

Machine Learning Foundations

Homework #3

1.

作業三
評分測驗 • 40 min

✓ 恭喜！您通過了！
通過條件：75% 或更高

堅持學習

成績
100%

作業三

最新提交作業的評分
100%

2.

當 $w^T x$ 與 y 同號時，代表 y 點分類正確，故不修正，梯度為 0。

當 $w^T x$ 與 y 異號時，則代表 y 點分類錯誤，需作修正，故對 $-y w^T x$ 作微分。

$$\text{即 } \frac{\partial -y w^T x}{\partial w} = -yx \Rightarrow w_{t+1} = w_t + \eta(-yx)$$

$\therefore \text{err}(w) = \max(0, -y w^T x)$ results in PLA

3.

一元的二階泰勒展開式：

$$f(x + \Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$$

當 Δx 趨近於 0，上式等價於：

$$\begin{aligned} f'(x) + f''(x)\Delta x &= 0 \\ \Rightarrow \Delta x &= -\frac{f'(x_n)}{f''(x_n)} \end{aligned}$$

推廣到二元，

$$(\Delta u, \Delta v) = -(\nabla^2 E(u, v))^{-1} \nabla E(u, v)$$

4.

根據課程投影片 10 中的第 11 頁，

$$\text{cross entropy error} = \min_w \left\{ \frac{1}{N} \sum_{n=1}^N -\ln \left(\theta(y_n w_y^T x_n) \right) \right\}$$

又根據課程投影片 10 中的第 10 頁， $h_y(x_n) = \theta(y_n w_y^T x_n)$

已知本題的 $h_y(x_n) = \frac{\exp(w_y^T x)}{\sum_{i=1}^K \exp(w_i^T x)}$ ，將此式代入 cross entropy error 的 $\theta(y_n w_y^T x_n)$

$$\Rightarrow \text{cross entropy error} = \frac{1}{N} \sum_{n=1}^N \left(-\ln \left(\frac{\exp(w_y^T x)}{\sum_{i=1}^K \exp(w_i^T x)} \right) \right)$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{n=1}^N \left(-\ln \left(\exp(w_{y_n}^T x_n) \right) + \ln \left(\sum_{i=1}^K \exp(w_i^T x_n) \right) \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\left(\ln \sum_{i=1}^K \exp(w_i^T x_n) \right) - w_{y_n}^T x_n \right)
\end{aligned}$$

5.

根據課程投影片9中的第7頁，

$$\begin{aligned}
E_{in} &= \frac{1}{N+K} \left(\sum_{n=1}^N (y_n - w^T x_n)^2 + \sum_{k=1}^K (\tilde{y}_k - w^T \tilde{x}_k)^2 \right) \\
&= \frac{1}{N+K} (\|Xw - y\|^2 + \|\tilde{X}w - \tilde{y}\|^2) \\
&= \frac{1}{N+K} \left((w^T X^T X w - 2w^T X^T y + y^T y) + (w^T \tilde{X}^T \tilde{X} w - 2w^T \tilde{X}^T \tilde{y} + \tilde{y}^T \tilde{y}) \right) \\
\nabla E_{in} &= \frac{1}{N+K} (2X^T X w - 2X^T y + 2\tilde{X}^T \tilde{X} w - 2\tilde{X}^T \tilde{y}) \\
&= \frac{2}{N+K} (X^T X w - X^T y + \tilde{X}^T \tilde{X} w - \tilde{X}^T \tilde{y}) \\
\nabla E_{in} &= 0 \\
&\Rightarrow \frac{2}{N+K} (X^T X w - X^T y + \tilde{X}^T \tilde{X} w - \tilde{X}^T \tilde{y}) = 0 \\
&\Rightarrow (X^T X + \tilde{X}^T \tilde{X}) w = X^T y + \tilde{X}^T \tilde{y} \\
&\Rightarrow w = (X^T X + \tilde{X}^T \tilde{X})^{-1} (X^T y + \tilde{X}^T \tilde{y})
\end{aligned}$$

6.

