

HTML HW4
B12901022 廖冠豪

5

By definition of the Hessian we have

$$A_E(\mathbf{w})_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} E_{in}(\mathbf{w})$$

$$\frac{\partial E_{in}(\mathbf{w})}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n) (-y_n x_{ni})}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}$$

$$\begin{aligned} \frac{\partial}{\partial w_j} \left(\frac{\partial E_{in}(\mathbf{w})}{\partial w_i} \right) &= \frac{1}{N} \sum_{n=1}^N \frac{(\exp(-y_n \mathbf{w}^T \mathbf{x}_n) (\mathbf{x}_{ni} \mathbf{x}_{nj})) (1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)) - \exp(-y_n \mathbf{w}^T \mathbf{x}_n)^2 (x_{ni} x_{nj})}{(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))^2} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))^2} x_{ni} x_{nj} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} \frac{1}{1 + \exp(y_n \mathbf{x}^T \mathbf{x}_n)} x_{ni} x_{nj} \\ &= \frac{1}{N} \sum_{n=1}^N h_t(\mathbf{x}_n) h_t(-\mathbf{x}_n) x_{ni} x_{nj} \end{aligned}$$

(x_{ni} is the i -th element of the \mathbf{x}_n , and is the same as X_{ji}) For a N by d matrix X and a N by N diagonal matrix D , we have

$$(X^T D X)_{ij} = \sum_{n=1}^N D_{nn} X_{ni} X_{nj}$$

By comparing the expressions, we can obtain

$$D_{ij} = \begin{cases} 0 & i \neq j \\ \frac{1}{N} h_t(\mathbf{x}_n) h_t(-\mathbf{x}_n) & i = j \end{cases}$$

6

In this problem, we use w_{j_i} to denote the i -th element of \mathbf{w}_j , which is the same as W_{ij} . Since E_{in} is minimized with SGD, we know that

$$V_{ij} = -\frac{\partial}{\partial w_{j_i}} \text{err}(W, \mathbf{x}, y)$$

$$\begin{aligned} \frac{\partial}{\partial w_{j_i}} \text{err}(W, \mathbf{x}, y) &= -\frac{1}{h_y(\mathbf{x})} \frac{\partial h_y(\mathbf{x})}{\partial w_{j_i}} \\ &= -\frac{1}{h_y(\mathbf{x})} \frac{\mathbb{I}[y = i] \exp(\mathbf{w}_y^T \mathbf{x}) (\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x})) x_j - \exp(\mathbf{w}_y^T \mathbf{x}) \exp(\mathbf{w}_i^T \mathbf{x}) x_j}{\left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}) \right)^2} \\ &= -\frac{1}{h_y(\mathbf{x})} \frac{(\mathbb{I}[y = i] - \exp(\mathbf{w}_i^T \mathbf{x})) \exp(\mathbf{w}_y^T \mathbf{x}) x_j}{\left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}) \right)^2} \\ &= -\left(\mathbb{I}[y = i] - \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x})} \right) x_j \\ &= (h_i(\mathbf{x}) - \mathbb{I}[y = i]) x_j \end{aligned}$$

Notice that

$$(\mathbf{x}_n \cdot \mathbf{u}^T)_{ij} = x_j u_i$$

Hence by comparing the expressions, \mathbf{u} is given by

$$u_i = -h_i(\mathbf{x}_n) + \mathbb{I}[y_n = i]$$

7

For MLR , we know that $\sum_{n=1}^N \text{err}(\mathbf{W}, \mathbf{x}, y)$ is minimized at its optimal solution, $(\mathbf{w}_1^*, \mathbf{w}_2^*)$. Also

