

HTML HW3
B12901022 廖冠豪

In this problem we use $\{-1, +1\}$ to denote boolean values (outcome of binary classification)

5

We consider hypotheses $h_1, h_2 : \mathbb{R} \rightarrow \{-1, +1\}$, where

$$\begin{aligned}h_1(x) &= \text{sign}(x) \\h_2(x) &= -\text{sign}(x)\end{aligned}$$

and hypotheses sets $\mathcal{H}_1 = \{h_1\}$ and $\mathcal{H}_2 = \{h_2\}$.

As both \mathcal{H}_1 and \mathcal{H}_2 only contain one hypothesis, their VC dimension are both 0.

$$d_{VC}(\mathcal{H}_1) = d_{VC}(\mathcal{H}_2) = 0$$

Now consider the hypotheses set $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2$, and $x = 1$. We have $h_1(x) = +1$ and $h_2(x) = -1$, so \mathcal{H} can shatter a set of one data vector and its VC dimension is hence greater than or equal to 1.

We conclude that

$$d_{VC}(\mathcal{H}_1 \cup \mathcal{H}_2) \geq 1 > d_{VC}(\mathcal{H}_1) + d_{VC}(\mathcal{H}_2) = 0$$

Therefore we have disproved the statement.

In this problem we use $\{-1, +1\}$ to denote boolean values (outcome of binary classification)

6

In the super-market case, we hope to find $f(\mathbf{x})$ that minimizes

$$\mathbb{E}_{\mathbf{x}, y \sim P(\mathbf{x}, y)} [10P(+1|\mathbf{x})\mathbb{I}[f(\mathbf{x}) = -1] + P(-1|\mathbf{x})\mathbb{I}[f(\mathbf{x}) = +1]]$$

We can see that for any \mathbf{x} exactly one out of $f(\mathbf{x}) = +1$ and $f(\mathbf{x}) = -1$ is true, so it would be optimal to choose an α such that the one that contributes to a smaller value in \mathbb{E} is always chosen (i.e. $f(\mathbf{x}) = -1$ if $10P(+1|\mathbf{x})$ is greater than $P(-1|\mathbf{x})$, and vice versa)

$$\begin{cases} P(+1|\mathbf{x}) + P(-1|\mathbf{x}) = 1 \\ 10P(+1|\mathbf{x}) \geq P(-1|\mathbf{x}) \end{cases} \implies P(+1|\mathbf{x}) \geq \frac{1}{11}$$

By the analysis above, we can see that the optimal choice of α is $\frac{1}{11}$, and we have the mini-target

$$f_{\text{MKT}}(\mathbf{x}) = \text{sign}(P(y = +1|\mathbf{x}) - \frac{1}{11})$$

In this problem we use $\{-1, +1\}$ to denote boolean values (outcome of binary classification)

7

$$\begin{aligned} E_{\text{out}}^{(1)}(h) &= \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x}) \llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket \\ E_{\text{out}}^{(2)}(h) &= \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x}) (P(-1|\mathbf{x}) \llbracket h(\mathbf{x}) = +1 \rrbracket + P(+1|\mathbf{x}) \llbracket h(\mathbf{x}) = -1 \rrbracket) \\ E_{\text{out}}^{(2)}(f) &= \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x}) (P(-1|\mathbf{x}) \llbracket f(\mathbf{x}) = +1 \rrbracket + P(+1|\mathbf{x}) \llbracket f(\mathbf{x}) = -1 \rrbracket) \end{aligned}$$

We see that

$$E_{\text{out}}^{(1)}(h) + E_{\text{out}}^{(2)}(f) = \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x}) (P(-1|\mathbf{x}) \llbracket f(\mathbf{x}) = +1 \rrbracket + P(+1|\mathbf{x}) \llbracket f(\mathbf{x}) = -1 \rrbracket + \llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket)$$

