

Deep Learning HW1

Yuanhe Guo

February 2024

1 Theory

1.1 Two-Layer Neural Nets

1.1.1 Regression Task

(a) 5 Steps:

| | |
|-------------------------------------|--|
| 1 $\tilde{y} = \text{model}(x)$; | // generate a prediction |
| 2 $L = F = C(\tilde{y}, y)$; | // compute the loss |
| 3 $\text{optimiser.zero_grad}()$; | // zero ∇params |
| 4 $L.\text{backward}()$; | // compute&accumulate ∇params |
| 5 $\text{optimizer.step}()$; | // step in towards $-\nabla \text{params}$ |

- (b) • layer1: $\text{input} = \mathbf{x}$
 $\text{output} = f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) = 3\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
- layer2: $\text{input} = 3\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
 $\text{output} = g(\mathbf{W}^{(2)}3\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}) = \mathbf{W}^{(2)}3\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$

- (c) • $\frac{\partial C}{\partial \mathbf{W}^{(2)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \frac{\partial \mathbf{s}_2}{\partial \mathbf{W}^{(2)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \mathbf{a}_1 = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} 3\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
- $\frac{\partial C}{\partial \mathbf{b}^{(2)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \frac{\partial \mathbf{s}_2}{\partial \mathbf{b}^{(2)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2}$
- $\frac{\partial C}{\partial \mathbf{W}^{(1)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \frac{\partial \mathbf{s}_2}{\partial \mathbf{a}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1} \frac{\partial \mathbf{s}_1}{\partial \mathbf{W}^{(1)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \mathbf{W}^{(2)} \frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1} \mathbf{x}$
- $\frac{\partial C}{\partial \mathbf{b}^{(1)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \frac{\partial \mathbf{s}_2}{\partial \mathbf{a}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1} \frac{\partial \mathbf{s}_1}{\partial \mathbf{b}^{(1)}} = \frac{\partial C}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} \mathbf{W}^{(2)} \frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1}$

- (d) • $\frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1} = \begin{bmatrix} \frac{\partial a_{11}^1}{\partial s_{11}^1} & \frac{\partial a_{12}^1}{\partial s_{12}^1} & \cdots & \frac{\partial a_{1j}^1}{\partial s_{1j}^1} \\ \frac{\partial a_{21}^1}{\partial s_{21}^1} & \frac{\partial a_{22}^1}{\partial s_{22}^1} & \cdots & \frac{\partial a_{2j}^1}{\partial s_{2j}^1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a_{i1}^1}{\partial s_{i1}^1} & \frac{\partial a_{i2}^1}{\partial s_{i2}^1} & \cdots & \frac{\partial a_{ij}^1}{\partial s_{ij}^1} \end{bmatrix}$

$$a_{ij}^1 = \begin{cases} 3 & \text{if } s_{ij}^1 > 0 \\ 0 & \text{if } s_{ij}^1 < 0 \end{cases}$$

$$\bullet \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$\bullet \frac{\partial C}{\partial \tilde{\mathbf{y}}} = (2\tilde{\mathbf{y}} - 2\mathbf{y})^T$$

1.1.2 Classification Task

- (a) • For (b), the equations will substitute f and g with other functions, resulting in layer1 output and layer2 input to be $\tanh(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$, layer2 output to be $(1 + \exp(-(\mathbf{W}^{(2)}\tanh(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})))^{-1}$
- For (c), we don't need to change any thing, but the value inside parameters may change accordingly.
- For (d),

$$\frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1} = \begin{bmatrix} 1 - \tanh^2(s_{11}^1) & 1 - \tanh^2(s_{12}^1) & \dots & 1 - \tanh^2(s_{1j}^1) \\ 1 - \tanh^2(s_{21}^1) & 1 - \tanh^2(s_{22}^1) & \dots & 1 - \tanh^2(s_{2j}^1) \\ \vdots & \vdots & \ddots & \vdots \\ 1 - \tanh^2(s_{i1}^1) & 1 - \tanh^2(s_{i2}^1) & \dots & 1 - \tanh^2(s_{ij}^1) \end{bmatrix}$$

$$\frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2} = \begin{bmatrix} \sigma(s_{11}^2) \cdot (1 - \sigma(s_{11}^2)) & \sigma(s_{12}^2) \cdot (1 - \sigma(s_{12}^2)) & \dots & \sigma(s_{1j}^2) \cdot (1 - \sigma(s_{1j}^2)) \\ \sigma(s_{21}^2) \cdot (1 - \sigma(s_{21}^2)) & \sigma(s_{22}^2) \cdot (1 - \sigma(s_{22}^2)) & \dots & \sigma(s_{2j}^2) \cdot (1 - \sigma(s_{2j}^2)) \\ \vdots & \vdots & \ddots & \vdots \\ \sigma(s_{i1}^2) \cdot (1 - \sigma(s_{i1}^2)) & \sigma(s_{i2}^2) \cdot (1 - \sigma(s_{i2}^2)) & \dots & \sigma(s_{ij}^2) \cdot (1 - \sigma(s_{ij}^2)) \end{bmatrix}$$

$\frac{\partial C}{\partial \tilde{\mathbf{y}}}$ will not change.

- (b) (a) (b) and (c) will be the same as 1.2(a).

- (b) For (d), $\frac{\partial \mathbf{a}_1}{\partial \mathbf{s}_1}, \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{s}_2}$ will not change.

$$\frac{\partial C}{\partial \tilde{\mathbf{y}}} = -\frac{1}{K} \left(\frac{\mathbf{y}}{\tilde{\mathbf{y}}} - \frac{1-\mathbf{y}}{1-\tilde{\mathbf{y}}} \right)$$

- (c) Compared with \tanh , ReLU function is easier to derive while maintaining non-linearity. When training a deep network, doing back propagation will be faster. Meanwhile, since ReLU won't saturate correct predictions, vanishing gradient is less likely to happen in a deep network.

