

5. Prove / disprove:

For any $M_1, M_2 \neq \emptyset$ (for binary classification, separate or same input space)
 $d_{VC}(M_1 \cup M_2) \leq d_{VC}(M_1) + d_{VC}(M_2)$

Pf:

Suppose $\exists M_1, M_2 \neq \emptyset$ s.t. $d_{VC}(M_1 \cup M_2) > d_{VC}(M_1) + d_{VC}(M_2)$

Let $d_{VC}(M_1) = k_1$, $d_{VC}(M_2) = k_2$, then

$$\begin{cases} m_{M_1}(k_1) = 2^{k_1}, & m_{M_1}(k_1+1) < 2^{k_1+1} \\ m_{M_2}(k_2) = 2^{k_2}, & m_{M_2}(k_2+1) < 2^{k_2+1} \end{cases}$$

$\Rightarrow d_{VC}(M_1 \cup M_2) > k_1 + k_2$, which means

\exists some $k_1 + k_2 + 1$ points that can be shattered by $M_1 \cup M_2$

Let $d_{VC}(M_1 \cup M_2) = k$, $k > k_1 + k_2$

$$m_{M_1 \cup M_2}(k_1 + k_2 + 1) \leq (k_1 + k_2 + 1)^k + 1 = 2^k$$

$$\therefore k_1, k_2 \in \mathbb{Z}^+ \cup \{0\}$$

$$\therefore (k_1, k_2) \in \{(0,1), (1,0)\} \Rightarrow (k_1 + k_2 + 1)^{k_1+1} = 2^{k_1+1} > 2^k \quad (\text{G})$$

Thus, $k_1 = k_2 = 0$, then $(k_1 + k_2 + 1)^{k_1+1} = 1 \neq 2$

function

By defn, $d_{VC}(M) = 0 \Leftrightarrow M = \{y = a, a \in \mathbb{R}\}$ (M has exactly one hypothesis, which is a constant)

Therefore, let $M_1 = \{y = a_1\}$, where $a_1 \in \mathbb{R}$, $M_2 = \{y = a_2\}$ where $a_2 \in \mathbb{R}$

$$\Rightarrow M_1 \cup M_2 = \{y = a_1, y = a_2\} \text{ where } a_1, a_2 \in \mathbb{R}.$$

$$\text{However, } d_{VC}(M_1 \cup M_2) = 0 \neq 0$$

Thus, $d_{VC}(M_1 \cup M_2) \leq d_{VC}(M_1) + d_{VC}(M_2) \quad \square$

6. $\therefore f_{\text{MKT}}(\vec{x}) = \text{sign}(P(y=+1 | \vec{x}) - \alpha)$

$\therefore \text{if } \begin{cases} P(y=+1 | \vec{x}) > \alpha, \text{ then } f_{\text{MKT}}(\vec{x}) = +1 \\ P(y=+1 | \vec{x}) \leq \alpha, \text{ then } f_{\text{MKT}}(\vec{x}) = -1 \end{cases}$

		\vec{y}	
		+1	-1
\vec{h}	+1	0	① FP.
	-1	② FN	0

If we consider the points on the decision boundary, the expected loss for assigning any such \vec{x} to either side should be the same (because the decision boundary is the place where we believe to have same probability to classify points to both sides)

Therefore, for \vec{x} on the boundary, we have that equality holds:

(FN) $\underbrace{1 \cdot P(y=+1 | \vec{x})}_{\text{the probability of having the point } \vec{x} \text{ (which is on the boundary) to have true label } y=+1} = 1 \cdot \underbrace{P(y=-1 | \vec{x})}_{\text{the probability of having the point } \vec{x} \text{ (which is on the boundary) to have true label } y=-1}$

$\Rightarrow 1 \cdot P(y=+1 | \vec{x}) = 1 \cdot (1 - P(y=+1 | \vec{x}))$

$\Rightarrow 1 \cdot P(y=+1 | \vec{x}) = 1 - P(y=+1 | \vec{x})$

$\Rightarrow 11 P(y=+1 | \vec{x}) = 1$

$\Rightarrow P(y=+1 | \vec{x}) = \frac{1}{11}$

Therefore, we assign α as $\alpha = \frac{1}{11}$, making $f_{\text{MKT}}(\vec{x}) = \text{sign}(P(y=+1 | \vec{x}) - \frac{1}{11})$ \square

because this means:

if $P(y=+1 | \vec{x}) > \frac{1}{11}$, then $f_{\text{MKT}}(\vec{x}) = +1$

\rightarrow we only need the probability to be more than $\frac{1}{11}$ (but not 0.5)

to make the model classify \vec{x} to be a member (+1)

\rightarrow we classify a customer to be a member more often because the cost of FN is greater.

