

Homework # 5

姓名	學號
操之晴	R13922A04

5.

$$\begin{aligned}
(P_R) \quad & \min \frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 + \frac{\lambda}{N} \sum_{i=0}^d \alpha_i w_i^2 \\
(P_v) \quad & \min \frac{1}{N+K} \left(\sum_{n=1}^N (w^T x_n - y_n)^2 + \sum_{k=1}^K (w^T \tilde{x}_k - \tilde{y}_k)^2 \right) \\
&= \frac{1}{N+K} \left(\sum_{n=1}^N (w^T x_n - y_n)^2 + \sum_{k=1}^K (\lambda \alpha_{k-1} w_{k-1}^2)^2 \right) \\
&= \frac{1}{N+K} \left(\sum_{n=1}^N (w^T x_n - y_n)^2 + \lambda \sum_{i=0}^d \alpha_i w_i^2 \right) \\
&= \frac{N}{N+K} (P_R)
\end{aligned}$$

Thus, solving (P_v) is the same as the optimal solution obtained by solving (P_R) .

6.

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

From the Taylor expansion provided around w^* , we have:

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

Substitute this approximation into $\tilde{E}_{aug}(w)$

$$\tilde{E}_{aug}(w) = E_{in}(w^*) + \frac{1}{2}(w - w^*)^T H (w - w^*) + \frac{\lambda}{N} ||w||^2$$

The minimizer of $\tilde{E}_{aug}(w)$ is found by taking the derivative with respect to w and setting it to zero:

$$\nabla_w \tilde{E}_{aug}(w) = H(w - w^*) + \frac{2\lambda}{N} w = 0$$

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

So, we get:

$$w = \left(H + \frac{2\lambda}{N} I\right)^{-1} H w^*$$

7.

$$\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - \tilde{y})^2 \right)$$

we can decompose the expression $(y_n - \tilde{y})^2$ as follows:

$$(y_n - \tilde{y})^2 = (y_n - 0)^2 - 2(y_n - 0)(\tilde{y} - 0) + (\tilde{y} - 0)^2$$

Taking the expectation over the validation set, we get: $\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - \tilde{y})^2 \right) =$

$$\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - 0)^2 \right) - 2 \mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - 0)(\tilde{y} - 0) \right) + \mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (\tilde{y} - 0)^2 \right)$$

$\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - 0)^2 \right) = \sigma^2$ as the labels y_n are independent and identically distributed with mean = 0 and variance σ^2

$\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (y_n - 0)(\tilde{y} - 0) \right) = 0$ as \tilde{y} and y_n are independent, and both have mean 0.

$$\mathbb{E} \left(\frac{1}{K} \sum_{n=N-K+1}^N (\tilde{y} - 0)^2 \right) = \frac{1}{K} * K * \frac{\sigma^2}{N - K} = \frac{\sigma^2}{N - K}$$

. Therefore, the answer is $\sigma^2 + \frac{\sigma^2}{N - K}$

8.

The given optimal solution $w_0^* = \frac{1}{N} \sum_{n=1}^N y_n$ minimizes

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

In LOOCV, we compute the prediction error for each sample when it is left out of the training set. For the n -th sample, let w_{-n} denote the mean calculated without y_n

$$w_{-n} = \frac{1}{N-1} \sum_{i \neq n} y_i$$

The LOOCV error for the n -th sample is $(w_{-n} - y_n)^2$

The total LOOCV error is the average over all N samples $E_{LOOCV} = \frac{1}{N} \sum_{n=1}^N (w_{-n} - y_n)^2$

Using the relationship between w_{-n} and w_0^* , we observe:

$$w_{-n} = \frac{1}{N-1} \left(\sum_{i=1}^N y_i - y_n \right) = \frac{N}{N-1} w_0^* - \frac{1}{N-1} y_n$$

$$\begin{aligned} w_{-n} - y_n &= \left(\frac{N}{N-1} w_0^* - \frac{1}{N-1} y_n \right) - y_n \\ &= \left(\frac{N}{N-1} w_0^* - \frac{N}{N-1} y_n \right) = \left(\frac{N}{N-1} (w_0^* - y_n) \right) \end{aligned}$$

