

1.

WEEK 4

第十二講: Nonlinear Transformation

Videos 58 min left

REQUIRED
 Quiz
 作業三
 40 min

GRADE
 100%

DUE
 Dec 24

2.

SGD:

$$err(W_t) \begin{cases} = 0 & (y_n W_t^T x_n > 0) \\ = -y_n W_t^T x_n & (y_n W_t^T x_n \leq 0) \end{cases}$$

$$\nabla err(W_t) \begin{cases} = 0 & (sign(W_t^T x_n) = y_n) \\ = -y_n x_n & (sign(W_t^T x_n) \neq y_n) \end{cases}$$

$$\begin{aligned} \therefore W_{t+1} &= W_t - \nabla err(W_t) \\ &= W_t - \mathbb{I}[sign(W_t^T x_n) \neq y_n](-y_n x_n) \\ &= W_t + \mathbb{I}[sign(W_t^T x_n) \neq y_n](y_n x_n) \end{aligned}$$

PLA:

$$W_{t+1} = W_t + \mathbb{I}[sign(W_t^T x_n) \neq y_n](y_n x_n)$$

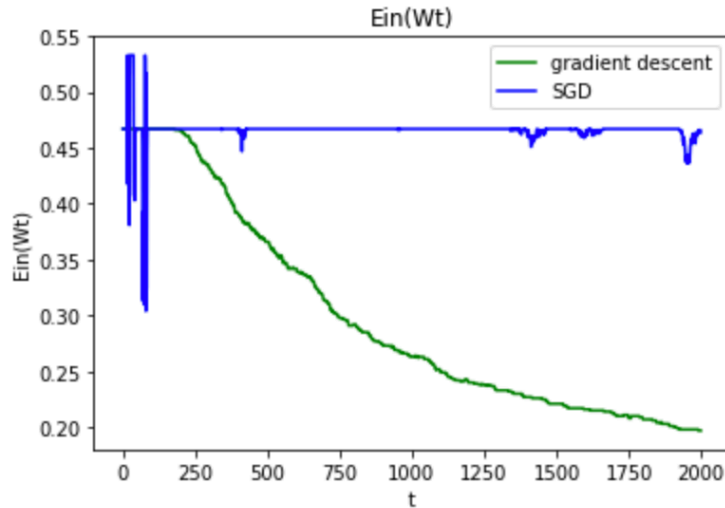
$$\therefore err(W_t) = \max(0, -y_n W_t^T x_n) \text{ results in PLA}$$

3.

$$\begin{aligned} \max_h \text{likelihood}(h) &\propto \prod_{n=1}^N h_{y_n}(x_n) \\ &= \max_h \ln \prod_{n=1}^N h_{y_n}(x_n) \\ &= \min_h \frac{1}{N} \sum_{n=1}^N -\ln(h_{y_n}(x_n)) \\ \therefore E_{in} &= \frac{1}{N} \sum_{n=1}^N -\ln(h_{y_n}(x_n)) \\ &= \frac{1}{N} \sum_{n=1}^N -\ln(\exp(W_{y_n}^T x_n) / (\sum_{k=1}^K \exp(W_k^T x_n))) \\ &= \frac{1}{N} \sum_{n=1}^N \left(\ln(\sum_{k=1}^K \exp(W_k^T x_n)) - \ln(\exp(W_{y_n}^T x_n)) \right) \\ &= \frac{1}{N} \sum_{n=1}^N (\ln(\sum_{k=1}^K \exp(W_k^T x_n)) - W_{y_n}^T x_n) \\ \therefore \frac{\partial E_{in}}{\partial W_i} &= \frac{1}{N} \sum_{n=1}^N \left[\left(x_n \cdot \exp(W_i^T x_n) / \left(\sum_{k=1}^K \exp(W_k^T x_n) \right) \right) - \mathbb{I}[y_n = i] x_n \right] \\ \therefore h_i(x_n) &= \exp(W_i^T x_n) / \left(\sum_{k=1}^K \exp(W_k^T x_n) \right) \end{aligned}$$

