

MLFoundation HW3

r07922100 楊力權

1.

此課程: 機器學習基石下 (Machine Learning Foundations)--Algorithmic Foundations



2.

SGD : $w_{t+1} \leftarrow w_t - \eta \nabla E_{in}$ PLA : $w_{t+1} \leftarrow w_t + (y_n x_n) [[y_n \neq \text{sign}(w_t^T x_n)]]$
 $\because \text{err}(w) = \max(0, -y w^T x) \therefore \nabla E_{in} = [[y w^T x < 0]](-yx) = [[\text{sign}(w^T x) \neq y]](-yx)$

SGD每次只使用一筆資料(x_n, y_n)更新 w 。

SGD : $w_{t+1} \leftarrow w_t - \eta [[\text{sign}(w^T x_n) \neq y_n]](-y_n x_n) = w_t + \eta [[\text{sign}(w^T x_n) \neq y_n]](y_n x_n)$

若learning rate $\eta = 1$ 則SGD在使用 $\text{err}(w) = \max(0, -y w^T x)$ 下等同於PLA。

3.

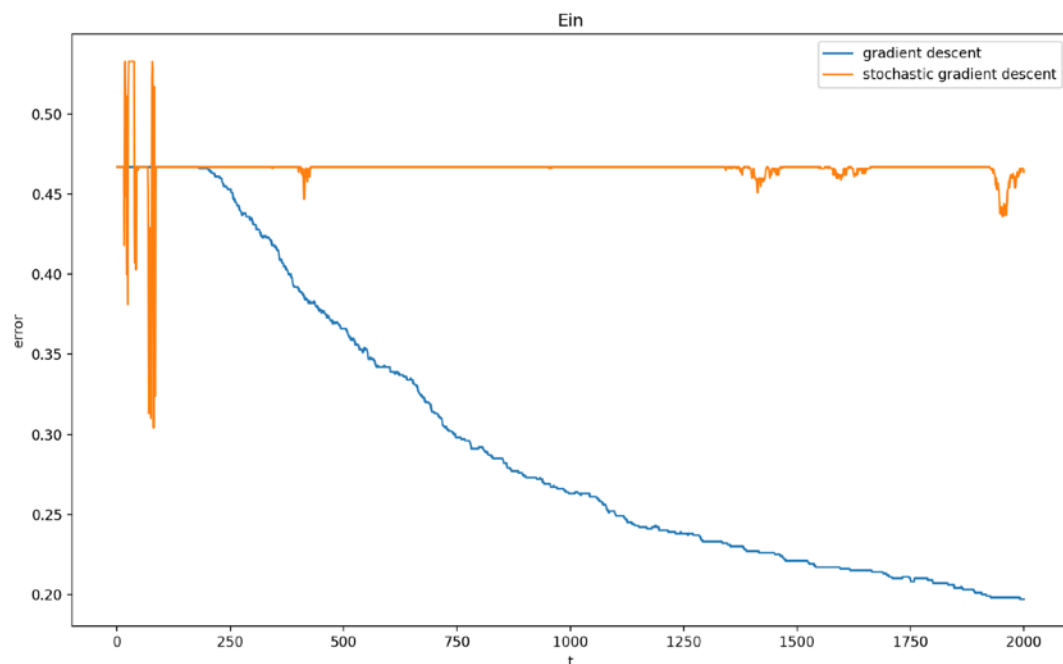
面對N筆資料，最大化likelihood，也就是最小化negative likelihood

$$\begin{aligned} \arg \max_w \text{likelihood}(w) &\propto \prod_{n=1}^N h_{y_n}(y_n w^T x_n) \Rightarrow \arg \max_w \ln \prod_{n=1}^N h_{y_n}(y_n w^T x_n) = \arg \min_w \sum_{n=1}^N -\ln \frac{e^{(w_{y_n}^T x_n)}}{\sum_{k=1}^K e^{(w_k^T x_n)}} \\ &= \arg \min_w \sum_{n=1}^N \ln \frac{\sum_{k=1}^K e^{(w_k^T x_n)}}{e^{(w_{y_n}^T x_n)}} = \arg \min_w \sum_{n=1}^N (\ln(\sum_{k=1}^K e^{(w_k^T x_n)}) - (w_{y_n}^T x_n)) \\ \text{error} &= \sum_{n=1}^N (\ln(\sum_{k=1}^K e^{(w_k^T x_n)}) - (w_{y_n}^T x_n)) \Rightarrow E_{in} = \frac{\text{error}}{N} = \frac{1}{N} \sum_{n=1}^N (\ln(\sum_{k=1}^K e^{(w_k^T x_n)}) - (w_{y_n}^T x_n)) \end{aligned}$$

