

HW 3

黃熙漢

## 5. Counter Exp:

$H_1 = \{h(x) = 0\}$ , can only assign a label of 0 to any value of input, I guess often use in benchmark for other algorithm

$H_2 = \{h(x) = 1\}$ , only assign a label of 1

Following the equation:

$$d_{vc}(H_1 \cup H_2) \leq d_{vc}(H_1) + d_{vc}(H_2)$$

$$d_{vc}(H_1 \cup H_2) = 1, \quad \because H_1 = 1 \& H_2 = 0, \quad \text{shattered if } N=1, \quad \text{all possibilities which include } x_0 = 0, x_0 = 1, \quad \{0, 1\} \text{ hypo}$$

$d_{vc}(H_1) = 0$ , can't shattered all points

$d_{vc}(H_2) = 0$ , same with  $H_1$

$$\therefore d_{vc}(H_1 \cup H_2) \leq d_{vc}(H_1) + d_{vc}(H_2)$$

$$1 \leq 0 \quad (\text{disprove})$$

6.

		predict	
		$h(x)$	
actual	$y(x)$	+1	-1
	+1	0	1 0
	-1	1	0

$$\text{see: } \mathbb{P} \text{ token} = \mathbb{T} \mathbb{P}$$

$$\mathbb{I} \mathbb{N} \text{ token} = \mathbb{T} \mathbb{I} \mathbb{N}$$

calculate that expected cost

$$\text{cost}(+) = \mathbb{T} \mathbb{P} \times \mathbb{P}(\gamma = -1 | x)$$

\* which  $\mathbb{P}(\gamma = -1 | x) + \mathbb{P}(\gamma = +1 | x) = 1$

$$\text{cost}(-) = \mathbb{T} \mathbb{I} \mathbb{N} \times \mathbb{P}(\gamma = +1 | x)$$

compare two expected cost:

$$\text{cost}(+) \leq \text{cost}(-)$$

$$\mathbb{T} \mathbb{P} \times \mathbb{P}(\gamma = -1 | x) \leq \mathbb{T} \mathbb{I} \mathbb{N} \times \mathbb{P}(\gamma = +1 | x)$$

$$\text{which } \mathbb{P}(\gamma = -1 | x) + \mathbb{P}(\gamma = +1 | x) = 1$$

$$\tau \models N = 10 \tau \models P$$

$$\frac{1}{10} \tau \models N \times (1 - P(Y = +1 | X)) \leq \tau \models P \times P(Y = +1 | X)$$

$$P(Y = +1 | X) \geq \frac{1}{11}$$

$\therefore$  Prove that  $\alpha = \frac{1}{11}$ , idea is sacrifice  $\tau \models P$  to reduce probability of error of  $\tau \models N$ .

7.

Define error for each  $x$

$$e(h, x) = P_{Y|X}(h(x) \neq Y)$$

$$e(f, x) = P_{Y|X}(f(x) \neq Y)$$

when  $h(x) = f(x)$ :

$$e(h, x) = e(f, x)$$

if  $h(x) \neq f(x)$ :

for given equation we know that,  $P(Y = f(x) | X) \geq \frac{1}{2}$ , so

$$P(Y = h(x) | X) \leq \frac{1}{2}$$

$$\text{Therefore } e(h, x) = P_{Y|X}(h(x) \neq Y) = 1 - P_{Y|X}(Y = h(x)) \geq \frac{1}{2}$$

$$e(f, x) = P_{Y|X}(f(x) \neq y) \leq \frac{1}{2}$$

$$e(h, x) = P_{Y|X}(y \neq f(x))$$

$$e(h, x) - e(f, x) = 2P_{Y|X}(y \neq f(x)) - 1 \geq 0$$

$$\text{Since } P_{Y|X}(y \neq f(x)) \geq \frac{1}{2}, \quad e(h, x) - e(f, x) \geq 0$$

upper bound:

$$e(h, x) \leq e(f, x) + \begin{cases} 1 & \text{when } h(x) \neq f(x) \\ 0 & \text{else} \end{cases}$$

so,

$$E_{out}(h) = E_x[e(h, x)] \leq E_x[e(f, x) + \mathbb{1}_{h(x) \neq f(x)}] = E_x[e(f, x)] + E_x[\mathbb{1}_{h(x) \neq f(x)}]$$

