

# Q2. Nueral Networks

Menghe JIN

EPFL — January 7, 2021

## 1 Derivation of gradients

To distinguish between matrices, vectors and scalar, we denote matrices as capital letters, vectors as bold lowercase letters and scalars as normal lowercase ones. The parameters involved in this derivation are  $\mathcal{L}$ ,  $\mathbf{x}$ ,  $W^{[1]}$ ,  $\mathbf{b}^{[1]}$ ,  $\mathbf{z}^{[1]}$ ,  $\mathbf{a}^{[1]}$ ,  $W^{[2]}$ ,  $\mathbf{b}^{[2]}$ ,  $\mathbf{z}^{[2]}$ ,  $\mathbf{a}^{[2]}$ ,  $W^{[3]}$ ,  $\mathbf{b}^{[3]}$ ,  $\mathbf{z}^{[3]}$ ,  $\mathbf{a}^{[3]}$ . For sigmoid function,  $\frac{d\delta(x)}{dx} = \delta(x)(1 - \delta(x))$ .

$$\mathcal{L} = - \left[ y \cdot \log(a^{[3]}) + (1 - y) \log(1 - a^{[3]}) \right]$$

$$\frac{d\mathcal{L}}{da^{[3]}} = -\frac{y}{a^{[3]}} + \frac{1-y}{1-a^{[3]}} = \frac{a^{[3]} - y}{a^{[3]}(1-a^{[3]})}$$

$$\frac{d\mathcal{L}}{dz^{[3]}} = \frac{d\mathcal{L}}{da^{[3]}} \frac{da^{[3]}}{dz^{[3]}} = \frac{d\mathcal{L}}{da^{[3]}} \cdot a^{[3]}(1-a^{[3]})$$

$$\nabla_{W^{[3]}} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial W_1^{[3]}}, \dots, \frac{\partial \mathcal{L}}{\partial W_{h_2}^{[3]}} \right) = \frac{d\mathcal{L}}{dz^{[3]}} \left( \frac{\partial z^{[3]}}{\partial W_1^{[3]}}, \dots, \frac{\partial z^{[3]}}{\partial W_{h_2}^{[3]}} \right) = \frac{d\mathcal{L}}{dz^{[3]}} (a_1^{[2]}, \dots, a_{h_2}^{[2]}) = \frac{d\mathcal{L}}{dz^{[3]}} \cdot (\mathbf{a}^{[2]})^T$$

$$\frac{\partial \mathcal{L}}{\partial b^{[3]}} = \frac{d\mathcal{L}}{dz^{[3]}} \frac{\partial z^{[3]}}{\partial b^{[3]}} = \frac{d\mathcal{L}}{dz^{[3]}} \cdot 1$$

$$\nabla_{\mathbf{a}^{[2]}} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial a_1^{[2]}}, \dots, \frac{\partial \mathcal{L}}{\partial a_{h_2}^{[2]}} \right)^T = \frac{d\mathcal{L}}{dz^{[3]}} \left( \frac{\partial z^{[3]}}{\partial a_1^{[2]}}, \dots, \frac{\partial z^{[3]}}{\partial a_{h_2}^{[2]}} \right)^T = \frac{d\mathcal{L}}{dz^{[3]}} (W_1^{[3]}, \dots, W_{h_2}^{[3]})^T = \frac{d\mathcal{L}}{dz^{[3]}} \cdot (W^{[3]})^T$$

$$\nabla_{\mathbf{z}^{[2]}} \mathcal{L} = \nabla_{\mathbf{a}^{[2]}} \mathcal{L} \cdot \nabla_{\mathbf{z}^{[2]}} \mathbf{a}^{[2]} = \nabla_{\mathbf{a}^{[2]}} \mathcal{L} \cdot \mathbf{a}^{[2]} (1 - \mathbf{a}^{[2]}), \quad (\text{element-wise multiplication})$$

$$\nabla_{W^{[2]}} \mathcal{L} = \left[ \frac{\partial \mathcal{L}}{\partial W_{ij}^{[2]}} \right] = \left[ \frac{\partial \mathcal{L}}{\partial z_i^{[2]}} \frac{\partial z_i^{[2]}}{\partial W_{ij}^{[2]}} \right] = \left[ \frac{\partial \mathcal{L}}{\partial z_i^{[2]}} a_j^{[1]} \right] = \nabla_{\mathbf{z}^{[2]}} \mathcal{L} \cdot (\mathbf{a}^{[1]})^T$$

