# analysis_chicago_type2

April 23, 2023

## 1 Analysis Template

### 1.1 Preprocess

```
[ ]: # resolve dependency
     # !pip install pmdarima
```

```
[ ]: import numpy as np
     import pandas as pd
     import matplotlib.pyplot as plt
     from statsmodels.tsa.stattools import adfuller
     from pandas.plotting import autocorrelation_plot
     from statsmodels.graphics.tsaplots import plot_acf,plot_pacf
     import statsmodels.api as sm
     from pmdarima.arima import ADFTest , auto_arima
     %matplotlib inline
```

```
[ ]: data_path = "../data/theft_occurrence_per_day.csv"
     crime = "type2"
     target = "Count"
     date = "Date"
     city = "chicago"
     fig_size = (20,5)
```

```
[ ]: df_by_day = pd.read_csv(data_path)
     df_by_day[date] = pd.to_datetime(df_by_day[date])
     df_by_day.set_index(date, inplace=True)
```

### 1.2 Profiling

#### 1.2.1 By day

```
[ ]: df_by_day.head()
```

```
[ ]:             Count
     Date
     2001-01-01    412
     2001-01-02    221
```
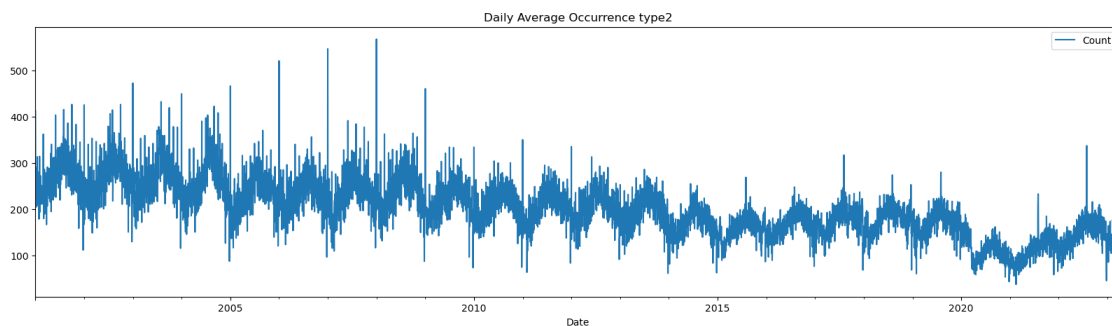
```
2001-01-03      226
2001-01-04      243
2001-01-05      265
```

```
[ ]: df_by_day.describe()
```

```
[ ]:            Count
     count  8132.000000
     mean    201.678554
     std      59.238356
     min      38.000000
     25%     160.000000
     50%     199.000000
     75%     241.000000
     max     567.000000
```

```
[ ]: df_by_day.plot(figsize=fig_size, title="Daily Average Occurrence " + crime)
     plt.show()
```
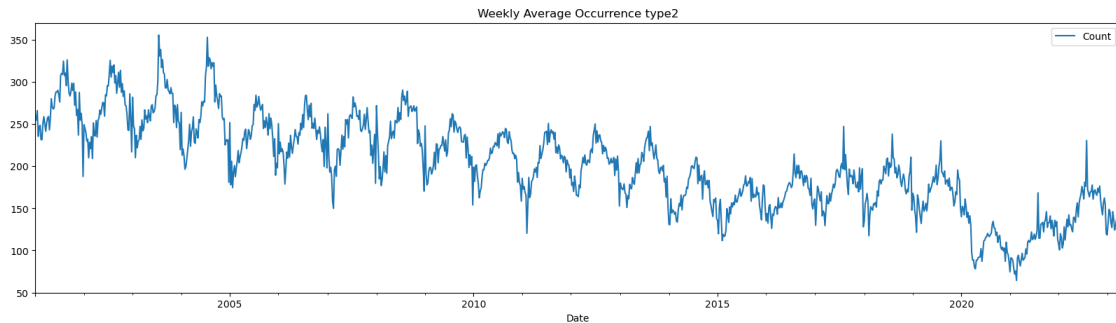


```
[ ]: df_by_day[target].sort_values(ascending=False).head()
```

```
[ ]: Date
     2008-01-01      567
     2007-01-01      546
     2006-01-01      520
     2003-01-01      472
     2005-01-01      466
     Name: Count, dtype: int64
```

