

Hossein Kargar, Maryam Fazeli

## 1 Short Questions

### 1.1

You're training a neural network and notice that the validation error is significantly lower than the training error. Name two possible reasons for this to happen and suggest two possible solutions.

### 1.2

Why do the layers in a deep architecture need to be non-linear?

### 1.3

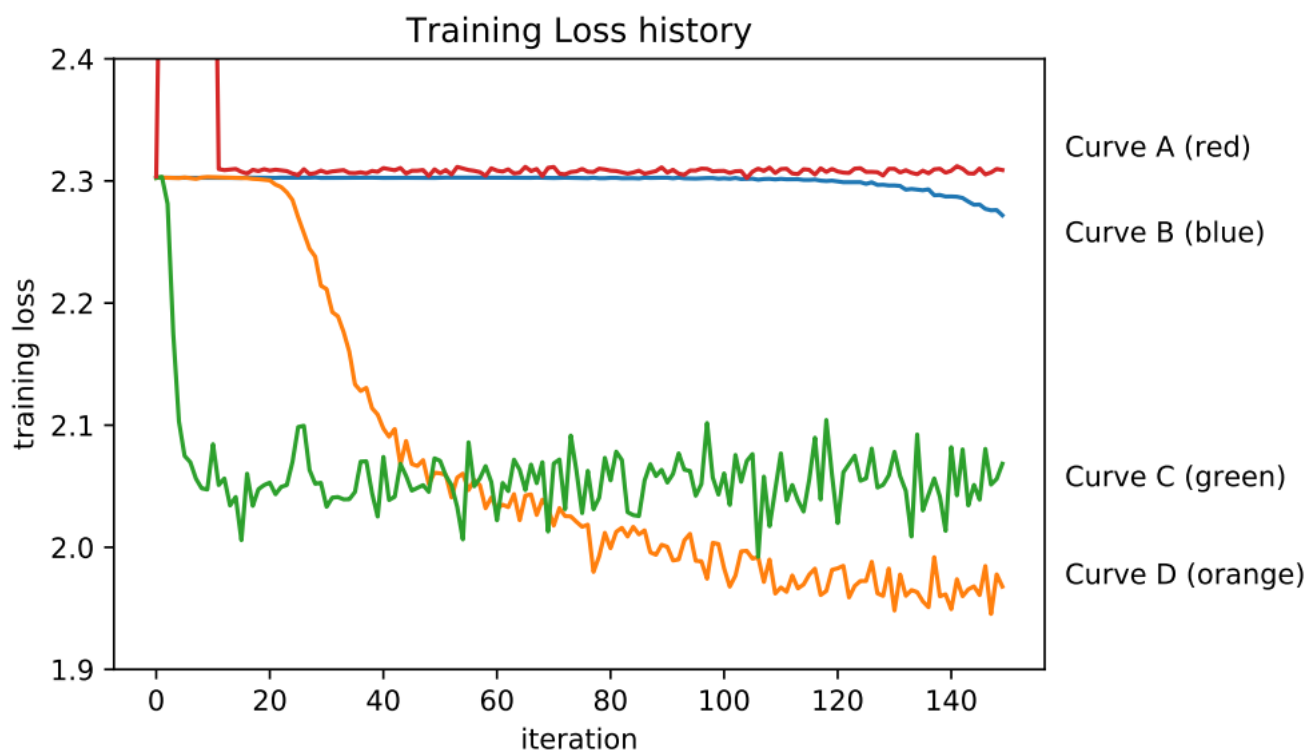
You are searching for the best learning rate ( $\alpha$ ) for your model. You decide to test the following values between 0.01 and 1:

$$\alpha \in \{0.01, 0.16, 0.21, 0.84, 0.94\}$$

Is that a good method? Explain why.

### 1.4

Your colleague trained a neural network using standard stochastic gradient descent and L2 weight regularization with four different learning rates ( $\alpha$ ) and plotted the corresponding loss curves. Unfortunately, he forgot which curve belongs to which learning rate. Please assign each of the learning rate values below to the curve (A/B/C/D) it probably belongs to and explain your thoughts.  $\alpha = [3e4, 4e1, 2e5, 8e3]$



1.5

Which of the following is true about the vanishing gradient problem?

- (i) Tanh is usually preferred over sigmoid because it doesn't suffer from vanishing gradients.
- (ii) Vanishing gradient causes deeper layers to learn more slowly than earlier layers.
- (iii) Leaky ReLU is less likely to suffer from vanishing gradients than sigmoid.
- (iv) None of the above.

1.6

Describe one advantage of using :

- (i) mini-batch gradient descent instead of full-batch gradient descent.
- (ii) mini-batch gradient descent instead of stochastic gradient descent with batch size 1.
- (iii) Adam optimizer instead of vanilla gradient descent.

## 2 Backpropagation

The softmax function has the desirable property that it outputs a probability distribution, and is often used as an activation function in many classification neural networks. Consider a 2-layer neural network for K-class classification using softmax activation and cross-entropy loss, as defined below:

$$\begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} \\ a^{[1]} &= \text{LeakyReLU}(z^{[1]}, \alpha = 0.01) \\ z^{[2]} &= W^{[2]}x + b^{[2]} \\ \hat{y} &= \text{softmax}(z^{[2]}) \\ L &= - \sum_{i=1}^K y_i \log(\hat{y}_i) \end{aligned}$$

where the model is given input  $x$  of shape  $D_x \times 1$ , and one-hot encoded label  $y \in \{0, 1\}^K$ . Assume that the hidden layer has  $D_a$  nodes, i.e.  $z^{[1]}$  is a vector of size  $D_a \times 1$ . Recall, the softmax function is computed as follows:

$$\hat{y} = [\frac{\exp(z_1^{[2]}}{Z}, \dots, \frac{\exp(z_K^{[2]}}{Z}]^T$$

where  $Z = \sum_{j=1}^K \exp(z_j^{[2]})$

2.1

What are the shapes of  $w^{[2]}, b^{[2]}$  ? If we were vectorizing across  $m$  examples, i.e. using a batch of samples  $X \in \mathbb{R}^{D_x \times m}$  as input, what would be the shape of the output of the hidden layer?