$$\begin{aligned} \sum_{n=1}^N \text{err}(\mathbf{W}, \mathbf{w}, y) &= \sum_{n=1}^N \ln \left(\frac{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)}{\exp(\mathbf{w}_1^T \mathbf{x}_n) + \exp(\mathbf{w}_2^T \mathbf{x}_n)} \right) \\ &= \begin{cases} -\sum_{n=1}^N \ln(1 + \exp((\mathbf{w}_2^T - \mathbf{w}_1^T) \mathbf{x}_n)) & \text{if } y_n = 1, \\ -\sum_{n=1}^N \ln(1 + \exp((\mathbf{w}_1^T - \mathbf{w}_2^T) \mathbf{x}_n)) & \text{if } y_n = 2 \end{cases} \\ &= -\sum_{n=1}^N \ln(1 + \exp(-y'_n (\mathbf{w}_2 - \mathbf{w}_1)^T \mathbf{x}_n)) \end{aligned}$$

For logistic regression

$$E_{\text{in}}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y'_n \mathbf{w}^T \mathbf{x}_n))$$

By comparing the expressions, we see that $E_{\text{in}}(\mathbf{w})$ is minimized with

$$\mathbf{w}_{lr} = \mathbf{w}_2^* - \mathbf{w}_1^*$$

Hence this is the optimal solution to the logistic regression.

8

The linear hypothesis that minimizes E_{in} is the line that passes $(x_1, f(x_1))$ and $(x_2, f(x_2))$.
Hence

$$\begin{aligned} g(x) &= f(x_1) + \frac{f(x_2) - f(x_1)}{x_2 - x_1}(x - x_1) \\ &= -2(x_1 + x_2)x + 2x_1x_2 + 1 \end{aligned}$$

For such a hypothesis, $E_{\text{in}} = 0$.

E_{out} is given by

$$\begin{aligned} E_{\text{out}} &= \int_0^1 (g(x) - f(x))^2 dx \\ &= \int_0^1 (2x_1x_2 - 2(x_1 + x_2)x + 2x^2)^2 dx \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}(|E_{\text{in}} - E_{\text{out}}|) &= \int_0^1 \int_0^1 |E_{\text{out}}| dx_1 dx_2 \quad (E_{\text{in}} = 0) \\ &= \int_0^1 \int_0^1 \int_0^1 (2x_1x_2 - 2(x_1 + x_2)x + 2x^2)^2 dx dx_1 dx_2 \\ &= \frac{2}{15} \end{aligned}$$

9

Let

$$\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}$$

Where

$$\mathbf{E} = \begin{bmatrix} | & \dots & | \\ \epsilon & \dots & \epsilon \\ | & \dots & | \end{bmatrix}$$

$$\begin{aligned} \mathbf{X}_h^T \mathbf{X}_h &= \begin{bmatrix} \mathbf{X}^T & \tilde{\mathbf{X}}^T \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{X}} \end{bmatrix} \\ &= \mathbf{X}^T \mathbf{X} + \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \\ &= 2\mathbf{X}^T \mathbf{X} + \mathbf{E}^T \mathbf{X} + \mathbf{X}^T \mathbf{E} + \mathbf{E}^T \mathbf{E} \end{aligned}$$

Therefore

$$\begin{aligned} \mathbb{E}[\mathbf{X}_h^T \mathbf{X}_h] &= \mathbb{E}[2\mathbf{X}^T \mathbf{X} + \mathbf{E}^T \mathbf{X} + \mathbf{X}^T \mathbf{E} + \mathbf{E}^T \mathbf{E}] \\ &= 2\mathbf{X}^T \mathbf{X} + \mathbb{E}[\mathbf{E}^T] \mathbf{X} + \mathbf{X}^T \mathbb{E}[\mathbf{E}] + \mathbb{E}[\mathbf{E}^T \mathbf{E}] \end{aligned}$$

Since ϵ is generated i.i.d. from a normal distribution with variance σ^2 , we have

$$\begin{aligned} \mathbb{E}[\mathbf{E}^T] &= \mathbf{O}_{(d+1) \times N} \\ \mathbb{E}[\mathbf{E}] &= \mathbf{O}_{N \times (d+1)} \\ \mathbb{E}[\mathbf{E}^T \mathbf{E}] &= N\sigma^2 \mathbf{I}_{d+1} \end{aligned}$$

