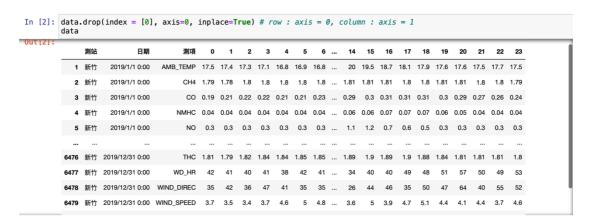
# 結果:

	將未來第一個小時 當預測目標,X只有 PM2.5	將未來第六個小時 當預測目標 X 只有 PM2.5	將未來第一個小時 當預測目標, x 取 所有 18 種屬性	將未來第六個小時 當預測目標 X 有所 有 18 種屬性
Linear Regression	2.6136558911044343	4.839609070011748	3.7329306223569962	3.7604398413161895
Random Forest Regression	3.285230352303523	5.560027285129604	3.5657181571815717	3.690995907230559

## 1.載入資料,因為資料是中文編碼,utf-8 會有亂碼,所以用 big5

In [1]:	impor from from from	t ni skle skle skle	earn.metr earn.ense		mean_ Rand	absol omFor	ute_e estCla	rror assifi	er												
	4	新竹	2019/1/1 0:00	NMHC	0.04	0.04	0.04	0.04	0.04	0.04	0.04	 0.06	0.06	0.07	0.07	0.07	0.06	0.05	0.04	0.04	0.0
	6476	新 竹	2019/12/31 0:00	THC	1.81	1.79	1.82	1.84	1.84	1.85	1.85	 1.89	1.9	1.89	1.9	1.88	1.84	1.81	1.81	1.81	1
	6477	新 竹	2019/12/31 0:00	WD_HR	42	41	40	41	38	42	41	 34	40	40	49	48	51	57	50	49	5
	6478	新 竹	2019/12/31 0:00	WIND_DIREC	35	42	36	47	41	35	35	 26	44	46	35	50	47	64	40	55	5
	6479	新 竹	2019/12/31 0:00	WIND_SPEED	3.7	3.5	3.4	3.7	4.6	5	4.8	 3.6	5	3.9	4.7	5.1	4.4	4.1	4.4	3.7	4
	6480	新 竹	2019/12/31 0:00	WS_HR	2.6	2.7	2.9	3.6	3.4	3.5	3.5	 3.5	3.6	3.4	3.6	3.5	3.7	3.5	3.1	3.2	3
	6481 r	ows	× 27 column	ıs																	

#### 2. 清理資料



#### (1) 處理前後空白字元

```
In [4]: data.columns

Out[4]: Index(['測站 ', '日期 ', '測項 ', '13', '14', '15', '6', '7', '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21', '22', '23'], dtype='object')

In [5]: data.columns = data.columns.str.strip() data.columns

Out[5]: Index(['測站', '日期', '測項', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21', '22', '23'], dtype='object')
```

# (2) 拿 10 月之後的資料

```
In [12]: data = data.loc[data.loc[data['日期'] == '2019/10/1 0:00'].index[0] : , :]
In [13]: data.head(30)
                             日期

        4825
        新竹
        2019/10/1 0:00
        AMB_TEMP
        24.7
        25.1
        25.4
        25.5
        25.3
        25.1
        24.6
        ...
        31.9
        30.4
        29.1
        28.7
        28.3
        28
        27.4
        27
        26.5

           4826 新竹 2019/10/1 0:00
                                         CH4 1.66 1.66 1.7 1.71 1.72 1.71 1.75 ... 1.69 1.72 1.74 1.74 1.78 1.82 1.82 1.83 1.93 1.96
                                       CO 0.05 0.13 0.15 0.17 0.16 0.16 0.22 ... 0.27 0.32 0.36 0.39 0.49 0.57 0.58 0.6 0.69 0.49
           4827 新竹 2019/10/1 0:00
                                        NMHC 0 0 0.02 0.02 0.01 0.03 ... 0.07 0.1 0.08 0.08 0.13 0.21 0.22 0.21 0.22 0.17
           4828 新竹 2019/10/1 0:00
           4829 新竹 2019/10/10:00 NO 0 0.3 0.3 0.3 0.3 1.2 ... 1.6 1.1 0.8 0.6 0.6 0.7 0.6 2 5.1 3.8
                                         NO2 1.3 1.1 1.3 1.8 2.1 2.3 6.1 ... 5.9 8.3 10.1 9.6 14.9 19.5 22.2 26.7 24 15.7
           4830 新竹 2019/10/1 0:00
                                        NOx 1.2 1.2 1.7 2.1 2.4 2.7 7.3 ... 7.4 9.4 10.8 10.3 15.5 20.2 22.8 28.6 29.1 19.5
           4832 新竹 2019/10/1 0:00
                                          O3 16.7 17.3 21.9 21.5 18.8 17.5 11.8 ...
                                                                                    41 54.1 57.7 49.3 39.1 30.8 23.3 13.1 7.5 5.5
           4833 新竹 2019/10/10:00 PM10 18 18 29 29 27 19 23 ... 17 23 47 53 47 35 33 21 22 14
           4834 新竹 2019/10/1 0:00
                                   PM2.5 2 4 6 9 10 7 4 ... 3 8 19 17 19 19 19 14 17
```

