

## HW2

1.

令  $P(+1|x) = \alpha$ ，則如果選  $y=+1$  當 mini-target 的話，有  $1-\alpha$  的機率會出現 false negative。Cost 是  $(1-\alpha)*1$ 。如果選  $y=-1$  當 mini-target 的話，有  $t$  的機率會出現 false positive。Cost 是  $\alpha*1000$ 。所以臨界值在於  $\alpha*1000=1-\alpha$ ，所以  $\alpha = \frac{1}{1001}$

2.

根據定義  $E_{\text{out}}(g) = E_{x \sim P(x)}[g(x) \neq f(x)]$ 。令總資料量為  $D$ ，則  $g(x) \neq f(x)$  的資料量為  $D * E_{\text{out}}(g)$ ，而  $g(x) = f(x)$  的資料量為  $D * (1 - E_{\text{out}}(g))$ 。這次的狀況是  $P(y = +f(x)|x) = 1 - \epsilon$ ，所以對於  $g(x) \neq f(x)$  的資料而言，維持錯誤的只剩  $(1 - \epsilon) * D * E_{\text{out}}(g)$  的資料量。變成正確的有  $\epsilon * D * E_{\text{out}}(g)$  的資料量。對於  $g(x) = f(x)$  的資料而言，維持正確的只剩  $(1 - \epsilon) * D * (1 - E_{\text{out}}(g))$  的資料量。變成錯誤的有  $\epsilon * D * (1 - E_{\text{out}}(g))$  的資料量。所以錯誤情況的期望值是  $((1 - \epsilon) * D * E_{\text{out}}(g) + \epsilon * D * (1 - E_{\text{out}}(g))) / D = (1 - \epsilon) * E_{\text{out}}(g) + \epsilon * (1 - E_{\text{out}}(g)) = E_{\text{out}}(g) - 2 * \epsilon * E_{\text{out}}(g) + \epsilon$ 。

3.

硬推。首先  $E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N (wx_n - y_n)^2$ ，所以對於  $w$ ，可以算偏微分  $\frac{\partial E_{\text{in}}(w)}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N 2 * (wx_n - y_n) * x_n$ 。偏微分等於，所以

$\sum_{n=1}^N (wx_n - y_n) * x_n = 0$ 。展開可得  $w \sum_{n=1}^N x_n^2 = \sum_{n=1}^N y_n x_n$ ，所以

$$w = \frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2}。$$

4.

沒有給點的數量，這裡我們用積分來模擬點的數量趨近無限的情況，也更能更精準得出  $\text{err}$ ， $\text{err} = \int_0^1 (w_0 + w_1 x - ax^2 - b)^2 dx = \int_0^1 w_0^2 + w_1^2 x^2 + a^2 x^4 + b^2 + 2w_0 w_1 x - 2w_0 ax^2 - 2w_0 b - 2w_1 ax^3 - 2w_1 bx - 2abx^2 dx = w_0^2 + \frac{w_1^2}{3} + \frac{a^2}{5} + b^2 + w_0 w_1 - \frac{2w_0 a}{3} - 2w_0 b - \frac{w_1 a}{2} - w_1 b - \frac{2ab}{3}$ 。首先對  $w_0$  做偏微分=0，可得  $2w_0 + w_1 = \frac{2a}{3} + 2b$ 。對  $w_1$  做偏微分=0，可得  $\frac{2w_1}{3} + w_0 = \frac{a}{2} + b$ 。解聯立可得  $w_0 = -\frac{a}{6} + b$ ， $w_1 = a$ 。

5.

