

Guolun Li

1.1

$$h_i = \phi(2(x_4 - x_{i_1})) = \begin{cases} 1, & \text{if } x_4 = x_{i_1} \\ 0, & \text{otherwise} \end{cases}$$

$$y = \phi(2(1 - \sum h_i)) = \begin{cases} 1, & \text{if } h_i = 1 \text{ for some } i \in \{1, 2, 3\} \\ 0, & \text{otherwise} \end{cases}$$

Thus,  $W^{(1)} \vec{x} + b^{(1)} = 2 \begin{pmatrix} x_4 - x_1 \\ x_4 - x_2 \\ x_4 - x_3 \end{pmatrix}$

$$\Rightarrow W^{(1)} = 2 \begin{pmatrix} -1 & 0 & 0 & 1 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 1 \end{pmatrix}, b^{(1)} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

and  $W^{(2)} \vec{h} + b^{(2)} = 2(1 - \sum h_i)$

$$\Rightarrow W^{(2)} = -2 \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}, b^{(2)} = -2$$

1.2 We use a brute force approach (list out all permutations) to achieve the goal.

We'll construct two hidden layers and one output layer.

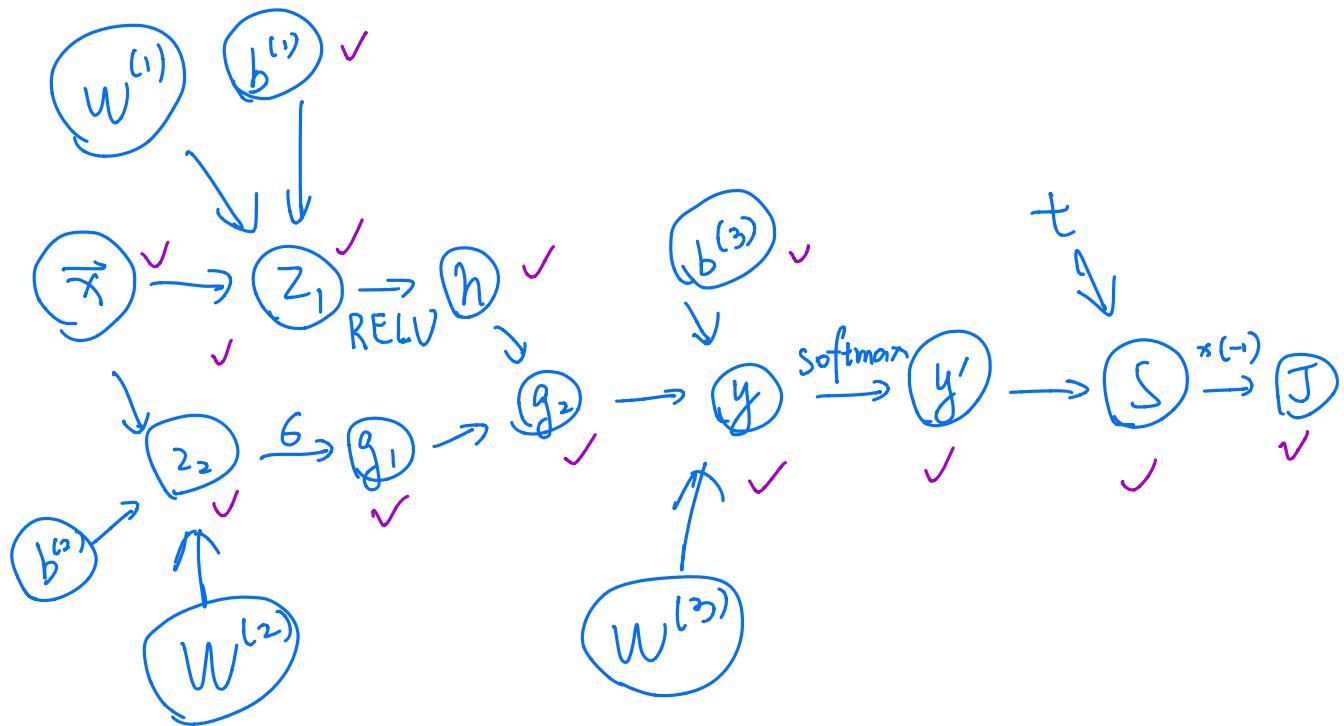
The first <sup>hidden</sup> layer compares, one by one, the first three nodes to the last three nodes. The second <sup>hidden</sup> layer enumerates all possibilities that the last three nodes are a permutation of the first three nodes.

