
DS-GA-1008: Deep Learning, Assignment 3

Jiayi Lu (jl6583)

Center for Data Science, New York University
726 Broadway, 7th floor, New York, NY 10012
jiayi.lu@nyu.edu

Abstract

1 This paper provides solutions to problems in assignment #3 of the 2016 spring
2 deep learning course (instructed by Prof. Zaid Harchaoui) at Courant Institute of
3 Mathematical Science

4 1 General Questions

5 (a)

6 If there are no non-linear activation between two linear modules, let's say $x_2 = W_1x_1 + b_1$, $x_3 =$
7 $W_2x_2 + b_2$, then using substitution we can simplify the two-layer structure by replacing the two
8 layers with a single layer module represented as $x_3 = W_2W_1x_1 + W_2b_1 + b_2$

9 (b)

10 The dictionary of sparse coding can be viewed as the weight matrix of the output layer in autoencoders,
11 however, in autoencoder, the weight matrix is learned from input data, while the sparse coding
12 dictionary is either hard-coded or learned from encoded inputs.

13 2 Softmax regression gradient calculation

14 (a)

15 Let $a_i = (Wx + b)_i$, according to multivariate chain rule,

$$\frac{\partial \ell}{\partial W_{ij}} = \sum_k \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial a_i} \frac{\partial a_i}{\partial W_{ij}} \quad (1)$$

16 We also have

$$\frac{\partial \ell}{\partial \hat{y}_k} = \frac{\partial - \sum_i y_i \log \hat{y}_i}{\partial \hat{y}_k} = -\frac{y_k}{\hat{y}_k} \quad (2)$$

(3)

17 Using the solution in Assignment 1,

$$\frac{\partial \hat{y}_k}{\partial a_i} = \begin{cases} \hat{y}_i(1 - \hat{y}_i), & k = i \\ -\hat{y}_k\hat{y}_i, & k \neq i \end{cases} \quad (4)$$

18 And $\frac{\partial a_i}{\partial W_{ij}} = x_j$, we could get:

$$\frac{\partial \ell}{\partial W_{ij}} = \sum_k \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial a_i} \frac{\partial a_i}{\partial W_{ij}} \quad (5)$$

$$= -y_i (1 - \hat{y}_i) x_j - \sum_{k \neq i} -y_k \hat{y}_i x_j \quad (6)$$

$$= x_j (\hat{y}_i - y_i) \quad (7)$$

19 **(b)**

20 In this case the loss function is inf, and the gradient is 0, x_j or $-x_j$. And since $\hat{y}_{c_2} = 1, \exp(a_j) = 0$,
21 a_j has to be $-\infty$.

22 3 Chain Rule

23 **(a)**

24 Let $m(x, y) = x^2 + \sigma(y), n(x, y) = 3x + y - \sigma(x)$, according to multivariate chain rule we have:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial m} \frac{\partial m}{\partial x} + \frac{\partial f}{\partial n} \frac{\partial n}{\partial x} \quad (8)$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial m} \frac{\partial m}{\partial y} + \frac{\partial f}{\partial n} \frac{\partial n}{\partial y}$$

25 Further, since we have:

$$\frac{\partial m}{\partial x} = \frac{\partial x^2}{\partial x} = 2x, \frac{\partial m}{\partial y} = \frac{\partial \sigma(y)}{\partial y} \quad (9)$$

$$\frac{\partial n}{\partial x} = 3 - \frac{\partial \sigma(x)}{\partial x}, \frac{\partial n}{\partial y} = 1$$

26 And $\frac{\partial f}{\partial m} = \frac{1}{n}, \frac{\partial f}{\partial n} = -\frac{m}{n^2}$ are gradOutputs of a division module with input m, n , we could compute
27 the gradient using only input and gradOutputs of given modules.

