

Linear Regression

1. ANS [d]

Consider a noisy target $y = \mathbf{w}_f^T \mathbf{x} + \epsilon$, where $\mathbf{x} \in \mathbb{R}^{d+1}$ (including the added coordinate $x_0 = 1$), $y \in \mathbb{R}$, $\mathbf{w}_f \in \mathbb{R}^{d+1}$ is an unknown vector, and ϵ is an i.i.d. noise term with zero mean and σ^2 variance. Assume that we run linear regression on a training data set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ generated i.i.d. from some $P(\mathbf{x})$ and the noise process above, and obtain the weight vector \mathbf{w}_{lin} . As briefly discussed in Lecture 5, it can be shown that the expected in-sample error $E_{\text{in}}(\mathbf{w}_{\text{lin}})$ with respect to \mathcal{D} is given by:

$$\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right).$$

For $\sigma = 0.1$ and $d = 19$, what is the smallest number of examples N such that $\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})]$ is no less than 0.005? Choose the correct answer; explain your answer.

- [a] 25
- [b] 30
- [c] 35
- [d] 40**
- [e] 45

$$\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right)$$

$$\sigma = 0.1 \quad d = 19$$

$$\Rightarrow \sigma^2 \left(1 - \frac{d+1}{N}\right) \leq 0.005$$

$$0.01 \left(1 - \frac{20}{N}\right) \leq 0.005$$

$$0.01 - \frac{0.2}{N} \leq 0.005$$

$$-0.2 \leq -0.005N$$

$$N \geq \frac{200}{5} = 40$$

~~不符合~~

smallest number of $N \geq 40$

~~这 d~~

2. ANS [c]

Consider the target function $f(x) = x^2$. Sample x uniformly from $[0, 1]$, and use all linear hypotheses $h(x) = w_0 + w_1 \cdot x$ to approximate the target function with respect to the squared error. What are the weights (w_0^*, w_1^*) of the optimal hypothesis? Choose the correct answer; explain your answer.

- [a] $(0, 1)$
- [b] $(\frac{1}{2}, \frac{1}{2})$
- [c] $(-\frac{1}{6}, 1)$**
- [d] $(-\frac{1}{4}, \frac{1}{4})$
- [e] $(\frac{1}{3}, 0)$

(Hint: The optimal hypothesis g^* must reach the minimum $E_{\text{out}}(g^*)$.)

$$f(x) = x^2$$

$$h(x) = w_0 + w_1 x$$

$$E_{\text{out}}(w) = \mathbb{E}_{(x,y) \sim P} (h(x) - f(x))^2$$

$$= \int_0^1 (w_0 + w_1 x - x^2)^2 dx$$

$$\int (w_0 + w_1 x - x^2)^2 dx$$

$$= \int (x^4 - w_1 x^3 - w_0 x^2 - w_1 x^3 + w_1^2 x^2 + w_0 w_1 x \\ - w_0 x^2 + w_0 w_1 x + w_0^2) dx$$

$$= \int (x^4 - 2w_1 x^3 + (w_1^2 - 2w_0) x^2 + 2w_0 w_1 x + w_0^2) dx$$

$$= \frac{1}{5} x^5 - \frac{2w_1}{4} x^4 + \frac{(w_1^2 - 2w_0)}{3} x^3 + \frac{2w_0 w_1}{2} x^2 + w_0^2 x$$

令各面積函数 $F(x)$

$$\begin{aligned}
 \int_0^1 (w_0 + w_1 x - x^2)^2 dx &= F(1) - F(0) = F(x) \Big|_0^1 \\
 &= \frac{1}{5} - \frac{w_1}{2} + \frac{(w_1^2 - 2w_0)}{3} + w_0 w_1 + w_0^2 \\
 &= \frac{1}{3} w_1^2 - \frac{1}{2} w_1 - \frac{2}{3} w_0 + w_0 w_1 + w_0^2 + \frac{1}{5}
 \end{aligned}$$

對 w_0 微分

$$\boxed{\frac{2}{3} + w_1 + 2w_0 = 0} \quad \left[\begin{array}{l} -\frac{2}{3} + w_1 + 2w_0 = 0 \\ \frac{2}{3} w_1 - \frac{1}{2} + w_0 = 0 \end{array} \right]$$

對 w_1 微分

$$\boxed{\frac{2}{3} w_1 - \frac{1}{2} + w_0 = 0} \quad \left[\begin{array}{l} -\frac{2}{3} + w_1 + 2w_0 = 0 \\ -1 + \frac{4}{3} w_1 + 2w_0 = 0 \end{array} \right]$$

$$\Rightarrow w_1 = 1, \quad w_0 = -\frac{1}{6}$$

$$\text{得最佳 } (w_0^*, w_1^*) = \left(-\frac{1}{6}, 1\right)$$

選 C

3. ANS [e]

Following the previous problem, assume that we sample two examples x_1 and x_2 uniformly from $[0, 1]$ to form the training set $\mathcal{D} = \{(x_1, f(x_1)), (x_2, f(x_2))\}$, and use linear regression to get g for approximating the target function with respect to the squared error. You can neglect the degenerate cases where x_1 and x_2 are the same. What is $\mathbb{E}_{\mathcal{D}}(|E_{\text{in}}(g) - E_{\text{out}}(g)|)$? Choose the correct answer; explain your answer.

- [a] $\frac{1}{60}$
- [b] $\frac{4}{15}$
- [c] $\frac{3}{20}$
- [d] $\frac{1}{25}$
- [e] $\frac{1}{30}$**



$$\text{求 } \mathbb{E}_{\mathcal{D}}(|E_{\text{in}}(g) - E_{\text{out}}(g)|)$$

由於 \mathcal{D} 只有 2 個，故 $E_{\text{in}}(g) = 0$

$$\text{剩下求 } \mathbb{E}_{\mathcal{D}}(E_{\text{out}}(g))$$

由第二題得出

$$\begin{aligned} g(x) &= \int_0^1 (w_0 + w_1 x - x^2)^2 dx \\ &= \frac{1}{3} w_1^2 - \frac{1}{2} w_1 - \frac{2}{3} w_0 + w_0 w_1 + w_0^2 + \frac{1}{5} \end{aligned}$$

用 x_1, x_2 替換 w_0, w_1

$$\text{由於 } \mathcal{D} = \{(x_1, f(x_1)), (x_2, f(x_2))\}$$

且 第二題 設 $f(x) = x^2$, $h(x) = w_0 + w_1 \cdot x$

$$\begin{aligned} \therefore \begin{cases} w_0 + w_1 x_1 = x_1^2 \\ w_0 + w_1 x_2 = x_2^2 \end{cases} \\ w_1 (x_1 - x_2) = x_1^2 - x_2^2 \\ \underline{w_1} = \frac{x_1^2 - x_2^2}{x_1 - x_2} = \underline{x_1 + x_2} \\ \underline{w_0} = x_1^2 - (x_1 + x_2)x_1 \\ = \underline{-x_1 x_2} \end{aligned}$$

$$T^* \lambda g(x) = \frac{1}{3} w_1^2 - \frac{1}{2} w_1 - \frac{2}{3} w_0 + w_0 w_1 + w_0^2 + \frac{1}{5}$$

$$\begin{aligned} g(x) &= \frac{1}{3}(x_1+x_2)^2 - \frac{1}{2}(x_1+x_2) - \frac{2}{3}(-x_1x_2) + (-x_1x_2)(x_1+x_2) \\ &\quad + x_1^2x_2^2 + \frac{1}{5} \\ &= \frac{1}{3}(x_1^2 + 2x_1x_2 + x_2^2) - \frac{1}{2}x_1 - \frac{1}{2}x_2 + \frac{2}{3}x_1x_2 - x_1^2x_2 - x_1x_2^2 \\ &\quad + x_1^2x_2^2 + \frac{1}{5} \\ &= \frac{1}{3}x_1^2 + \frac{2}{3}\underline{x_1x_2} + \frac{1}{3}x_2^2 - \frac{1}{2}x_1 - \frac{1}{2}x_2 + \frac{2}{3}\underline{x_1x_2} - x_1^2x_2 - x_1x_2^2 \\ &\quad + x_1^2x_2^2 + \frac{1}{5} \\ &= \frac{1}{3}x_1^2 + \frac{1}{3}x_2^2 - \frac{1}{2}x_1 - \frac{1}{2}x_2 + \frac{4}{3}x_1x_2 - x_1^2x_2 - x_1x_2^2 \\ &\quad + x_1^2x_2^2 + \frac{1}{5} \end{aligned}$$

