

5. Let $\tilde{w} = [w_0 \dots w_d]^T$

$$\begin{aligned}
 & \sum_{k=1}^K (\tilde{w}^T \tilde{x}_k - \tilde{y}_k)^2 \\
 &= \sum_{k=1}^{d+1} ([w_0 \ w_1 \ \dots \ w_d] \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \sqrt{\lambda \alpha_{k-1}} \\ 0 \\ \vdots \\ 0 \end{bmatrix} - 0)^2 \\
 &= \sum_{k=1}^{d+1} (\sqrt{\lambda \alpha_{k-1}} w_{k-1})^2 \\
 &= \sum_{k=1}^{d+1} \lambda \alpha_{k-1} w_{k-1}^2 = \lambda \alpha_0 w_0^2 + \lambda \alpha_1 w_1^2 + \dots + \lambda \alpha_d w_d^2 = \lambda \sum_{i=0}^d \alpha_i w_i^2
 \end{aligned}$$

Pr of solve:

$$\begin{aligned}
 & \min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N+K} \left(\sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2 + \sum_{k=1}^K (\tilde{w}^T \tilde{x}_k - \tilde{y}_k)^2 \right) \\
 &= \min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N+K} \left(\sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2 + \lambda \sum_{i=0}^d \alpha_i w_i^2 \right) \quad \text{K=d+1} \\
 &= \min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N+d+1} \left(\sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2 + \lambda \sum_{i=0}^d \alpha_i w_i^2 \right)
 \end{aligned}$$

Pr is to solve:

$$\min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2 + \frac{\lambda}{N} \sum_{i=0}^d \alpha_i w_i^2$$

$$= \min_{\tilde{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \left(\sum_{n=1}^N (\tilde{w}^T \tilde{x}_n - y_n)^2 + \lambda \sum_{i=0}^d \alpha_i w_i^2 \right)$$

equivalence

" $\frac{1}{N}, \frac{1}{N+d+1} > 0$

□

6. L2 regularization, twice differentiable, convex $E_{in}(\vec{w})$

$$\min_{\vec{w} \in \mathbb{R}^{d+1}} E_{aug}(\vec{w})$$

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}) + \frac{\lambda}{N} \|\vec{w}\|_2^2$$

$$\nabla E_{in}(\vec{w}^*) = 0$$

$$\tilde{E}_{in}(\vec{w}) = E_{in}(\vec{w}^*) + \underbrace{(\vec{w} - \vec{w}^*)^T \nabla E_{in}(\vec{w}^*)}_0 + \frac{1}{2} (\vec{w} - \vec{w}^*)^T H (\vec{w} - \vec{w}^*)$$

$H \in \mathbb{R}^{(d+1) \times (d+1)}$: Hessian, positive semi-definite ("convex")

approx $E_{aug}(\vec{w})$: $\tilde{E}_{aug}(\vec{w}) = \tilde{E}_{in}(\vec{w}) + \frac{\lambda}{N} \|\vec{w}\|_2^2$

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}^*) + \cancel{(\vec{w} - \vec{w}^*)^T \nabla E_{in}(\vec{w}^*)} + \frac{1}{2} (\vec{w} - \vec{w}^*)^T H (\vec{w} - \vec{w}^*) + \frac{\lambda}{N} \vec{w}^T \vec{w}$$

$$\frac{\partial}{\partial \vec{w}} E_{aug}(\vec{w}) = \frac{\partial}{\partial \vec{w}} E_{in}(\vec{w}^*) + \frac{1}{2} \frac{\partial}{\partial \vec{w}} (\vec{w} - \vec{w}^*)^T H (\vec{w} - \vec{w}^*) + \frac{\partial}{\partial \vec{w}} \left(\frac{\lambda}{N} \vec{w}^T \vec{w} \right)$$

$$= 0 + \frac{1}{2} (H^T + H) (\vec{w} - \vec{w}^*) + \frac{\lambda}{N} \cdot 2\vec{w}$$

$$= H \cdot (\vec{w} - \vec{w}^*) + \frac{2\lambda}{N} \vec{w}$$

$$\frac{d}{d\vec{w}} \vec{w}^T A \vec{w} = (A^T + A) \vec{w}$$

$$\vec{w}^T \vec{w} = \vec{w}^T I \vec{w} \quad \therefore \frac{\partial}{\partial \vec{w}} = (I^T + I) \vec{w} = 2\vec{w}$$

" $E_{in}(\vec{w})$ twice differentiable

\therefore the second partial derivatives \rightarrow continuous

\rightarrow Hessian: symmetric

Set $\frac{\partial}{\partial \vec{w}} E_{aug}(\vec{w}) = 0$

$$\Rightarrow \frac{1}{2} (H^T + H) (\vec{w} - \vec{w}^*) = \frac{1}{2} \cdot 2H \cdot (\vec{w} - \vec{w}^*) = H (\vec{w} - \vec{w}^*)$$

$$H \cdot (\vec{w} - \vec{w}^*) + \frac{2\lambda}{N} \vec{w} = 0$$

$$\Rightarrow H \vec{w} - H \vec{w}^* + \frac{2\lambda}{N} \vec{w} = 0$$

$$\Rightarrow (H + \frac{2\lambda}{N} I) \vec{w} = H \vec{w}^*$$

$$\Rightarrow \vec{w} = (H + \frac{2\lambda}{N} I)^{-1} H \vec{w}^* \quad \square$$

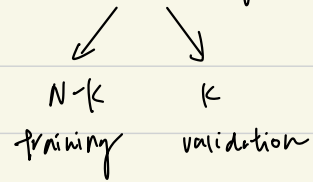
$$\{h_1, h_2, \dots, h_{d+1}\} \\ \parallel \text{vec}$$

since H : positive semi-definite, $\forall h_i \in \lambda(1)$, $h_i \geq 0$
also, $\frac{2\lambda}{N} > 0$

$$\text{Therefore, } \lambda(H + \frac{2\lambda}{N} I) = \{h_1 + \frac{2\lambda}{N}, \dots, h_{d+1} + \frac{2\lambda}{N}\}$$

