

## Week 6 Deep learning

Introduction: This hourly data set contains the PM2.5 data in Beijing, Shanghai, Guangzhou, Chengdu and Shenyang. Meanwhile, meteorological data for each city are also included.

<https://archive.ics.uci.edu/ml/datasets/PM2.5+Data+of+Five+Chinese+Cities>  
(<https://archive.ics.uci.edu/ml/datasets/PM2.5+Data+of+Five+Chinese+Cities>)

## Objective

I used recurrent neural networks and LSTM model to predict the PM2.5 level in Beijing, Shanghai, Guangzhou, Chengdu and Shenyang using time series data.

## Imports

In [1]:

```
1 import sys, os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import warnings
5 warnings.simplefilter(action='ignore')
6 import seaborn as sns
7 import pandas as pd
8 from datetime import datetime
```

In [2]:

```
1 import tensorflow
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, SimpleRNN, LSTM, Activation, Dropout
```

In [3]:

```
1 import math
2 from sklearn.preprocessing import MinMaxScaler
3 from sklearn.metrics import mean_squared_error
```

## Data description

In [4]:

```
1 os.getcwd()
```

Out[4]:

```
'/Users/yanzhenlei'
```

In [5]:

```

1 df_Beijing = pd.read_csv('/Users/yanzhenlei/Beijing.csv')
2 df_Beijing = df_Beijing[df_Beijing.year >= 2015]
3 df_Beijing.head(10)

```

Out[5]:

	No	year	month	day	hour	season	PM_Dongsi	PM_Dongsihuan	PM_Nongzhangua
<b>43824</b>	43825	2015	1	1	0	4	5.0	32.0	8.
<b>43825</b>	43826	2015	1	1	1	4	4.0	12.0	7.
<b>43826</b>	43827	2015	1	1	2	4	3.0	19.0	7.
<b>43827</b>	43828	2015	1	1	3	4	4.0	9.0	11.
<b>43828</b>	43829	2015	1	1	4	4	3.0	11.0	5.
<b>43829</b>	43830	2015	1	1	5	4	3.0	18.0	3.
<b>43830</b>	43831	2015	1	1	6	4	3.0	20.0	6.
<b>43831</b>	43832	2015	1	1	7	4	3.0	22.0	7.
<b>43832</b>	43833	2015	1	1	8	4	NaN	NaN	NaN
<b>43833</b>	43834	2015	1	1	9	4	5.0	37.0	11.

In [67]:

```
1 df_Beijing.info()
```

```

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 8760 entries, 2015-01-01 00:00:00 to 2015-12-31 23:00:00
0
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   No                     8760 non-null  int64
1   year                  8760 non-null  int64
2   month                 8760 non-null  int64
3   day                   8760 non-null  int64
4   hour                  8760 non-null  int64
5   season                8760 non-null  int64
6   PM_Dongsi             8760 non-null  float64
7   PM_Dongsihuan         5465 non-null  float64
8   PM_Nongzhanguan       8760 non-null  float64
9   PM_US Post            8631 non-null  float64
10  DEWP                  8755 non-null  float64
11  HUMI                  8421 non-null  float64
12  PRES                  8421 non-null  float64

```

There are 8760 units for analysis. We need to transform a variable into a datetime type.

In [6]:

```
1 df_Beijing.isnull().sum()
```

Out[6]:

```
No          0
year         0
month        0
day          0
hour         0
season       0
PM_Dongsi    164
PM_Dongsihuan 3295
PM_Nongzhanguan 287
PM_US Post   129
DEWP         5
HUMI         339
PRES         339
TEMP         5
cbwd         5
Iws          5
precipitation 459
Iprec        459
dtype: int64
```

PM\_Dongsi had the the least missing values for PM values so I used PM\_Dongsi as example.

