

Deep Learning: Homework 1: Backpropagation

Due on 02/25/2021

Xinyi Zhao

Problem 1**1.2 Regression Task**(a) **Answer:**

- 1 Feed data forward to the model to get the logits. Namely, put input \mathbf{x} to the model and get the corresponding output $\hat{\mathbf{y}}$.
- 2 Compute the loss using loss function and accuracy of the model. Namely, calculate $\ell_{MSE}(\hat{\mathbf{y}}, \mathbf{y})$ in this example. And calculate the ratio of the number of accurate $\hat{\mathbf{y}}$ over the number of all outputs.
- 3 Zero the gradients before running the backward pass, namely, clean up the gradient calculations, so that they are not accumulated for the next pass.
- 4 Backward pass to compute the gradient of loss w.r.t our learnable parameters. Namely, do back-propagation to calculate $\frac{\partial \ell}{\partial \mathbf{W}}$.
- 5 Use SGD to update parameters. Namely, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial \ell}{\partial \mathbf{W}}$.

(b) **Answer:**

Layer	Input	Output
$Linear_1$	\mathbf{x}	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
f	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$	$(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+$
$Linear_2$	$(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+$	$\mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$
g	$\mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$	$\mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$
Loss	$\mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$	$\ \mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)} - \mathbf{y}\ ^2$

(c) **Answer:**

Parameter	Gradient
$\mathbf{W}^{(1)}$	$\mathbf{x} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \mathbf{W}^{(2)} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$
$\mathbf{b}^{(1)}$	$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \mathbf{W}^{(2)} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$
$\mathbf{W}^{(2)}$	$(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})^+ \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$
$\mathbf{b}^{(2)}$	$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$

(d) **Answer:**Assume $\mathbf{z}_1 \in \mathbb{R}^d$, then $\mathbf{z}_2 \in \mathbb{R}^d, \mathbf{z}_3 \in \mathbb{R}^k$.

$$(1) \text{ Since } \mathbf{z}_2 = (\mathbf{z}_1)^+, \text{ then } \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} \frac{\partial z_{21}}{\partial z_{11}} & \frac{\partial z_{21}}{\partial z_{12}} & \cdots & \frac{\partial z_{21}}{\partial z_{1d}} \\ \frac{\partial z_{22}}{\partial z_{11}} & \frac{\partial z_{22}}{\partial z_{12}} & \cdots & \frac{\partial z_{22}}{\partial z_{1d}} \\ \cdots & & & \\ \frac{\partial z_{2d}}{\partial z_{11}} & \frac{\partial z_{2d}}{\partial z_{12}} & \cdots & \frac{\partial z_{2d}}{\partial z_{1d}} \end{bmatrix}.$$

$$\frac{\partial z_{2i}}{\partial z_{1j}} = \begin{cases} 1, & \text{if } i = j \text{ and } z_{1j} \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$(2) \text{ Since } \hat{\mathbf{y}} = g(\mathbf{z}_3), \text{ then } \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial z_{31}} & \frac{\partial \hat{y}_1}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_1}{\partial z_{3k}} \\ \frac{\partial \hat{y}_2}{\partial z_{31}} & \frac{\partial \hat{y}_2}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_2}{\partial z_{3k}} \\ \cdots & & & \\ \frac{\partial \hat{y}_k}{\partial z_{31}} & \frac{\partial \hat{y}_k}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_k}{\partial z_{3k}} \end{bmatrix}.$$

$$\frac{\partial \hat{y}_i}{\partial z_{3j}} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

Namely, $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} =$

$$\begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & & & & \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

$$(3) \quad \frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \begin{bmatrix} \frac{\partial \ell}{\partial \hat{y}_1} & \frac{\partial \ell}{\partial \hat{y}_2} & \dots & \frac{\partial \ell}{\partial \hat{y}_k} \end{bmatrix}, \text{ where } \frac{\partial \ell}{\partial \hat{y}_i} = 2(\hat{y}_i - y_i).$$

Problem 2

1.3 Classification Task

(a) **Answer:**

The equations of $\mathbf{z}_2, \mathbf{z}_3, \hat{\mathbf{y}}, \ell(\hat{\mathbf{y}}, \mathbf{y}), \frac{\partial \ell}{\partial \mathbf{W}^{(2)}}, \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}, \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$ need to change.

