

$$\delta L / \delta b_j^2 = \sum_k \delta L / \delta y_{jk}^2 \quad \text{bias gradient}$$

$$\delta L / \delta w_{ipq}^2 = \sum_k \delta L / \delta y_{jk}^2 \hat{n}_{i+p, q+k}^1$$

where $\hat{n}_{i+p, q+k}^1$ ~~corresponds to~~ ^{is} output of maxpool 1 corresponding
to window containing element at position $(i+p, q+k)$

4) Layer 1 (Conv 1 sigmoid 1 Maxpool 1)

$$\delta L / \delta y_{pq}^1 = \sum_j \sum_k \delta L / \delta y_{jk}^2 \times w_{jp(k+u)q+v}^2 \times \sigma'(y_{pq}^1)$$

following above rules.

$$\left. \begin{aligned} \delta L / \delta w &= \delta L / \delta \hat{y} \times \delta \hat{y} / \delta w = (\hat{y} - y) \times (x)^T \\ \delta L / \delta b &= \delta L / \delta \hat{y} \times \delta \hat{y} / \delta b = (\hat{y} - y) \end{aligned} \right\} \text{[Ans.]}$$

$$\delta L / \delta b_p^1 = \sum_q \sum_i \sum_j \left(\frac{\delta L}{\delta y_{jq}^2} \frac{\delta y_{jq}^2}{\delta b_p^1} \right)$$

$$\delta L / \delta w_{pq}^1 = \sum_i \sum_j \left(\frac{\delta L}{\delta y_{jq}^2} \frac{\delta y_{jq}^2}{\delta w_{pq}^1} \right)$$

3) Given

a) 2 K-D Gaussian distributions $N(\mu_0, \Sigma_0), N(\mu_1, \Sigma_1)$

R.T.P

$$D_{KL}[N(\mu_0, \Sigma_0) || N(\mu_1, \Sigma_1)] = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) \right) - K + \log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right)$$

Proof: We know that the KL divergence between 2 PDFs are given

$$\text{as } D_{KL}(P_1 || P_2) = \sum_{n \in X} P_1(n) \log \left(\frac{P_1(n)}{P_2(n)} \right)$$

$$= - \sum_{n \in X} P_1(n) \log \left(\frac{P_2(n)}{P_1(n)} \right)$$

$$= + E_{P_1} \log \left(\frac{P_1}{P_2} \right)$$

~~Pdf~~ Pdf of a multivariate Gaussian distribution is given as

$$p(n) = \frac{1}{(2\pi)^{N/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2} (n - \mu)^T \Sigma^{-1} (n - \mu)}$$

$$\text{In this case } D_{KL}[N(\mu_0, \Sigma_0) || N(\mu_1, \Sigma_1)] = E_{N_0} (\log N_0 - \log N_1)$$

$$= \frac{1}{2} E_{N_0} \left[-\log(\det(\Sigma_0)) - (n - \mu_0)^T \Sigma_0^{-1} (n - \mu_0) \right. \\ \left. + \log(\det(\Sigma_1)) - (n - \mu_1)^T \Sigma_1^{-1} (n - \mu_1) \right]$$

$$= \frac{1}{2} \log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) + \frac{1}{2} E_{N_0} \left[- (n - \mu_0)^T \Sigma_0^{-1} (n - \mu_0) \right. \\ \left. + (n - \mu_1)^T \Sigma_1^{-1} (n - \mu_1) \right]$$

$$= \frac{1}{2} \log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) + \frac{1}{2} E_{N_0} \left[-\text{tr}(\Sigma_0^{-1} (n - \mu_0)(n - \mu_0)^T) \right. \\ \left. + \text{tr}(\Sigma_1^{-1} (n - \mu_1)(n - \mu_1)^T) \right]$$

$$= \frac{1}{2} \log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) - \frac{1}{2} K + \frac{1}{2} \text{tr}(\Sigma_1^{-1} (\Sigma_0 + \mu_0 \mu_0^T - 2\mu_1 \mu_0^T + \mu_1 \mu_1^T))$$

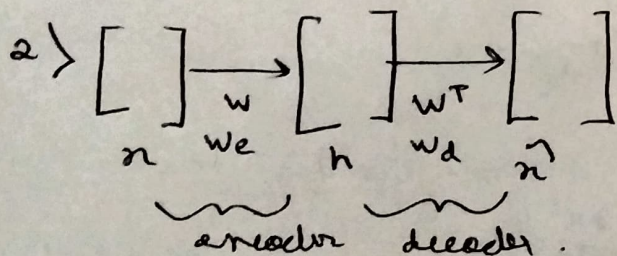
$$\begin{aligned}
&= \frac{1}{2} \left(\log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) - K + \text{tr}(\Sigma_1^{-1} \Sigma_0) + \text{tr}(\mu_0^T \Sigma_1^{-1} \mu_0 - 2 \mu_0^T \mu_1 + \mu_1^T \Sigma_1^{-1} \mu_1) \right) \\
&= \frac{1}{2} \left(\log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) - K + \text{tr}(\Sigma_1^{-1} \Sigma_0) - (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) \right) \\
&= \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - K + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right) \\
&\quad \text{[Proved]}
\end{aligned}$$

b) KL-Divergence is not a metric loss

∴ it does not follow the symmetry property of metric space

$$\therefore D_{KL}(N(\mu_0, \Sigma_0) || N(\mu_1, \Sigma_1)) \neq D_{KL}(N(\mu_1, \Sigma_1) || N(\mu_0, \Sigma_0))$$

[Proved]



$$h = w_e x + b$$

$$\hat{n} = w_d h + c$$

$$= w_d w_e x + w_d b + c$$

$$L = (n - \hat{n})^2 \times \frac{1}{2N}$$

$$= (n - w_d w_e x - w_d b - c)^2 \times \frac{1}{2N}$$

$$\frac{\partial L}{\partial w_d} = (0 - w_e x - b - c) \times \frac{1}{N} (n - w_d w_e x - w_d b - c)$$

$$\frac{\partial L}{\partial w_e} = \frac{1}{N} (n - w_d w_e x - w_d b - c) (0 - (w_d x)^T - 0 - 0)$$

~~$\frac{\partial L}{\partial W}$~~

~~$W^T W^T$~~ $W_d = W^T$

$W_e = W$

$$L = (x - W^T W x - W^T b - c)^2 \approx \frac{1}{2N}$$

$$\frac{\partial L}{\partial W} = (0 - (x + x^T)W - b) (x - W^T W (x - W^T b - c)) \frac{1}{N}$$
$$= (-(x + x^T)W - b) L'(x, \hat{n})$$

$$\frac{\partial L}{\partial W_d} = (-W x - b) L'(x, \hat{n})$$

$$\frac{\partial L}{\partial W_e} = (-(W^T x)^T) L'(x, \hat{n})$$
$$= (-(W x^T)) L'(x, \hat{n})$$

$$\therefore \frac{\partial L}{\partial W_d} + \frac{\partial L}{\partial W_e} = \frac{\partial L}{\partial W} \quad \square \text{ Proved.}$$