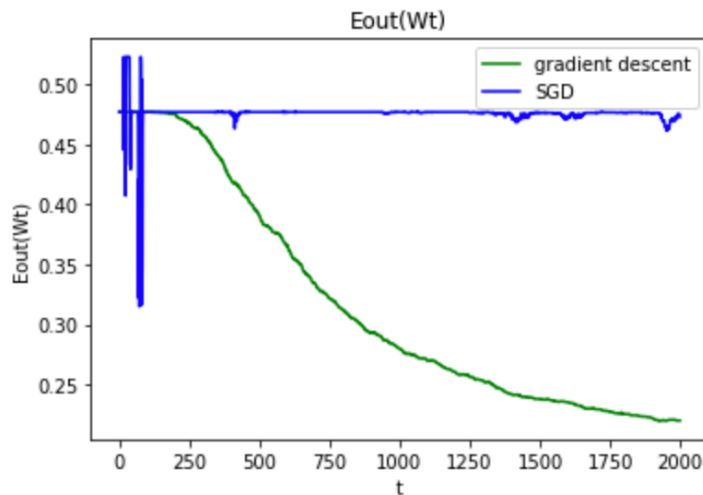
$$\begin{aligned}\therefore \frac{\partial E_{in}}{\partial W_i} &= \frac{1}{N} \sum_{n=1}^N [x_n \cdot h_i(x_n) - \mathbb{I}[y_n = i] x_n] \\ &= \frac{1}{N} \sum_{n=1}^N [(h_i(x_n) - \mathbb{I}[y_n = i]) x_n]\end{aligned}$$

4.



上圖中藍色和綠色的曲線分別為 gradient descent(GD)和 SGD 的 E_{in} ，GD 的 E_{in} 隨 t 逐步下降，但是 SGD 的 E_{in} 卻有明顯的上下震盪。SGD 每次僅利用一個 data point 的梯度更新參數，梯度方向不一定是 loss function 最小的方向，這讓 SGD 的梯度稍稍偏離「正軌」，導致 E_{in} 並不會隨著 t 穩定下降。因為 GD 和 SGD 的 learning rate 分別是 0.01 和 0.001，learning rate 不同導致他們最後的 E_{in} 也不同，較大的 learning rate 可以讓最後的 E_{in} 更小。

5.



上圖中藍色和綠色的曲線分別為 gradient descent(GD)和 SGD 的 E_{out} ，在 VC bound 的保證下，無論 SGD 還是 GD，與之對應的 E_{in} 和 E_{out} 都很相近。GD 的 E_{out} 隨 t 逐步下降，最後的 E_{out} 約為 0.22，略高於其 E_{in} 。SGD 的 E_{out} 有明顯的上下震盪，與其 E_{in} 的振幅相似。

6.

設定：

$$h_k(X) = \begin{bmatrix} h_k(x_1) \\ h_k(x_2) \\ \vdots \\ h_k(x_N) \end{bmatrix} \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$A = \begin{bmatrix} h_1(x_1) & \dots & h_K(x_1) \\ h_1(x_2) & \dots & h_K(x_2) \\ \vdots & & \vdots \\ h_1(x_N) & \dots & h_K(x_N) \end{bmatrix} = [h_1(X) \quad h_2(X) \quad \dots \quad h_K(X)]$$

$$H(X) = [h_1(X) \quad h_2(X) \quad \dots \quad h_K(X)] \cdot W = AW$$

由題意可知：

$$\textcircled{1} e_0^2 = \frac{1}{N} \|Y\|^2$$

$$= \frac{1}{N} Y^T Y$$

$$\textcircled{2} e_k^2 = \frac{1}{N} \|Y - h_k(X)\|^2$$

$$= \frac{1}{N} (Y^T Y + h_k^T(X) h_k(X) - 2h_k^T(X) Y)$$

$$\textcircled{3} RMSE^2(H) = E^2 = \frac{1}{N} \|Y - H(X)\|^2$$

$$= \frac{1}{N} (Y^T Y + H^T(X) H(X) - 2H^T(X) Y)$$

將①帶入②得：

$$Ne_k^2 = Ne_0^2 + h_k^T(X) h_k(X) - 2h_k^T(X) Y$$

$$h_k^T(X) Y = \frac{1}{2} (Ne_0^2 + h_k^T(X) h_k(X) - Ne_k^2) \dots\dots\dots \textcircled{4}$$