$\Rightarrow (H + \frac{2\lambda}{N} I)$: invertible

1. N labels: $y_1, \dots, y_N \sim \text{i.i.d. } \mathcal{M}(0, \sigma^2)$



Expected validation error: $E\left(\frac{1}{k} \sum_{n=N-k+1}^N (y_n - 0)^2\right) = \sigma^2 \Rightarrow \frac{1}{k} E\left[\sum_{n=N-k+1}^N y_n^2\right] = \sigma^2 \Rightarrow E\left[\sum_{n=N-k+1}^N y_n^2\right] = k\sigma^2$

\uparrow
 $E[(X - E[X])^2] = \sigma^2$

$N-k$ samples estimate mean: $\bar{y} = \frac{1}{N-k} \sum_{n=1}^{N-k} y_n$

Q: expected validation error: $E\left(\frac{1}{k} \sum_{n=N-k+1}^N (y_n - \bar{y})^2\right)$

expansion \downarrow

$$\begin{aligned}
 E\left(\frac{1}{k} \sum_{n=N-k+1}^N (y_n - \bar{y})^2\right) &= \frac{1}{k} E\left[\sum_{n=N-k+1}^N (y_n^2 - 2y_n\bar{y} + \bar{y}^2)\right] \quad \text{linearity of expectation} \\
 &= \frac{1}{k} \left(E\left[\sum_{n=N-k+1}^N y_n^2\right] - 2E\left[\sum_{n=N-k+1}^N y_n\bar{y}\right] + E\left[\sum_{n=N-k+1}^N \bar{y}^2\right] \right) \\
 &= \frac{1}{k} \left(k\sigma^2 - 2 \sum_{n=N-k+1}^N E[y_n\bar{y}] + E\left[\sum_{n=N-k+1}^N \bar{y}^2\right] \right) \quad \text{linearity of expectation} \\
 &= \frac{1}{k} \left(k\sigma^2 - 2 \sum_{n=N-k+1}^N \cancel{E[y_n]E[\bar{y}]} + E\left[\sum_{n=N-k+1}^N \bar{y}^2\right] \right) \quad \begin{array}{l} \text{if } X, Y: \text{independent} \\ \text{then } E[XY] = E[X]E[Y] \end{array} \\
 &\quad \begin{array}{l} \because E[y_n] = 0 \\ \text{(labels generated from distribution with mean } = 0) \end{array} \\
 &\quad \begin{array}{l} \because \text{values of } y_n \text{ range from } N-k+1 \text{ to } N \\ \text{and } \bar{y} = \frac{1}{N-k} \sum_{n=1}^{N-k} y_n \text{ depends on } y_n \\ \text{that range from } 1 \text{ to } N-k \end{array} \\
 &\quad \therefore y_n, \bar{y}: \text{independent} \\
 &\quad \text{Thus, } E[y_n\bar{y}] = E[y_n]E[\bar{y}] \\
 &= \frac{1}{k} \left(k\sigma^2 + E\left[\sum_{n=N-k+1}^N \bar{y}^2\right] \right)
 \end{aligned}$$

$$\begin{aligned}
 E[\bar{y}] &= E\left[\frac{1}{N-k} \sum_{n=1}^{N-k} y_n\right] = \frac{1}{N-k} \sum_{n=1}^{N-k} E[y_n] = 0 \quad \because E[y_n] = 0 \\
 E[\bar{y}^2] &= E\left[\sum_{n=N-k+1}^N \bar{y}^2\right] \quad \begin{array}{l} E[\bar{y}] = 0 \\ \text{def of variance,} \\ \text{expected value: constant} \end{array} \\
 &= \sum_{n=N-k+1}^N E[(\bar{y} - E[\bar{y}])^2] \quad \rightarrow \sum_{n=N-k+1}^N c = kc \\
 &= k \cdot \text{Var}(\bar{y}) \\
 &= k \cdot \text{Var}\left(\frac{1}{N-k} \sum_{n=1}^{N-k} y_n\right) \\
 &= k \cdot \left(\frac{1}{N-k}\right)^2 \text{Var}\left(\sum_{n=1}^{N-k} y_n\right) \quad \text{Binomial's identity} \\
 &= k \cdot \left(\frac{1}{N-k}\right)^2 \left[\sum_{n=1}^{N-k} \text{Var}(y_n) + \sum_{i \neq j} \text{Cov}(y_i, y_j) \right]
 \end{aligned}$$

$$\begin{aligned}
 &= K \cdot \left(\frac{1}{N-K} \right)^2 \sum_{n=1}^{N-K} \text{Var}(y_n) \quad \swarrow \text{ ' Cov}(y_i, y_j) = 0 \quad \forall i, j \in \{1, \dots, N-K\} \\
 &\quad \text{since } y_1, \dots, y_{N-K} \text{ drawn independently} \\
 &= K \cdot \left(\frac{1}{N-K} \right)^2 \sum_{n=1}^{N-K} \sigma^2 \quad \swarrow \text{ given by the problem} \\
 &= K \cdot \left(\frac{1}{N-K} \right)^2 (N-K) \sigma^2 \\
 &= \frac{K \sigma^2}{N-K}
 \end{aligned}$$

↓

$$\begin{aligned}
 &= \frac{1}{K} \left(K \sigma^2 + \frac{K \sigma^2}{N-K} \right) \\
 &= \sigma^2 + \frac{\sigma^2}{N-K} \\
 &= \sigma^2 \left(1 + \frac{1}{N-K} \right) \quad \square
 \end{aligned}$$

i.i.d.