28 **(b)**

29 We compute partial derivatives described in (a):

$$\frac{\partial f}{\partial m} = \frac{1}{3x+y-\sigma(x)} = \frac{1}{3-\frac{1}{1+e^{-1}}} \quad (10)$$

$$\frac{\partial f}{\partial n} = -\frac{x^2+\sigma(y)}{(3x+y-\sigma(x))^2} = -\frac{3/2}{\left(3-\frac{1}{1+e^{-1}}\right)^2}$$

30 Combining all previous evaluations we could get gradients at $x = 1, y = 0$ as follows:

$$\frac{\partial f}{\partial x} = \frac{3+5e-e^2}{2(3+2e)^2} \approx 0.065 \quad (11)$$

$$\frac{\partial f}{\partial y} = -\frac{(1+e)(3+4e)}{4(3+2e)^2} \approx -0.181$$

31 4 Variants of Pooling

32 **(a)**

- 33 1. Max Pooling: its corresponding Torch7 implementation is SpatialMaxPooling
- 34 2. Average Pooling: its corresponding Torch module is SpatialAveragePooling
- 35 3. LP(p -norm) Pooling: in Torch it's implementation is SpatialLPPooling

36 **(b)**

37 Let P be the pooling region (a sub-matrix of input), r is the representative value of that region, we
38 have:

39 1. Max Pooling: $r = \max_{1 \leq i \leq kW, 1 \leq j \leq kH} (P_{ij})$

40 2. Average Pooling: $r = \frac{1}{n} \sum_{i=1}^{kW} \sum_{j=1}^{kH} P_{ij}$

41 3. LP Pooling: $r = \left(\sum_{i=1}^{kW} \sum_{j=1}^{kH} \|P_{ij}\|^p \right)^{1/p}$

42 **(c)**

43 One of the major reason for using pooling is to capture transition invariant features, take max pooling
44 as an example, suppose we have the two following input matrices with "L" shaped feature at different
45 locations:

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \text{ and } \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix}$$

46 After pooling on these inputs in the group of 2×2 we have:

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \text{ and } \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}$$

47 As we can see, the "L" shaped feature is captured in a smaller region, that is to say, pooling makes a
48 kind of transitional invariant down sampling.

49 **5 Convolution**

50 **(a)**

51 Assume that there are no paddings and stride = 1, we would end up with 9 (3×3) values after forward
52 propagation

53 **(b)**

The values generated are listed below,

109	92	72
108	85	74
110	74	79

Table 1: Output of convolution

54

55 **(c)**

56 Let E be the output of the layer above the convolutional layer, Y be the output matrix of the
57 convolutional layer, X be the input matrix, K be the kernel matrix, according to multivariate chain

rule, we could write $\frac{\partial E}{\partial X}$ as follows:

$$\frac{\partial E}{\partial x_{ij}} = \frac{\partial E}{\partial y_{11}} \frac{\partial y_{11}}{\partial x_{ij}} + \dots + \frac{\partial E}{\partial y_{33}} \frac{\partial y_{33}}{\partial x_{ij}} \quad (12)$$

$$= \sum_{m,n=1}^3 \frac{\partial E}{\partial y_{mn}} \frac{\partial y_{mn}}{\partial x_{ij}} \quad (13)$$

Note that $\frac{\partial E}{\partial y_{mn}} = 1, \forall m, n = 1, 2, 3$, and,

$$y_{mn} = k_{11}x_{mn} + k_{12}x_{m,n+1} + k_{13}x_{m,n+2} + \quad (14)$$

$$\dots \quad (15)$$

$$+ k_{31}x_{m+2,n} + k_{32}x_{m+2,n+1} + k_{33}x_{m+2,n+2} \quad (16)$$

Hence,

$$\frac{\partial E}{\partial x_{11}} = k_{11} \quad (17)$$

$$\frac{\partial E}{\partial x_{12}} = k_{11} + k_{12} \quad (18)$$

$$\frac{\partial E}{\partial x_{13}} = k_{11} + k_{12} + k_{13} \quad (19)$$

$$\dots \quad (20)$$

The gradient backpropagated through this convolutional layer is therefore a 5×5 matrix looks like this:

$$\frac{\partial E}{\partial X} = \begin{pmatrix} k_{11} & k_{11} + k_{12} & k_{11} + k_{12} + k_{13} & k_{12} + k_{13} & k_{13} \\ k_{11} + k_{21} & \dots & \dots & \dots & k_{13} + k_{23} \\ \vdots & & & & \vdots \\ \vdots & & & & \vdots \\ k_{31} & \dots & \dots & \dots & k_{33} \end{pmatrix} \quad (21)$$