Claim:

↙ naive

$$2 \quad \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq y)] \leq \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq f(\vec{x}))] + \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(f(\vec{x}) \neq y)]$$

$$\text{Suppose } \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq y)] > \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq f(\vec{x}))] + \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(f(\vec{x}) \neq y)]$$

$\downarrow S$
 $\downarrow S_1$
 $\downarrow S_2$

Let S be the set that contains all \vec{x}_i s.t. $h(\vec{x}_i) \neq y_i$

S_1 be the set that contains all \vec{x}_j s.t. $h(\vec{x}_j) \neq f(\vec{x}_j)$

S_2 be the set that contains all \vec{x}_k s.t. $f(\vec{x}_k) \neq y_k$

Suppose $\exists \vec{x}_t \in S$, then $h(\vec{x}_t) \neq y_t$ $\left\{ \begin{array}{l} \text{if } \vec{x}_t \in S_1, \text{ then } f(\vec{x}_t) = y_t \rightarrow \vec{x}_t \notin S_2 \\ \text{if } \vec{x}_t \in S_2, \text{ then } f(\vec{x}_t) \neq y_t \rightarrow \vec{x}_t \notin S_1 \end{array} \right.$

Thus, $S_1 \cap S_2 = \emptyset$

If $\vec{x}_t \in S$ and $\vec{x}_t \notin S_1 \cap S_2$

$\rightarrow h(\vec{x}_t) \neq y_t \cap h(\vec{x}_t) \neq f(\vec{x}_t) \cap f(\vec{x}_t) \neq y_t \rightarrow$ \nexists binary classification

Therefore, $\forall \vec{x}_t \in S$, either $\vec{x}_t \in S_1$ or $\vec{x}_t \in S_2$

$$\therefore S \subseteq S_1 \cup S_2 \Rightarrow |S| \leq |S_1 \cup S_2| = |S_1| + |S_2|$$

\uparrow
 $\because S_1 \cap S_2 = \emptyset$

Which leads to a contradiction.

$$\text{Thus, } \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq y)] \leq \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x})}} [\mathbb{I}(h(\vec{x}) \neq f(\vec{x}))] + \mathbb{E}_{\substack{\vec{x} \sim p(\vec{x}) \\ y \sim p(y|\vec{x})}} [\mathbb{I}(f(\vec{x}) \neq y)] \quad \square$$

8. linear regression on dataset: $\{(\tilde{x}_n, y_n)\}_{n=1}^N$ $\tilde{x}_n \in \mathbb{R}^{d+1}$ with $\tilde{x}_{n0} = 1 \forall n \in \{1, \dots, N\}$

Assume: $X^T X$: invertible

unique sol \vec{w}_{LSR} acquired on $\{(\tilde{x}_n, y_n)\}_{n=1}^N$

+ change $\tilde{x}_{n0} = 1 \forall n \in \{1, \dots, N\}$

Proof: $\vec{w}_{LSR} = D \vec{w}_{LSRy}$ where D : diagonal matrix

" $X^T X$: invertible $\therefore \vec{w}_{LSR} = \overbrace{(X^T X)^{-1}}^{X^+} X^T y$

Write:

$$X' = \begin{bmatrix} 1\tilde{x}_0 & \tilde{x}_{11} & \tilde{x}_{12} & \dots & \tilde{x}_{1d} \\ 1\tilde{x}_0 & \tilde{x}_{21} & \tilde{x}_{22} & \dots & \tilde{x}_{2d} \\ \vdots & \vdots & & \ddots & \vdots \\ 1\tilde{x}_0 & \tilde{x}_{N1} & \tilde{x}_{N2} & \dots & \tilde{x}_{Nd} \end{bmatrix} = \begin{bmatrix} 1 & \tilde{x}_{11} & \tilde{x}_{12} & \dots & \tilde{x}_{1d} \\ 1 & \tilde{x}_{21} & \tilde{x}_{22} & \dots & \tilde{x}_{2d} \\ \vdots & \vdots & & \ddots & \vdots \\ 1 & \tilde{x}_{N1} & \tilde{x}_{N2} & \dots & \tilde{x}_{Nd} \end{bmatrix} \begin{bmatrix} 1\tilde{x}_0 & & & & 0 \\ & 1 & & & \\ & 0 & \ddots & & \\ & & & \ddots & 1 \end{bmatrix} = X' D'$$

$N \times (d+1)$ $N \times (d+1)$ $\parallel D'$ $\det(d_{11} \times d_{dd})$

$$\rightarrow X'^T X' = (X' D')^T (X' D') = D'^T X^T X D' = D'^T X^T X D'$$

\uparrow
" D' diagonal
 $\therefore D'^T = D'$

" $\det(D') = 1\tilde{x}_0$
and $\exists (X^T X)^{-1} \therefore \det(X^T X) \neq 0$

\downarrow

$$\det(X'^T X') = \det(D'^T X^T X D') = \det(D') \det(X^T X) \det(D') = \det(D')^2 \det(X^T X) \neq 0$$