\mathcal{D} : N labels, $y_1, \dots, y_N \sim M < \infty$

w_0 : estimate mean

$$E_{in}[w] = \frac{1}{N} \sum_{n=1}^N (w_0 - y_n)^2 \quad w_0^* = \frac{1}{N} \sum_{n=1}^N y_n$$

$$E_{in}[w_0^*] = \frac{1}{N} \sum_{n=1}^N (w_0^* - y_n)^2$$

Claim:

for $N \geq 2$

$$E_{out}(A_{avg}) = \left(\frac{N}{N-1} \right)^2 E_{in}(w_0^*)$$

$$= \frac{1}{N} \sum_{n=1}^N \left[(w_0^*)^2 - 2w_0^* y_n + y_n^2 \right] \quad \text{---} \quad *$$

Proof: for $N \geq 2$, each time we can partition the N examples into $\begin{bmatrix} N-1 & : & \text{training} \\ 1 & : & \text{validation} \end{bmatrix}$

Let D_n be the new training set, which is:

$$D_n = \{(\tilde{x}_1, y_1), (\tilde{x}_2, y_2), \dots, (\tilde{x}_n, y_n), \dots, (\tilde{x}_N, y_N)\}$$

Also, let g_n be the hypothesis learned from D_n

\rightarrow Let error on validation set $\{(\tilde{x}_n, y_n)\}$: e_n

$$\text{where } e_n = E_{val}[g_n] = e[g_n(\tilde{x}_n), y_n]$$

$$E_{out}(A_{avg}) = \frac{1}{N} \sum_{n=1}^N e_n$$

$$= \frac{1}{N} \sum_{n=1}^N \left[\left(\frac{1}{N-1} \sum_{i \neq n}^N y_i \right) - y_n \right]^2$$

$$= \frac{1}{N} \sum_{n=1}^N \left(\frac{Nw_0^* - y_n}{N-1} - y_n \right)^2$$

$$= \frac{1}{N} \sum_{n=1}^N \left[\frac{Nw_0^* - y_n - y_n(N-1)}{N-1} \right]^2$$

$$= \frac{1}{N(N-1)^2} \sum_{n=1}^N [Nw_0^* - y_n - y_n(N-1)]^2$$

$$= \frac{N^2}{N(N-1)^2} \sum_{n=1}^N \left[(w_0^*)^2 - 2w_0^* y_n + y_n^2 \right] \quad \text{---} \quad *$$

$$= \frac{N^2}{N(N-1)^2} \cdot N \cdot E_{in}[w_0^*]$$

$$= \frac{N^2}{(N-1)^2} E_{in}[w_0^*] \quad \square$$

$$w_0^* = \frac{1}{N} \sum_{i=1}^N y_i$$

$$\rightarrow Nw_0^* = \sum_{i=1}^N y_i$$

$$\rightarrow Nw_0^* - y_n = \sum_{i=1}^N y_i - y_n = \sum_{i \neq n}^N y_i$$

$$\rightarrow \frac{Nw_0^* - y_n}{N-1} = \frac{1}{N-1} \sum_{i \neq n}^N y_i$$

$$[Nw_0^* - y_n - y_n(N-1)]^2$$

$$= [(Nw_0^* - y_n) - y_n(N-1)]^2$$

$$= (Nw_0^* - y_n)^2 - 2(Nw_0^* - y_n)y_n(N-1) + y_n^2(N-1)^2$$

$$= N(w_0^*)^2 - 2Nw_0^*y_n + 1 - 2N(N-1)w_0^*y_n + 2(N-1) + (N-1)^2$$

$$= N(w_0^*)^2 - 2Nw_0^*y_n [1 + (N-1)] + N^2$$

$$= N(w_0^*)^2 - 2Nw_0^*y_n + N^2$$

$$= N^2 \left[(w_0^*)^2 - 2w_0^*y_n + 1 \right]$$

$$(2Nw_0^* - 2y_n)y_n(N-1)$$

$$= (2Nw_0^*y_n - 2y_n^2)(N-1)$$

$$= 2N(N-1)w_0^*y_n - 2(N-1)y_n^2$$

$$[w_0^* + 1]^2 = N^2$$

9. binary classifier g : $P(g(\vec{x}) = -1 \mid y = +1) = \epsilon_+$
 $P(g(\vec{x}) = +1 \mid y = -1) = \epsilon_-$

distribution: $P(y = +1) = P(y = -1) = \frac{1}{2}$

$\rightarrow E_{out}(g) = \frac{1}{2}\epsilon_+ + \frac{1}{2}\epsilon_-$

distribution: $P(y = -1) = p$

constant classifier $g_c \rightarrow E_{out}(g_c) = p$
 (always predict +1)

$E_{out}(g)$ over distribution $P(y = -1) = p$:
 $\rightarrow P(y = +1) = 1 - p$

$$E_{out}(g) = P(g(\vec{x}) = -1 \mid y = +1) P(y = +1) + P(g(\vec{x}) = +1 \mid y = -1) P(y = -1)$$

$$= (1-p)\epsilon_+ + p\epsilon_-$$

$E_{out}(g) = E_{out}(g_c)$:

$$(1-p)\epsilon_+ + p\epsilon_- = p$$

$$\rightarrow \epsilon_+ - p\epsilon_+ + p\epsilon_- = p$$

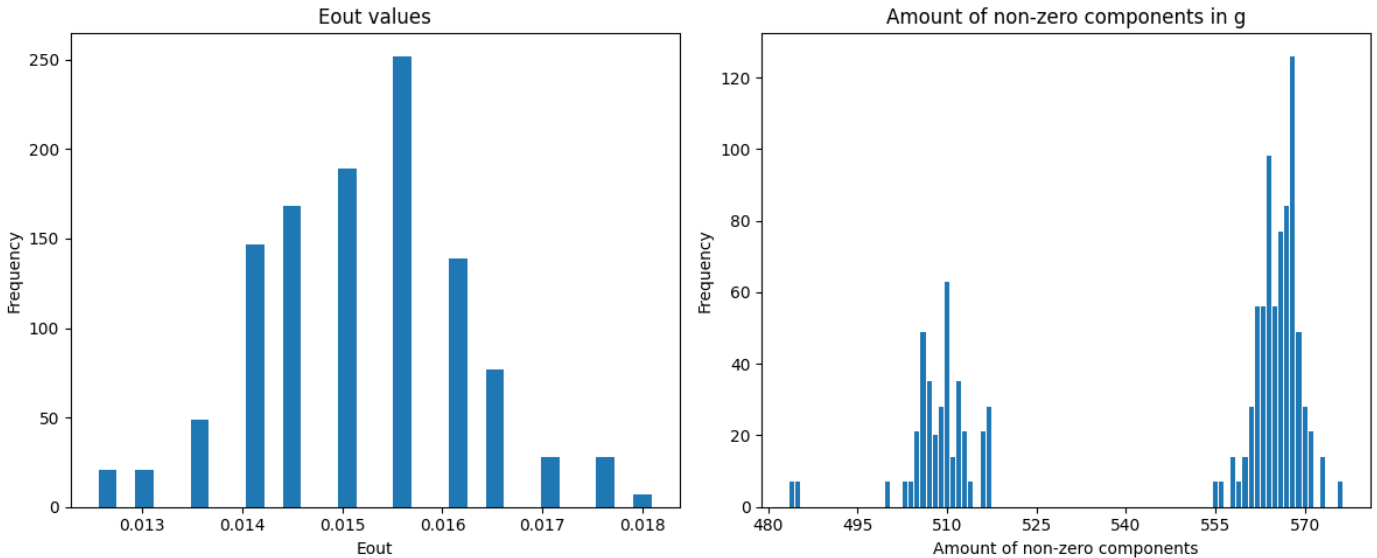
$$\rightarrow \epsilon_+ = p + p\epsilon_+ - p\epsilon_-$$

$$\rightarrow \epsilon_+ = p(1 + \epsilon_+ - \epsilon_-)$$

$$\rightarrow p = \frac{\epsilon_+}{1 + \epsilon_+ - \epsilon_-} \quad \square$$

HTML homework 5: q10 report

The resulting plot is as follows:



Some additional information that I printed in order to do the comparison between this question and question 11:

```
Statistics for Eout:
Mean: 0.0152
Median: 0.0151
Standard Deviation: 0.0010
Min: 0.0126
Max: 0.0181

Statistics for Non-zero Components:
Mean: 546.5
Median: 563.0
Standard Deviation: 27.2
Min: 484
Max: 576
```

Code¹:

Preprocessing the dataset:

```
1 from liblinear.liblinearutil import *
2 from itertools import combinations_with_replacement
3 import numpy as np
4 import scipy.sparse
5 import matplotlib.pyplot as plt
6 from tqdm import tqdm
7 from joblib import Parallel, delayed
8
[28] ✓ 0.4s
```

Dataset

- `mnist.scale` for training
- `mnist.scale.t` for testing

```
1 train_data = 'mnist.scale'
2 test_data = 'mnist.scale.t'
```

[2]

Empty markdown cell, double-click or press enter to edit.

Read and save data

The training data is saved in the `X_train` list, with each element being a dictionary of features. example format: {153: 0.0117647, 154: 0.0705882, 155: 0.0705882, ...}

The corresponding label list is saved in `y_train`.

Similarly, the testing data is saved in `X_test` and the corresponding label list is saved in `y_test`.

```
1 def read_linear_format(file_path):
2     X, y = [], []
3     with open(file_path, 'r') as f:
4         for line in f:
5             parts = line.strip().split()
6             if parts[0] == '2' or parts[0] == '6':
7                 y.append(int(parts[0]))
8                 features = {}
9                 for item in parts[1:]:
10                     index, value = item.split(":")
11                     features[int(index)] = float(value)
12                 X.append(features)
13     return X, y
14
15 X_train, y_train = read_linear_format(train_data)
16 X_test, y_test = read_linear_format(test_data)
17
```

Error function:

0/1 error

```
1 def ZeroOneError(predictions, y):
2     return sum(p_i != y_i for p_i, y_i in zip(predictions, y)) / len(predictions)
```

[6]

¹ Most of the screenshots here are from my .ipynb file, since I wrote the explanations here, however, I ran the code in a .py file (which has equivalent contents), and part of the screenshot is from this file.

Explanation of deriving the relationship between C, λ :

Regularized logistic regression

Original Formulation in README.md

As the explanation in the liblinear README file, for `L1-regularized_logistic_regression` (s 6), we solve

$$\min_{\vec{w}} \sum |w_j| + C \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

Components of this equation

In this equation, the `L1 regularizer` is defined as:

$$\Omega(\vec{w}) = \sum |w_j| = \|\vec{w}\|_1$$

info in Lec14 slide 18

- Note: using L1 regularization encourages sparsity, which means fewer parameters.

And the `in-sample error` for logistic regression is:

$$E_{in}(\vec{w}) = \frac{1}{N} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

in textbook p.91

Derivation of C

Since we knew that:

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}) + \frac{\lambda}{N} \Omega(\vec{w})$$

We can derive that:

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}) + \frac{\lambda}{N} \sum |w_j|$$

For minimization:

$$\begin{aligned} \arg \min_{\vec{w}} E_{aug}(\vec{w}) &= \arg \min_{\vec{w}} (E_{in}(\vec{w}) + \frac{\lambda}{N} \sum |w_j|) \\ &= \arg \min_{\vec{w}} \frac{N}{\lambda} \left(\frac{1}{N} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i)) \right) + \frac{\lambda}{N} \sum |w_j| \\ &= \arg \min_{\vec{w}} \frac{1}{\lambda} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i)) + \sum |w_j| \end{aligned}$$

Compare to the original equation marked with (*), we can see that:

$$C = \frac{1}{\lambda}$$

Some additional information:

Select the best λ^*

Selecting optimal λ^* using the following equation:

$$\lambda^* = \arg \min_{\log_{10}(\lambda) \in \{-2, -1, 0, 1, 2, 3\}} E_{in}(\vec{w}_{\lambda})$$

get amount of non-zeros

To get the amount of non-zero components in the model, we can use the function `.get_decfun()`:

```
[W, b] = model_.get_decfun()
```