$$\nabla_{\mathbf{b}^{[2]}} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial b_1^{[2]}}, \dots, \frac{\partial \mathcal{L}}{\partial b_{h_2}^{[2]}} \right)^T = \nabla_{\mathbf{z}^{[2]}} \mathcal{L} \cdot \left( \frac{\partial z_1^{[2]}}{\partial b_1^{[2]}}, \dots, \frac{\partial z_{h_2}^{[2]}}{\partial b_{h_2}^{[2]}} \right)^T = \nabla_{\mathbf{z}^{[2]}} \mathcal{L}, \quad (\text{element-wise multiplication})$$

$$\nabla_{\mathbf{a}^{[1]}} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial a_1^{[1]}}, \dots, \frac{\partial \mathcal{L}}{\partial a_{h_2}^{[1]}} \right)^T = \left[ \frac{\partial z_i^{[2]}}{\partial a_j^{[1]}} \right] \cdot \nabla_{\mathbf{z}^{[2]}} \mathcal{L} = (W^{[2]})^T \cdot \nabla_{\mathbf{z}^{[2]}} \mathcal{L}$$

$$\nabla_{\mathbf{z}^{[1]}} \mathcal{L} = \nabla_{\mathbf{a}^{[1]}} \mathcal{L} \cdot \nabla_{\mathbf{z}^{[1]}} \mathbf{a}^{[1]} = \nabla_{\mathbf{a}^{[1]}} \mathcal{L} \cdot \mathbf{a}^{[1]} (1 - \mathbf{a}^{[1]}), \quad (\text{element-wise multiplication})$$

$$\nabla_{W^{[1]}} \mathcal{L} = \left[ \frac{\partial \mathcal{L}}{\partial W_{ij}^{[1]}} \right] = \left[ \frac{\partial \mathcal{L}}{\partial z_i^{[1]}} \frac{\partial z_i^{[1]}}{\partial W_{ij}^{[1]}} \right] = \left[ \frac{\partial \mathcal{L}}{\partial z_i^{[1]}} x_j \right] = \nabla_{\mathbf{z}^{[1]}} \mathcal{L} \cdot \mathbf{x}^T$$

$$\nabla_{\mathbf{b}^{[1]}} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial b_1^{[1]}}, \dots, \frac{\partial \mathcal{L}}{\partial b_{h_1}^{[1]}} \right)^T = \nabla_{\mathbf{z}^{[1]}} \mathcal{L} \cdot \left( \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}}, \dots, \frac{\partial z_{h_1}^{[1]}}{\partial b_{h_1}^{[1]}} \right)^T = \nabla_{\mathbf{z}^{[1]}} \mathcal{L}, \quad (\text{element-wise multiplication})$$

## 2 Implementation and results

The settings are: learning rate is  $1e-3$  and for each batch this learning rate multiplies with 0.95; the batch size is 16; for each batch, the program iterates 150 times. For the training part, we average the loss over all samples in the batch and we compute the accuracy on the batch as well. Figure 1 and 2 are two examples of loss and accuracy during training. The final accuracy on the test data is 0.7097.

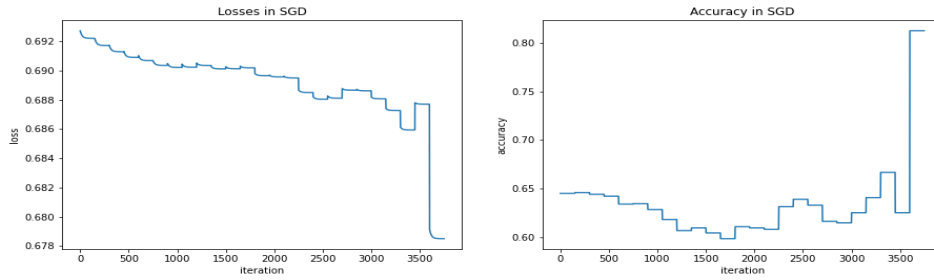


Figure 1: Loss and accuracy during training (1).

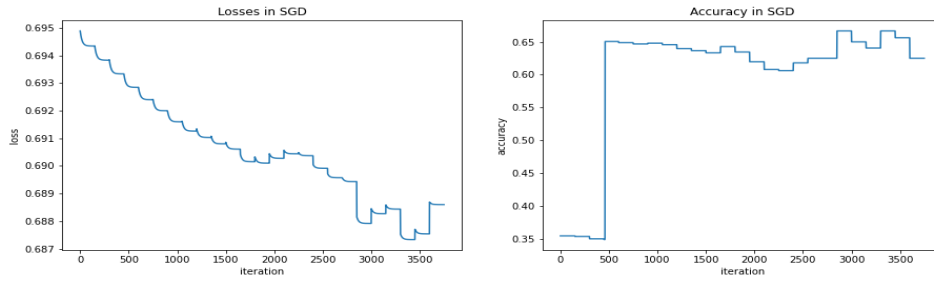


Figure 2: Loss and accuracy during training (2).