

$$W^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \phi^{(1)} = \begin{cases} z & \text{if } z > 0 \\ -z & \text{if } z \leq 0 \end{cases} = |z|$$

$$W^{(2)} = \begin{bmatrix} 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix} \quad b^{(2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \phi^{(2)} = z \quad (\text{linear})$$

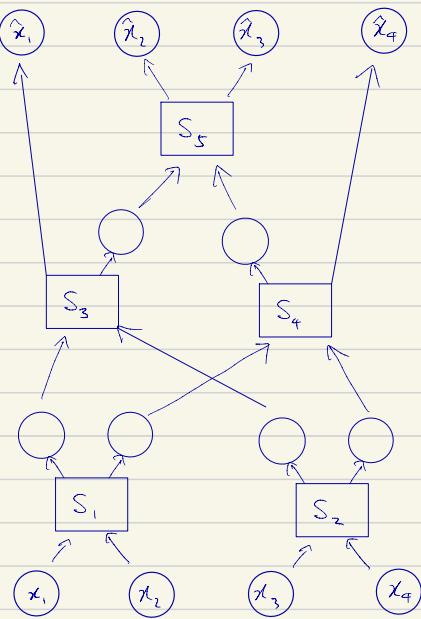
$$h_1 = x_1 + x_2$$

$$h_2 = x_1 - x_2$$

$$y_1 = 0.5 h_1 + 0.5 h_2$$

$$y_2 = 0.5 h_1 - 0.5 h_2$$

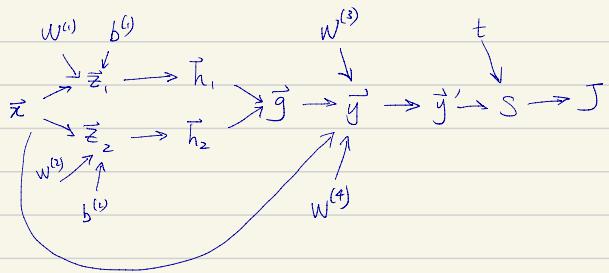
1.2) using merge sort



1.3)

$$W^{(3)} = [0 \ 0 \ | \ 0]$$

2.1.1)



2.1.2)

Dm $\bar{y} = 1$
 $\bar{s} = \bar{g} \frac{\partial \bar{y}}{\partial \bar{s}} = -\bar{f}$

$\bar{y}' = \bar{s} \frac{\partial \bar{s}}{\partial \bar{y}'} = \bar{s} \left[\frac{\mathbb{I}(t=1)}{y'_1} \frac{\mathbb{I}(t=2)}{y'_2} \dots \frac{\mathbb{I}(t=N)}{y'_N} \right]^T = \bar{s} \frac{1}{\bar{y}'}^T \circ \text{one-hot}(t) = \bar{s} \frac{1}{\bar{y}'}^T \underbrace{\left[0 \ 0 \ 0 \dots 0 \ 1 0 \dots 0 \right]^T}_{\text{where the } t^{\text{th}} \text{ index is } 1}$

$\bar{y} = \text{softmax}'(\bar{y}) \bar{y}'$

$\bar{g} = w^{(3)T} \bar{y}$

$\bar{h}_1 = \begin{bmatrix} h_{21} & 0 \\ 0 & h_{22} \\ \vdots & \vdots \\ 0 & h_{2m} \end{bmatrix} \begin{bmatrix} 1 \\ \bar{g} \end{bmatrix} = \bar{g} \circ h_2$

$\bar{h}_2 = \begin{bmatrix} h_{11} & h_{12} & 0 \\ 0 & \ddots & h_{1m} \end{bmatrix} \begin{bmatrix} 1 \\ \bar{g} \end{bmatrix} = \bar{g} \circ h_1$