2.2

What is  $\frac{\partial \hat{y}_k}{\partial z_k^{[2]}}$  ? Simplify your answer in terms of element(s) of  $\hat{y}$ .

2.3

What is  $\frac{\partial \hat{y}_k}{\partial z_i^{[2]}}$ , for  $i \neq k$ ? Simplify your answer in terms of element(s) of  $\hat{y}$ .

2.4

Assume that the label  $y$  has 1 at its  $k^{th}$  entry, and 0 elsewhere. What is  $\frac{\partial L}{\partial z_i^{[2]}}$ ? Simplify your answer in terms of  $\hat{y}$ . (Consider both cases where  $i = k$  and  $i \neq k$ ).

2.5

What is  $\frac{\partial z^{[2]}}{\partial \alpha^{[1]}}$ ? Refer to this result as  $\delta_1$ .

2.6

What is  $\frac{\partial \alpha^{[1]}}{\partial z^{[1]}}$ ? Refer to this result as  $\delta_2$ .

2.7

Denote  $\frac{\partial L}{\partial z^{[2]}}$  with  $\delta_0$ . What is  $\frac{\partial L}{\partial W^{[1]}}$  and  $\frac{\partial L}{\partial b^{[1]}}$ ? You can reuse notations from previous parts.

2.8

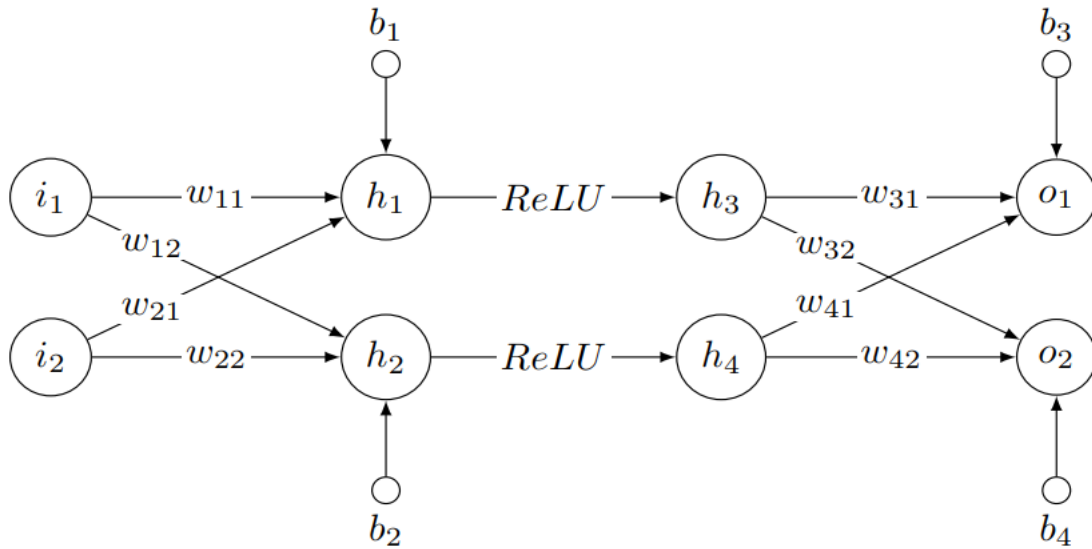
To avoid running into issues with numerical stability, one trick may be used when implementing the softmax function. Let  $m = \max_{i=1}^K z_i$  be the max of  $z_i$ , then

$$\hat{y}_i = \frac{\exp(z_i^{[2]} - m)}{\sum_{j=1}^K \exp(z_j^{[2]} - m)}$$

What is the numerical problem with the initial softmax computation? Why the modified formula would help to resolve that problem?

### 3 Two Layer Neural Network

Given the following neural network with fully connected layers and ReLU activations, including two input units ( $i_1, i_2$ ), four hidden units ( $h_1, h_2$ ) and ( $h_3, h_4$ ). The output units are indicated as ( $o_1, o_2$ ) and their targets are indicated as ( $t_1, t_2$ ). The weights and biases of the fully connected layer are called w and b with specific sub-descriptors.



The values of the variables are given in the following table:

Variable	$i_1$	$i_2$	$w_{11}$	$w_{12}$	$w_{21}$	$w_{22}$	$w_{31}$	$w_{32}$	$w_{41}$	$w_{42}$	$b_1$	$b_2$	$b_3$	$b_4$	$t_1$	$t_2$
Value	2.0	-1.0	1.0	-0.5	0.5	-1.0	0.5	-1.0	-0.5	1.0	0.5	-0.5	-1.0	0.5	1.0	0.5

3.1

Compute the output ( $o_1, o_2$ ) with the input ( $i_1, i_2$ ) and network parameters as specified above. Write down all calculations, including intermediate layer results.

3.2

Compute the mean squared error of the output ( $o_1, o_2$ ) calculated above and the target ( $t_1, t_2$ ).

### 3.3

Update the weight  $w_{21}$  using gradient descent with a learning rate of 0.1 as well as the loss computed previously. (Please write down all your computations.)