

# Deep Learning : Theory and Practice

## Midterm Exam

Date: March 8, 2019 (6:15PM)

### Instructions

1. This exam is open notes. However, computers, mobile phones and other handheld devices are not allowed.
2. Notation - bold symbols are vectors, capital bold symbols are matrices and regular symbols are scalars.
3. All questions even if they are multiple choice need justification. No marks will be provided for answers without justification.
4. **Online exam instructions** - Download question paper, write name and question number and the answer in A4 sheets, scan and send the answer sheets back in email to *deeplearning.cce2019@gmail.com*. **No external online resources can be accessed during the exam.** Any notes, books or slide material from the class can be used.
5. Total Duration In Class - **60 minutes**
6. Total Duration Online Exam - **90 minutes**
7. Total Marks - **50 points**

Name - .....

Organization - .....

1. The sigmoid function is defined as

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

and the ReLU function is defined as

$$\text{ReLU}(x) = \begin{cases} 0 & ; x \leq 0 \\ x & ; x > 0 \end{cases}$$

The derivative of  $\sigma(\text{ReLU}(a))$  with respect to  $a$  is ....

(a)  $\begin{cases} \sigma(a) - \sigma(a)^2 & a \geq 0 \\ 0.5 & a < 0 \end{cases}$

(b)  $\begin{cases} \sigma(a) - \sigma(a)^2 & a > 0 \\ 0 & a \leq 0 \end{cases}$

(c)  $\begin{cases} \sigma(a) - \sigma(a)^2 & a \geq 0 \\ 0 & a < 0 \end{cases}$

(d)  $\begin{cases} \sigma(a) - \sigma(a)^2 & a > 0 \\ 0.5 & a \leq 0 \end{cases}$

(e)  $\sigma(a) - \sigma(a)^2$

(Points 5)

**Solution:**

$$\sigma(\text{ReLU}(a)) = \begin{cases} 0.5 & ; a \leq 0 \\ \sigma(a) & ; a > 0 \end{cases} \quad (1)$$

$$\frac{\partial}{\partial a} \sigma(\text{ReLU}(a)) = \begin{cases} 0 & ; a < 0 \\ \sigma(a)(1 - \sigma(a)) & ; a > 0 \end{cases} \quad (2)$$

As the derivative is discontinuous at  $a = 0$ , both options (b) and (c) are accepted.

2. If A is constrained to be a symmetric matrix, the derivative  $\frac{\partial}{\partial A} \mathbf{x}^T A \mathbf{x} = \dots$

- (a)  $2\mathbf{x}\mathbf{x}^T$
- (b)  $\mathbf{x}\mathbf{x}^T$
- (c)  $\mathbf{x}\mathbf{x}^T - \text{diag}(\mathbf{x}\mathbf{x}^T)$
- (d)  $2\mathbf{x}\mathbf{x}^T - \text{diag}(\mathbf{x}\mathbf{x}^T)$

(Points 5)

**Solution:** Assuming  $\mathbf{x}$  is a  $D$  dimensional vector and A is a symmetric matrix of dimension  $D \times D$ ,

$$\mathbf{x}^T A \mathbf{x} = \sum_{m=1}^D \sum_{n=1}^D x_m A_{mn} x_n$$

$$\frac{\partial}{\partial A_{ij}} (\mathbf{x}^T A \mathbf{x}) = \begin{cases} x_i x_j + x_j x_i & \text{if } i \neq j \\ x_i x_i & \text{if } i = j \end{cases} \quad \because A_{ij} = A_{ji}$$

$$= \begin{cases} 2x_i x_j & \text{if } i \neq j \\ x_i x_i & \text{if } i = j \end{cases}$$

$$\frac{\partial}{\partial A} (\mathbf{x}^T A \mathbf{x}) = \begin{pmatrix} x_1 x_1 & 2x_1 x_2 & \dots & 2x_1 x_D \\ 2x_2 x_1 & x_2 x_2 & \dots & 2x_2 x_D \\ \vdots & \vdots & \ddots & \vdots \\ 2x_D x_1 & 2x_D x_2 & \dots & x_D x_D \end{pmatrix}$$

