Problem 1

(a) Supervised Learning learns from labeled examples to make predictions, while Unsupervised Learning discovers patterns in data without any guidance from labels.

(b) 2) True

(c) data matrix $X \in \mathbb{R}^{n \times d}$ has full column rank.

→ Explain why $X^T X$ is PD (Hint: linear independence).

↳ for any matrix $A \neq 0$,

$$A^T \underline{X^T X} A > 0 .$$

$XA = 0$, only when $A = 0$.

Proof: $(X^T X)^T = X^T X$, thus $X$ is symmetric.

Let matrix $A \in \mathbb{R}^d$, $A \neq 0$.

Because $X \in \mathbb{R}^{n \times d}$ has full column rank.

So $X$ is linear independent,

so $XA \neq 0$, $\|XA\|_2^2 > 0$

$$\|XA\|_2^2 = (XA)^T(XA) = A^T X^T X A > 0$$

so $X^T X$ is Positive Definite.

(c)

Problem 2

(2.a)

Problem 2.

(a) $\quad \min_{\theta \in \mathbb{R}^d} \| X\theta - y \|_2^2$

$= \min_{\theta \in \mathbb{R}^d} \| V\Sigma_1 U_1^T \theta - y \|_2^2$

Let $A := V\Sigma_1$ which is square matrix of full rank.

Let $z = U_1^T \theta$.

Solve: $\mathcal{L}(z) = \min_z \| Az - y \|_2^2$

$= (Az - y)^T (Az - y) = z^T A^T A z - 2z^T A^T y + y^T y$

$\nabla_z \mathcal{L}(z) = 2A^T(Az - y)$

set $\nabla_z \mathcal{L}(z) = 0 \quad \Rightarrow$ the optimal solution $\hat{z}$

satisfies: $A^T A \hat{z} = A^T y$

Because $A \in \mathbb{R}^{n \times n}$ is a square matrix of full rank.

$\Rightarrow A^T A \in \mathbb{R}^{n \times n}$ is full rank and invertible

We have $\hat{z} = (A^T A)^{-1} A^T y$

So, $\hat{z} = A^\dagger y \quad \Rightarrow U_1^T \hat{\theta} = (V\Sigma_1)^\dagger y = (\Sigma_1^T V^T V\Sigma_1)^{-1} \Sigma_1^T V^T y$

$= (\Sigma_1^T \Sigma_1)^{-1} \Sigma_1^T (\Sigma_1 \Sigma_1^{-1}) V^T y$

$= \Sigma_1^{-1} V^T y$

we want to solve: $\boxed{U_1^T \hat{\theta} = \Sigma_1^{-1} V^T y}$

this is a system of $\underline{n}$ equations in $\underline{d}$ unknowns $(d > n)$.

so, there is infinite many solutions.

$\Rightarrow$ one particular solution $\theta_p$:

$\quad$ Try: $\theta_p = U_1 \hat{z} = U_1 \Sigma_1^{-1} V^T y$

$\quad\quad U_1^T \theta_p = (U_1^T U_1) \Sigma_1^{-1} V^T y = \Sigma_1^{-1} V^T y \quad \checkmark$

$\Rightarrow$ homogeneous solution $\theta_h$:

$\quad$ solve: $U_1^T \theta_h = 0$

$\quad$ we know $U_1^T U_2 = 0$

$\quad$ So: $\theta_h = U_2 \alpha \quad (\text{for any } \alpha \in \mathbb{R}^{d-n})$.