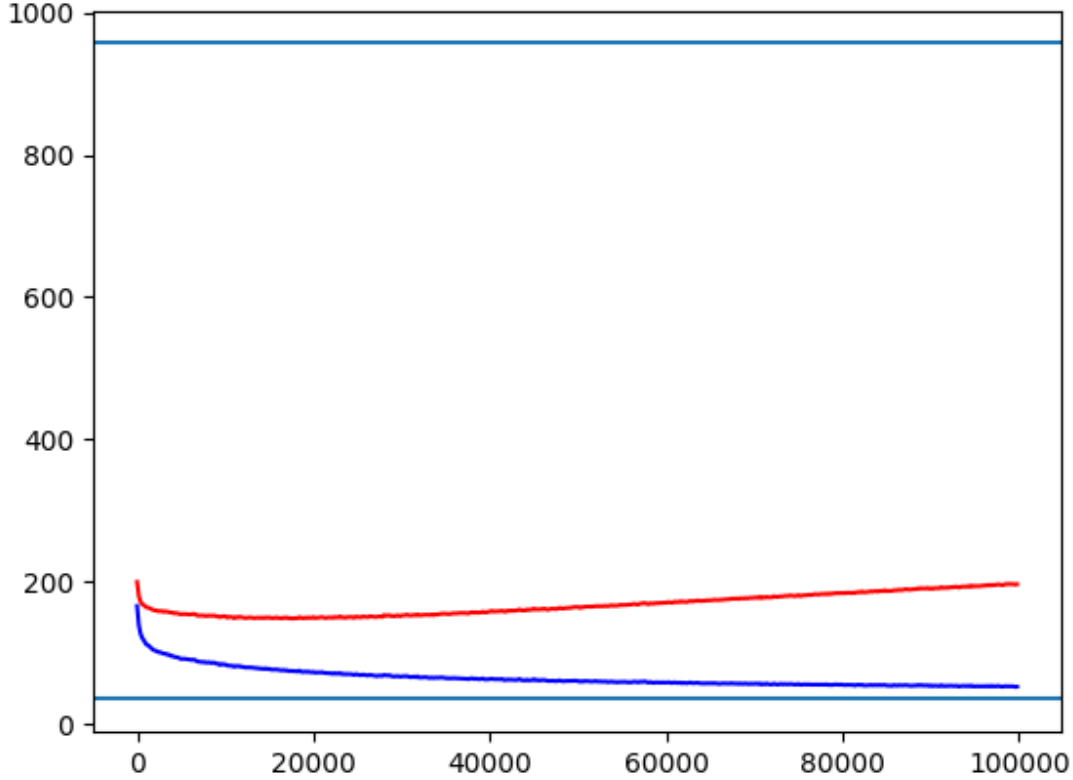
Hence

$$\mathbb{E}[\mathbf{X}_h^T \mathbf{X}_h] = 2\mathbf{X}^T \mathbf{X} + N\sigma^2 \mathbf{I}_{d+1}$$

So $(\alpha, \beta) = (2, N)$

10

Figure:



The horizontal light-blue lines at top and bottom are the average $E_{\text{out}}(\mathbf{w}_{lin})$ and $E_{\text{in}}(\mathbf{w}_{lin})$ respectively.

The red curve is the average $E_{\text{out}}(\mathbf{w}_t)$, and the dark-blue curve is the average $E_{\text{in}}(\mathbf{w}_t)$.

The horizontal axis represents t in the SGD process.

Findings:

First we can see that compared with the average values of linear regression, the E_{in} of SGD is slightly larger, and the E_{out} of SGD is significantly lower. This implies that SGD is probably a better approach compared with directly computing the weight vector in this case.

