

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

	Public	Private	Both (RMS)
All feature	7.46237	5.53562	6.57001
PM 2.5	7.44013	5.62719	6.59624

Private 全部污染源優於只取 PM2.5，Public 則相反，整體而言全取較佳。

2. (1%)將 feature 從抽前 9 小時改成抽前 5 小時，討論其變化

	Public	Private	Both (RMS)
All feature	7.66119	5.44024	6.64417
PM 2.5	7.57651	5.79427	6.74452

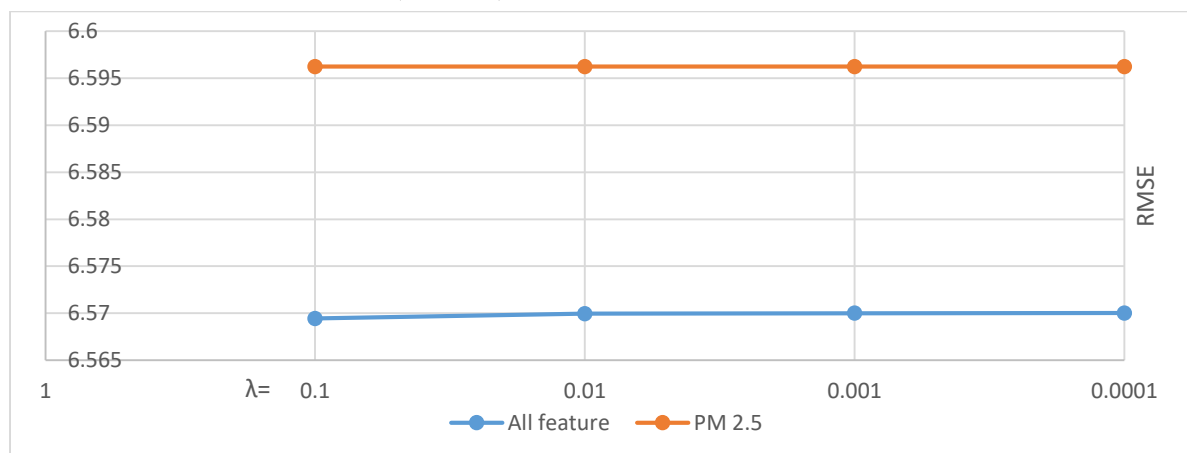
Private 全部污染源優於只取 PM2.5，Public 則相反，整體而言全取較佳。

除了全取的 Private 變小，其餘誤差 5 小時相較於 9 小時皆變大。

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

All feature	Public	Private	Both (RMS)
$\lambda=0.1$	7.46198	5.53477	6.56943
$\lambda=0.01$	7.46233	5.53553	6.56995
$\lambda=0.001$	7.46236	5.53561	6.57000
$\lambda=0.0001$	7.46237	5.53562	6.57001
PM 2.5	Public	Private	Both (RMS)
$\lambda=0.1$	7.44012	5.62720	6.59624
$\lambda=0.01$	7.44013	5.62719	6.59624
$\lambda=0.001$	7.44013	5.62719	6.59624
$\lambda=0.0001$	7.44013	5.62719	6.59624

在此範圍 λ 影響不大(見下圖)。



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)

(c) $(X^T X)^{-1} X^T y$

$$L(w) = \|y - Xw\|^2 = (y - Xw)^T (y - Xw) = y^T y - w^T X^T y - y^T Xw + w^T X^T Xw$$

$(w^T X^T y)^T = y^T Xw$ 為 1×1 的矩陣

$$L(w) = y^T y - 2w^T X^T y + w^T X^T Xw$$

對 w 微分並令它為零

$$-X^T y + (X^T X)w = 0$$

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

Source:

[https://en.wikipedia.org/wiki/Linear_least_squares_\(mathematics\)#Derivation_directly_in_terms_of_matrices](https://en.wikipedia.org/wiki/Linear_least_squares_(mathematics)#Derivation_directly_in_terms_of_matrices)