Main code:

```
1 def run_single_experiment(experiment):
2     np.random.seed(experiment)
3     min_Ein = np.inf
4     opt_log10_lambda = 0
5
6     for log10_lambda in (-2, -1, 0, 1, 2, 3):
7         c = 1 / (10 ** log10_lambda)
8         prob = problem(y_train, X_train)
9         param = parameter('-s 6 -c ' + str(c))
10        model = train(prob, param)
11
12        train_label, __, _ = predict(y_train, X_train, model)
13        Ein = ZeroOneError(train_label, y_train)
14        if Ein == min_Ein:
15            opt_log10_lambda = max(opt_log10_lambda, log10_lambda)
16            if opt_log10_lambda == log10_lambda:
17                opt_model = model
18        elif Ein < min_Ein:
19            min_Ein = Ein
20            opt_log10_lambda = log10_lambda
21            opt_model = model
22
23        test_label, __, _ = predict(y_test, X_test, opt_model)
24        Eout = ZeroOneError(test_label, y_test)
25
26        W = np.array(opt_model.get_decfun()[0])
27        non_zero_count = np.count_nonzero(W)
28
29        return Eout, non_zero_count
30
31 experiment_amount = 10
32 results = Parallel(n_jobs=-1)(
33     delayed(run_single_experiment)(i) for i in tqdm(range(experiment_amount))
34 )
35
36 Eouts, non_zero_count_list = zip(*results)
37 Eouts = list(Eouts)
38 non_zero_count_list = list(non_zero_count_list)
```

Plotting:

```
# aim: plotting
# subaim: First subplot (Eout values)
ax1.hist(Eouts, bins=30)
ax1.set_title('Eout values')
ax1.set_xlabel('Eout')
ax1.set_ylabel('Frequency')

# subaim: Second subplot (non-zero components)
min_val = int(min(non_zero_count_list))
max_val = int(max(non_zero_count_list))
integer_bins = np.arange(min_val, max_val + 2) - 0.5 # +2 to include max_val, -0.5 for bin edges

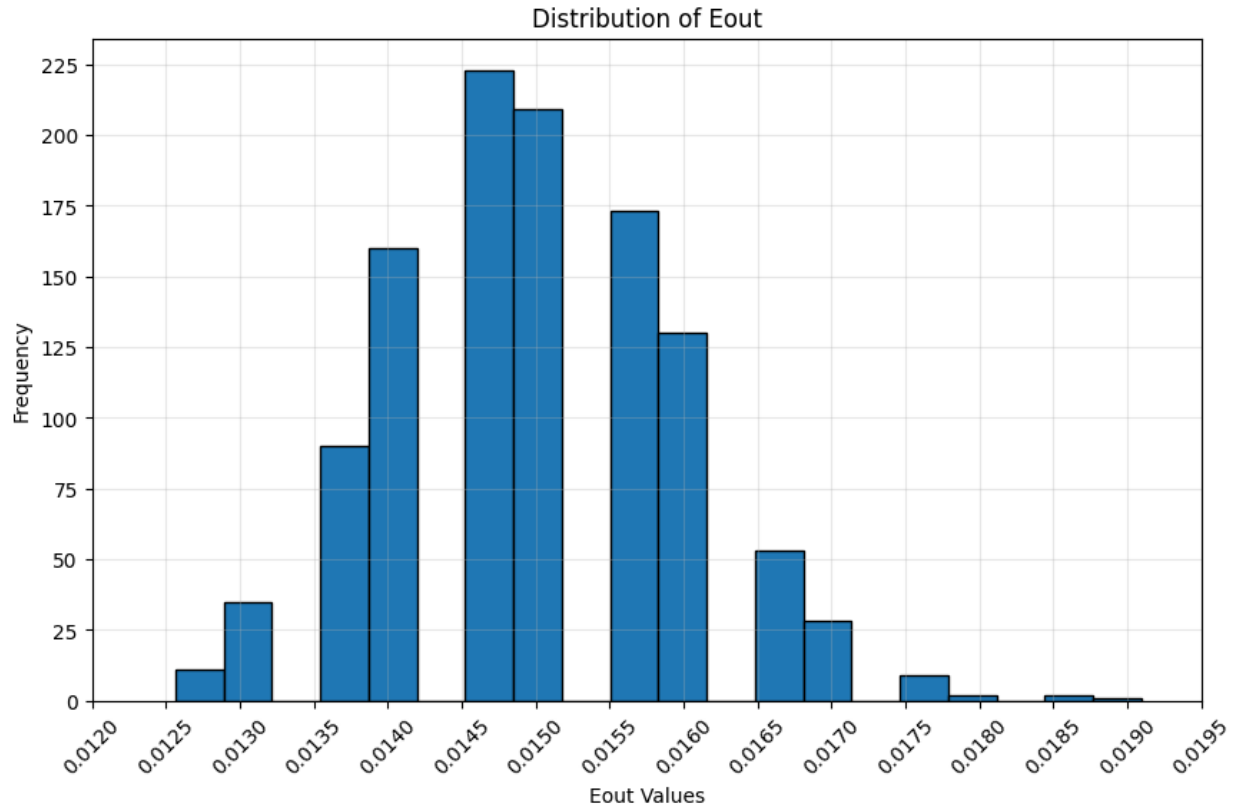
ax2.hist(non_zero_count_list,
         bins=integer_bins,
         align='mid',
         rwidth=0.8)

ax2.xaxis.set_major_locator(plt.MaxNLocator(integer=True))
ax2.set_title('Amount of non-zero components in g')
ax2.set_xlabel('Amount of non-zero components')
ax2.set_ylabel('Frequency')

plt.tight_layout()
plt.show()
```

HTML homework 5: q11 report

The resulting picture is as follows:



In order to compare with the previous question, some statistics are:

```
Statistics for Eout:
Mean: 0.0150
Median: 0.0151
Standard Deviation: 0.0010
Min: 0.0126
Max: 0.0191
```

```
Statistics for Eout:
Mean: 0.0152
Median: 0.0151
Standard Deviation: 0.0010
Min: 0.0126
Max: 0.0181

Statistics for Non-zero Components:
Mean: 546.5
Median: 563.0
Standard Deviation: 27.2
Min: 484
Max: 576
```

¹ Left: statistics of q11 result; Right: statistics of q10 result

In this problem, the difference from the previous is that we further divide the training set into a subtrain set (with 8000 examples), and a validation set (3876 examples).

