

1.1

$$J = \frac{1}{n} \|X\hat{w} - t\|_2^2 = \frac{1}{n} (X\hat{w} - t)^T \underbrace{(X\hat{w} - t)}_y$$

using chain Rule: $\frac{dJ}{d\hat{w}} = \frac{dJ}{dy} \cdot \frac{dy}{d\hat{w}}$

$$\left. \begin{aligned} \frac{dJ}{dy} &= \frac{d}{dy} \left(\frac{1}{n} y^T y \right) = \frac{2}{n} y \\ \frac{dy}{d\hat{w}} &= \frac{d}{d\hat{w}} (X\hat{w} - t) = X \end{aligned} \right\} \frac{dJ}{d\hat{w}} = \frac{2}{n} y \cdot X = \frac{2}{n} (X\hat{w} - t) \cdot X = \frac{2}{n} X^T (X\hat{w} - t)$$

Underparameterized Model

1.2.1

$$J = \frac{1}{n} \sum_{i=1}^n (\hat{w}^T x_i - t_i)^2 = \frac{1}{n} \|X\hat{w} - t\|_2^2$$

Gradient descent $\hat{w} \leftarrow \hat{w} - \alpha \frac{\partial J}{\partial \hat{w}}$ This converges once $\frac{dJ}{d\hat{w}} = 0$

set $\frac{dJ}{d\hat{w}} = 0$ to find the minimum from 1.1 $\frac{dJ}{d\hat{w}} = \frac{2}{n} X^T (X\hat{w} - t)$

$$\frac{2}{n} X^T (X\hat{w} - t) = 0 \rightarrow X^T (X\hat{w} - t) = 0 \rightarrow X^T X \hat{w} - X^T t = 0 \rightarrow X^T X \hat{w} = X^T t \rightarrow \hat{w} = (X^T X)^{-1} X^T t$$

\uparrow
 $X^T X$ is invertible since $d > n$

1.2.2

- Show $\mathcal{L}_{\text{loss}} = \frac{1}{n} \|X\hat{w} - t\|_2^2 = \frac{1}{n} \|X(X^T X)^{-1} X^T t - t\|_2^2$

from 1.2.1 $\rightarrow \hat{w} = (X^T X)^{-1} X^T t$ sub into Loss $\rightarrow \frac{1}{n} \|X(X^T X)^{-1} X^T t - t\|_2^2$

$t_i = w^* x_i + \varepsilon_i$ $t = \sum w^* x_i + \varepsilon_i = Xw^* + \varepsilon$ sub this in for $t \rightarrow \mathcal{L}_{\text{loss}} = \frac{1}{n} \|X(X^T X)^{-1} X^T (Xw^* + \varepsilon) - (Xw^* + \varepsilon)\|_2^2$

$$\mathcal{L}_{\text{loss}} = \frac{1}{n} \left\| \underbrace{X(X^T X)^{-1} X^T}_{XIX^{-1} = I} Xw^* + X(X^T X)^{-1} X^T \varepsilon - (Xw^* + \varepsilon) \right\|_2^2 = \frac{1}{n} \left\| \cancel{Xw^*} + X(X^T X)^{-1} X^T \varepsilon - \cancel{Xw^*} - \varepsilon \right\|_2^2$$

$\mathcal{L}_{\text{loss}} = \frac{1}{n} \|X(X^T X)^{-1} X^T \varepsilon - \varepsilon\|_2^2$ factor out ε $\mathcal{L}_{\text{loss}} = \frac{1}{n} \|(X(X^T X)^{-1} X^T - I) \varepsilon\|_2^2$

- find $E \left[\frac{1}{n} \|(X(X^T X)^{-1} X^T - I) \varepsilon\|_2^2 \right] = \frac{1}{n} E \left[\|(X(X^T X)^{-1} X^T - I) \varepsilon\|_2^2 \right]$

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

Proof that $(X(X^T X)^{-1} X^T - I)$ is a symmetric matrix.

