

314652021 羅瑋翔

Assignment 2

1. Read [Deep Learning: An Introduction for Applied Mathematicians](#). Consider a network as defined in (3.1) and (3.2). Assume that $n_L = 1$, find an algorithm to calculate $\nabla a^{[L]}(x)$.
2. There are unanswered questions during the lecture, and there are likely more questions we haven't covered. Take a moment to think about them and write them down here.

$$(3.1) \quad a^{[1]} = x \in \mathbb{R}^{n_1}$$

$$(3.2) \quad a^{[i]} = \sigma(W^{[i]} a^{[i-1]} + b^{[i]}) \in \mathbb{R}^{n_i}, i = 1, 2, \dots, L$$

$$1. \text{ Define } z^{[i]} = W^{[i]} a^{[i-1]} + b^{[i]}, \quad a^{[i]} = \sigma(z^{[i]}).$$

$$\text{Then, } \frac{\partial a^{[i]}}{\partial a^{[i-1]}} = \frac{\partial a^{[i]}}{\partial z^{[i]}} \cdot \frac{\partial z^{[i]}}{\partial a^{[i-1]}} = \sigma'(z^{[i]}) \cdot W^{[i]}$$

$$\text{and } \sigma'(z^{[i]}) = \begin{bmatrix} \sigma'(z_1^{[i]}) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma'(z_{n_i}^{[i]}) \end{bmatrix} \in \mathbb{R}^{n_i \times n_i}$$

$$\therefore \sigma(z_i^{[i]}) \text{ relies on } i\text{-th input, and } W \in \mathbb{R}^{n_i \times n_{i-1}}$$

\Rightarrow Hence, we have

$$\begin{aligned} \frac{\partial a^{[L]}}{\partial a^{[1]}} &= \frac{\partial a^{[L]}}{\partial a^{[L-1]}} \cdot \frac{\partial a^{[L-1]}}{\partial a^{[L-2]}} \cdot \dots \cdot \frac{\partial a^{[2]}}{\partial a^{[1]}} \\ &= (\sigma'(z^{[L]}) \cdot W^{[L]}) \cdot (\sigma'(z^{[L-1]}) \cdot W^{[L-1]}) \dots (\sigma'(z^{[2]}) \cdot W^{[2]}) \\ &\in \mathbb{R}^{n_L \times n_1} \end{aligned}$$

★ Since we need to obtain $z^{[i]}$ and $a^{[i]}$ to calculate $\nabla a^{[L]}(x)$, we can use "Forward passing and back propagation" to solve this

Algo:

```
1   $a^{[1]} = x$ 
2  for  $i = 2$  to  $L$ 
3       $z^{[i]} = W^{[i]} a^{[i-1]} + b^{[i]}$ 
4       $a^{[i]} = \sigma(z^{[i]})$  //  $\sigma$  is the activation function.
5       $S^{[i]} = I_{n_i}^T \cdot \text{diag } \sigma'(z^{[i]})$  // saved this for back propagation
6   $p^{[L]} = a^{[L]}$ 
7  for  $i = L$  to  $2$ .
8       $p^{[i-1]} = (p^{[i]} \odot S^{[i]}) W^{[i]}$  //  $\odot$  is elementwise product
```

\Rightarrow By this, we obtain $p^{[1]} (\nabla a^{[L]}(x)) \neq$.

2. During the class, it says that using a tanh neural network with $3n+2$ neurons in the layer and $2n+1$ neurons in the output layers, we can approximate the functions $x^p = 0, \dots, 2n$ to any desired accuracy.

\Rightarrow The question is, if the domain of x become larger, such as $x \in [-10, 10]$, is the lemma still holds?

And is tanh still an effective way to approximate x^p ?

Is there exists any way to approximate x^p for $x \in (-\infty, \infty)$?

Assignment 2

The definition of the neural network

$$f \in FC(W = 64, L = 2, d_{in} = 1, d_{out} = 1)$$

- The choice of activation function is $\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$. (Since it performs a good approximation of polynomials on the interval $x \in [-1, 1]$)
- I initialize W (the weight) with Glorot/Xavier initialization. (Maintain a moderate variance for the outputs of each layer.) And set the bias $b^{[l]}$ to $0_{1 \times n_l}$.
- Then, I use **Adam optimizer** to maintain a good choice of learning rate α . (Since adam optimizer can choose suitable momentum and velocity overtime, it might better approaches minimum than SGD.)