$$E_D(\pi_{\text{out}}(q))$$

$$\begin{aligned} &= \int_0^1 \int_0^1 \left[\frac{1}{3}x_1^2 + \frac{1}{3}x_2^2 - \frac{1}{2}x_1 - \frac{1}{2}x_2 + \frac{4}{3}x_1x_2 - x_1^2x_2 - x_1x_2^2 + x_1^2x_2^2 + \frac{1}{5} \right] dx_1 dx_2 \\ &= \int_0^1 \left[\frac{1}{9}x_1^3 + \frac{1}{3}x_1x_2^2 - \frac{1}{4}x_1^2 - \frac{1}{2}x_1x_2 + \frac{2}{3}x_1^2x_2 - \frac{1}{3}x_1^3x_2 \right. \\ &\quad \left. - \frac{1}{2}x_1^2x_2^2 + \frac{1}{3}x_1^3x_2^2 + \frac{1}{5}x_1 \right] \Big|_0^1 dx_2 \\ &= \int_0^1 \left[\frac{1}{9} + \frac{1}{3}x_2^2 - \frac{1}{4} - \frac{1}{2}x_2 + \frac{2}{3}x_2 - \frac{1}{3}x_2^2 - \frac{1}{2}x_2^2 + \frac{1}{3}x_2^2 + \frac{1}{5} \right] dx_2 \\ &= \left[\frac{1}{9}x_2 + \frac{1}{9}x_2^3 - \frac{1}{4}x_2 - \frac{1}{4}x_2^2 + \frac{1}{3}x_2^2 - \frac{1}{6}x_2^2 - \frac{1}{6}x_2^3 + \frac{1}{9}x_2^3 + \frac{1}{5}x_2 \right] \Big|_0^1 \\ &= \frac{1}{9} + \frac{1}{9} - \frac{1}{4} - \frac{1}{4} + \frac{1}{3} - \frac{1}{6} - \frac{1}{6} + \frac{1}{9} + \frac{1}{5} \\ &= \frac{3}{9} - \frac{2}{4} + \frac{1}{3} - \frac{2}{6} + \frac{1}{5} = \frac{1}{3} - \frac{1}{2} + \frac{1}{3} - \frac{1}{3} + \frac{1}{5} = \frac{10 - 15 + b}{30} \\ &= \frac{1}{30} \end{aligned}$$

$$\boxed{E_D(\pi_{\text{out}}(q)) = \frac{1}{30} \text{ 得證}}$$

Cross-Entropy Error

4. ANS [b]

In class, we introduced our version of the cross-entropy error function

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N -\ln \theta(y_n \mathbf{w}^T \mathbf{x}_n).$$

based on the definition of $y_n \in \{-1, +1\}$. If we transform y_n to $y'_n \in \{0, 1\}$ by $y'_n = \frac{y_n+1}{2}$, which of the following error function is equivalent to E_{in} above? Choose the correct answer; explain your answer.

- [a] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(\mathbf{w}^T \mathbf{x}_n) + (1 - y'_n) \ln(\theta(-\mathbf{w}^T \mathbf{x}_n)))$
- [b] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(-\mathbf{w}^T \mathbf{x}_n) + (1 - y'_n) \ln(\theta(\mathbf{w}^T \mathbf{x}_n)))$
- [c] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(\mathbf{w}^T \mathbf{x}_n) - (1 - y'_n) \ln(\theta(-\mathbf{w}^T \mathbf{x}_n)))$
- [d] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(-\mathbf{w}^T \mathbf{x}_n) - (1 - y'_n) \ln(\theta(\mathbf{w}^T \mathbf{x}_n)))$
- [e] none of the other choices

Linear Models Logistic Regression Error

Likelihood

target function $f(\mathbf{x}) = P(+1|\mathbf{x})$ \Leftrightarrow $P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1 \\ 1 - f(\mathbf{x}) & \text{for } y = -1 \end{cases}$

consider $\mathcal{D} = \{(\mathbf{x}_1, \circ), (\mathbf{x}_2, \times), \dots, (\mathbf{x}_N, \times)\}$

probability that f generates \mathcal{D} likelihood that h generates \mathcal{D}

$P(\mathbf{x}_1)f(\mathbf{x}_1) \times$	$P(\mathbf{x}_1)h(\mathbf{x}_1) \times$
$P(\mathbf{x}_2)(1 - f(\mathbf{x}_2)) \times$	$P(\mathbf{x}_2)(1 - h(\mathbf{x}_2)) \times$
\dots	\dots
$P(\mathbf{x}_N)(1 - f(\mathbf{x}_N))$	$P(\mathbf{x}_N)(1 - h(\mathbf{x}_N))$

if $h \approx f$, then likelihood(h) \approx probability using f

- probability using f usually **large**

Hsuan-Tien Lin (NTU CSIE) Machine Learning 25/52

Linear Models Logistic Regression Error

Likelihood of Logistic Hypothesis

likelihood(h) \approx (probability using f) \approx **large**

$$g = \underset{h}{\operatorname{argmax}} \text{ likelihood}(h)$$

when logistic: $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$ $1 - h(\mathbf{x}) = h(-\mathbf{x})$

likelihood(h) = $P(\mathbf{x}_1)h(\mathbf{x}_1) \times P(\mathbf{x}_2)(1 - h(\mathbf{x}_2)) \times \dots \times P(\mathbf{x}_N)(1 - h(\mathbf{x}_N))$

likelihood(logistic h) $\propto \prod_{n=1}^N h(y_n \mathbf{x}_n)$

Hsuan-Tien Lin (NTU CSIE) Machine Learning 26/52

根據投影片與題意

$$\begin{aligned} \text{可假設 } f(x) &= P(y_n = 1 | x) \\ &= P(y_n' = 1 | x) \text{ 为 target function} \end{aligned}$$

且此題只有 2 種情況: $y_n' = 1$ or 0

$$\begin{aligned} \Rightarrow P(y_n' | x) &= \begin{cases} f(x) & \text{for } y_n' = 1 \\ 1-f(x) & \text{for } y_n' = 0 \end{cases} \\ &= \underbrace{f(x)^{y_n'} (1-f(x))^{1-y_n'}}_{\text{if } h \approx f, \text{ then likelihood}(h) \approx \text{probability using } f} \end{aligned}$$

且 probability using f is large

故 likelihood(h) 也大

\Rightarrow likelihood(h) \approx (probability using f) \approx large

$$g = \underset{h}{\operatorname{argmax}} \text{ likelihood}(h)$$

但 $\text{likelihood}(h) = L(h)$

且取 log 時 $h(x) = \Theta(w^T x)$, $1-h(x) = \Theta(-w^T x)$

$$L(\log h) = P(x_1) h(x_1) \times P(x_2) h(x_2) \times \dots \times P(x_N) (1-h(x_N))$$

$$\Rightarrow \underset{h}{\operatorname{argmax}} L(\log h) \propto \prod_{n=1}^N \underbrace{h(x)^{y_n'} (1-h(x))^{1-y_n'}}_{\text{if } h \approx f, \text{ then likelihood}(h) \approx \text{probability using } f}$$

$$\max_w L(w) \propto \prod_{n=1}^N \underbrace{\Theta(w^T x)^{y_n'} (1-\Theta(w^T x))^{1-y_n'}}_{\text{if } h \approx f, \text{ then likelihood}(h) \approx \text{probability using } f}$$

$$\propto \prod_{n=1}^N \Theta(w^T x)^{y_n'} \underbrace{\Theta(-w^T x)^{1-y_n'}}_{\text{if } h \approx f, \text{ then likelihood}(h) \approx \text{probability using } f}$$

$$\text{由於 } y_n' = 0 \text{ or } 1 \Rightarrow \propto \prod_{n=1}^N (\Theta(z y_n - 1) w^T x)$$

\Rightarrow 求 $\max_w L(w)$ 相當於求 $\max_w \prod_{n=1}^N (\theta(z_n y_n - 1) w^T x)$
 且相當於求 $\max_w \prod_{n=1}^N \ln(\theta(z_n y_n - 1) w^T x)$ 取 ln 大小關係 不變

$$\max_w \prod_{n=1}^N \ln(\theta(z_n y_n - 1) w^T x) = \min_w \frac{1}{N} \sum_{n=1}^N -\ln(\theta(z_n y_n - 1) w^T x)$$