根據公式， $w_{\text{lin}} = (X^T X)^{-1} X^T y$ 。其中  $y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix}$ 。根據題目，對

所有  $y_n$ ，有  $y'_n = ay_n + b$ ，所以  $y' = \begin{bmatrix} ay_1 + b \\ ay_2 + b \\ \dots \\ ay_N + b \end{bmatrix} = \begin{bmatrix} ay_1 \\ ay_2 \\ \dots \\ ay_N \end{bmatrix} + b \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} =$

$ay + b \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$ 。因此， $w'_{\text{lin}} = (X^T X)^{-1} X^T y' = (X^T X)^{-1} X^T \left( ay + \right.$

$$b \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} = a w_{\text{lin}} + b (X^T X)^{-1} X^T \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} = a w_{\text{lin}} + b \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}。這裡我們觀察$$

一下  $(X^T X)^{-1} X^T \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$  這個東西。這東西相當於 data 的 y 都是 1 的時

候求  $w_{\text{lin}}$ 。因為在這種情況，對所有資料  $x_n$ ，其第一維度，也就是

$x_0$ ，都是 1。所有有個非常直觀的解，就是  $w_{\text{lin}} = \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}$ 。這時，對所有

$x_n$ ， $w_{\text{lin}}^T x_n = [1, 0, 0, \dots, 0] \begin{bmatrix} 1 \\ x_1 \\ \dots \\ x_d \end{bmatrix} = 1$ ，所以  $\text{err} = (w_{\text{lin}}^T x_n - y)^2 = 0$ 。因為

是 square error，所以 err 必然大於等於 0。所以  $\begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}$  必然是最優解，

因為不可能有 err 比 err=0 還小。題目說  $w_{\text{lin}}$  有唯一性，所以  $\begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}$  是

唯一的，所以符合  $(X^T X)^{-1} X^T \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}$ 。因此  $a w_{\text{lin}} +$

$$b (X^T X)^{-1} X^T \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} = a w_{\text{lin}} + b \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}。$$

6.

根據定義， $E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n w^T x_n))$ 。對於任意的  $w_i$ ，上課有說過， $\frac{\partial E_{in}(w)}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(y_n w^T x_n)} (-y_n x_{ni}) = \frac{1}{N} \sum_{n=1}^N h_t(y_n x_n) (-y_n x_{ni})$ 。這時，任意再挑一個  $w_j$ ，則  $\frac{\partial^2 E_{in}(w)}{\partial w_i \partial w_j} = \frac{1}{N} \sum_{n=1}^N \frac{y_n x_{ni}}{((1 + \exp(y_n w^T x_n))^2} * \exp(y_n w^T x_n) * y_n x_{nj} = \frac{1}{N} \sum_{n=1}^N \frac{x_{ni} x_{nj}}{((1 + \exp(y_n w^T x_n))^2} * \exp(y_n w^T x_n) \circ \frac{1}{N} \sum_{n=1}^N \frac{x_{ni} x_{nj}}{((1 + \exp(y_n w^T x_n))^2} * \exp(y_n w^T x_n)$  就是  $A_E(w)$  的第  $ij$  項。所以當  $X = \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_N \end{bmatrix}$  時 ( $X$  是個  $N \times (d+1)$

矩陣)， $X^T X$  的第  $ij$  項就是  $\sum_{n=1}^N x_{ni} x_{nj}$ 。給定  $w_t$ ， $A_E(w_t)$  可以表示為

$X^T D X$ ，因為  $D$  是 diagonal，所以  $D$  的對角線上第  $i$  個值是

$$\frac{1}{N} \frac{\exp(y_i w_t^T x_i)}{(1 + \exp(y_i w_t^T x_i))^2} = \frac{1}{N} \frac{1}{1 + \exp(y_i w_t^T x_i)} * \frac{\exp(y_i w_t^T x_i)}{1 + \exp(y_i w_t^T x_i)} = \frac{1}{N} \frac{1}{1 + \exp(y_i w_t^T x_i)} *$$

$$\frac{1}{1 + \exp(-y_i w_t^T x_i)} = \frac{1}{N} h_t(y_i x_i) * h_t(-y_i x_i) \circ$$

總共有  $N$  個，所以  $D$  是個  $N \times N$  矩陣

7.