The first hidden layer has 9 nodes, each node checks if  $x_i = x_j$ , for some  $i \in \{1, 2, 3\}$ ,  $j \in \{4, 5, 6\}$ . This is similar to Q1.1, and the activation function would be  $\phi_1(z) = I(z \in [-1, 1])$ .

The second hidden layer has 6 nodes, each node checks if  $(x_1, x_2, x_3) = (x_i, x_j, x_k)$ , where  $(x_i, x_j, x_k)$  is a permutation of  $(x_4, x_5, x_6)$ . This can be done by putting weight 1 on the nodes that check  $x_1 = x_i, x_2 = x_j, x_3 = x_k$ , and weight 0 on other nodes, then using an activation function  $\phi_2(z) = I(z = 3)$ . Then, the output layer checks if any node in the second hidden layer is 1. This is also similar to Q1.1, we simply put weight 1 on each node and use  $\phi(z) = I(z = 1)$ .

### 2.1.1 Computational Graph [0pt]

Draw the computation graph relating  $\mathbf{x}$ ,  $t$ ,  $\mathbf{z}_1$ ,  $\mathbf{h}$ ,  $\mathbf{z}_2$ ,  $\mathbf{g}_1$ ,  $\mathbf{g}_2$ ,  $\mathbf{y}$ ,  $\mathbf{y}'$ ,  $\mathcal{S}$  and  $\mathcal{J}$ .



### 2.1.2 Backward Pass [1pt]

Derive the backprop equations for computing  $\bar{\mathbf{x}} = \frac{\partial \mathcal{J}}{\partial \mathbf{x}}$ , one variable at a time, similar to the vectorized backward pass derived in Lec 2.

$$\bar{\mathcal{J}} = 1$$

$$\bar{\mathcal{S}} = -\bar{\mathcal{J}}$$

$$\bar{\mathbf{y}'} = \bar{\mathcal{S}} \nabla_{\mathbf{y}'} \mathcal{S} = -\bar{\mathcal{J}} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ y_t \\ \vdots \end{pmatrix}_{t^{\text{th}} \text{ position}}, \text{ since } \frac{\partial \mathcal{S}}{\partial y_k} = \frac{\partial}{\partial y_k} (\log y_t) = \begin{cases} \frac{1}{y_t}, & \text{if } k=t \\ 0, & \text{o.w.} \end{cases}$$

$\bar{y} = (\frac{\partial y'}{\partial y})^T \bar{y}'$ , where  $(\frac{\partial y'}{\partial y})_{i,j} = I(i=j) \text{softmax}(y_i) - \text{softmax}(y_i) \cdot \text{softmax}(y_j)$   
 (see lemma 1)

$$\bar{g}_2 = (W^{(2)})^T \bar{y}$$

$$\bar{h} = \bar{g}_2 \odot g_1$$

$$\bar{g}_1 = \bar{g}_2 \odot h$$

$$\bar{z}_2 = g'(z_2) \odot \bar{g}_1, \text{ where } g'(x) = \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right) = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\bar{z}_1 = \text{RELU}'(z_1) \odot \bar{h}, \text{ where } \text{RELU}'(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{if } x < 0 \end{cases}$$

(Define  $\text{RELU}'(0) = 1$  for computation purpose)

$$\bar{x} = W^{(1)^T} \bar{z}_1 + W^{(2)^T} \bar{z}_2$$

Lemma 1:

$$\begin{aligned} \left( \frac{\partial y'}{\partial y} \right)_{i,j} &= \frac{\partial y'_i}{\partial y_j} = \frac{\partial}{\partial y_j} \left( \frac{e^{y_i}}{\sum e^{y_k}} \right) \\ &= \frac{\partial}{\partial y_j} \frac{e^{y_i} \sum e^{y_k} - e^{y_j} \cdot e^{y_i}}{(\sum e^{y_k})^2} \\ &= \frac{1}{(\sum e^{y_k})^2} \left[ I(i=j) e^{y_i} \cdot \sum e^{y_k} - e^{y_i+y_j} \right] \\ &= I(i=j) \text{softmax}(y_i) - \text{softmax}(y_i) \cdot \text{softmax}(y_j) \end{aligned}$$

## 2.2.2 Computation Cost [1pt]

What is the number of scalar multiplications and memory cost of computing the Hessian  $\mathbf{H}$  in terms of  $n$ ?