From the statistics, we can see that most of them are quite similar, I think part of the reason is that, even though 3876 examples are taken out from the original training set, so one may expect that as we train on a smaller subtrain set, the out of sample error should be bigger than in q10 while we have the whole training set used, we got an expected value slightly smaller because we retrain the whole training set after the optimal λ is chosen. We did not use g^- to predict the test set, but g instead.

Apart from this reason, I think this result may also due to the precision, if more digits are taken into consideration, maybe we can observe the difference.

Another cause may be the limited number of possible λ , this constraints the variability of our model, even though there are infinitely many possible models, most of them are alike, while only slightly different by the λ value chosen.

Code:

```
1  Eouts = []
2
3  for experiment in tqdm(range(1126)):
4
5      np.random.seed(experiment)
6
7      total_size = len(X_train)
8      indices = np.random.permutation(total_size)
9
10     subtrain_indices = indices[:8000]
11     validation_indices = indices[8000:]
12
13     X_subtrain = [X_train[i] for i in subtrain_indices]
14     X_validation = [X_train[i] for i in validation_indices]
15     y_subtrain = [y_train[i] for i in subtrain_indices]
16     y_validation = [y_train[i] for i in validation_indices]
17
18     subtrain_prob = problem(y_subtrain, X_subtrain)
19
20     min_validation_err = np.inf
21     opt_log10_lambda = 0
22     for log10_lambda in (-2, -1, 0, 1, 2, 3):
23         subtrain_pred_res = []
24         c = 1 / (10 ** log10_lambda)
25         param = parameter('-s 6 -c ' + str(c))
26         model_by_subtrain = train(subtrain_prob, param)
27
28         validation_label, _, _ = predict(y_validation, X_validation, model_by_subtrain)
29         validation_err = ZeroOneError(validation_label, y_validation)
30         if validation_err == min_validation_err:
31             opt_log10_lambda = max(opt_log10_lambda, log10_lambda) # break tie by choosing the larger lambda
32         elif validation_err < min_validation_err:
33             min_validation_err = validation_err
34             opt_log10_lambda = log10_lambda
35
36     # rerun the model with the best lambda on the whole training set
37     whole_train_prob = problem(y_train, X_train)
38     param_with_opt_lambda = parameter('-s 6 -c ' + str(1 / (10 ** opt_log10_lambda)))
39     whole_train_model = train(whole_train_prob, param_with_opt_lambda)
40
41     test_label, _, _ = predict(y_test, X_test, whole_train_model)
42     Eout = ZeroOneError(test_label, y_test)
43
44     Eouts.append(Eout)
```

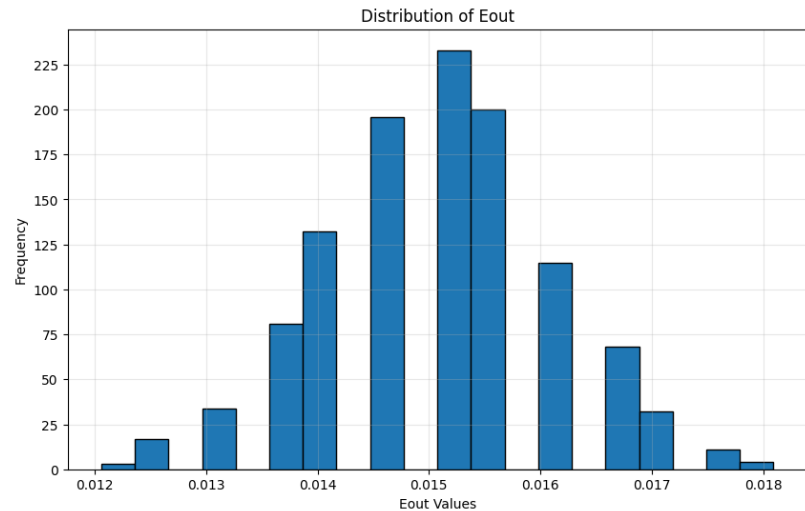
✓ 61m 42.5s

2

² Other part of codes that are similar to problem 10 (like reading in the data, split the data to get proper form of X and y, the ZeroOneError function...) are omitted in this report, can check for problem 10 for details.

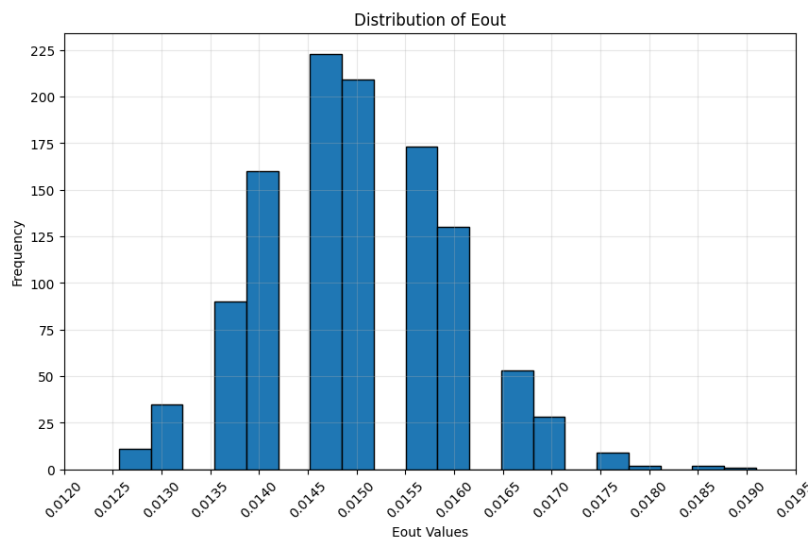
HTML homework 5: q12 report

The resulting picture is as follows:



(q12)

To find the difference, the plot below is the result of the previous problem:



(q11)

From the two plots we can see that the subtrain / validation split is more dispersed, and the variability of the 3-fold cross validation is smaller.