(3) NR 表示無降雨,以0取代

```
In [14]: data = data.replace(['NR'], [0.0])
    data.head(30)
```

(4) 把怪怪的字元拿掉,用前後兩個值的平均

```
In [18]:

symbols = ['#', '*', 'x', 'A']

for i in range(3) (en(data.columns)): # column

for j in range(3) (en(data.columns)): # column

for symbol in symbols:

if symbol in data.iloc[i, j]:

front_j = j + 1

front_i = i

back_j = j + 1

front_j = = 2: # (代表取到了測項,要往前一天的最後一小時取資料

front_j = len(data.columns) - 1

front_i = i - 18 # 前一天的 row index 在前面18個

if back_j == len(data.columns): # (代表取到超過今天的數據了

back_j = 3 # 隔天的第一舉數據

back_j = i + 18

# 下一個的symbol 不一定是上一個index 的symbol ,所以要再重测。

# 因為我是用row by row, column by colum,所以不需要去测前面的

while True:

s = 0 # 用來記錄是否都不是symbol

for symbol in symbols:

if symbol in data.iloc[back_i, back_j]:

back_j + 1

else:

s + 1

if s = 4: # (代表4個symbol也都不在back裡面了

break

# 因為可以讓中間都是異常值的數直接被平均數替代掉

data.iloc[i, j:back_j] = str((float(data.iloc[front_i, front_j]) + float(data.iloc[back_i, back_j]))
```

怪字元拿掉了!!

掃描 data 的方式由左至右、由上至下

#### 這邊是分三種狀況:

- a. 如果取到了' 測項'那欄,那就要回前 18 個列取前一天 24 時的資料。
- b. 如果取到了超過 data 欄位的資料,那就要去之後 18 個列取下一天 0 時的資料。
- c. 剩下的狀況就是目前有異常值,但還要繼續檢查下一個是不是異常值。
- (5) 將資料切割成訓練集(10.11 月)以及測試集(12 月)

```
In [20]: data.loc[data['日期'] == '2019/12/31 0:00'].index[0]
Out[20]: 6463
training_set
Out[21]:
                                 測項 0 1 2 3 4 5 6 ... 14 15 16 17 18 19 20
         4825 新竹 2019/10/1 0:00 AMB_TEMP 24.7 25.1 25.4 25.5 25.3 25.1 24.6 ... 32.4 31.9 30.4 29.1 28.7 28.3 28 27.4 27 26.5
         4826 新竹 2019/10/1 0:00
                                CH4 1.66 1.66 1.7 1.71 1.72 1.71 1.75 ... 1.69 1.72 1.74 1.74 1.78 1.82 1.82 1.83 1.93 1.96
         4827 新竹 2019/10/1 0:00 CO 0.05 0.13 0.15 0.17 0.16 0.16 0.22 ... 0.27 0.32 0.36 0.39 0.49 0.57 0.58 0.6 0.69 0.49
                                NMHC 0 0 0 0.02 0.02 0.01 0.03 ... 0.07 0.1 0.08 0.08 0.13 0.21 0.22 0.21 0.22 0.17
         4828 新竹 2019/10/1 0:00
         4829 新竹 2019/10/1 0:00 NO 0 0.3 0.3 0.3 0.3 0.3 1.2 ... 1.6 1.1 0.8 0.6 0.6 0.7 0.6 2 5.1 3.8
         5918 新竹 2019/11/30 0:00 THC 1.84 1.82 1.84 1.84 1.91 1.97 1.9 ... 1.83 1.94 1.95 1.94 2.04 2.23 2.2 2 2.05 1.94
                               WD_HR 56 65 67 66 95 96 41 ... 281 271 257 246 213 165 198 162 167 168
         5920 新竹 2019/11/30 0:00 WIND_DIREC 53 70 98 20 94 83 39 ... 258 273 244 167 178 162 225 170 115 179
         5921 新竹 2019/11/30 0:00 WIND_SPEED 1.7 1.7 0.8 0.7 0.9 1.3 1.3 ... 2.1 1.7 1.3 0.8 0.9 1.2 0.8 0.9 0.9
         5922 新竹 2019/11/30 0:00 WS_HR 1.8 1.3 0.9 0.4 0.7 0.8 0.9 ... 1.5 1.6 1.5 1.3 0.6 0.9 0.5 0.6 0.6 1
        1098 rows x 27 columns
```