由於 $\theta(-s) = 1 - \theta(s)$

$$\min_w -\ln \theta(-s) \Leftrightarrow \min_w -\ln(1 - \theta(s))$$

$$\Leftrightarrow \max_w \ln(1 - \theta(s))$$

$$\Leftrightarrow \min_w \ln(\theta(s) - 1)$$

$$\Leftrightarrow \min_w \ln \theta(s)$$

$$\Leftrightarrow \min_w \frac{1}{N} \sum_{n=1}^N -\ln(\theta(z_n y_n - 1) w^T x) \Leftrightarrow \min_w \frac{1}{N} \sum_{n=1}^N \ln(\theta(1 - z_n y_n) w^T x)$$

$$\min_w \frac{1}{N} \sum_{n=1}^N \ln(\theta(1 - z_n y_n) w^T x) = \min_w \frac{1}{N} \sum_{n=1}^N \ln(\theta(1 - y_n' - y_n) w^T x)$$

$$= \min_w \frac{1}{N} \sum_{n=1}^N \ln \left[\theta[(1 - y_n') w^T x - y_n' w^T x] \right]$$

假設 $E_{in_new}(w)$

$$E_{in_new}(w) = \frac{1}{N} \sum_{n=1}^N \ln \left[\theta[(1 - y_n') w^T x - y_n' w^T x] \right]$$

$$y_n' = 0, E_{in_new}(w) = \frac{1}{N} \sum_{n=1}^N \ln \theta(w^T x)$$

$$y_n' = 1, E_{in_new}(w) = \frac{1}{N} \sum_{n=1}^N \ln \theta(-w^T x)$$

$$\Leftrightarrow E_{in}'(w) = \frac{1}{N} \sum_{n=1}^N \left(\underbrace{y_n' \ln \theta(-w^T x)}_{得證} + \underbrace{(1 - y_n') \ln \theta(w^T x)}_b \right)$$

5. ANS [e]

Consider a coin with an unknown head probability μ . Independently flip this coin N times to get y_1, y_2, \dots, y_N , where $y_n = 1$ if the n -th flipping results in head, and 0 otherwise. Define $\nu = \frac{1}{N} \sum_{n=1}^N y_n$. How many of the following statements about ν are true? Choose the correct answer; explain your answer by briefly illustrating why those statements are true.

(Note: μ is similar to the role of the “target function” and $\hat{\mu}$ is similar to the role of the “hypothesis” in our machine learning framework.)

- [a] 0
- [b] 1
- [c] 2
- [d] 3
- [e] 4**

全部都對，選 e。

- With probability more than $1 - \delta$,

$$\mu \leq \nu + \sqrt{\frac{1}{2N} \ln \frac{2}{\delta}}$$

for all $N \in \mathbb{N}$ and $0 < \delta < 1$.

- in big sample (N large), ν is probably close to μ (within ϵ)

$$\mathbb{P}[|\nu - \mu| > \epsilon] \leq 2 \exp(-2\epsilon^2 N)$$

- called **Hoeffding's Inequality**, for marbles, coin, polling, ...

⊕ Hoeffding's Inequality

$$\Rightarrow \mathbb{P}[|\nu - \mu| > \epsilon] \leq \underbrace{2 \exp(-2\epsilon^2 N)}_{\delta}$$

With probability $\geq 1 - \delta$ (Good: $|\nu - \mu| \leq \epsilon$)

$$\therefore \delta = 2 \exp(-2\epsilon^2 N)$$

$$\frac{\delta}{2} = \exp(-2\epsilon^2 N)$$

取 ln

$$\ln\left(\frac{\delta}{2}\right) = -2\epsilon^2 N$$

$$\epsilon^2 = -\frac{1}{2N} \ln\left(\frac{\delta}{2}\right)$$

$$= \frac{1}{2N} \ln\left(\frac{2}{\delta}\right)$$

$$\epsilon = \sqrt{\frac{1}{2N} \ln\left(\frac{2}{\delta}\right)}$$

$$\Rightarrow |\nu - \mu| \leq \epsilon = \sqrt{\frac{1}{2N} \ln\left(\frac{2}{\delta}\right)}$$

$$\nu - \sqrt{\frac{1}{2N} \ln\left(\frac{2}{\delta}\right)} \leq \mu \leq \nu + \sqrt{\frac{1}{2N} \ln\left(\frac{2}{\delta}\right)}$$

得證

- ν maximizes likelihood($\hat{\mu}$) over all $\hat{\mu} \in [0, 1]$.

Coin Flip MLE 伯努利分佈, $y_n = 0$ or 1

$$\text{likelihood}(\hat{\mu}) = L_{y_n}(\hat{\mu}) = P(y_n | \hat{\mu}) = \hat{\mu}^{y_n} (1 - \hat{\mu})^{1-y_n}$$

由於拋硬幣為 i.i.d., 故 likelihood = 所有可能性連乘

$$\text{likelihood}(\hat{\mu}) = L_{y_n}(\hat{\mu}) = \prod_{y_n \in y} \hat{\mu}^{y_n} (1 - \hat{\mu})^{1-y_n}$$

$$\text{EX: } y = \begin{matrix} h & t & t & t & h & h \\ 1 & 0 & 0 & 0 & 1 & 1 \end{matrix} \quad (\text{head, tail})$$

$$\text{則 } L_{y_n}(\hat{\mu}) = \hat{\mu}^h \cdot (1 - \hat{\mu})^{n-h}$$

$$\Leftrightarrow L_{y_n}(\hat{\mu}) = \hat{\mu}^h \cdot (1 - \hat{\mu})^{n-h} \quad \left(\begin{array}{l} n \text{ 次丟的總次數} \\ h \text{ 是 head 次數} \end{array} \right)$$

想找出 maximizes of $L_{y_n}(\hat{\mu})$

取 log-likelihood (" $\hat{\mu}$ that maximize $L_{y_n}(\hat{\mu})$)
 等價於 $\hat{\mu}$ that maximize log-Lyn($\hat{\mu}$)

$$\log-Lyn(\hat{\mu}) = h \cdot \log(\hat{\mu}) + (n-h) \cdot \log(1 - \hat{\mu})$$

對公偏微分

$$\Leftrightarrow \log-Lyn(\hat{\mu})' = \frac{h}{\hat{\mu}} - \frac{n-h}{1-\hat{\mu}} = 0$$

$$\Leftrightarrow \hat{\mu} = \frac{h}{n}, \text{ 而且 } v = \frac{1}{n} \sum_{n=1}^N y_n$$

也就是 $\frac{\text{head 次數}}{\text{總次數}}$

$$\Leftrightarrow v = \frac{h}{n} = \hat{\mu}$$

故 v 可以 maximizes likelihood($\hat{\mu}$)

- v minimizes the squared error

$$E^{\text{sqr}}(\hat{y}) = \frac{1}{N} \sum_{n=1}^N (\hat{y} - y_n)^2$$

over all $\hat{y} \in \mathbb{R}$.