Note that:

$$E_x[e(f, x)] = E_{out}^{(2)}(f)$$

$$E_x[\mathbb{1}_{h(x) \neq f(x)}] = P_x[h(x) \neq f(x)] = E_{out}^{(1)}(h)$$

$\therefore$  prove that  $E_{out}^{(2)}(h) \leq E_{out}^{(2)}(f) + E_{out}^{(1)}(h)$ ,  
means for any hypothesis, its error under the noisy distribution  
 $E_{out}(h)$  is bounded by its error to the target function

$\|e_{\text{one}}(h)\|$  plus the irreducible error  $\|e_{\text{one}}^*(f)\|$ .

8.

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \end{bmatrix} \quad X' = \begin{bmatrix} 1126 & x_1 \\ 1126 & x_2 \\ 1126 & x_3 \end{bmatrix}$$

relationship between  $X$  &  $X'$

$X = X' C$ ,  $C$  is diagonal matrix

$$\text{Define } C : \begin{bmatrix} 1126 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

we know  $w_{\text{lin}} = (X^T X)^{-1} X^T y$  (unique solution)

$$\begin{aligned} w_{\text{luck}} &= (X'^T X')^{-1} X'^T y \\ &= (S X^T S X)^{-1} S X^T y \\ &= S^{-1} (X^T X)^{-1} S S^{-1} y \\ &= S^{-1} (X^T X)^{-1} y \end{aligned}$$

$$w_{\text{luck}} = S^{-1} w_{\text{lin}}$$

$\therefore w_{\text{lin}} = S w_{\text{luck}}$ , which  $D = S$

$w_{\text{lin}} = D w_{\text{luck}}$ , I learned how the weight change by scale features.

9. likelihood function:  $\prod_{n=1}^N (h(x))$

$$\min \bar{L}(w) = \frac{1}{N} \sum_{n=1}^N -\ln\left(\frac{1}{2} \left( \frac{w^T x}{\sqrt{1+(w^T x)^2}} + 1 \right)\right)$$

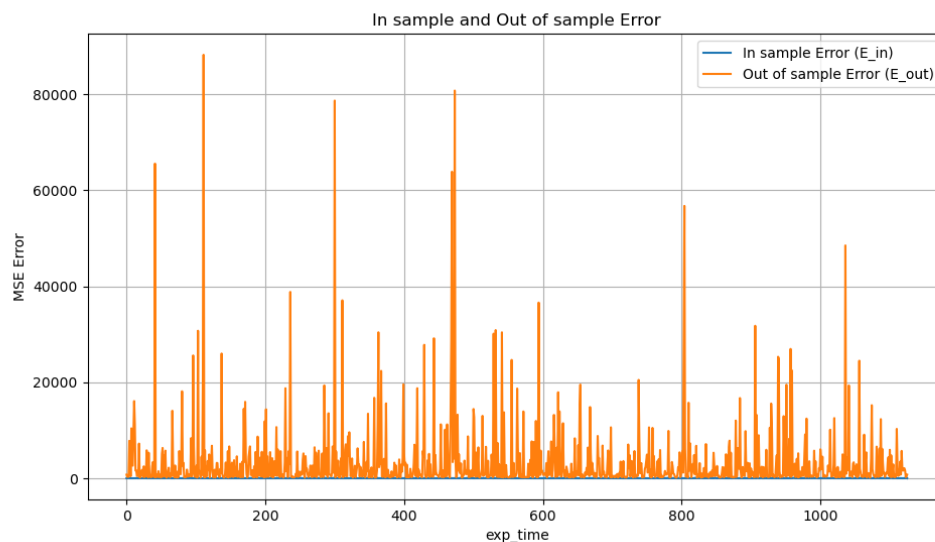
$$\nabla \bar{L}(w) = \frac{1}{N} \sum_{n=1}^N -\ln\left(\frac{1}{2} \frac{w^T x}{\sqrt{1+(w^T x)^2}} + \frac{1}{2}\right)$$