(6) 再把每小時的資料從 object transfer to float, testing\_set 也做一樣的事

```
測站
Out[24]:
               object
        日期
               object
object
             float64
             float64
             float64
             float64
             float64
float64
             float64
             float64
             float64
        11
             float64
             float64
float64
float64
        12
             float64
```

```
In [25]: testing_set = data.loc[data.loc[data['日期'] == '2019/12/1 0:00'].index[0]: , :]
           testing_set
           5924 新竹 2019/12/1 0:00
                                           CH4 1.86 1.91 1.89 1.87 1.96 1.86 1.86 ... 1.73 1.74 1.75 1.78 1.77 1.75 1.77 1.81 1.79 1.79

        5925
        新竹
        2019/12/1 0:00
        CO
        0.28
        0.29
        0.24
        0.2
        0.23
        0.2
        0.3
        0.26
        0.28
        0.36
        0.37
        0.32
        0.29
        0.32
        0.27
        0.29

           5926 新竹 2019/12/1 0:00
                                         NMHC 0.08 0.1 0.09 0.09 0.1 0.08 0.11 ... 0.04 0.05 0.07 0.12 0.12 0.08 0.09 0.12 0.09 0.09
                                      NO 0.6 0.3 0.3 0.6 0.3 1.95 1.95 ... 0.3 0.5 0.3 0.3 0.3 0.3 0.3 0.3 0.4 0.5
           5927 新竹 2019/12/1 0:00
                                          THC 1.81 1.79 1.82 1.84 1.84 1.85 1.85 ... 1.89 1.9 1.89 1.9 1.88 1.84 1.81 1.81 1.81 1.8
           6476 新竹 2019/12/31 0:00
           6477 新竹 2019/12/31 0:00
                                       WD HR 42 41 40 41 38 42 41 ... 34 40 40 49 48 51 57 50 49
           6478 新竹 2019/12/31 0:00 WIND_DIREC 35 42 36 47 41 35 35 ... 26 44 46 35 50 47 64 40 55 52
           6479 新竹 2019/12/31 0:00 WIND SPEED 3.7 3.5 3.4 3.7 4.6 5 4.8 ... 3.6 5 3.9 4.7 5.1 4.4 4.1 4.4 3.7 4.6
           6480 新竹 2019/12/31 0:00 WS_HR 2.6 2.7 2.9 3.6 3.4 3.5 3.5 ... 3.5 3.6 3.4 3.6 3.5 3.7 3.5 3.1 3.2 3.6
           558 rows x 27 columns
```

```
日期
                object
        測項
                object
              float64
float64
              float64
              float64
float64
              float64
              float64
              float64
float64
              float64
              float64
float64
        12
13
              float64
              float64
              float64
              float64
```

(7) 製作時序資料: 將資料形式轉換為行(row)代表 18 種屬性,欄(column)代表逐時數據資料 hint: 將訓練集每 18 行合併,轉換成維度為(18,61*24)的 DataFrame*(每個屬性都有 61 天 24 小時共 1464 筆資料)

#### Training\_set

```
In [29]: # 有61筆data・i處理到 len(data) - 2 = 59就好
for i in range(int(len(training_set) / 18) - 1):
    new_training_set = pd.merge(new_training_set, training_set.iloc[(i + 1) * 18 : (i + 2) * 18, 2 : ],on='測項')
          new_training_set
 12 新 2019/10/1
竹 0:00
                        SO2 1.10 1.30 1.40 1.50 1.50 1.60 1.80 ... 1.90 4.60 3.90 2.30 3.30 3.50 3.00 3.10 2.40
                        THC 1.65 1.65 1.70 1.73 1.74 1.72 1.78 ... 1.83 1.94 1.95 1.94 2.04 2.23 2.20 2.00 2.05
                                                                                                                                              1 94
                      WD_HR 292.00 286.00 291.00 286.00 268.00 277.00 247.00 ... 281.00 271.00 257.00 246.00 213.00 165.00 198.00 162.00 167.00 168.00
    新 2019/10/1
竹 0:00 WIND_DIREC 283.00 289.00 276.00 289.00 250.00 262.00 239.00 ... 258.00 273.00 244.00 167.00 178.00 162.00 225.00 170.00 115.00 179.00
     新 2019/10/1
竹 0:00 WIND_SPEED 2.10 1.40
                                                               2.40 2.00 ...
                                                                                 2.10
                                                                                             1.30
                                                                                                                                              1.00
                                           2.00
                                                  2.50
                                                          1.80
                                                                                      1.70
                                                                                                      0.80
                                                                                                            0.90
                                                                                                                   1.20
                                                                                                                          0.80
                                                                                                                                0.90
 17 新 2019/10/1 0:00
                      WS_HR 0.90 1.10 1.80 2.20 1.80 1.90 2.00 ... 1.50 1.60 1.50 1.30 0.60 0.90 0.50 0.60 0.60
                                                                                                                                            1.00
 18 rows × 1467 columns
```