$\bar{z}_1 = \begin{bmatrix} \mathbb{I}(z_{11} > 0) & 0 \\ 0 & \mathbb{I}(z_{12} > 0) \\ \vdots & \vdots \\ 0 & \mathbb{I}(z_{1m} > 0) \end{bmatrix} \begin{bmatrix} 1 \\ \bar{h}_1 \end{bmatrix} = \bar{h}_1 \circ \mathbb{I}(\bar{z}_1 > 0)$, where $\mathbb{I}(\bar{z}_1 > 0) = \begin{bmatrix} \mathbb{I}(z_{11} > 0) \\ \mathbb{I}(z_{12} > 0) \\ \vdots \\ \mathbb{I}(z_{1m} > 0) \end{bmatrix}$

$\bar{z}_2 = \begin{bmatrix} \sigma(z_{21})(1-\sigma(z_{21})) & 0 \\ \sigma(z_{22})(1-\sigma(z_{22})) & \ddots \\ \vdots & \vdots \\ 0 & \sigma(z_{2m})(1-\sigma(z_{2m})) \end{bmatrix} \begin{bmatrix} 1 \\ \bar{h}_2 \end{bmatrix} = \bar{h}_2 \circ \sigma(\bar{z}_2) \circ (1-\sigma(\bar{z}_2))$

Sinc: $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\begin{aligned} \sigma'(z) &= \frac{d}{dz} (1+e^{-z})^{-1} \\ &= -(1+e^{-z})^{-2} (-e^{-z}) \\ &= \frac{-e^{-z}}{(1+e^{-z})^2} \\ &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ &= \sigma(z) (1 - \sigma(z)) \end{aligned}$$

$\mathcal{D} \times 1 \quad \bar{x} = W^{(1)T} \bar{z}_1 + W^{(2)T} \bar{z}_2 + W^{(4)T} \bar{y}$

$$2.2.1) \quad \mathbf{x} = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \quad \bar{\mathbf{y}} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \xrightarrow{\mathbf{z} = \mathbf{x}^T \bar{\mathbf{y}}} \xrightarrow{\mathbf{h} = \text{ReLU}(\mathbf{z})} \xrightarrow{\mathbf{y} = \text{ReLU}(\mathbf{h})}$$

forward pass:

$$\mathbf{z} = \begin{bmatrix} 1 & 2 & 1 \\ -2 & 1 & 0 \\ 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 8 \\ -6 \end{bmatrix}$$

$$\mathbf{h} = \text{ReLU} \begin{bmatrix} 8 \\ -6 \end{bmatrix} = \begin{bmatrix} 8 \\ 0 \end{bmatrix}$$

backward pass:

$$\bar{\mathbf{h}} = \mathbf{W}^{(2)^T} \bar{\mathbf{y}} = \begin{bmatrix} -2 & 4 & 1 \\ 1 & -2 & -3 \\ -3 & 4 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ -4 \\ 7 \end{bmatrix}$$

$$\bar{\mathbf{z}} = \bar{\mathbf{h}} \circ \mathbb{I}(\mathbf{z} > 0) = \begin{bmatrix} 3 \\ -4 \\ 7 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 3 \\ -4 \\ 0 \end{bmatrix}$$

$$\frac{\partial J}{\partial \mathbf{W}^{(1)}} = \mathbf{x} \bar{\mathbf{z}}^T = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} [3 \ -4 \ 0] = \begin{bmatrix} 3 & -4 & 0 \\ 9 & -12 & 0 \\ 3 & -4 & 0 \end{bmatrix}$$

$$\frac{\partial J}{\partial \mathbf{W}^{(2)}} = \mathbf{h} \bar{\mathbf{y}}^T = \begin{bmatrix} 8 \\ 0 \\ 1 \end{bmatrix} [1 \ 1 \ 1] = \begin{bmatrix} 8 & 8 & 8 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\left\| \frac{\partial J}{\partial \mathbf{W}^{(1)}} \right\|_F^2 = \left\| \begin{bmatrix} 3 & -4 & 0 \\ 9 & -12 & 0 \\ 3 & -4 & 0 \end{bmatrix} \right\|_F^2 = 9 + 16 + 81 + 144 + 9 + 16$$

$$= \boxed{275}$$

$$\left\| \frac{\partial J}{\partial \mathbf{W}^{(2)}} \right\|_F^2 = \left\| \begin{bmatrix} 8 & 8 & 8 \\ 0 & 0 & 0 \end{bmatrix} \right\|_F^2 = 64 + 64 + 64 + 1 + 1 + 1 = \boxed{195}$$