In [66]:

```
1 df_Beijing['PM_Dongsi'].describe()
```

Out[66]:

```
count      8760.000000
mean        87.145776
std         91.899063
min          3.000000
25%         22.000000
50%         58.000000
75%        117.000000
max         685.000000
Name: PM_Dongsi, dtype: float64
```

In [8]:

```
1 df_Beijing['PM_Dongsi'] = df_Beijing['PM_Dongsi'].interpolate()
```

In [9]:

```
1 def make_date(row):
2     return datetime(year = row['year'], month = row['month'], day = row['day'],
3 df_Beijing['date'] = df_Beijing.apply(make_date,axis=1)
```

In [10]:

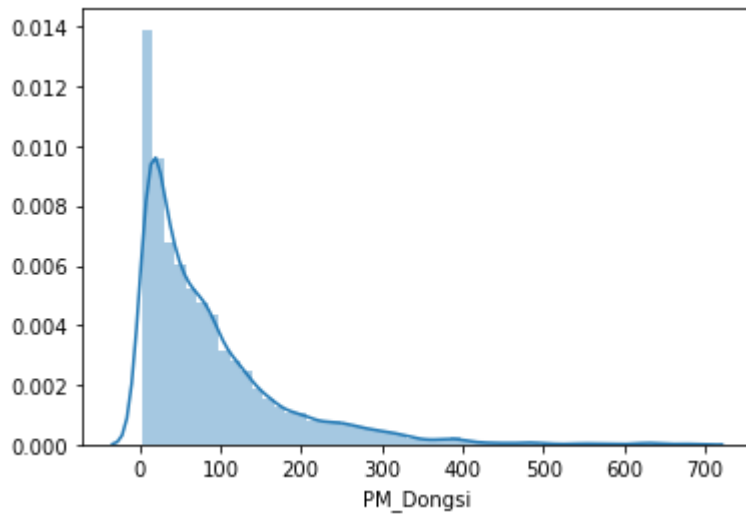
```
1 df_Beijing.set_index(df_Beijing.date,inplace=True)
```

In [68]:

```
1 sns.distplot(df_Beijing['PM_Dongsi'])
```

Out[68]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f98be25ae80>



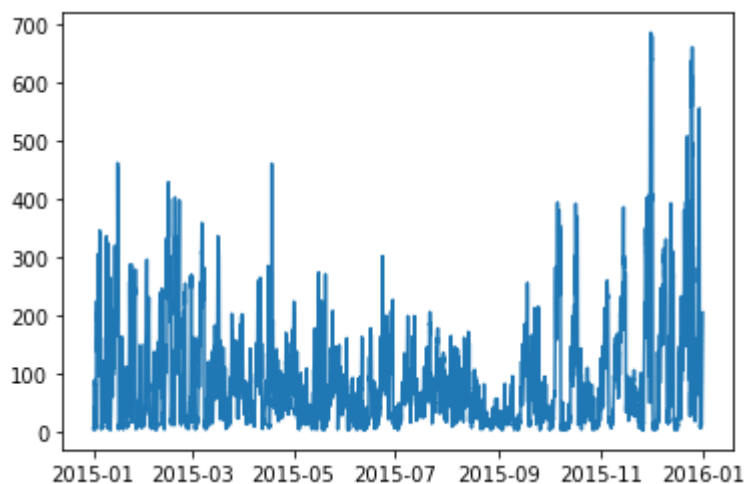
Most of the PM\_Dongsi level concentrated in the zone less than 100 but the extreme large level is also abundant

In [11]:

```
1 plt.plot(df_Beijing['PM_Dongsi'])
```

Out[11]:

[<matplotlib.lines.Line2D at 0x7f9908674438>]



## Data engineering

We use the last 56 days of the PM\_Dongsi series, and will train a model that takes in 12 time steps in order to predict the next time step. We use the last day of data for visually testing the model.