#### 再把不想留的 column 拿掉

```
In [30]: new_training_set.drop(['測站', '日期'], axis=1, inplace=True) new_training_set
                                                        16.70 17.30 21.90 21.50 18.80 17.50
                                                                                                                                                                                                                        8.20
                                                                                                                                                                                                                                                                           64.10 52.10 48.10
                                                                                                                                                                                                                                                                                                                                            51.80 37.00 21.00
                                                                                                                                                                                                                                                                                                                                                                                                                    15.40
                                 PM10 18.00 18.00 29.00 29.00 27.00 19.00 23.00 13.00 10.00 ... 30.00 51.00 41.00 38.00 42.00 56.00 56.00 51.00 41.00 47.00
 10
                       {\sf RAINFALL} \quad 0.00 
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         0.00
                                       RH
                                                       89.00 87.00
                                                                                                 85.00 84.00 85.00 86.00 88.00 90.00 85.00 ... 54.00 54.00 65.00
                                                                                                                                                                                                                                                                                                                                              75.00 78.00
                                                                                                                                                                                                                                                                                                                                                                                           79.00
                                                                                                                                                                                                                                                                                                                                                                                                                  81.00
                                                                                                                                                                                                                                                                                                                                                                                                                                         83.00
                                                                                                                                                                                                                                                                                                                                                                                                                                                              84.00
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       85.00
                                                                                                                                                                                                                                                                                                                      3.90
                                                                                                                                                                                                                                                                                                                                                                                             3.50
                                                                                                                                                                                                                                                                                                                                                                                                                    3.00
                                    SO2
                                                      1.10
                                                                           1.30 1.40 1.50 1.50 1.60
                                                                                                                                                                                            1.80 2.00 1.90 ... 1.90 4.60
                                                                                                                                                                                                                                                                                                                                              2.30 3.30
                                                                                                                                                                                                                                                                                                                                                                                                                                         3.10
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         2.20
 13
                                    THC
                                                          1.65
                                                                                 1.65
                                                                                                        1.70
                                                                                                                             1.73
                                                                                                                                                    1.74
                                                                                                                                                                          1.72
                                                                                                                                                                                                 1.78
                                                                                                                                                                                                                        1.84
                                                                                                                                                                                                                                              1.81 ...
                                                                                                                                                                                                                                                                              1.83
                                                                                                                                                                                                                                                                                                   1.94
                                                                                                                                                                                                                                                                                                                            1.95
                                                                                                                                                                                                                                                                                                                                                   1.94
                                                                                                                                                                                                                                                                                                                                                                      2.04
                                                                                                                                                                                                                                                                                                                                                                                               2.23
                                                                                                                                                                                                                                                                                                                                                                                                                     2.20
                                                                                                                                                                                                                                                                                                                                                                                                                                           2.00
                                                                                                                                                                                                                                                                                                                                                                                                                                                                 2.05
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         1.94
                            WD_HR 292.00 286.00 291.00 286.00 277.00 247.00 203.00 193.00 ... 281.00 271.00 257.00 246.00 213.00 165.00 198.00 162.00 167.00 168.00
 15 WIND_DIREC 283.00 289.00 276.00 289.00 20.00 250.00 260.00 260.00 289.00 250.00 260.00 289.00 183.00 187.00 ... 258.00 273.00 244.00 167.00 178.00 162.00 225.00 170.00 115.00 179.00
 16 WIND_SPEED 2.10 1.40 2.00 2.50 1.80 2.40 2.00 1.70 1.30 0.80 0.90 1.20 0.80 0.90 0.90
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      1.00
                                                                                                                                                                                                2.00
                                                                                                                                                                                                                                                                             1.50
                                                                                                                                                                                                                                                                                                                           1.50
                                                                                                                          2.20
                                                                                                                                                                                                                      1.00
                                                                                                                                                                                                                                              1.50 ...
                                                                                                                                                                                                                                                                                                                                                                                                                     0.50
                                                                                                                                                                                                                                                                                                                                                                                                                                            0.60
18 rows × 1465 columns
```