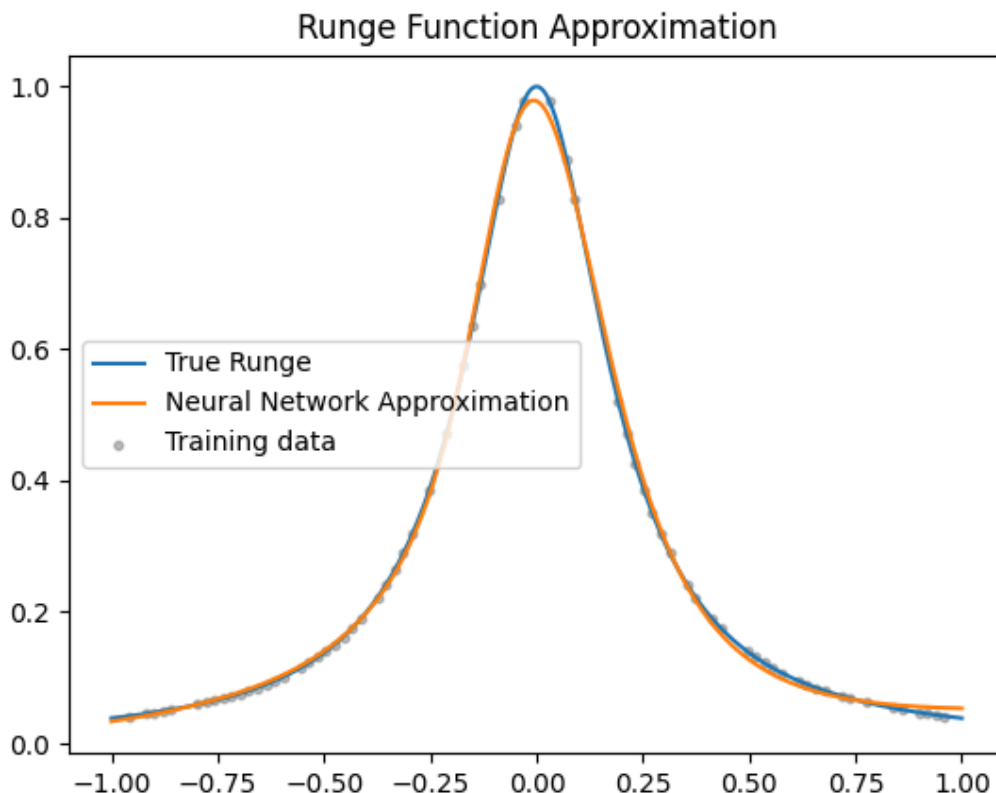


Figure 1. Runge Function Approximation with Neural Network

One can observe is that when x approaches 0, it deviates away.

- Xavier initialization would let most of the inputs not lies near 0
- And N might not be big enough to spread on the interval.

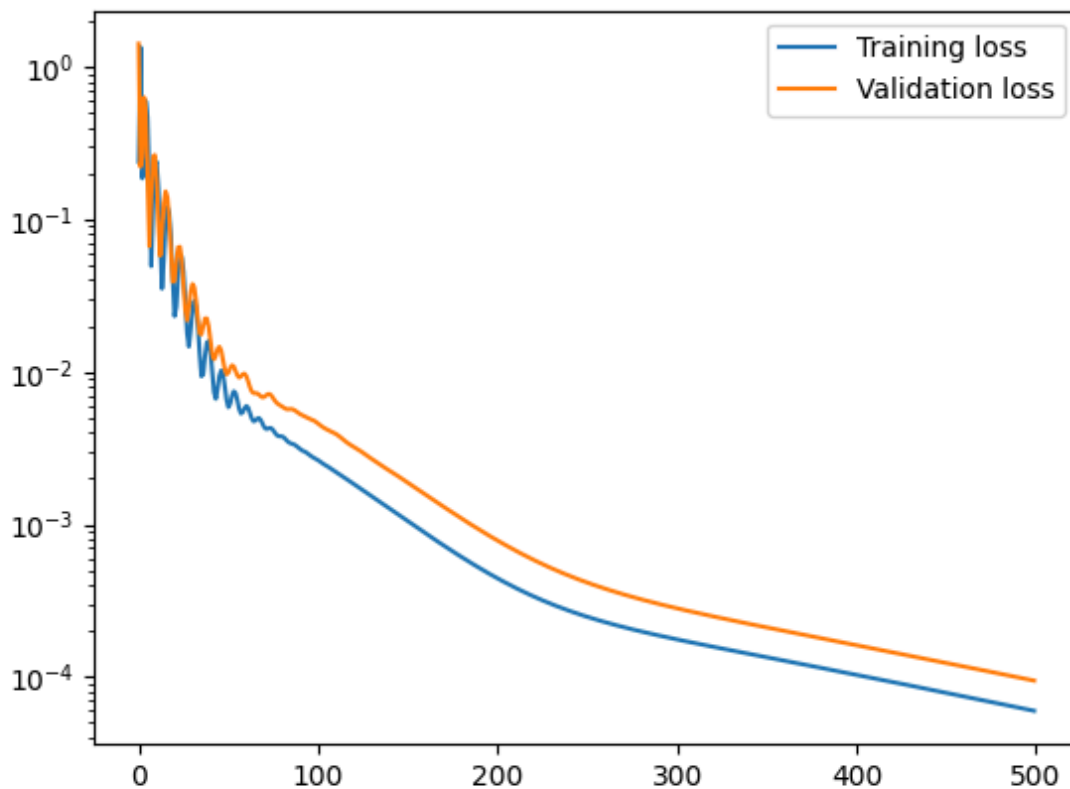


Figure 2. Training and Validation Loss Curves vs. Epoch

At the start of training, the loss is higher and oscillating. Since the model just started training and has not learned the characteristics of the function. As the epoch grows, the model converges gradually, and the loss decreases significantly. At around epoch 500, the validation loss becomes around 10^{-4} .