Actually, $\frac{\partial E}{\partial X}$ can be represented as the convolution of $\frac{\partial E}{\partial Y}$ added with zero padding of width 2 and K rotated by 180 degrees:

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} - \text{convolve} - \begin{pmatrix} k_{33} & k_{32} & k_{31} \\ k_{23} & k_{22} & k_{21} \\ k_{13} & k_{12} & k_{11} \end{pmatrix} \quad (22)$$

The result of this convolution is:

$$\frac{\partial E}{\partial X} = \begin{pmatrix} 4 & 7 & 10 & 6 & 3 \\ 9 & 17 & 25 & 16 & 8 \\ 11 & 23 & 34 & 23 & 11 \\ 7 & 16 & 24 & 17 & 8 \\ 2 & 6 & 9 & 7 & 3 \end{pmatrix} \quad (23)$$

6 Optimization

(a)

According to the definition of autoencoder on Wikipedia, the reconstruction loss can be formulated as:

$$\ell(x, \hat{x}) = \|x - \hat{x}\|^2 \quad (24)$$

where $\hat{x} = \sigma_2(W_2\sigma_1(W_1x + b_1) + b_2)$

71 **(b)**

72 We can use chain rules to derive the derivatives, for simplicity we skip b_1, b_2 here by letting both of
 73 them equal to 0, let $h_1(x) = \sigma_1(W_1x)$, $h_2(x) = \sigma_2(W_2x)$:

$$\frac{\partial \ell}{\partial W_1^{(ij)}} = \sum_k \frac{\partial \ell}{\partial \hat{x}_k} \frac{\partial \hat{x}_k}{\partial h_1^{(i)}} \frac{\partial h_1^{(i)}}{\partial W_1^{(ij)}} \quad (25)$$

$$\frac{\partial \ell}{\partial W_2^{(ij)}} = \frac{\partial \ell}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial W_2^{(ij)}} \quad (26)$$

74 **(c)**

$$W_1(t+1) \leftarrow W_1(t) - \alpha \frac{\partial \ell}{\partial W_1(t)} \quad (27)$$

$$W_2(t+1) \leftarrow W_2(t) - \alpha \frac{\partial \ell}{\partial W_2(t)} \quad (28)$$

75 **(d)**

76 Let $v_X = \eta v_X + \alpha \frac{\partial \ell}{\partial X(t)}$, $\eta \in (0, 1]$, we have

$$W_1(t+1) \leftarrow W_1(t) - v_{W_1} \quad (29)$$

$$W_2(t+1) \leftarrow W_2(t) - v_{W_2} \quad (30)$$

77 **7 Top-k error**

78 Top-k error is the fraction of test images for which the correct label is not among the five labels
 79 considered most probable by the model. The motivation for ImageNet to use top-5 and top-1 error
 80 is to better evaluate model performance with the large number of class labels to predict and the
 81 ambiguity between classes. In addition to the percentage of predictions that generated the exact target
 82 class label (which might be fairly small), we also want to know the performance of the model in
 83 generating predictions "similar" to the target class.

84 **8 t-SNE**

85 **(a)**

86 The crowding problem refers to the fact that when mapping clusters embedded in high-dimensional
 87 space to 2-d space, the volume of the spherical space around a certain point increases slower in 2-d
 88 space when distance increases, therefore, if we want to allocate enough space to visualize small
 89 distances, the moderate distant points has to be placed very far away, on the other way, if moderate
 90 distance are kept closer, small distance points will be crowded together.