$(X^T X)^T = X^T X$ hence $X^T X$ is symmetric therefore $(X^T X)^{-1}$ is symmetric.

$(X(X^T X)^{-1} X^T)^T = X(X^T X)^{-1} X^T$ hence $X(X^T X)^{-1} X^T$ is symmetric.

The Identity matrix is also symmetric. The difference of two symmetric matrix is also symmetric.

→ $(X(X^T X)^{-1} X^T - I)$ is symmetric.

$$\mathcal{E}^T (X(X^T X)^{-1} X^T - I)^2 \mathcal{E}$$

From matrixcookbook $E[X^T A X] = \text{Tr}(A \text{Var}[X]) + E[X]^T A E[X]$ if A is symmetric

Since $\text{Var}[\mathcal{E}] = \sigma^2$ and $E[\mathcal{E}] = 0$

$$E[\mathcal{E}^T (X(X^T X)^{-1} X^T - I)^2 \mathcal{E}] = \text{Tr}((X(X^T X)^{-1} X^T - I)^2 \sigma^2) = \sigma^2 \text{Tr}((X(X^T X)^{-1} X^T - I)^2)$$

$$(X(X^T X)^{-1} X^T - I)^2 = \underbrace{(X(X^T X)^{-1} X^T)^2}_{\substack{\text{I} \\ \text{I}}} - 2(X(X^T X)^{-1} X^T) + \underbrace{I}_{\text{I}}$$

$$\rightarrow (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T) = X \underbrace{(X^T X)^{-1} X^T X}_{I} (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T$$

$$(X(X^T X)^{-1} X^T - I)^2 = X(X^T X)^{-1} X^T - 2X(X^T X)^{-1} X^T + I = I - X(X^T X)^{-1} X^T$$

$$E[\mathcal{E}^T (X(X^T X)^{-1} X^T - I)^2 \mathcal{E}] = \sigma^2 \text{Tr}(I - X(X^T X)^{-1} X^T) = \sigma^2 \text{Tr}(I) - \text{Tr}(X(X^T X)^{-1} X^T)$$

$$\left. \begin{array}{l} \text{using cyclic property of trace} \rightarrow \text{Tr}(X(X^T X)^{-1} X^T) = \text{Tr}(X^T X)^{-1} X^T X = \text{Tr}(I) \\ \text{Since } X \text{ is } n \times d \quad (X^T X)^{-1} X^T X \text{ is } d \times d \end{array} \right\} \text{Tr}(X(X^T X)^{-1} X^T) = d$$

Dimension of I is $n \times n$ Since $X(X^T X)^{-1} X^T$ is an $n \times n$ matrix → $\text{Tr}(I) = n$

Hence

$$E[\mathcal{E}^T (X(X^T X)^{-1} X^T - I)^2 \mathcal{E}] = \sigma^2 (n - d)$$

$$Loss = \frac{1}{n} \mathcal{E}^T (X(X^T X)^{-1} X^T - I)^2 \mathcal{E} \rightarrow E[Loss] = \sigma^2 \left(1 - \frac{d}{n}\right)$$

Overparameterized Model

1.3.1

$$\hat{w} = [w_1, w_2]^T \quad x_1 = [1, 1] \quad t = 3$$

$$\hat{w} x_1 = t \quad \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} [1, 1] = 3 \quad w_1 + w_2 = 3 \leftarrow \text{equation of the line}$$

$$w = \begin{bmatrix} w_1 \\ 3 - w_1 \end{bmatrix} \text{ this shows that } w_1 \text{ and } w_2 \text{ can be any value}$$

