

## HW3

### 1. Lemma:

#### Lemma 3.1

Let  $k \in \mathbb{N}_0$  and  $s \in 2\mathbb{N} - 1$ . Then it holds that for all  $\varepsilon > 0$  there exists a shallow tanh neural network  $\Psi_{s,\varepsilon} : [-M, M] \rightarrow \mathbb{R}^{\frac{s+1}{2}}$  of width  $\frac{s+1}{2}$  such that

$$\max_{\substack{p \leq s, \\ p \text{ odd}}} \left\| f_p - (\Psi_{s,\varepsilon})_{\frac{p+1}{2}} \right\|_{W^{k,\infty}} \leq \varepsilon, \quad (17)$$

Moreover, the weights of  $\Psi_{s,\varepsilon}$  scale as  $O\left(\varepsilon^{-s/2} \left(2(s+2)\sqrt{2M}\right)^{s(s+3)/2}\right)$  for small  $\varepsilon$  and large  $s$ .

Introduction: 任意選一個小於  $s$  且大於 0 的奇數  $p$ , 有一個 neural network  $\Psi_{s,\varepsilon} : \mathbb{R}^{\frac{s+1}{2}} \rightarrow \mathbb{R}$ , 其第  $\frac{p+1}{2}$  個輸出值會近似一個 monomial function  $f_p$ .

Why  $\Psi_{s,\varepsilon} : [-M, M] \rightarrow \mathbb{R}^{\frac{s+1}{2}}$ :

因我們要逼近所有  $f_p$  ( $p \leq s$  且  $p$  is odd), 附帶條件的  $p$  有  $\frac{s+1}{2}$  個, 所以

以  $\Psi_{s,\varepsilon}$  會 output 一個  $\frac{s+1}{2}$ -dimension 之向量, 且第  $i$  項將近似  $f_i$  之值  $i \in [1, \dots, \frac{s+1}{2}]$

Why width =  $\frac{s+1}{2}$ :

首先  $f_p = x^p$ , 且  $p$  是 odd,  $f_p$  是奇函數

而  $\tanh(x) = x - \frac{x^3}{3} + \frac{3x^5}{15} - \frac{17x^7}{315} \dots = \sum_{n=1}^{\infty} \frac{2^{2n}(2^{2n}-1)B_{2n}x^{2n}}{(2n)!}$  where  $B_n(n) = \sum_{k=0}^n \binom{n}{k} (-1)^k \left(\frac{k}{2}\right)^{n-1} \frac{(n!)^{n-1}}{1+k}$

By Taylor series,  $\tanh(x)$  只含奇數次方的項數  $\rightarrow$  也是奇函數

為了逼近  $f_p$ , 我們至少取  $p$ -th Taylor polynomial of  $\tanh$  為訓練的函數, 稱其為  $\phi$

由  $\phi$ ,  $\phi$  只含  $\frac{s+1}{2}$  個項數, 靠調整各項之權重取得近似  $f_p$  值

所以 hidden layer 要有  $\frac{s+1}{2}$  個 neurons  $\rightarrow$  width =  $\frac{s+1}{2}$

What is  $\max_{\substack{p \leq s, \\ p \text{ odd}}} \|f_p - (\Psi_{s,\varepsilon})_{\frac{p+1}{2}}\|_{W^{k,\infty}} \leq \varepsilon$ :

讓兩個函數在 Sobolev Spaces 上的距離, 當任意選擇一個小於  $s$  的奇數  $p$ , 皆小於  $\varepsilon$

The weight of  $\Psi_{s,\varepsilon}$  scale as  $O\left(\varepsilon^{-s/2} \left(2(s+2)\sqrt{2M}\right)^{s(s+3)/2}\right)$  for small  $\varepsilon$  and large  $s$ :

This shows the cost of approximation: as  $\varepsilon \rightarrow 0$  (very small error), weights blow up like  $\varepsilon^{-s/2}$

As  $s$  increases, the blow-up is even more severe.

### Lemma 3.2

Let  $k \in \mathbb{N}_0$ ,  $s \in 2\mathbb{N} - 1$  and  $M > 0$ . For every  $\varepsilon > 0$ , there exists a shallow tanh neural network  $\psi_{s,\varepsilon} : [-M, M] \rightarrow \mathbb{R}^s$  of width  $\frac{3(s+1)}{2}$  such that

$$\max_{p \leq s} \|f_p - (\psi_{s,\varepsilon})_p\|_{W^{k,\infty}} \leq \varepsilon. \quad (26)$$

Furthermore, the weights scale as  $O\left(\varepsilon^{-s/2} \left(\sqrt{M}(s+2)\right)^{3s(s+3)/2}\right)$  for small  $\varepsilon$  and large  $s$ .

