


Machine Learning Foundation (NTU, Fall 2018) Homework #3

陳熙 R07922151

1.



此课程: 機器學習基石下 (Machine Learning Foundations)---Algorithmic Foundati...

测验

作業三

20 个问题

您的分数

100.00%

我们会保留您的最高分数。

[查看最新提交内容](#)

再次参加

2. 在 SGD 使用這個 error function 即，

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \nabla \text{err}(\mathbf{w}_t) \\ &= \mathbf{w}_t - \nabla \max(0, -y\mathbf{w}_t^T \mathbf{x}) \\ &= \begin{cases} \mathbf{w}_t, y\mathbf{w}_t^T \mathbf{x} \geq 0 \\ \mathbf{w}_t - y\mathbf{x}, y\mathbf{w}_t^T \mathbf{x} < 0 \end{cases}\end{aligned}$$

其中， $y\mathbf{w}_t^T \mathbf{x} \geq 0$ 表示 y 和 $\mathbf{w}_t^T \mathbf{x}$ 同號，即 $\text{sign}(\mathbf{w}_t^T \mathbf{x}) = y$ ，此時 \mathbf{w} 不變。當

$\text{sign}(\mathbf{w}_t^T \mathbf{x}) \neq y$ 時， $\mathbf{w}_{t+1} = \mathbf{w}_t - y\mathbf{x}$ 。因此，這一算法和 PLA 的結果相同。

3.

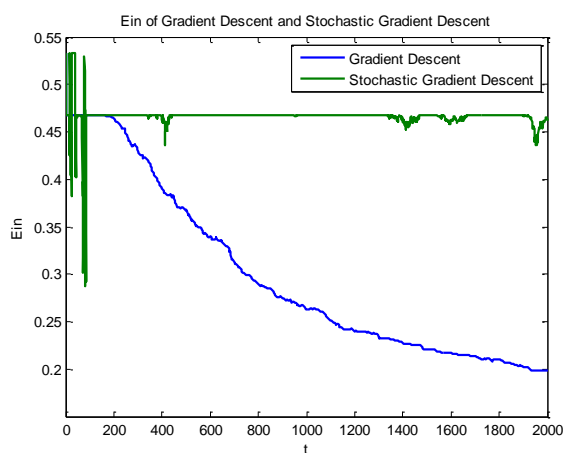
$$\begin{aligned}E_{in} &= \frac{1}{N} \sum_{n=1}^N -\ln(h_{y_n}(\mathbf{x}_n)) \\ &= \frac{1}{N} \sum_{n=1}^N -\ln \left(\frac{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)} \right)\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)}{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)} \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\ln \left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n) \right) - \ln(\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)) \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\ln \left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n) \right) - \mathbf{w}_{y_n}^T \mathbf{x}_n \right)
\end{aligned}$$

以上可以知道， $E_{in} = \frac{1}{N} \sum_{n=1}^N \left(\ln \left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n) \right) - \mathbf{w}_{y_n}^T \mathbf{x}_n \right)$ 。

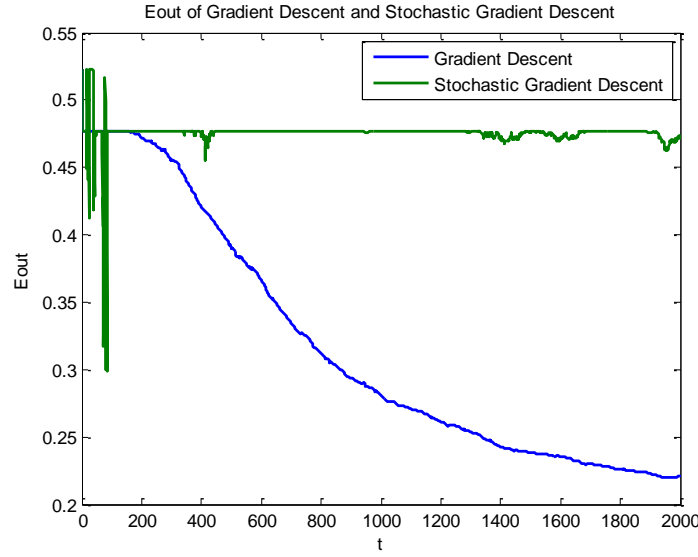
$$\begin{aligned}
\frac{\partial E_{in}}{\partial w_i} &= \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)} \frac{\partial \sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)}{\partial w_i} - y_n = i \mathbf{x}_n \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\frac{\exp(\mathbf{w}_i^T \mathbf{x}_n)}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)} \frac{\partial (\mathbf{w}_i^T \mathbf{x}_n)}{\partial w_i} - y_n = i \mathbf{x}_n \right) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\frac{\exp(\mathbf{w}_i^T \mathbf{x}_n)}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n)} \mathbf{x}_n - y_n = i \mathbf{x}_n \right) \\
&= \frac{1}{N} \sum_{n=1}^N ((h_i(\mathbf{x}_n) - y_n = i) \mathbf{x}_n)
\end{aligned}$$