1.3.2

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

$$\text{Loss} = \frac{1}{n} \|X \hat{w} - t\|_2^2 \quad \text{gradient decent: } \hat{w} \leftarrow \hat{w} - \alpha \frac{\partial \text{loss}}{\partial \hat{w}}$$

$$\text{As calculate in 1.1 } \frac{d \text{loss}}{d \hat{w}} = \frac{2}{n} X^T (X \hat{w} - t)$$

$$\hat{w}' \leftarrow \hat{w}(0) - \alpha \frac{2}{n} X^T (X \hat{w} - t)$$

The gradient decent converges at a point where $X^T (X \hat{w} - t) = 0$

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

Given that we start from $w(0) = 0$ which is in the row space of X

and

$\alpha \frac{2}{n} X^T (X \hat{w} - t)$ is also in the row space of X (due to X^T multiplication)

} Any resultant w from subtracting these two when updating the \hat{w} is also in the row space of X . Therefore any \hat{w} can be represented via $\hat{w} = X^T \alpha$

$$X^T X \hat{w} = X^T t \quad \hat{w} = X^T \alpha \rightarrow X^T X X^T \alpha = X^T t \quad \xrightarrow{X X^T \text{ is invertible when } d > n} X^T \alpha = X^T (X X^T)^{-1} t$$

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

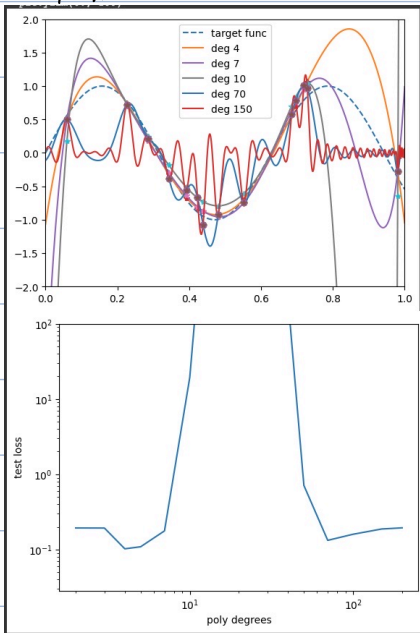
1.3.4

Code snippet:

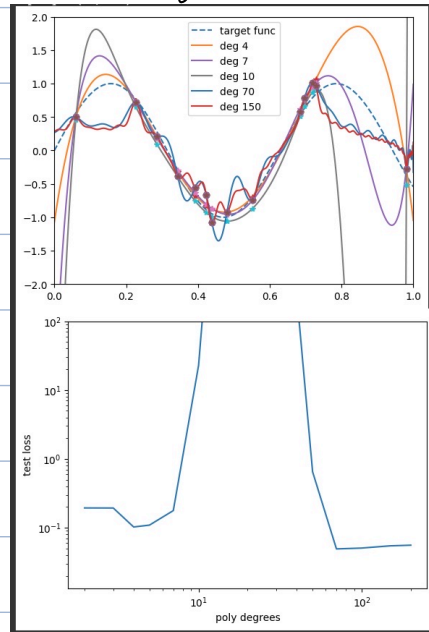
```
[6] # to be implemented; fill in the derived solution for the underparameterized (d<n) and overparameterized (d>n) problem

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
        ## W = ... (Your solution for Part 1.3.2)
    else:
        W = np.linalg.inv(X_expand.T @ X_expand) @ X_expand.T @ t
        ## W = ... (Your solution for Part 1.2.1)
    return W
```

Poly type = Chebyshev



poly type = Legendre



Actual $n = 14$

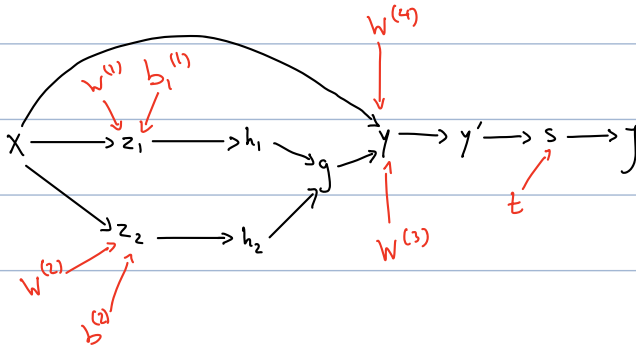
It can be seen that overparameterization does not always lead to overfitting.