要 minimizes the squared error

\hookrightarrow 微分

$$E^{\text{sqr}}(\hat{y}) = \frac{1}{N} \sum_{n=1}^N (\hat{y} - y_n)^2 = \frac{1}{N} \sum_{n=1}^N (\hat{y}^2 - 2\hat{y}y_n + y_n^2)$$

$$\frac{\partial E^{\text{sqr}}(\hat{y})}{\partial \hat{y}} = \frac{1}{N} \sum_{n=1}^N (2\hat{y} - 2y_n) = 0$$

$$\hookrightarrow \frac{1}{N} \sum_{n=1}^N (\hat{y} - y_n) = 0$$

$\Rightarrow v = \bar{y}$

$$\begin{aligned} \nabla E^{\text{sqr}}(v) &= \frac{1}{N} \sum_{n=1}^N (v - y_n) = 0 \\ &= \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{N} \sum_{n=1}^N y_n - y_n \right) = 0 \end{aligned}$$

$$\hookrightarrow \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{N} \sum_{n=1}^N y_n \right) - \frac{1}{N} \sum_{n=1}^N y_n = 0$$

$$\frac{1}{N} \sum_{n=1}^N y_n - \frac{1}{N} \sum_{n=1}^N y_n = 0$$

成立，得證

v 以最小化 squared error

- When $0 < \nu < 1$, it minimizes the cross-entropy error (which is similar to the cross-entropy error for logistic regression)

$$E^{\text{ce}}(\hat{y}) = -\frac{1}{N} \sum_{n=1}^N \left(y_n \ln \hat{y} + (1-y_n) \ln(1-\hat{y}) \right)$$

over all $\hat{y} \in (0, 1)$.

minimizes the cross-entropy error

$$\hat{y} \in (0, 1)$$

$$E^{\text{ce}}(\hat{y}) = -\frac{1}{N} \sum_{n=1}^N (y_n \ln \hat{y} + (1-y_n) \ln(1-\hat{y}))$$

$$\frac{\partial E^{\text{ce}}(\hat{y})}{\partial \hat{y}} = -\frac{1}{N} \sum_{n=1}^N \left[y_n \cdot \frac{\partial \ln(\square)}{\partial \square} \cdot \frac{\partial \hat{y}}{\partial \hat{y}} + (1-y_n) \cdot \frac{\partial \ln(\blacksquare)}{\partial \blacksquare} \cdot \frac{\partial (1-\hat{y})}{\partial \hat{y}} \right]$$

$$= -\frac{1}{N} \sum_{n=1}^N \left[y_n \cdot \frac{1}{\hat{y}} \cdot 1 + (1-y_n) \cdot \frac{1}{1-\hat{y}} \cdot (-1) \right]$$

$$= -\frac{1}{N} \sum_{n=1}^N \left[\frac{y_n}{\hat{y}} - \frac{1-y_n}{1-\hat{y}} \right]$$

$$= -\frac{1}{N} \sum_{n=1}^N \frac{y_n}{\hat{y}} + \frac{1}{N} \sum_{n=1}^N \frac{1-y_n}{1-\hat{y}}$$

$$-\frac{1}{N} \sum_{n=1}^N \frac{y_n}{\hat{y}} + \frac{1}{N} \sum_{n=1}^N \frac{1-y_n}{1-\hat{y}} = 0$$

$$[\hat{y} \text{ 是常数}]$$

$$-\cancel{\frac{1}{N} \cdot \frac{1}{\hat{y}} \sum_{n=1}^N y_n} + \cancel{\frac{1}{N} \cdot \frac{1}{1-\hat{y}} \sum_{n=1}^N (1-y_n)} = 0$$

$$-\frac{1}{\hat{y}} \sum_{n=1}^N y_n + \frac{1}{1-\hat{y}} \sum_{n=1}^N (1-y_n) = 0$$

$$-(1-\hat{y}) \sum_{n=1}^N y_n + \underbrace{\hat{y} \sum_{n=1}^N (1-y_n)}_{= 0} = 0$$

$$-\sum_{n=1}^N y_n + \cancel{\hat{y} \sum_{n=1}^N y_n} + \hat{y} \sum_{n=1}^N 1 - \cancel{\hat{y} \sum_{n=1}^N y_n} = 0$$

$$-\sum_{n=1}^N y_n + \hat{y} N = 0$$

$$\hat{y} = \frac{1}{N} \sum_{n=1}^N y_n = \nu$$

成立，得证

ν 最小化 cross-entropy error

Stochastic Gradient Descent

6. ANS [a]

In the perceptron learning algorithm, we find one example $(\mathbf{x}_{n(t)}, y_{n(t)})$ that the current weight vector \mathbf{w}_t mis-classifies, and then update \mathbf{w}_t by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}.$$

The algorithm can be viewed as optimizing some $E_{in}(\mathbf{w})$ that is composed of one of the following point-wise error functions with stochastic gradient descent (neglecting any non-differentiable points of the error function). What is the error function? Choose the correct answer; explain your answer.

- [a] $\text{err}(\mathbf{w}, \mathbf{x}, y) = \max(0, -y\mathbf{w}^T \mathbf{x})$
- [b] $\text{err}(\mathbf{w}, \mathbf{x}, y) = -\max(0, -y\mathbf{w}^T \mathbf{x})$
- [c] $\text{err}(\mathbf{w}, \mathbf{x}, y) = \max(y\mathbf{w}^T \mathbf{x}, -y\mathbf{w}^T \mathbf{x})$
- [d] $\text{err}(\mathbf{w}, \mathbf{x}, y) = -\max(y\mathbf{w}^T \mathbf{x}, -y\mathbf{w}^T \mathbf{x})$
- [e] none of the other choices

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \underbrace{y_{n(t)} \mathbf{x}_{n(t)}}_{-\nabla \text{err}(\mathbf{w}, \mathbf{x}, y)}$$

[a] $\text{err}(\mathbf{w}, \mathbf{x}, y)$ 有二種情況

① $\text{err}(\mathbf{w}, \mathbf{x}, y) = 0$

$$\nabla \text{err}(\mathbf{w}, \mathbf{x}, y) = 0$$

$\Rightarrow \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + 0$ 代表沒有錯不需要更新的情形

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

$$\nabla \text{err}(\mathbf{w}, \mathbf{x}, y) = \frac{\partial -y\mathbf{w}^T \mathbf{x}}{\partial \mathbf{w}} = -y\mathbf{x}$$

$$\Rightarrow \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - (-y\mathbf{x})$$

$\Rightarrow \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y\mathbf{x}$ 與題目的 update function
相同，得證

Multinomial Logistic Regression

7. ANS [a]

In Lecture 6, we solve multiclass classification by OVA or OVO decompositions. One alternative to deal with multiclass classification is to extend the original logistic regression model to Multinomial Logistic Regression (MLR). For a K -class classification problem, we will denote the output space $\mathcal{Y} = \{1, 2, \dots, K\}$. The hypotheses considered by MLR can be indexed by a matrix

$$W = \begin{bmatrix} | & | & \cdots & | & \cdots & | \\ w_1 & w_2 & \cdots & w_k & \cdots & w_K \\ | & | & \cdots & | & \cdots & | \end{bmatrix}_{(d+1) \times K},$$

that contains weight vectors (w_1, \dots, w_K) , each of length $d+1$. The matrix represents a hypothesis

$$h_y(\mathbf{x}) = \frac{\exp(w_y^T \mathbf{x})}{\sum_{i=1}^K \exp(w_i^T \mathbf{x})}$$

that can be used to approximate the target distribution $P(y|\mathbf{x})$ for any (\mathbf{x}, y) . MLR then seeks for the maximum likelihood solution over all such hypotheses. For a given data set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ generated i.i.d. from some $P(\mathbf{x})$ and target distribution $P(y|\mathbf{x})$, the likelihood of $h_y(\mathbf{x})$ is proportional to $\prod_{n=1}^N h_{y_n}(\mathbf{x}_n)$. That is, minimizing the negative log likelihood is equivalent to **minimizing** **an $E_{\text{in}}(W)$** that is composed of the following error function

$$\text{err}(W, \mathbf{x}, y) = -\ln h_y(\mathbf{x}) = -\sum_{k=1}^K \llbracket y = k \rrbracket \ln h_k(\mathbf{x}).$$

When minimizing $E_{\text{in}}(W)$ with SGD, we update the $W^{(t)}$ at the t -th iteration to $W^{(t+1)}$ by

$$W^{(t+1)} \leftarrow W^{(t)} + \eta \cdot V,$$

where V is a $(d+1) \times K$ matrix whose k -th column is an update direction for the k -th weight vector. Assume that an example (\mathbf{x}_n, y_n) is used for the SGD update above. What is the y_n -th column of V ? Choose the correct answer; explain your answer.