# Testing\_set

```
In [32]: for i in range(int(len(testing_set) / 18) - 1):
         new_testing_set = pd.merge(new_testing_set, testing_set.iloc[(i + 1) * 18 : (i + 2) * 18, 2 : ],on='測項')
new_testing_set
          12 新 2019/12/1 0:00
                                SO2 2:10 1:80 2:00
                                                       1.60 1.60 1.60 1.70 ... 1.00 0.50 0.50 1.00 1.20 1.40 1.40 1.70 1.60 1.40
                               THC 1.94 2.01 1.98 1.96 2.06 1.94 1.97 ... 1.89 1.90 1.89 1.90 1.88 1.84 1.81 1.81 1.81 1.8
                              WD_HR 171.00 181.00 169.00 179.00 187.00 187.00 156.00 ... 34.00 40.00 40.00 48.00 51.00 57.00 50.00 49.00 53.00
                2019/12/1 0:00 WIND_DIREC 176.00 181.00 174.00 177.00 189.00 167.00 92.00 ... 26.00 44.00 46.00 35.00 50.00 47.00 64.00 40.00 55.00 52.00
                2019/12/1
0:00 WIND_SPEED
                                            0.70
                                                  0.90
                                                        0.80
                                                              0.70
                                                                     0.80
                                                                          0.60 ... 3.60 5.00 3.90 4.70 5.10 4.40
                                                                                                                   4.10
              新 2019/12/1
                              WS_HR 0.90 0.50 0.90 0.80 0.50 0.60 0.40 ... 3.50 3.60 3.40 3.60 3.50 3.70 3.50 3.10 3.20 3.61
         18 rows x 747 columns
```

```
In [33]: new_testing_set.drop(['測站', '日期'], axis=1, inplace=True) new_testing_set
                                                                   O3 21.50
                                                                                                  18.40 19.20 20.90 18.90 20.10 14.40 9.00 16.60 ... 33.40 31.40 30.60 31.20 33.30 34.10 34.10 34.10 33.70
                                                            PM10 41.00 41.00 39.00 37.00 41.00 39.00 38.00 41.00 39.00 38.00 41.00 ... 22.00 27.00 32.00 26.00 24.00 22.00 24.00 27.00 28.00 30.00
                                                            PM2.5
                                                                              18.00 18.00 18.00
                                                                                                                                     16.00 18.00 18.00 15.00 26.00 21.00 ... 14.00 18.00 20.00 17.00 18.00 12.00 13.00 10.00 12.00 11.00
                                                 \mathsf{RH} \quad 86.00 \quad 88.00 \quad 89.00 \quad 89.00 \quad 88.00 \quad 88.00 \quad 88.00 \quad 87.00 \quad 84.00 \quad 78.00 \quad \dots \quad 73.00 \quad 74.00 \quad 76.00 \quad 74.00 \quad 73.00 \quad 72.00 \quad 73.00 \quad 72.00 \quad 72.00 \quad 73.00 \quad 72.00 \quad 73.00 \quad 72.00 \quad 73.00 \quad 
                                  12
                                                             SO2 2.10 1.80 2.00 1.60 1.60 1.60 1.70 1.80 2.00 ... 1.00 0.50 0.50 1.00 1.20 1.40 1.40 1.70 1.60 1.40
                                                                                                                        1.98
                                                                                                                                           1.96
                                                                                                                                                           2.06
                                                                                                                                                                              1.94
                                                                                                                                                                                                     1.97
                                                                                                                                                                                                                  2.05
                                                                                                                                                                                                                                   1.93 ... 1.89
                                                                                                                                                                                                                                                                            1.90 1.89 1.90 1.88 1.84 1.81
                                                         WD HR 171.00 181.00 169.00 179.00 187.00 187.00 156.00 54.00 49.00 ... 34.00 40.00 40.00 49.00 48.00 51.00 57.00 50.00 49.00 53.00
                                 15 WIND_DIREC 176.00 181.00 174.00 177.00 189.00 167.00 92.00 35.00 67.00 ... 26.00 44.00 46.00 35.00 50.00 47.00 64.00 40.00 55.00 52.00
                                  WS_HR 0.90 0.50 0.90 0.80
                                                                                                                                                          0.50
                                                                                                                                                                                                 0.40 0.70 0.90 ... 3.50 3.60 3.40 3.60 3.50 3.70 3.50 3.10 3.20 3.60
                                                                                                                                                                             0.60
                               18 rows × 745 columns
```

- 3. 將未來第一個小時當預測目標, X 只有 PM2.5 (e.g. X[0]會有 6 個特徵, 即第 0~5 小時的 PM2.5 數值)
- (1) 這裡整理 training set 的 x, y 資料

```
In [35]: training x = new_training_set.iloc[new_training_set.loc[new_training_set['NUT'] == 'PM2.5'].index[0], 1:]

print(training_x, shape)
print(len(training_x))

0 x 2
1_x 4
2_x 6
3_x 9
4_x 10
...
19 34
20 36
21 32
22 27
23 22
Name: 9, Length: 1464, dtype: object
(1464,)
1464

In [36]: training_x_np = np.zeros((len(training_x) - 6, 6))
print(training_x_np.shape[0])

[[0. 0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0.]
```

```
In [38]: training_y_np = np.zeros(len(training_x_np))
    def get_PM25_data_y_from(data_df, data_np, amount_in_group):
        for i in range(len(data_df) - amount_in_group):
            data_np[i] = data_df.iloc[i + amount_in_group,]
        return data_np
        training_y_np = get_PM25_data_y_from(training_x, training_y_np, 6)
        print(training_y_np)
        print(training_y_np)
        [ 4. 6. 7. ... 32. 27. 22.]
        1458
```