The LOOCV error is the average of the squared errors over all n :

$$\begin{aligned} E_{LOOCV} &= \frac{1}{N} \sum_{n=1}^N (w_{-n} - y_n)^2 = \frac{1}{N} \sum_{n=1}^N \left(\frac{N}{N-1} (w_0^* - y_n) \right)^2 \\ &= \left(\frac{N}{N-1} \right)^2 \frac{1}{N} \sum_{n=1}^N (w_0^* - y_n)^2 \\ &= \left(\frac{N}{N-1} \right)^2 E_{in}(w_0^*) \end{aligned}$$

$$\begin{aligned}
9. \quad E_{out} &= \frac{1}{N} \sum [g(x) \neq y] = \frac{1}{N} \sum [g(x) = -1, y = 1] + \frac{1}{N} \sum [g(x) = 1, y = -1] \\
&= \frac{1}{N} \sum \mathbb{P}(g(x) = -1 | y = 1) \mathbb{P}(y = 1) + \frac{1}{N} \sum \mathbb{P}(g(x) = 1 | y = -1) \mathbb{P}(y = -1) \\
&= \frac{1}{N} N \left(\frac{1}{2} \epsilon_+ \right) + \frac{1}{N} N \left(\frac{1}{2} \epsilon_- \right) = \frac{1}{2} \epsilon_+ + \frac{1}{2} \epsilon_-
\end{aligned}$$

$$\mathbb{P}(y = -1) = p, \quad \mathbb{P}(y = 1) = 1 - p$$

$$\begin{aligned}
E_{out}(g_c) &= \mathbb{P}(g(x) = -1 | y = 1) \mathbb{P}(y = 1) + \mathbb{P}(g(x) = 1 | y = -1) \mathbb{P}(y = -1) \\
&= \epsilon_+ (1 - p) + \epsilon_- (p) = p
\end{aligned}$$

$$\text{So, } p = \frac{\epsilon_+}{1 + \epsilon_+ - \epsilon_-}$$

10.

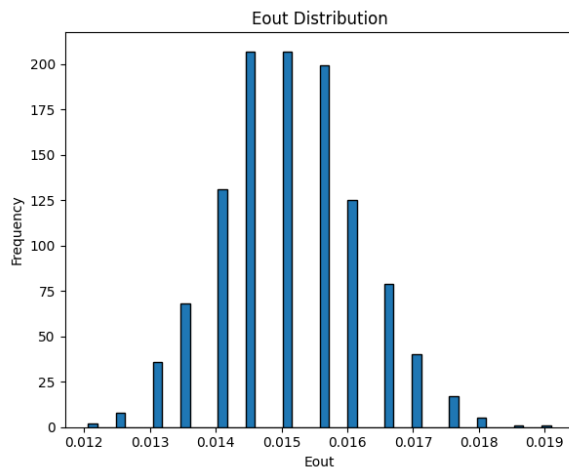
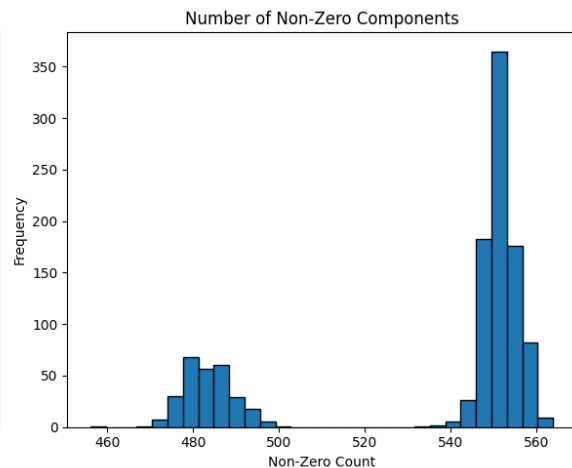
Figure 10a. histogram of $E_{out}(g)$ 