In [12]:

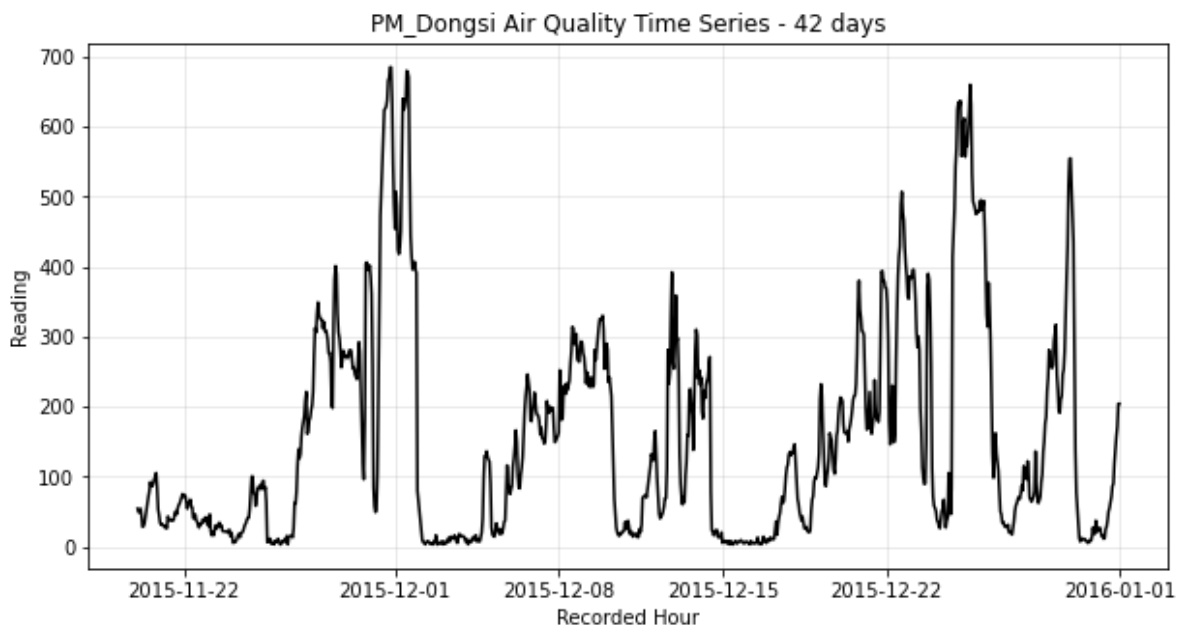
```
1 def get_n_last_days(df,series_name,n_days):
2     return df[series_name][-(24*n_days):]
```

In [13]:

```
1 def plot_n_last_days(df,series_name,n_days):
2     plt.figure(figsize=(10,5))
3     plt.plot(get_n_last_days(df,series_name,n_days),'k-')
4     plt.title('{0} Air Quality Time Series - {1} days'
5             .format(series_name, n_days))
6     plt.xlabel('Recorded Hour')
7     plt.ylabel('Reading')
8     plt.grid(alpha=0.3)
```

In [14]:

```
1 plot_n_last_days(df_Beijing, 'PM_Dongsi', 42)
```



In [43]:

```
1 def get_keras_format_series(series):
2     series = np.array(series)
3     return series.reshape(series.shape[0], series.shape[1], 1)
```

In [44]:

```
1 def get_train_test_data(df, series_name, series_days, input_hours,
2                         test_hours, sample_gap=3):
3
4     forecast_series = get_n_last_days(df, series_name, series_days).values
5
6     train = forecast_series[:-test_hours]
7     test = forecast_series[-test_hours:]
8
9     train_X, train_y = [], []
10
11     for i in range(0, train.shape[0]-input_hours, sample_gap):
12         train_X.append(train[i:i+input_hours])
13         train_y.append(train[i+input_hours])
14
15     train_X = get_keras_format_series(train_X)
16     train_y = np.array(train_y)
17     test_X_init = test[:input_hours]
18     test_y = test[input_hours:]
19
20     return train_X, test_X_init, train_y, test_y
```

In [45]:

```
1 series_days = 56
2 input_hours = 12
3 test_hours = 24
4
5 train_X, test_X_init, train_y, test_y = \
6     (get_train_test_data(df_Beijing, 'PM_Dongsi', series_days,
7                          input_hours, test_hours))
```