From the test error vs poly degree it can be seen that as the poly degree gets closer to the target n of 14 the test error increases. However, around a poly degree of 70 the test error drops, showing that increase in poly degree does not always cause overfitting.

Back propagation

Automatic Differentiation

2.1.1



2.1.2

$$\bullet \quad \bar{y} = -1$$

$$\bullet \quad \bar{s} = \frac{\partial j}{\partial s} \bar{y} = -s$$

$$\bullet \quad \bar{y}' = \frac{ds}{dy'} \bar{s}$$

Since $\mathbb{I}(t=k) = 1$ only when $k=t$
and $\frac{d \log(y'_k)}{dy'_k} = \frac{1}{y'_k}$

$$\bar{y} = [0 \dots \frac{1}{y'_t} \dots 0]^T \bar{s}$$

$$\bullet \quad \bar{y} = \frac{dy'}{dy} \bar{y}' = (sofmax'(y)) \bar{y}'$$

$$\bullet \quad \bar{g} = \frac{dy}{dg} \bar{y} = w^{(3)T} \bar{y}$$

$$\bullet \quad \bar{h}_2 = \frac{\partial g^T}{\partial h_2} \bar{g} = \frac{\partial g_i^T}{\partial h_{2i}} = \begin{bmatrix} h_{1,1} & 0 & 0 & \dots \\ 0 & h_{2,2} & 0 & \dots \\ 0 & 0 & \dots & h_{1,n} \end{bmatrix} \text{ when multiplied by } \bar{g} = h_1 \circ \bar{g} \quad \bar{h}_2 = h_1 \circ \bar{g}$$

$$\bullet \quad \bar{h}_1 = \frac{\partial g^T}{\partial h_1} \bar{g} = h_2 \circ \bar{g} \quad \text{same reasoning as above}$$

$$\bullet \quad \bar{z}_2 = \frac{\partial h_2^T}{\partial z_2} \bar{h}_2 = \sigma'(z_2) \circ \bar{h}_2 \quad \bullet \quad \bar{z}_1 = \frac{\partial h_1^T}{\partial z_1} \bar{h}_1 \quad \frac{\text{Relu}(z_1)^T}{\partial z_1} = \text{Relu}'(z_1) \circ \bar{h}_1 \quad \text{Relu}'(z_1) = \begin{cases} 1 & \text{if } z_1 > 0 \\ 0 & \text{if } z_1 < 0 \end{cases}$$

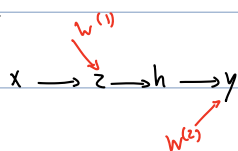
$$\sigma'(z_2) = \frac{d}{dz} (1 + e^{-z})^{-1} = -(1 + e^{-z})^{-2} \times \frac{d}{dz} (1 + e^{-z})$$

$$\bar{z}_2 = (- (1 + e^{-z})^2 \cdot e^{-z}) \circ \bar{h}_2$$

$$\bullet \quad \bar{x} = \frac{\partial z_1^T}{\partial x} \bar{z}_1 + \frac{\partial z_2^T}{\partial x} \bar{z}_2 + \frac{\partial y^T}{\partial x} \bar{y} = w^{(1)T} \bar{z}_1 + w^{(2)T} \bar{z}_2 + w^{(3)T} \bar{y}$$