- [a] $(1 - h_{y_n}(\mathbf{x}_n))\mathbf{x}_n$
- [b] $(h_{y_n}(\mathbf{x}_n) - 1)\mathbf{x}_n$
- [c] $(-h_{y_n}(\mathbf{x}_n))\mathbf{x}_n$
- [d] $(h_{y_n}(\mathbf{x}_n))\mathbf{x}_n$
- [e] none of the other choices

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \cdot \underline{V}$$

$-\nabla \text{err}$

minimizing $E_{in}(w)$

$$h_y(x) = \frac{\exp(w_y^T x)}{\sum_{k=1}^K \exp(w_k^T x)}$$

$$\text{err}(w, x, y) = -\ln h_y(x) = -\sum_{k=1}^K [y=k] \ln h_k(x)$$

$$W^{(t+1)} = \begin{bmatrix} | & | & | & | \\ w_1^{(t+1)} & w_2^{(t+1)} & \dots & w_K^{(t+1)} \\ | & | & | & | \end{bmatrix}$$

$$V = \begin{bmatrix} | & | & | \\ v_1 & \dots & v_K \\ | & | & | \end{bmatrix}$$

$$W^{(t)} = \begin{bmatrix} | & | & | \\ w_1^{(t)} & w_2^{(t)} & \dots & w_K^{(t)} \\ | & | & | \end{bmatrix}$$

共有 K 個 classifier 要去跑 $w^{(t+1)} \leftarrow w^{(t)} + \eta \cdot \underline{V}$

focus 在第 k 個 classifier 跑的 update rule 上

$$\rightarrow w_k^{(t+1)} \leftarrow w_k^{(t)} + \eta \cdot \underline{v_k}$$

$$v_k = -\nabla \text{err}(w, x, k), \quad (y \text{ 用 } k \text{ 代})$$

$$-\nabla \text{err}(w, x, k) = -\frac{\partial \text{err}(w, x, k)}{\partial w_k} = -\frac{\partial -\ln h_k(x)}{\partial w_k}$$

$$h_k(x) = \frac{\exp(w_k^T x)}{\sum_{i=1}^K \exp(w_i^T x)}$$

$$= -\nabla \left(-\ln \left(\frac{\exp(w_k^T x)}{\sum_{i=1}^K \exp(w_i^T x)} \right) \right)$$

$$= \nabla \left(\ln \exp(w_k^T x) - \ln \sum_{i=1}^K \exp(w_i^T x) \right)$$

$$= \frac{\partial \ln(\square)}{\partial \square} \cdot \frac{\partial \exp(\square)}{\partial \square} \cdot \frac{\partial w_k^T x}{\partial w_k}$$

$$= \frac{\partial \ln(\square)}{\partial \square} \cdot \frac{\partial \sum_{i=1}^K \exp(w_i^T x)}{\partial w_k^T x} \cdot \frac{\partial w_k^T x}{\partial w_k}$$

$$\begin{aligned}
&= \frac{1}{D} \cdot \exp(0) \cdot X \\
&\quad - \frac{1}{D} \cdot \underbrace{\left(\exp(w_1^T x) + \exp(w_2^T x) + \dots + \exp(w_K^T x) \right)}_{\sum w_k^T x} \cdot X \\
&= \frac{\exp(w_K^T x)}{\sum_{k=1}^K \exp(w_k^T x)} \cdot X - \frac{1}{D} \cdot \exp(w_K^T x) \cdot X \\
&= X - \frac{\exp(w_K^T x)}{\sum_{k=1}^K \exp(w_k^T x)} \cdot X = X - h_K(x) \cdot X \\
&= (1 - h_K(x)) X \\
&\text{k} = y_n \Rightarrow v = \underbrace{(1 - h_{y_n}(x_n)) X_n}_{\text{得證. 這裏}} \quad \text{*}
\end{aligned}$$

Nonlinear Transformation

8. ANS [e]

Given the following training data set:

$$\mathbf{x}_1 = (0, 1), y_1 = -1 \quad \mathbf{x}_2 = (0, -1), y_2 = -1 \quad \mathbf{x}_3 = (-1, 0), y_3 = +1 \quad \mathbf{x}_4 = (1, 0), y_4 = +1$$

Use the quadratic transform $\Phi_2(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$ and take $\text{sign}(0) = 1$. Which of the following weight vector $\tilde{\mathbf{w}}$ represents a linear classifier in the \mathcal{Z} -space that can separate all the transformed examples perfectly? Choose the correct answer; explain your answer.

- [a] $(0, -1, 0, 0, 0, 0)$
- [b] $(0, 0, -1, 0, 0, 0)$
- [c] $(0, 0, 0, -1, 0, 0)$
- [d] $(0, 0, 0, 0, -1, 0)$
- [e] $(0, 0, 0, 0, 0, -1)$**

$$\Phi_2(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$$

選最為對應的 $\tilde{\mathbf{w}}$

[a] $(0, -1, 0, 0, 0, 0)$

$$\Rightarrow h(\mathbf{x}) = -x_1$$

$$x_1 = (0, 1), y_1 = -1$$

$$h(x_1) = \text{sign}(0) = 1 \neq y_1 \text{ 錯}$$

[b] $(0, 0, -1, 0, 0, 0)$

$$\Rightarrow h(\mathbf{x}) = -x_2$$

$$x_1 = (0, 1), y_1 = -1 : h(x_1) = -1 = y_1$$

$$x_2 = (0, -1), y_2 = -1 : h(x_2) = 1 \neq y_2 = -1 \text{ 錯}$$

[c] $(0, 0, 0, -1, 0, 0)$

$$\Rightarrow h(\mathbf{x}) = -x_1^2$$

$$x_1 = (0, 1), y_1 = -1 : h(x_1) = \text{sign}(0) = 1 \neq y_1 \text{ 錯}$$

[d] $(0, 0, 0, 0, -1, 0)$

$$\Rightarrow h(\mathbf{x}) = -x_1 x_2$$

$$x_1 = (0, 1), y_1 = -1 : h(x_1) = \text{sign}(0) = 1 \neq y_1 \text{ 錯}$$

(e) $(0, 0, 0, 0, 0, -1)$

$$\exists h(x) = -x_2^2$$

$$x_1 = (0, 1), y_1 = -1 : h(x_1) = -1 = y_1$$

$$x_2 = (0, -1), y_2 = -1 : h(x_2) = -1 = y_2$$

$$x_3 = (-1, 0), y_3 = 1 : h(x_3) = \text{sign}(0) = 1 = y_3$$

$$x_4 = (1, 0), y_4 = 1 : h(x_4) = \text{sign}(0) = 1 = y_4$$

都符合，選 e ✗

9. ANS [b]

Consider a feature transform $\Phi(x) = \Gamma x$ where Γ is a $(d+1)$ by $(d+1)$ invertible matrix. For a training data set $\{(x_n, y_n)\}_{n=1}^N$, run linear regression on the original data set, and get w_{lin} . Then, run linear regression on the Φ -transformed data, and get \tilde{w} . For simplicity, assume that the matrix X (with every row being x_n^T) satisfies that $X^T X$ is invertible. What is the relationship between w_{lin} and \tilde{w} ? Choose the correct answer; explain your answer.

- [a] $w_{\text{lin}} = \Gamma \tilde{w}$
- [b] $w_{\text{lin}} = \Gamma^T \tilde{w}$**
- [c] $w_{\text{lin}} = (\Gamma^{-1})^T \tilde{w}$
- [d] $w_{\text{lin}} = \Gamma^{-1} \tilde{w}$
- [e] none of the other choices

Linear Regression Algorithm

- ➊ from \mathcal{D} , construct **input matrix X** and **output vector y** by

$$X = \underbrace{\begin{bmatrix} \cdots \\ x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}}_{N \times (d+1)} \quad y = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{N \times 1}$$

FUN TIME

- ➋ calculate pseudo-inverse $\underbrace{X^\dagger}_{(d+1) \times N}$
- ➌ return $\underbrace{w_{\text{LIN}}}_{(d+1) \times 1} = X^\dagger y$

After getting w_{LIN} , we can calculate the predictions $\hat{y}_n = w_{\text{LIN}}^T x_n$. If all are collected in a vector \hat{y} similar to how we form y , what is the matrix formula of \hat{y} ?