In [46]:

```
1 train_X.shape
```

Out[46]:

```
(436, 12, 1)
```

In [47]:

```
1 train_y.shape
```

Out[47]:

```
(436,)
```

In [48]:

```
1 test_X_init.shape
```

Out[48]:

```
(12,)
```

In [49]:

```
1 test_y.shape
```

Out[49]:

```
(12,)
```

## Data Modeling

In [50]:

```
1 def fit_SimpleRNN(train_X, train_y, cell_units, epochs):
2     model = Sequential()
3     model.add(SimpleRNN(cell_units, input_shape=(train_X.shape[1],1)))
4     model.add(Dense(1))
5     model.compile(loss='mean_squared_error', optimizer='adam')
6     model.fit(train_X, train_y, epochs=epochs, batch_size=64, verbose=0)
7     return model
```

In [51]:

```
1 model = fit_SimpleRNN(train_X, train_y, cell_units=10, epochs=10)
```

In [52]:

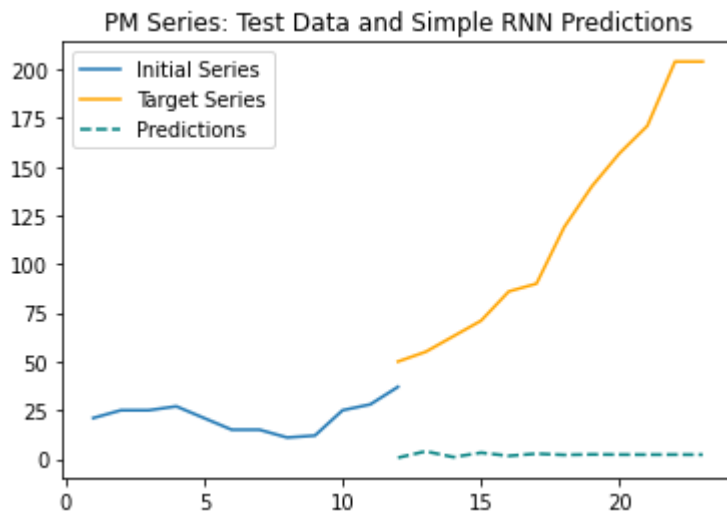
```
1 def predict(X_init, n_steps, model):
2
3     X_init = X_init.copy().reshape(1,-1,1)
4     preds = []
5
6     for _ in range(n_steps):
7         pred = model.predict(X_init)
8         preds.append(pred)
9         X_init[:, :-1, :] = X_init[:, 1:, :]
10        X_init[:, -1, :] = pred
11
12    preds = np.array(preds).reshape(-1,1)
13
14    return preds
```

In [53]:

```
1 def predict_and_plot(X_init, y, model, title):
2
3     y_preds = predict(test_X_init, n_steps=len(y), model=model)
4     start_range = range(1, test_X_init.shape[0]+1)
5     predict_range = range(test_X_init.shape[0], test_hours)
6     plt.plot(start_range, test_X_init)
7     plt.plot(predict_range, test_y, 'orange')
8     plt.plot(predict_range, y_preds, 'teal', linestyle='--')
9
10    plt.title(title)
11    plt.legend(['Initial Series', 'Target Series', 'Predictions'])
```

In [41]:

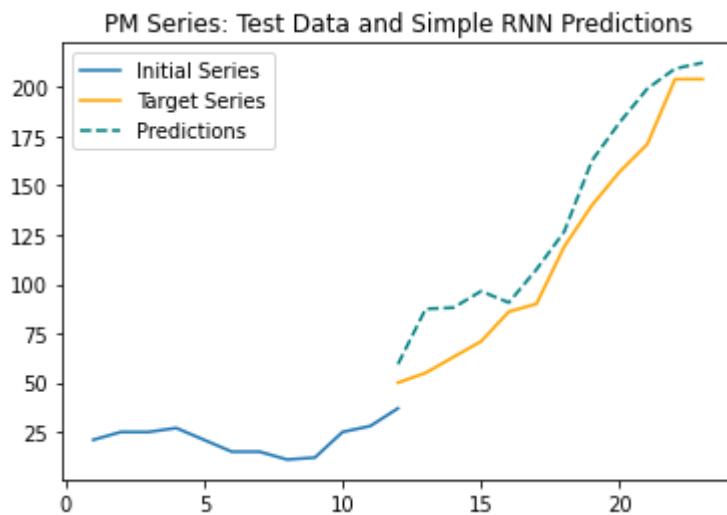
```
1 predict_and_plot(test_X_init,test_y,model,'PM Series: Test Data and Simple RNN P
```



The predicts is bad. We can also pass over the training data many more times, increasing epochs, giving the model more opportunity to learn the patterns in the data.