2.2.2

$$L = \mathbf{x}^T \mathbf{V} \mathbf{V}^T \mathbf{x}$$

$$\nabla L = \mathbf{V} \mathbf{V}^T \mathbf{x}$$

$\mathbf{H} = \nabla^2 L = 2\mathbf{V} \mathbf{V}^T$  which requires  $2n^2$  multiplications and

$n^2 + n$  memory cost

$$\Downarrow \\ O(n^2)$$

$$O(n^2)$$

### 2.3 Vector-Hessian Products [1pt]

Compute  $\mathbf{z} = \mathbf{H}\mathbf{y} = \mathbf{v}\mathbf{v}^\top \mathbf{y}$  where  $n = 3$ ,  $\mathbf{v}^\top = [1, 2, 3]$ ,  $\mathbf{y}^\top = [1, 1, 1]$  using two algorithms: reverse-mode and forward-mode autodiff.

In backpropagation (also known as reverse-mode autodiff), you will compute  $\mathbf{M} = \mathbf{v}^\top \mathbf{y}$  first, then compute  $\mathbf{v}\mathbf{M}$ . Whereas, in forward-mode, you will compute  $\mathbf{H} = \mathbf{v}\mathbf{v}^\top$  then compute  $\mathbf{H}\mathbf{y}$ .

Write down the numerical values of  $\mathbf{z}^\top = [z_1, z_2, z_3]$  for the given  $\mathbf{v}$  and  $\mathbf{y}$ . What is the time and memory cost of evaluating  $\mathbf{z}$  with backpropagation (reverse-mode) in terms of  $n$ ? What about forward-mode?

For Q 2.3 and Q 2.4, we use the lemma that computing the matrix product  $MN$ , where  $M$  is  $a \times b$ ,  $N$  is  $b \times c$ , takes in total  $abc$  scalar multiplications. This is easy to see, as each entry of  $MN$  is the dot product of one row from  $M$  and one column from  $N$  (takes  $b$  steps), and there are  $ac$  such entries.

$$\mathbf{z}^\top = (6 \ 12 \ 18)$$

Back Prop:

$$\text{Time: } M = \mathbf{v}^\top \mathbf{y}, \quad \mathbf{z} = \mathbf{v}M$$

$$\text{Total time} = n + n = 2n$$

$$\text{Memory: } \mathbf{y}, \mathbf{v}, M, \mathbf{z}$$

$$\text{Total Memory} = n + n + 1 + n = 3n + 1$$

Forward :

$$\text{Time: } H = \mathbf{v}\mathbf{v}^\top, \quad \mathbf{z} = H\mathbf{y}$$

$$\text{Total time} = n^2 + n^2 = 2n^2$$

$$\text{Memory: } \mathbf{v}, H, \mathbf{y}, \mathbf{z}$$

$$\text{Total Memory} = n + n^2 + n + n = n^2 + 3n$$

$$n \times 1 \times n \times n \times 1 \times m$$

## 2.4 Trade-off of Reverse- and Forward-mode Autodiff [1pt] $Z = \mathbf{V} \mathbf{V}^T \mathbf{y}_1 \mathbf{y}_2^T$ $H = \mathbf{V} \mathbf{V}^T$

Consider computing  $\mathbf{Z} = \mathbf{H} \mathbf{y}_1 \mathbf{y}_2^T$  where  $\mathbf{v} \in \mathbb{R}^{n \times 1}$ ,  $\mathbf{y}_1 \in \mathbb{R}^{n \times 1}$  and  $\mathbf{y}_2 \in \mathbb{R}^{m \times 1}$ . What are the time and memory cost of evaluating  $\mathbf{Z}$  with reverse-mode in terms of  $n$  and  $m$ ? What about forward-mode? When is forward-mode a better choice? (Hint: Think about the shape of  $\mathbf{Z}$ , "tall" v.s. "wide".)

Back Prop:

$$\text{Time: } \mathbf{B}_1 = \mathbf{y}_1 \mathbf{y}_2^T, \quad \mathbf{B}_2 = \mathbf{V}^T \mathbf{B}_1, \quad \mathbf{Z} = \mathbf{V} \mathbf{B}_2$$

$$\text{Total time} = nm + nm + nm = 3mn$$

$$\text{Memory: } \mathbf{y}_1, \mathbf{y}_2, \mathbf{B}_1, \mathbf{V}, \mathbf{B}_2, \mathbf{Z}$$

$$\text{Total Memory} = n + m + nm + n + m + nm = 2(m+n+mn)$$

Forward:

$$\text{Time: } \mathbf{A}_1 = \mathbf{V} \mathbf{V}^T, \quad \mathbf{A}_2 = \mathbf{A}_1 \mathbf{y}_1, \quad \mathbf{Z} = \mathbf{A}_2 \mathbf{y}_2^T$$

$$\text{Total time} = n^2 + n^2 + nm = n(2n+m)$$

$$\text{Memory: } \mathbf{V}, \mathbf{A}_1, \mathbf{y}_1, \mathbf{A}_2, \mathbf{y}_2, \mathbf{Z}$$

$$\text{Total Memory} = n + n^2 + n + n + m + mn = 3n + m + n^2 + mn$$

$$\frac{\text{Back Prop time}}{\text{Forward time}} = n(2m - 2n)$$

$$\frac{\text{Back Prop memory}}{\text{Forward memory}} = (n+1)(m-n)$$

Thus, given the dimension of  $\mathbf{Z}$  being  $n \times m$ , using back propagation is better when  $m < n$ , using forward mode is better when  $m > n$ .