4. 通過圖形可以發現，GD 得到的 E_{in} 逐漸下降，而 SGD 不是十分穩定，也可能是 η 設置過小，或是 T 不夠大，使其在很小的區間內波動。



5. 通過觀察可以發現， E_{out} 圖形的趨勢和 E_{in} 圖形相似，說明這個參數下的 GD

此時表現得更好，更為穩定。



6. 根據題意，令 $\mathbf{y} = [y_1, \dots, y_n]^T$ ， $\mathbf{w} = [w_1, \dots, w_k]^T$ ， $\mathbf{h}(\mathbf{x}_n) = [h_1(\mathbf{x}_n), \dots, h_k(\mathbf{x}_n)]$ ，

$\mathbf{h}(X) = [\mathbf{h}(\mathbf{x}_1), \dots, \mathbf{h}(\mathbf{x}_n)]^T$ ，有

$$\begin{aligned}
 \text{RMSE}(H) &= \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - H(\mathbf{x}))^2} \\
 &= \sqrt{\frac{1}{N} \sum_{n=1}^N \left(\sum_{k=1}^K w_k h_k(\mathbf{x}_n) - y_n \right)^2} \\
 &= \sqrt{\frac{1}{N} \sum_{n=1}^N (\mathbf{h}^T(\mathbf{x}_n) \mathbf{w} - y_n)^2} \\
 &= \sqrt{\frac{1}{N} (\mathbf{h}(X) \mathbf{w} - \mathbf{y})^T (\mathbf{h}(X) \mathbf{w} - \mathbf{y})} \\
 &= \sqrt{\frac{1}{N} (\mathbf{w}^T \mathbf{h}^T(X) \mathbf{h}(X) \mathbf{w} - \mathbf{w}^T \mathbf{h}^T(X) \mathbf{y} - \mathbf{y}^T \mathbf{h}(X) \mathbf{w} + \mathbf{y}^T \mathbf{y})}
 \end{aligned}$$

令 $f(\mathbf{w}) = \mathbf{w}^T \mathbf{h}^T(X) \mathbf{h}(X) \mathbf{w} - \mathbf{w}^T \mathbf{h}^T(X) \mathbf{y} - \mathbf{y}^T \mathbf{h}(X) \mathbf{w}$ ，則

$$\frac{df(\mathbf{w})}{d\mathbf{w}} = 2(\mathbf{h}^T(X) \mathbf{h}(X) \mathbf{w} - \mathbf{h}^T(X) \mathbf{y})，\text{ 令 } \frac{df(\mathbf{w})}{d\mathbf{w}} = 0，\text{ 有}$$

$$\mathbf{w} = (\mathbf{h}^T(X) \mathbf{h}(X))^{-1} \mathbf{h}^T(X) \mathbf{y}$$

$$e_k^2 = \frac{1}{N} \sum_{n=1}^N (y_n - h_k(\mathbf{x}))^2 = \frac{1}{N} \|\mathbf{y} - \mathbf{h}_k(X)\|^2 = \frac{1}{N} (\mathbf{y}^T \mathbf{y} + \mathbf{h}_k^T(X) \mathbf{h}_k(X) - 2\mathbf{h}_k^T(X) \mathbf{y})$$

$e_k^2 = \frac{1}{N} \left(e_0 + \mathbf{h}_k^T(X) \mathbf{h}_k(X) - 2\mathbf{h}_k^T(X) \mathbf{y} \right)$ ，則

$$\mathbf{h}_k^T(X) \mathbf{y} = \frac{1}{2} \left(e_0^2 + \mathbf{h}_k^T(X) \mathbf{h}_k(X) - N e_k^2 \right)$$

令 $\mathbf{e} = [e_1^2, \dots, e_K^2]^T$ ， $\mathbf{e}_0 = [e_0^2, \dots, e_0^2]^T$ ， $\tilde{\mathbf{h}}(X) = [\mathbf{h}_1^T(X) \mathbf{h}_1(X), \dots, \mathbf{h}_K^T(X) \mathbf{h}_K(X)]^T$

則， $\mathbf{h}^T(X) \mathbf{y} = \frac{1}{2} (\mathbf{e}_0^2 + \tilde{\mathbf{h}}(X) - N \mathbf{e})$ 。

代入 $\mathbf{w} = (\mathbf{h}^T(X) \mathbf{h}(X))^{-1} \mathbf{h}^T(X) \mathbf{y}$ 中，有

$$\mathbf{w} = \frac{1}{2} (\mathbf{h}^T(X) \mathbf{h}(X))^{-1} (\mathbf{e}_0^2 - N \mathbf{e} + \tilde{\mathbf{h}}(X))$$