We see that starting from about 200 at $t = 0$, $E_{\text{in}}(\mathbf{w}_t)$ decreases monotonically as t increases. This is as expected, since the model can fit the $N = 64$ training data vectors better after more iterations. We also see that as t increases, $E_{\text{out}}(\mathbf{w}_t)$ first decreases slightly, after reaching its minimum at small t (around 0), it begins to increase slowly. At the end of iterations $t = 100000$, $E_{\text{out}}(\mathbf{w}_t)$ is approximately at the same level as when the iteration started. A possible cause for this phenomenon is that at the very beginning of the SGD process, E_{out} decreases because some information about the

data set is learned by the algorithm. But as t increases, the weight vector gets more biased by the $N = 64$ training data vectors, and its ability to classify the entire data set is compromised, hence the increasing E_{out} . This also means that a large number of iterations may not be meaningful or beneficial in this learning problem, since most iterations result in the increase of E_{out} . Moreover, there is no improve in E_{out} at the end of the 100000 iterations compared with at the beginning, when little "learning" is done.

Code snapshot:

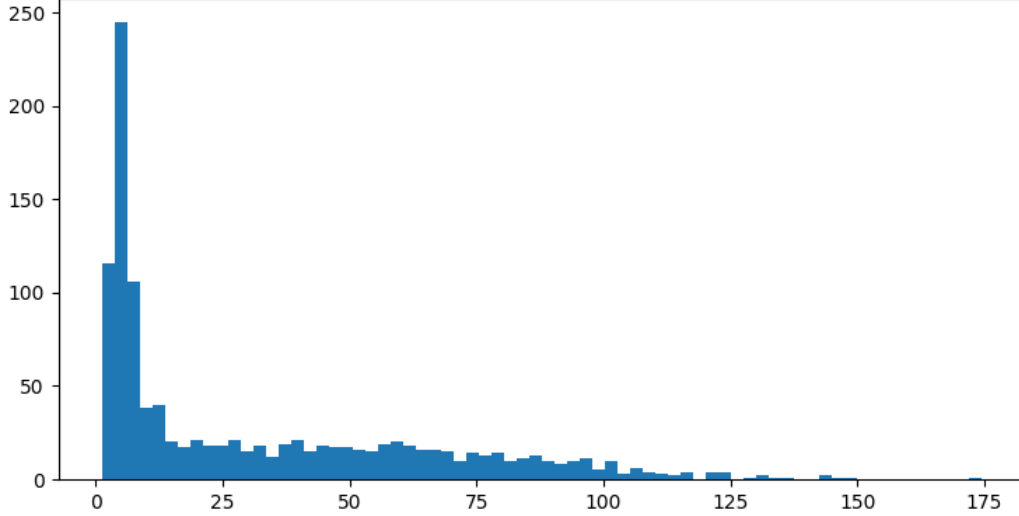
```

import numpy as np
import matplotlib.pyplot as plt

2
3
4 ### parse data
5 def parse_data_line(line):
6     parts = line.strip().split()
7     y = int(parts[0])
8     vector = []
9     vector.append(1)
10    for part in parts[1:]:
11        _, value = part.split(":")
12        vector.append(float(value))
13    return y, np.array(vector)
14
15
16 def parse_data_file():
17     file_path = "data.txt"
18     matrix = []
19     y_values = []
20     with open(file_path, "r") as file:
21         for line in file:
22             y, vector = parse_data_line(line)
23             y_values.append(y)
24             matrix.append(vector)
25     return np.array(y_values), np.array(matrix)
26
27
28 y, matrix = parse_data_file()
29 ### end parse
30
31 e_lin_in = []
32 e_lin_out = []
33 e_t_in = np.zeros(500)
34 e_t_out = np.zeros(500)
35
36 def sgd(random_indices):
37     w = np.zeros(13)
38     for j in range(1, 100001):
39         ind = random_indices[np.random.randint(0, 64)]
40         vec = matrix[ind, :].transpose()
41         w = np.subtract(w, 0.02 * (np.dot(vec, w) - y[ind]) * vec)
42         if j % 200 == 0:
43             err = np.subtract(y, np.matmul(matrix, w))
44             e_in = 0
45             for k in random_indices:
46                 e_in += err[k] ** 2
47             e_out = np.linalg.norm(err) ** 2
48             # print(e_out)
49             e_out += e_in
50             # print(e_out)
51             # for j in range(err.shape[0]):
52                 if j in random_indices:

```

Figure:



The average value of $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$ is **32.45148423777838**

This figure is the histogram for $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$, the horizontal axis is the value of $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$, and the vertical axis is the number of times a value happens in the 1126 experiments (call this "frequency" in the following context).