PLA 就是每一步挑錯誤的 1 個點，做  $w_{t+1} = w_t + y_n x_n$ 。題目要求的演算法是每一步挑 1 個點，做  $w_{t+1} = w_t - \eta \nabla err = w_t -$

$\eta \nabla \max(0, 1 - y w^T x)$ 。而  $\max(0, 1 - y w^T x)$  就是

$$\begin{cases} 0 & \text{if } y w^T x \geq 1 \\ 1 - y w^T x & \text{if } y w^T x < 1 \end{cases}, \text{ 所以對任意 } w_i, \text{ 偏微分是}$$

$\begin{cases} 0 & \text{if } yw^T x \geq 1 \\ -yx_i & \text{if } yw^T x < 1 \end{cases}$ ，所以梯度是 $\begin{cases} 0 & \text{if } yw^T x \geq 1 \\ -yx & \text{if } yw^T x < 1 \end{cases}$ ，也就是 $\max(0, 1 - yx)$ 。所以更新 $w_{t+1}$ 的算式就是 $w_{t+1} = w_t + \eta y_n x_n$  if  $yw^T x < 1$ 。兩者比較，PLA 是 $yw^T x < 0$  的點才會更新。SGD 是 $yw^T x < 1$  的點就會更新。條件比較寬，所以 SGD 可以停下來，PLA 就可以停下來。但 PLA 可以停下來，SGD 不一定可以，可能所有 data 都符合 $0 < yw^T x < 1$ ，這樣 SGD 就停不下來。根據這個式子，可以發現  $x$  的大小會對 SGD 的判斷產生影響，所以 normalized 會對 SGD 能不能停下有影響。而且 SGD 有 $\eta$ 這個參數，會決定 SGD 會跑多快。

8.

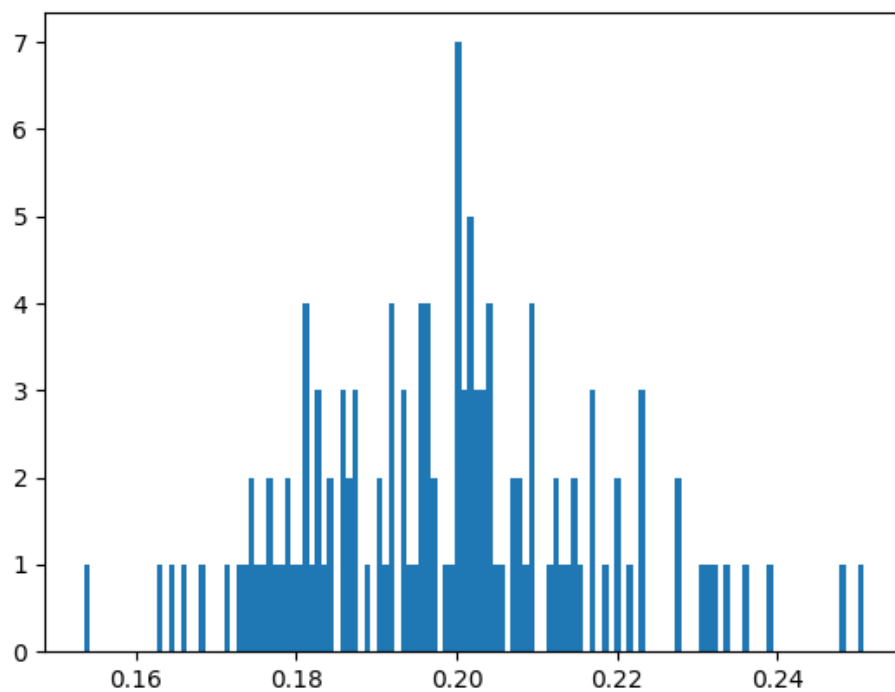
$E_{in}(W) = \frac{-1}{N} \sum_{n=1}^N \ln \frac{\exp(w_{y_n}^T x_n)}{\sum_{i=1}^K \exp(w_i^T x_n)} = \frac{-1}{N} \sum_{n=1}^N (\ln \exp(w_{y_n}^T x_n) - \ln(\sum_{i=1}^K \exp(w_i^T x_n))) = \frac{-1}{N} \sum_{n=1}^N (w_{y_n}^T x_n - \ln(\sum_{i=1}^K \exp(w_i^T x_n)))$ 。對於  $W$  中任意的 $w_{kj}$ ，做偏微分。當 $y_n = k$ 時， $w_{kj}$ 的偏微分如下。