91 To alleviate this problem, t-SNE used two techniques:

- 92 1. t-SNE minimizes a KL-divergence from high-dimensional edge strength to low-dimensional
 93 distance, by using KL-divergence, t-SNE imposes a bigger penalty for similar datapoints
 94 that are placed far apart than dissimilar points placed closely together.
- 95 2. t-SNE employs a heavy tail distribution (student-t) in low dimensional to avoid the situation
 96 when moderately dissimilar points are repulsed together.

97 **(b)**

98 Let $d_{ij} = -\|y_i - y_j\|^2$, we first compute:

$$\frac{\partial C}{\partial y_t} = - \sum_i \sum_j p_{ij} \frac{\partial \log q_{ij}}{\partial y_t} \quad (31)$$

$$= - \sum_{i,j} p_{ij} \frac{\partial d_{ij} - \partial \log \sum_{k \neq l} e^{d_{kl}}}{\partial y_t} \quad (32)$$

$$= - \sum_{i,j} p_{ij} \frac{\partial d_{ij}}{\partial y_t} - \frac{1}{\sum_{k \neq l} e^{d_{kl}}} \frac{\partial \sum_{k \neq l} e^{d_{kl}}}{\partial y_t} \quad (33)$$

99 And since

$$\frac{\partial d_{ij}}{\partial y_t} = \begin{cases} -2(y_i - y_j), & \text{when } t = i \\ -2(y_j - y_i), & \text{when } t = j \\ 0, & \text{otherwise} \end{cases} \quad (34)$$

$$d_{ij} = d_{ji} \quad (35)$$

$$\frac{\partial \sum_{k \neq l} e^{d_{kl}}}{\partial y_t} = \frac{\partial \left(2 \sum_m e^{d_{tm}} \right)}{\partial y_t} = -4 \sum_m (y_t - y_m) e^{d_{tm}} \quad (36)$$

100 Combine the derivatives above, we get

$$\frac{\partial C}{\partial y_t} = - \sum_{i,j} p_{ij} \frac{\partial d_{ij}}{\partial y_t} - \frac{1}{\sum_{k \neq l} e^{d_{kl}}} \frac{\partial \sum_{k \neq l} e^{d_{kl}}}{\partial y_t} \quad (37)$$

$$= -2 \sum_m p_{tm} (-2(y_t - y_m)) - 4 \sum_{i,j} p_{ij} \sum_m (y_t - y_m) e^{d_{tm}} \quad (38)$$

$$= 4 \sum_m (p_{tm} - q_{tm}) (y_t - y_m) \quad (39)$$

101 **9 Proximal gradient descent**

102 **(a)**

$$\text{prox}_{h,t}(x) = \arg \min_z \frac{1}{2} \sum_i (z_i - x_i)^2 + t |z_i| \quad (40)$$

$$\frac{\partial \frac{1}{2} \sum_i (z_i - x_i)^2 + t |z_i|}{\partial z_i} = \begin{cases} z_i - x_i + t, & z_i \geq 0 \\ z_i - x_i - t, & z_i < 0 \end{cases} \quad (41)$$

103 Note that if $x_i - t < 0$, $z_i \geq 0$, or $x_i + t \geq 0$, $z_i < 0$, the derivative does not equals 0, hence the
104 minimizer of this function is 0, which is the property of soft thresholding function.

105 **(b)**

106 ISTA minimizes $\frac{1}{2} \|z - Dx\|_2^2 + t \|z\|$, note that the first term is convex differentiable, and second
107 term is convex, simple, non-differentiable, so it is a proximal gradient descent method.

108 **(c)**

109 let $f = \frac{1}{2} \|z - x\|_2^2 + th(z)$, we have

$$\partial f = z - x + t \partial h(z) \quad (42)$$

110 Since u minimizes f , we have $0 \in u - x + t \partial h(u)$, hence $\frac{x-u}{t} \in \partial h(u)$

111 **References**

- 112 [1] Lapin, M., Hein, M., & Schiele, B. (2015). *Loss Functions for Top-k Error: Analysis and Insights*. arXiv
113 preprint arXiv:1512.00486.