使用連鎖率求 E_{in} 對 w_i 的微分

$$\begin{aligned} \frac{\partial E_{in}}{\partial w_i} &= \frac{1}{N} \sum_{n=1}^N \left(\frac{\partial}{\partial w_i} \ln(\sum_{k=1}^K e^{(w_k^T x_n)}) - \frac{\partial}{\partial w_i} (w_{y_n}^T x_n) \right) = \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{\sum_{k=1}^K e^{(w_k^T x_n)}} \frac{\partial}{\partial w_i} (\sum_{k=1}^K e^{(w_k^T x_n)}) - [[y_n = i]] x_n \right) \\ &= \frac{1}{N} \sum_{n=1}^N \left(\frac{e^{(w_i^T x_n)}}{\sum_{k=1}^K e^{(w_k^T x_n)}} x_n - [[y_n = i]] x_n \right) = \frac{1}{N} \sum_{n=1}^N (h_i(x_n) x_n - [[y_n = i]] x_n) = \frac{1}{N} \sum_{n=1}^N (h_i(x_n) - [[y_n = i]]) x_n \end{aligned}$$

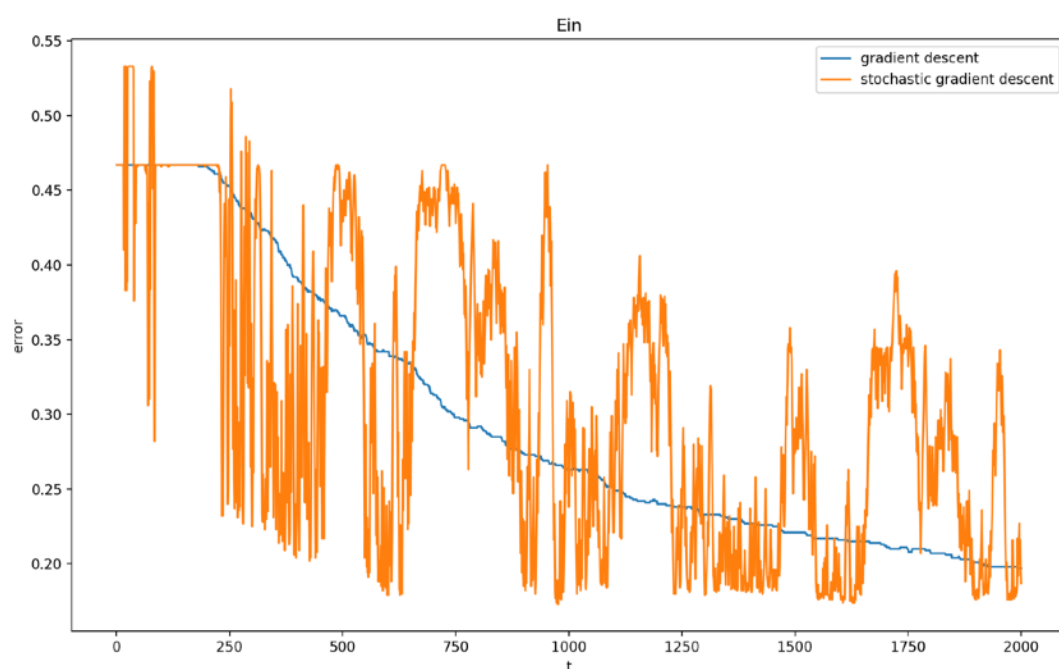
4.