### 3.2 Underparameterized Model [1pt]

First consider the underparameterized  $d < n$  case. Write down the solution obtained by gradient descent assuming training converges. Show your work. Is the solution unique?

$$\begin{aligned} L &= \frac{1}{n} \|A\|^2, A = X\hat{w} - t \\ \Rightarrow \bar{A} &= \frac{2}{n} A \\ \Rightarrow \bar{w} &= X^T \bar{A} = X^T \cdot \frac{2}{n} A = \frac{2}{n} X^T (X\hat{w} - t) \\ \text{Gradient Descent will result in } \bar{w} &= 0 \Rightarrow X^T X \hat{w} = X^T t \\ &\Rightarrow \hat{w}^* = (X^T X)^{-1} X^T t. \end{aligned}$$

Since  $d < n$ , then  $X^T X$  is invertible, so  $\hat{w}^*$  has an explicit formula so it is unique.

#### 3.3.2 [1pt]

Now, let's generalize the previous 2D case to the general  $d > n$ . Show that gradient descent from zero initialization i.e.  $\hat{w}(0) = 0$  finds a unique minimizer if it converges. Write down the solution and show your work.

See next page.

$$\text{Solution: } \hat{w} = x^T(x^T x)^{-1} t$$

Proof:

$$L = \frac{1}{n} \|Ax - t\|^2, A = x \hat{w} - t.$$

$$\Rightarrow \bar{A} = \frac{1}{n} A \Rightarrow \bar{w} = x^T \bar{A} = x^T \cdot \frac{1}{n} A = \frac{1}{n} x^T (x \hat{w} - t) \text{ i.e. gradient}$$

$$\text{Thus, } \hat{w}_{k+1} = \hat{w}_k - \eta \cdot \frac{1}{n} x^T (x \hat{w}_k - t) \text{ i.e. gradient descent}$$

where  $\eta > 0$  is learning rate

We'll show two things:

1) Given  $\hat{w}_0 = \vec{0}$ ,  $\forall k \in N, \hat{w}_k = x^T c_k$  for some  $c_k \in \mathbb{R}^d$ .

2) Given that gradient descent converges, there will be  $x^T \hat{w} - t = 0$

We'll show 1) first by induction:

• Obviously  $\hat{w}_0 = x^T \vec{0}$  satisfies the requirement

• Assume that for  $k \in N$ ,  $\hat{w}_k = x^T c_k$  for some  $c_k \in \mathbb{R}^n$ .

$$\begin{aligned} \text{Then } \hat{w}_{k+1} &= x^T c_k - \eta \cdot \frac{1}{n} x^T (x \hat{w}_k - t) \\ &= x^T \left( c_k - \underbrace{\eta \cdot \frac{1}{n} x^T x \hat{w}_k + t}_{= c_{k+1} \in \mathbb{R}^n} \right) \end{aligned}$$

So 1) has been shown.

For 2), since gradient descent converges, then

$$\lim_{k \rightarrow \infty} \hat{w}_{k+1} = \lim_{k \rightarrow \infty} \hat{w}_k = \hat{w}, \text{ where } \hat{w} \text{ is the solution.}$$

$$\text{Thus, } \lim_{k \rightarrow \infty} \hat{w}_{k+1} = \lim_{k \rightarrow \infty} \hat{w}_k - \eta \cdot \frac{2}{n} X^T(X \hat{w}_k - t)$$

$$\Rightarrow \hat{w} = \hat{w} - \eta \cdot \frac{2}{n} X^T(X \hat{w} - t)$$

$$\Rightarrow X^T(X \hat{w} - t) = 0$$

$$\Rightarrow X X^T(X \hat{w} - t) = 0$$

In the case of  $d > n$ ,  $X X^T$  is  $n \times n$  thus invertible.