$$= \frac{1}{N} \sum \left( \frac{1}{2 \sqrt{1+(w^T x)^2}} + \frac{1}{2} \right) \cdot \frac{1}{2} \cdot \frac{x \frac{w^T x}{\sqrt{1+(w^T x)^2}} - w^T x \cdot (1+(w^T x)^2)^{-\frac{1}{2}} \cdot 2w^T x}{1+(w^T x)^2}$$

#

10. I noticed a significant gap between  $E_{out}$  and  $E_{in}$ . I believe this gap occurs because the training data size is very small, with only 32 samples, which has likely caused the model to overfit. Overfitting happens when the model performs exceptionally well on the training data but struggles to generalize to unseen data. As a result,  $E_{in}$  is much lower than  $E_{out}$ , indicating that the model has learned specific patterns or noise in the training set that don't hold in the test set.

MSE Result :



Code :

```

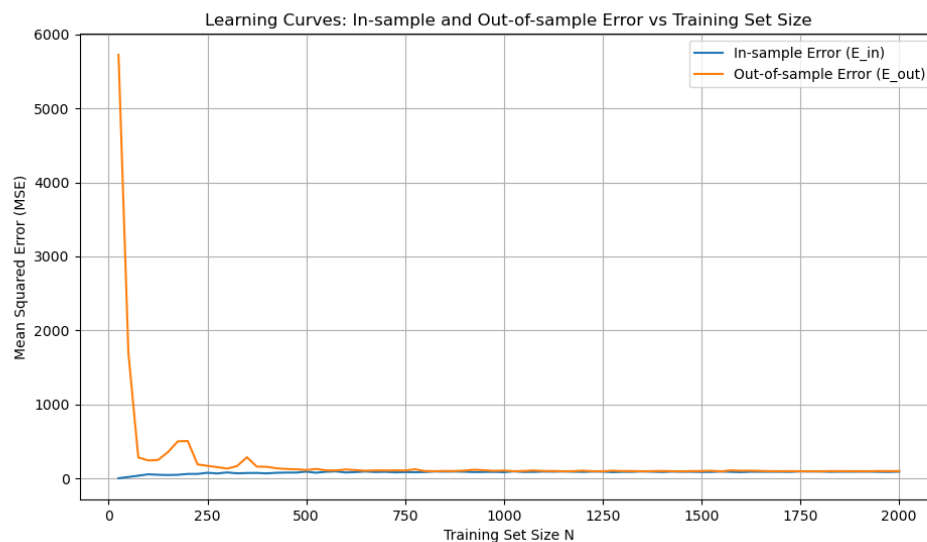
21 for i in range(data_num) :
22     x[i,0] = 1 #設置x0
23     for j in range(1,feat + 1) :
24         x[i,j] = data[i][j] #data[第幾個字典][第幾個鍵值]
25
26 for t in range(experiment_time) :
27     splice_data = np.random.permutation(data_num)
28     train_set = splice_data[:N] #N_exp 個 traindata
29     test_set = splice_data[N:]
30
31     xtrain_set = x[train_set] # train set 是一組數字 index, x[train_set] 就會對著他要的行來取。
32     xtest_set = x[test_set]
33
34     ytrain_set = arr_label[train_set]
35     ytest_set = arr_label[test_set]
36
37     #W_Lin = x^pseudoinv @ y
38     x_pseudo_inv = np.linalg.pinv(xtrain_set)
39     w_Lin = x_pseudo_inv@ytrain_set
40
41     # y = wx (outsample)
42     y_ht_train = xtrain_set @ w_Lin #don't forget mtrx rule
43     y_ht_test = xtest_set @ w_Lin
44
45     #E_in = (1/N)(y - y^)^2 , y^ = w_Lin@X ,
46     per_E_in = np.mean((ytrain_set - y_ht_train)**(2))
47     per_E_out = np.mean((ytest_set - y_ht_test)**(2))