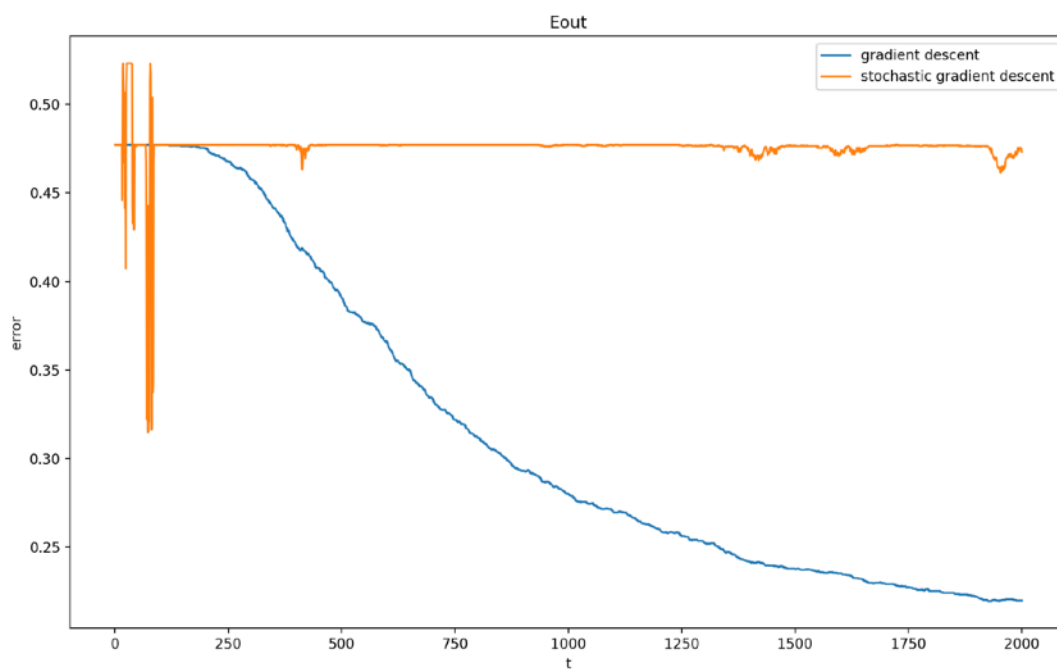
GD使用 $lr=0.01$ ，而SGD使用 $lr=0.001$ 。

gradient descent每個iteration使用所有的數據來計算平均梯度方向，因此訓練穩定，而SGD跳動劇烈且最後幾乎沒有降低Ein。

我認為這應該是出題時的問題，因為當SGD也使用 $lr=0.01$ 時，產生下圖，GD與SGD的Ein皆有下降，而明顯不同之處是SGD非常不穩定的震盪，因為每筆數據有不同程度雜訊的緣故，GD能夠平均掉雜訊對梯度帶來的極端影響，而SGD只使用一筆數據對 W 的梯度。但是即便如此兩者都能成功降低Ein。

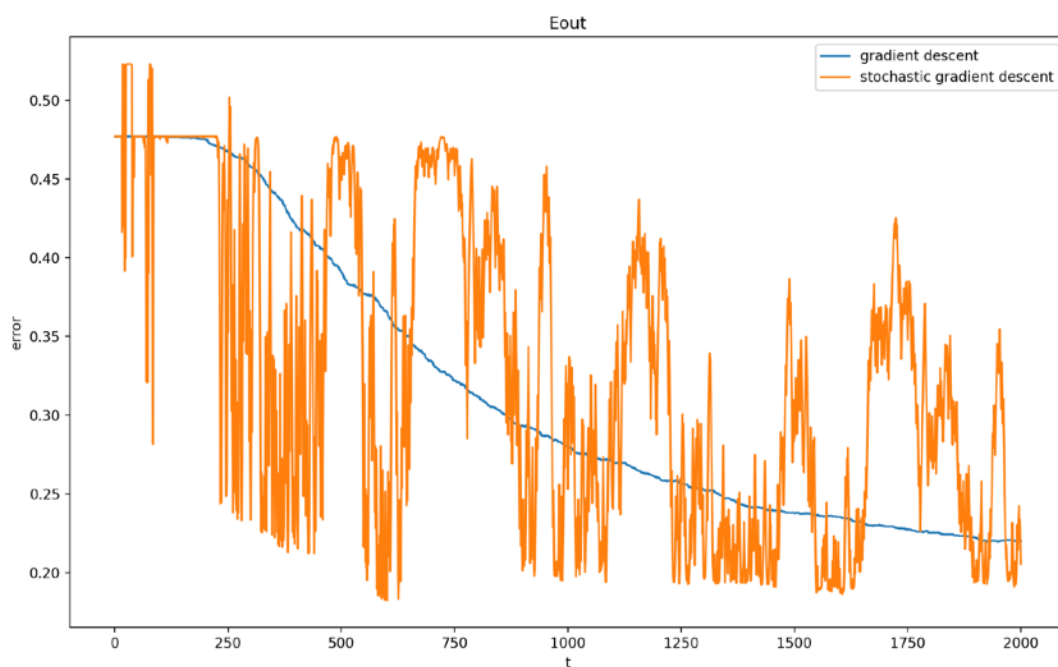


5.