The reason why the subtrain / validation split is more dispersed may due to the fact that the validation set size is not that big compared to the 3-fold version. In the previous question, the validation set size is 3876, and in this question, using 3-fold splits the training set into 3 sets,

with each about 3959 examples, which is slightly higher. Also, E_{CV} is calculated by the mean error over using each of the 3 sets as the validation set, therefore these 2 reasons introduce stability in the 3-fold version.

On the other hand, the subtrain / validation split relies only on one split, if the validation set is not that representative, our selection of λ may be biased, this would magnify the effect of the unstableness in the previous question.

Code:

The markdown screenshot below describes the modification of this problem, other parts like reading in the data and the error function..., are the same as problem 10:

3 folds cross validation

textbook p.150

In this problem, we're asked to conduct 3-folds cross validation on the training data, which means we do the following steps:

1. partition the training data into 3 disjoint subsets (folds): D_1, D_2, D_3 , each of size $\approx \frac{N}{3}$
2. for each fold D_i , we train the model using the remaining 2 folds, and evaluate the model on D_i

Each set serves as a validation set to compute the validation error for the hypothesis g learned on a training set which is the complement of the validation set $D \setminus D_i$

→ we select the optimal λ^* by this validation error E_{CV}

3. rerun the model with the best λ^* on the whole training set
4. evaluate the model on the testing data

```

1  def run_single_experiment(experiment):
2      np.random.seed(experiment)
3
4      total_size = len(X_train)
5      indices = np.random.permutation(total_size)
6
7      fold_size = total_size // 3
8
9      fold1_indices = indices[:fold_size]
10     fold2_indices = indices[fold_size:2*fold_size]
11     fold3_indices = indices[2*fold_size:]
12
13     X_fold1 = [X_train[i] for i in fold1_indices]
14     X_fold2 = [X_train[i] for i in fold2_indices]
15     X_fold3 = [X_train[i] for i in fold3_indices]
16
17     y_fold1 = [y_train[i] for i in fold1_indices]
18     y_fold2 = [y_train[i] for i in fold2_indices]
19     y_fold3 = [y_train[i] for i in fold3_indices]
20
21     min_Ecv = np.inf
22     opt_log10_lambda = 0
23     for log10_lambda in (-2, -1, 0, 1, 2, 3):
24         each_fold_as_valid_err = []
25         c = 1 / (10 ** log10_lambda)
26
27         for i in range(1,4):
28             if i == 1:
29                 X_subtrain = X_fold2 + X_fold3
30                 y_subtrain = y_fold2 + y_fold3
31                 X_validation = X_fold1
32                 y_validation = y_fold1
33             elif i == 2:
34                 X_subtrain = X_fold1 + X_fold3
35                 y_subtrain = y_fold1 + y_fold3
36                 X_validation = X_fold2
37                 y_validation = y_fold2
38             else:
39                 X_subtrain = X_fold1 + X_fold2
40                 y_subtrain = y_fold1 + y_fold2
41                 X_validation = X_fold3
42                 y_validation = y_fold3
43
44             subtrain_prob = problem(y_subtrain, X_subtrain)
45             param = parameter('-s 6 -c ' + str(c))
46             model_by_subtrain = train(subtrain_prob, param)
47
48             validation_label, _, _ = predict(y_validation, X_validation, model_by_subtrain)
49             validation_err = ZeroOneError(validation_label, y_validation)
50             each_fold_as_valid_err.append(validation_err)
51
52     Ecv_each_lambda = np.mean(each_fold_as_valid_err)

```



```

53
54     if Ecv_each_lambda == min_Ecv:
55         opt_log10_lambda = max(opt_log10_lambda, log10_lambda) # break tie by choosing larger lambda
56     elif Ecv_each_lambda < min_Ecv:
57         min_Ecv = Ecv_each_lambda
58         opt_log10_lambda = log10_lambda
59
60     whole_train_prob = problem(y_train, X_train)
61     param_with_opt_lambda = parameter('-s 6 -c ' + str(1 / (10 ** opt_log10_lambda)))
62     whole_train_model = train(whole_train_prob, param_with_opt_lambda)
63
64     test_label, _, _ = predict(y_test, X_test, whole_train_model)
65     return ZeroOneError(test_label, y_test)
66
67 n_jobs = -1
68 Eouts = Parallel(n_jobs=n_jobs)(
69     delayed(run_single_experiment)(experiment)
70     for experiment in tqdm(range(1126))
71 )

```

✓ 166m 7.7s

13. elastic net solves:

$$\min_{\vec{w} \in \mathbb{R}^n} \frac{1}{N} \|\vec{y} - X\vec{w}\|_2^2 + \underbrace{\frac{\lambda_1}{N} \|\vec{w}\|_1}_{\text{not smooth}} + \frac{\lambda_2}{N} \|\vec{w}\|_2^2$$

\therefore use coordinate descent (instead of gradient descent):

$$w_i^{(tr+1)} = \arg\min_{w_i \in \mathbb{R}} \frac{1}{N} \sum_{n=1}^N (y_n - \sum_{j \neq i} w_j^{(t)} x_{n,j} - w_i x_{n,i})^2 + \underbrace{\frac{\lambda_1}{N} \left(\sum_{j \neq i} |w_j^{(t)}| + |w_i| \right)}_{\text{constant}} + \underbrace{\frac{\lambda_2}{N} \left(\sum_{j \neq i} (w_j^{(t)})^2 + w_i^2 \right)}_{\text{constant}}$$