$\therefore X'^T X'$: invertible

$$\text{Thus, } \bar{w}_{wcy} = (X^T X)^{-1} X^T \bar{y} = \underbrace{(D^T X^T X D)^{-1}}_{\substack{A \\ B \\ C}} (X D)^T \bar{y} = D^{-1} (X^T X)^{-1} \underbrace{D^T D^{-1}}_{\substack{I \\ D^T D = I}} X^T \bar{y}$$

$\because D^T = D$
 $\therefore D^{-1} D^T = D^{-1} D = I$

$\because \det(D) \neq 0 \therefore D: \text{invertible}$
 $(ABC)^{-1} = C^{-1} B^{-1} A^{-1}$

$$= D^{-1} \underbrace{(X^T X)^{-1} X^T \bar{y}}_{\substack{\parallel \\ \bar{w}_{wcy}}}$$

Therefore, $\bar{w}_{wcy} = D^{-1} \bar{w}_{wcy}$

$$\rightarrow D = D^{-1} = \begin{bmatrix} 1 & 0 & & 0 \\ 0 & 1 & & 0 \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix} \quad \square$$

9. Using a new family of sigmoid hypotheses, $\tilde{h}(\vec{x}) = \frac{1}{2} \left(\frac{\tilde{w}^T \vec{x}}{\sqrt{1 + (\tilde{w}^T \vec{x})^2}} + 1 \right)$

We can write

$$\tilde{h}(\vec{x}) = \theta(\tilde{w}^T \vec{x}) \quad , \quad \text{where } \theta(s) = \frac{1}{2} \left(\frac{s}{\sqrt{1+s^2}} + 1 \right)$$

Also, we will have likelihood:

$$P(y|\vec{x}) = \begin{cases} \tilde{h}(\vec{x}) & \text{for } y=+1 \\ 1 - \tilde{h}(\vec{x}) & \text{for } y=-1 \end{cases}$$

substitute $\tilde{h}(\vec{x})$ by $\theta(\tilde{w}^T \vec{x})$:

$$P(y|\vec{x}) = \begin{cases} \theta(\tilde{w}^T \vec{x}) & \text{for } y=+1 \\ 1 - \theta(\tilde{w}^T \vec{x}) & \text{for } y=-1 \end{cases}$$

We can verify that:

$$\begin{aligned} \boxed{1 - \theta(s)} &= 1 - \frac{1}{2} \left(\frac{s}{\sqrt{1+s^2}} + 1 \right) = 1 - \frac{s}{2\sqrt{1+s^2}} - \frac{1}{2} = \frac{1}{2} - \frac{s}{2\sqrt{1+s^2}} \\ &= \frac{1}{2} \left(\frac{-s}{\sqrt{1+(-s)^2}} + 1 \right) = \boxed{\theta(-s)} \end{aligned}$$

Thus,

$$P(y|\vec{x}) = \begin{cases} \theta(\tilde{w}^T \vec{x}) & \text{for } y=+1 \\ \theta(-\tilde{w}^T \vec{x}) & \text{for } y=-1 \end{cases}$$

$$\Rightarrow P(y|\vec{x}) = \theta(y \cdot \tilde{w}^T \vec{x})$$

We can then derive the log likelihood:

$$\begin{aligned}\ln \prod_{i=1}^N P(y_i | \vec{x}_i) &= \sum_{i=1}^N \ln(P(y_i | \vec{x}_i)) \\ &= \sum_{i=1}^N \ln(\theta(y_i \cdot \vec{w}^T \vec{x}_i))\end{aligned}$$

To find the maximum log likelihood is equivalent to find the minimum cross entropy error:

$$\max_{\vec{w}} \sum_{i=1}^N \ln(\theta(y_i \cdot \vec{w}^T \vec{x}_i)) \geq \min_{\vec{w}} - \sum_{i=1}^N \ln(\theta(y_i \cdot \vec{w}^T \vec{x}_i))$$

By our definition of $\theta(s) = \frac{1}{2} \left(\frac{s}{\sqrt{1+s^2}} + 1 \right)$,

$$\begin{aligned}\ln(\theta(y \vec{w}^T \vec{x})) &= \ln \left(\frac{1}{2} \left(\frac{y \vec{w}^T \vec{x}}{\sqrt{1 + (y \vec{w}^T \vec{x})^2}} + 1 \right) \right) \\ &= \ln \frac{y \vec{w}^T \vec{x} + \sqrt{1 + (y \vec{w}^T \vec{x})^2}}{2 \sqrt{1 + (y \vec{w}^T \vec{x})^2}}\end{aligned}$$

\therefore the EM we want to minimize is:

$$\tilde{E}_M(\vec{w}) = \frac{1}{N} \left[- \sum_{i=1}^N \ln(\theta(y_i \vec{w}^T \vec{x}_i)) \right] = \frac{1}{N} \sum_{i=1}^N \ln \frac{2 \sqrt{1 + (y \vec{w}^T \vec{x})^2}}{y \vec{w}^T \vec{x} + \sqrt{1 + (y \vec{w}^T \vec{x})^2}}$$