```



11. Description: As shown in the graph, we can observe that as we use more training data, both  $E_{in}$  and  $E_{out}$  begin to converge. Initially, there is a noticeable gap between  $E_{in}$  and  $E_{out}$ , which indicates overfitting when the model is trained on smaller datasets. However, as the amount of training data increases, the gap narrows, and the errors converge, indicating that the model is generalizing better and overfitting is reduced.

Result:



Code:

```
label, data = svm_read_problem('ml_hw_2_feat')

arr_label = np.array(label)

N_Range = list(range(25, 2001, 25))
data_num = 8192
experiment_time = 16
feat = 12

E_in = []
E_out = []

np.random.seed(1)

x = np.zeros((data_num, feat + 1))

for i in range(data_num):
    x[i, 0] = 1 #設置x0
    for j in range(1, feat + 1):
        x[i, j] = data[i][j] #data[第幾個字典][第幾個鍵值]

for N in N_Range:
    E_in_avg = []
    E_out_avg = []
    for t in range(experiment_time):
        splice_data = np.random.permutation(data_num)
        train_set = splice_data[:N] #N exp 個 traindata
```

```

E_out_avg = []
for t in range(experiment_time) :
    splice_data = np.random.permutation(data_num)
    train_set = splice_data[:N] #N_exp 個 traindata
    test_set = splice_data[N:]

    xtrain_set = x[train_set] # train set 是一組數字 index, x
    xtest_set = x[test_set]

    ytrain_set = arr_label[train_set]
    ytest_set = arr_label[test_set]

    #W_Lin = x^pseudoinv @ y
    x_pseudo_inv = np.linalg.pinv(xtrain_set)
    w_Lin = x_pseudo_inv@ytrain_set

    # y = wx (outsample)
    y_ht_train = xtrain_set @ w_Lin #don't forget mtrx rule
    y_ht_test = xtest_set @ w_Lin

    #E_in = (1/N)(y - y^)^2 , y^ = wLin@X ,
    per_E_in = np.mean((ytrain_set - y_ht_train)**(2))
    per_E_out = np.mean((ytest_set - y_ht_test)**(2))

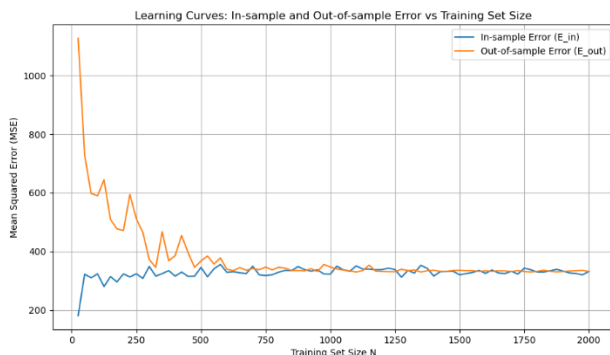
    E_out_avg.append(per_E_out) , E_in_avg.append(per_E_in)
E_out.append(np.mean(E_out_avg))
E_in.append(np.mean(E_in_avg))

```

12.

Description: We could see that final Eout values of the two models are very close, which indicates that the noise level has a greater impact on the model's performance, and increasing the number of features did not significantly reduce the generalization error. I guess it could be because most of the useful information in the data has already been captured by the first two features, and adding more features did not provide much additional valuable information.

Result:



Code: `feat = 2`