Given $h(\mathbf{x}), f(\mathbf{x})$, define

$$\begin{aligned} a(\mathbf{x}) &= P(-1|\mathbf{x}) \llbracket f(\mathbf{x}) = +1 \rrbracket + P(+1|\mathbf{x}) \llbracket f(\mathbf{x}) = -1 \rrbracket + \llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket \\ b(\mathbf{x}) &= P(-1|\mathbf{x}) \llbracket h(\mathbf{x}) = +1 \rrbracket + P(+1|\mathbf{x}) \llbracket h(\mathbf{x}) = -1 \rrbracket \end{aligned}$$

We then consider the four different cases for $h(\mathbf{x}), f(\mathbf{x})$.

Case 1: $h(\mathbf{x}) = +1, f(\mathbf{x}) = +1$

$$a(\mathbf{x}) = P(-1|\mathbf{x})$$

$$b(\mathbf{x}) = P(-1|\mathbf{x})$$

Case 2: $h(\mathbf{x}) = +1, f(\mathbf{x}) = -1$

$$\begin{aligned}a(\mathbf{x}) &= P(+1|\mathbf{x}) + 1 \\b(\mathbf{x}) &= P(-1|\mathbf{x})\end{aligned}$$

Case 3: $h(\mathbf{x}) = -1, f(\mathbf{x}) = -1$

$$\begin{aligned}a(\mathbf{x}) &= P(+1|\mathbf{x}) \\b(\mathbf{x}) &= P(+1|\mathbf{x})\end{aligned}$$

Case 4: $h(\mathbf{x}) = -1, f(\mathbf{x}) = +1$

$$\begin{aligned}a(\mathbf{x}) &= P(-1|\mathbf{x}) + 1 \\b(\mathbf{x}) &= P(+1|\mathbf{x})\end{aligned}$$

We can see that for all possible cases we have $a(\mathbf{x}) \geq b(\mathbf{x})$.

Hence

$$E_{\text{out}}^{(1)}(h) + E_{\text{out}}^{(2)}(f) = \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x})a(\mathbf{x}) \geq E_{\text{out}}^{(2)}(h) = \sum_{\mathbf{x} \sim P(\mathbf{x})} P(\mathbf{x})b(\mathbf{x})$$

Q.E.D.

Note that when dealing with a continuous distribution of \mathbf{x} , we only need to change the summation to integral, and the following argument applies.

Assuming that $X^T X$ is invertible, we can express \mathbf{w}_{LIN} in the following form

$$\mathbf{w}_{\text{LIN}} = (X^T X)^{-1} X^T \mathbf{y}$$

Now if we replace x_0 with 1126, we get the new X matrix

$$X' = X D'$$

Where D' is the diagonal matrix with

$$D'_{ij} = \begin{cases} 1126 & i = j = 0 \\ 1 & i = j \neq 0 \\ 0 & \text{else} \end{cases}$$

Clearly D'^{-1} exists, and D' is symmetric, so

$$(X'^T X')(D'^{-1}(X^T X)^{-1} D'^{-1}) = D'(X^T X) D' D'^{-1} (X^T X)^{-1} D'^{-1} = I$$

So $(X'^T X')$ is invertible with

$$(X'^T X')^{-1} = D'^{-1} (X^T X)^{-1} D'^{-1}$$

Hence we can repeat the linear regression procedure and obtain

$$\begin{aligned} \mathbf{w}_{\text{LUCKY}} &= (X'^T X')^{-1} X'^T \mathbf{y} \\ &= D'^{-1} (X^T X)^{-1} D'^{-1} D' X^T \mathbf{y} \\ &= D'^{-1} (X^T X)^{-1} D'^{-1} X^T \mathbf{y} \\ &= D'^{-1} \mathbf{w}_{\text{LIN}} \end{aligned}$$

Hence we have proved the statement and found the diagonal matrix $D = D'^{-1}$

$$D_{ij} = \begin{cases} \frac{1}{1126} & i = j = 0 \\ 1 & i = j \neq 0 \\ 0 & \text{else} \end{cases}$$

In this problem we use $\{-1, +1\}$ to denote boolean values (outcome of binary classification)

9

Let $f(\mathbf{x})$ be the target function we want to approximate with $\tilde{h}(\mathbf{x})$.

