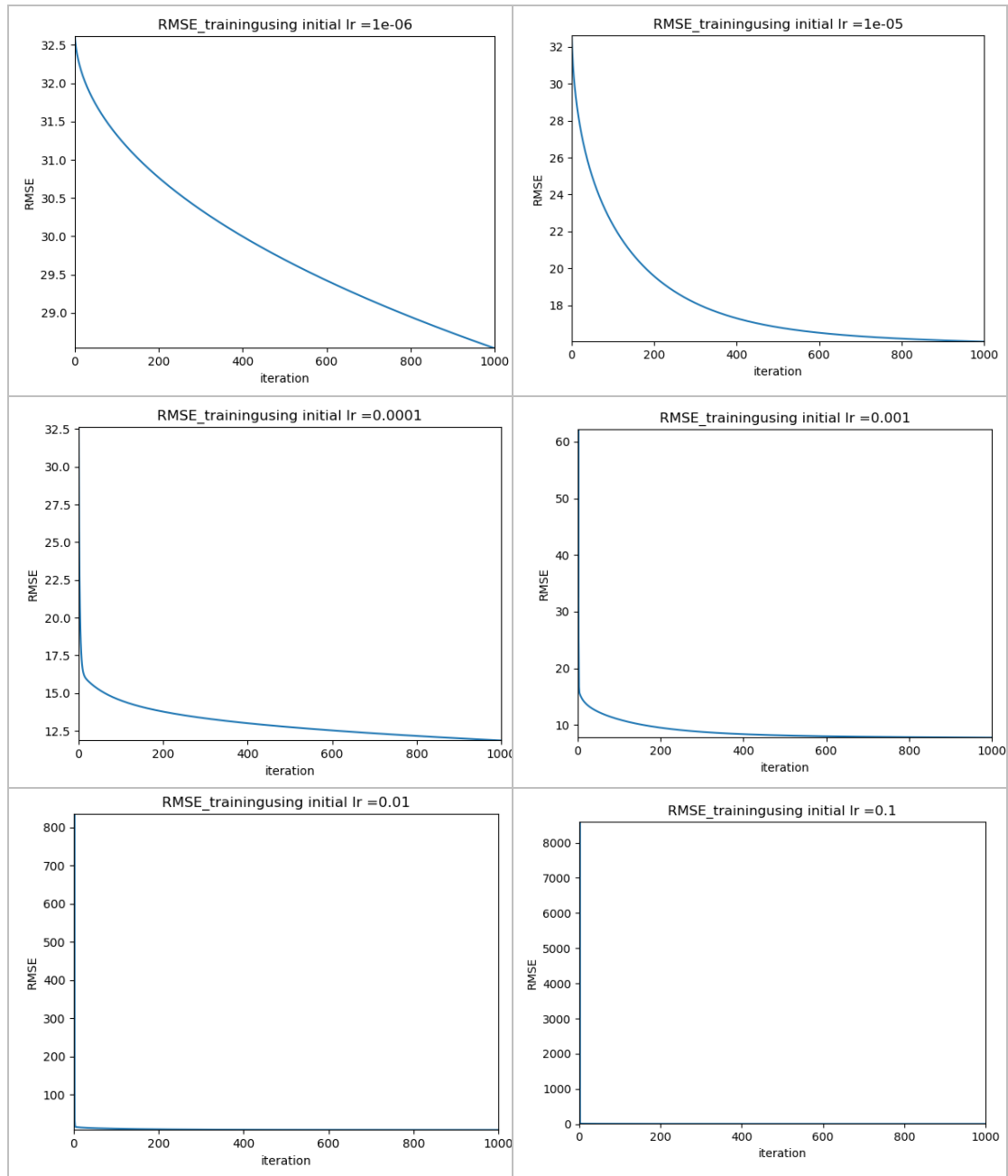


Homework 1 Report - PM2.5 Prediction

學號：r06631035 系級：生機所碩二 姓名：王凱陞

1.



由左至右，上至下的順序 learning rate 分別為 0.000001、0.00001、0.0001、0.001、0.01、0.1 的 RMSE-iteration 變化圖，固定 iteration 為 1000 且加入 regularization 固定參數為 0.05，可從上述順序看出，learning rate 的上升會使整個 training 的收斂加快，也就是如同課堂所說的 learning rate 的提高會愈快趨近到低處，但這部分看不到的是 learning rate 過大產生的震盪效果，是由於這部分我已經採用 adagrad，所以實際 training 進去的 learning rate 會隨著權重和 iteration 的進行而被拉低。

2.

