
Deep Learning Assignment 3

Anonymous Author(s)

Affiliation

Address

email

1 General Questions

(a)

(b)

2 Softmax regression gradient calculation

(a) By chain rule, we have:

$$\frac{\partial l}{\partial W_{i,j}} = \frac{\partial l}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial W_{i,j}} \quad (1)$$

$$\frac{\partial l}{\partial \hat{y}_i} = \frac{\partial}{\partial \hat{y}_i} \left(- \sum_j y_j \log \hat{y}_j \right) \quad (2)$$

When $i \neq j$, the elements in the loss function are constant compared to y_i . We can just consider the element that $i = j$.

$$\frac{\partial l}{\partial \hat{y}_i} = \frac{\partial}{\partial \hat{y}_i} (-y_i \log \hat{y}_i) = \frac{-y_i}{(\ln 10) \hat{y}_i} = \frac{-y_i}{\hat{y}_i} \quad (3)$$

From assignment 1, we have:

$$(X_{out})_i = \frac{\exp(\beta(X_{in})_i)}{\sum_j \exp(\beta(X_{in})_j)} \quad (4)$$

and also

$$\frac{\partial (X_{out})_i}{\partial (X_{in})_i} = \frac{\partial}{\partial (X_{in})_i} \frac{\exp(\beta(X_{in})_i)}{\sum_j \exp(\beta(X_{in})_j)} = \beta(X_{out})_i (1 - (X_{out})_i) \quad (5)$$

So

$$\frac{\partial \hat{y}_i}{\partial W_{i,j}} = \frac{\partial}{\partial W_{i,j}} \frac{\exp(W_i x + b_i)}{\sum_k \exp(W_k x + b_k)} = x_j \hat{y}_i (1 - \hat{y}_i) \quad (6)$$

Now we can calculate $\frac{\partial l}{\partial W_{i,j}}$

$$\frac{\partial l}{\partial W_{i,j}} = \frac{\partial l}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial W_{i,j}} = \frac{-y_i}{\hat{y}_i} x_j \hat{y}_i (1 - \hat{y}_i) = -y_i x_j (1 - \hat{y}_i) \quad (7)$$

12 **(b)** When $y_{c1} = 1$ and $\hat{y}_{c2} = 1$ and $c1 \neq c2$, the model has wrong output with confidence 100%
 13 for example:
 14 model output = [1, 0]
 15 true label = [0, 1]
 16

17 From the loss function in part(a), we will get a very very large loss at $y_{c1} \log \hat{y}_{c1}$ and get 0 error in
 18 other positions. As for the gradient, since

$$\frac{\partial l}{\partial W_{i,j}} = -y_i x_j (1 - \hat{y}_i) \quad (8)$$

19 when $i = c1$, it will increase the W value that contribute to the correct label, otherwise the W value
 20 remain unchanged. With the softmax function, the \hat{y}_i for the wrong label will be reduced and the
 21 correct the \hat{y}_i will increase in next forward propagation.

22 **3 Chain rule**

23 **(a)** Let

$$f = \frac{a}{b} \quad (9)$$

$$a = x^2 + \sigma(y) \quad (10)$$

$$b = 3x + y - \sigma(x) \quad (11)$$

$$\frac{\partial f}{\partial x} = \frac{\frac{\partial a}{\partial x} \cdot b - \frac{\partial b}{\partial x} \cdot a}{b^2} = \frac{2x \cdot b - (3 - \frac{\partial \sigma(x)}{\partial x}) \cdot a}{b^2} \quad (12)$$

$$\frac{\partial f}{\partial y} = \frac{\frac{\partial a}{\partial y} \cdot b - \frac{\partial b}{\partial y} \cdot a}{b^2} = \frac{\frac{\partial \sigma(y)}{\partial y} \cdot b - a}{b^2} \quad (13)$$

24 Where

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x) \cdot (1 - \sigma(x)) \quad (14)$$

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

$$\sigma(x) = \sigma(1) = 0.731 \quad (15)$$

$$\sigma(y) = \sigma(0) = 0.5 \quad (16)$$

$$a = 1 + \sigma(0) = 1.5 \quad (17)$$

$$b = 3 * 1 + 0 - \sigma(1) = 2.269 \quad (18)$$

26 So we can calculate derivative of f by:

$$\frac{\partial f}{\partial x} = \frac{2 * 1 * 2.269 - (3 - 0.731 * (1 - 0.731)) * 1.5}{2.269^2} = 0.0646 \quad (19)$$

$$\frac{\partial f}{\partial y} = \frac{0.5 * (1 - 0.5) * 2.269 - 1.5}{2.269^2} = -0.181 \quad (20)$$

27 **4 Variants of pooling**

28 **(a)** SpatialMaxPooling SpatialAveragePooling SpatialAdaptivePooling

29 **(b)**

30 **(c)**

31 5 Convolution

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

34 **(b)** Let X be a 3×3 matrix on Image Matrix and W is the 3×3 convolution filter.
35 According to the definition of convolution, each element in the output is the point product of these
36 two matrix.

$$F = \sum W \cdot * X \quad (21)$$

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

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

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

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

$$\frac{\partial E}{\partial X_{ij}} = \frac{\partial}{\partial X_{ij}} \dots = \sum W_{ij} \frac{\partial E}{\partial X(i-1)} \quad (22)$$

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

42 6 Optimization

43 **(a)**

44 **(b)**

45 **(c)**

46 **(d)**

47 7 Top-k error

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

$$E = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{R_i > k} \quad (23)$$

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

$$R = \sum_{c \in C} \mathbb{1}_{\hat{y}_c > \hat{y}_L} \quad (24)$$

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

57 **8 t-SNE**

58 **(a)** The crowding problem happens when we try to project high dimensional dataset into lower
 59 dimensional space. The distance between records in lower dimensional space will be much smaller
 60 than in high dimensional space. In some clustering algorithms, the records will form a large group in
 61 the center and fail to produce nature clustering. t-SNE use heavy-tailed distribution to reduce this
 62 problem. It convert distances into probabilities using a Gaussian distribution in high-dimensional
 63 space. Then use a much heavier tails probability distribution in the low-dimensional map. This allows
 64 a moderate distance in the high-dimensional space to be modeled by a much larger distance in low
 65 dimension.

(b)

$$\frac{\partial C}{\partial y_i} = \frac{\partial}{\partial X_{ij}} \dots = \sum W_{ij} \frac{\partial E}{\partial X^{(i-1)}} \quad (25)$$

66 **9 Proximal gradient decent**

67 **(a)**

68 **(b)**

69 **(c)**

70 **(d)**