Consider $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N), \}$, $y_i \in \{-1, +1\}$.

The probability that f generates \mathcal{D} is

$$P(\mathbf{x}_1)f(\mathbf{x}_1) \times P(\mathbf{x}_2)(1 - f(\mathbf{x}_2)) \times \dots \times P(\mathbf{x}_N)(1 - f(\mathbf{x}_N))$$

The likelihood that \tilde{h} generates \mathcal{D} is

$$P(\mathbf{x}_1)h(\mathbf{x}_1) \times P(\mathbf{x}_2)(1 - h(\mathbf{x}_2)) \times \dots \times P(\mathbf{x}_N)(1 - h(\mathbf{x}_N))$$

Also,

$$\begin{aligned} 1 - \tilde{h}(\mathbf{x}) &= 1 - \frac{1}{2} \left(\frac{\mathbf{w}^T \mathbf{x}}{\sqrt{1 + (\mathbf{w}^T \mathbf{x})^2}} + 1 \right) \\ &= \frac{1}{2} \left(\frac{-\mathbf{w}^T \mathbf{x}}{\sqrt{1 + (\mathbf{w}^T \mathbf{x})^2}} + 1 \right) \\ &= \tilde{h}(-\mathbf{x}) \end{aligned}$$

Therefore the likelihood for some hypothesis \tilde{h} is proportional to

$$\prod_{n=1}^N \tilde{h}(y_n \mathbf{x}_n)$$

We hope to find

$$\begin{aligned}
\arg \max_{\mathbf{w}} \prod_{n=1}^N \tilde{h}(y_n \mathbf{x}_n) &= \arg \max_{\mathbf{w}} \prod_{n=1}^N \left(\frac{1}{2} \left(\frac{y_n (\mathbf{w}^T \mathbf{x}_n)}{\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} + 1 \right) \right) \\
&= \arg \max_{\mathbf{w}} \ln \left(\prod_{n=1}^N \frac{1}{2} \left(\frac{y_n (\mathbf{w}^T \mathbf{x}_n)}{\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} + 1 \right) \right) \\
&= \arg \max_{\mathbf{w}} \sum_{n=1}^N \ln \left(\frac{1}{2} \left(\frac{y_n (\mathbf{w}^T \mathbf{x}_n)}{\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} + 1 \right) \right) \\
&= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N -\ln \left(\frac{1}{2} \left(\frac{y_n (\mathbf{w}^T \mathbf{x}_n)}{\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} + 1 \right) \right) \\
&= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \text{err}(\mathbf{w}, \mathbf{x}_n, y_n)
\end{aligned}$$

Hence we've found the error function

$$\begin{aligned}
\text{err}(\mathbf{w}, \mathbf{x}_n, y_n) &= -\ln \left(\frac{1}{2} \left(\frac{y_n \mathbf{w}^T \mathbf{x}_n}{\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} + 1 \right) \right) \\
&= \ln \frac{2\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}}{y_n \mathbf{w}^T \mathbf{x}_n + \sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}}
\end{aligned}$$