① y
 ② $XX^T y$
 ③ $XX^\dagger y$
 ④ $XX^\dagger XX^T y$

Reference Answer: ③

根據投影片可知 X

$$\begin{aligned} \Phi(x) &= \Gamma x \quad (X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}) \\ \tilde{x} &= \begin{bmatrix} -(\Gamma x_1)^T \\ \vdots \\ -(\Gamma x_N)^T \end{bmatrix} = \begin{bmatrix} -x_1^T \Gamma^T \\ \vdots \\ -x_N^T \Gamma^T \end{bmatrix} \\ &= \begin{bmatrix} -x_1^T \\ \vdots \\ -x_N^T \end{bmatrix} \Gamma^T = X \Gamma^T \end{aligned}$$

$$\begin{aligned}
\tilde{w} &= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y \\
&= ((X^T)^T (X^T)^{-1})^{-1} (X^T)^T y \\
&= (I X^T \quad X^T)^{-1} (I X^T) y \\
&= (I^T)^{-1} X^{-1} (X^T)^{-1} \underbrace{I^{-1} I}_{\text{In}} X^T y \\
&= (I^T)^{-1} X^{-1} (X^T)^{-1} X^T y \\
&= (I^T)^{-1} \underbrace{(X^T X)^{-1}}_{W_{LIN}} X^T y \\
&= (I^T)^{-1} W_{LIN}
\end{aligned}$$

$$\tilde{w} = (I^T)^{-1} W_{LIN}$$

$$\text{得意 } w_{LIN} = I^T \tilde{w} \neq b$$

10. ANS [c]

After “visualizing” the data and noticing that all $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are distinct, Dr. Trans magically decides the following transform

$$\Phi(\mathbf{x}) = ([\mathbf{x} = \mathbf{x}_1], [\mathbf{x} = \mathbf{x}_2], \dots, [\mathbf{x} = \mathbf{x}_N]).$$

That is, $\Phi(\mathbf{x})$ is a N -dimensional vector whose n -th component is 1 if and only if $\mathbf{x} = \mathbf{x}_n$. If we run linear regression after applying this transform, what is the optimal $\tilde{\mathbf{w}}$? Choose the correct answer; explain your answer.

[a] $\mathbf{1}$, the vector of all 1s.

[b] $\mathbf{0}$, the vector of all 0s.

[c] \mathbf{y}

[d] $-\mathbf{y}$

[e] none of the other choices

(Note: Be sure to also check what $E_{\text{in}}(\tilde{\mathbf{w}})$ is!)

$$\Phi(\mathbf{x}) = ([\mathbf{x} = \mathbf{x}_1], [\mathbf{x} = \mathbf{x}_2], \dots, [\mathbf{x} = \mathbf{x}_N])$$

$$\begin{aligned} & N \text{ 單資料 經過 } \Phi(\mathbf{x}) \\ \Leftrightarrow \Phi(\mathbf{x}_1) &= (\underbrace{1, 0, 0, \dots, 0}_N) \\ \Phi(\mathbf{x}_2) &= (0, 1, 0, \dots, 0) \\ & \vdots \\ \Phi(\mathbf{x}_N) &= (0, 0, \dots, 0, 1) \end{aligned}$$

$$\Rightarrow \text{transform } \Phi(\mathbf{x}) = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} = I_N$$

Linear Regression 的 optimal \mathbf{w} 是 \mathbf{w}_{LIN}

$$\tilde{\mathbf{w}}_{\text{LIN}} = \mathbf{X}^T \mathbf{y} \quad \text{pseudo-inverse}$$

$$\tilde{\mathbf{w}}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = (I_N^T I_N)^{-1} I_N = I_N$$

$$\tilde{\mathbf{w}}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = I_N \mathbf{y} = \mathbf{y}$$

選 c。

Linear Regression

11. ANS [c]

Assume that we couple linear regression with one-versus-all decomposition for multi-class classification, and get K weight vectors $\mathbf{w}_{[k]}^*$. Assume that the squared error $E_{in}^{sqr}(\mathbf{w}_{[k]}^*)$ for the k -th binary classification problem is e_k . What is the tightest upper bound of $E_{in}^{0/1}(g)$, where g is the multi-class classifier formed by the one-versus-all decomposition? Choose the correct answer; explain your answer.

- [a] $2 \sum_{k=1}^K e_k$
- [b] $\sum_{k=1}^K e_k$
- [c] $\frac{1}{2} \sum_{k=1}^K e_k$**
- [d] $\frac{1}{K} \sum_{k=1}^K e_k$
- [e] $\frac{1}{2K} \sum_{k=1}^K e_k$

For one binary classifier

$$\text{if } E_{in}^{sqr}(\mathbf{w}_{[k]}^*) = e_k$$

$$e_k = \left[\mathbf{w}_{[k]}^{*\top} \mathbf{x} - (2 \llbracket y=k \rrbracket - 1) \right]^2$$

$$\text{if } E_{in}^{0/1}(\mathbf{w}_{[k]}^*) = \left[\text{sign}(\mathbf{w}_{[k]}^{*\top} \mathbf{x}) = \llbracket y=k \rrbracket \right]$$

假設 $k=4$: $\Delta, O, \square, \star$

有4個 binary classifier

Δ : Δ or not Δ

O : O or not O

\square : \square or not \square

\star : \star or not \star

假設 Label: Δ

傳出來

Δ : $! \Delta$ 錯

O : O 錯

\square : $! \square$

\star : $! \star$

binary
最少 2 個 classifier

發生錯誤
就會使結果出錯

要找 $E_{in}^{0/1}(g)$ 的 tightest upper bound

M : 至少多少個 classifier 判斷錯誤才會出現 mistake

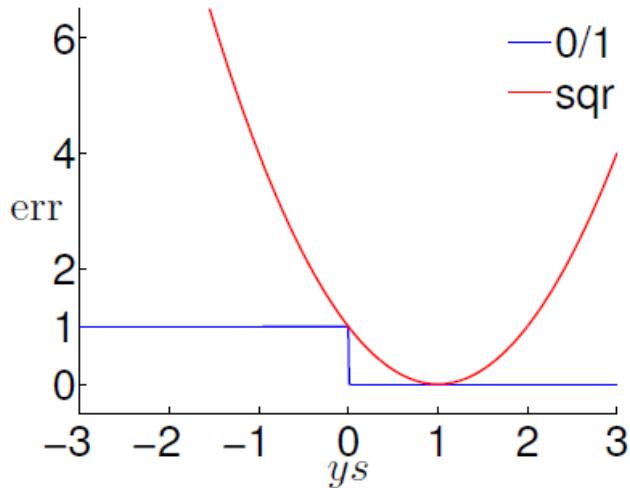
$$\underline{M=2}$$

$$\begin{aligned} E_{in}^{0/1}(g) &= \frac{1}{2} \sum_{k=1}^K E_{in}^{0/1}(w_k^*) \leq \frac{1}{2} \sum_{k=1}^K E_{in}^{sqr}(w_k^*) \\ &\leq \frac{1}{2} \sum_{k=1}^K e_k \end{aligned}$$

假設 g 有 10 個犯錯

代表在 K 個 w_k^* 中至少有 2×10 個 binary classifier 犯錯

$E_{in}^{0/1}$ 會一直在 E_{in}^{sqr} 下面 (如下圖)



選 C。

Experiments with Linear and Nonlinear Models

Next, we will play with transform + linear regression for binary classification. Please use the following set for training:

https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat

and the following set for testing (estimating E_{out}):

https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat

Each line of the data set contains one (\mathbf{x}_n, y_n) with $\mathbf{x}_n \in \mathbb{R}^{10}$. The first 10 numbers of the line contains the components of \mathbf{x}_n orderly, the last number is y_n , which belongs to $\{-1, +1\} \subseteq \mathbb{R}$. That is, we can use those y_n for either binary classification or regression.

12. ANS [b]

(*) Consider the following *homogeneous* order- Q polynomial transform

$$\Phi(\mathbf{x}) = (1, x_1, x_2, \dots, x_{10}, x_1^2, x_2^2, \dots, x_{10}^2, \dots, x_1^Q, x_2^Q, \dots, x_{10}^Q).$$

Transform the training and testing data according to $\Phi(\mathbf{x})$ with $Q = 2$, and implement the linear regression algorithm on the transformed data. What is $|E_{\text{in}}^{0/1}(g) - E_{\text{out}}^{0/1}(g)|$, where g is the hypothesis returned by the transform + linear regression procedure? Choose the closest answer; provide your code.

