

Softmax with temperature from the paper

"Distilling the Knowledge in a Neural Network"

Neural networks typically produce class probabilities by using a "softmax" output layer that converts the **logit,  $z_i$** , computed for each class into a probability,  $q_i$ , by comparing  $z_i$  with the other logits.

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)} \quad (1)$$

2

where  **$T$  is a temperature** that is normally set to 1. Using a higher value for  $T$  produces a softer probability distribution over classes.

\* Key Concept

**Model** learns **probability distribution**

**Loss** measures **difference between model distribution**  
**and data distribution**

**Optimizer** trains model to **decrease loss**

## Details about "Loss"

Training dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

Key assumption: independent and identically distributed (i.i.d.)

$$(x_i, y_i) \perp (x_j, y_j) \quad (x_i, y_i) \sim p(x, y)$$

$$\begin{aligned} \Rightarrow p(D) &= \prod p(x_i, y_i) \\ &= \prod p(x_i) p(y_i | x_i) \end{aligned}$$

$$p(D) = \prod p(x_i) p_\theta(y_i | x_i) \quad \text{model of true } p(y_i | x_i)$$

choose  $\theta$  such that  $p(D)$  is maximized

$$\log p(D) = \sum \cancel{\log p(x_i)} + \log p_\theta(y_i | x_i)$$

$$\theta^* \leftarrow \operatorname{argmax}_{\theta} \sum \log p_\theta(y_i | x_i) \quad \text{maximum likelihood estimation (MLE)}$$

$$\theta^* \leftarrow \operatorname{argmin}_{\theta} - \sum \log p_\theta(y_i | x_i) \quad \text{negative log-likelihood (NLL)} \\ \text{our loss function}$$

Information of an event  $E$

$$I(E) = -\log_2 p(E)$$

Entropy of random variable  $X$

$$H(X) = E(I(X))$$

$$= -\sum p(x_i) \log p(x_i)$$

Cross entropy of the distribution  $Q$  relative to a distribution  $P$

$$\begin{aligned} H(P, Q) &= E_P(I(Q)) \\ &= E_P(-\log q(x)) \\ &= -\sum p(x) \log q(x) \end{aligned}$$

KL Divergence (Kullback - Leibler divergence)

relative entropy

statistical distance measuring how <sup>model</sup> probability distribution  $Q$  is different from <sub>data</sub> reference probability distribution  $P$

$$\begin{aligned} D_{KL}(P \parallel Q) &= E_P(I(Q) - I(P)) \\ &= \sum p(x) \log \frac{p(x)}{q(x)} \\ &= -\sum p(x) \log \frac{q(x)}{p(x)} \end{aligned}$$

JS Divergence (Jensen - Shannon divergence)

$$JSD(P \parallel Q) = \frac{1}{2} D_{KL}(P \parallel \frac{P+Q}{2}) + \frac{1}{2} D_{KL}(Q \parallel \frac{P+Q}{2})$$

Summary

for 1 data point  $(x_i, y_i)$

- cross entropy :  $-\sum_y p(y|x_i) \log p_\theta(y|x_i)$
- negative log-likelihood :  $-\log p_\theta(y_i|x_i)$
- KL divergence :  $-\sum_y p(y|x_i) \log \frac{p_\theta(y|x_i)}{p(y|x_i)}$

KL Divergence  $\longleftrightarrow$  Cross Entropy

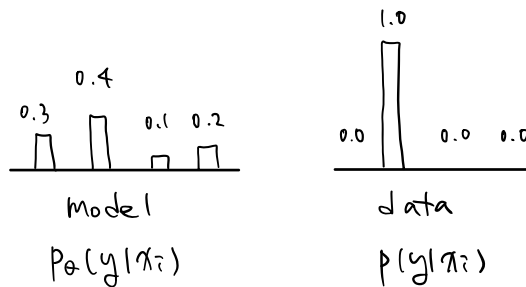
$$\begin{aligned}
 & -\sum p(x) \log q(x) + \cancel{\sum p(x) \log p(x)} & -\sum p(x) \log q(x) \\
 & -\sum p(y|x_i) \log p_\theta(y|x_i) + \cancel{\sum p(y|x_i) \log p(y|x_i)} & -\sum p(y|x_i) \log p_\theta(y|x_i)
 \end{aligned}$$

Cross Entropy  $\longleftrightarrow$  NLL

Consider only  $(x_i, y_i) \rightarrow$  model output is probability distribution

$$H(p, p_\theta) = -\sum_y p(y|x_i) \log p_\theta(y|x_i) - \log p_\theta(y_i|x_i)$$

if  $p(y|x_i)$  is one-hot encoded ...



$$H(p, p_\theta) = -\sum_y p(y|x_i) \log p_\theta(y|x_i)$$

$$\begin{aligned}
 &= -p(y_1|x_i) \log p_\theta(y_1|x_i) - p(y_2|x_i) \log p_\theta(y_2|x_i) \\
 &\quad \dots - p(y_n|x_i) \log p_\theta(y_n|x_i)
 \end{aligned}
 \quad \left. \begin{array}{l} \text{1개의 class} \\ y_i \text{에 해당하는 class 뿐} \\ \text{남은 나머지는 전부 0} \end{array} \right\}$$

$$= -p(y_i|x_i) \log p_\theta(y_i|x_i)$$

$$= -\log p_\theta(y_i|x_i)$$

$$= \text{NLL}$$