Therefore

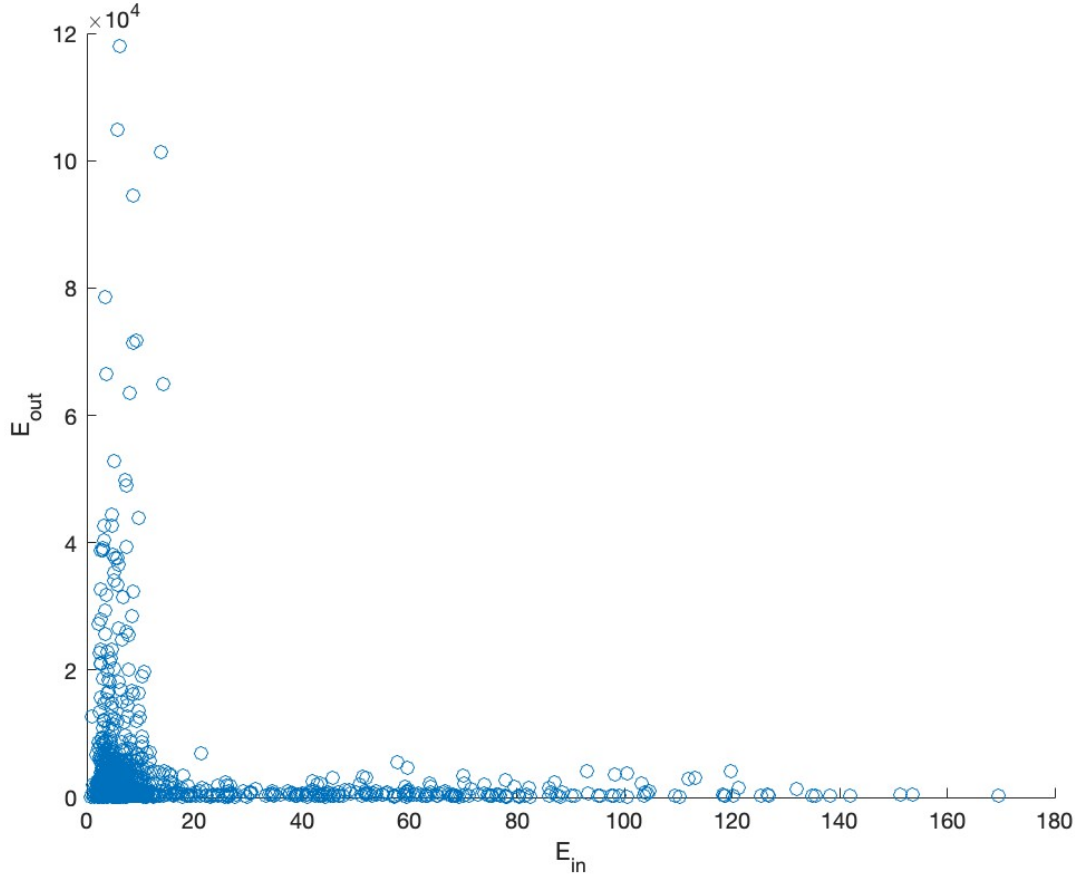
$$\tilde{E}_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{2\sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}}{y_n \mathbf{w}^T \mathbf{x}_n + \sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2}} \right)$$

We can then compute $\nabla \tilde{E}_{\text{in}}(\mathbf{w})$

$$\begin{aligned}
(\nabla \tilde{E}_{\text{in}}(\mathbf{w}))_i &= \frac{\partial \tilde{E}_{\text{in}}(\mathbf{w})}{\partial \mathbf{w}_i} \\
&= \sum_{n=1}^N \frac{\partial \ln \circ}{\partial \circ} \frac{\partial \circ}{\partial \square} \frac{\partial \square}{\partial \mathbf{w}_i} \quad \left(\circ = \left(\frac{2\sqrt{1 + \square^2}}{y_n \square + \sqrt{1 + \square^2}} \right), \square = \mathbf{w}^T \mathbf{x}_n \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{\circ} \right) \left(-\frac{2y_n}{\sqrt{1 + \square^2} (y_n \square + \sqrt{1 + \square^2})^2} \right) (\mathbf{x}_{n,i}) \\
&= \frac{1}{N} \sum_{n=1}^N \left(-\frac{y_n \mathbf{x}_{n,i}}{(1 + (\mathbf{w}^T \mathbf{x}_n)^2)(y_n \mathbf{w}^T \mathbf{x}_n + \sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2})} \right)
\end{aligned}$$

