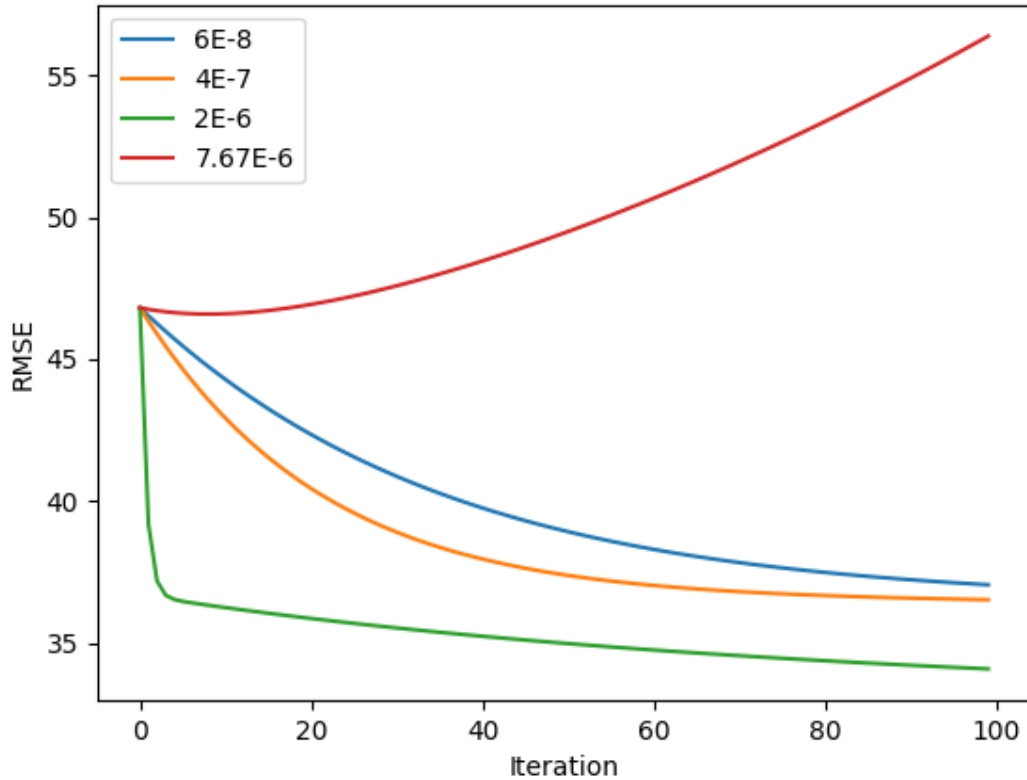


Homework 1 Report - PM2.5 Prediction

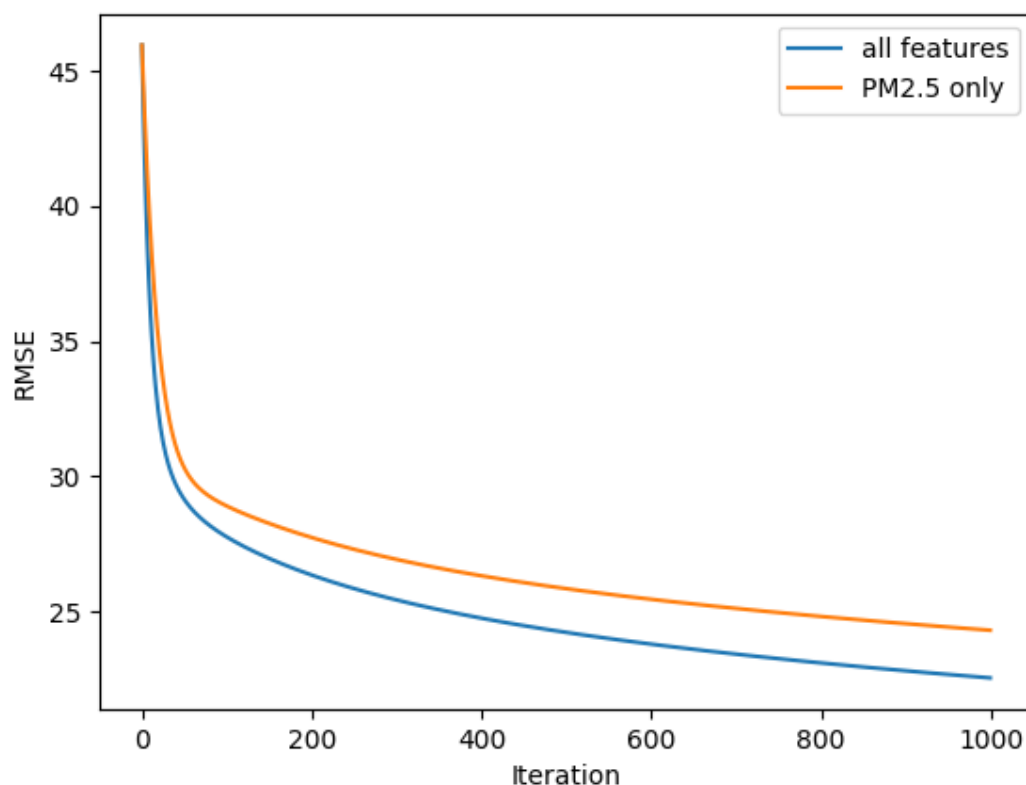
學號: B06901031 系級: 電機二 姓名: 楊宗賢

1. (1%) 請分別使用至少 4 種不同數值的 learning rate 進行 training (其他參數需一致), 對其作圖, 並且討論其收斂過程差異。



$\eta=6 \times 10^{-8}$ 和 $\eta=4 \times 10^{-7}$ 都太小, 以致於雖然藍線與黃線皆嚴格遞減, 斜率仍然太小; $\eta=7.67 \times 10^{-6}$ 太大了, 因此很快就向外發散; $\eta=2 \times 10^{-6}$ 一下就把 RMSE 從初值 46 減至 36, 且在 iteration 100 處亦有不錯的斜率, 是較剛好的 learning rate。

2. (1%) 請分別使用每筆 data9 小時內所有 feature 的一次項（含 bias 項）以及每筆 data9 小時內 PM2.5 的一次項（含 bias 項）進行 training，比較並討論這兩種模型的 root mean-square error（根據 kaggle 上的 public/private score）。



所有 feature: training RMSE=22.537, public score=10.32572, private score=10.37691;
只用 PM2.5: training RMSE=24.302, public score=12.30281, private score=12.51373。

兩種 model 的學習曲線形狀類似，但只用 PM2.5 的 model 總是略遜一籌。我認為使用較多的 feature 才能得到最佳的預測結果。

3. (1%)請分別使用至少四種不同數值的 regularization parameter λ 進行 training (其他參數需一至), 討論及討論其 RMSE(traning, testing) (testing 根據 kaggle 上的 public/private score) 以及參數 weight 的 L2 norm。

4~6 (3%) 請參考數學題目，將作答過程以各種形式 (latex 尤佳) 清楚地呈現在 pdf 檔中 (手寫再拍照也可以，但請注意解析度)。

4-b