Introduction: 任意選一個小於 $s$ 且大於0的數  $p$ , 有一個 neural network  $\psi_{s,\varepsilon} : \mathbb{R}^s \rightarrow \mathbb{R}$ , 其第  $p$  個輸出值會近似一個 monomial function  $f_p$ .

Why  $\psi_{s,\varepsilon} : [-M, M] \rightarrow \mathbb{R}^s$ :

因我們要逼近所有  $f_p (p \leq s)$ , 總共  $s$  個, 所以

以  $\psi_{s,\varepsilon}$  會 output 一個  $s$ -dimension 之向量, 且第  $i$  項將近似  $f_i$  之值

Why width  $\frac{3(s+1)}{2}$ :

若  $p$  是奇數, 在 lemma 3.1 中, 可直接用 tanh 逼近 (需大約 2 neurons)

而  $p-1$  是偶數:

因 tanh 的展開式只包含奇數次方, 至少取  $p$  次 Taylor polynomial of tanh

我們可用奇函數點對稱性質將最高次的  $x^p$  消掉, 例  $\tanh(x+a)$  和  $\tanh(x-a)$

$\Rightarrow$  偶數次方的項數留下, 而我們訓練這偶函數, 使其逼近  $f_p$

因可重複用奇數的 neural  $\Rightarrow$  需約 1 neuron

$\Rightarrow$  所以一對奇偶函數平均需 3 neurons

因  $s$  是奇數  $\Rightarrow \frac{s+1}{2} \times 3 \rightarrow$  平均一對奇偶需 3 neurons

Why  $\max_{p \leq s} \|f_p - (\psi_{s,\varepsilon})_p\|_{W^{k,\infty}} \leq \varepsilon$ :

表示確定每個  $(\psi_{s,\varepsilon})_p$  逼近  $f_p$ , 在 Sobolev space 中, 所有  $p \leq s$ ,  $(\psi_{s,\varepsilon})_p$  和  $f_p$  的距離皆小於

一開始任意給定的值  $\varepsilon$ .

The weight scale as  $O\left(\varepsilon^{-s/2} \left(\sqrt{M}(s+2)\right)^{3s(s+3)/2}\right)$  for small  $\varepsilon$  and large  $s$ :

當  $\varepsilon \rightarrow 0$  或  $s \rightarrow \infty$  時, 權重會增加

## 2. Program:

### (1) Introduction:

Use a neural network to approximate both the Runge function and its derivative.

### (2) Method:

(a) Dataset: Random uniform sampling in  $[-1, 1]$ . Training: Validation = 80:20.

(b) Model architecture:

- a. Numbers of layer: 4
- b. Hidden size: [64, 64]
- c. Activation: tanh

(c) Training:

a. Loss function:

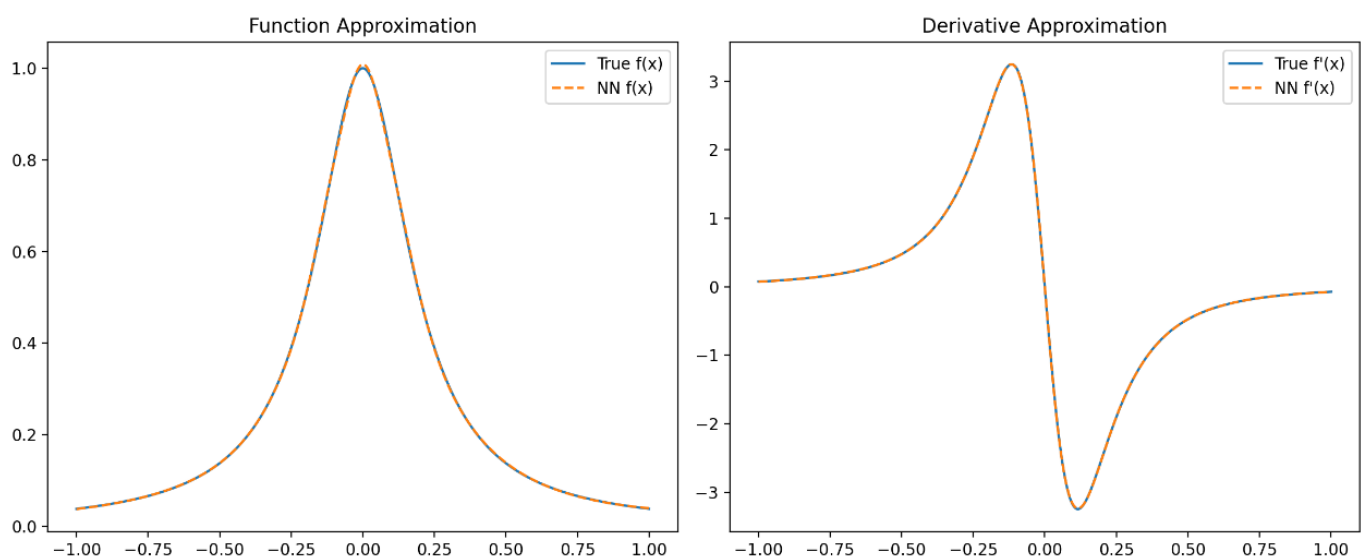
$$\text{loss function: } \frac{1}{N} \sum_{i=1}^N \|f_{\text{true}}(x_i) - f_{\text{predict}}(x_i)\|^2 + \|f'_{\text{true}}(x_i) - f'_{\text{predict}}(x_i)\|^2$$