\rightarrow closed form sol: $w_i^{(tr+1)} \leftarrow \underset{\substack{\text{sign flip} \\ \{1,1\}}}{\underset{\mathbb{R}}{\alpha \cdot \max(\beta, 0)}}$ if $\beta \leq 0$, $w_i^{(tr+1)} = 0 \rightarrow \text{sparsity}$

$$\begin{aligned} & (y_n - \sum_{j \neq i} w_j^{(t)} x_{n,j} - w_i x_{n,i})^2 \\ &= \left(y_n - \sum_{j \neq i} w_j^{(t)} x_{n,j} \right) - w_i x_{n,i} \Big)^2 \\ &= \left(y_n - \sum_{j \neq i} w_j^{(t)} x_{n,j} \right)^2 - 2 \left(y_n - \sum_{j \neq i} w_j^{(t)} x_{n,j} \right) w_i x_{n,i} + (w_i x_{n,i})^2 \\ &= y_n^2 - 2 y_n \sum_{j \neq i} w_j^{(t)} x_{n,j} + \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} \right)^2 - 2 y_n w_i x_{n,i} + 2 \sum_{j \neq i} w_j^{(t)} x_{n,j} w_i x_{n,i} + (w_i x_{n,i})^2 \end{aligned}$$

expand the update equation:

$$w_i^{(tr+1)} = \arg\min_{w_i \in \mathbb{R}} \left\{ \frac{1}{N} \sum_{n=1}^N \left[y_n^2 - 2 y_n \sum_{j \neq i} w_j^{(t)} x_{n,j} + \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} \right)^2 - 2 y_n w_i x_{n,i} + 2 \sum_{j \neq i} w_j^{(t)} x_{n,j} w_i x_{n,i} + (w_i x_{n,i})^2 \right] + \frac{\lambda_1}{N} \sum_{j \neq i} |w_j^{(t)}| + \frac{\lambda_1}{N} |w_i| + \frac{\lambda_2}{N} \sum_{j \neq i} (w_j^{(t)})^2 + \frac{\lambda_2}{N} w_i^2 \right\}$$

$$= \arg\min_{w_i \in \mathbb{R}} \left[\cancel{\frac{1}{N} \sum_{n=1}^N y_n^2} + \cancel{\frac{1}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} \right)^2} - \frac{2 w_i}{N} \sum_{n=1}^N x_{n,i} y_n + \frac{2 w_i}{N} \sum_{n=1}^N \sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} + \frac{w_i^2}{N} \sum_{n=1}^N (x_{n,i})^2 + \frac{\lambda_1}{N} \sum_{j \neq i} |w_j^{(t)}| + \frac{\lambda_1}{N} |w_i| + \frac{\lambda_2}{N} \sum_{j \neq i} (w_j^{(t)})^2 + \frac{\lambda_2}{N} w_i^2 \right]$$

$$= \arg\min_{w_i \in \mathbb{R}} \frac{2 w_i}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} y_n \right) + \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) w_i^2 + \frac{\lambda_1}{N} |w_i|$$

Case 1:

if $w_i > 0$, then we get:

$$\arg\min_{w_i \in \mathbb{R}} \frac{2 w_i}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} y_n \right) + \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) w_i^2 + \frac{\lambda_1}{N} w_i$$

$$\Rightarrow \frac{\partial}{\partial w_i} \left(\frac{2 w_i}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} y_n \right) + \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) w_i^2 + \frac{\lambda_1}{N} w_i \right) = 0$$

$$\rightarrow \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] + w_i \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) + \frac{\lambda_1}{N} = 0$$

$$\rightarrow w_i \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) = - \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] - \frac{\lambda_1}{N}$$

$$\rightarrow w_i = \frac{- \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] - \frac{\lambda_1}{N}}{2 \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right)}$$

Case 2:
What if $w_i < 0$, then we get:

$$\arg \min_{w_i \in \mathbb{R}} \quad \frac{2w_i}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) + \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) w_i^2 - \frac{\lambda_1}{N} w_i$$

$$\rightarrow \frac{\partial}{\partial w_i} \left(\frac{2w_i}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) + \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) w_i^2 - \frac{\lambda_1}{N} w_i \right) = 0$$

$$\rightarrow \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] + w_i \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right) - \frac{\lambda_1}{N} = 0$$

$$\rightarrow w_i = \frac{- \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] + \frac{\lambda_1}{N}}{2 \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right)}$$

\therefore For case 1:

$$- \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] - \frac{\lambda_1}{N} > 0 \quad (\because \text{this case is } w_i > 0)$$

$$\rightarrow - \frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] > \frac{\lambda_1}{N} > 0$$

$$\rightarrow \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) < - \frac{\lambda_1}{2} < 0$$

For case 2:

$$-\frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] + \frac{\lambda_1}{N} < 0 \quad (\text{if the case is } w_i < 0)$$

$$\Rightarrow -\frac{2}{N} \left[\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right] > -\frac{\lambda_1}{N} < 0$$

$$\Rightarrow \underline{\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) > \frac{\lambda_1}{2} > 0} \quad \triangleright$$

Δ :

The two results show that, if

$$\left| \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right| > \frac{\lambda_1}{2},$$

then we should update w_i as

$$\text{sign} \left(\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right) \frac{\left| -\frac{2}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right| - \frac{\lambda_1}{N}}{2 \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right)}$$

Therefore, we have $\alpha = \text{sign} \left(\sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right)$

$$\beta = \frac{\left| -\frac{2}{N} \sum_{n=1}^N \left(\sum_{j \neq i} w_j^{(t)} x_{n,j} x_{n,i} - x_{n,i} \right) \right| - \frac{\lambda_1}{N}}{2 \left(\frac{\sum_{n=1}^N (x_{n,i})^2}{N} + \frac{\lambda_2}{N} \right)}$$

□