2.2.2)

$$\begin{aligned}\left\| \frac{\partial Y}{\partial W^{(1)}} \right\|_F^2 &= \text{trace} \left(\frac{\partial Y^T}{\partial W^{(1)}} \frac{\partial Y}{\partial W^{(1)}} \right) \\ &= \text{trace} (\bar{z} x^T x \bar{z}^T) \\ &= \text{trace} (x^T x \bar{z}^T \bar{z}) \quad (\text{Cyclic Property of Trace}) \\ &= (x^T x)(\bar{z}^T \bar{z}) \quad (\text{Scalar Multiplication}) \\ &= \|x\|_2^2 \|\bar{z}\|_2^2 \\ &= \left\| \begin{bmatrix} 1 \\ 3 \\ 7 \end{bmatrix} \right\|_2^2 \left\| \begin{bmatrix} 3 \\ -4 \\ 0 \end{bmatrix} \right\|_2^2 \\ &= (1+9+1)(9+16+0) \\ &= 11 \cdot 25 \\ &= \boxed{275} \quad \text{same as previous}\end{aligned}$$

$$\begin{aligned}\left\| \frac{\partial Y}{\partial W^{(2)}} \right\|_F^2 &= \text{trace} \left(\frac{\partial Y^T}{\partial W^{(2)}} \frac{\partial Y}{\partial W^{(2)}} \right) \\ &= \text{trace} (\bar{y} h^T h \bar{y}^T) \\ &= \text{trace} (h^T h \bar{y}^T \bar{y}) \quad (\text{Cyclic Property of Trace}) \\ &= (h^T h)(\bar{y}^T \bar{y}) \quad (\text{Scalar Multiplication}) \\ &= \|h\|_2^2 \|\bar{y}\|_2^2 \\ &= \left\| \begin{bmatrix} 8 \\ 1 \\ 0 \end{bmatrix} \right\|_2^2 \left\| \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\|_2^2 \\ &= (64+1+0)(1+1+1) = 65 \cdot 3 = \boxed{195} \quad \text{same as previous}\end{aligned}$$

2-2-3)

	T (Naive)	T (Efficient)	M (Naive)	M (Efficient)
Forward Pass	$D^2 N(k-1) \textcircled{①}$	$D^2 N(k-1) \textcircled{①}$	$O(D^2 k + DNk)$	$O(D^2 k + DNk)$
Backward Pass	$D^2 N(2k-1) \textcircled{②}$	$D^2 N(k-1) \textcircled{③}$	$O(D^2 Nk)$	$O(D^2 k + DNk)$
Gradient Norm Computation	$D^2 Nk$	$Nk(2D+1)$	$O(D^2 Nk)$	$O(DNk)$

$$\text{model 1: } h_1 = w_1 x$$

$$h_2 = w_2 h_1$$

:

$$h_{k-1} = w_{k-1} h_{k-2}$$

$$y = h_k = w_k h_{k-1}$$

general assumptions: \bar{y} is given and we only need to compute what is needed to calculate the gradient norms of all the weights.