$$= 2 \begin{pmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_D \\ x_2 x_1 & x_2 x_2 & \dots & x_2 x_D \\ \vdots & \vdots & \ddots & \vdots \\ x_D x_1 & 2x_D x_2 & \dots & x_D x_D \end{pmatrix} - \begin{pmatrix} x_1 x_1 & 0 & \dots & 0 \\ 0 & x_2 x_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x_D x_D \end{pmatrix}$$

$$= 2\mathbf{x}\mathbf{x}^T - \text{diag}(\mathbf{x}\mathbf{x}^T)$$

Option (d) is the right answer.

3. **Neural Networks - Cost Function** - Let us define a NN with  $\{\mathbf{x}_i\}_{i=1}^N$  and  $\{\mathbf{t}_i\}_{i=1}^N$  denote the input and targets to the NN for  $N$  training data points. The task is classification with hard targets  $\mathbf{t} \in \mathcal{B}^{C \times 1}$ , where  $\mathcal{B}$  denotes boolean variable (0 or 1), and  $C$  is the number of classes. Note that every data point  $\mathbf{x}_i$  is associated with only one class label  $c_i$  where  $c_i \in \{1, 2, \dots, C\}$  classes. The output of the network is denoted as  $\mathbf{y}_i$  where  $\mathbf{y}_i = \{y_i^1, y_i^2, \dots, y_i^C\} \in \mathcal{R}^{C \times 1}$  and  $0 < y_i < 1$ . The NN cost function is defined using mean square error (MSE).

$$E_{MSE} = \sum_{i=1}^N \|\mathbf{t}_i - \mathbf{y}_i\|^2$$

Which of the following statements are true regarding the MSE loss function (justify your answer).

- (a)  $\sum_{i=1}^N [1 - y_i^{c_i}] > E_{MSE}$
- (b)  $\sum_{i=1}^N [1 - y_i^{c_i}]^2 \leq E_{MSE}$
- (c)  $\sum_{i=1}^N [1 - y_i^{c_i}]^2 \geq E_{MSE}$
- (d) None of the above

(Points 10)

**Solution:**

Option (a) **need not always be true**. Here is a counter example:

Let  $N = 1$ ,  $C = 2$ ,  $c_1 = 1$ ,  $c_2 = 0$ ,  $y_1^1 = 0.2$ ,  $y_1^2 = 0.8$ .

$$\begin{aligned} 1 - y_1^{c_1} &= 1 - 0.2 \\ &= 0.8 \\ E_{MSE} &= (1 - 0.2)^2 + (0 - 0.8)^2 \\ &= 1.28 \end{aligned}$$

Option (b) **is always true**. Proof:

$$\begin{aligned} E_{MSE} &= \sum_{i=1}^N \|\mathbf{t}_i - \mathbf{y}_i\|^2 \\ &= \sum_{i=1}^N \sum_{c=1}^C (t_i^c - y_i^c)^2 \\ &= \sum_{i=1}^N \left[ (1 - y_i^{c_i})^2 + \underbrace{\sum_{\substack{c=1 \\ c \neq c_i}}^C (-y_i^c)^2}_{\text{A non-negative quantity}} \right] \\ &\geq \sum_{i=1}^N [1 - y_i^{c_i}]^2 \end{aligned}$$

Options (c) and (d) **are false** by contradiction.

4. A binary linear classifier is trained using mean square error loss and is given by  $y = \mathbf{w}^T \mathbf{x} + b$ , where  $y > 0$  is associated with class  $C_1$  and  $y \leq 0$  with class  $C_2$ . Separately a logistic regression classifier is trained using maximum likelihood criterion and it outputs  $y = \sigma(\mathbf{w}^T \mathbf{x} + b)$  with  $y > 0.5$  associated with class  $C_1$  and  $y \leq 0.5$  associated with class  $C_2$ . Check whether the following statements are true or false (justify your choice in each case).