$\frac{\partial E_{in}(W)}{\partial w_{kj}} = \frac{-1}{N} \sum_{n=1}^N (x_{nj} - \frac{x_{nj} \exp(w_k^T x_n)}{\sum_{i=1}^K \exp(w_i^T x_n)}) = \frac{-1}{N} \sum_{n=1}^N x_{nj} (1 - \frac{\exp(w_{y_n}^T x_n)}{\sum_{i=1}^K \exp(w_i^T x_n)}) = \frac{-1}{N} \sum_{n=1}^N x_{nj} (1 - h_{y_n}(x_n))$ 。當 $y_n \neq k$ 時， $w_{kj}$ 的偏微分如下。 $\frac{\partial E_{in}(W)}{\partial w_{kj}} = \frac{-1}{N} \sum_{n=1}^N (-\frac{x_{nj} \exp(w_k^T x_n)}{\sum_{i=1}^K \exp(w_i^T x_n)}) = \frac{-1}{N} \sum_{n=1}^N x_{nj} (-h_{y_n}(x_n))$ 。當 $y_n = k$ 時，集合 $w_{kj}$ 所有  $j$ ，則 $\nabla w_{y_n} = \frac{-1}{N} \sum_{n=1}^N x_n (1 - h_{y_n}(x_n))$ ，是個  $(d+1)*1$  的矩陣。當 $y_n \neq k$ 時，集合 $w_{kj}$ 所有  $j$ ，則 $\nabla w_k =$

$\frac{-1}{N} \sum_{n=1}^N x_n (-h_k(x_n))$ ，是個 $(d+1)*1$  的矩陣。把 $\nabla w_1 \sim \nabla w_K$ 這  $K$  個向量橫著排起來，就是一個 $(d+1)*K$  的矩陣，也就是 $\nabla E_{in}(W)$ 。

9.

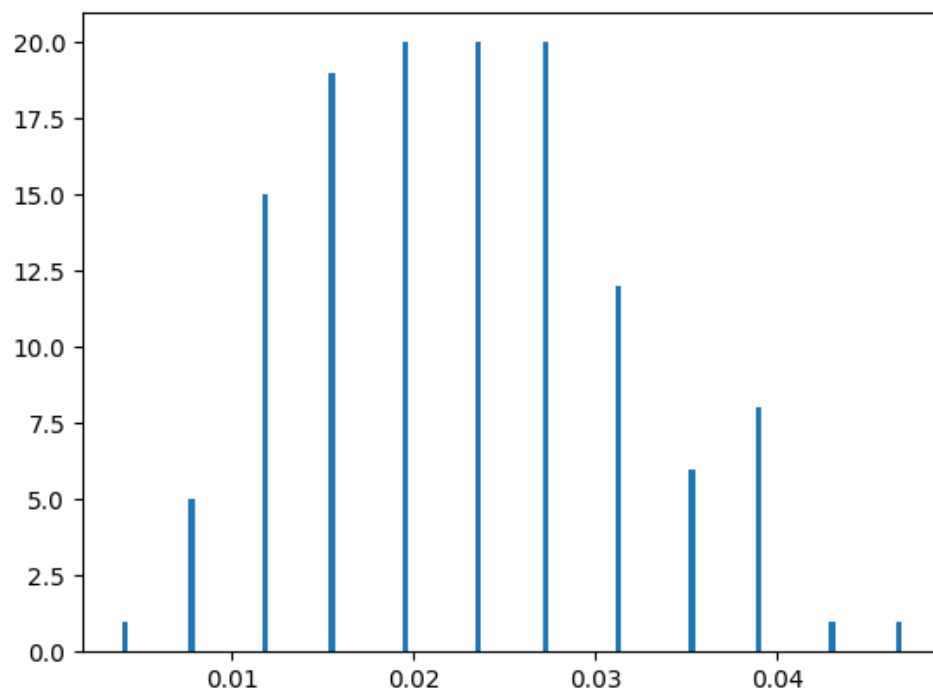
這是我畫的圖:



err 中位數=0.199

10

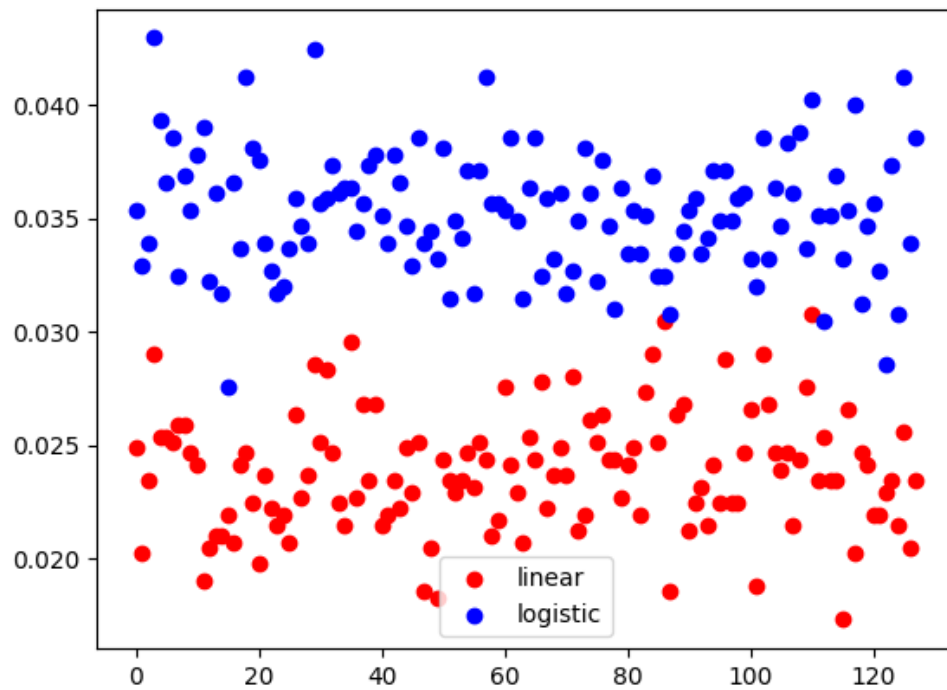
這是我畫的圖:



Err 中位數是 0.023。也就是 256 個平均約 5 個出錯。如上課所說  $\text{sqr err}$  比 0/1 err 大。

11.