Submission and Description	Private Score	Public Score
submission.csv a minute ago by r06631035_王凱陞 For Question 2, all feature.	8.79604	9.07504

上圖為所有 feature 的一次項所得到的結果，iteration 為 5000，learning rate 為 0.01，未加入 regularization，僅做將小於零的資料轉為 0 的資料處理所跑出的結果。

Submission and Description	Private Score	Public Score
submission.csv just now by r06631035_王凱陞 For Question 2, only PM2.5 as feature.	10.10747	10.03291

此為僅有 PM2.5 的一次項，其餘參數和資料處理皆和上圖一致。

可由此二結果看到，若單一只用 PM2.5 的一次項作為 feature，明顯會造成結果不理想，雖然並非所有 feature 都跟預測 PM2.5 具一定的相關性，但使用全部的 feature 會增加更多的資料，且有部分 feature 對 PM2.5 的預測有其相關性，故預測的結果會較好一些，但 PM2.5 的一次項在 public/private score 上較為接近一些，這有可能是資料分佈的情形不同，而使用單一使用 PM2.5 預測結果較差以外，也僅用單一 feature，故受到其餘 feature 的資料分部不均影響較小，所以才導致此情況。

3.

	Private	Public	Weight L2 norm	RMSE for training
$\lambda = 0.001$	7.63446	7.21778	14.05656	6.38875
$\lambda = 0.01$	7.64086	7.20220	13.73623	6.38943
$\lambda = 0.05$	7.66619	7.14815	12.32904	6.40561
$\lambda = 0.1$	7.74193	7.11676	10.62326	6.45660
$\lambda = 0.5$	10.54464	9.09383	8.236014	8.40138

上表為 regularization parameter λ 由上到下分別為 0.001、0.01、0.05、0.1、0.5 所輸出的結果，此結果有經資料處理且 iteration 固定為 6617，learning rate 設定在 0.001，可由結果看到當 λ 從 0.001 遞減至 0.1 間，public score 都是呈現變好的情況，但 private score 則是有持續變差，到了 0.5 則是兩個都變得很差，而 weight 都是下降的情況，表示 regularization 對權重確實有 bound 住的情況，這點在 training 的 RMSE 也可以看到，結果都是愈來愈差，在合理的情況下確實可以避免 overfitting 的狀況產生，在 testing score 上呼應了老師上課所說的 regularization 的參數設定上，需要經過 validation set 多方的嘗試，不可以太過相信在 public 的結果，不然很容易在 private 產生過多的限制或是 overfitting，像上表的結果明顯 regularization 太大之後就會產生 private 比較差的情況。

下圖為 kaggle 輸出結果，僅供參考：

Submission and Description	Private Score	Public Score
submission.csv just now by r06631035_王凱陞 Question 3, lambda = 0.5, L2 norm: 8.23601370871	10.54464	9.09383
submission.csv 4 minutes ago by r06631035_王凱陞 Question 3, lambda = 0.1, L2 norm: 10.6232589836	7.74193	7.11676
submission.csv 6 minutes ago by r06631035_王凱陞 Question 3, lambda = 0.05, L2 norm: 12.3290443083	7.66619	7.14815
submission.csv 8 minutes ago by r06631035_王凱陞 Question 3, lambda = 0.01, L2 norm: 13.7360522893	7.64086	7.20220
submission.csv 9 minutes ago by r06631035_王凱陞 Question 3, lambda = 0.001, L2 norm: 14.05656328	7.63446	7.21778

4 (a). 根據題意:

$$\text{SSE}(\mathbf{w}) = E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n (t_n - \mathbf{w}^T \mathbf{x}_n)^2$$

Let

$$\mathbf{x} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N] = \begin{bmatrix} x_{11} & \cdots & x_{1N} \\ \vdots & \ddots & \vdots \\ x_{(M+1)1} & \cdots & x_{(M+1)N} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_{M+1} \end{bmatrix}$$

$$\mathbf{r} = \begin{bmatrix} r_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & r_N \end{bmatrix} \quad (\text{主對角線為 } r_1, r_2, \dots, r_N)$$

$$\rightarrow \text{則 } \mathbf{t} - \mathbf{x}^T \mathbf{w} = [t_1 \ t_2 \ \cdots \ t_N] - [\mathbf{w}^T \mathbf{x}_1 \ \mathbf{w}^T \mathbf{x}_2 \ \cdots \ \mathbf{w}^T \mathbf{x}_N]$$

$$\begin{aligned} \rightarrow \text{SSE} &= [t_1 - \mathbf{w}^T \mathbf{x}_1 \ t_2 - \mathbf{w}^T \mathbf{x}_2 \ \cdots \ t_N - \mathbf{w}^T \mathbf{x}_N] \begin{bmatrix} r_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & r_N \end{bmatrix} \begin{bmatrix} t_1 - \mathbf{w}^T \mathbf{x}_1 \\ \vdots \\ t_N - \mathbf{w}^T \mathbf{x}_N \end{bmatrix} \\ &= (\mathbf{t} - \mathbf{w}^T \mathbf{x}) \mathbf{r} (\mathbf{t} - \mathbf{w}^T \mathbf{x})^T = \mathbf{t} \mathbf{r} \mathbf{t}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T \mathbf{w} - \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{t}^T + \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \mathbf{w} \end{aligned}$$