- (a) If the weight  $\mathbf{w}$  and the bias  $b$  are the same for both classifiers, the logistic regression performs better classification than the linear classifier.
- (b) When trained using the same training data, both the linear classifier and the logistic regression classifier have the same parameters  $\mathbf{w}$  and  $b$ .

(Points 10)

**Answer:**

(a)

$$\begin{array}{llll}
 \sigma(\mathbf{w}^T \mathbf{x} + b) = 0.5 & \iff & \mathbf{w}^T \mathbf{x} + b = 0 & \text{(Boundary condition)} \\
 \sigma(\mathbf{w}^T \mathbf{x} + b) < 0.5 & \iff & \mathbf{w}^T \mathbf{x} + b < 0 & \text{(Decided class } C_1\text{)} \\
 \underbrace{\sigma(\mathbf{w}^T \mathbf{x} + b) > 0.5} & \iff & \underbrace{\mathbf{w}^T \mathbf{x} + b > 0} & \text{(Decided class } C_2\text{)} \\
 \text{Logistic regression decision rule} & & \text{Linear classifier decision rule} & 
 \end{array}$$

If the weight  $\mathbf{w}$  and the bias  $b$  are the same for both classifiers, the decision made on a test feature vector  $\mathbf{x}$  is the same by both logistic regression and linear classifiers. Hence, statement (a) is **FALSE**.

(b) The linear classifier minimizes the mean squared error function

$$\begin{aligned}
 E_{MSE} &= \sum_{i=1}^N (y_i - t_i)^2 \\
 &= \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + b - t_i)^2
 \end{aligned}$$

and the logistic regression classifier maximizes the likelihood

$$\begin{aligned}
 L &= \prod_{i=1}^N p(t_i = 1 | \mathbf{x}_i) \\
 &= \prod_{i=1}^N (y_i)^{t_i} (1 - y_i)^{1-t_i} \\
 &= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i + b)^{t_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b))^{1-t_i}
 \end{aligned}$$

The two classifier are optimized on two different functions, which do not necessarily have the same optimal point. Hence, statement (b) is **FALSE**.

Moreover, the MSE of a linear classifier is a convex function in  $\mathbf{w}$  and  $b$ , whereas the likelihood function is non convex. It could have several local minima. When trained using the same training data, it can not be guaranteed that the linear classifier and the logistic regression classifier have the same parameters  $\mathbf{w}$  and  $b$ .

5. A deep neural network with 1 hidden layer is described as follows. The input is denoted as  $\mathbf{x}$ , and the output of layer 1 as  $\mathbf{z}^{(1)} = \phi(W^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$ , where  $\phi$  is the non-linearity applied. The output activation of the second layer is given as  $\mathbf{a}^{(2)} = (W^{(2)}\mathbf{z}^{(1)} + \mathbf{b}^{(2)})$ , and the network output is  $\mathbf{y} = \psi(\mathbf{a}^{(2)})$  where  $\psi$  can be of two types, linear activation ( $\psi(\mathbf{a}^{(2)}) = \mathbf{a}^{(2)}$ ) and soft-max activation ( $\psi(\mathbf{a}^{(2)}) = \text{softmax}(\mathbf{a}^{(2)})$ ). The error function in the case of linear activation is the mean square loss  $E_{MSE} = \|\mathbf{t} - \mathbf{y}\|^2$ , while we use a cross entropy loss  $E_{CE} = -\sum_{i=1}^C t_i \log y_i$  for the soft-max activation. Here  $\mathbf{t}$  is the target vector.

Check whether the following statement is true or false.

“The derivative of cross entropy loss  $\left(\frac{\partial E_{CE}}{\partial \mathbf{a}^{(2)}}\right)$  with softmax target activation is equivalent to derivative of MSE loss  $\left(\frac{\partial E_{MSE}}{\partial \mathbf{a}^{(2)}}\right)$  with linear target activation.”