$\Rightarrow$ general solution :

$\quad \hat{\theta} = \theta_p + \theta_h = U_1 \Sigma_1^{-1} V^T y + U_2 \alpha \quad (\forall \alpha \in \mathbb{R}^{d-n})$

• The optimal function value for $\min_{\theta \in \mathbb{R}^d} \| X\theta - y \|_2^2$

Let $\theta = \theta_p$,

Optimal function value $= \| V\Sigma_1 U_1^T \theta_p - y \|_2^2$

$= \| V\Sigma_1 U_1^T U_1 \Sigma_1^{-1} V^T y - y \|_2^2$

$= \| y - y \|_2^2 = 0$

KOKUYO

(2.b)

(b). Solve $\min\limits_{\theta \in \mathbb{R}^d} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2$

Let $J(\theta) = \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2$

$$= (X\theta - y)^T (X\theta - y) + \lambda \theta^T \theta$$

$$= \theta^T X^T X \theta - 2\theta^T X^T y + y^T y + \lambda \theta^T \theta$$

$$= \theta^T (X^T X + \lambda I)\theta - 2\theta^T X^T y + y^T y$$

$\nabla_\theta J(\theta) = 2(X^T X + \lambda I)\theta - 2X^T y$

Set $\nabla_\theta J(\theta) = 0 \Rightarrow \boxed{(X^T X + \lambda I)\theta = X^T y}$

Because $X \in \mathbb{R}^{n \times d}$, $n < d$

$\Rightarrow X^T X$ is ~~s.e.~~ positive semi-definite.

note that $\lambda > 0$.

therefore $(\lambda I + X^T X)$ is positive definite.

and invertible.

so. $\hat{\theta} = (X^T X + \lambda I)^{-1} X^T y$

## Problem 3

### (3.a)

**Problem 3.**

(a) $y = X\theta^* + \epsilon$ $\Rightarrow$ $\epsilon = y - X\theta^*$

$\Rightarrow$ $\epsilon_i = y_i - x_i^T\theta$

$\epsilon_i \sim L(0, b)$ , $p(\epsilon_i) = \frac{1}{2b} e^{-\frac{|\epsilon_i|}{b}}$ , $b > 0$.

So, $p(y_i | x_i, \theta) = \frac{1}{2b} \exp\left(-\frac{|y_i - x_i^T\theta|}{b}\right)$

since the data are independent:

$$P(Y | X, \theta) = \prod_{i=1}^{n} P(y_i | x_i, \theta) = \prod_{i=1}^{n} \frac{1}{2b} \exp\left(-\frac{|y_i - x_i^T\theta|}{b}\right)$$

$$L(\theta) = \log P(Y | X, \theta) = \sum_{i=1}^{n} \log\left(\frac{1}{2b} \exp\left(-\frac{|y_i - x_i^T\theta|}{b}\right)\right)$$

$$= \sum_{i=1}^{n} \left(-\log(2b) - \frac{|y_i - x_i^T\theta|}{b}\right)$$

$$= -n \log(2b) - \frac{1}{b}\sum_{i=1}^{n} |y_i - x_i^T\theta|$$

$\hat{\theta} = \underset{\theta}{\text{argmax }} L(\theta)$

$= \underset{\theta}{\text{argmax}} \left(-\frac{1}{b}\sum_{i=1}^{n} |y_i - x_i^T\theta|\right)$

$= \underset{\theta}{\text{argmin}} \sum_{i=1}^{n} |y_i - x_i^T\theta| = \underset{\theta}{\text{argmin}} \|y - X\theta\|_1$

### (3.b)

(b) $\nabla_\theta L(\theta) = \nabla_\theta H_\mu(X\theta - y)$

Let $r = X\theta - y$ $\Rightarrow$ $\nabla_\theta L(\theta) = X^T \nabla_r H_\mu(r) \big|_{r = X\theta - y}$

$$= X^T \nabla_r \left(\sum_{j=1}^{n} h_\mu(r_j)\right)$$

$$\nabla_r H_\mu(r) = \begin{bmatrix} \frac{d}{dr_1} h_\mu(r_1) \\ \frac{d}{dr_2} h_\mu(r_2) \\ \vdots \\ \frac{d}{dr_j} h_\mu(r_j) \end{bmatrix}$$

$$h_\mu(z) = \begin{cases} |z| & , \ |z| \ge \mu \\ \dfrac{z^2}{2\mu} + \dfrac{\mu}{2} & , \ |z| \le \mu. \end{cases}$$