We see that $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$ mostly fall between 0 and 10, as the frequency is high in this range. The frequency drops rapidly at about $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly}) = 10$, after the rapid drop, $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$ decreases slowly, the largest value is slightly less than 150.

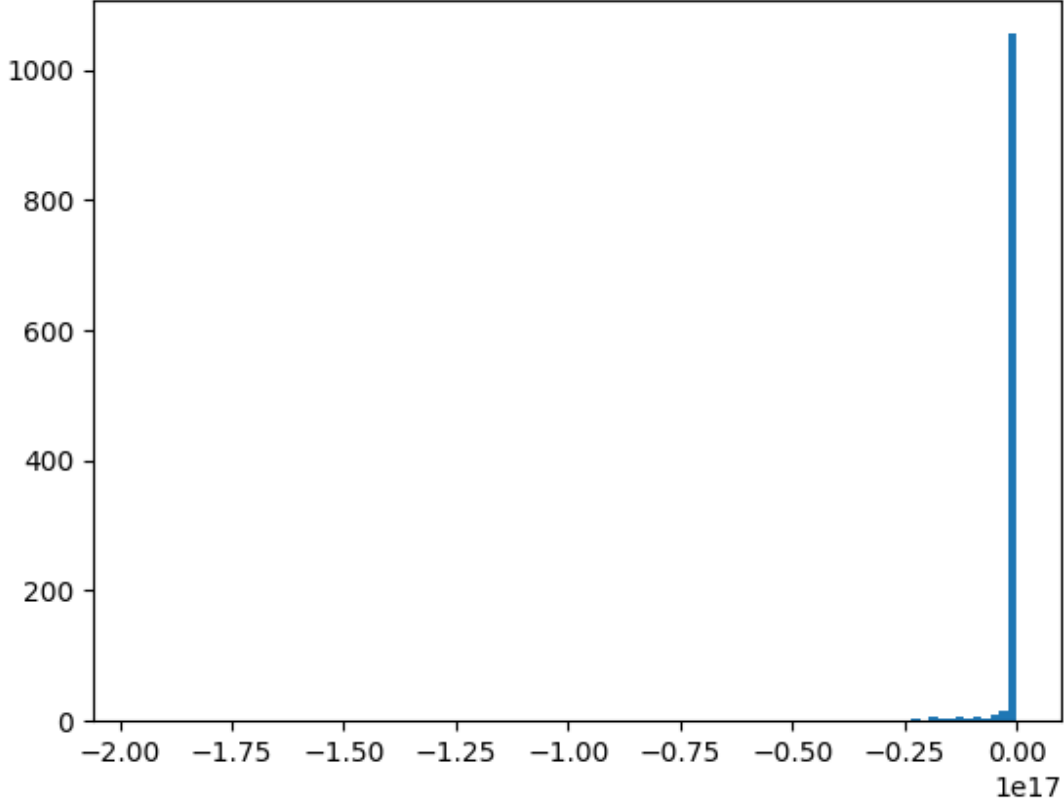
The histogram and the average value of $E_{in}^{sqr}(\mathbf{w}_{lin}) - E_{in}^{sqr}(\mathbf{w}_{poly})$ implies that the homogeneous polynomial transform does indeed result in a E_{in} gain. (which is in fact a decrease in E_{in} , but can be seen as a gain in the ability to fit the training set judging by E_{in}). Furthermore, compared with the average value of $E_{in}(\mathbf{w}_{lin})$ in P10, which is about 35 (see the light-blue horizontal line at the bottom), the E_{in} gain (≈ 32.5) is quite significant. This further shows that the polynomial transform reduces E_{in} effectively. This is because by allowing a more complicated model, the classifier can fit the $N = 64$ training data vectors better.