$$\therefore \textcircled{3} E^2 = \frac{1}{N} (Y^T Y + H^T(X) H(X) - 2H^T(X) Y)$$

$$= \frac{1}{N} \left(Y^T Y + H^T(X) H(X) - 2W^T \begin{bmatrix} \text{---} & h_1^T(X) & \text{---} \\ & \vdots & \\ \text{---} & h_K^T(X) & \text{---} \end{bmatrix} Y \right)$$

$$= \frac{1}{N} \left(Y^T Y + H^T(X) H(X) - 2W^T \begin{bmatrix} h_1^T(X) Y \\ \vdots \\ h_K^T(X) Y \end{bmatrix} \right)$$

代入④得：

$$E^2 = \frac{1}{N} \left(Y^T Y + H^T(X) H(X) - 2W^T \begin{bmatrix} \frac{1}{2} (Ne_0^2 + h_1^T(X) h_1(X) - Ne_1^2) \\ \vdots \\ \frac{1}{2} (Ne_0^2 + h_K^T(X) h_K(X) - Ne_K^2) \end{bmatrix} \right)$$

$$= \frac{1}{N} \left(Y^T Y + W^T A^T A W - W^T \begin{bmatrix} Ne_0^2 + h_1^T(X) h_1(X) - Ne_1^2 \\ \vdots \\ Ne_0^2 + h_K^T(X) h_K(X) - Ne_K^2 \end{bmatrix} \right)$$

$$\therefore \nabla E^2(W) = \frac{1}{N} \left(2A^T A W - \begin{bmatrix} N e_0^2 + h_1^T(X) h_1(X) - N e_1^2 \\ \vdots \\ N e_0^2 + h_K^T(X) h_K(X) - N e_K^2 \end{bmatrix} \right)$$

$$\therefore W_{min} = \underset{W}{\operatorname{argmin}} E^2 = \underset{W}{\operatorname{argmin}} E$$

$$\nabla E^2(W_{min}) = 0$$

$$\therefore 2A^T A W_{min} - \begin{bmatrix} N e_0^2 + h_1^T(X) h_1(X) - N e_1^2 \\ \vdots \\ N e_0^2 + h_K^T(X) h_K(X) - N e_K^2 \end{bmatrix} = 0$$

$$W_{min} = \frac{1}{2} (A^T A)^{-1} \begin{bmatrix} N e_0^2 + h_1^T(X) h_1(X) - N e_1^2 \\ \vdots \\ N e_0^2 + h_K^T(X) h_K(X) - N e_K^2 \end{bmatrix}$$

$$= \frac{1}{2} \left(\begin{bmatrix} \text{---} & h_1^T(X) & \text{---} \\ & \vdots & \\ \text{---} & h_K^T(X) & \text{---} \end{bmatrix} \begin{bmatrix} h_1(X) & \dots & h_K(X) \end{bmatrix} \right)^{-1} \begin{bmatrix} N e_0^2 + h_1^T(X) h_1(X) - N e_1^2 \\ \vdots \\ N e_0^2 + h_K^T(X) h_K(X) - N e_K^2 \end{bmatrix}$$

$$= \frac{1}{2} \left(\begin{bmatrix} h_1^T(X) h_1(X) & h_1^T(X) h_2(X) & \dots & h_1^T(X) h_K(X) \\ h_2^T(X) h_1(X) & h_2^T(X) h_2(X) & \dots & h_2^T(X) h_K(X) \\ \vdots & \vdots & \dots & \vdots \\ h_K^T(X) h_1(X) & h_K^T(X) h_2(X) & \dots & h_K^T(X) h_K(X) \end{bmatrix} \right)^{-1} \left(\begin{bmatrix} h_1^T(X) h_1(X) \\ h_2^T(X) h_2(X) \\ \vdots \\ h_K^T(X) h_K(X) \end{bmatrix} + N \begin{bmatrix} e_0^2 - e_1^2 \\ e_0^2 - e_2^2 \\ \vdots \\ e_0^2 - e_K^2 \end{bmatrix} \right)$$