

# Time Series Regression

## 資料探勘 HW3

Department : IMM, 洪翊誠, ID : 310653005

November 2021

### Introduction

本次作業主題，做回歸分析 (使用線性回歸模型) 問題: 新竹地區 2020 年 10 12 月之空氣品質資料，進行時間序列分析 & 迴歸預測 pm2.5 之值。  
資料來源: [https://airtw.epa.gov.tw/CHT/Query/His\\_Data.aspx](https://airtw.epa.gov.tw/CHT/Query/His_Data.aspx)

### Load data

一般測站資料註記說明：

# 表示儀器檢核為無效值

\* 表示程式檢核為無效值

x 表示人工檢核為無效值

A 係指因儀器疑似故障警報所產生的無效值

空白表示缺值

流程:

1. 閱覽資料:  
整理表格，處理讀檔問題 (*utf8*) 編碼，(可能在不同界面上讀檔會有問題，colab 無法順利執行，在 jupyter lab 可以成功)
2. 處理上述符號型資料:  
因為是時間序列資料，利用前後小時的資料補齊缺失值，若前一小時也是缺失值，則在往前找非缺失值的數值取平均 (需要非常小心，要先把表格重新排序表示)
3. 整理資料架構:  
因為是要利用前 6 小時該第 17 種欄位屬性去預測第 18 種資料 PM2.5，故需要先調整資料架構。

#### 4. 放入模型:

- LinearRegression 訓練線性回歸模型
- XGBoost : XGBRegressor

#### 5. 預測結果: 計算 MAE(mean absolute error)

### 閱覽資料

- 讀資料

	測站	日期	測項	00	01	02	03	04	05	06	...	14	15	16	17	18	19	20	21	22	23
1	新竹	2020/01/01 00:00:00	AMB_TEMP	15.2	15.2	15.3	15.3	15.3	15.4	15.5	...	18.1	18.2	17.9	17.3	16.7	16.4	16.2	16.1	16	15.8
2	新竹	2020/01/01 00:00:00	CH4	1.74	1.74	1.77	1.78	1.77	1.77	1.77	...	1.78	1.78	1.77	1.8	1.81	1.82	1.85	1.83	1.92	1.94
3	新竹	2020/01/01 00:00:00	CO	0.28	0.25	0.24	0.22	0.2	0.19	0.2	...	0.28	0.29	0.28	0.34	0.39	0.41	0.46	0.49	0.58	0.52
4	新竹	2020/01/01 00:00:00	NMHC	0.06	0.07	0.05	0.05	0.05	0.05	0.07	...	0.09	0.09	0.07	0.08	0.12	0.12	0.16	0.14	0.17	0.2
5	新竹	2020/01/01 00:00:00	NO	0.3	0.6	0.6	0.6	0.3	0.3	0.5	...	1.6	1.6	1.2	0.7	0.9	1.1	1.1	1.7	1.8	1.4

Figure 1: DataFrame

- 取出我們所需保留資料並且拉直

```
1 # df_train, df_test 都純取值
2 month = ['10', '11']
3 tmp = df['日期'].apply(lambda x: True if x[5:7] in month else
4                          False)
5 df_train = df[tmp]
6 df_train = df_train.iloc[:,3:]
7
8 month = ['12']
9 tmp = df['日期'].apply(lambda x: True if x[5:7] in month else
10                        False)
11 df_test = df[tmp]
12 df_test = df_test.iloc[:,3:]
```

	0	1	2	3	4	5	6	7	8	9	...	1454	1455	1456	1457	1458	1459	1460	1461	1462	1463
0	23.7	23.8	23.8	23.9	23.9	23.8	24.1	24.7	26	27.2	...	21.6	21.5	20.4	20	20.1	19.9	19.4	18.9	18.9	18.7
1	1.97	1.95	1.96	1.96	1.95	1.96	1.97	1.97	1.96	1.98	...	1.93	1.94	1.93	1.94	1.94	1.95	1.95	1.95	1.95	1.95
2	0.23	0.22	0.21	0.2	0.2	0.22	0.24	0.29	0.27	0.33	...	0.26	0.27	0.27	0.29	0.29	0.31	0.25	0.22	0.2	0.18
3	0.06	0.05	0.03	0.03	0.03	0.04	0.04	0.05	0.06	0.07	...	0.05	0.06	0.06	0.09	0.07	0.09	0.07	0.07	0.07	0.06
4	1.2	0.7	0.5	0.7	0.5	0.3	0.7	0.9	1	1.8	...	2.5	2.4	2	1.8	1.6	1.6	1.8	1.7	1.6	1.6