To minimize $\tilde{E}_m(\vec{w})$, we need to find the place where $\nabla \tilde{E}_m(\vec{w}) = 0$

First we calculate $\nabla \mu(\theta(s))$

$$\nabla \mu(\theta(y_i \vec{w}^T \vec{x}_i)) = \frac{1}{\theta(y_i \vec{w}^T \vec{x}_i)} \cdot \nabla \theta(y_i \vec{w}^T \vec{x}_i)$$

calculate $\nabla \theta(s)$: $\frac{d}{ds} \left(\frac{s}{\sqrt{1+s^2}} \right) = \frac{\sqrt{1+s^2} - \frac{s^2}{\sqrt{1+s^2}}}{1+s^2} = \frac{(1+s^2) - s^2}{(1+s^2)^{3/2}} = \frac{1}{(1+s^2)^{3/2}}$

$$\frac{d}{ds} \left(\frac{s}{\sqrt{1+s^2}} \right) = \frac{\sqrt{1+s^2} - \frac{s^2}{\sqrt{1+s^2}}}{1+s^2} = \frac{(1+s^2) - s^2}{(1+s^2)^{3/2}} = \frac{1}{(1+s^2)^{3/2}}$$

$$\therefore \frac{d}{ds} \theta(s) = \frac{1}{2} \cdot (1+s^2)^{-\frac{3}{2}}$$

$$\rightarrow \nabla \theta(y_i \vec{w}^T \vec{x}_i) = \frac{1}{2} \frac{1}{[1 + (y_i \vec{w}^T \vec{x}_i)^2]^{\frac{3}{2}}} y_i \vec{x}_i = \frac{y_i \vec{x}_i}{2[1 + (\vec{w}^T \vec{x}_i)^2]^{\frac{3}{2}}}$$

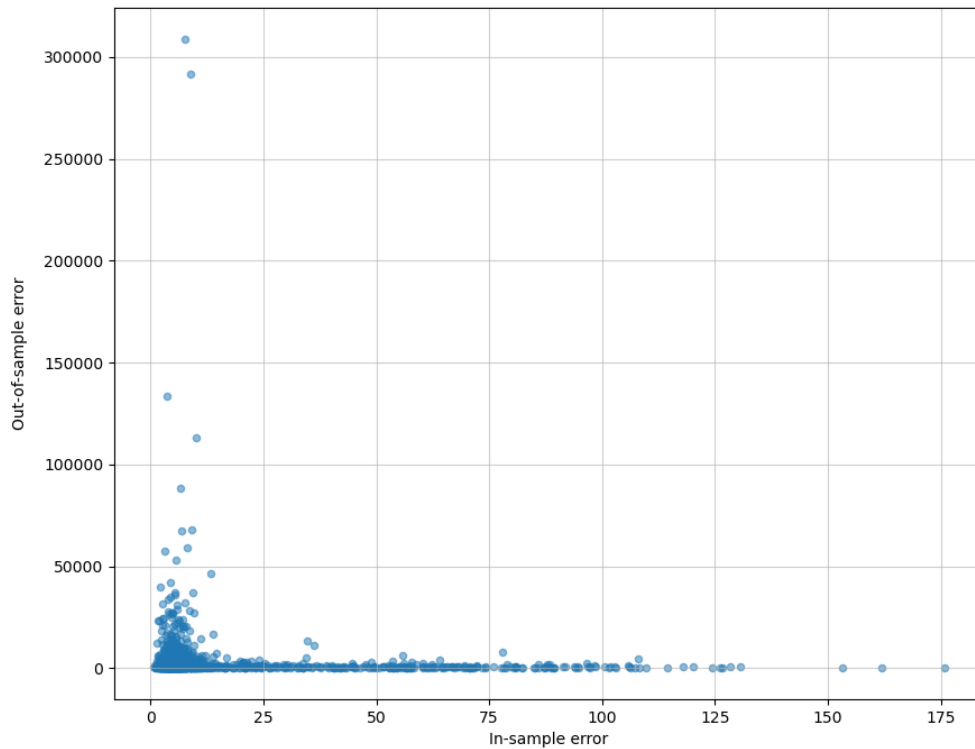
$$\rightarrow \nabla \mu(\theta(y_i \vec{w}^T \vec{x}_i)) = \frac{1}{\theta(y_i \vec{w}^T \vec{x}_i)} \cdot \frac{y_i \vec{x}_i}{2[1 + (\vec{w}^T \vec{x}_i)^2]^{\frac{3}{2}}}$$

$$\rightarrow \nabla \tilde{E}_m(\vec{w}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\theta(y_i \vec{w}^T \vec{x}_i)} \frac{y_i \vec{x}_i}{[1 + (\vec{w}^T \vec{x}_i)^2]^{\frac{3}{2}}}$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{y_i \vec{x}_i}{\left(\frac{y_i \vec{w}^T \vec{x}_i}{\sqrt{1 + (y_i \vec{w}^T \vec{x}_i)^2}} + 1 \right) [1 + (\vec{w}^T \vec{x}_i)^2]^{\frac{3}{2}}}$$

$$\theta(y_i \vec{w}^T \vec{x}_i) = \frac{1}{2} \left(\frac{y_i \vec{w}^T \vec{x}_i}{\sqrt{1 + (y_i \vec{w}^T \vec{x}_i)^2}} + 1 \right)$$