再解 $\text{SSE}(\mathbf{w} + \Delta \mathbf{w}) - \text{SSE}(\mathbf{w})$

$$\begin{aligned} &= \frac{1}{2} [\mathbf{t} \mathbf{r} \mathbf{t}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T (\mathbf{w} + \Delta \mathbf{w}) - (\mathbf{w} + \Delta \mathbf{w})^T \mathbf{x} \mathbf{r} \mathbf{t}^T + (\mathbf{w} + \Delta \mathbf{w})^T \mathbf{x} \mathbf{r} \mathbf{x}^T (\mathbf{w} + \Delta \mathbf{w}) - \\ &\quad \mathbf{t} \mathbf{r} \mathbf{t}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T \mathbf{w} - \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{t}^T + \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \mathbf{w}] \end{aligned}$$

$$= \frac{1}{2} [-\mathbf{t} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} - \Delta \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{t}^T + \Delta \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \mathbf{w} + \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} + \Delta \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w}]$$

由於 $\mathbf{x} \mathbf{r} \mathbf{t}^T$ 、 $\mathbf{x} \mathbf{r} \mathbf{x}^T \mathbf{w}$ 為 scalar，則

$$\rightarrow \frac{1}{2} [2 \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} - 2 \mathbf{t} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} + \Delta \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w}]$$

又因為 $\Delta \mathbf{w}$ 為 infinitesimal 故可消掉，則

$$= \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} - \mathbf{t} \mathbf{r} \mathbf{x}^T \Delta \mathbf{w} = (\mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T) \Delta \mathbf{w}$$

此 form 即符合 $f(\mathbf{w} + \Delta \mathbf{w}) - f(\mathbf{w}) = \nabla_{\mathbf{w}} f(\mathbf{w})(\Delta \mathbf{w})$ ，

故 $\nabla_{\mathbf{w}} E_D(\mathbf{w}) = (\mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T)$ ，而最佳解 \mathbf{w}^* ，即為 $\nabla_{\mathbf{w}} E_D(\mathbf{w}) = 0$ 時的 \mathbf{w} ，故

$$\nabla_{\mathbf{w}} E_D(\mathbf{w}) = (\mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T - \mathbf{t} \mathbf{r} \mathbf{x}^T) = \mathbf{0}$$

$$\rightarrow \mathbf{w}^T \mathbf{x} \mathbf{r} \mathbf{x}^T = \mathbf{t} \mathbf{r} \mathbf{x}^T$$

$$\rightarrow \mathbf{w}^T = \mathbf{t} \mathbf{r} \mathbf{x}^T (\mathbf{x} \mathbf{r} \mathbf{x}^T)^{-1}$$

$$\rightarrow \mathbf{w}^T = \mathbf{t} \mathbf{r} \mathbf{x}^T (\mathbf{x} \mathbf{r} \mathbf{x}^T)^{-1} \quad (\because ((\mathbf{x} \mathbf{r} \mathbf{x}^T)^{-1})^T = (\mathbf{x} \mathbf{r} \mathbf{x}^T)^{-1}, \mathbf{r}^T = \mathbf{r})$$

$$\rightarrow \mathbf{w}^* = (\mathbf{x} \mathbf{r} \mathbf{x}^T)^{-1} \mathbf{x} \mathbf{r} \mathbf{t}^T$$

4 (b).

根據上述， $\mathbf{w}^* = (\mathbf{x}\mathbf{r}\mathbf{x}^T)^{-1}\mathbf{x}\mathbf{r}\mathbf{t}^T$

代入題目參數:

$$\begin{aligned}\mathbf{w}^* &= \begin{pmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{pmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 5 & 1 \\ 5 & 6 \end{bmatrix}^{-1} \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 10 & 5 \end{bmatrix}^T \\&= \begin{pmatrix} 108 & 107 \\ 107 & 127 \end{pmatrix}^{-1} \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 10 & 5 \end{bmatrix}^T \\&= \begin{pmatrix} \frac{127}{2267} & \frac{-107}{2267} \\ \frac{-107}{2267} & \frac{108}{2267} \end{pmatrix} \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 10 & 5 \end{bmatrix}^T \\&= \begin{pmatrix} \frac{127}{2267} & \frac{-107}{2267} \\ \frac{-107}{2267} & \frac{108}{2267} \end{pmatrix} \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 10 & 5 \end{bmatrix}^T \\&= \begin{bmatrix} \frac{5175}{2267} \\ \frac{-2575}{2267} \end{bmatrix}\end{aligned}$$