Therefore we have

$$\nabla \tilde{E}_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \left(-\frac{y_n \mathbf{x}_n}{(1 + (\mathbf{w}^T \mathbf{x}_n)^2)(y_n \mathbf{w}^T \mathbf{x}_n + \sqrt{1 + (\mathbf{w}^T \mathbf{x}_n)^2})} \right)$$



From the scatter plot we can see E_{out} is significantly larger than E_{in} , for the majority of data points E_{in} ranges between 0 to 2×10^4 and the maxima for E_{in} is about 1.2×10^5 . On the other hand, for most data points E_{out} is less than 40, and the maxima of E_{out} is no less than 180.

This difference between E_{out} and E_{in} is consistent with theory, in lecture we're introduced with the equations

$$\overline{E_{out}} = \text{noise level} \cdot \left(1 + \frac{d+1}{N}\right)$$

$$\overline{E_{in}} = \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right)$$

The expected generalization error is $\frac{2(d+1)}{N}$. In this case $N = 32$ is relatively small, so the difference between E_{out} and E_{in} is quite large.


```

0  data_vec = zeros(8192, 13);
1  y = zeros(1, 8192);
2  E_in = zeros(1126, 1);
3  E_out = zeros(1126, 1);
4  file = fopen('./data.txt');
5
6  line = fgetl(file);
7  line_count = 0;
8  while ischar(line)
9      line_string = string(line);
10     parts = strsplit(line_string, ' ');
11     %disp(parts);
12     y(line_count + 1) = str2double(parts(1));
13     x = zeros(13, 1);
14     x(1) = 1;
15     %disp(length(parts));
16     for j = 2:13
17         foo = strsplit(parts(j), ':');
18         x(j) = str2double(foo(2));
19     end
20     data_vec(line_count + 1, :) = x;
21     line_count = line_count + 1;
22     disp(line_count);
23     line = fgetl(file);
24 end
25 %disp(data_vec(8192, :));
26
27
28 count = 1;
29 while count <= 1126
30     in = 0;
31     out = 0;
32     indices = zeros(1, 32);
33     y_vec = zeros(32, 1);
34     mat = zeros(32, 13);
35     for i = 1:32
36         r = randi([1, 8192]);
37         indices(i, 1) = r;
38         mat(i, :) = data_vec(r, :);
39         y_vec(i, 1) = y(1, r);
40     end
41     w = pinv(mat) * y_vec;
42     %disp(w);
43     for i = 1:8192
44         err = (dot(w, data_vec(i, :)) - y(1, i)) ^ 2;
45         if ismember(i, indices(1, :))
46             in = in + err;
47         else
48             out = out + err;
49         end
50     end
51     in = in / 32;
52     out = out / 8160;

```

NORMAL main P10.m utf-8 matlab Top 1:5

Code snapshot:



From the figure we see that for small N , as in problem 10, the difference between E_{out} and E_{in} is large, with E_{out} reaching almost 7000 and E_{in} approximately 27 for $N = 25$.

As N increases, E_{out} decreases rapidly until it reaches the same level as E_{in} (around 100). During the process, E_{in} also increases slightly. At approximately $N = 400$, E_{out} and E_{in} become almost identical. For larger N , E_{out} and E_{in} are very stable, and their value doesn't change much as N increases to larger values.

This result is also consistent with theory, from the two equations

$$\overline{E_{out}} = \text{noise level} \cdot \left(1 + \frac{d+1}{N}\right)$$

$$\overline{E_{in}} = \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right)$$

we see that as N becomes larger, both E_{out} and E_{in} converges to the noise level σ^2 . This phenomenon is clearly reflected in the figure, as E_{out} and E_{in} approaches the same value for large N .