$$\Rightarrow \quad \frac{d}{dz} h_\mu(z) = \begin{cases} \text{sign}(z), & |z| \ge \mu. \\ \dfrac{z}{\mu}, & |z| \le \mu \end{cases}$$

$$\text{Let} \quad g_j = \frac{d}{dr_j} h_\mu(r_j) = \begin{cases} \text{sign}(r_j), & |r_j| \le \mu \\ \dfrac{r_j}{\mu}, & |r_j| \ge \mu. \end{cases}$$
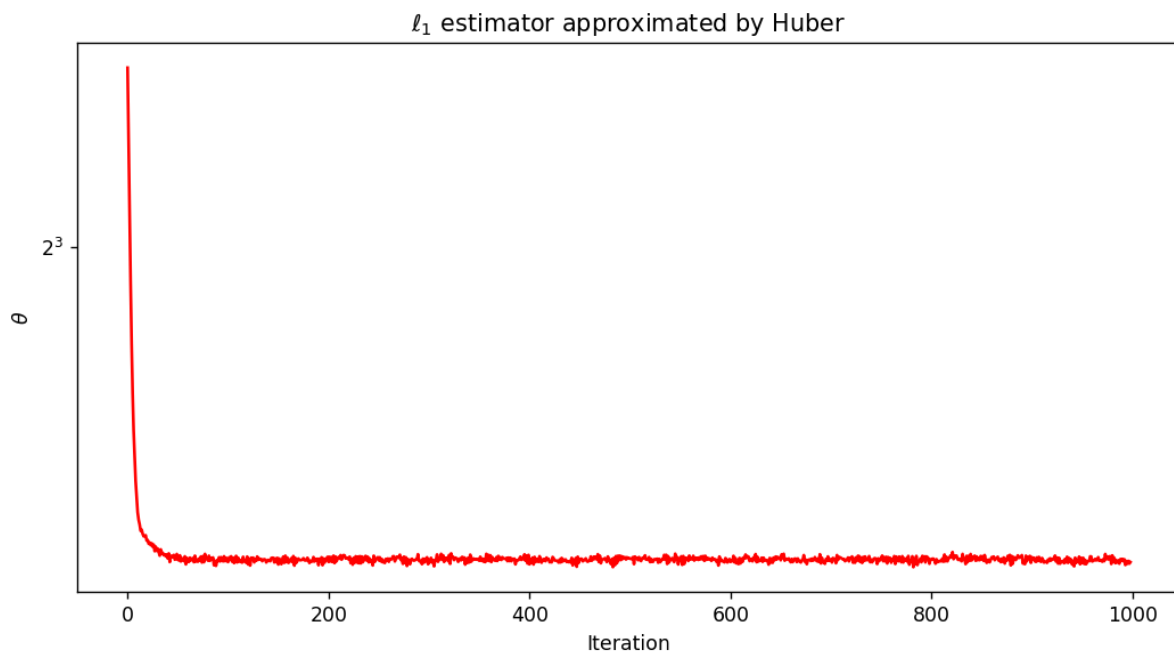
$\bullet \quad \nabla_r H_\mu(r) = g.$

Now, $\nabla_\theta L(\theta) = X^T \nabla_r H_\mu(r) = X^T g$

$$\text{where} \quad g_j = \begin{cases} \text{sign}(r_j), & |r_j| \le \mu \\ \dfrac{r_j}{\mu}, & |r_j| \ge \mu. \end{cases}$$

$$r = X\theta - y$$

(3.c)

Code: p3.py



$\ell_1$ estimator approximated by Huber

# Problem 4

## (a). proof:

since the data is linearly separable,

we have: $y_i(\theta^{*T}x_i) > 0$ for all $i = 1, 2, \ldots n$.

Therefore, $\rho = \min\limits_{1 \le i \le n} y_i(\theta^{*T}x_i) > 0$.