note ①: assuming \bar{y} is given, forward pass only needs to compute h_1, h_2, \dots, h_{k-1} since we do not need the output $y = h_k$ for computing the gradient norms of the weights. (hence the $k-1$)

note ②: assuming \bar{y} is given, then we only need to compute $\bar{h}_{k-1}, \bar{h}_{k-2}, \dots, \bar{h}_1$ and $\frac{\partial y}{\partial w^{(1)}}, \frac{\partial y}{\partial w^{(2)}} \dots, \frac{\partial y}{\partial w^{(k)}}, \frac{\partial y}{\partial w^{(k)}}$ for computing the gradient norms of the weights. (hence the $2k-1$)

note ③: assuming \bar{y} is given, then we only need to compute $\bar{h}_{k-1}, \bar{h}_{k-2}, \dots, \bar{h}_1$ for computing the gradient norms of the weights. (hence the $k-1$)

$$3) \quad X = n \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} \quad t = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_n \quad \hat{w} = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_d \quad \varepsilon = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_n$$

$$3.1) \quad \mathcal{L} = \frac{1}{n} \| X \hat{w} - t \|_2^2$$

$$\frac{\partial \mathcal{L}}{\partial \hat{w}} = \frac{2}{n} X^T (X \hat{w} - t) = \frac{2}{n} (X^T X \hat{w} - X^T t)$$

3.2.1)

$$\text{set } \frac{\partial \mathcal{L}}{\partial \hat{w}} = \frac{2}{n} (X^T X \hat{w} - X^T t) = 0$$

$$\Leftrightarrow X^T X \hat{w} - X^T t = 0$$

$$\Leftrightarrow X^T X \hat{w} = X^T t$$

$$\Leftrightarrow \hat{w} = (X^T X)^{-1} X^T t$$

invertible since $n > d$

3.2.2)

$$\text{error} = \frac{1}{n} \| X \hat{w} - t \|_2^2, \quad \hat{w} = (X^T X)^{-1} X^T t \quad t = X w^* + \varepsilon$$

$$= \frac{1}{n} \| X (X^T X)^{-1} X^T t - t \|_2^2$$

$$= \frac{1}{n} \| X (X^T X)^{-1} X^T (X w^* + \varepsilon) - (X w^* + \varepsilon) \|_2^2$$

$$= \frac{1}{n} \| X (X^T X)^{-1} X^T X w^* + X (X^T X)^{-1} X^T \varepsilon - X w^* - \varepsilon \|_2^2$$

$$= \frac{1}{n} \| X w^* + X (X^T X)^{-1} X^T \varepsilon - X w^* - \varepsilon \|_2^2$$

$$= \frac{1}{n} \| X (X^T X)^{-1} X^T \varepsilon - \varepsilon \|_2^2$$

$$= \frac{1}{n} \| (X (X^T X)^{-1} X^T - I) \varepsilon \|_2^2 \quad \text{as wanted}$$

expectation of error:

$$\text{Error} = \frac{1}{n} \| (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I}) \boldsymbol{\varepsilon} \|^2$$
$$= \frac{1}{n} [(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I}) \boldsymbol{\varepsilon}]^T (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I}) \boldsymbol{\varepsilon}$$

$$\text{where } A = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$
$$A^T = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$
$$= \frac{1}{n} ([A - \mathbf{I}] \boldsymbol{\varepsilon}]^T (A - \mathbf{I}) \boldsymbol{\varepsilon})$$
$$= \frac{1}{n} (\boldsymbol{\varepsilon}^T (A - \mathbf{I})^T (A - \mathbf{I}) \boldsymbol{\varepsilon})$$
$$= \frac{1}{n} \boldsymbol{\varepsilon}^T (AA^T - A - A^T \mathbf{I}) \boldsymbol{\varepsilon}$$

$$= \frac{1}{n} \boldsymbol{\varepsilon}^T (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{I}) \boldsymbol{\varepsilon}$$
$$= \frac{1}{n} \boldsymbol{\varepsilon}^T (\mathbf{A}^T - \mathbf{A} - \mathbf{A}^T + \mathbf{I}) \boldsymbol{\varepsilon}$$
$$= \frac{1}{n} \boldsymbol{\varepsilon}^T (\mathbf{I} - \mathbf{A}) \boldsymbol{\varepsilon}$$
$$= \frac{1}{n} (\boldsymbol{\varepsilon}^T - \boldsymbol{\varepsilon}^T \mathbf{A}) \boldsymbol{\varepsilon}$$
$$= \frac{1}{n} (\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T \mathbf{A} \boldsymbol{\varepsilon})$$