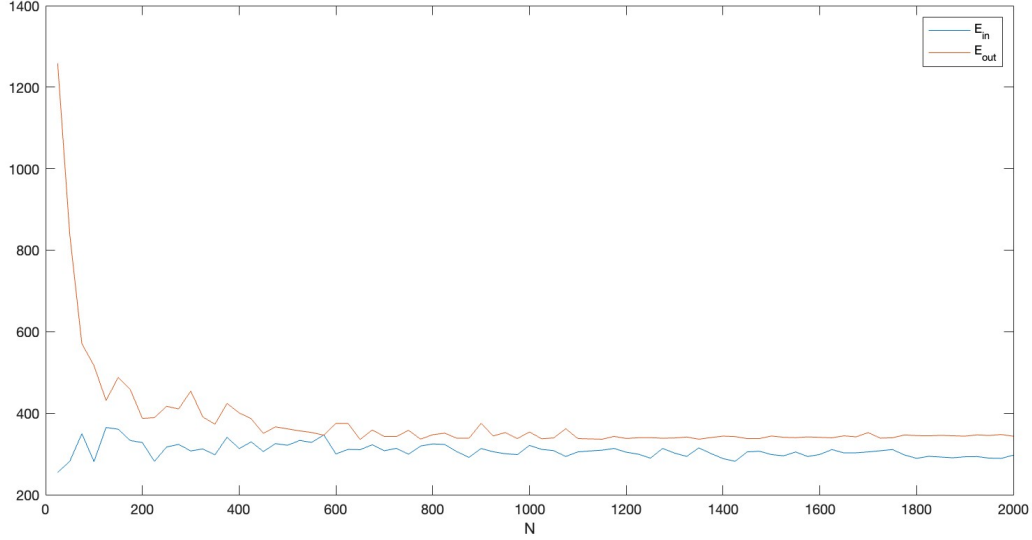
```

5     disp(line_count);
6     line = fgetl(file);
7 end
8 %disp(data_vec(8192, :));
9
10
11
12 for N = 25:25:2000
13     count = 1;
14     iin = 0;
15     oout = 0;
16     while count <= 16
17         in = 0;
18         out = 0;
19         indices = zeros(1, N);
20         y_vec = zeros(N, 1);
21         mat = zeros(N, 13);
22         for l = 1:N
23             r = randi([1, 8192]);
24             indices(l, 1) = r;
25             mat(l, :) = data_vec(r, :);
26             y_vec(l, 1) = y(1, r);
27         end
28         w = pinv(mat) * y_vec;
29         %disp(w);
30         for l = 1:8192
31             err = (dot(w, data_vec(l, :)) - y(1, l)) ^ 2;
32             if ismember(l, indices(1, :))
33                 in = in + err;
34             else
35                 out = out + err;
36             end
37         end
38         in = in / N;
39         out = out / (8192 - N);
40         iin = iin + in;
41         oout = oout + out;
42         count = count + 1;
43     end
44     oout = oout / 16;
45     iin = iin / 16;
46     E_out(N/25, 1) = oout;
47     E_in(N/25, 1) = iin;
48     disp(N);
49 end
50 vec = 25:25:2000;
51 plot(vec, E_in);
52 hold on;
53 plot(vec, E_out);
54 legend('E_in', 'E_out');
55 xlabel('N');
56 hold off;

```

NORMAL main P11.m utf-8 36% 28:1

Code snapshot:



The figure above is mostly similar to that in problem 10 in their trend of growth(increase/decrease). But there are some differences.

Firstly we see that the value to which E_{out} and E_{in} converge is different. This may be due to the noise level term σ^2 in the equation. For this set the noise level for each dimension of \mathbf{x} may be different, so training with different dimensions of \mathbf{x} may lead to different σ^2 values. Secondly, the difference between E_{in} and E_{out} for small N and large N (stable value) is larger in problem 11. We can see that in problem 11 E_{out} decreased by more than 99% ($\approx 7000 \rightarrow \approx 100$), and E_{in} increased by about 300% ($\approx 27 \rightarrow \approx 100$), whereas in this problem E_{out} only decreased by approximately 70% ($\approx 1250 \rightarrow \approx 400$) and E_{in} only increased by about 100% ($\approx 200 \rightarrow \approx 400$). This is due to the difference in d .