$$\text{take } w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \Rightarrow \frac{\partial E_D(w)}{\partial w_k} = \frac{1}{2} \sum_{n=1}^3 2r_n(t_n - w^T x_n)(-x_{nk})$$

$$x_n = \begin{bmatrix} x_{n1} \\ x_{n2} \end{bmatrix} = \sum_{n=1}^3 r_n(w^T x_n - t_n)(x_{nk})$$

the minimum $E_D(w)$ locates at where $\frac{\partial E_D(w)}{\partial w_1} = \frac{\partial E_D(w)}{\partial w_2} = 0$

$$\Rightarrow \begin{cases} [w_1 \ w_2] \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 \times 2 \\ 1 \times 5 \\ 3 \times 5 \end{bmatrix} = [0 \ 10 \ 5] \begin{bmatrix} 2 \times 2 \\ 1 \times 5 \\ 3 \times 5 \end{bmatrix} \\ [w_1 \ w_2] \begin{bmatrix} 2 & 5 & 5 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} 2 \times 3 \\ 1 \times 1 \\ 3 \times 6 \end{bmatrix} = [0 \ 10 \ 5] \begin{bmatrix} 2 \times 3 \\ 1 \times 1 \\ 3 \times 6 \end{bmatrix} \end{cases} \Rightarrow \begin{cases} 108w_1 + 107w_2 = 125 \\ 107w_1 + 127w_2 = 100 \end{cases}$$

$$\Rightarrow w_1 = \frac{5195}{2267}, \quad w_2 = \frac{-2595}{2267} \Rightarrow w^* = \begin{bmatrix} \frac{5195}{2267} \\ \frac{-2595}{2267} \end{bmatrix}$$

