

2.

分以下兩種情形來討論：

(i) $\text{err} = 0$: 不用做更新

(ii) $\text{err} \neq 0$:

其梯度為 $\nabla \text{err}(\mathbf{w}) = -y\mathbf{x}$ 。故若是今天有一筆資料 (x_n, y_n) 要拿其去做 SGD 進行修正，結果會是如下：

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta(-y_n \mathbf{x}_n)$$

而若是做 PLA 的修正，結果會是：

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_n \mathbf{x}_n$$

可以看出兩者修正的方向都是相同的。

3.

首先，由 16 題可知：

$$E_{in} = \frac{1}{N} \sum_{n=1}^N (\ln(\sum_{k=1}^K e^{w_k^T x_n}) - w_{y_n}^T x_n)$$

對 sigma 內前項做微分，可得：(由鏈鎖律)

$$\frac{1}{\sum_{k=1}^K e^{w_k^T x_n}} (x_n e^{w_k^T x_n}) = h_i(x_n)$$

對 sigma 內後項做微分：

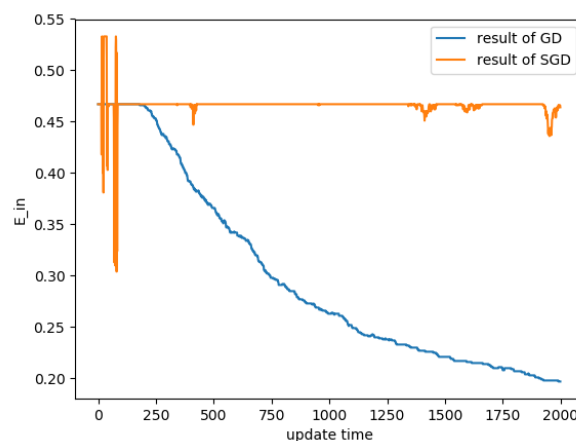
(i) 若 $w_{y_n} = w_i$ ，對 w_i 微分則為 x_n

(ii) 若 $w_{y_n} \neq w_i$ ，對 w_i 微分則為 0

結合以上可得整個 E_{in} 對 w_i 的微分即為：

$$\frac{1}{N} \sum_{n=1}^N (h_i(x_n) - [y_n = i] x_n)$$

4.

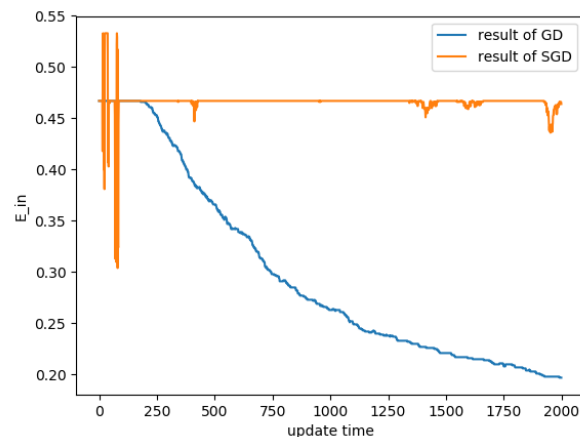


(i) 使用 stochastic gradient descent(SGD)在一開始做更新時，變動幅度很大。而

且在做 2000 次以內的更新對於降低 E_{in} 的效果相當有限， E_{in} 大約都落在 0.46 附近。

(ii) 使用 gradient descent(GD)能夠穩定使 E_{in} 下降，可以可降低製 0.2 附近。

5.



(i) E_{out} 的趨勢和 E_{in} 相同，故不再贅述。

(ii) 大體而言， E_{out} 會大於 E_{in} 。

*註：兩題程式我用同一個 code(45.py)來跑，並製出兩張圖。一張是 E_{in} 對更新次數，另一張是 E_{out} 對更新次數。

6.

令 $Y = [y_1, y_2, \dots, y_n]^T$, $\bar{h}_i = [h_i(x_1), h_i(x_2), \dots, h_i(x_n)]^T$

故原本的函式可以表成： $(Y - \bar{h}_i)^2 = N e_i^2$

展開可得 $Y^T Y - 2\bar{h}_i^T Y + \bar{h}_i^2 = N e_i^2$ 。又因 $(Y - \bar{h}_0)^2 = Y^T Y = N e_0^2$ ，故可得出：

$$\bar{h}_i^T Y = N(e_0^2 - e_i^2) + \bar{h}_i^2$$

再令：(便於表示)

$$H_j' = \begin{bmatrix} h_1(x_j) \\ h_2(x_j) \\ \vdots \\ h_k(x_j) \end{bmatrix}, H' = \begin{bmatrix} h_1(x_1) & h_1(x_2) & \dots & h_1(x_n) \\ h_2(x_1) & h_2(x_2) & \dots & h_2(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ h_k(x_1) & h_k(x_2) & \dots & h_k(x_n) \end{bmatrix}, W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}$$

而題目所求為 $RMSE(H)$ 的最小值，其極值發生點會和 $(RMSE(H))^2$ 相同，即要找以下式子：

$$\min_w \frac{1}{N} \sum_{j=1}^N (y_j - W^T H_j')^2$$

將其對 w_1 做偏微可得：

$$\frac{1}{N} \sum_{j=1}^N -2h_i(x_j) (y_j - w_i h_i(x_j)) = \frac{2}{N} (\bar{h}_i^T \bar{h}_i w_i - \bar{h}_i^T Y)$$

把全部的 w_i 和在一起，可得其梯度：

$$\nabla(\text{RMSE}(H))^2 = \frac{2}{N} (H' H'^T W - H' Y) \triangleq 0$$

於是可得：

$$W = (H' H'^T)^{-1} H' Y$$

又因為前面曾算出：

$$\bar{h}_i^T Y = N(e_0^2 - e_i^2) + \bar{h}_i^2$$

故可得：

$$H' Y = \begin{bmatrix} N(e_0^2 - e_1^2) + \bar{h}_1^2 \\ N(e_0^2 - e_2^2) + \bar{h}_2^2 \\ \vdots \\ N(e_0^2 - e_n^2) + \bar{h}_n^2 \end{bmatrix}$$

結合以上可得：

$$W = (H' H'^T)^{-1} \begin{bmatrix} N(e_0^2 - e_1^2) + \bar{h}_1^2 \\ N(e_0^2 - e_2^2) + \bar{h}_2^2 \\ \vdots \\ N(e_0^2 - e_n^2) + \bar{h}_n^2 \end{bmatrix}, \text{其中 } H' = \begin{bmatrix} h_1(x_1) & h_1(x_2) & \dots & h_1(x_n) \\ h_2(x_1) & h_2(x_2) & \dots & h_2(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ h_k(x_1) & h_k(x_2) & \dots & h_k(x_n) \end{bmatrix}$$