b. Optimizer: Adam optimizer

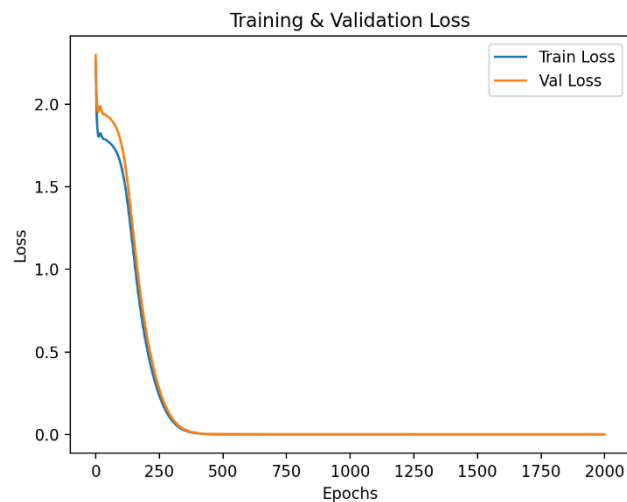
c. Learning rate:  $1e-3$

### (3) Result:

(a) The true function and the neural network prediction:



(b) The training/Validation loss curves:



(c) Compute and report errors (MSE or max error):

```
Epoch 1/2000 train MSE=2.254505 val MSE=2.297833
Epoch 200/2000 train MSE=0.541207 val MSE=0.618725
Epoch 400/2000 train MSE=0.006625 val MSE=0.008052
Epoch 600/2000 train MSE=0.000185 val MSE=0.000195
Epoch 800/2000 train MSE=0.000073 val MSE=0.000072
Epoch 1000/2000 train MSE=0.000054 val MSE=0.000054
Epoch 1200/2000 train MSE=0.000044 val MSE=0.000044
Epoch 1400/2000 train MSE=0.000036 val MSE=0.000037
Epoch 1600/2000 train MSE=0.000030 val MSE=0.000031
Epoch 1800/2000 train MSE=0.000024 val MSE=0.000025
Epoch 2000/2000 train MSE=0.000020 val MSE=0.000020

--- Error Metrics ---
Function MSE: 6.752685e-06, Max Error: 0.009458
Derivative MSE: 1.309726e-05, Max Error: 0.013237
```

(4) Discussion:

The result has shown that the neural network has completely approximate both function and its derivative. In the first picture, the prediction function curve is similar to the true function curve, and also there is no obvious difference between the prediction derivative function curve and the true derivative function curve. In the loss curve, since the validation loss closely tracked the training loss, there is no obvious overfitting. In the last picture,  $MSE=1.309726e-05$ , indicating that the network achieves a good fit.

