

1. 請簡明扼要地闡述你如何抽取模型的輸入特徵 (feature)

答：

我總共寫了兩個版本：

hw1_best.py 取前 9 個小時的 pm2.5 指標做一維和二維的 feature。

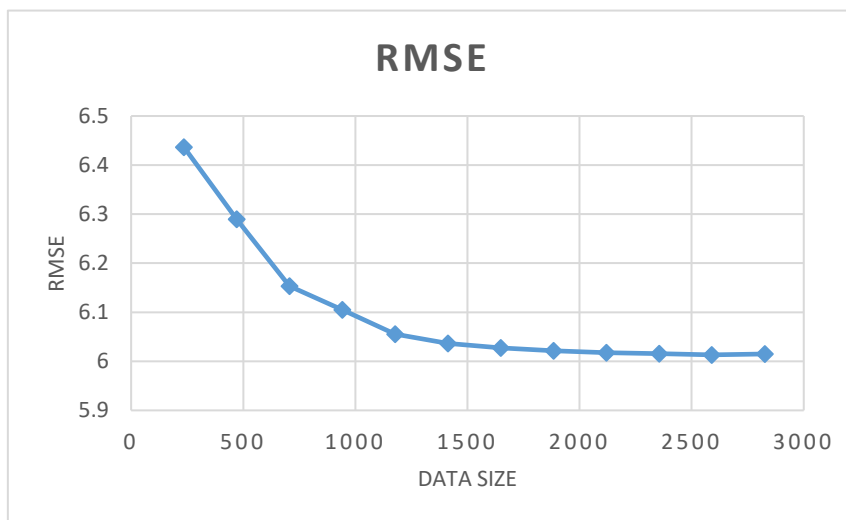
hw1.py 則是取前 9 個小時的 pm10、pm2.5、RAINFALL 指標做一維和二維的 feature。

2. 請作圖比較不同訓練資料量對於 PM2.5 預測準確率的影響

答：

我把整個 data set 分一半，一半只用來計算最後結果的 loss，另一半則是用來 train 出結果。而 train 的時候，我只拿這一半資料中的 $\frac{1}{12} \sim \frac{12}{12}$ 共 12 種不同的資料量去 train 出結果。除此之外，為了怕結果因為資料量太小導致誤差太大，我又多取了不同的資料 train 出不同結果，再算出 RMSE 並平均起來。

得到的結果如下圖：



從圖中可以發現：在資料量等距增加的時候，RMSE 會下降得越來越慢。看起來很像是和資料量成反比的圖形。

而我在網路上查到簡易推導誤差的教學，發現這個 RMSE 以數學的公式推導出來，也會是和資料量成反比。（公式： $E_{out} = noise_level \cdot \left(1 - \frac{d+1}{N}\right)$ ，其中 N 是資料量、 d 是 feature 的元素數量）

3. 請比較不同複雜度的模型對於 PM2.5 預測準確率的影響

答：

我先把資料 shuffle 一遍，再把資料切成兩塊，一塊拿來 train，另一塊則是當作 validation set。接著針對同一筆變數同時用 1 次、2 次、3 次的複雜度來 train，算出 RMSE，再分別對維度取平均。結果如下：

表一 RMSE 先標準化再平均

dim	RMSE 先標準化再平均
1	-0.815379555
2	-0.07284098
3	0.888220534

表二 RMSE 直接平均

dim	RMSE 直接平均
1	11.17740882
2	11.2917164
3	16.28339887

在表一裡面，可以看出來 1 次的預測效果最好，2 次的效果其次，3 次的效果最差。

在表二裡面，可以看出來 1 次和 2 次的預測效果很接近，但是 3 次的效果明顯比較差，很像是 overfitting 的現象。

4. 請討論正規化(regularization)對於 PM2.5 預測準確率的影響

答：

我把資料 shuffle 一遍，再把資料切成兩塊，一塊拿來 train，另一塊則是當作 validation set。而下表就是對於不同的 dim 和 regular 數值，算出的 RMSE 的結果。

表三 對於不同的 dim 和 regular 數值，RMSE 的結果

regular \ dim	0	10	100	1000	10000	100000	1000000
1	5.904	5.902	5.896	5.908	6.531	48.957	227.289
2	5.916	5.915	5.910	5.937	6.512	10.175	180.401
3			5.941	5.965	7.087	12.865	64.950

從表格中可以發現：在 regular 很小的時候，對資料不會有影響，但是當 regular 比 1000 大之後，不管維度是多少，都會讓 loss 變大。

5. 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 \mathbf{x}^n ，其標註(label)為一存量 y^n ，模型參數為一向量 \mathbf{w} (此處忽略偏權值 \mathbf{b})，則線性回歸的損失函數(loss function)為 $\sum_{n=1}^N (y^n - \mathbf{w} \cdot \mathbf{x}^n)^2$ 。若將所有訓練資料的特徵值以矩陣 $\mathbf{X} = [\mathbf{x}^1 \mathbf{x}^2 \dots \mathbf{x}^N]$ 表示，所有訓練資料的標註以向量 $\mathbf{y} = [y^1 y^2 \dots y^N]^T$ 表示，請以 \mathbf{X} 和 \mathbf{y} 表示可以最小化損失函數的向量 \mathbf{w} 。

答：

設 loss function : $E_{in} = \sum_{n=1}^N (y^n - \mathbf{w} \cdot \mathbf{x}^n)^2$

則 $E_{in} = \|\mathbf{X} \cdot \mathbf{w}^T - \mathbf{y}\|^2 = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$

故 $\nabla E_{in} = 2(\mathbf{X}^T \mathbf{X} \mathbf{w}^T - \mathbf{X}^T \mathbf{y})$

設 $0 = \nabla E_{in} = 2(\mathbf{X}^T \mathbf{X} \mathbf{w}^T - \mathbf{X}^T \mathbf{y})$ ，則

$\mathbf{X}^T \mathbf{X} \mathbf{w}^T = \mathbf{X}^T \mathbf{y}$

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

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

另外，如果 $(\mathbf{X}^T \mathbf{X})$ 或 $(\mathbf{X} \mathbf{X}^T)$ 是不可逆的，那就無法做這個運算。