- [a] 0.28
- [b] 0.32** 21
Ein: 0.181
- [c] 0.36 Eout: 0.5073333333333333
- [d] 0.40 0.3263333333333333
- [e] 0.44

Colab :

https://colab.research.google.com/drive/18Z_WwMk1o4neFzO4tIG21bEd18GpWKu9?usp=sharing

```
import urllib.request
import numpy as np
import random
import math
```

```

def getData(url):
    #get data
    contents = urllib.request.urlopen(url).read()
    text = str(contents, 'utf-8')
    # data preprocessing
    sptext = text.split('\n')

X = []
Y = []

for i in range(len(sptext)-1):
    data = sptext[i].split('\t')
    # print(data)
    label = data[-1]
    # print(type(data))
    data.pop()
    # print(len(data))
    X += [data]
    Y += [label]

X = np.asarray(X, dtype = float)
Y = np.asarray(Y, dtype = float)
# print(X)
# print(Y)
return X, Y

```

```

def transform(X, Q):
    trans_X = []
    for x in X:
        x_Q = []
        for q in range(2, Q+1):
            x_q = np.zeros(len(x))
            x_q = np.power(x, q)
            x_Q = np.concatenate([x_Q, x_q])
        # print(x_Q)
        x = np.insert(x, 0, 1)
        x_Q = np.asarray(x_Q)
        x = np.concatenate([x, x_Q])
        trans_X.append(x)
    # print(trans_X[0])
    return trans_X

```

```

if __name__ == '__main__':
    train_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat"
    test_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat"
    train_X, train_Y = getData(train_data)
    test_X, test_Y = getData(test_data)
    Q = 2
    trans_train_X = transform(train_X, Q)
    # print(type(trans_train_X[0]))
    trans_test_X = transform(test_X, Q)
    # print(len(trans_test_X[0]))

    # Calculate pseudo-inverse X_cross
    trans_train_X_cross = np.linalg.pinv(trans_train_X)
    # print(trainX_cross, type(trainX_cross))

    # Calculate WLIN = X_cross . Y
    WLIN = np.matmul(trans_train_X_cross, train_Y)
    print(len(WLIN))

    # Calculate Ein
    EinN = 0
    for i in range(len(trans_train_X)):
        h_xn = np.matmul(WLIN.T, trans_train_X[i])
        if np.sign(h_xn) != train_Y[i]:
            EinN += 1
    # print("EinN: ", EinN)
    Ein = EinN/len(trans_train_X)
    print("Ein: ", Ein)

    # Calculate Eout
    EoutN = 0
    for i in range(len(trans_test_X)):
        h_xn = np.matmul(WLIN.T, trans_test_X[i])
        if np.sign(h_xn) != test_Y[i]:
            EoutN += 1
    # print("EoutN: ", EoutN)
    Eout = EoutN/len(trans_test_X)
    print("Eout: ", Eout)

    result = abs(Ein - Eout)
    print(result)

```

結果： 0.3263333333333333 . 選 b 。

13. ANS [d]

(*) Repeat the previous problem, but with $Q = 8$ instead. What is $|E_{\text{in}}^{0/1}(g) - E_{\text{out}}^{0/1}(g)|$, where g is the hypothesis returned by the transform + linear regression procedure? Choose the closest answer; provide your code.

- [a] 0.30 81
- [b] 0.35 Ein: 0.052
- [c] 0.40 Eout: 0.5096666666666667
- [d] 0.45** 0.4576666666666667
- [e] 0.50

Colab :

[https://colab.research.google.com/drive/15zNujr4Yuvh9ChRQ0bCMdUIbBm0j69I ?usp=sharing](https://colab.research.google.com/drive/15zNujr4Yuvh9ChRQ0bCMdUIbBm0j69I?usp=sharing)

```
import urllib.request
import numpy as np
import random
import math
```

```

def getData(url):
    #get data
    contents = urllib.request.urlopen(url).read()
    text = str(contents, 'utf-8')
    # data preprocessing
    sptext = text.split('\n')

X = []
Y = []

for i in range(len(sptext)-1):
    data = sptext[i].split('\t')
    # print(data)
    label = data[-1]
    # print(type(data))
    data.pop()
    # print(len(data))
    X += [data]
    Y += [label]

X = np.asarray(X, dtype = float)
Y = np.asarray(Y, dtype = float)
# print(X)
# print(Y)
return X, Y

```

```

def transform(X, Q):
    trans_X = []
    for x in X:
        x_Q = []
        for q in range(2, Q+1):
            x_q = np.zeros(len(x))
            x_q = np.power(x, q)
            x_Q = np.concatenate([x_Q, x_q])
        # print(x_Q)
        x = np.insert(x, 0, 1)
        x_Q = np.asarray(x_Q)
        x = np.concatenate([x, x_Q])
        trans_X.append(x)
    # print(trans_X[0])
    return trans_X

```

```

if __name__ == '__main__':
    train_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat"
    test_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat"
    train_X, train_Y = getData(train_data)
    test_X, test_Y = getData(test_data)
    Q = 8
    trans_train_X = transform(train_X, Q)
    # print(type(trans_train_X[0]))
    trans_test_X = transform(test_X, Q)
    # print(len(trans_test_X[0]))

    # Calculate pseudo-inverse X_cross
    trans_train_X_cross = np.linalg.pinv(trans_train_X)
    # print(trainX_cross, type(trainX_cross))

    # Calculate WLIN = X_cross . Y
    WLIN = np.matmul(trans_train_X_cross, train_Y)
    print(len(WLIN))

    # Calculate Ein
    EinN = 0
    for i in range(len(trans_train_X)):
        h_xn = np.matmul(WLIN.T, trans_train_X[i])
        if np.sign(h_xn) != train_Y[i]:
            EinN += 1
    # print("EinN: ", EinN)
    Ein = EinN/len(trans_train_X)
    print("Ein: ", Ein)

    # Calculate Eout
    EoutN = 0
    for i in range(len(trans_test_X)):
        h_xn = np.matmul(WLIN.T, trans_test_X[i])
        if np.sign(h_xn) != test_Y[i]:
            EoutN += 1
    # print("EoutN: ", EoutN)
    Eout = EoutN/len(trans_test_X)
    print("Eout: ", Eout)

    result = abs(Ein - Eout)
    print(result)

```

結果：0.4576666666666667，選 d。

14. ANS [a]

(*) Repeat the previous problem, but with Φ_2 (the full order-2 polynomial transform introduced in the lecture, which is of $1 + 10 + 45 + 10$ dimensions) instead. What is $|E_{\text{in}}^{0/1}(g) - E_{\text{out}}^{0/1}(g)|$, where g is the hypothesis returned by the transform + linear regression procedure? Choose the closest answer; provide your code.

- [a] 0.33 66
- [b] 0.41 Ein: 0.168
- [c] 0.49 Eout: 0.5066666666666667
- [d] 0.57 0.3386666666666667
- [e] 0.65

Colab :

<https://colab.research.google.com/drive/1FV1xQNaU9aZW-1zMVmUhwhK7fiNIkKEz?usp=sharing>

```
import urllib.request
import numpy as np
import random
import math

def getData(url):
    #get data
    contents = urllib.request.urlopen(url).read()
    text = str(contents,'utf-8')
    # data preprocessing
    sptext = text.split('\n')

    X = []
    Y = []

    for i in range(len(sptext)-1):
        data = sptext[i].split('\t')
        # print(data)
        label = data[-1]
        # print(type(data))
        data.pop()
        # print(len(data))
        X += [data]
        Y += [label]

    X = np.asarray(X, dtype = float)
    Y = np.asarray(Y, dtype = float)
    # print(X)
    # print(Y)
    return X, Y
```

```
def transform_all(X, Q):
    trans_X = []
    for x in X:
        x_Q = []
        for i in range(len(x)):
            for j in range(i+1):
                x_Q.append(x[i]*x[j])
        # print(len(x_Q))
        x = np.insert(x, 0, 1)
        x_Q = np.asarray(x_Q)
        x = np.concatenate([x, x_Q])
        trans_X.append(x)
    # print(trans_X[0])

    return trans_X
```