這是我畫的圖：



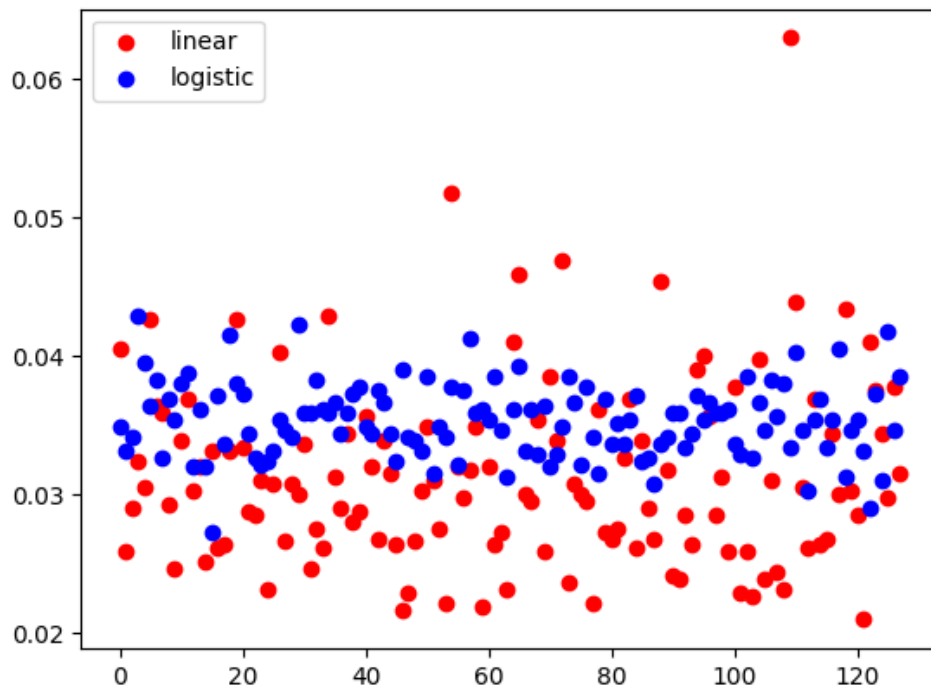
Linear regression 的 median err 是 0.023

Logistic regression 的 median err 是 0.034

12.

這是我畫的圖:





Linear regression 的 median err 是 0.030

Logistic regression 的 median err 是 0.035

linear regression 的 err 變化幅度較大，而 Logistic regression 的 err 變化幅度較小。

可以發現說 linear regression 比較受雜訊的影響，而 Logistic regression 受到的影響比較小。

13.

對於 $\hat{X}$ 而言，因為我第 6 題算出來的 D 其對角第 i 項是

$\frac{1}{N} h_t(y_i x_i) * h_t(-y_i x_i)$ ，把這個值叫做 $D_i$ 好了。只要把每個 $D_i$ 開根號，乘到 X 裡面每第 i 個 row 就好了。因為 $h_t$ 必為正，所以 $D_i$ 開根

號必為實數。

對於 $\hat{y}$ 而言，考慮第 6 題的 $E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n w^T x_n))$ ，

$\frac{\partial E_{in}(w)}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(y_n w^T x_n)} (-y_n x_{ni}) = \frac{1}{N} \sum_{n=1}^N h_t(y_n x_n) (-y_n x_{ni})$ ，所以

$$-\nabla E_{in}(w_t) = \frac{-1}{N} \sum_{n=1}^N h_t(y_n x_n) (-y_n x_n) = X^T *$$

$$\begin{bmatrix} \frac{1}{N} h_t(y_1 x_1)(y_1) \\ \frac{1}{N} h_t(y_2 x_2)(y_2) \\ \dots \\ \frac{1}{N} h_t(y_N x_N)(y_N) \end{bmatrix} = \hat{X}^T \cdot \begin{bmatrix} \frac{\frac{1}{N} h_t(y_1 x_1)(y_1)}{\text{sqrt}\left(\frac{1}{N} h_t(y_1 x_1) * h_t(-y_1 x_1)\right)} \\ \frac{\frac{1}{N} h_t(y_2 x_2)(y_2)}{\text{sqrt}\left(\frac{1}{N} h_t(y_2 x_2) * h_t(-y_2 x_2)\right)} \\ \dots \\ \frac{\frac{1}{N} h_t(y_N x_N)(y_N)}{\text{sqrt}\left(\frac{1}{N} h_t(y_N x_N) * h_t(-y_N x_N)\right)} \end{bmatrix} = \hat{X}^T .$$

$$\begin{bmatrix} \text{sqrt}\left(\frac{h_t(y_1 x_1)}{h_t(-y_1 x_1)}\right)(y_1) \\ \text{sqrt}\left(\frac{h_t(y_2 x_2)}{h_t(-y_2 x_2)}\right)(y_2) \\ \dots \\ \text{sqrt}\left(\frac{h_t(y_N x_N)}{h_t(-y_N x_N)}\right)(y_N) \end{bmatrix} \circ \text{因此} \hat{y} = \begin{bmatrix} \text{sqrt}\left(\frac{h_t(y_1 x_1)}{h_t(-y_1 x_1)}\right)(y_1) \\ \text{sqrt}\left(\frac{h_t(y_2 x_2)}{h_t(-y_2 x_2)}\right)(y_2) \\ \dots \\ \text{sqrt}\left(\frac{h_t(y_N x_N)}{h_t(-y_N x_N)}\right)(y_N) \end{bmatrix}$$