Question 10: Report



Explanation:

The first thing we would probably notice is that, there's 2 points in the upper left corner of the plot, which indicates a small in-sample error, and a relatively large out of sample error. I think this is the result that I would expect, since we only choose 32 from all the 8192 examples, using such a small subset of the dataset would make the model vary greatly, depending on the subset we choose, and if the subset contains outliers, or the examples chosen are not representative enough for the entire dataset, then overfitting to such examples would cause high out of sample error, meaning that the model is not well generalized. In addition to these 2 points, we can also see that in about 10 experiments, we have out of sample error greater than 50000, which is also

quite high. Most of the other points in the lower left corner still shows that the out of sample error is not that small, this result verifies our intuition.

Next, we can see that in the lower middle part, and the lower right corner, there are experiments having nearly 0 out of sample error but higher in sample errors (compared to those having only in sample errors less than 25.) For this part, this result actually matches that in theory, higher in sample error (due to underfitting) would often end up with lower out of sample error, because models like this generalize better than those that overfit. However, I still think that it's surprising to use 32 points to generalize well on the remaining points, but the amount of experiments like this is small, so I think maybe it's reasonable due to probability.

Last, observe the whole plot, we can see that the distribution varies a lot over all of the experiments. Even though we don't know the correlations between the 8192 examples, randomly choosing 32 from them still results in a high variability (since there would be $C(8192, 32)$ possibilities.) As mentioned before, this would cause our distribution to depend highly on the 32 samples we choose, so I don't think we can tell more by only conducting 1126 experiments.

Code:

Read in the data:

```

5  def read_libsvm_format(file_path):
6      X, y = [], []
7      with open(file_path, 'r') as f:
8          for line in f:
9              parts = line.strip().split()
10             y.append(float(parts[0]))
11             features = {}
12             for item in parts[1:]:
13                 index, value = item.split(":")
14                 features[int(index)] = float(value)
15             X.append(features)
16         return X, y
17
18     # aim: read the data from the file 'cpusmall_scale'
19     X, y = read_libsvm_format('cpusmall_scale')
20
21     # explain: the data X returned is in the form of a list of dictionaries, where each dictionary contains the features of a data point
22     # ex: X[0] = {1: -0.993496, 2: -0.993043, 3: -0.850291, 4: -0.963479, 5: -0.960727, 6: -0.900596, 7: -0.96642, 8: -0.863996, 9: -0.606175, 10: -0.999291, 11: 0.08118}
23
24     # explain: the label y is a list of floats
25     # ex: y[0] = 90.0
26

```

Calculate w_{LIN} and E_{in} :

```

27 N = 32
28 feature_amount = 12
29 in_out_sample_error = []
30 # aim: perform linear regression for 1126 times, calculate the in-sample error and estimate the out of sample error
31 for experiment in tqdm(range(1126)):
32     np.random.seed(experiment)
33
34     random_sample_indices = np.random.choice(len(X), N, replace=False)
35     X_sample = [X[i] for i in random_sample_indices]
36     y_sample = [y[i] for i in random_sample_indices]
37
38     # subaim: convert each sample data points (which is a dictionary) to a ndarray and save in X_sample_array
39     X_sample_array = []
40     for sample in X_sample:
41         input_vector = np.concatenate((np.array([1]), np.zeros(feature_amount)))
42         # explain: originally we need subtract 1 because indices in libsvm format start from 1, but now we add the 0-th element (1) to each input vector
43         for index, value in sample.items():
44             input_vector[index] = value
45         X_sample_array.append(input_vector)
46
47     # subaim: convert X_sample_array to a matrix, and y_sample to an array
48     X_mat = np.array(X_sample_array)
49     y_array = np.array(y_sample)
50
51     # subaim: perform linear regression
52     # explain: we do linear regression by using the normal equation ( $w = (X^T X)^{-1} X^T y$ )
53     # explain: we use np.linalg.inv() calculate the inverse of  $X^T X$ , and use @ to denote matrix multiplication
54     w = np.linalg.inv(X_mat.T @ X_mat) @ X_mat.T @ y_array
55
56     # subaim: calculate the in-sample error
57     # explain: the in-sample error is the mean squared error of the N points
58     in_sample_error = np.mean((X_mat @ w - y_array) ** 2)
59