**Q.E.D.**

## (b). proof:

$$\theta_k = \theta_{k-1} + y_{k-1}x_{k-1}$$

$$\Rightarrow \theta_k^T\theta^* = (\theta_{k-1} + y_{k-1}x_{k-1})^T\theta^*$$

$$= \theta_{k-1}^T\theta^* + y_{k-1}x_{k-1}^T\theta^* = \theta_{k-1}^T\theta^* + y_{k-1}\theta^{*T}x_{k-1}$$

because, $y_{k-1}\theta^{*T}x_{k-1} \ge \min\limits_{1 \le k \le n} y_k(\theta^{*T}x_k) = \rho$

Therefore, $\cancel{\theta_k^T\theta^* \ge \theta_{k-1}^T\theta^*}$

$$\theta_k^T\theta^* = \theta_{k-1}^T\theta^* + y_{k-1}\theta^{*T}x_{k-1} \ge \theta_{k-1}^T\theta^* + \rho$$

**Q.E.D.**

We can know: $\theta_1^T\theta^* \ge \theta_0^T\theta^* + \rho = 0 + \rho = \rho$

$$\theta_2^T\theta^* \ge \theta_1^T\theta^* + \rho \ge 2\rho$$

$$\vdots$$

By induction:

$$\theta_k^T\theta^* \ge k\rho.$$

**Q.E.D.**

## (c)

$$\|\theta_k\|^2 = \|\theta_{k-1} + y_{k-1}x_{k-1}\|^2 = (\theta_{k-1} + y_{k-1}x_{k-1})^T(\theta_{k-1} + y_{k-1}x_{k-1})$$

$$= (\theta_{k-1}^T + y_{k-1}x_{k-1}^T)(\theta_{k-1} + y_{k-1}x_{k-1})$$

$$= \theta_{k-1}^T\theta_{k-1} + 2y_{k-1}(\theta_{k-1}^Tx_{k-1}) + y_{k-1}^2 x_{k-1}^T x_{k-1}$$

$$= \|\theta_{k-1}\|^2 + 2y_{k-1}(\theta_{k-1}^Tx_{k-1}) + y_{k-1}^2\|x_{k-1}\|^2$$

Note: $y_{k-1}^2 = 1$, since $y_i \in \{-1, 1\}$

$y_{k-1}\theta_{k-1}^Tx_{k-1} < 0$, because $x_{k-1}$ is misclassified.

Therefore $\|\theta_k\|^2 \le \|\theta_{k-1}\|^2 + \|x_{k-1}\|^2$

**Q.E.D.**

(d). from question (4.c):

$$\|\theta_k\|^2 \leq \|\theta_{k-1}\|^2 + \|x_{k-1}\|^2 \leq \|\theta_{k-1}\|^2 + R^2$$

$$(R = \max_{1 \leq i \leq n} \|x_i\|)$$

Apply recursively: $\|\theta_1\|^2 \leq \|\theta_0\|^2 + R^2 = 0 + R^2$

$$\|\theta_2\|^2 \leq \|\theta_1^2\| + R^2 \leq 2R^2$$

$$\vdots$$

$$\|\theta_k\|^2 \leq kR^2 \qquad Q.E.D.$$

(e). from question (4.b): $\theta_k^T \theta^* \geq k\rho$

from (4.d): $\|\theta_k\|^2 \leq kR^2$

$$\Rightarrow \frac{\theta_k^T \theta^*}{\|\theta_k\|} \geq \frac{k\rho}{\sqrt{kR^2}} = \frac{k\rho}{\sqrt{k}R} = \sqrt{k} \cdot \frac{\rho}{R} \qquad Q.E.D.$$

recall the Cauchy-Schwarz inequality:

$$\theta_k^T \theta^* \leq \|\theta_k\| \cdot \|\theta^*\| \qquad \Rightarrow \frac{\theta_k^T \theta^*}{\|\theta_k\| \|\theta^*\|} \leq 1$$

$$\Rightarrow \frac{\theta_k^T \theta^*}{\|\theta_k\|} \leq \|\theta^*\|$$

So, $\sqrt{k} \frac{\rho}{R} \leq \frac{\theta_k^T \theta^*}{\|\theta_k\|} \leq \|\theta^*\|$