Figure 10b. histogram of non-zero components

```

HWS > src > q10.py > ...
1 from liblinear.liblinearutil import *
2 import numpy as np
3 from sklearn.datasets import load_svmlight_file
4 import matplotlib.pyplot as plt
5 from tqdm import tqdm
6
7 # parameters
8 EXPERIMENTS = 1126
9 lambdas = [-2, -1, 0, 1, 2, 3]
10
11 # loading data
12 train_x, train_y = load_svmlight_file('mnist.scale')
13 test_x, test_y = load_svmlight_file('mnist.scale.t')
14
15 train_x = train_x.astype(np.float64)
16 test_x = test_x.astype(np.float64)
17
18 # filter class 2 & 6
19 def filter_data(x, y, class1=2, class2=6):
20     mask = (y == class1) | (y == class2)
21     x_filtered = x[mask]
22     y_filtered = np.where(y[mask] == class1, 1, -1) # class 2 : 1, class 6 : -1
23     return x_filtered, y_filtered
24
25 train_x, train_y = filter_data(train_x, train_y)
26 test_x, test_y = filter_data(test_x, test_y)
27
28 # finding best lambda
29 Ein_histogram = []
30 Eout_histogram = []
31 non_zero_histogram = []
32
33 for i in tqdm(range(EXPERIMENTS), desc="Experiments"):
34     best_lambda = None
35     best_Ein = float('inf')
36
37     for log_lambda in lambdas:
38         current_lambda = 10 ** log_lambda
39         c = 1 / (2 * current_lambda)
40
41         model = train(train_y, train_x, f'~s 6 -c (c) -0.1 -q')
42
43         Ein_acc, _ = predict(train_y, train_x, model, '-q')
44         Ein_error = 1 - (Ein_acc[0] / 100) # 0/1 error

```

Figure 10c. screenshots of my code

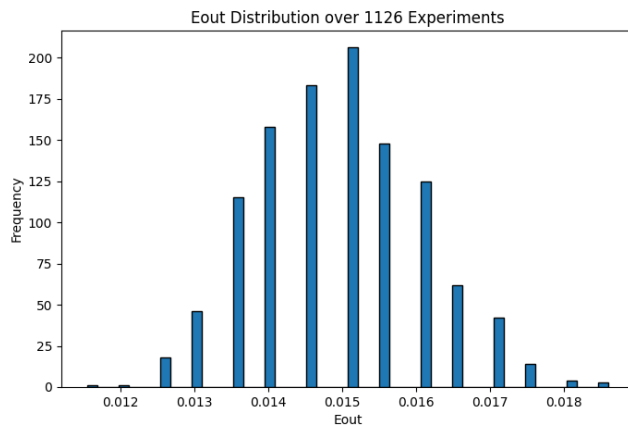
Eout Distribution

The test error E_{out} is mainly concentrated between $[0.014, 0.017]$, and the peak is around 0.015, showing very stable model performance. It looks like a normal distribution, indicating that different random seeds have minimal impact.

Number of Non-Zero Components

The number of non-zero weights shows two clusters, one around 540 and another around 480. This could be because the best λ values chosen in different experiments mainly fall into two categories: Larger λ (0.1) stronger regularization, shrinking more weights, leading to fewer non-zero components. Smaller λ (0.01) relaxes regularization, allowing more weights to remain non-zero.

11.

Figure 11a. histogram of $E_{out}(g)$

```

HW5 > # 0 q11.py > ...
1 from sklearn.linear_model import LogisticRegression
2 import numpy as np
3 from sklearn.datasets import load_svmlight_file
4 import matplotlib.pyplot as plt
5 from tqdm import tqdm
6 from sklearn.model_selection import train_test_split
7
8 # parameters
9 EXPERIMENTS = 1126 # 实验次数
10 lambdas = [-2, -1, 0, 1, 2, 3]
11 SUB_TRAIN_SIZE = 8000 # 子训练集大小
12
13 # loading data
14 train_x, train_y = load_svmlight_file('mnist.scale')
15 test_x, test_y = load_svmlight_file('mnist.scale.t')
16
17 train_x = train_x.astype(np.float64)
18 test_x = test_x.astype(np.float64)
19
20 # filter class 2 & 6
21 def filter_data(x, y, class1=2, class2=6):
22     mask = (y == class1) | (y == class2)
23     x_filtered = x[mask]
24     y_filtered = np.where(mask == class1, 1, -1) # class 2 : 1, class 6 : -1
25     return x_filtered, y_filtered
26
27 train_x, train_y = filter_data(train_x, train_y)
28 test_x, test_y = filter_data(test_x, test_y)
29
30 Eout_histogram = {}
31
32 for i in tqdm(range(EXPERIMENTS), desc="Experiments"):
33     # splitting sub train set
34     sub_train_x, val_x, sub_train_y, val_y = train_test_split(
35         train_x, train_y, train_size=SUB_TRAIN_SIZE)
36
37     # finding best lambda on subset
38     best_lambda = None
39     best_eval = float("inf")
40
41     for log_lambda in lambdas:
42         current_lambda = 10 ** log_lambda
43         c = 1 / (2 * current_lambda)

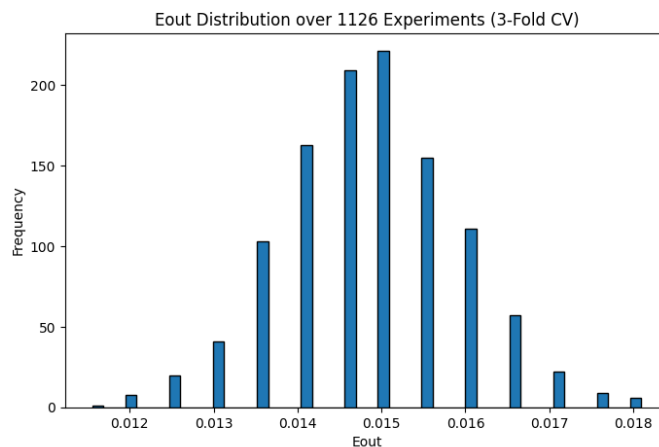
```