## (2) 這裡整理 testing\_set 的 x, y 資料

```
In [39]: testing_x = new_testing_set.iloc[new_testing_set.loc[new_testing_set['渊项'] == 'PM2.5'].index[0], 1:] testing_x_np = np.zeros((len(testing_x) - 6, 6)) testing_x_np = get_PM25_data_x_from(testing_x, testing_x_np, 6) print(testing_x_np) print(len(testing_x_np))
                           [[18. 18. 18. 16. 18. 18.]
[18. 18. 16. 18. 18. 15.]
[18. 16. 18. 18. 15. 26.]
                             [18. 20. 17. 18. 12. 13.]
[20. 17. 18. 12. 13. 10.]
[17. 18. 12. 13. 10. 12.]]
In [40]:
testing_y_np = np.zeros(len(testing_x_np))
testing_y_np = get_PM25_data_y_from(testing_x, testing_y_np, 6)
print(testing_y_np)
                           print(len(testing_y_np))
                                                                       16.
                                                                                     15.
                                                                                                                                 11.
                                                                                                                                               17.
                                                                                                                                                              17.
                                                                                                                                                                                                      9. 14.
9. 9.
8. 8.
15. 13.
13. 9.
                             12.
10.
11.
10.
                                          11. 19. 22. 14. 16. 19. 17. 18. 15. 13. 7. 3. 4. 4. 5. 6. 4. 3. 5. 4. 10. 6. 6. 8. 11. 8. 9. 8. 13. 11. 13. 12. 13. 12. 16. 12. 11. 18. 15.
                                                                                                                                                                                          10.
15.
                                                      12. 13. 13. 12. 16. 12. 11. 9. 17. 13. 12. 9. 10. 9. 8. 8. 13. 15. 19. 10. 21. 22. 24. 30. 30. 11. 19. 11. 14. 18. 14. 11. 7. 9. 9. 11. 23. 24. 23. 18. 20. 17. 9. 8. 7. 7. 8. 5. 11. 13. 12. 9. 7. 8. 21. 6. 15. 16. 15. 13. 12. 10. 12. 15. 12. 13. 9. 6. 5. 7. 9. 10. 11. 10.
                                                                                                                                                            8.
15.
32.
21.
                                                                                                                                                                            11.
13.
                                                                                                                                                                                          11.
15.
                             13.
9.
15.
29.
8.
21.
12.
12.
24.
15.
                                           14.
10.
20.
31.
13.
22.
12.
8.
21.
                                                                                                                                30.
18.
11.
                                                                                                                                                                           30.
19.
                                                                                                                                                                                         31.
15.
                                                                                                                                               31.
21.
                                                                                                                                               9.
                                                                                                                                                             9.
                                                                                                                                                                            19.
18.
                                                                                                                                                                                         20.
13.
7.
19.
16.
                                                                                                                                             19. 16.
6. 6.
12. 14.
13. 14.
7. 3.
3. 5.
8. 8.
                                                                                                                                                                           6.
17.
14.
3.
5.
7.
                                                                                    13.
9.
18.
                                           6. 15.
6. 5.
18. 20.
                                                                     12.
7.
17.
                                                                                                                  6.
11.
13.
                                                                                                                                5.
10.
10.
                                1.
7.
                                                                                                    10.
12.
```

#### (3) 做 Linear Regression

"' The best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of y, disregarding the input features, would get a R^2 score of 0.0. "

```
In [41]:

def linear_regression(training_set_x, training_set_y, testing_set_x, testing_set_y):
    x, y = training_set_x, training_set_y
    reg = LinearRegression().fit(x, y)
    predict_y = reg.predict(testing_set_x)
    y_true = testing_set_y
    y_pred = predict_y
    return mean_absolute_error(y_true, y_pred)
```

"" 用於評估預測結果和真實資料集的接近程度的程度,其其值越小說明擬合效果越好。 MAE output is non-negative floating point. The best value is 0.0. ""

```
In [42]: linear_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[42]: 2.6136558911044343
```

# (4) 做 Random Forest Regression

```
In [43]: def random_forest_regression(training_set_x, training_set_y, testing_set_x, testing_set_y):
    clf = RandomForestClassifier(max_depth=10, random_state=1)
    x, y = training_set_x, training_set_y
    y = np_array(y, dtype=int)
    clf.fit(x, y)
    predict_y = clf.predict(testing_set_x)
    y_true = testing_set_y
    y_pred = predict_y
    return mean_absolute_error(y_true, y_pred)
In [44]: random_forest_regression(training_x_np, training_y_np, testing_x_np, testing_y_np)
```

Out[44]: 3.285230352303523

- 4. 將未來第六個小時當預測目標 X 只有 PM2.5 (e.g. X[0]會有 6 個特徵,即第 0~5 小時的 PM2.5 數值)
- (1) 這裡整理 training\_set 的 x, y 資料

(2) 這裡整理 testing\_set 的 x, y 資料

(3) 做 Linear Regression

```
In [49]: linear_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[49]: 4.839609070011748
```

(4) 做 Random Forest Regression

```
In [50]: random_forest_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[50]: 5.560027285129604
```

- 5. 將未來第一個小時當預測目標, x 取所有 18 種屬性 (e.g. X[0]會有 18\*6 個特徵, 即第 0~5 小時的所有 18 種屬性數值)
- (1) 這裡整理 training\_set 的 x, y 資料