5. 根據題意，加入高斯雜訊，則 y 可表示為 $y(x, \mathbf{w}, \epsilon) = w_o + \sum_{i=1}^D w_i(x_i + \epsilon_i)$
則 SSE 為

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w},) - t_n)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i(x_{n,i} + \epsilon_i) - t_n)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n + \sum_{i=1}^D w_i \epsilon_i)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n)^2 + \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n) (\sum_{i=1}^D w_i \epsilon_i) + \\ &\quad \frac{1}{2} \sum_{n=1}^N (\sum_{i=1}^D w_i \epsilon_i)^2 \end{aligned}$$

則 $\frac{1}{2} \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n)^2$ ，即為加高斯雜訊的 Loss function，

根據題意取期望值處理其他項 (即題意的平均):

$\sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n) (\sum_{i=1}^D w_i \epsilon_i)$ 的期望值，會因為 $\sum_{i=1}^D w_i \epsilon_i$ 項為線性函數，其期望值可表示為 $\sum_{i=1}^D w_i E[\epsilon_i]$ ，又 $E[\epsilon_i] = 0$ ，故此項為 0。

$\frac{1}{2} \sum_{n=1}^N (\sum_{i=1}^D w_i \epsilon_i)^2$ 項， $(\sum_{i=1}^D w_i \epsilon_i)^2$ 期望值，由題意可知每一項加入的高斯雜訊 ϵ_i 皆為獨立，且 $\sum_{i=1}^D w_i \epsilon_i$ 為線性函數， w_i 可做為 scalar 提出，則可寫為

$$(\sum_{i=1}^D w_i E[\epsilon_i])(\sum_{j=1}^D w_j E[\epsilon_j]) = \sum_{i=1}^D \sum_{j=1}^D w_i w_j E[\epsilon_i] E[\epsilon_j] = \sum_{i=1}^D \sum_{j=1}^D w_i w_j E[\epsilon_i \epsilon_j]$$

又 $E[\epsilon_i \epsilon_j] = \delta_{ij} \sigma^2$ ，故

$$\sum_{i=1}^D \sum_{j=1}^D w_i w_j E[\epsilon_i \epsilon_j] = \sum_{i=1}^D \sum_{j=1}^D w_i w_j \delta_{ij} \sigma^2 = \sigma^2 \sum_{i=1}^D w_i^2$$

此即為 regularization 的形式，令 $\sigma^2 = \lambda$ ，和前面的 Loss function 做結合，則可得下式，如同課堂所教的 regularization:

$$\frac{1}{2} \sum_{n=1}^N (w_o + \sum_{i=1}^D w_i x_{n,i} - t_n)^2 + \lambda \sum_{i=1}^D w_i^2$$

6.

先解等號左邊部分，已知 $|A|$ 即為 A 的特徵值 $\lambda_1, \lambda_2, \dots, \lambda_n$ 的乘積，取 \ln 後即為連加故可寫為 $\frac{d}{d\alpha} \sum_{i=1}^n \ln(\lambda_i)$ ，則

$$\frac{d}{d\alpha} \sum_{i=1}^n \ln(\lambda_i) = \sum_{i=1}^n \frac{1}{\lambda_i} \frac{d\lambda_i}{d\alpha}$$

等號右邊部分，已知特徵值系統 $Au = \lambda u$ ， A 為可逆矩陣，則 $A^{-1}u = \frac{1}{\lambda}u$ ：

$Au = \lambda u$ ，對 α 微分

$$\rightarrow \frac{d}{d\alpha} Au = \frac{d}{d\alpha} \lambda u, \text{ 左乘 } A^{-1}$$

$$\rightarrow A^{-1} \frac{d}{d\alpha} Au = A^{-1} \frac{d}{d\alpha} \lambda u = \frac{d\lambda}{d\alpha} A^{-1}u = \frac{d\lambda}{d\alpha} \frac{1}{\lambda} u, \text{ 由於 } A \text{ 為實對稱矩陣，可利用}$$

eigenvector 進行對角化，故 α 會和 λ 有關，故可寫為 $\frac{d\lambda}{d\alpha}$ ，

$$\text{則可得 } A^{-1} \frac{d}{d\alpha} A = \frac{d\lambda}{d\alpha} \frac{1}{\lambda} I, I \text{ 為 } n \times n \text{ 的單位矩陣，}$$

根據題意求 $A^{-1} \frac{d}{d\alpha} A$ 的 trace，即為 $\text{Trace}\left(\frac{d\lambda}{d\alpha} \frac{1}{\lambda} I\right)$

$$\text{故 } \text{Trace}\left(\frac{d\lambda}{d\alpha} \frac{1}{\lambda} I\right) = \sum_{i=1}^n \frac{d\lambda_i}{d\alpha} \frac{1}{\lambda_i}$$

則左式等於右式，等號成立，故得證。