
Deep Learning Assignment 3

Anonymous Author(s)

Affiliation

Address

email

1 General Questions

2 Softmax regression gradient calculation

Given

$$\hat{y} = \sigma(Wx + b), \text{ where } x \in \mathbb{R}^d, W \in \mathbb{R}^{k \times d}, b \in \mathbb{R}^k \quad (1)$$

where d is the input dimension, k is the number of classes, σ is the softmax function:

$$\sigma(a)_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)} \quad (2)$$

Which means a given input x will output y with probability of each class

2.1 Derive $\frac{\partial l}{\partial W_{ij}}$

If the given cross-entropy loss defined as followed:

$$l(y, \hat{y}) = - \sum_i y_i \log \hat{y}_i \quad (3)$$

As W_{ij} will affect the prediction of class i by multiplying index j in x , therefore we can derive:

$$\frac{\partial l}{\partial W_{ij}} = \frac{\partial l}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial W_{ij}} \quad (4)$$

where:

$$l(y, \hat{y}) = - \sum_i y_i \log \hat{y}_i = -(y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + \dots + y_i \log \hat{y}_i + \dots) \quad (5)$$

and therefore

$$\frac{\partial l}{\partial \hat{y}_i} = \frac{-y_i}{\hat{y}_i} \quad (6)$$

And we can rewrite (1) and (2) and care the value only for \hat{y}_i :

$$\hat{y}_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)} = \frac{\exp(a_i)}{C + \exp(a_i)}, \text{ where } C = \sum_{k \neq i} \exp(a_k) \quad (7)$$

Since

$$\frac{\partial \exp(a_i)}{\partial W_{ij}} = W_{ij} \exp(a_i) \quad (8)$$

Therefore

$$\frac{\partial \hat{y}_i}{\partial W_{ij}} = W_{ij} \hat{y}_i (1 - \hat{y}_i) \quad (9)$$

Combining (6) and (9) to (4), and we will get the result:

$$\frac{\partial l}{\partial W_{ij}} = \frac{\partial l}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial W_{ij}} = -X_j y_i (1 - \hat{y}_i) \quad (10)$$

15 **2.2 What happen when** $y_{c_1} = 1, \hat{y}_{c_2} = 1, c_1 \neq c_2$

16 This means something like $y = [1, 0, 0]^T$ and $\hat{y} = [0, 0, 1]^T$, and the predict is far different from true
17 lable. This will cause the log part in loss (3) become negative infinity. We may not need to worry this
18 because before one of the class predicted close to 1 and everything else close to 0, it will generate a
19 great positive loss the the class that is miss-predicted trying to make the predict right to true label.

20 **3 Chain rule**

21 **4 Variants of pooling**

22 **5 Convolution**

23 (a) As it is using 3x3 kernal along x and y axis of input, which is 5 and 5 respectively. The output of
24 this layer will be $(5 - 3 + 1) \times (5 - 3 + 1)$ which is 3x3.

25 (b) Assuming the kernel operation is point-point multiplication and summation, then the output of
26 this layer is:

27
$$\begin{pmatrix} 109 & 92 & 72 \\ 108 & 85 & 74 \\ 110 & 74 & 79 \end{pmatrix}$$

28 (c)
$$\begin{pmatrix} 4 & 7 & 10 & 6 & 3 \\ 9 & 17 & 25 & 16 & 8 \\ 11 & 23 & 34 & 23 & 11 \\ 7 & 16 & 24 & 17 & 8 \\ 2 & 6 & 9 & 7 & 3 \end{pmatrix}$$

29

30 **6 Optimization**

31 **7 Top-k error**

32 **8 t-SNE**

33 **9 Proximal gradient descent**