Figure 11b. screenshots of my code

The E_{out} range is roughly between $[0.012, 0.018]$, which is slightly wider than in Problem 10. The distribution is still bell-shaped, with a peak around 0.015. However, the tails are slightly longer. This could be because problem 10 is more stable since the entire training set was used to find the best λ^* . Problem 11 shows more instability due to data splitting.

In conclusion, the model in Problem 10 is more stable, with concentrated test errors, making it suitable for evaluating the best model performance. Problem 11 is closer to real world, as data splitting simulates actual training/validation processes, but it introduces some instability.

12.

Figure 12a. histogram of $E_{out}(g)$

```

HMS > src > 000.py
1 from sklearn.linear_model import *
2 import numpy as np
3 from sklearn.datasets import load_svmlight_file
4 from sklearn.model_selection import KFold
5 import matplotlib.pyplot as plt
6 from tqdm import tqdm
7
8 # parameters
9 EXPERIMENTS = 1126
10 lambdas = [-2, -1, 0, 1, 2, 3]
11 K_FOLDS = 3
12
13 # loading data
14 train_x, train_y = load_svmlight_file('mist-scale')
15 test_x, test_y = load_svmlight_file('mist-scale.t')
16
17 train_x = train_x.astype(np.float64)
18 test_x = test_x.astype(np.float64)
19
20 # filter class 2 & 8
21 def filter_data(x, y, class1=2, class2=8):
22     mask = (y == class1) | (y == class2)
23     x_filtered = x[mask]
24     y_filtered = np.where(mask == class1, 1, -1) # class 2 : 1-class 8 : -1
25     return x_filtered, y_filtered
26
27 train_x, train_y = filter_data(train_x, train_y)
28 test_x, test_y = filter_data(test_x, test_y)
29
30 Eout_histogram = []
31
32 for i in tqdm(range(EXPERIMENTS), desc="Experiments"):
33     # 3-fold cross validation
34     kf = KFold(n_splits=K_FOLDS, shuffle=True)
35     best_lambda = None
36     best_ecv = float('inf')
37
38     for log_lambda in lambdas:
39         current_lambda = 10 ** log_lambda
40         c = 1 / (2 * current_lambda)
41         fold_errors = []
42
43         # split current fold (train & valid set)
44         for train_idx, val_idx in kf.split(train_x):
45             sub_train_x, val_x = train_x[train_idx], train_x[val_idx]
46             sub_train_y, val_y = train_y[train_idx], train_y[val_idx]

```

Figure 12b. screenshots of my code

Problem 12 is very similar to problem 11. Both have their error distribution centers around 0.015 with similar shapes. However, there's still some slightly difference. In Problem 12, E_{out} is mainly concentrated between $[0.014, 0.016]$, making it more focused than in Problem 11. Extreme values also appear less frequently compared to Problem 11, and the tails are shorter. It indicated that problem 12 is more stable, with higher error concentration and almost no extreme values.

This might be because problem 11 uses a single random split, and the distribution of the sub-training set and validation set might be uneven. This causes more variation in the selected λ^* , leading to greater fluctuations in the results. Problem 12 uses 3-fold cross-validation, which reduces the impact of uneven distributions by validating in turns. This makes the selected λ^* more representative and the errors more concentrated.

For problems 10 to 12, the full code has been made available on [GitHub](#).

Feel free to check it out if required or interested!