

Hand-written Part

学号: t11902210 姓名: 张凡 (交换生)

code question report 在后面

Problem 1

Consider a binary classification data set $\{(x_n, y_n)\}_{n=1}^N$ with $x_n \in \mathbb{R}^d$ and $y_n \in \{-1, 1\}$. For any weight vector w within a linear model, define an error function

$$\text{err}(w^T x, y) = (\max(1 - y w^T x, 0))^2$$

That is, $E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (\max(1 - y_n w^T x_n, 0))^2$

Running gradient descent to optimize $E_{in}(w)$ requires calculating its gradient direction $\nabla E_{in}(w)$ (and then move opposite to that direction). What is $\nabla E_{in}(w)$?

$$\nabla E_{in}(w) = \frac{\partial E_{in}(w)}{\partial w} = \frac{1}{N} \sum_{n=1}^N \max(1 - y_n w^T x_n, 0) \cdot \frac{\partial f(w)}{\partial w}$$

$$f(w) = [\max(1 - y_n w^T x_n, 0)]^2 \\ = \begin{cases} (1 - y_n w^T x_n)^2 & (y_n w^T x_n \leq 1) \\ 0 & (y_n w^T x_n > 1) \end{cases}$$

$$\frac{\partial f(w)}{\partial w} = \begin{cases} -2 y_n x_n (1 - y_n w^T x_n) & 0 \quad (y_n w^T x_n > 1) \\ (y_n w^T x_n \leq 1) \\ 0 & (y_n w^T x_n > 1) \end{cases}$$

$$\nabla E_{in}(w) = -2 y_n x_n (\max(1 - y_n w^T x_n, 0))$$

$$\nabla E_{in}(w) = \begin{cases} -\frac{2}{N} \sum_{n=1}^N y_n x_n (1 - y_n w^T x_n) & (y_n w^T x_n \leq 1) \\ 0 & (y_n w^T x_n > 1) \end{cases}$$

$$= -\frac{2}{N} \sum_{n=1}^N y_n x_n (\max(1 - y_n w^T x_n, 0))$$

Problem 2.

Consider a process that generates d -dimensional vectors x_1, x_2, \dots, x_N independently from a multivariate Gaussian distribution $\mathcal{N}(\mu, I)$, where $\mu \in \mathbb{R}^d$ is an unknown parameter vector and $I \in \mathbb{R}^{d \times d}$ is an identity matrix. The maximum likelihood estimate of μ is

$$\mu^* = \operatorname{argmax}_{\mu \in \mathbb{R}^d} \prod_{n=1}^N p_{\mu}(x_n),$$

where p_{μ} is the probability density function of $\mathcal{N}(\mu, I)$. Prove that

$$\mu^* = \frac{1}{N} \sum_{n=1}^N x_n.$$

Gaussian distribution (d)

$$N(\mu, \Sigma) \text{ definition} \Rightarrow p_{\mu}(x_n) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x_n - \mu)^T \Sigma^{-1}(x_n - \mu)}$$

maximum likelihood:

$$\left(x \sim N(\mu, \Sigma) \text{ for } f(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)} \right)$$

$$\mu^* = \arg \max \prod_{n=1}^N p_{\mu}(x_n) \quad \mu_0 = \arg \prod_{n=1}^N p_{\mu}(x_n)$$

turn "multiply" to "add" via "log"

$$\log \mu_0 = \sum_{n=1}^N \log p_{\mu}(x_n) \rightarrow \text{find max}$$

$$\log p_{\mu}(x_n) = \log \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} + \left(-\frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) \right)$$

$$= -\frac{1}{2} \log (2\pi)^d |\Sigma| - \frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

$$= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

\rightarrow find max

上述问题等价于 $\rightarrow -\log p_{\mu}(x_n) \rightarrow \min.$

$$\rightarrow P_{\mu}(x_n) = \frac{1}{2} (\log 2\pi + \frac{1}{2} \log |\Sigma| + \frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu))$$

$$\sum \log p_{\mu}(x_n) = \frac{dN}{2} \log 2\pi + \frac{N}{2} \log |\Sigma| + \frac{1}{2} \sum (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

