

# Homework 1

PENG Qiheng

Student ID 225040065

September 27, 2025

Problem 1.

(a) In supervised learning, we have a dataset of input data and labels, and the goal is to learn a mapping from inputs to labels.

In unsupervised learning, we only have input data and the goal is to discover the hidden patterns or structures in the data.

(b) Answer: 2

1. Regression is used to fit continuous labels.

3. For a given training dataset, perceptrons with different initialization may converge to different linear classifiers.

4. Least squares is a maximum likelihood estimator only if error follows a Gaussian Distribution.

(c) 1. Since we have

$$r(\mathbf{X}^T \mathbf{X}) = r(\mathbf{X}) = d \quad (1)$$

$$\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{d \times d} \quad (2)$$

thus,  $\mathbf{X}^T \mathbf{X}$  is invertible

2. As

$$(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X} \quad (3)$$

we know that  $\mathbf{X}^T \mathbf{X}$  is Symmetric Matrices

3. Thus, we can conclude that  $\mathbf{X}^T \mathbf{X}$  is Positive Definite

Problem 2.

(a) Let  $\mathbf{X} = \mathbf{V}\Sigma_1\mathbf{U}_1^T$ ,  $\mathbf{z} = \mathbf{U}_1^T\boldsymbol{\theta}$  and  $\mathbf{A} := \mathbf{V}\Sigma_1$  we have

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{V}\Sigma_1\mathbf{U}_1^T\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2^2 \quad (4)$$

Since  $\mathbf{A}$  is of full rank, we have optimal  $\mathbf{z}$

$$\mathbf{z}^* = \mathbf{A}^{-1}\mathbf{y} \quad (5)$$

And we have

$$\mathbf{U}_1^T\boldsymbol{\theta}^* = \mathbf{A}^{-1}\mathbf{y} \quad (6)$$

$$\mathbf{U}_1^T\boldsymbol{\theta}^* - \mathbf{A}^{-1}\mathbf{y} = \mathbf{0} \quad (7)$$

Considering that  $\mathbf{U}_1^T\mathbf{U}_2 = \mathbf{0}$

$$\mathbf{U}_1^T\boldsymbol{\theta}^* - \mathbf{A}^{-1}\mathbf{y} = \mathbf{U}_1^T\mathbf{U}_2 \quad (8)$$

$$\boldsymbol{\theta}^* = \mathbf{U}_1(\mathbf{V}\Sigma_1)^{-1}\mathbf{y} + \mathbf{U}_2\mathbf{c}, \mathbf{c} \in \mathbb{R}^{d-n} \quad (9)$$

And the optimal function value is 0, if  $\mathbf{c} = \mathbf{0}$

(b) Let

$$\mathcal{L}(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\theta}\|_2^2 \quad (10)$$

Take the gradient

$$\nabla\mathcal{L}(\boldsymbol{\theta}) = 2\mathbf{X}^T(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) + 2\lambda\boldsymbol{\theta} \quad (11)$$

$$(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})\boldsymbol{\theta}^* = \mathbf{X}^T\mathbf{y} \quad (12)$$

Since  $\mathbf{X}^T\mathbf{X}$  is semi-positive definite and  $\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}$  is positive definite.

Thus, we have

$$\boldsymbol{\theta}^* = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} \quad (13)$$

Problem 3.

(a) As the assumption, we have

$$\epsilon_i = y_i - \mathbf{X}_i \boldsymbol{\theta} \quad (14)$$

$$P(D|\boldsymbol{\theta}) = \prod_{i=1}^n \left( \frac{1}{2b} e^{-\frac{|y_i - \mathbf{X}_i \boldsymbol{\theta}|}{b}} \right) = 2b^{-n} e^{-\frac{|\mathbf{y} - \mathbf{X} \boldsymbol{\theta}|}{b}} \quad (15)$$

Let  $\mathcal{L}(\boldsymbol{\theta}) = \log P(D|\boldsymbol{\theta})$ , we have

$$\mathcal{L}(\boldsymbol{\theta}) = -n \log(2b) - \frac{1}{b} |\mathbf{y} - \mathbf{X} \boldsymbol{\theta}| \quad (16)$$

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} |\mathbf{y} - \mathbf{X} \boldsymbol{\theta}| \quad (17)$$

(b) According to the definition of  $h_\mu(z)$ , we have

$$\frac{\partial h_\mu(z)}{\partial z} = \begin{cases} \text{sign}(z), & |z| \geq \mu \\ \frac{z}{\mu}, & |z| \leq \mu \end{cases} \quad (18)$$