1.2 Conceptual Questions

- (a) Softmax function outputs a probability distribution controlled by variable β aka "coldness". When $\beta \rightarrow \infty$, the output will become an argmax function. Therefore softmax is actually softargmax.
- (b) See fig. 1

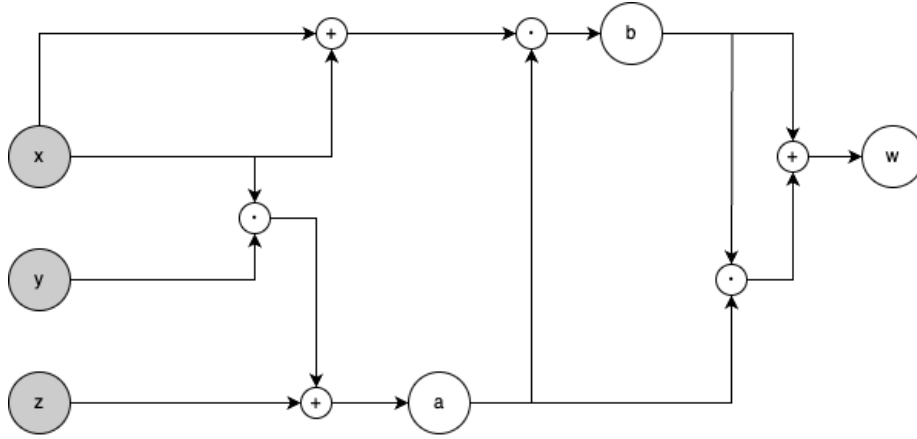


Figure 1: Computational Graph

(c) See fig. 2

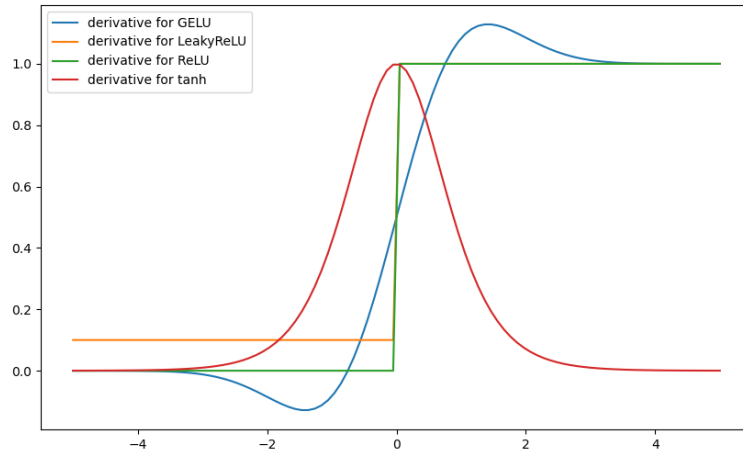


Figure 2: Derivatives for GELU, LeakyReLU, ReLU, tanh

- (d) (a) $J_f = \frac{\partial \mathbf{W}_1 \mathbf{x}}{\partial \mathbf{x}} = \mathbf{W}_1$
 $J_g = \frac{\partial \mathbf{W}_2 \mathbf{x}}{\partial \mathbf{x}} = \mathbf{W}_2$
 (b) $J_h = \frac{\partial f(\mathbf{x}) + g(\mathbf{x})}{\partial \mathbf{x}} = J_f + J_g = \mathbf{W}_1 + \mathbf{W}_2$
 (c) $J_h = \mathbf{W}_1 + \mathbf{W}_2 = 2\mathbf{W}_1 = 2\mathbf{W}_2$

- (e) (a) $J_f = \frac{\partial \mathbf{W}_1 \mathbf{x}}{\partial \mathbf{x}} = \mathbf{W}_1$
 $J_g = \frac{\partial \mathbf{W}_2 \mathbf{x}}{\partial \mathbf{x}} = \mathbf{W}_2$
(b) $J_h = \frac{\partial g(f(\mathbf{x}))}{\partial f(\mathbf{x})} \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{W}_2 \mathbf{W}_1$
(c) $J_h = \mathbf{W}_2 \mathbf{W}_1 = \mathbf{W}_1^2 = \mathbf{W}_2^2$

1.3 Deriving Loss Functions

1. For Perceptron, \tilde{y} is either -1 or 1 . Notice that $(y - \tilde{y}x_i)$ is always a non-positive number. By add a negative sign and summing up all cases, we get the loss function $-(y - \tilde{y})\sum_{i=1}^d w_i x_i$
2. For Adaline / Least Mean Squares, we want to minimize the distance between \mathbf{y} and $\tilde{\mathbf{y}}$, which is $\|y - \tilde{y}\|_2$. For sake of mathematical convenience, we set the loss function to be $\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{y}})^2$.
3. For Logistic Regression, we want to maximize the case where $\tanh(b + w_i x_i) = y_i$. Since y_i is either 1 or -1 , and \tanh function is odd, we know the previous function is the same as

$$\begin{aligned} \tanh(y_i(b + w_i x_i)) &= 1 \\ \frac{2}{1 + \exp(-y_i(b + w_i x_i))} &= 1 \\ 1 + \exp(-y_i(b + w_i x_i)) &= 1 \\ \log(1 + \exp(-y_i(b + w_i x_i))) &= 0 \end{aligned}$$

Here we can also ignore b since it won't affect the gradient. We sum up all cases, and multiply it by -2 for mathematical convenience, we get the loss function $-2\log(1 + \exp(-y\sum_{i=1}^d w_i x_i))$.