Figure 2: DataFrame(FLATTEN)

## Function

- **Fill the average value (FILLAVE)**

將符號位置填入數直，利用前後一小時取平均，若前後是符號，就在往更前面或更後面搜尋（若末項是符號，則取與前項相同）。

```

1 symbol_list = ['#', '*', 'x', 'A']
2 def FILLAVE(df_FLATTEN):
3     for a in range(0,len(df_FLATTEN)):
4         for i in range(1,len(df_FLATTEN.columns)-1):
5             j = i-1
6             k = i+1
7             if sum([str(df_FLATTEN.iloc[a,i]).rfind(sym) for sym in
8                 symbol_list])!=-4:
9                 while sum([str(df_FLATTEN.iloc[a,j]).rfind(sym) for
10                     sym in symbol_list])!=-4:
11                     j = j-1
12                 while sum([str(df_FLATTEN.iloc[a,k]).rfind(sym) for
13                     sym in symbol_list])!=-4:
14                     k = k+1
15                 if k>=len(df_FLATTEN.columns)-1:
16                     k=j
17                 df_FLATTEN.iloc[a,i] =
18                     str((float(df_FLATTEN.iloc[a,j]) +
19                         float(df_FLATTEN.iloc[a,k]))/2)
20             if sum([1 for sym in symbol_list if sym in
21                 df_FLATTEN.iloc[a,-1]])!=0:
22                 df_FLATTEN.iloc[a,-1] = df_FLATTEN.iloc[a,-2]
23     return df_FLATTEN
24 df_Train = FILLAVE(df_Tr)
25 df_Test = FILLAVE(df_Te)