固定值

只需求 $\rightarrow \sum (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$ 取最小, (Σ 为单位阵)

$$\begin{aligned} g(\mu) &= \sum (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) = \sum \cancel{\frac{1}{2} \|x_n - \mu\|^2} = \sum \|x_n - \mu\|^2 \\ &\stackrel{\partial \sum \|x_n - \mu\|^2}{\cancel{\partial \mu}} = \cancel{\sum 2(x_n - \mu)} = 0 \rightarrow \cancel{x_n = b}. \end{aligned}$$

$$\frac{\partial g(\mu)}{\partial \mu} = -2 \sum (x_n - \mu) = 0 \quad \begin{aligned} &\sum_{n=1}^N (x_n - \mu) \\ &= \sum_{n=1}^N x_n - N\mu = 0 \end{aligned}$$

$$\text{if } \mu = \frac{1}{N} \sum_{n=1}^N x_n \Rightarrow \mu^* = \frac{1}{N} \sum_{n=1}^N x_n$$

Problem 3.

A classic binary classification data set that cannot be separated by any line is called the XOR data set, with

$$\overbrace{x = [x_1, x_2]}^{x_1 = [+1, +1]} \quad y$$

$$\begin{array}{ll} x_1 = [+1, +1] & y_1 = -1 \\ x_2 = [-1, +1] & y_2 = +1 \\ x_3 = [-1, -1] & y_3 = +1 \\ x_4 = [+1, -1] & y_4 = -1 \end{array}$$

You can see why it is called XOR by

interpreting +1 as a boolean value value of

True and -1 as FALSE.

Consider a second-order feature transform $\tilde{x}_2(x) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$ that converts the data set to

$z = \Phi_2(x)$	y
$z_1 = \Phi_2(x_1)$	$y_1 = -1$
$z_2 = \Phi_2(x_2)$	$y_2 = +1$
$z_3 = \Phi_2(x_3)$	$y_3 = -1$
$z_4 = \Phi_2(x_4)$	$y_4 = +1$

Show a perception \tilde{w} in the z -space that separates the data. That is $y_n = \text{sign}(\tilde{w}^T z_n)$ for $n=1, 2, 3, 4$. Then, plot the classification boundary $\tilde{w}^T \Phi_2(x)$ in the x -space. Your boundary should look like a quadratic curve that classifies x_1, x_2, x_3, x_4 perfectly.

解: 已知 x_1, \dots, x_4 代入 $z_i = z$, (特征转换)

$$\begin{cases} z_1 = \Phi_2([+1, +1]) = (1, 1, 1, 1, 1, 1)^T \Rightarrow y_1 = -1 \\ z_2 = \Phi_2([-1, +1]) = (1, -1, 1, 1, -1, 1)^T \Rightarrow y_2 = +1 \\ z_3 = \Phi_2([-1, -1]) = (1, -1, -1, 1, 1, 1)^T \Rightarrow y_3 = -1 \\ z_4 = \Phi_2([-1, +1, -1]) = (1, 1, -1, 1, -1, 1)^T \Rightarrow y_4 = +1 \end{cases}$$

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

先设 $\tilde{w} = (0, 0, 0, 0, 0, 0)^T$

① $\tilde{w}^T z_1 = 0 \neq (-1) \Rightarrow \tilde{w} = \tilde{w} + y_1 z_1 = (-1, -1, -1, -1, -1, -1)^T$

② $\tilde{w}^T z_2 = -1 + 1 - 1 - 1 + 1 - 1 = -2 \Rightarrow \tilde{w} = \tilde{w} + y_2 z_2 = (0, -2, 0, 0, -2, 0)^T$