$$\Rightarrow E(\text{error}) = E\left[\frac{1}{n} (\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T \mathbf{A} \boldsymbol{\varepsilon})\right]$$
$$= \frac{1}{n} E[\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T \mathbf{A} \boldsymbol{\varepsilon}]$$
$$= \frac{1}{n} (E[\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}] - E[\boldsymbol{\varepsilon}^T \mathbf{A} \boldsymbol{\varepsilon}]) \text{ by linearity of expectation}$$

$$= \frac{1}{n} (n\sigma^2 - \sigma^2 d) \text{ see next page}$$

$$= \sigma^2 - \frac{\sigma^2 d}{n}$$
$$= \boxed{\sigma^2 \left(1 - \frac{d}{n}\right)}$$

$$\star E[\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}]$$

$$= E[\varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2]$$

$$= E[\varepsilon_1^2] + E[\varepsilon_2^2] + \dots + E[\varepsilon_n^2] \quad \text{by linearity of Expectation}$$

$$= n \sigma^2 \quad \text{since } \text{var}(\varepsilon) = E(\varepsilon^2) - [E(\varepsilon)]^2 \text{ and } E(\varepsilon) = 0.$$

$$\star E[\boldsymbol{\varepsilon}^T \mathbf{A} \boldsymbol{\varepsilon}] = E\left[\begin{bmatrix} \varepsilon_1 & \varepsilon_2 & \dots & \varepsilon_n \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & \ddots & & \\ \vdots & & \ddots & \\ A_{n1} & & & A_{nn} \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}\right]$$

$$= E\left[\begin{bmatrix} \sum_{i=1}^n \varepsilon_i A_{1i} & \sum_{i=1}^n \varepsilon_i A_{i2} & \dots & \sum_{i=1}^n \varepsilon_i A_{in} \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}\right]$$

$$= E\left[\sum_{j=1}^n \sum_{i=1}^n \varepsilon_i \varepsilon_j A_{ij}\right]$$

$$= E\left[\sum_{i=1}^n A_{ii} \varepsilon_i^2 + \sum_{i,j \neq i}^n A_{ij} \varepsilon_i \varepsilon_j\right]$$

$$= \sum_{i=1}^n E[A_{ii} \varepsilon_i^2] + \sum_{i,j \neq i}^n E[A_{ij} \varepsilon_i \varepsilon_j] \quad \text{by linearity of expectation}$$

$$= \sum_{i=1}^n A_{ii} E[\varepsilon_i^2] + \sum_{i,j \neq i}^n A_{ij} E[\varepsilon_i] E[\varepsilon_j] \quad \text{since } \varepsilon_i \text{ and } \varepsilon_j \text{ are independent when } i \neq j.$$

$$= \sum_{i=1}^n A_{ii} \sigma^2 + \sum_{i,j \neq i}^n A_{ij} (0)(0)$$

$$= \sigma^2 \sum_{i=1}^n A_{ii}$$

$$= \sigma^2 \text{tr}(\mathbf{A}) \quad , \text{ recall } \mathbf{A} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$= \sigma^2 \text{tr}(\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)$$

$$= \sigma^2 \text{tr}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}) \quad \text{by cyclic property of trace.}$$

$$= \sigma^2 \text{tr}(\mathbf{I}) \quad , \mathbf{I} \text{ is dimension } d \times d$$

$$= \sigma^2 d$$

3.3.2.)

$$\hat{\mathbf{w}} = \mathbf{X}^T \mathbf{a} \quad \text{for some } \mathbf{a} \in \mathbb{R}^n$$