Gradient Norm Computation

2.2.1



$$\bullet \quad \bar{y} = 1 \quad \bullet \quad \bar{y} = \frac{dy^T}{dy} \bar{y} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \bullet \quad \bar{h} = \frac{dy^T}{dh} \bar{y} = (W^{(2)})^T \bar{y} = \begin{bmatrix} -2 & 1 & -3 \\ 4 & -2 & 4 \\ 1 & -3 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -4 \\ 6 \\ 4 \end{bmatrix}$$

$$\bullet \quad \bar{z} = \frac{dh^T}{dz} \bar{h} = \text{Relu}'(z) \circ \bar{h} \rightarrow h = \text{Relu}(z) = \begin{bmatrix} 8 \\ 1 \\ 0 \end{bmatrix} \quad \text{Relu}'(z) = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \rightarrow \bar{z} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 \\ 6 \\ 4 \end{bmatrix} = \begin{bmatrix} -4 \\ 6 \\ 0 \end{bmatrix}$$

$$\bullet \quad z = W^{(1)} x = \begin{bmatrix} 1 & 2 & 1 \\ -2 & 1 & 0 \\ 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 8 \\ 1 \\ -6 \end{bmatrix}$$

$$\bullet \quad \frac{\partial J}{\partial W^{(1)}} = \bar{z} \frac{dz^T}{dW^{(1)}} = \bar{z} x^T = \begin{bmatrix} -4 \\ 6 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 3 & 1 \end{bmatrix} = \begin{bmatrix} -4 & -12 & -4 \\ 6 & 18 & 6 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\bullet \quad \left\| \frac{\partial J}{\partial W^{(1)}} \right\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 = 4^2 + 12^2 + 4^2 + 6^2 + 18^2 + 6^2 = 572$$

$$\bullet \quad \frac{\partial J}{\partial W^{(2)}} = \bar{y} \frac{dy^T}{dW^{(2)}} = \bar{y} h^T = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 8 & 1 & 0 \\ 8 & 1 & 0 \\ 8 & 1 & 0 \end{bmatrix}$$

$$\bullet \quad \left\| \frac{\partial J}{\partial W^{(2)}} \right\|_F^2 = 8^2 + 8^2 + 8^2 + 1^2 + 1^2 + 1^2 = 195$$

2.2.2

$$\left\| \frac{\partial J}{\partial W^{(1)}} \right\|_F^2 = \|x\|_2^2 \|z\|_2^2 = (1^2 + 3^2 + 1^2) \times ((-4)^2 + 6^2 + 0^2) = 11 \times 52 = 572$$

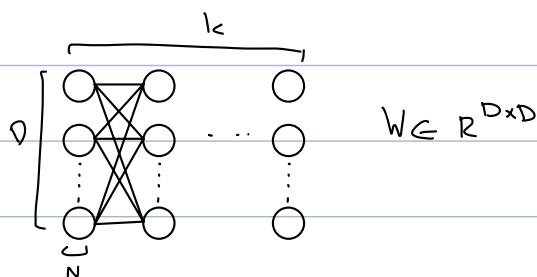
Same as the answer in the prev. question

$$\left\| \frac{\partial J}{\partial W^{(2)}} \right\|_F^2 = \|h\|_2^2 \|\bar{y}\|_2^2 = (8^2 + 1^2 + 0^2) \times (1^2 + 1^2 + 1^2) = 65 \times 3 = 195$$

Same as the answer in the prev. question

2.2.3

	$T(\text{Naive})$	$T(\text{Efficient})$	$M(\text{Naive})$	$M(\text{Efficient})$
Forward Pass	$N k D^2$	$N k D^2$	$O(kD^2 + kND)$	$O(kD^2 + kND)$
Backward Pass	$2 N k D^2$	$N k D^2$	$O(kND^2)$	$O(kD^2 + kND)$
Gradient Norm Computation	$N k D^2$	$(2D+1) N k$	$O(kND^2)$	$O(kND)$



Forward pass:

$T(\text{Naive})$: k layer with $D \times D$ matrix each D neurons and N training sample $\underline{N k D^2}$