```

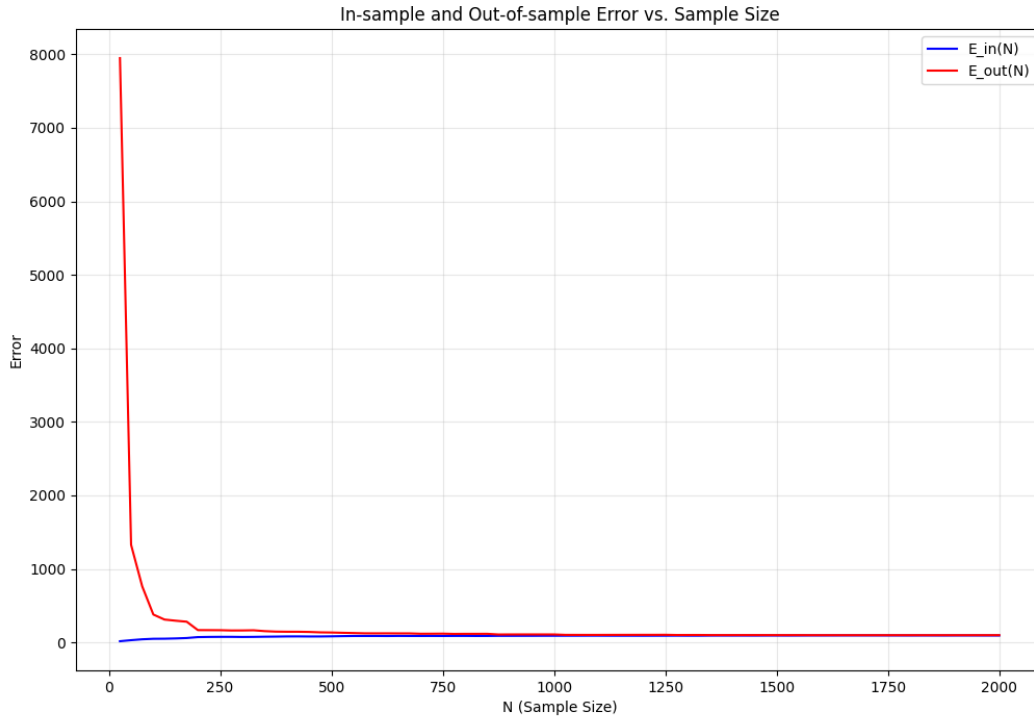
Estimate E_{out} :

```

60 # subaim: estimate the out of sample error
61 # explain: the out of sample error is estimated by averaging the squared error of the rest of the 8192 - N data points
62 out_of_sample_error = 0
63 for i in range(len(X)):
64     if i not in random_sample_indices:
65         unseen_input_vector = np.concatenate((np.array([1]), np.zeros(feature_amount)))
66         for index, value in X[i].items():
67             unseen_input_vector[index] = value
68         out_of_sample_error += (unseen_input_vector @ w - y[i]) ** 2
69 out_of_sample_error /= (len(X) - N)
70 in_out_sample_error.append((in_sample_error, out_of_sample_error))

```

Question 11: Report



Explanations:

In this problem, we gradually increase the sample size (until 2000, which is nearly $\frac{1}{4}$ of our dataset), and calculate the average in sample error and out of sample error over only 16 experiments.

From the graph we can see that when the sample size is under about 100, we cannot generalize well since this is represented by the high out of sample error. But after we met a certain sample size (which is maybe 100 or 150 or some amount around them), we have enough data to represent the whole dataset by the sample, so that we have enough information to construct our model.

Also, if we really look at the examples, we could find that the number of features is 12, which is enough to provide sufficient degrees of freedom to fit the training data. In a geometrical view, we can think of finding a 12 dimensional hyperplane to split the N data points, in a 13 dimensional space (since we add 1 to the input vector.) In lower dimensional spaces, like 2D or 3D spaces, the degree of freedom is smaller, therefore it is hard to find a line or a 2D plane that can pass through all the points, but for a 12 dimensional hyperplane, it is easier to do so. Thus, it is possible to find the “optimal” hyperplane to fit best to the sample points, producing nearly 0 in sample error in all cases.

Code:

(The reading data part is ignored in this report since it's the same in the previous question.)

The first thing to do is to define a list of all possible N s that we need to iterate through (25, 50, 75,...,2000), then for each value of N , we conduct the experiment 16 times.

```

21 N_list = list(range(25, 2025, 25))
22 feature_amount = 12
23 seed = 0
24
25 # aim: for each sample size N, perform linear regression for 16 times, calculate the average in-sample error and estimate the average out of sample error
26 avg_in_error_for_each_N = []
27 avg_out_error_for_each_N = []
28 for N in tqdm(N_list):
29     in_out_sample_error = []
30     # subaim: perform linear regression for 16 times, calculate the in-sample error and estimate the out of sample error
31     for experiment in range(16):

```

The random choosing process is done without replacement. We save the chosen examples in X_sample and their corresponding y in y_sample :

```

31     for experiment in range(16):
32         seed += 1
33         np.random.seed(seed)
34
35         random_sample_indices = np.random.choice(len(X), N, replace=False)
36         X_sample = [X[i] for i in random_sample_indices]
37         y_sample = [y[i] for i in random_sample_indices]