Therefore, $\sqrt{k} \leq \|\theta^*\| \frac{R}{\rho}$
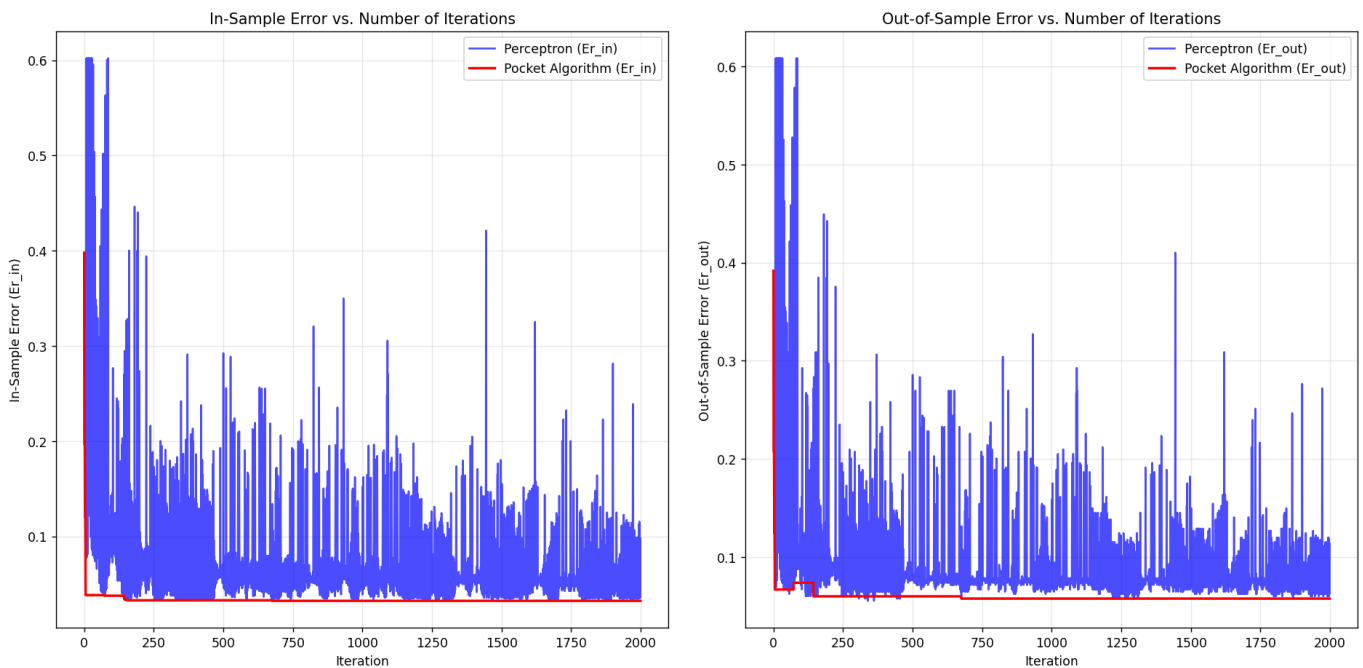
$$\Rightarrow k \leq \frac{R^2}{\rho^2} \|\theta^*\|^2$$

So, $\bar{k} \leq \frac{R^2 \|\theta^*\|^2}{\rho^2} \qquad Q.E.D.$

Problem 5

(5.1) code: p5.py

(5.2)



Discussion of the results:

The Perceptron's errors go up and down a lot during training. This is because the data (digits "1" and "6") cannot be perfectly separated by a straight line. The Pocket Algorithm is more stable. It remembers the best solution it has found so far. Even if the Perceptron's current weights get worse, the Pocket keeps the best ones. So, its training error only gets better or stays the same.

(5.3)