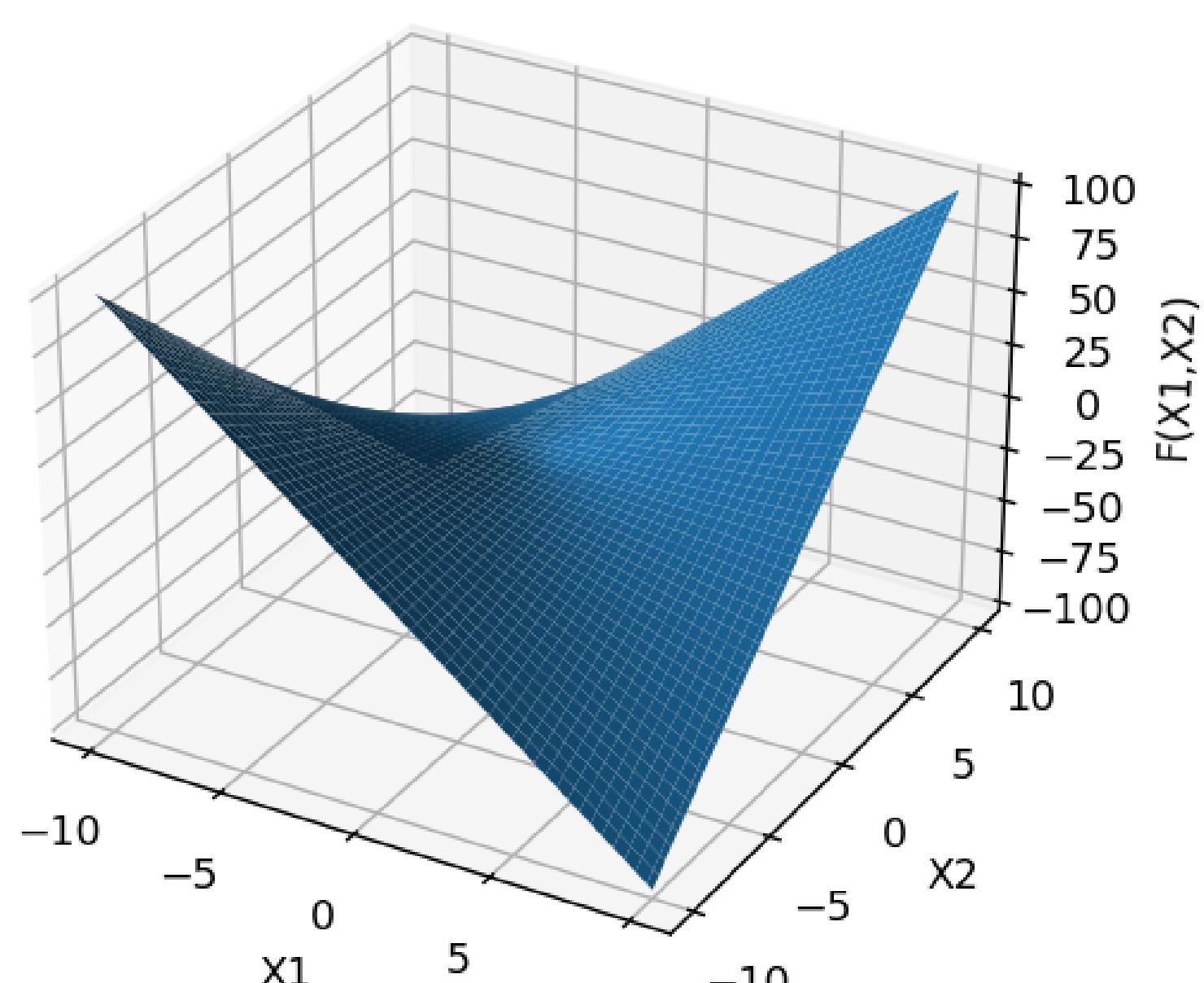
③ $\tilde{w}^T z_3 = 0 + 0 + 0 + 0 + 0 = 0 \Rightarrow \tilde{w} = \tilde{w} + y_3 z_3 = (-1, -1, 1, -1, 1, -1)^T$

④ $\tilde{w}^T z_4 = -1 - 3 - 1 - 1 + 4 - 1 = -3 \Rightarrow \tilde{w} = \tilde{w} + y_4 z_4 = (0, 0, 0, 0, 0, -3)^T$

权值向量 $\tilde{w} = (0, 0, 0, 0, 0, -3)^T$

$\tilde{w}^T \Phi_2(x) = -3x_1 - 3x_2 = 0 \Rightarrow x_1 = x_2 = 0$

$x_1 = 0 \Rightarrow x_2 = 0$



The left side shows the final image I drew using Python and related drawing libraries.

The mathematical expression for the image is below.

$$f(x_1, x_2) = x_1 * x_2$$

Problem. 4

Consider building binary classification with AdaBoost algorithm, a weak classifier $g_t(x)$ is trained on a data set $\{(x_n, y_n)\}_{n=1}^N$, with weights $\{w_n\}_{n=1}^N$, at the time step t .