Substitute this into

$$\frac{\partial L}{\partial \hat{\mathbf{w}}} = \frac{2}{n} \mathbf{X}^T (\mathbf{X} \hat{\mathbf{w}} - \mathbf{t}) = 0$$

$$\Rightarrow \frac{2}{n} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T \mathbf{a} - \mathbf{t}) = 0$$

$$\Rightarrow \mathbf{X}^T (\mathbf{X} \mathbf{X}^T \mathbf{a} - \mathbf{t}) = 0$$

$$\Rightarrow \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T \mathbf{a} - \mathbf{t}) = \mathbf{X} \cdot 0$$

$$\Rightarrow (\mathbf{X} \mathbf{X}^T \mathbf{a} - \mathbf{t}) = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \cdot 0,$$

$\mathbf{X} \mathbf{X}^T$ is invertible since $d > n$.

$$\Rightarrow \mathbf{X} \mathbf{X}^T \mathbf{a} - \mathbf{t} = 0$$

$$\Rightarrow \mathbf{X} \mathbf{X}^T \mathbf{a} = \mathbf{t}$$

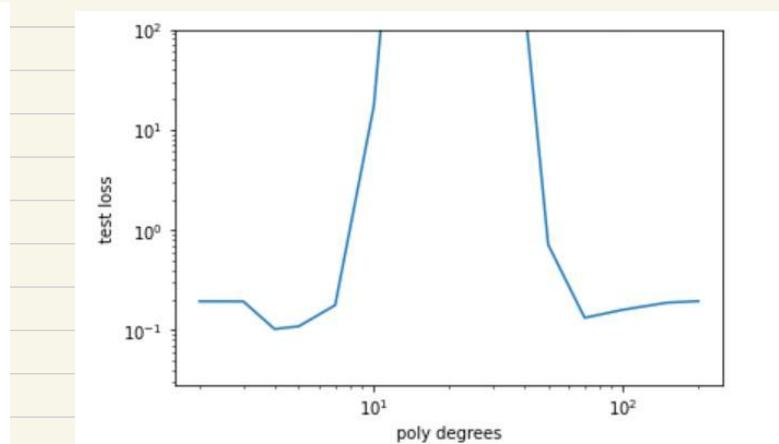
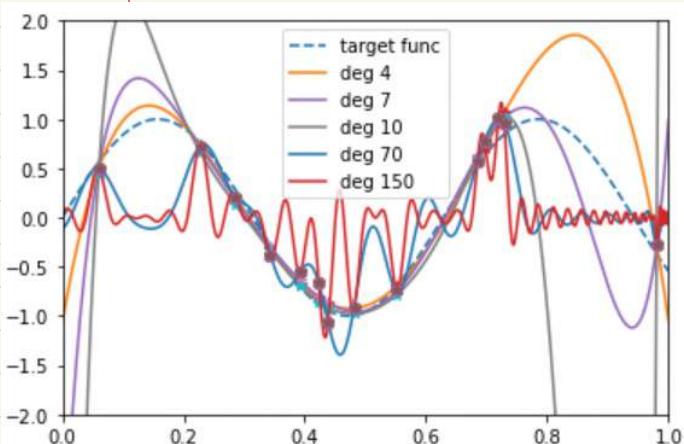
$$\Rightarrow \mathbf{a} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{t}$$

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

3.3.4)

to be implemented; fill in the derived solution for the underparameterized ($d < n$) and overparameterized ($d > n$) problem

```
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 @ np.linalg.inv(X_expand @ X_expand.T)) @ t
    else:
        W = (np.linalg.inv(X_expand.T @ X_expand) @ X_expand.T) @ t
    return W
```



Conclusion:

Overparametrization does not always lead to overfitting. From the second graph above, we see that test loss for polynomial degree 10^2 is about the same as test loss for degree ≈ 4 , and much less than test loss for degrees between $10 - 50$.