```

Then convert data type:

```

37 # subsubaim: convert each sample data points (which is a dictionary) to a ndarray and save in X_sample_array
38 X_sample_array = []
39 for sample in X_sample:
40     input_vector = np.concatenate((np.array([1.0]), np.zeros(feature_amount, dtype=float)))
41     for index, value in sample.items():
42         input_vector[index] = float(value)
43     X_sample_array.append(input_vector)
44
45 # subsubaim: convert X_sample_array to a matrix, and y_sample to an array
46 X_mat = np.array(X_sample_array, dtype=float)
47 y_array = np.array(y_sample, dtype=float)

```

Calculate in sample error and out of sample error:

```

51 # subsubaim: perform linear regression
52 # explain: we do linear regression by using the normal equation (w = (X^TX)^(-1)X^Ty)
53 # explain: we use np.linalg.inv() calculate the inverse of X^TX, and use @ to denote matrix multiplication
54 w = np.linalg.inv(X_mat.T @ X_mat) @ X_mat.T @ y_array
55
56 # subsubaim: calculate the in-sample error
57 # explain: the in-sample error is the mean squared error of the N points
58 in_sample_error = np.mean((X_mat @ w - y_array) ** 2)
59
60 # subsubaim: estimate the out of sample error
61 # explain: the out of sample error is estimated by averaging the squared error of the rest of the 8192 - N data points
62 out_of_sample_error = 0.0
63 for i in range(len(X)):
64     if i not in random_sample_indices:
65         unseen_input_vector = np.concatenate((np.array([1.0]), np.zeros(feature_amount, dtype=float)))
66         for index, value in X[i].items():
67             unseen_input_vector[index] = float(value)
68         out_of_sample_error += (unseen_input_vector @ w - y[i]) ** 2
69 out_of_sample_error /= (len(X) - N)
70 in_out_sample_error.append((in_sample_error, out_of_sample_error))

```

The in_out_sample_error is a list containing the result of the 16 experiments for any given N.

We then calculate $\bar{E}_{in}(N)$ and $\bar{E}_{out}(N)$:

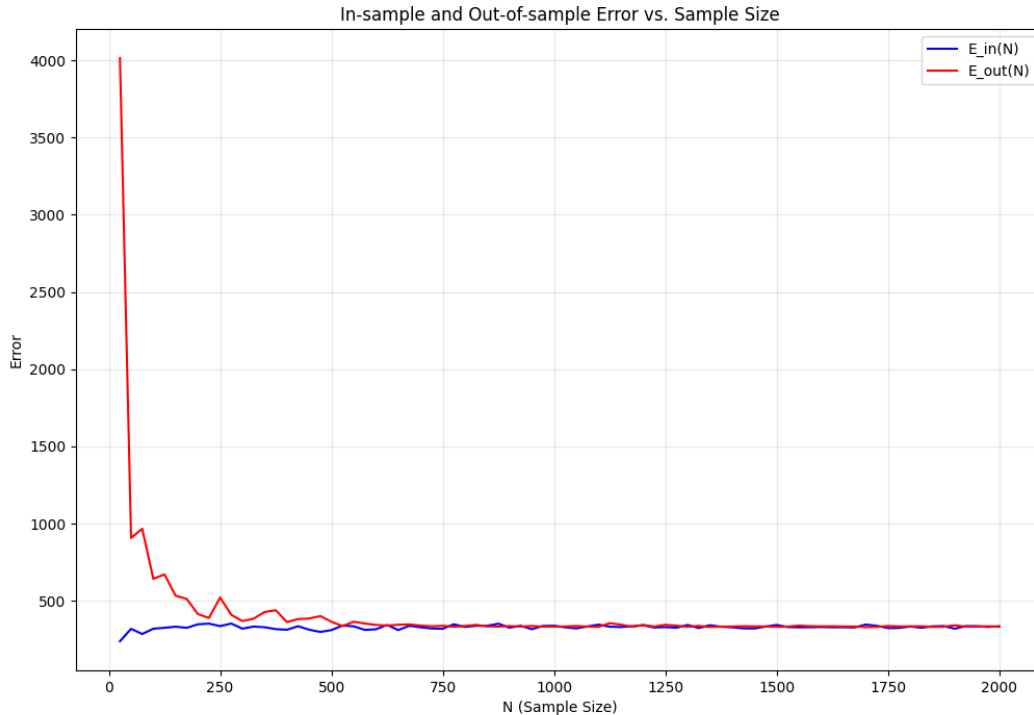
```

72 # subaim: calculate the average in-sample error and out of sample error for each N
73 avg_in_sample_error = 0.0
74 avg_out_of_sample_error = 0.0
75 for each_exp_error_result in in_out_sample_error:
76     avg_in_sample_error += each_exp_error_result[0]
77     avg_out_of_sample_error += each_exp_error_result[1]
78 avg_in_sample_error /= 16
79 avg_out_of_sample_error /= 16
80 avg_in_error_for_each_N.append(avg_in_sample_error)
81 avg_out_error_for_each_N.append(avg_out_of_sample_error)

```

The last part is to plot the result, but it's similar to the previous question.

