

Jinchao Xu

Deep Learning Algorithms and Analysis

Summer 2020

Contributors:

This set of notes are based on contributions from many of graduate students, post-doctoral fellows and other collaborators. Here is a partial list:

Juncai He, Qingguo Hong, Li Jiang.....

Contents

0.1 KL divergence, cross-entropy and logistic regression	3
--	---

0.1 KL divergence, cross-entropy and logistic regression

KL divergence and cross-entropy

Given two discrete probability distributions,

$$p = \begin{pmatrix} p_1 \\ \vdots \\ p_k \end{pmatrix}, q = \begin{pmatrix} q_1 \\ \vdots \\ q_k \end{pmatrix}$$

namely $0 \leq p_i, q_i \leq 1$ and $\sum_{i=1}^k p_i = \sum_{i=1}^k q_i = 1$. The KL divergence defines a special distance between p and q :

$$(0.1) \quad D_{\text{KL}}(q, p) = \sum_{i=1}^k q_i \log \frac{q_i}{p_i}.$$

$D_{\text{KL}}(q, p)$ works like a “distance” without the symmetry:

Lemma 1.

1. $D_{\text{KL}}(q, p) \geq 0$;
2. $D_{\text{KL}}(q, p) = 0$ if and only if $p = q$;

Proof. We first note that the elementart inequality

$$(0.2) \quad \log x \leq x - 1, \quad \text{for any } x \geq 0,$$

and the equality holds if and only if $x = 1$.

$$(0.3) \quad -D_{\text{KL}}(q, p) = -\sum_{i=1}^c q_i \log \frac{q_i}{p_i} = \sum_{i=1}^k q_i \log \frac{p_i}{q_i} \leq \sum_{i=1}^k q_i \left(\frac{p_i}{q_i} - 1 \right) = 0.$$

And the equality holds if and only if

$$(0.4) \quad \frac{p_i}{q_i} = 1 \quad \forall i = 1 : k.$$

□

Note that

$$(0.5) \quad D_{\text{KL}}(q, p) = \sum_{i=1}^k q_i \log \frac{q_i}{p_i} = \sum_{i=1}^k q_i \log q_i - \sum_{i=1}^k q_i \log p_i$$

We write

$$(0.6) \quad H(q, p) = H(q) + D_{\text{KL}}(q, p),$$

where

$$(0.7) \quad H(q) = -\sum_{i=1}^k q_i \log q_i,$$

which is called entropy for distribution p and

$$(0.8) \quad H(q, p) = -\sum_{i=1}^k q_i \log p_i.$$

which is called cross-entropy for distribution p and q .

It follows from the relation (0.6) that

$$(0.9) \quad \arg \min_p D_{\text{KL}}(q, p) = \arg \min_p H(q, p).$$

Cross-entropy

In §0.1 we introduced the concept of cross-entropy, which can be used to define a loss function in machine learning and optimization. Let us assume y_i is the true label for x_i , for example $y_i = e_{k_i}$ if $x_i \in A_{k_i}$. Then, consider the predicted distribution

$$(0.10) \quad p(x; \theta) = \frac{1}{\sum_{i=1}^k e^{w_i x + b_i}} \begin{pmatrix} e^{w_1 x + b_1} \\ e^{w_2 x + b_2} \\ \vdots \\ e^{w_k x + b_k} \end{pmatrix} = \begin{pmatrix} p_1(x; \theta) \\ p_2(x; \theta) \\ \vdots \\ p_k(x; \theta) \end{pmatrix}$$

for any data $x \in A$. By (0.9), we have

$$(0.11) \quad \arg \min_{\theta} \sum_{i=1}^N D_{\text{KL}}(y_i, \mathbf{p}(x_i; \theta)) = \arg \min_{\theta} \sum_{i=1}^N H(y_i, \mathbf{p}(x_i; \theta)),$$

Recall that we have all data $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$. Then, it is natural to consider the loss function as following:

$$(0.12) \quad \sum_{j=1}^N H(y_j, \mathbf{p}(x_j; \theta)),$$

which measures the distance between the real label and predicted one for all data. In the meantime, we can check that

$$\begin{aligned} \sum_{j=1}^N H(y_j, \mathbf{p}(x_j; \theta)) &= - \sum_{j=1}^N y_j \cdot \log p(x_j; \theta) \\ &= - \sum_{j=1}^N \log p_{i_j}(x_j; \theta) \quad (\text{because } y_j = e_{i_j} \text{ for } x_j \in A_{i_j}) \\ (0.13) \quad &= - \sum_{i=1}^k \sum_{x \in A_i} \log p_i(x; \theta) \\ &= - \log \prod_{i=1}^k \prod_{x \in A_i} p_i(x; \theta) \\ &= L(\theta) \end{aligned}$$

That is to say, the logistic regression loss function defined by likelihood is exact the loss function defined by measuring the distance between real label and predicted one via cross-entropy. Or we can note as

$$(0.14) \quad \min_{\theta} L_{\lambda}(\theta) \Leftrightarrow \min_{\theta} \sum_{j=1}^N H(y_j, \mathbf{p}(x_j; \theta)) + \lambda R(\|\theta\|) \Leftrightarrow \min_{\theta} \sum_{j=1}^N D_{\text{KL}}(y_j, \mathbf{p}(x_j; \theta)) + \lambda R(\|\theta\|).$$