```

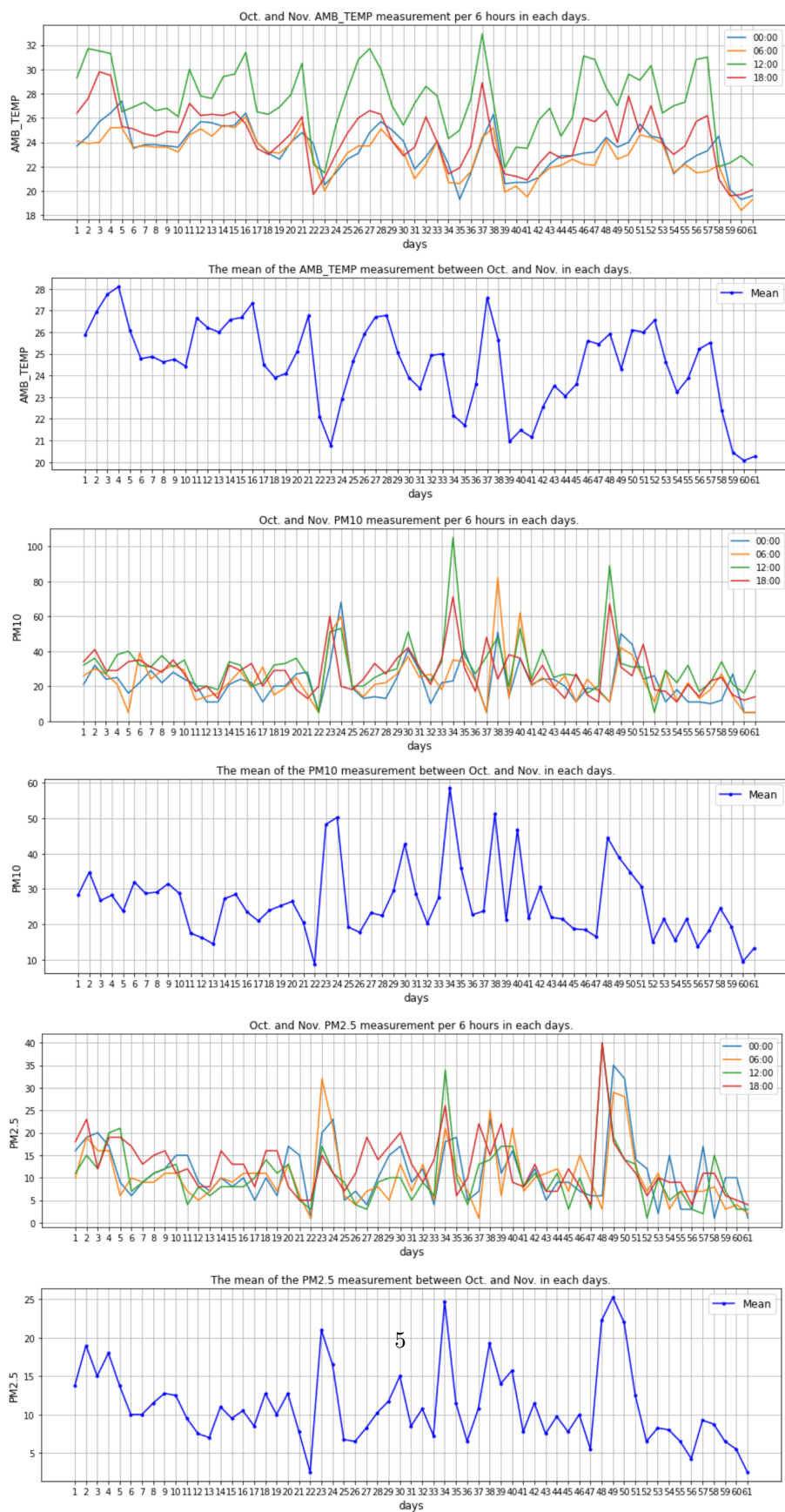
	0	1	2	3	4	5	6	7	8	9	...	734	735	736	737	738	739	740	741	742	743
11	85	86	87	87	86	86	87	86	85	80	...	58	60	61	62	62	61	61	60	60	62
12	0.5	1.9	1.8	1.8	1.8	2	1.8	2.5	3.1	3.3	...	#	#	#	#	#	#	#	#	#	#
13	2.01	2	2.01	2.01	1.99	2.01	1.97	2.03	2.04	2.06	...	2.03	2.07	2.07	2.1	2.1	2.07	2.07	2.05	2.04	2.07
14	43	45	49	43	50	35	36	34	34	37	...	54	50	52	45	47	42	42	47	45	44
15	31	38	40	47	55	37	30	41	35	35	...	48	43	44	33	50	40	46	46	51	38
16	4.3	4.2	3.8	3.9	4.7	5.3	4.3	4	4.6	5.3	...	4.5	4.4	4.2	3.8	3.7	4.7	4.5	4.4	3.9	3.9
17	3.5	3	3	3.1	3	4	3.6	3.3	3.7	4	...	3.7	3.1	3.3	3.1	2.9	3.3	3.1	2.9	2.8	2.6

Figure 3: DataFrame(goal:Filled it)

- Plot the Oct. and Nov. data figure  
藉由視覺化了解數值變化在每天的特定時間走向是如何。

```
1 num = 9
2 Mean = 0
3 per_hour = 6
4 figsize = (15,4)
5 fontsize = 12
6
7 for num in range(18):
8     fig = plt.figure(figsize=figsize)
9     for hours in range(0, 24, per_hour):
10         x1 = np.arange(hours, df_Train.shape[1],24)
11         x = np.arange(1,len(df_Train.iloc[num,x1])+1)
12
13         plt.plot(x, np.float64(df_Tr.iloc[num,x1]))
14         Mean += np.float64(df_Tr.iloc[num,x1])
15     plt.title(f'Oct. and Nov. {indexlist[num]} measurement per
16             {per_hour} hours in each days.',{"fontsize" : fontsize})
17     plt.ylabel(indexlist[num],{"fontsize" : fontsize})
18     plt.xlabel('days',{"fontsize" : fontsize})
19     plt.xticks(x)
20     plt.legend([f'{int(e):02d}:00' for e in range(0, 24, per_hour)])
21     plt.grid()
22     Mean = Mean/(24//per_hour)
23
24     fig = plt.figure(figsize=figsize)
25     plt.title(f'The mean of the {indexlist[num]} measurement
26             between Oct. and Nov. in each days.',{"fontsize" :
27             fontsize})
28     plt.ylabel(indexlist[num],{"fontsize" : fontsize})
29     plt.xlabel('days',{"fontsize" : fontsize})
30     plt.plot(x, Mean, '.-', color='blue')
31     plt.xticks(x)
32     plt.grid()
33     plt.legend(['Mean'],fontsize=12)
34     Mean=0
```