Question 12: Report



Explanations:

The difference between this question and the previous question 11 is that we only preserve the first 2 features instead of using all the 12 features. We could see that in the resulting figure, the highest out of sample value is lower comparing to the previous question. Although they both gradually decrease as the sample size increases.

This is because when we're in a higher dimensional space, we would have to face the problem of generalization. Although using all the 12 features results in a space that can form more flexible models, meaning that is more possible to find a hyperplane that can separate our training data, making the in sample error to be about 0, when the model is used to evaluate the

unseen data, over-parameterized models may cause overfitting. In this question we decrease the amount of features to 2, indicates a choice toward generalization than approximation.

However, we can see that the in sample error is a nonzero value, even though the value is low for over all possible N . This is also the result as we reduce the dimension of the space to 3 (including the fixed zeroth coordinate), our free parameter is reduced to 2, so that it is hard to make the training data linearly separable.

Nonetheless, reducing the number of features would make d_{VC} small, and as the amount of data needed is proportional to the VC dimension, we could reach the same level of generalization error using less examples, and this is reflected in the drastic decrease in the out of sample error in a low value of N .

Code:

The initialization part (getting the dataset, and generate random indices is the same as the previous 2 questions), so the first difference occurs when we save the input vectors, as we only save the first 2 features instead of 12:

```
38 # subsubaim: convert each sample data points (which is a dictionary) to a ndarray and save in X_sample_array
39 X_sample_array = []
40 for sample in X_sample:
41     input_vector = [1, sample.get(1, 0), sample.get(2, 0)]
42     X_sample_array.append(input_vector)
```

The following part is almost the same as the previous question, except modifying it to be more clear and simple, but the concept behind is all the same:


```
44     # subsubaim: convert X_sample_array to a matrix, and y_sample to an array
45     X_mat = np.array(X_sample_array)
46     y_array = np.array(y_sample)
47
48     # subsubaim: perform linear regression
49     w = np.linalg.inv(X_mat.T @ X_mat) @ X_mat.T @ y_array
50
51     # subsubaim: calculate the in-sample error
52     in_sample_error = np.mean((X_mat @ w - y_array) ** 2)
53
54     # subsubaim: estimate the out of sample error
55     out_sample_indices = list(set(range(len(X)) - set(random_sample_indices))
56     X_out = np.array([[1, X[i].get(1, 0), X[i].get(2, 0)] for i in out_sample_indices])
57     y_out = np.array([y[i] for i in out_sample_indices])
58     out_of_sample_error = np.mean((X_out @ w - y_out) ** 2)
59     in_out_sample_error.append((in_sample_error, out_of_sample_error))
60
61     # subaim: calculate the average in-sample error and out of sample error for each N
62     avg_in_sample_error = np.mean([error[0] for error in in_out_sample_error])
63     avg_out_of_sample_error = np.mean([error[1] for error in in_out_sample_error])
64
65     avg_in_error_for_each_N.append(avg_in_sample_error)
66     avg_out_error_for_each_N.append(avg_out_of_sample_error)
```

13. By the definition of $B(N, k)$, it is the max number of dichotomies on N points,
 \times no subset of size k of the N points can be shattered
 by these dichotomies

If we consider the case that we have a set S , containing dichotomies of N points
 that have number of $\{-1\}$'s $\leq k-1$

x_1	x_2	x_3	x_4	...	x_{N-1}	x_N	containing # $\{-1\}$
+1	+1	+1	+1		+1	+1	0
-1	+1	+1	+1		+1	+1	} $1 \rightarrow \binom{N}{1}$ dichotomies
+1	-1	+1	+1		+1	+1	
:		-1			:	:	
:					:	:	
:					-1		
+1	+1					-1	} $2 \rightarrow \binom{N}{2}$ dichotomies
-1	-1						
-1		-1					
:							:

Similarly, it is easy to see that we will have $\sum_{i=0}^{k-1} \binom{N}{i}$ dichotomies in S

Suppose \exists subset of size k of the N points:

$$X' = \{\tilde{x}_j, \tilde{x}_{j+1}, \dots, \tilde{x}_{j+k-1}\} \subseteq \{\tilde{x}_1, \dots, \tilde{x}_N\}$$

\downarrow
size k

\nexists X' can be shattered by these $\sum_{i=0}^N \binom{N}{i}$ dichotomies,
 \rightarrow then $[\tilde{x}_j, \tilde{x}_{j+1}, \dots, \tilde{x}_{j+k-1}] = [\underbrace{-1, -1, \dots, -1}_{k \text{ } \{-1\}'s}] \in S$ ~~X~~

" By our assumption, S has
dichotomies with number of $\{-1\}'s$
 $\leq k-1$

Thus, we obtain the result:

$$B(n, k) \geq \sum_{i=0}^{k-1} \binom{N}{i} \quad \square$$