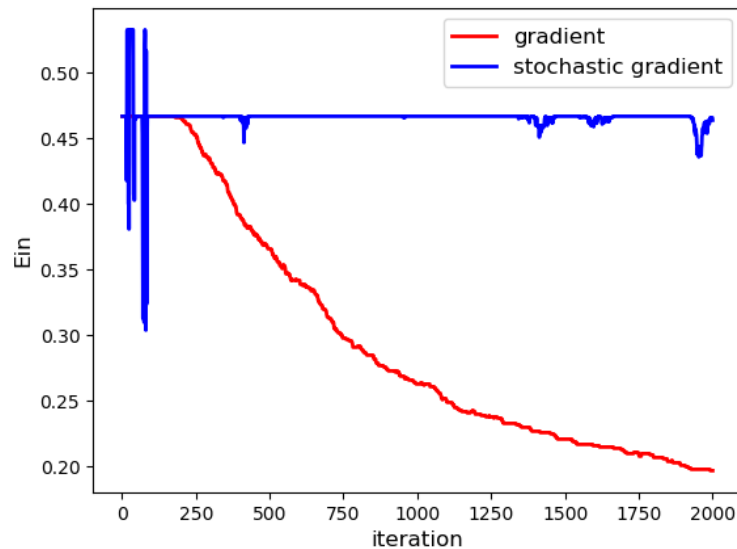
本題的w和上一題的w等價，

所以將 $\tilde{X} = \sqrt{\lambda}I, \tilde{y} = 0$ 代入上一題的 $w = (X^T X + \tilde{X}^T \tilde{X})^{-1} (X^T y + \tilde{X}^T \tilde{y})$

$$\Rightarrow w = (X^T X + \sqrt{\lambda}I \sqrt{\lambda}I)^{-1} (X^T y + (\sqrt{\lambda}I)^T 0) = (X^T X + \lambda I)^{-1} X^T y$$

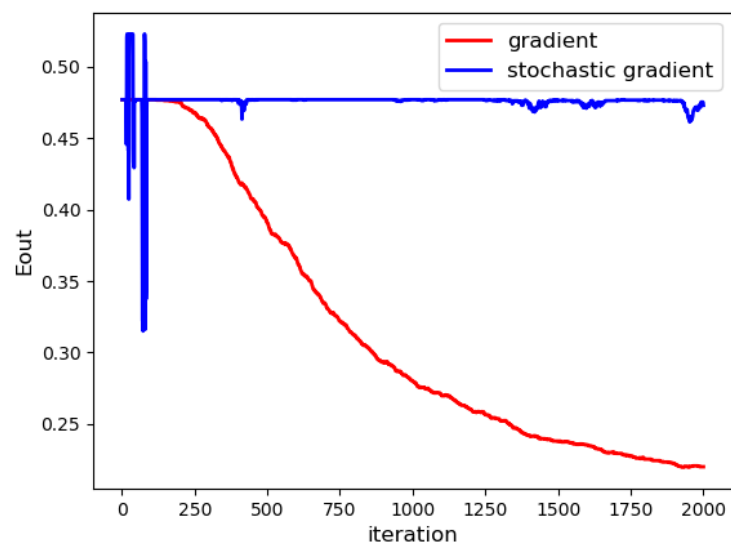
與課程投影片14的第10頁中Augmented Error的optimal solution相符

7.



根據上圖的結果，可以發現用gradient descent去更新 w 的話， E_{in} 確實會隨著越來越小。然而，用stochastic gradient descent更新 w 在前幾次讓 E_{in} 十分不穩，一直到更新2000次的時候 E_{in} 都沒有下降。我認為stochastic gradient的效果會不如預期可能是因為learning rate較小，且iteration的次數只有2000次，或許對stochastic gradient來說不夠多次，可能更新次數要再多一點才會有較明顯的效果。

8.



根據上圖的結果，可以發現跟上一題的 E_{in} 相比，無論是gradient descent或是stochastic gradient descent所造成的 E_{out} 走向都和上一題差不多，唯獨 E_{out} 都會比 E_{in} 高一點點。其實這樣的效果蠻正常的，因為是使用training data去計算gradient，所以同樣用training data去算出Error的話本來就會比較低一點。

9.

(a)

已知 $X^T X w_{lin} = X^T y$,

又已知 $X = U \Gamma V^T$, $w_{lin} = V \Gamma^{-1} U^T y$, 代入左式

$$\text{左式} = (U \Gamma V^T)^T (U \Gamma V^T) V \Gamma^{-1} U^T y = V \Gamma^T U^T U \Gamma V^T V \Gamma^{-1} U^T y$$

又已知 $U^T U = I$, $V^T V = I$,

$$\therefore \text{左式} = V \Gamma^T \Gamma \Gamma^{-1} U^T y = V \Gamma^T U^T y = \text{右式}$$

故得證

(b)

若 $X^T X w = X^T y$, 則 $\|X w - y\|$ 會是最小值

$$\Leftrightarrow X w = \text{Proj}_{R(X)} y$$

$$\Leftrightarrow (y - X w) \perp R(X)$$

$$\Leftrightarrow \langle y - X w, X b \rangle = 0, \forall b : (d+1) \times 1$$

$$\Leftrightarrow (X b)^T (y - X w) = 0, \forall b : (d+1) \times 1$$

$$\Leftrightarrow b^T X^T y - b^T X^T X w = 0, \forall b : (d+1) \times 1$$

$$\Leftrightarrow b^T (X^T y - X^T X w) = 0, \forall b : (d+1) \times 1$$

$$\Leftrightarrow \langle X^T y - X^T X w, b \rangle = 0, \forall b : (d+1) \times 1$$

$$\Leftrightarrow X^T y - X^T X w = 0$$

$$\Leftrightarrow X^T X w = X^T y$$