In [54]:

```
1 model = fit_SimpleRNN(train_X, train_y, cell_units=30, epochs=1200)
2 predict_and_plot(test_X_init,test_y,model,'PM Series: Test Data and Simple RNN P
```





In [55]:

```
1 model.summary()
```

Model: "sequential\_6"

Layer (type)	Output Shape	Param #
=====		
simple_rnn_6 (SimpleRNN)	(None, 30)	960
=====		
dense_6 (Dense)	(None, 1)	31
=====		
Total params: 991		
Trainable params: 991		
Non-trainable params: 0		

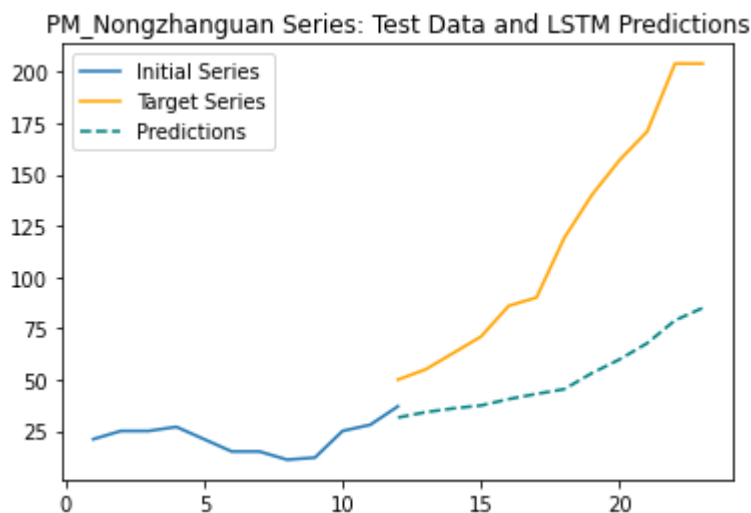
## LSTM Modeling

In [60]:

```
1 def fit_LSTM(train_X, train_y, cell_units, epochs):
2
3     model = Sequential()
4
5     model.add(LSTM(cell_units, input_shape=(train_X.shape[1],1)))
6
7     model.add(Dense(1))
8
9     model.compile(loss='mean_squared_error', optimizer='adam')
10    model.fit(train_X, train_y, epochs=epochs, batch_size=64, verbose=0)
11
12    return model
```

In [65]:

```
1 series_days = 56
2 input_hours = 12
3 test_hours = 24
4
5 train_X, test_X_init, train_y, test_y = \
6     (get_train_test_data(df_Beijing, 'PM_Dongsi', series_days,
7                           input_hours, test_hours))
8
9 model = fit_LSTM(train_X, train_y, cell_units=30, epochs=1200)
10
11 predict_and_plot(test_X_init, test_y, model,
12                  'PM_Nongzhanguan Series: Test Data and LSTM Predictions')
```



LSTM seems not a good way to predict the PM2.5

## Insights and key findings

I used two deep learning models(LSTM and RNN) to predict PM2.5 in time series data. The RNNs model perform better than LSTM model to predict data in my attempt. The parameters(more cell units and more epochs) should be tuned for more attempts. Better result than before and it fits the curve more accurately. There are 991 parameters to run so the computation would be slow.

## Summary

## Flaws and next step

The attempts seems too simple since the time-limited attempts. I should try more parameters in series\_days,input\_hours,test\_hours to figure out the optimal models. Maybe write a loop function to run as much as parameters to find the optimal parameters combination. I can increase the cell\_Units and try more training epochs.