$$\Rightarrow (X X^T)^{-1}(X X^T)(X \hat{w} - t) = (X X^T)^{-1} \vec{0}$$

$$\Rightarrow X \hat{w} - t = 0$$

Combining 1) and 2): Let  $\hat{w} = X^T c$ , where

$$\hat{w} = \lim_{k \rightarrow \infty} \hat{w}_k \text{ and } c = \lim_{k \rightarrow \infty} c_k, \text{ then:}$$

$$X(X^T c) - t = 0$$

$$\Rightarrow X X^T c = t$$

$$\Rightarrow c = (X X^T)^{-1} t$$

$$\Rightarrow \hat{w} = X^T (X X^T)^{-1} t \text{ is unique}$$

### 3.3.3 [1pt]

Visualize and compare underparameterized with overparameterized polynomial regression: <https://colab.research.google.com/drive/1Atkk9hjUaXV-bDttCxAv9WiMsB6EJTKU>. Include your code snippets for the `fit_poly` function in the write-up. Does overparameterization (higher degree polynomial) always lead to overfitting, i.e. larger test error?

Function Fill In for `fit_poly()`:

```
def fit_poly(X, d, t):
    X_expand = poly_expand(X, d=d, poly_type = poly_type)
    n = X.shape[0]
    if d > n:
        W = X_expand.T.dot(np.linalg.inv(X_expand.dot(X_expand.T))).dot(t)
    else:
        W = np.linalg.inv(X_expand.T.dot(X_expand)).dot(X_expand.T).dot(t)
    return W
```

Loss Values Output:

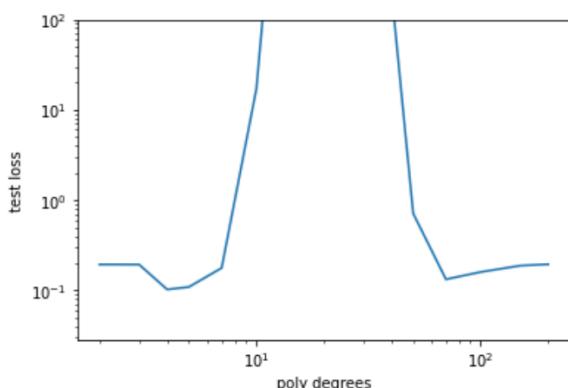
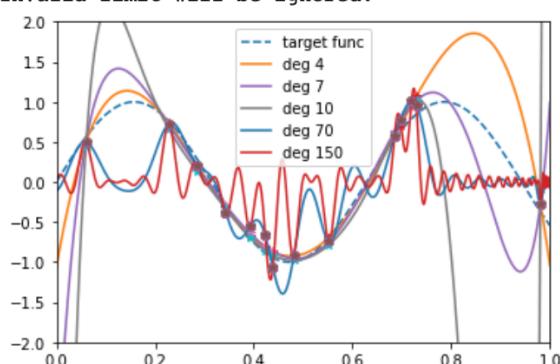
```
poly_type = 'chebyshev' # try legendre or chebyshev
loss_val_list = []

poly_degrees = [2, 3, 4, 5, 7, 10, 15, 20, 30, 50, 70, 100, 150, 200]
plot_poly_degrees = [4, 7, 10, 70, 150] ## only plot these polynomials
## 
plot_target_func()

for d in poly_degrees:
    W = fit_poly(X, d, t)
    plot_flag = True if d in plot_poly_degrees else False
    loss_val = plot_prediction(X, W, d, domain, plot_flag)
    loss_val_list.append(loss_val)
plt.legend()

plot_val_loss(poly_degrees, loss_val_list)
```

↳ /usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:43: UserWarning: Attempted to set no  
Invalid limit will be ignored.



```
[7] for i in range(len(poly_degrees)):
    print(f"degree = {poly_degrees[i]}, loss = {loss_val_list[i]}")
```

degree = 2, loss = 0.19386411705081336  
degree = 3, loss = 0.193432748835684  
degree = 4, loss = 0.10255222991025262  
degree = 5, loss = 0.10915309162008865  
degree = 7, loss = 0.177356881829107  
degree = 10, loss = 17.094862480309498  
degree = 15, loss = 1323020.060699293  
degree = 20, loss = 15596274997.211277  
degree = 30, loss = 317604.6227527557  
degree = 50, loss = 0.7120003515525691  
degree = 70, loss = 0.13281043237871068  
degree = 100, loss = 0.1599430277105052  
degree = 150, loss = 0.18800282052817618  
degree = 200, loss = 0.19468029146824442

As can be seen, the loss function increases with polynomial degree first, then decreases back to the original level. By comparing degree 7 and degree 70, we see that degree 70 has a smaller loss. This shows that overparametrization ( $n=14, d=71$ ) doesn't necessarily cause overfitting.