主要就是差在要產生交叉的項，利用雙層 for 迴圈去兩兩相乘。

```

if __name__ == '__main__':
    train_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat"
    test_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat"
    train_X, train_Y = getData(train_data)
    test_X, test_Y = getData(test_data)
    Q = 8
    trans_train_X = transform_all(train_X, Q)
    # print(type(trans_train_X[0]))
    trans_test_X = transform_all(test_X, Q)
    # print(len(trans_test_X[0]))

    # Calculate pseudo-inverse X_cross
    trans_train_X_cross = np.linalg.pinv(trans_train_X)
    # print(trainX_cross, type(trainX_cross))

    # Calculate WLIN = X_cross . Y
    WLIN = np.matmul(trans_train_X_cross, train_Y)
    print(len(WLIN))

    # Calculate Ein
    EinN = 0
    for i in range(len(trans_train_X)):
        h_xn = np.matmul(WLIN.T, trans_train_X[i])
        if np.sign(h_xn) != train_Y[i]:
            EinN += 1
    # print("EinN: ", EinN)
    Ein = EinN/len(trans_train_X)
    print("Ein: ", Ein)

    # Calculate Eout
    EoutN = 0
    for i in range(len(trans_test_X)):
        h_xn = np.matmul(WLIN.T, trans_test_X[i])
        if np.sign(h_xn) != test_Y[i]:
            EoutN += 1
    # print("EoutN: ", EoutN)
    Eout = EoutN/len(trans_test_X)
    print("Eout: ", Eout)

    result = abs(Ein - Eout)
    print(result)

```

結果： 0.33866666666666667，選 a。

15. ANS [c]

(*) Instead of transforming to a higher dimensional space, we can also transform to a lower dimensional space. Consider the following 10 transforms:

$$\begin{aligned}\Phi^{(1)}(\mathbf{x}) &= (x_0, x_1) \\ \Phi^{(2)}(\mathbf{x}) &= (x_0, x_1, x_2) \\ &\dots \\ \Phi^{(10)}(\mathbf{x}) &= (x_0, x_1, x_2, \dots, x_{10})\end{aligned}$$

Run $\Phi^{(i)}$ + linear regression to get a hypothesis g_i . Which i leads to a $\Phi^{(i)}$ that reaches the minimum $|E_{\text{in}}^{0/1}(g_i) - E_{\text{out}}^{0/1}(g_i)|$? Choose the closest answer; provide your code.

- [a] 1
- [b] 2
- [c] 3** 3
- [d] 5
- [e] 8

Colab :

<https://colab.research.google.com/drive/1dsdq4NWMBq0R6YL0orOvoMyh30TwAoZn?usp=sharing>

```
import urllib.request
import numpy as np
import random
import math
```

```

def getData(url):
    #get data
    contents = urllib.request.urlopen(url).read()
    text = str(contents, 'utf-8')
    # data preprocessing
    sptext = text.split('\n')

    X = []
    Y = []

    for i in range(len(sptext)-1):
        data = sptext[i].split('\t')
        # print(data)
        label = data[-1]
        # print(type(data))
        data.pop()
        # print(len(data))
        X += [data]
        Y += [label]

    X = np.asarray(X, dtype = float)
    Y = np.asarray(Y, dtype = float)
    # print(X)
    # print(Y)
    return X, Y

def transform_low(X, i):
    trans_X = []
    for x in X:
        x = x[:i+1]
        x = np.insert(x, 0, 1)
        trans_X.append(x)
    # print(trans_X[0])
    return trans_X

```

每次只抓出前 $i+1$ 項(i 從 0 開始) · 再去補 1。

```

if __name__ == '__main__':
    train_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat"
    test_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat"
    train_X, train_Y = getData(train_data)
    test_X, test_Y = getData(test_data)
    I = 10
    E = []
    for i in range(I):
        trans_train_X = transform_low(train_X, i)
        # print(type(trans_train_X[0]))
        trans_test_X = transform_low(test_X, i)
        # print(len(trans_test_X[0]))

        # Calculate pseudo-inverse X_cross
        trans_train_X_cross = np.linalg.pinv(trans_train_X)
        # print(trans_train_X_cross, type(trans_train_X_cross))

        # Calculate WLIN = X_cross * Y
        WLIN = np.matmul(trans_train_X_cross, train_Y)
        print(len(WLIN))

        # Calculate Ein
        EinN = 0
        for i in range(len(trans_train_X)):
            h_xn = np.matmul(WLIN.T, trans_train_X[i])
            if np.sign(h_xn) != train_Y[i]:
                EinN += 1
        # print("EinN: ", EinN)
        Ein = EinN/len(trans_train_X)
        print("Ein: ", Ein)

        # Calculate Eout
        EoutN = 0
        for i in range(len(trans_test_X)):
            h_xn = np.matmul(WLIN.T, trans_test_X[i])
            if np.sign(h_xn) != test_Y[i]:
                EoutN += 1
        # print("EoutN: ", EoutN)
        Eout = EoutN/len(trans_test_X)
        print("Eout: ", Eout)

        result = abs(Ein - Eout)
        print(result)
        E.append(result)

    print(E)
    print(E.index(min(E))+1)
    print(min(E))

```

3

結果：0.1323333333333333，選 C。

16. ANS [d]

(*) Consider a transform that randomly chooses 5 out of 10 dimensions. That is, $\Phi(\mathbf{x}) = (x_0, x_{i_1}, x_{i_2}, x_{i_3}, x_{i_4}, x_{i_5})$, where i_1 to i_5 are distinct random integers uniformly and independently generated within $\{1, 2, \dots, 10\}$. Run Φ + linear regression to get a hypothesis g . What is the average $|E_{\text{in}}^{0/1}(g_i) - E_{\text{out}}^{0/1}(g_i)|$ over 200 experiments, each generating Φ with a different random seed? Choose the closest answer; provide your code.

- [a] 0.06
- [b] 0.11
- [c] 0.16
- [d] 0.21**
- [e] 0.26

Colab :

<https://colab.research.google.com/drive/1VSDBvI55SRdAiN53wn4tOIBPqlOltgjN?usp=sharing>

```
import urllib.request
import numpy as np
import random
import math
```

```

def getData(url):
    #get data
    contents = urllib.request.urlopen(url).read()
    text = str(contents, 'utf-8')
    # data preprocessing
    sptext = text.split('\n')

    X = []
    Y = []

    for i in range(len(sptext)-1):
        data = sptext[i].split('\t')
        # print(data)
        label = data[-1]
        # print(type(data))
        data.pop()
        # print(len(data))
        X += [data]
        Y += [label]

    X = np.asarray(X, dtype = float)
    Y = np.asarray(Y, dtype = float)
    # print(X)
    # print(Y)
    return X, Y

```

```

def transform_ran(X, x_i):
    trans_X = []
    for x in X:
        xl = []
        for i in x_i:
            xl.append(x[i-1])
        # print(xl)
        x = np.asarray(xl)
        x = np.insert(x, 0, 1)
        trans_X.append(x)
    # print(trans_X[0])
    return trans_X

```

在主程式產生好每次 iteration 要隨機抓取的五個 index，再進入 transform_ran

function，依照 random index list : x_i 來抓出每一筆 x。

```

if __name__ == '__main__':
    train_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_train.dat"
    test_data = "https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw3/hw3_test.dat"
    train_X, train_Y = getData(train_data)
    test_X, test_Y = getData(test_data)
    iteration = 200
    E = []
    x_index = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
    for iter in range(iteration):
        random.seed(iter)
        x_i = random.sample(x_index, 5)
        print(x_i)
        trans_train_X = transform_ran(train_X, x_i)
        # print(type(trans_train_X[0]))
        trans_test_X = transform_ran(test_X, x_i)
        # print(len(trans_test_X[0]))

        # Calculate pseudo-inverse X_cross
        trans_train_X_cross = np.linalg.pinv(trans_train_X)
        # print(trainX_cross, type(trainX_cross))

        # Calculate WLIN = X_cross * Y
        WLIN = np.matmul(trans_train_X_cross, train_Y)
        # print(len(WLIN))

        # Calculate Ein
        EinN = 0
        for i in range(len(trans_train_X)):
            h_xn = np.matmul(WLIN.T, trans_train_X[i])
            if np.sign(h_xn) != train_Y[i]:
                EinN += 1
        # print("EinN: ", EinN)
        Ein = EinN/len(trans_train_X)
        print("Ein: ", Ein)

        # Calculate Eout
        EoutN = 0
        for i in range(len(trans_test_X)):
            h_xn = np.matmul(WLIN.T, trans_test_X[i])
            if np.sign(h_xn) != test_Y[i]:
                EoutN += 1
        # print("EoutN: ", EoutN)
        Eout = EoutN/len(trans_test_X)
        print("Eout: ", Eout)

        result = abs(Ein - Eout)
        print("Ein - Eout: ", result)
        E.append(result)

    print(len(E))
    print(np.mean(E))

```