
Deep Learning Assignment 3

Anonymous Author(s)

Affiliation

Address

email

1 General Questions

(a)

(b)

2 Softmax regression gradient calculation

(a)

(b)

3 Chain rule

(a) Let

$$f = \frac{a}{b} \tag{1}$$

$$a = x^2 + \sigma(y) \tag{2}$$

$$b = 3x + y - \sigma(x) \tag{3}$$

$$\frac{\partial f}{\partial x} = \frac{\frac{\partial a}{\partial x} \cdot b - \frac{\partial b}{\partial x} \cdot a}{b^2} = \frac{2x \cdot b - (3 - \frac{\partial \sigma(x)}{\partial x}) \cdot a}{b^2} \tag{4}$$

$$\frac{\partial f}{\partial y} = \frac{\frac{\partial a}{\partial y} \cdot b - \frac{\partial b}{\partial y} \cdot a}{b^2} = \frac{\frac{\partial \sigma(y)}{\partial y} \cdot b - a}{b^2} \tag{5}$$

Where

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x) \cdot (1 - \sigma(x)) \tag{6}$$

(b) For $x = 1, y = 0$,

$$\sigma(x) = \sigma(1) = 0.269 \tag{7}$$

$$\sigma(y) = \sigma(0) = 0.5 \tag{8}$$

$$a = 1 + \sigma(0) = 1.5 \quad (9)$$

$$b = 3 * 1 + 0 - \sigma(1) = 2.731 \quad (10)$$

11 So we can calculate derivative of f by:

$$\frac{\partial f}{\partial x} = \frac{2 * 1 * 2.731 - (3 - 0.269 * (1 - 0.269)) * 1.5}{2.731^2} = 1.633 \quad (11)$$

$$\frac{\partial f}{\partial y} = \frac{0.5 * (1 - 0.5) * 2.731 - 1.5}{2.731^2} = -0.1096 \quad (12)$$

12 4 Variants of pooling

13 (a) SpatialMaxPooling SpatialAveragePooling SpatialAdaptivePooling

14 (b)

15 (c)

16 5 Convolution

17 (a) Assume we use zero padding and step size = 1, then: $(5 - 3 + 1) * (5 - 3 + 1) = 9$ values will
18 be generated.

19 (b) Let X be a 3x3 matrix on Image Matrix and W is the 3x3 convolution filter.
20 According to the definition of convolution, each element in the output is the point product of these
21 two matrix.

$$F = \sum W \cdot * X \quad (13)$$

22 For example: $F_11 = 4 * 4 + 3 * 5 + 3 * 2 + 5 * 3 + 5 * 3 + 5 * 2 + 2 * 4 + 4 * 3 + 3 * 4 = 109$

23 So we have the output F = $\begin{pmatrix} 109 & 92 & 72 \\ 108 & 85 & 74 \\ 110 & 74 & 79 \end{pmatrix}$

24 (c) Let $\frac{\partial E}{\partial X(i-1)} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$

25 The definition of F = $\begin{pmatrix} \sum_{i=1}^3 \sum_{j=1}^3 X_{ij} W_{ij} & \sum_{i=2}^4 \sum_{j=1}^3 X_{ij} W_{ij} & \sum_{i=3}^5 \sum_{j=1}^3 X_{ij} W_{ij} \\ \sum_{i=1}^3 \sum_{j=2}^4 X_{ij} W_{ij} & \sum_{i=2}^4 \sum_{j=2}^4 X_{ij} W_{ij} & \sum_{i=3}^5 \sum_{j=2}^4 X_{ij} W_{ij} \\ \sum_{i=1}^3 \sum_{j=3}^5 X_{ij} W_{ij} & \sum_{i=2}^4 \sum_{j=3}^5 X_{ij} W_{ij} & \sum_{i=3}^5 \sum_{j=3}^5 X_{ij} W_{ij} \end{pmatrix}$

$$\frac{\partial E}{\partial X_{ij}} = \frac{\partial}{\partial X_{ij}} \dots = \sum W_{ij} \frac{\partial E}{\partial X(i-1)} \quad (14)$$

26 result = $\begin{pmatrix} 4 & 7 & 10 & 6 & 3 \\ 9 & 17 & 25 & 16 & 8 \\ 11 & 23 & 34 & 23 & 11 \\ 7 & 16 & 24 & 17 & 8 \\ 2 & 6 & 9 & 7 & 3 \end{pmatrix}$

27 **6 Optimization**

28 **(a)**

29 **(b)**

30 **(c)**

31 **(d)**

32 **7 Top-k error**

33 Top-k error: the fraction of test images for which the correct label is not among the k labels considered
34 most probable by the model.

$$E = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{R_i > k} \quad (15)$$

35 Where function R counts the number of probability for prediction \hat{y}_c that's larger than the probability
36 of true label \hat{y}_L for each test image.

$$R = \sum_{c \in C} \mathbb{1}_{\hat{y}_c > \hat{y}_L} \quad (16)$$

37 The top-1 errors represents the error rate that the prediction is not the same with true label. This is
38 useful to compare the performance of different models. However, it cannot show how good the
39 model is in general. Sometimes the categories in a image may be ambiguous and it may be described
40 as multiple labels. The order of several possible labels is ambiguous. Using top-5 errors provides a
41 more general understanding of learning ability of the model.

42 **8 t-SNE**

43 **(a)**

44 **(b)**

45 **9 Proximal gradient decent**

46 **(a)**

47 **(b)**

48 **(c)**

49 **(d)**