And we can rewrite  $h'_\mu(z)$  as

$$h'_\mu(z) = \frac{z}{\max(|z|, \mu)} \quad (19)$$

Thus, we have

$$H'_\mu(\mathbf{z}) = \frac{\partial \sum_{j=1}^n h_\mu(z_j)}{\partial \mathbf{z}} = \begin{bmatrix} h'_\mu(z_1) \\ \vdots \\ h'_\mu(z_n) \end{bmatrix} = h'_\mu(\mathbf{z}) \quad (20)$$

Finally, we have

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = \frac{\partial(\mathbf{X} \boldsymbol{\theta} - \mathbf{y})}{\partial \boldsymbol{\theta}} \frac{\partial H_\mu(\mathbf{z})}{\partial \mathbf{z}} = \mathbf{X}^T h'_\mu(\mathbf{X} \boldsymbol{\theta} - \mathbf{y}) \quad (21)$$

(c)

After calculating, output  $\|\hat{\boldsymbol{\theta}}_{LS} - \boldsymbol{\theta}^*\| = 59.5136$

code as follows:

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  d = 50 # feature dimension
4  X = np.load("data/X.npy")
5  y = np.load("data/y.npy")
6  print("data_shape:", X.shape, y.shape)
7  theta_star = np.load("data/theta_star.npy")
8  ##### part (1): least square estimator #####
9  theta_hat = np.linalg.inv(X.T @ X) @ X.T @ y
10 Error_LS = np.linalg.norm(theta_hat - theta_star, 2)
11 print("Estimator_approximated_by_LS:", Error_LS)
12 ##### part (2): L1 estimator #####
13 mu = 1e-5 # smoothing parameter
14 alpha = 0.001 # stepsize
15 T = 1000 # iteration number
16 # random initialization
17 theta = np.random.randn(d, 1)
18 Error_huber = []
19 for _ in range(1, T):
20     # calculate the l2 error of the current iteration
21     Error_huber.append(np.linalg.norm(theta - theta_star, 2))
22     # calculate gradient
23     z = X @ theta - y
24     h_prime = z / np.maximum(np.abs(z), mu)
25     grad = X.T @ h_prime
26     # gradient descent update
27     theta = theta - alpha * grad
28     ##### plot the figure #####
29     plt.figure(figsize=(10, 5))
30     plt.yscale("log", base=2)
31     plt.plot(Error_huber, "r-")
32     plt.title(r"$\ell_1$ estimator approximated by Huber")
33     plt.ylabel(r"$\theta$") # set the label for the y axis
34     plt.xlabel("Iteration") # set the label for the x axis
35     plt.show()
36

```

The Huber smoothing error is shown in Figure 1.

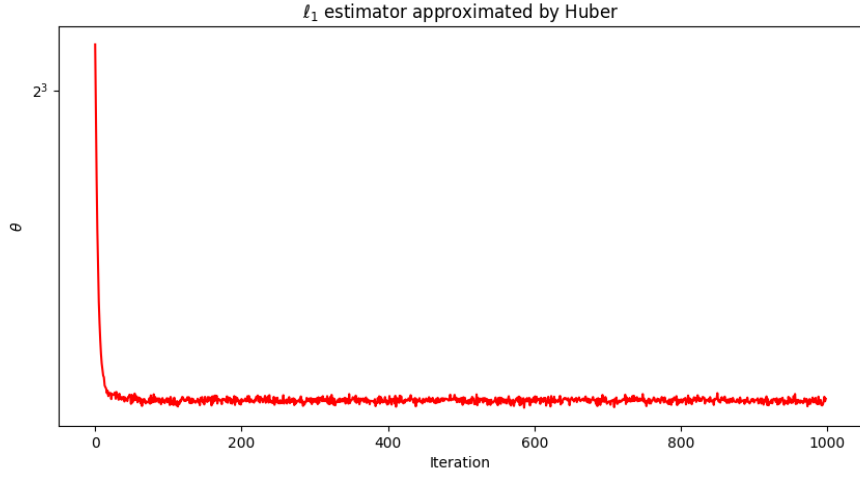


Figure 1: Huber smoothing error

Problem 4.

- (a) Since  $\theta^*$  is a classifier that correctly separates all the data points, we can get

$$y_i(\theta^{*T} \mathbf{x}_i) > 0, \forall i = 1, \dots, n$$

We can easily have

$$\rho \geq y_i(\theta^{*T} \mathbf{x}_i) > 0$$

- (b) According to update (1), we have

$$\theta_k^T \theta^* = \theta_{k-1}^T \theta^* + y_{k-1} \mathbf{x}_{k-1}^T \theta^* \geq \theta_{k-1}^T \theta^* + \rho$$

After recursion, we have

$$\theta_k^T \theta^* \geq k\rho$$

- (c) According to update (1), we have

$$\begin{aligned} \|\theta_k\|^2 &= \|\theta_{k-1} + y_{k-1} \mathbf{x}_{k-1}^T\|^2 \\ &= \|\theta_{k-1}\|^2 + \|\mathbf{x}_{k-1}\|^2 + 2y_{k-1} \theta_{k-1}^T \mathbf{x}_{k-1} \end{aligned}$$

Since  $\mathbf{x}_{k-1}$  is misclassified by  $\boldsymbol{\theta}_{k-1}$ , we have  $y_{k-1}\boldsymbol{\theta}_{k-1}^T\mathbf{x}_{k-1} \leq 0$

$$\|\boldsymbol{\theta}_k\|^2 \leq \|\boldsymbol{\theta}_{k-1}\|^2 + \|\mathbf{x}_{k-1}\|^2$$

(d) According to the problem, we have  $R \geq \|\mathbf{x}_i\|, \forall i = 1, \dots, n$  And we have

$$\|\boldsymbol{\theta}_k\|^2 \leq \|\boldsymbol{\theta}_{k-1}\|^2 + \|\mathbf{x}_{k-1}\|^2 \leq \|\boldsymbol{\theta}_{k-1}\|^2 + R^2$$

After recursion, we have

$$\|\boldsymbol{\theta}_k\|^2 \leq kR^2$$

(e) According to steps 2 and 4, we have  $\boldsymbol{\theta}_k^T\boldsymbol{\theta}^* \geq k\rho$  and  $\|\boldsymbol{\theta}_k\|^2 \leq kR^2$  we can get

$$\frac{\boldsymbol{\theta}_k^T\boldsymbol{\theta}^*}{\|\boldsymbol{\theta}_k\|} \geq \sqrt{k}\frac{\rho}{R}$$

By using Cauchy-Schwarz inequality, we have

$$k^2\rho^2 \leq \|\boldsymbol{\theta}_k\|^2\|\boldsymbol{\theta}^*\|^2 \leq kR^2\|\boldsymbol{\theta}^*\|^2$$

Thus, we have

$$\bar{k} \leq \frac{R^2\|\boldsymbol{\theta}^*\|^2}{\rho^2}$$

Problem 5. Main code as follows:

```
1      # train
2      iters = 2000
3      d = 2
4      num_sample = X.shape[0]
5      threshold = 1e-4
6      theta = np.zeros((d + 1, 1))
7
8      X = np.hstack([X, np.ones((num_sample, 1))])
9      X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
10     Er_in_perceptron = []
11     Er_out_perceptron = []
12     Er_in_pocket = []
13     Er_out_pocket = []
14     pocket = theta.copy()
15     best_error = cal_error(theta, X, y)
16
17     for iterate in range(iters):
18         for i in random.sample(range(num_sample), num_sample):
19             if np.sign(X[i] @ theta)[0] != y[i]:
20                 theta += (y[i] * X[i]).reshape(-1, 1)
21                 break
22
23         cur_error = cal_error(theta, X, y)
24         if cur_error < best_error:
25             best_error = cur_error
26             pocket = theta.copy()
27
28         Er_in_perceptron.append(cur_error)
29         Er_out_perceptron.append(cal_error(theta, X_test, y_test))
30         Er_in_pocket.append(best_error)
31         Er_out_pocket.append(cal_error(pocket, X_test, y_test))
32
33     # print(f"theta (perceptron): {theta}")
34     # print(f"pocket (pocket): {pocket}")
35
36     # plot Er_in and Er_out
37     fig, axs = plt.subplots(1, 2, figsize=(14, 6))
38
39     # perceptron
40     axs[0].plot(Er_in_perceptron, label="Perceptron_Train_Error")
```