下方列幾張圖示描述，取每天的 00:00、06:00、12:00、18:00 時刻的對應屬性資料，連續繪製觀察其每日變化。有圖中可以看出在一日當中 PM2.5 波動變化其實不大，且在天空氣很糟的情況下，附近天數的空氣也是很糟的。(要看更多圖片可以點選附件的.html)



- Splitting Data 分割訓練集 10~11 月和測試集 12 月

```
1 def SPDA(df_Train, period, pred, feature=0):
2     length = len(df_Train.columns) - pred
3     train_x = [ [] for i in range(length)]
4     if feature!=1:
5         Index = df_Train.index
6     else:
7         Index = [9]
8     for i in range(length):
9         for j in Index:
10            train_x[i].append(df_Train.loc[j, i:i+period-1])
11     train_x = np.float64(train_x).reshape((len(train_x),-1))
12     train_y =
13         np.float64(list(df_Train.loc[9,pred:])).reshape(len(df_Train.loc[9,pred:]),1)
14     return train_x, train_y
```

- Training and Predict

```
1 namelst = ['predict 6-th(1 feature)','predict 6-th(18 features)',
2            'predict 11-th(1 feature)','predict 11-th(18 features)']
3 loss = []
4 for idx,lst in enumerate(data_list):
5     train_X, train_Y, test_X, test_Y = lst
6     lr_model = LinearRegression()
7     lr_model.fit(train_X, train_Y)
8     pred_Y = lr_model.predict(test_X)
9     MAE = mean_absolute_error(test_Y, pred_Y)
10    Plot_the_result(pred_Y, test_Y, hour=10, name=namelst[idx])
11    loss.append(MAE)
12 print('The MAE(mean absolute error) is')
13 print(loss)
```

```
1 import xgboost as xgb
2 loss = []
3 for idx,lst in enumerate(data_list):
4     train_X, train_Y, test_X, test_Y = lst
5     xgb_model = xgb.XGBRegressor(objective='reg:squarederror',
6                                   verbosity=2)
7     xgb_model.fit(train_X, train_Y)
8     pred_Y = xgb_model.predict(test_X)
9     MAE = mean_absolute_error(test_Y, pred_Y)
10    Plot_the_result(pred_Y, test_Y, hour=10, name=namelst[idx])
11    loss.append(MAE)
12 print('The MAE(mean absolute error) is')
13 print(loss)
```

MAE (Mean Absolute Error) Table:

第一橫列說明在使用幾個特徵下預測第幾時刻時間

(6: 往後推 1 小時，11: 往後推 6 小時)

(feature(s), time)	1-6	1-11	18-6	18-11
LinearRegression	2.5224	4.5794	2.6959	6.0882
XGBoost	3.0083	5.0085	2.9761	4.6741

- Show the Result Function Plot

```
1 def Plot_the_result(pred_Y, test_Y, hour=10,  
    name='LinearRegression'):  
2     figsize = (12,4)  
3     fontsize = 12  
4     mean_pred_Y = [np.mean(pred_Y[24*i:24*(i+1)]) for i in  
        range(31)]  
5     mean_test_Y = [np.mean(test_Y[24*i:24*(i+1)]) for i in  
        range(31)]  
6     x = np.arange(1,len(pred_Y[hour:-1:24])+1)  
7     fig = plt.figure(figsize=figsize)  
8     plt.title(f'The prediction of the PM2.5 measurement in Dec. in  
        each days at time {hour:02d}:00.',{"fontsize" : fontsize})  
9     plt.ylabel('Measure of PM2.5',{'fontsize' : fontsize})  
10    plt.xlabel(f'Method : {name}' ,{"fontsize" : fontsize})  
11    plt.xticks(x)  
12    plt.plot(x, pred_Y[hour:-1:24], color='blue', marker='*',  
        linestyle = 'dashed')  
13    plt.plot(x, test_Y[hour:-1:24], color='red', marker='o')  
14    plt.plot(x, mean_pred_Y, color='green', marker='*', linestyle =  
        'dashed')  
15    plt.plot(x, mean_test_Y, color='orange', marker='o')  
16    plt.legend(['pred', 'real', 'M_pred', 'M_real'],fontsize=10)  
17    plt.grid()
```