Code snapshot:

```
1044: main: H066.tex (-25%) N1 vim: main P10.py (-25%) 302
50 def poly_regression(random_indices):
51     sub_x = poly_matrix(random_indices, :)
52     sub_y = y[random_indices]
53     x_dagger = np.linalg.pinv(sub_x)
54     w_lin = np.matmul(x_dagger, sub_y)
55     err = np.subtract(y, np.matmul(poly_matrix, w_lin))
56     e_in = 0
57     for j in random_indices:
58         e_in += err[j] ** 2
59     return e_in / 64
60
61 def linear_regression(random_indices):
62     sub_x = matrix[random_indices, :]
63     sub_y = y[random_indices]
64     x_dagger = np.linalg.pinv(sub_x)
65     w_lin = np.matmul(x_dagger, sub_y)
66     err = np.subtract(y, np.matmul(matrix, w_lin))
67     e_in = 0
68     for j in random_indices:
69         e_in += err[j] ** 2
70     return e_in / 64
71
72 for i in range(1126):
73     print(i)
74     np.random.seed(i)
75     random_indices = np.random.choice(matrix.shape[0], size = 64, replace = False)
76     e_in_lin = linear_regression(random_indices)
77     e_in_poly = poly_regression(random_indices)
78     e_lin_in.append(e_in_lin)
79     e_poly_in.append(e_in_poly)
80 #
81 # random_indices = np.random.choice(matrix.shape[0], size = 64, replace = False)
82 # e_in, e_out = linear_regression(random_indices)
83 # e_lin_in.append(e_in)
84 # e_lin_out.append(e_out)
85 # sgd(random_indices)
86 #
87 # print(e_t_in)
88 # print(e_t_out)
89 #
90 # print(e_poly_in)
91 # print(e_lin_in)
92 #
93 de = []
94 for i in range(0, len(e_poly_in)):
95     de.append(e_lin_in[i] - e_poly_in[i])
96     plt.hist(de, bins = 30)
97     print(np.mean(de))
98     plt.show()
99 #
100 NORMAL main P11.py utf-8 < python Bot 94:1
```

12

Figure:



The average value of $E_{out}^{sqr}(\mathbf{w}_{lin}) - E_{out}^{sqr}(\mathbf{w}_{poly})$ is -1051028513192239.2

This figure is the histogram of $E_{out}^{sqr}(\mathbf{w}_{lin}) - E_{out}^{sqr}(\mathbf{w}_{poly})$, the horizontal axis is the value of $E_{out}^{sqr}(\mathbf{w}_{lin}) - E_{out}^{sqr}(\mathbf{w}_{poly})$, and the vertical axis is the number of times a value happens in the 1126 experiments. We see that $E_{out}^{sqr}(\mathbf{w}_{poly}) - E_{out}^{sqr}(\mathbf{w}_{lin})$ is huge, which implies that $E_{out}^{sqr}(\mathbf{w}_{poly})$ is far greater than $E_{out}^{sqr}(\mathbf{w}_{lin})$.

Together with the result in P11, we can see that adopting the polynomial transform in this problem makes classifying performance a lot worse. This is because the polynomial transform allows a more complicated model, and although E_{in} is decreased, the classifying performance is worse due to the hazard of overfitting.