[ 4. 6. 7. ... 32. 27. 22.]

```
In [51]: training_x = new_training_set.iloc[:, 1:]
training_x
              7 16.70 17.30 21.90 21.50 18.80 17.50 11.80 8.20 10.40 13.20 ... 64.10 52.10 48.10 51.80 37.00 21.00 15.40 23.00 16.40
             8 18.00 18.00 29.00 29.00 27.00 19.00 23.00 13.00 10.00 13.00 ... 30.00 51.00 41.00 38.00 42.00 56.00 56.00 51.00 41.00
                 2.00
                                    9.00 10.00
                                               7.00 4.00 6.00
                                                                  7.00
                                                                         6.00 ... 16.00 35.00
                                                                                             17.00 20.00 22.00 34.00 36.00 32.00 27.00
                        4.00
                              6.00
             11 89.00 87.00 85.00 84.00 85.00 86.00 88.00 90.00 85.00 75.00 ... 54.00 54.00 65.00 75.00 78.00 79.00 81.00 83.00 84.00
             13 1.65 1.65
                              1.70
                                   1.73 1.74 1.72 1.78 1.84 1.81
                                                                         1.81 ... 1.83 1.94
                                                                                             1.95 1.94 2.04 2.23
                                                                                                                      2.20
                                                                                                                             2.00 2.05
             14 292.00 286.00 291.00 286.00 268.00 277.00 247.00 203.00 193.00 200.00 ... 281.00 271.00 257.00 246.00 213.00 165.00 198.00 162.00 167.00 1
             15 283.00 289.00 276.00 289.00 250.00 262.00 239.00 183.00 187.00 196.00 ... 258.00 273.00 244.00 167.00 178.00 162.00 225.00 170.00 115.00 1
             16 2.10 1.40 2.00 2.50 1.80 2.40 2.00 1.70 2.30 2.80 ... 2.10 1.70 1.30 0.80 0.90 1.20 0.80 0.90 0.90
             17 0.90 1.10 1.80 2.20 1.80 1.90 2.00 1.00 1.50 2.00 ... 1.50 1.60 1.50 1.30 0.60 0.90
                                                                                                                       0.50
             18 rows x 1464 columns
    In [52]: print(training_x.shape[1])
    In [53]: training_x_np = np.zeros((training_x.shape[1] - 6, 18 * 6))
print(len(training_x_np))
              print(training_x_np.shape)
              (1458, 108)
In [54]: def get_PM25_Alldata_x_from(data_df, data_np, amount_in_group):
    # data_df: 原本的資料
    # data_np: 我要存的格式
               # data_ip: 我要评的语式
# amount_in_group: 幾個資料一組
for i in range(data_df.shape[1] - amount_in_group):
    count = 0
    for j in range(18):
        k = i
                       j in range(i, i + 6):
k = i
for k in range(i, i + 6):
    data_np[i][count] = data_df.iloc[j, k]
    count += 1
    k += 1
          return data_np
training_x_np = get_PM25_Alldata_x_from(training_x, training_x_np, 6)
          print(training_x_np)
          [[24.7 25.1 25.4 ... 2.2 1.8 1.9]
[25.1 25.4 25.5 ... 1.8 1.9 2.]
[25.4 25.5 25.3 ... 1.9 2. 1.]
           [25.1 23.6 22.4 ... 0.6 0.9 0.5]
[23.6 22.4 22.1 ... 0.9 0.5 0.6]
[22.4 22.1 21.7 ... 0.5 0.6 0.6]]
return data_np
training_y_np = get_PM25_Alldata_y_from(training_x, training_y_np, 6)
          print(training_y_np)
```

## (2) 再處理 testing set 的 x, y

```
In [56]: testing_x = new_testing_set.iloc[:, 1:]
    testing_x_np = pn.zeros((testing_x.shape[1] - 6, 18 * 6))
    testing_x_np = get_PMES_Alldata_x_from(testing_x, testing_x_np, 6)
    print(testing_x_np)
    print(testing_x_np)

[[20.5 20.1 19.9 ... 0.8 0.5 0.6]
    [20.1 19.9 19.8 ... 0.5 0.6 0.4]
    [19.9 19.8 20.1 ... 0.6 0.4 0.7]
    ...
    [15.9 15.4 15.3 ... 3.5 3.7 3.5]
    [15.4 15.3 15.3 ... 3.5 3.7 3.5]
    [15.4 15.3 15.1 ... 3.5 3.1 3.2]]
    738

In [57]: testing_y_np = np.zeros(len(testing_x_np))
    print(len(testing_y_np))
    testing_y_np = get_PMES_Alldata_y_from(testing_x, testing_y_np, 6)
    print(testing_y_np)

    5. 10. 6. 8. 7. 7. 12. 9. 6. 6. 8. 7. 6. 7.
    8. 6. 6. 16. 15. 17. 8. 11. 17. 17. 8. 11. 10. 14.
    12. 11. 19. 22. 14. 16. 19. 17. 18. 15. 13. 8. 9. 14.
    10. 7. 3. 4. 4. 18. 19. 19. 18. 18. 19. 14.
    10. 7. 3. 4. 4. 18. 19. 19. 18. 18. 19. 14.
    10. 10. 6. 6. 8. 11. 0. 9. 3. 13. 11. 10. 8. 8.
    10. 13. 12. 13. 13. 12. 16. 12. 11. 18. 15. 15. 18. 8.
    10. 13. 12. 13. 13. 12. 16. 12. 11. 18. 15. 15. 15. 13.
    13. 14. 11. 9. 17. 13. 12. 9. 7. 8. 11. 11. 13. 9.
    9. 10. 10. 9. 8. 8. 13. 15. 13. 15. 13. 16.
    15. 20. 19. 10. 21. 22. 24. 30. 31. 32. 30. 31. 26. 28.
    29. 31. 30. 11. 19. 11. 14. 18. 20. 17. 19. 19. 19. 20. 20. 20.
    21. 22. 23. 24. 23. 18. 20. 17. 19. 16. 18. 31. 19. 19. 20. 20. 20.
    21. 22. 23. 24. 23. 18. 20. 17. 19. 19. 19. 25. 27.
    24. 21. 21. 6. 15. 15. 16. 15. 13. 13. 14. 16. 15. 14.
    15. 12. 13. 12. 9. 7. 8. 15. 6. 6. 6. 7. 6. 8.
    7. 6. 6. 7. 9. 10. 11. 10. 10. 14. 16. 15. 14.
    16. 15. 12. 13. 9. 6. 8. 7. 7. 9. 8. 12. 14. 17. 19. 25. 27.
    24. 21. 21. 6. 15. 16. 15. 13. 13. 14. 14. 16. 15. 14.
    16. 15. 12. 13. 9. 6. 8. 7. 5. 5. 6. 8. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9. 10. 11. 10. 8. 8. 7. 7.
    7. 6. 6. 7. 9.
```