(Points 10)

**Answer:**

For the first case, we have

$$\begin{aligned} E_{MSE} &= \|\mathbf{t} - \mathbf{y}\|^2 \\ &= (\mathbf{t} - \mathbf{y})^T (\mathbf{t} - \mathbf{y}) \end{aligned}$$

Taking derivative wrt  $\mathbf{a}^{(2)}$

$$\begin{aligned} \frac{\partial E_{MSE}}{\partial \mathbf{a}^{(2)}} &= \frac{\partial \mathbf{y}}{\partial \mathbf{a}^{(2)}} \frac{\partial E_{MSE}}{\partial \mathbf{y}} \\ &= I \{-2(\mathbf{t} - \mathbf{y})\} \quad [\text{since } \mathbf{y} = \mathbf{a}^{(2)}] \\ &= 2(\mathbf{y} - \mathbf{t}) \end{aligned}$$

For the second case,

$$\begin{aligned} E_{CE} &= -\sum_{j=1}^C t_j \log y_j \\ &= -\sum_{j=1}^C t_j \log \frac{e^{a_j^{(2)}}}{\sum_{c=1}^C e^{a_c^{(2)}}} \\ &= -\sum_{j=1}^C t_j a_j^{(2)} + \sum_{j=1}^C t_j \log \sum_{c=1}^C e^{a_c^{(2)}} \end{aligned}$$

Taking derivative wrt  $a_i^{(2)}$

$$\begin{aligned} \frac{\partial E_{CE}}{\partial a_i^{(2)}} &= -t_i + \sum_{j=1}^C t_j \underbrace{\frac{e^{a_i^{(2)}}}{\sum_{c=1}^C e^{a_c^{(2)}}}}_{\text{softmax}(a_i^{(2)})} \\ &= y_i - t_i \end{aligned}$$

$$\frac{\partial E_{CE}}{\partial \mathbf{a}^{(2)}} = \mathbf{y} - \mathbf{t}$$

Though the derivatives are equivalent in terms of the network output  $\mathbf{y}$ , the network outputs itself are different in both cases. (The first network has a linear activation and the second network has softmax activation. Hence they are not equivalent.)

6. Assume a neural network with 1 hidden layer, where all the weights and bias parameters of the model are initialized to 0 with *Tanh* as the activation function in the hidden layer and linear activation at the output layer. Check whether the following statement is true or false.

“The weight matrix in the first layer  $W^1$  does not get updated during the backpropagation”

(Points 10)

**Answer:**

The neural network can be expressed as follows:

$$\mathbf{x} \xrightarrow{W^{(1)}\mathbf{x} + \mathbf{b}^{(2)}} \mathbf{a}^{(1)} \xrightarrow{\tanh(\mathbf{a}^{(1)})} \mathbf{y}^{(1)} \xrightarrow{W^{(2)}\mathbf{y}^{(1)} + \mathbf{b}^{(2)}} \mathbf{a}^{(2)} \xrightarrow{\text{(linear activation)}} \mathbf{y}^{(2)}$$

The backpropogated derivatives of the error function wrt  $W^{(1)}$  and  $W^{(2)}$  are given by

$$\frac{\partial E}{\partial W^{(2)}} = \underbrace{\frac{\partial \mathbf{a}^{(2)}}{\partial W^{(2)}}}_{=y^{(1)T} = \mathbf{0}} \frac{\partial \mathbf{y}^{(2)}}{\partial \mathbf{a}^{(2)}} \frac{\partial E}{\partial \mathbf{y}^{(2)}}$$

$$\frac{\partial E}{\partial W^{(1)}} = \frac{\partial \mathbf{a}^{(1)}}{\partial W^{(1)}} \frac{\partial \mathbf{y}^{(1)}}{\partial \mathbf{a}^{(1)}} \underbrace{\frac{\partial \mathbf{a}^{(2)}}{\partial \mathbf{y}^{(1)}}}_{=W^{(2)}} \frac{\partial \mathbf{y}^{(2)}}{\partial \mathbf{a}^{(2)}} \frac{\partial E}{\partial \mathbf{y}^{(2)}}$$

When the weights and biases are initialized to zero, the gradients wrt  $W^{(2)}$  is zero, and hence  $W^{(2)}$  do not get updated. The derivative wrt  $W^{(1)}$  will remain zero, and hence  $W^{(1)}$  will also never get updated.

So the statement is **TRUE**.