(5) Code:

```
1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from sklearn.model_selection import train_test_split
7  from sklearn.metrics import mean_squared_error
8
9  # --- Runge function and derivative ---
10 def runge(x):
11     return 1.0/(1.0+25.0*x**2)
12
13 def Drunge(x):
14     return (-50.0*x)/((1.0+25.0*x**2)**2)
15
16 # --- Neural network ---
17 class MLP(nn.Module):
18     def __init__(self, hidden_sizes=[64, 64]):
19         super().__init__()
20         layers = []
21         in_dim = 1
22         for h in hidden_sizes:
23             layers.append(nn.Linear(in_dim, h))
24             layers.append(nn.Tanh())
25             in_dim = h
26         layers.append(nn.Linear(in_dim, 2))
27         self.net = nn.Sequential(*layers)
28     def forward(self, x):
29         return self.net(x)
30
31 # --- Data ---
32 N_total=2000
33 x_all = np.random.uniform(-1, 1, size=(N_total, 1)).astype(np.float32)
34 y_all = runge(x_all).astype(np.float32)
35 y_allD = Drunge(x_all).astype(np.float32)
```



```

36
37 # Split into train/val
38 x_train, x_val, y_train, y_val, yD_train, yD_val = train_test_split(x_all, y_all, y_allD, test_size=0.2, random_state=1)#訓練集和測試集
39
40 # Convert to torch tensors
41 X_train=torch.from_numpy(x_train)
42 Y_train=torch.from_numpy(y_train)
43 yD_train=torch.from_numpy(yD_train)
44 X_val=torch.from_numpy(x_val)
45 Y_val=torch.from_numpy(y_val)
46 yD_val=torch.from_numpy(yD_val)
47
48 # --- Model, loss, optimizer ---
49 lr=1e-3
50 epochs=2000
51 model =MLP(hidden_sizes=[64, 64])
52 criterion=nn.MSELoss()
53 optimizer=optim.Adam(model.parameters(), lr=lr)
54 train_losses=[]
55 val_losses=[]
56
57 # --- Training loop ---
58 for ep in range(epochs):
59     model.train()#start training by using training data
60     running_loss=0.0
61     optimizer.zero_grad()
62     out=model(X_train)
63     y_pred=out[:, 0]
64     yD_pred=out[:, 1]
65     y_loss=criterion(Y_train.squeeze(), y_pred)
66     yD_loss=criterion(yD_train.squeeze(), yD_pred)
67     loss=y_loss+yD_loss
68     loss.backward()
69     optimizer.step()
70

```

```

71     # Validation loss
72     model.eval()
73     with torch.no_grad():
74         out_val=model(X_val)
75         f_pred, fD_pred=out_val[:, 0], out_val[:, 1]
76         loss_val=criterion(f_pred, Y_val.squeeze())+criterion(fD_pred, yD_val.squeeze())
77     train_losses.append(loss.item())
78     val_losses.append(loss_val.item())
79
80     if (ep+1)%200 == 0 or ep == 0:
81         print(f'Epoch {ep+1}/{epochs} train MSE={loss.item():.6f} val MSE={loss_val.item():.6f}')
82
83 # --- Predictions on all points ---
84 model.eval()
85 with torch.no_grad():
86     out_all = model(torch.from_numpy(x_all))
87     f_pred_all = out_all[:,0].numpy()
88     df_pred_all = out_all[:,1].numpy()
89
90 # --- Error metrics ---
91 mse_f = mean_squared_error(y_all.squeeze(), f_pred_all)
92 mse_df = mean_squared_error(y_allD.squeeze(), df_pred_all)
93 maxerr_f = np.max(np.abs(y_all.squeeze() - f_pred_all))
94 maxerr_df = np.max(np.abs(y_allD.squeeze() - df_pred_all))
95
96 print("\n--- Error Metrics ---")
97 print(f"Function MSE: {mse_f:.6e}, Max Error: {maxerr_f:.6f}")
98 print(f"Derivative MSE: {mse_df:.6e}, Max Error: {maxerr_df:.6f}")
99
100 # --- Plots ---
101 #grid
102 x_grid=np.linspace(-1, 1, 500, dtype=np.float32).reshape(-1, 1)#convert to an array with 1 column
103 y_true_grid=runge(x_grid)
104 yD_true_grid=Drunge(x_grid)
105 model.eval()

```

```
106 with torch.no_grad():
107     out_grid = model(torch.from_numpy(x_grid))
108     f_pred_grid = out_grid[:,0].numpy()
109     df_pred_grid = out_grid[:,1].numpy()
110
111 plt.figure(figsize=(12,5))
112
113 # Function
114 plt.subplot(1,2,1)
115 plt.plot(x_grid, y_true_grid, label="True f(x)")
116 plt.plot(x_grid, f_pred_grid, "--", label="NN f(x)")
117 plt.legend()
118 plt.title("Function Approximation")
119
120 # Derivative
121 plt.subplot(1,2,2)
122 plt.plot(x_grid, yD_true_grid, label="True f'(x)")
123 plt.plot(x_grid, df_pred_grid, "--", label="NN f'(x)")
124 plt.legend()
125 plt.title("Derivative Approximation")
126
127 plt.tight_layout()
128 plt.savefig("function_and_derivative.png", dpi=200)
129 plt.show()
130
131 # Loss curves
132 plt.figure()
133 plt.plot(train_losses, label="Train Loss")
134 plt.plot(val_losses, label="Val Loss")
135 plt.legend()
136 plt.title("Training & Validation Loss")
137 plt.xlabel("Epochs")
138 plt.ylabel("Loss")
139 plt.savefig("loss_curves.png", dpi=200)
140 plt.show()
```