Layer	Input	Output
$Linear_1$	\mathbf{x}	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
σ	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
$Linear_2$	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$
σ	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
Loss	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$	$\ \sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}) - \mathbf{y}\ ^2$

$$\frac{\partial \ell}{\partial \mathbf{W}^{(2)}} = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$$

Assume $\mathbf{z}_1 \in \mathbb{R}^d$, then $\mathbf{z}_2 \in \mathbb{R}^d, \mathbf{z}_3 \in \mathbb{R}^k$.

$$(1) \text{ Since } \mathbf{z}_2 = \sigma(\mathbf{z}_1), \text{ then } \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} \frac{\partial z_{21}}{\partial z_{11}} & \frac{\partial z_{21}}{\partial z_{12}} & \cdots & \frac{\partial z_{21}}{\partial z_{1d}} \\ \frac{\partial z_{22}}{\partial z_{11}} & \frac{\partial z_{22}}{\partial z_{12}} & \cdots & \frac{\partial z_{22}}{\partial z_{1d}} \\ \cdots & & & \\ \frac{\partial z_{2d}}{\partial z_{11}} & \frac{\partial z_{2d}}{\partial z_{12}} & \cdots & \frac{\partial z_{2d}}{\partial z_{1d}} \end{bmatrix}.$$

$$\frac{\partial z_{2i}}{\partial z_{1j}} = \begin{cases} \frac{\exp(-z_{1i})}{(1+\exp(-z_{1i}))^2}, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

(2) Since $\hat{\mathbf{y}} = \sigma(\mathbf{z}_3)$, then $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial z_{31}} & \frac{\partial \hat{y}_1}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_1}{\partial z_{3k}} \\ \frac{\partial \hat{y}_2}{\partial z_{31}} & \frac{\partial \hat{y}_2}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_2}{\partial z_{3k}} \\ \cdots & & & \\ \frac{\partial \hat{y}_k}{\partial z_{31}} & \frac{\partial \hat{y}_k}{\partial z_{32}} & \cdots & \frac{\partial \hat{y}_k}{\partial z_{3k}} \end{bmatrix}.$

$$\frac{\partial \hat{y}_i}{\partial z_{3j}} = \begin{cases} \frac{\exp(-z_{3i})}{(1+\exp(-z_{3i}))^2}, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

(b) **Answer:**

The equations of $\ell(\hat{\mathbf{y}}, \mathbf{y})$, $\frac{\partial \ell}{\partial \hat{\mathbf{y}}}$ need to change.

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \begin{bmatrix} \frac{\partial \ell}{\partial \hat{y}_1} & \frac{\partial \ell}{\partial \hat{y}_2} & \cdots & \frac{\partial \ell}{\partial \hat{y}_k} \end{bmatrix}, \text{ where } \frac{\partial \ell}{\partial \hat{y}_i} = -\frac{1}{k} \left(\frac{y_i}{\hat{y}_i} - \frac{1-y_i}{1-\hat{y}_i} \right).$$

(c) **Answer:**

Because the use of binary activation function always leaves the network susceptible to vanishing gradients. Vanishing gradient prevents weights downstream from being modified by the neural network, which may completely stop the neural network from further training. For example, the logistic sigmoid function $\sigma(s)$ shows that when s is large, $\sigma(s)$ is 1, and when s is small, $\sigma(s)$ is 0. Because the sigmoid function is flat at $\sigma(s) = 1$ and $\sigma(s) = 0$, the gradient is 0, which results in a vanishing gradient. This tends to lose more information during training. Instead, ReLU can be used for intermediate layers to keep more useful information to improve accuracy. So in classification, we can use sigmoid function only on the output of the last layer instead of ReLU. In general, final layers of the networks usually have different kinds of activation functions to accommodate for different types of output, and ReLU is used for intermediate layers. This choice can be beneficial for training a (deeper) network.