```

41     axs[0].plot(Er_in_pocket, label="Pocket_Train_Error")
42     axs[0].set_xlabel("Iteration")
43     axs[0].set_ylabel("Error_Rate")
44     axs[0].set_title("Training_Error")
45     axs[0].legend()
46     # pocket
47     axs[1].plot(Er_out_perceptron, label="Perceptron_Test_Error")
48     axs[1].plot(Er_out_pocket, label="Pocket_Test_Error")
49     axs[1].set_xlabel("Iteration")
50     axs[1].set_ylabel("Error_Rate")
51     axs[1].set_title("Test_Error")
52     axs[1].legend()
53
54     plt.tight_layout()
55     plt.show()
56
57     # plot decision boundary
58     fig, axs = plt.subplots(1, 2, figsize=(14, 6))
59
60     ax0 = plot_feature(X, y, plot_num=500, ax=axs[0], classes=np.unique(y))
61     x1 = np.linspace(X[:, 0].min(), X[:, 0].max(), 100)
62     x2_perceptron = -(theta[0] * x1 + theta[2]) / theta[1]
63     x2_pocket = -(pocket[0] * x1 + pocket[2]) / pocket[1]
64     ax0.plot(x1, x2_perceptron, "r--", label="Perceptron_Boundary")
65     ax0.plot(x1, x2_pocket, "b-", label="Pocket_Boundary")
66     ax0.set_title("Decision_Boundaries_(Train)")
67     ax0.legend()
68
69     ax1 = plot_feature(X_test, y_test, plot_num=500, ax=axs[1], classes=np.unique(y_test))
70     x1_test = np.linspace(X_test[:, 0].min(), X_test[:, 0].max(), 100)
71     x2_perceptron_test = -(theta[0] * x1_test + theta[2]) / theta[1]
72     x2_pocket_test = -(pocket[0] * x1_test + pocket[2]) / pocket[1]
73     ax1.plot(x1_test, x2_perceptron_test, "r--", label="Perceptron_Boundary")
74     ax1.plot(x1_test, x2_pocket_test, "b-", label="Pocket_Boundary")
75     ax1.set_title("Decision_Boundaries_(Test)")
76     ax1.legend()
77
78     plt.tight_layout()
79     plt.show()
80

```



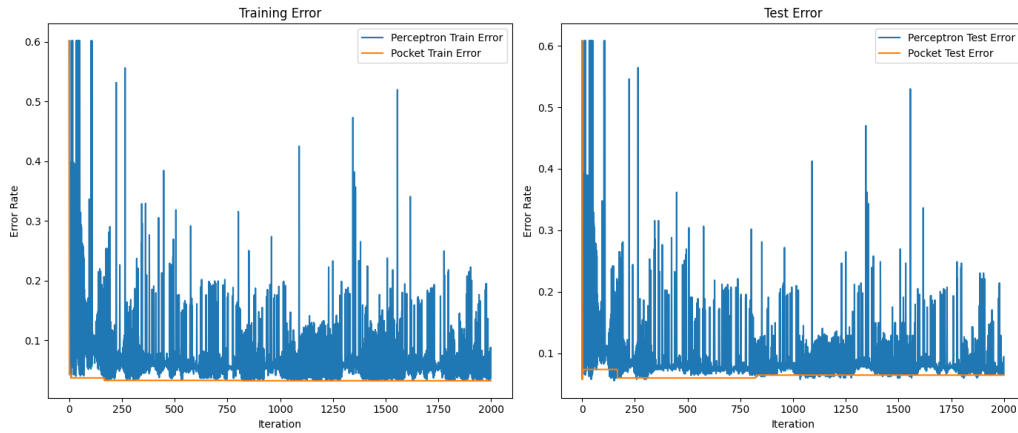


Figure 2: Training and Test Error Rate

From the Figure 2, we observe that:

The error of the Perceptron algorithm fluctuates more, showing that the Perceptron is sensitive to noisy or non-separable samples.

The Pocket algorithm maintains and uses the best solution found so far, so its test error curve is smoother and lower than that of the Perceptron. The error of pocket algorithms decreases as the number of iterations increases.

And the out-of-sample error is generally higher than the in-sample error as we study in lectures.

Overall, the Pocket algorithm is more robust to outliers and non-separable data.

The decision boundaries of the two algorithms are shown in Figure 3.

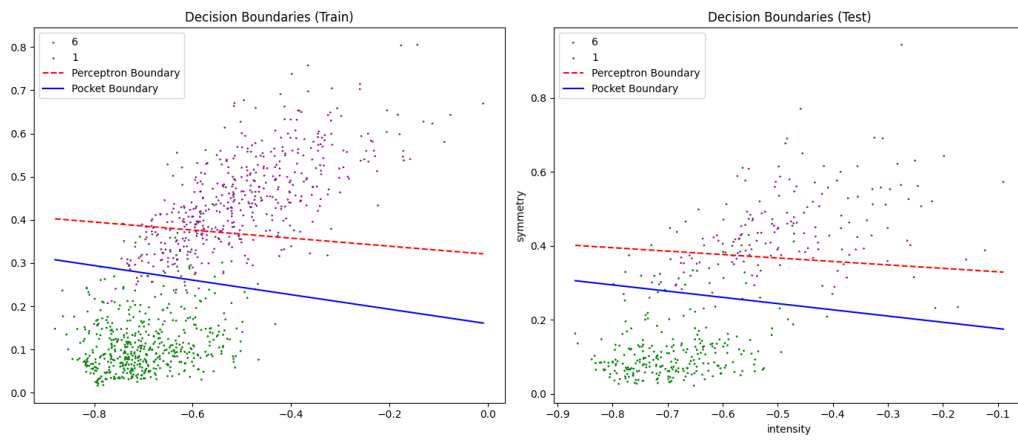


Figure 3: Decision Boundaries