上圖為GD $\text{lr}=0.01$ ，SGD $\text{lr}=0.001$ ；下圖為GD $\text{lr}=0.01$ ，SGD $\text{lr}=0.01$ 。

結果與第4題非常接近，因此觀察討論與第4.題一樣。而Eout比Ein高出一點點但差距不多，可以說明20維的資料在資料夠多時，可以做到Ein與Eout差不多的結果。



6.

$$\text{設 } x = [d_1, d_2, \dots] \in R^d, X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \in R^{N \times d}, h = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \end{bmatrix} \in R^d, H = [h_1, h_2, \dots] \in R^{d \times K}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \end{bmatrix} \in R^K, Y \in R^N, XHW \in R^N$$

$$\text{RMSE}(H) = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - \sum_{k=1}^K w_k h_k(x_n))^2} = \sqrt{\frac{1}{N} \|Y - XHW\|^2}$$

$$\text{Optimal} : \nabla E_{in} = 0 \Rightarrow (XH)^T(XH)W = (XH)^TY \Rightarrow W = ((XH)^T(XH))^{-1}(XH)^TY$$

$$Y \text{ 未知，但我們知道每個 } e_k = \text{RMSE}(h_k) = \sqrt{\frac{1}{N} \|Y - Xh_k\|^2} = \sqrt{\frac{1}{N} (Y - Xh_k)^T(Y - Xh_k)}$$

$$\Rightarrow Ne_k^2 = Y^TY + (Xh_k)^T(Xh_k) - 2(Xh_k)^TY$$

$$\Rightarrow N(e_k^2 - e_m^2) = 2(h_m^T - h_k^T)X^TY + (Xh_k)^T(Xh_k) - (Xh_m)^T(Xh_m)$$

且我們知道 $h_0(x) = 0$ ，因此把 m 用 0 代入

$$\Rightarrow N(e_k^2 - e_0^2) = -2h_k^TX^TY + (Xh_k)^T(Xh_k) \Rightarrow (Xh_k)^TY = \frac{(Xh_k)^T(Xh_k) - N(e_k^2 - e_0^2)}{2}$$

我們可以做一個矩陣如下 (Y 左方的矩陣必須 $\text{Rank}=N$ 否則 $\det=0$ 無反矩陣，因此若 Xh_k 與 $\{Xh_m | m < k\}$ 平行則必須跳過 k)

$$\begin{bmatrix} (Xh_1)^T \\ (Xh_2)^T \\ \vdots \\ (Xh_N)^T \end{bmatrix} Y = \begin{bmatrix} \frac{(Xh_1)^T(Xh_1) - N(e_1^2 - e_0^2)}{2} \\ \frac{(Xh_2)^T(Xh_2) - N(e_2^2 - e_0^2)}{2} \\ \vdots \\ \frac{(Xh_N)^T(Xh_N) - N(e_N^2 - e_0^2)}{2} \end{bmatrix} \Rightarrow Y = \begin{bmatrix} (Xh_1)^T \\ (Xh_2)^T \\ \vdots \\ (Xh_N)^T \end{bmatrix}^{-1} \begin{bmatrix} \frac{(Xh_1)^T(Xh_1) - N(e_1^2 - e_0^2)}{2} \\ \frac{(Xh_2)^T(Xh_2) - N(e_2^2 - e_0^2)}{2} \\ \vdots \\ \frac{(Xh_N)^T(Xh_N) - N(e_N^2 - e_0^2)}{2} \end{bmatrix}$$

把 Y 代回

$$W = ((XH)^T(XH))^{-1}(XH)^TY = ((XH)^T(XH))^{-1}(XH)^T \begin{bmatrix} (Xh_1)^T \\ (Xh_2)^T \\ \vdots \\ (Xh_N)^T \end{bmatrix}^{-1} \begin{bmatrix} \frac{(Xh_1)^T(Xh_1) - N(e_1^2 - e_0^2)}{2} \\ \frac{(Xh_2)^T(Xh_2) - N(e_2^2 - e_0^2)}{2} \\ \vdots \\ \frac{(Xh_N)^T(Xh_N) - N(e_N^2 - e_0^2)}{2} \end{bmatrix}$$

$$\text{其中 } X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \in R^{N \times d}, H = [h_1, h_2, \dots] \in R^{d \times K}$$