From the equations

$$\begin{aligned}\overline{E_{out}} &= \text{noise level} \cdot \left(1 + \frac{d+1}{N}\right) \\ \overline{E_{in}} &= \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right)\end{aligned}$$

we see that the difference between E_{in} and E_{out} at small N and their value at large N is $\frac{d+1}{N}$, hence for larger d and same N , E_{in} and E_{out} differs more with their asymptotic value. This is consistent with the cases for problem 11 and 12, as in problem 11 $d = 13$ and in problem 12 $d = 3$.

```
HW3: nvim HW3.tex (-zsh) 361 code: nvim P12.m (-zsh) 362 +
0 data_vec = zeros(8192, 3);
1 y = zeros(1, 8192);
2 E_in = zeros(80, 1);
3 E_out = zeros(80, 1);
4 file = fopen("./data.txt");
5
6 line = fgetl(file);
7 line_count = 0;
8 while ischar(line)
9     line_string = string(line);
10    parts = strsplit(line_string, ' ');
11    %disp(parts);
12    y(line_count + 1) = str2double(parts(1));
13    x = zeros(3, 1);
14    x(1) = 1;
15    %disp(length(parts));
16    for j = 2:3
17        foo = strsplit(parts(j), ':');
18        x(j) = str2double(foo(2));
19    end
20    data_vec(line_count + 1, :) = x;
21    line_count = line_count + 1;
22    disp(line_count);
23    line = fgetl(file);
24 end
25 %disp(data_vec(8192, :));
26
27
28
29 for N = 25:25:2000
30     count = 1;
31     iin = 0;
32     oout = 0;
33     while count <= 16
NORMAL  main P12.m utf-8 < matlab Top 1:1
```

Code snapshot:

In this problem we use $\{-1, +1\}$ to denote boolean values (outcome of binary classification)

13

This problem is done in collaboration with B12901035 鄭宇彥

Let \mathcal{S} be the set that contains all 2^N dichotomies of some N data vectors $\mathbf{x} \in \mathcal{D}$, that is, \mathcal{S} contains all 2^N distinct boolean vectors of length N .

Let $\mathcal{S}_0, \mathcal{S}_1, \mathcal{S}_2 \dots \mathcal{S}_N$ be subsets of \mathcal{S} . The subset \mathcal{S}_i contains all the dichotomies where exactly i data vectors are classified as $+1$.

From the definition above, we can see that \mathcal{S}_i has $\binom{N}{i}$ elements, and any two of these subsets are mutually exclusive.

Now consider the set $\mathcal{S} \setminus (\mathcal{S}_k \cup \mathcal{S}_{k+1} \cup \dots \cup \mathcal{S}_N)$ for $k \in \{0, 1, 2, \dots, N\}$. For the dichotomies in this set, no k data vectors are scattered. This is true because for any k data vectors, this set does not contain any dichotomy where all these k data vectors are classified as $+1$, since any such dichotomy must be in one of $\mathcal{S}_k, \mathcal{S}_{k+1}, \dots, \mathcal{S}_N$.

The number of elements in $\mathcal{S} \setminus (\mathcal{S}_k \cup \mathcal{S}_{k+1} \cup \dots \cup \mathcal{S}_N)$ is

$$2^N - \left(\sum_{i=k}^N \binom{N}{i} \right) = \sum_{i=0}^{k-1} \binom{N}{i}$$

Hence we have

$$B(N, k) \geq \sum_{i=0}^{k-1} \binom{N}{i} \quad \forall k \in \{0, 1, 2, \dots, N\}$$

Lastly, for $k \in \{N+1, N+2, N+3 \dots\}$ obviously $B(N, k) = 2^N = \sum_{i=0}^{k-1} \binom{N}{i}$, as these N data vectors can be shattered.

Therefore, we have proven

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