4-a

$$\text{take } w = \begin{bmatrix} w_1 \\ \vdots \\ w_k \end{bmatrix}, \quad X = [x_1 \cdots x_n] = \begin{bmatrix} x_{11} & \cdots & x_{n1} \\ \vdots & & \vdots \\ x_{1k} & \cdots & x_{nk} \end{bmatrix}, \quad r = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}, \quad t = [t_1 \cdots t_n]$$

then w^* is the solution of $w^T X \left(\underbrace{[r \ r \cdots r]}_k \circ X^T \right) = t \left(\underbrace{[r \ r \cdots r]}_k \circ X^T \right)$,

where \circ is the entrywise product.

5.

after adding the noise, $y'(x_n, w) = w_0 + \sum_{i=1}^D w_i(x_i + \epsilon_i) = y(x_n, w) + \sum_{i=1}^D \epsilon_i w_i$,

$$\begin{aligned} \Rightarrow E'(w) &= \frac{1}{2} \sum_{n=1}^N (y'(x_n, w)^2 - 2t_n y'(x_n, w) + t_n^2) \\ &= \frac{1}{2} \sum_{n=1}^N \left[(y(x_n, w) - t_n)^2 + 2(y(x_n, w) - t_n) \left(\sum_{i=1}^D \epsilon_i w_i \right) + \left(\sum_{i=1}^D \epsilon_i w_i \right)^2 \right] \\ &= E(w) + \frac{1}{2} \sum_{n=1}^N \left[\sum_{i=1}^D (\epsilon_i \cdot 2w_i (y(x_n, w) - t_n)) + \sum_{i=1}^D \sum_{j=i+1}^D (\epsilon_i \epsilon_j \cdot 2w_i w_j) \right. \\ &\quad \left. + \sum_{i=1}^D (\epsilon_i^2 w_i^2) \right], \end{aligned}$$

because ϵ_i is stochastically independent to x , w and t ,

$$\begin{aligned} E(E'(w)) &= E(E(w)) + \frac{1}{2} \sum_{n=1}^N \left[\overset{=0}{E(\epsilon_i)} E(2w_i (y(x_n, w) - t_n)) \right. \\ &\quad \left. + \overset{=0}{E(\epsilon_i \epsilon_j (i \neq j))} E(2w_i w_j) \right. \\ &\quad \left. + E(\epsilon_i^2) E(w_i^2) \right] \\ &= E \left(\frac{1}{2} \sum_{n=1}^N \left[(y(x_n, w) - t_n)^2 + \sigma^2 \sum_{i=1}^D w_i^2 \right] \right), \end{aligned}$$

which is equivalent to adding

the L2 regularization term $\lambda \sum (w_i^2)$.

6.

the matrix A has n eigenvalues $\lambda_1 \sim \lambda_n$,

their corresponding eigenvectors are $x_1 \sim x_n$,

because A is real, symmetric and non-singular,

we can apply diagonalization on A :

$$A = P \Lambda P^{-1}, \text{ where } P = [x_1 \cdots x_n] \text{ and } \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}.$$

$$\Rightarrow A^{-1} = P^{-1} \begin{bmatrix} \lambda_1^{-1} & & \\ & \ddots & \\ & & \lambda_n^{-1} \end{bmatrix} P, \text{ thus, } A^{-1} \text{ has } n \text{ eigenvalues } \lambda_1^{-1} \sim \lambda_n^{-1}; \quad \textcircled{1}$$

$$|A| = \prod_n \lambda_i \Rightarrow \ln |A| = \sum_n \ln(\lambda_i) \Rightarrow \frac{d}{d\alpha} \ln |A| = \sum_n \lambda_i^{-1} \left(\frac{d}{d\alpha} \lambda_i \right) \quad \textcircled{2};$$

$$\frac{d}{d\alpha} A \text{ has } n \text{ eigenvalues } \frac{d}{d\alpha} \lambda_1 \sim \frac{d}{d\alpha} \lambda_n \quad \textcircled{3};$$

.....

This equation is equivalent to "Jacobi's formula".