$T(\text{Efficient})$: same as Naive

$M(\text{Naive})$: For the k , $D \times D$ matrices $k D^2$, Also need to save value of neurons (values needed for back propagation) so k layer with N training sample and D neurons kND . $\rightarrow \underline{O(kD^2 + kND)}$

$M(\text{Efficient})$: Same as Naive

Backward Pass:

$T(\text{Naive})$: Error vector computation $D \times D$ matrix multiplied by D dimension error vector. (For N training sample and k layers) Gradient $D \times D$ gradient matrix and D neurons for k layers and N sample. $\rightarrow \underline{2 N k D^2}$

$T(\text{Efficient})$: Error vector computation same as Naive. For Gradient we don't store all of the $D \times D$ gradient matrices just add them after each step. $\rightarrow \underline{N k D^2}$

$M(\text{Naive})$: We need to store Error vector $[O(N k D)]$, parameters/weights $[O(k D^2)]$, and gradient $[O(k N D^2)]$ $O(k N D^2)$ dominates $\rightarrow \underline{O(k N D^2)}$

$M(\text{Efficient})$: Same as Naive but gradient is just $O(k D^2)$ as we are updating the sum each time.

$\rightarrow \underline{O(k D^2 + N k D)}$

Gradient Norm Computation:

$T(\text{Naive})$: Square all the gradient D^2 scalar computation for k layers and N training sample $\rightarrow NkD^2$

$T(\text{Efficient})$: $\|x\|_2^2 \|\bar{y}\|_2^2 \rightarrow$ Norm of a D dimension vector has D scalar computation. This is done twice.

Then they are multiplied by each other. $(2D+1)$ scalar computation for k layers and N samples $\rightarrow (2D+1)Nk$

$M(\text{Naive})$: Storing $D \times D$ matrices for k layer and N training sample. $\rightarrow O(kND^2)$

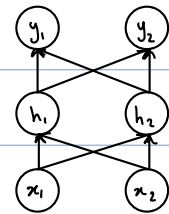
$M(\text{Efficient})$: Stores D dimension vector for k layer and N samples $\rightarrow O(kND)$

Hard Coding Networks

Sort two numbers

3.1 $\max(x_1, x_2) = \frac{1}{2}(x_1 + x_2) + \frac{1}{2}|x_1 - x_2|$ $\min(x_1, x_2) = \frac{1}{2}(x_1 + x_2) - \frac{1}{2}|x_1 - x_2|$

$$W^{(2)} W^{(1)} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}(x_1 + x_2) - \frac{1}{2}|x_1 - x_2| \\ \frac{1}{2}(x_1 + x_2) + \frac{1}{2}|x_1 - x_2| \end{bmatrix}$$



$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} ex_1 + fx_2 \\ gx_1 + hx_2 \end{bmatrix} = \begin{bmatrix} a(ex_1 + fx_2) + b(gx_1 + hx_2) \\ c(ex_1 + fx_2) + d(gx_1 + hx_2) \end{bmatrix}$$

$$\begin{bmatrix} a(ex_1 + fx_2) + b(gx_1 + hx_2) \\ c(ex_1 + fx_2) + d(gx_1 + hx_2) \end{bmatrix} = \begin{bmatrix} \frac{1}{2}(x_1 + x_2) - \frac{1}{2}|x_1 - x_2| \\ \frac{1}{2}(x_1 + x_2) + \frac{1}{2}|x_1 - x_2| \end{bmatrix} \rightarrow \begin{matrix} a = \frac{1}{2} & b = -\frac{1}{2} & c = \frac{1}{2} & d = \frac{1}{2} \\ e = 1 & f = 1 & g = 1 & h = -1 \end{matrix}$$

$$\phi^{(2)}(z) = z$$

$$\phi^{(1)}(z) = |z|$$

$$W^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$b^{(2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\phi^{(1)}(z) = |z|$$

$$\phi^{(2)}(z) = z$$

Perform Sort

3.2 Merge Sort