## Conclusion

從下方預測結果 (每日取 10:00 時刻) 和平均 (每天計算平均) 結果來看，用越多欄位反而導致跟實際資料差異大；而預測的時間越久越遠，也會導致預測偏差比較大。在做預測的過程中，理應要透過預測的數值疊加進輸入值，進一步預測後面項目 (模擬沒有未來資料的情況下，預測後面資料的過程中不會出現沒有資料能輸入的問題，但這就會如同颱風軌跡預測一樣，後面會隨之誤差越來越大)

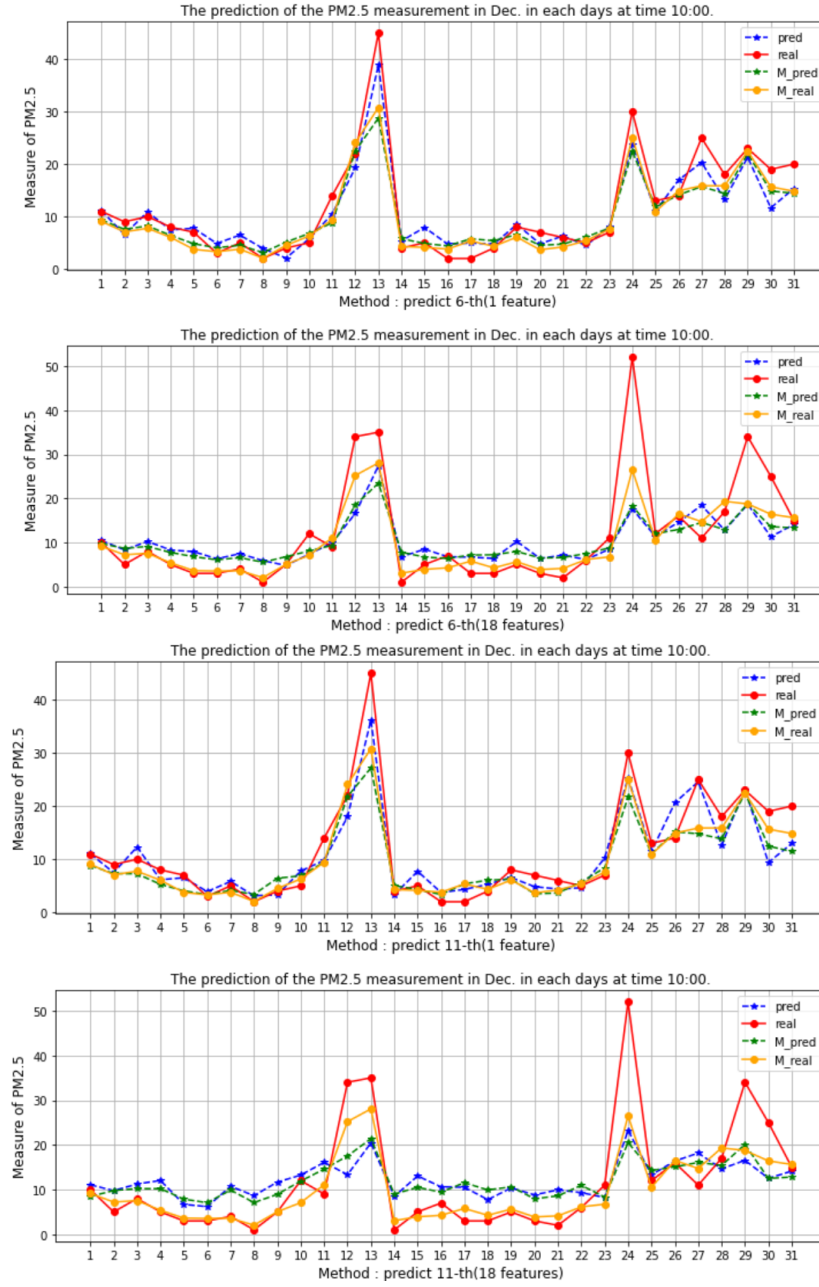


Figure 5: Result (Linear Regression)



