

Homework 1: Theory

Yuxiang Chai

September 27, 2022

1.1 Two-Layer Neural Nets

1.2 Regression Task

- (a) i. $\hat{\mathbf{y}} = \text{model}(\mathbf{x})$. First we need to generate the prediction of the given model and input.
ii. $\ell = \ell_{\text{MSE}} = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$. Second we compute the loss between the prediction and the ground truth.
iii. `optimizer.zero_grad()`. Then we need to set the optimizer with zero gradient. The optimizer can be SGD, Adam or others.
iv. `ℓ.backward()`. Then we compute the gradient.
v. `optimizer.step()`. Last we update the weight of each parameter.
- (b) i. For Linear_1 , the input is \mathbf{x} , and the output is $\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$.
ii. For f , the input is $\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$, and the output is $5 \cdot (\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+$.
iii. For Linear_2 , the input is $5 \cdot (\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+$, and the output is $5 \cdot \mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$.
iv. For g , the input is $5 \cdot \mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$, and the output is $5 \cdot \mathbf{W}^{(2)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+ + \mathbf{b}^{(2)}$.

- (c) i.

$$\frac{\partial \ell}{\partial \mathbf{b}^{(2)}} = \frac{\partial \ell}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{b}^{(2)}} = \frac{\partial \ell}{\partial \mathbf{z}_3} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$$

- ii.

$$\frac{\partial \ell}{\partial \mathbf{W}^{(2)}} = \frac{\partial \ell}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{W}^{(2)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{W}^{(2)}} = 5 \cdot \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} ((\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})^+)^{\top}$$

- iii.

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

- iv.

$$\frac{\partial \ell}{\partial \mathbf{W}^{(1)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}^{(1)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \mathbf{W}^{(2)} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \mathbf{x}^{\top}$$

- (d) i. $\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$ is a $m \times m$ matrix \mathbf{M} if \mathbf{z}_1 and \mathbf{z}_2 is \mathbb{R}^m , where $\mathbf{M}_{ij} = 0$ if $i \neq j$ and $\mathbf{M}_{ii} = 5$ if $\mathbf{z}_1[i] > 0$. Assume that $\text{ReLU}()$ gradient at 0 is 0.
ii. $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$ is a $K \times K$ identity matrix \mathbf{I} .
iii. $\frac{\partial \ell}{\partial \hat{\mathbf{y}}}$ is a K dimension row vector \mathbf{V} , where $\mathbf{V} = 2(\hat{\mathbf{y}} - \mathbf{y})^{\top}$.

1.3 Classification Task

- (a) (b) i. For Linear_1 , the input is \mathbf{x} , and the output is $\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$.

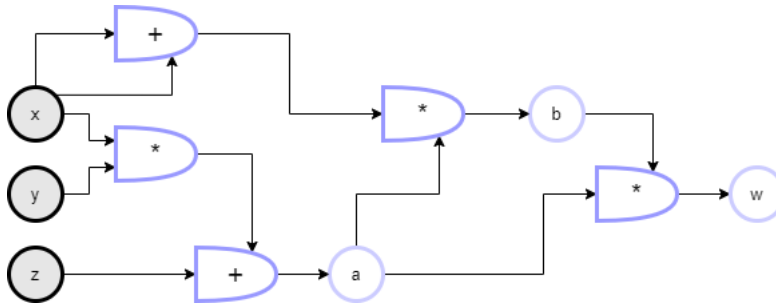
- (c) First, using ReLU can prevent vanishing gradient problem in sigmoid or tanh because of the gradient of ReLU is 1 or 0.

And using ReLU is faster than using sigmoid or tanh. As we can see in the previous parts, the gradient computation of ReLU is easier than those of sigmoid or tanh.

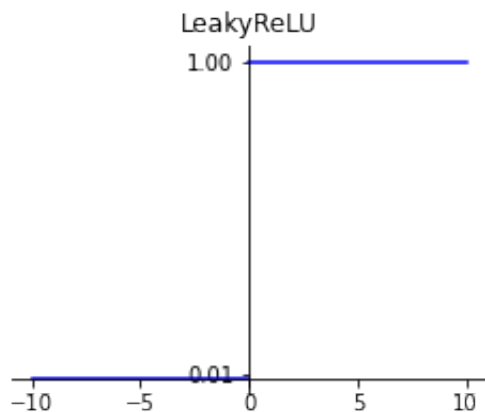
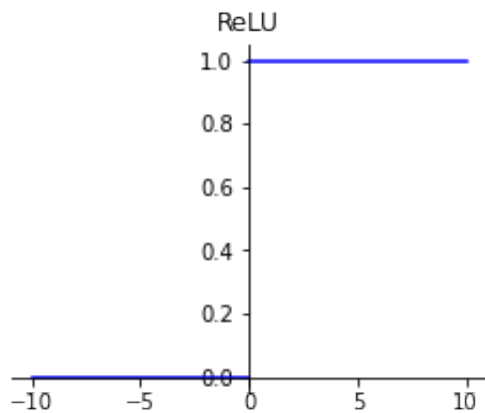
1.4 Conceptual Questions

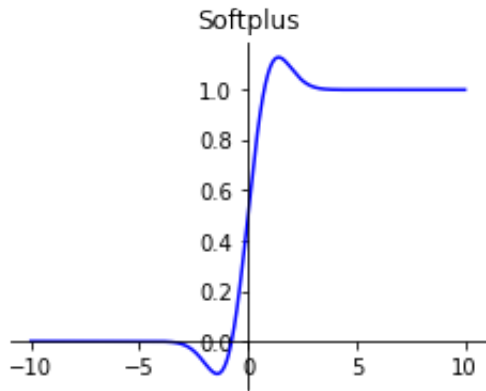
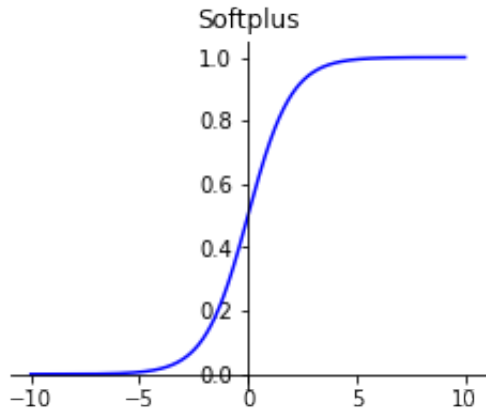
- (a) Because the definition of pytorch softmax is $z_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$, where the denominator for every element is the same. So the larger the numerator the larger z_i . So softmax is actually choosing the index that has the largest probability.

- (b)



- (c)





(d) Four types: Scaling, Rotation, Translation and Shearing.

Linear transformation can separate the current space.

Non-linear transformation is to prevent the model from becoming a linear transformation and non-linearly separate space.

(e) For every iteration, do the following:

- i. $\hat{Y} = F_{\theta}(D)$, where $\hat{Y} = \hat{y}_1, \hat{y}_2, \dots, \hat{y}_N$
- ii. $\ell = \ell_{\text{MSE}} = \|\hat{Y} - Y\|^2$
- iii. Compute gradients of θ , $\frac{\partial \ell}{\partial \theta}$
- iv. Update θ , $\theta = \theta - \eta \frac{\partial \ell}{\partial \theta}$