The error rate is defined as $q_t = \frac{\sum_{n=1}^N w_n \cdot \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n}$,

if $g_t(x_n) \neq y_n$ and $\delta(g_t(x_n), y_n) = 0$ otherwise. For the next time step, the data set is reweighted to emphasize on misclassified samples through the following rules

$$w_n^{t+1} = \begin{cases} w_n^t \cdot dt & \text{if } g_t(x_n) \neq y_n \\ w_n^t / dt & \text{if } g_t(x_n) = y_n \end{cases}$$

Show that $dt = \sqrt{1 - \epsilon_t} / \epsilon_t$ can degrade the previous classifier $g_t(x)$ and make its error rate become 0.5 with new weights $\{w_n^{t+1}\}_{n=1}^N$:

$$\frac{\sum_{n=1}^N w_n^{t+1} \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^t} = 0.5$$

if:

$$\epsilon_t = \frac{\sum_{n=1}^N w_n^t \cdot \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^t}$$

$$w_n^{t+1} = \begin{cases} w_n^t \cdot dt & (g_t(x_n) \neq y_n) \\ w_n^t / dt & (g_t(x_n) = y_n) \end{cases}$$

$$dt = \sqrt{(1 - \epsilon_t) / \epsilon_t} = \sqrt{\frac{1}{\epsilon_t} - 1}$$

$$w_n^{t+1} = w^t \cdot \epsilon$$

$$= \sqrt{\frac{\sum_{n=1}^N w_n^t}{\sum_{n=1}^N w_n^t \delta(g_t(x_n), y_n)}} - 1$$

$$= \sqrt{\frac{\sum_{n=1}^N w_n^t (1 - \delta(g_t(x_n), y_n))}{\sum_{n=1}^N w_n^t \delta(g_t(x_n), y_n)}} = \cancel{\sqrt{\dots}}$$

$$\frac{\sum_{n=1}^N w_n^{t+1} \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^{t+1}} = \frac{\sum_{n=1}^N w_n^t \cdot dt \cdot \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^{t+1}} = \frac{\sum_{n=1}^N w_n^t \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^{t+1}} \cdot dt$$

$$= \frac{\sum_{n=1}^N w_n^t \delta(g_t(x_n), y_n)}{\sum_{n=1}^N w_n^t} \cdot \frac{\sum_{n=1}^N w_n^t}{\sum_{n=1}^{t+1} w_n^t} \cdot dt \cdot \epsilon_t \cdot \frac{\sum_{n=1}^N w_n^t}{\sum_{n=1}^{t+1} w_n^t}$$

$$= \epsilon_t \cdot \frac{\sum_{n=1}^N w_n^t}{\sum_{n=1}^{t+1} w_n^t} \cdot \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \sqrt{1 - \epsilon_t} \cdot \frac{\sum_{n=1}^N w_n^t}{\sum_{n=1}^{t+1} w_n^t}$$

$$\sum_{n=1}^N w_n^{t+1} dt = \sum_{n=1}^N \left(w_n^t \cdot dt \cdot \delta(g_t(x_n), y_n) + \frac{w_n^t dt}{dt} (1 - \delta(g_t(x_n), y_n)) \right)$$

$$= \sum_{n=1}^N w_n^t \cdot \frac{(dt^2 - 1) \delta(g_t(x_n), y_n)}{dt} + \sum_{n=1}^N w_n^t \cdot$$

$$= (\cancel{dt^2 - 1}) \cdot \epsilon_t \cdot \frac{\sum_{n=1}^N w_n^t}{\cancel{dt}} + \frac{\sum_{n=1}^N w_n^t}{dt}$$

$$\frac{\sum_{n=1}^N w_n^{t+1}}{\sum_{n=1}^N w_n^t} = \frac{(\cancel{dt^2 - 1}) \cdot \epsilon_t + 1}{dt} = \frac{1 + \frac{1 - \epsilon_t}{\epsilon_t} \cdot \epsilon_t}{dt} = \frac{2 - \epsilon_t}{\sqrt{1 - \epsilon_t}}$$

$$\cancel{dt} = \sqrt{1 - \epsilon_t} \cdot \frac{dt \cdot \epsilon_t}{\cancel{dt^2 - 1} \cdot \epsilon_t + 1} = \frac{dt \cdot \epsilon_t}{(\cancel{dt^2 - 1}) \cdot \epsilon_t + 1} = \frac{dt \cdot \epsilon_t}{dt^2 - \cancel{dt} + \cancel{dt} + 1} = \frac{dt \cdot \epsilon_t}{dt^2 - dt + dt + 1} = \frac{dt \cdot \epsilon_t}{dt^2 - dt + dt} = \frac{dt \cdot \epsilon_t}{dt^2} = \frac{\epsilon_t}{dt} = \frac{1 - \epsilon_t}{2 - 2\epsilon_t}$$

$$= \frac{1 - \epsilon_t}{2 - 2\epsilon_t} = \frac{1}{2} \quad \star$$

Report
见后