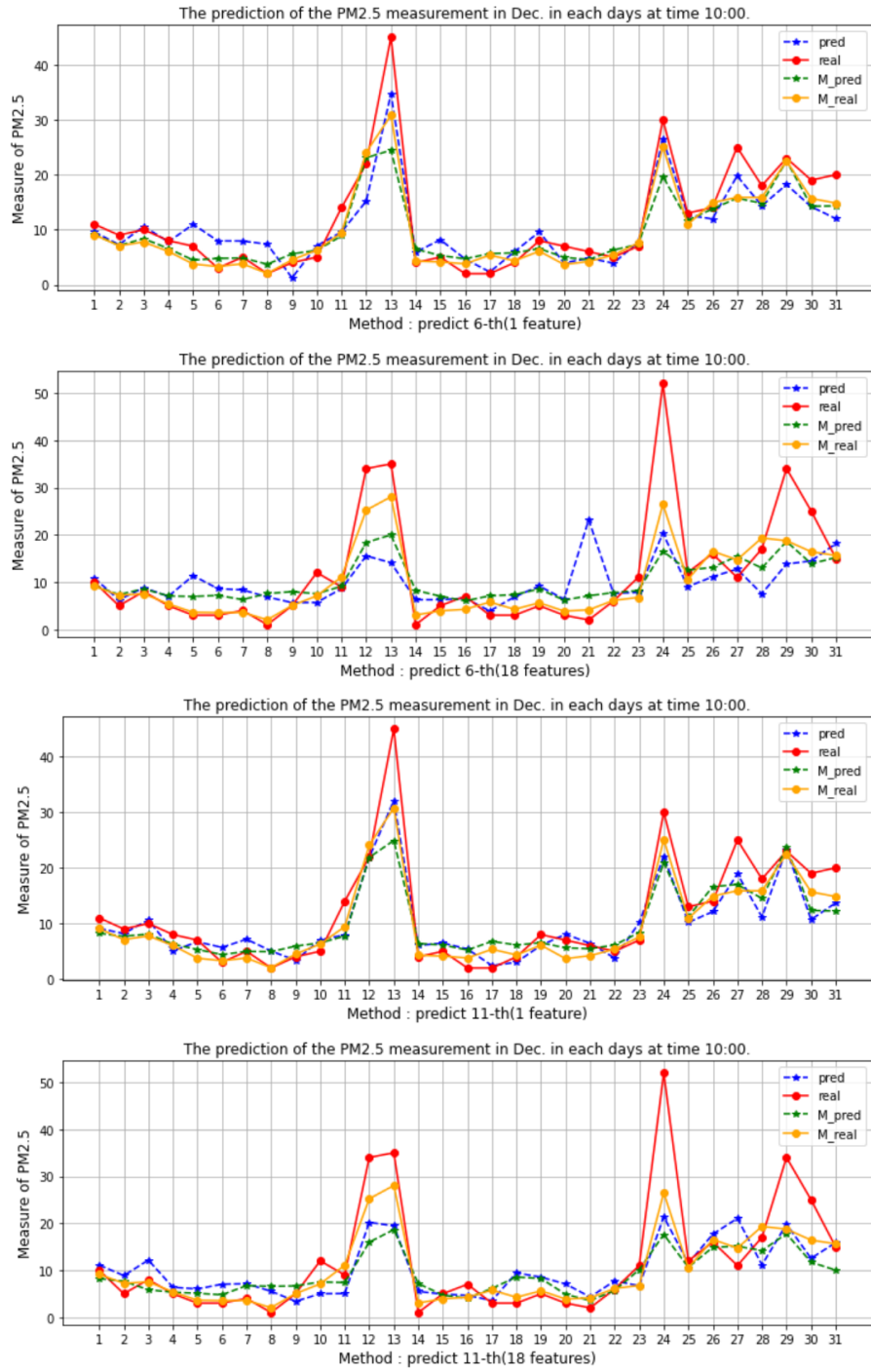


Figure 6: Result (Linear Regression)