

請實做以下兩種不同feature的模型，回答第(1)~(3)題：

- (1) 抽全部9小時內的污染源feature的一次項(加bias)
- (2) 抽全部9小時內pm2.5的一次項當作feature(加bias)

備註：

- a. NR請皆設為0，其他的數值不要做任何更動
- b. 所有 advanced 的 gradient descent 技術(如: adam, adagrad 等) 都是可以用的

1. (2%)記錄誤差值 (RMSE)(根據kaggle public+private分數)，討論兩種feature的影響

- (1) RMSE \approx 6.475 (public:7.46631 private:5.30105)
- (2) RMSE \approx 6.596 (public:7.44013 private:5.62719)

可以發現兩次的public分數都差不多，對於public的那份data，大概抽取多少feature跟只抽，PM2.5是差不多的，但對於private的那份data，可以發現只抽取pm2.5的表現下降不少，因此以這份資料而言，或許資料抽取全部data的表現會優於只抽pm2.5的表現。

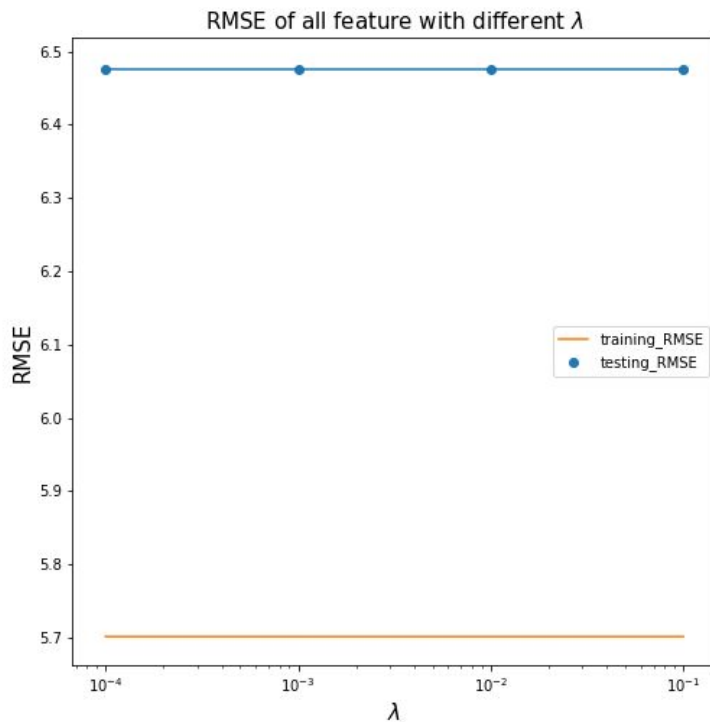
2. (1%)將feature從抽前9小時改成抽前5小時，討論其變化

- (1) RMSE \approx 6.601(public:7.66477 private:5.32990)
- (2) RMSE \approx 6.775(public:7.57904 private:5.79187)

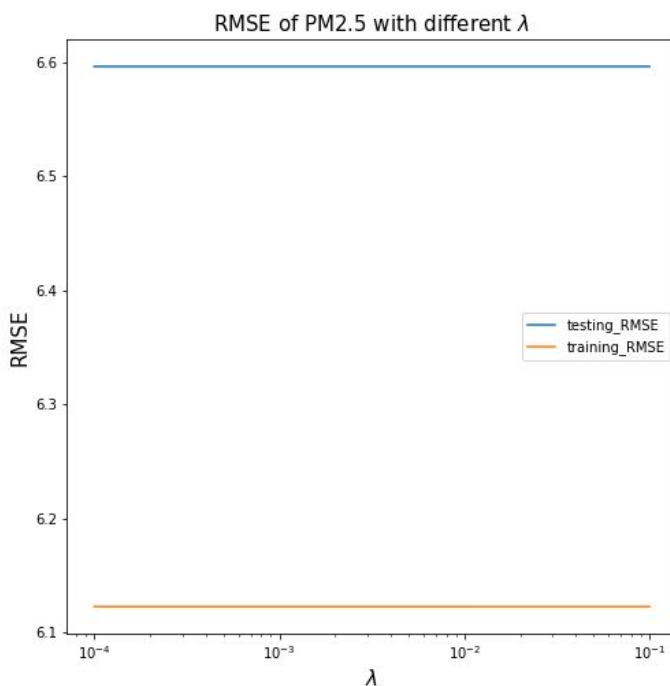
與上面比較，發現不管public，或是private的表現都差了不少，顯然以這份資料而言，預測未來的pm2.5我們抽取9小時前的資料，會比只抽取5小時來的好。

3. (1%) Regularization on all the weight with $\lambda=0.1$ 、 0.01 、 0.001 、 0.0001 ，並作圖

(1) 做了 Regularization 之後發現，用這四個大小的 λ ，RMSE 都沒有差多少，表示這個函數可能對於這組資料，並沒有 Overfit 的情況發生。



(2) 同上，四個大小的 λ 幾乎都沒有讓 RMSE 發生變化，表示這函數對於這組資料可能沒有 Overfitting 的狀況發生。



4. (1%) 在線性回歸問題中，假設有 N 筆訓練資料，每筆訓練資料的特徵 (feature) 為一向量 x^n ，其標註 (label) 為一存量 y^n ，模型參數為一向量 w (此處忽略偏權值 b)，則線性回歸的損失函數 (loss function) 為 $\sum_{n=1}^N (y^n - x^n \cdot w)^2$ 。若將所有訓練資料的特徵值以矩陣 $X = [x^1 \ x^2 \ \dots \ x^N]^T$ 表示，所有訓練資料的標註以向量 $y = [y^1 \ y^2 \ \dots \ y^N]^T$ 表示，請問如何以 X 和 y 表示可以最小化損失函數的向量 w ？請寫下算式並選出正確答案。(其中 $X^T X$ 為 invertible)

- (a) $(X^T X) X^T y$
- (b) $(X^T X)^0 X^T y$
- (c) $(X^T X)^{-1} X^T y$
- (d) $(X^T X)^{-2} X^T y$

(C)

Handwritten derivation of the linear regression solution:

$$X = [x^1 \ x^2 \ \dots \ x^N]^T \quad y = [y^1 \ y^2 \ \dots \ y^N]^T$$

goal: minimize mean squared error

$$\Rightarrow \min \left(\frac{1}{2} \|y - Xw\|^2 \right)$$

so, $L = \frac{1}{2} \|y - Xw\|^2$

$$= \frac{1}{2} (y - Xw)^T (y - Xw) = \frac{1}{2} (y^T y - 2w^T X^T y + w^T X^T X w)$$

take derivative w.r.t. w and set it to zero:

$$\frac{\partial}{\partial w} \frac{1}{2} (y^T y - 2w^T X^T y + w^T X^T X w) = 0$$

$$\frac{1}{2} (2X^T y - 2X^T X w) = 0$$

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

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

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

so the answer is (C)