To conclude, the linear regression with the polynomial transform in this problem results in a bad learning model. (compared with direct linear regression)

Code snapshot:

```

52
51
50 ##### end parse
49
48 e_lin_in = []
47 e_poly_in = []
46
45 def poly_regression(random_indices):
44     sub_x = poly_matrix(random_indices, :)
43     sub_y = y[random_indices]
42     x_dagger = np.linalg.pinv(sub_x)
41     w_lin = np.matmul(x_dagger, sub_y)
40     err = np.subtract(y, np.matmul(poly_matrix, w_lin))
39     e_in = 0
38     e_out = np.linalg.norm(err) ** 2
37     print(e_out / 8128)
36     for j in random_indices:
35         e_in += err[j] ** 2
34     e_out -= e_in
33     return e_out / 8128
32
31
30 def linear_regression(random_indices):
29     sub_x = matrix[random_indices, :]
28     sub_y = y[random_indices]
27     x_dagger = np.linalg.pinv(sub_x)
26     w_lin = np.matmul(x_dagger, sub_y)
25     err = np.subtract(y, np.matmul(matrix, w_lin))
24     e_in = 0
23     e_out = np.linalg.norm(err) ** 2
22     print(e_out / 8128)
21     for j in random_indices:
20         e_in += err[j] ** 2
19     e_out -= e_in
18     return e_out / 8128
17
16
15 for i in range(1126):
14     print(i)
13     np.random.seed(i ** 2)
12     random_indices = np.random.choice(matrix.shape[0], size = 64, replace = False)
11     e_out_lin = linear_regression(random_indices)
10     e_out_poly = poly_regression(random_indices)
9     e_lin_in.append(e_out_lin)
8     e_poly_in.append(e_out_poly)
7
6
5 de = []
4 for i in range(0, len(e_poly_in)):
3     de.append(e_lin_in[i] - e_poly_in[i])
2 plt.hist(de, bins = 100)
1 print(np.mean(de))
0 plt.show()

```

NORMAL
main
P12.py
utf-8
python
Bot
87:1

13

Notice that

$$\begin{aligned}\operatorname{sign}\left(w_0 + \sum_{i=1}^d w_i x_i\right) &= \operatorname{sign}\left(\exp\left(w_0 + \sum_{i=1}^d w_i x_i\right) - 1\right) \\ &= \operatorname{sign}\left(-e^{-w_0} + \prod_{i=1}^d (\exp(x_i))^{w_i}\right)\end{aligned}$$

Therefore we see that $h_{\mathbf{w}}(\mathbf{x}) = \tilde{h}_{\mathbf{u}}(\mathbf{x}')$ if we let

$$\begin{aligned}u_0 &= -e^{-w_0} \\ u_i &= w_i \quad i = 1, 2, \dots, d \\ x'_i &= \exp(x_i) - 1\end{aligned}$$

Hence if there exist n data vectors shattered by \mathcal{H}_1 , we can find a corresponding set of n data vectors that is shattered by \mathcal{H}_2 , which implies that $d_{vc}(\mathcal{H}_1) \leq d_{vc}(\mathcal{H}_2)$.

We can also see that $h_{\mathbf{w}}(\mathbf{x}') = \tilde{h}_{\mathbf{u}}(\mathbf{x})$ if we let

$$\begin{aligned}x'_i &= \ln(1 + |x_i|) \\ w_0 &= \begin{cases} -\ln(-u_0) & u_0 < 0 \\ 1 & u_0 \geq 0 \end{cases} \\ w_i &= \begin{cases} u_i & u_0 < 0 \\ 0 & u_0 \geq 0 \end{cases} \quad i = 1, 2, \dots, d\end{aligned}$$

(Observe that $\tilde{h}_{\mathbf{u}}(\mathbf{x}) = +1$ if $u_0 \geq 0$)

Hence if there exist n data vectors shattered by \mathcal{H}_2 , we can find a corresponding set of n data vectors that is shattered by \mathcal{H}_1 , which implies that $d_{vc}(\mathcal{H}_2) \leq d_{vc}(\mathcal{H}_1)$.

Since $d_{vc}(\mathcal{H}_1) \leq d_{vc}(\mathcal{H}_2)$ and $d_{vc}(\mathcal{H}_2) \leq d_{vc}(\mathcal{H}_1)$, $d_{vc}(\mathcal{H}_2) = d_{vc}(\mathcal{H}_1)$, and the statement is proved.