#### (3) 做 Linear Regression

```
In [58]: linear_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[58]: 3.7329306223569962
```

#### (4) 做 Random Forest Regression

```
In [59]: random_forest_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[59]: 3.565718157187
```

6. 將未來第六個小時當預測目標 X 有所有 18 種屬性 (e.g. X[0]會有 18\*6 個特徵,即第 0~5 小時的所有 18 種屬性數值)

## (1) 這裡整理 training set 的 x, y 資料

```
In [60]: training_x = new_training_set.iloc[:, 1:]
         training x
                                                      8.20 10.40 13.20 ... 64.10
         8 18.00 18.00 29.00 29.00 27.00 19.00 23.00 13.00 10.00 13.00 ... 30.00 51.00 41.00 38.00 42.00 56.00 56.00 51.00 41.00
         9 2.00
                  4.00
                         6.00 9.00 10.00 7.00 4.00
                                                     6.00
                                                           7.00
                                                                 6.00 ... 16.00 35.00 17.00 20.00 22.00 34.00 36.00 32.00 27.00
            0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
         11 89.00 87.00 85.00 84.00 85.00 86.00 88.00 90.00 85.00 75.00 ... 54.00 54.00 65.00 75.00 78.00 79.00 81.00 83.00 84.00
         12 1.10 1.30 1.40 1.50 1.50 1.60 1.80 2.00 1.90 1.90 ... 1.90 4.60 3.90 2.30 3.30 3.50 3.00 3.10 2.40
                   1.65 1.70 1.73 1.74 1.72 1.78 1.84 1.81 1.81 ... 1.83 1.94 1.95 1.94 2.04 2.23 2.20
         14 292.00 286.00 291.00 286.00 260.00 277.00 247.00 203.00 193.00 200.00 ... 281.00 271.00 257.00 246.00 213.00 165.00 198.00 162.00 167.00
         15 283.00 289.00 276.00 289.00 250.00 262.00 239.00 183.00 187.00 196.00 ... 258.00 273.00 244.00 167.00 178.00 162.00 225.00 170.00
         16 2.10 1.40 2.00 2.50 1.80 2.40 2.00 1.70 2.30 2.80 ... 2.10 1.70 1.30 0.80 0.90 1.20 0.80 0.90 0.90
            0.90 1.10 1.80 2.20 1.80 1.90 2.00 1.00 1.50 2.00 ... 1.50 1.60 1.50 1.30 0.60 0.90 0.50 0.60
                                                                                                                        0.60
         18 rows × 1464 columns
```

```
In [61]: training_x_np = np.zeros((training_x.shape[1] - 11, 18 * 6))
    training_x_np = get_PM25_Alldata_x_from(training_x, training_x_np, 11)
    print(training_x_np)
    print(len(training_x_np))

[[24.7 25.1 25.4 ... 2.2 1.8 1.9]
        [25.1 25.4 25.5 ... 1.8 1.9 2. ]
        [25.4 25.5 25.3 ... 1.9 2. 1. ]
        ...
        [23.7 24.9 25.1 ... 1.7 1.5 1.6]
        [24.9 25.1 25.1 ... 1.5 1.6 1.5]
        [25.1 25.1 25.2 ... 1.6 1.5 1.3]]

In [62]: training_y_np = np.zeros(len(training_x_np))
        training_y_np = get_PM25_Alldata_y_from(training_x, training_y_np, 11)
        print(training_y_np)
        print(len(training_y_np))

        [4. 6. 7. ... 17. 20. 22.]
        [453]
```

#### (2)再處理 testing set 的 x, y

```
In [63]: testing_x = new_testing_set.iloc[:, 1:]
    testing_x_np = np.zeros((testing_x.shape[1] - 11, 18 * 6))
    testing_x_np = get_PM25_Alldata_x_from(testing_x, testing_x_np, 11)
    print(testing_x_np)
    print(len(testing_x_np))

[[20.5 20.1 19.9 ... 0.8 0.5 0.6]
    [20.1 19.9 19.8 ... 0.5 0.6 0.4]
    [19.9 19.8 20.1 ... 0.6 0.4 0.7]
    ...
    [16.2 16.8 16.9 ... 3.8 3.5 3.6]
    [16.8 16.9 17.1 ... 3.5 3.6 3.4]
    [16.9 17.1 16.3 ... 3.6 3.4 3.6]]

733

In [64]: testing_y_np = np.zeros(len(testing_x_np))
    testing_y_np = get_PM25_Alldata_y_from(testing_x, testing_y_np, 11)
    print(training_y_np)

[ 4.  6.  7.  ... 17. 20. 22.]
```

## (3)做 Linear Regression

```
In [65]: linear_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[65]: 3.7604398413161895
```

## (4)做 Random Forest Regression

```
In [66]: random_forest_regression(training_x_np, training_y_np,testing_x_np, testing_y_np)
Out[66]: 3.690995907230559
```