### 1.2.2 By week

```
[ ]: df_by_week = pd.DataFrame(df_by_day[target].resample('W').mean())
```
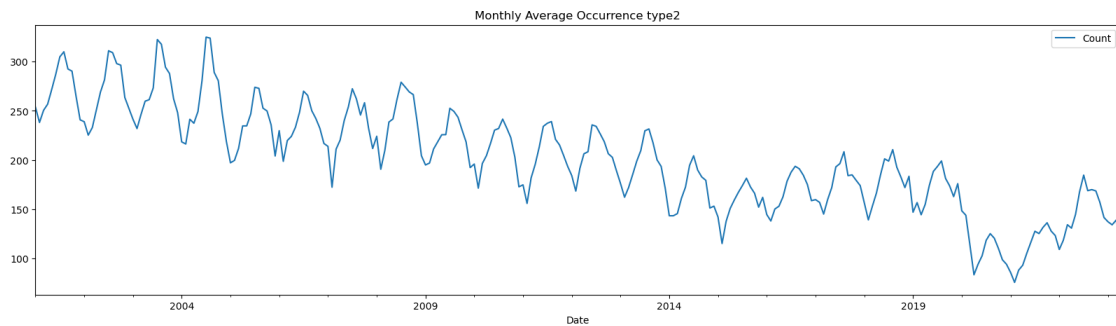
```
df_by_week.plot(
    figsize=fig_size,
    title="Weekly Average Occurrence " + crime)
plt.show()
```



### 1.2.3  By month

```
df_by_month = pd.DataFrame(df_by_day[target].resample('M').mean())
```

```
df_by_month.plot(
    figsize=fig_size,
    title="Monthly Average Occurrence " + crime)
plt.show()
```



## 1.3  Analysis

```
#Ho: It is non stationary
#H1: It is stationary

def adfuller_test(count):
    result=adfuller(count)
```

```
    labels = ['ADF Test Statistic','p-value','#Lags Used','Number of␣
 ↪Observations Used']
    for value,label in zip(result,labels):
        print(label+' : '+str(value) )
    if result[1] <= 0.05:
        print("strong evidence against the null hypothesis(Ho), reject the null␣
 ↪hypothesis. Data has no unit root and is stationary")
    else:
        print("weak evidence against null hypothesis, time series has a unit␣
 ↪root, indicating it is non-stationary ")
```

### 1.3.1 Checking stationary

```
[ ]: adfuller_test(df_by_month[target])
```

```
ADF Test Statistic : -1.2545143770290907
p-value : 0.6497168328079318
#Lags Used : 13
Number of Observations Used : 254
weak evidence against null hypothesis, time series has a unit root, indicating
it is non-stationary
```

### 1.3.2 Checking seasonality

```
[ ]: df_by_month['seasonal_first_difference'] = df_by_month[target] -␣
 ↪df_by_month[target].shift(12)
```

```
[ ]: adfuller_test(df_by_month['seasonal_first_difference'].dropna())
```
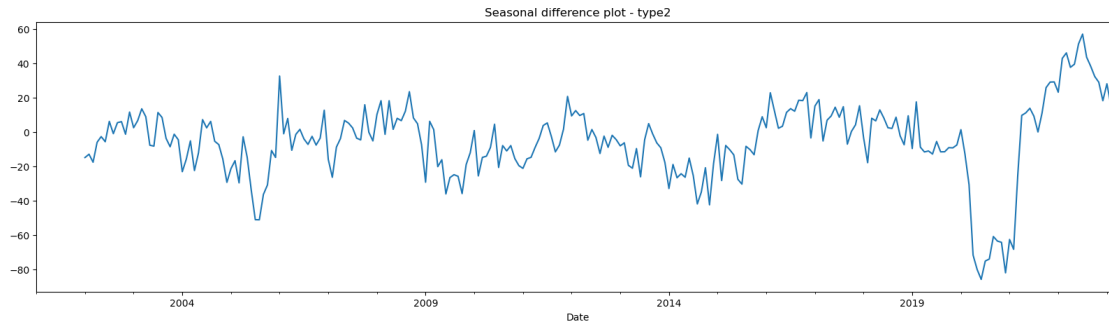
```
ADF Test Statistic : -4.324409680804535
p-value : 0.0004029082855870488
#Lags Used : 12
Number of Observations Used : 243
strong evidence against the null hypothesis(Ho), reject the null hypothesis.
Data has no unit root and is stationary
```

```
[ ]: df_by_month['seasonal_first_difference'].plot(figsize=fig_size, title='Seasonal␣
 ↪difference plot - ' + crime)
```

```
[ ]: <Axes: title={'center': 'Seasonal difference plot - type2'}, xlabel='Date'>
```
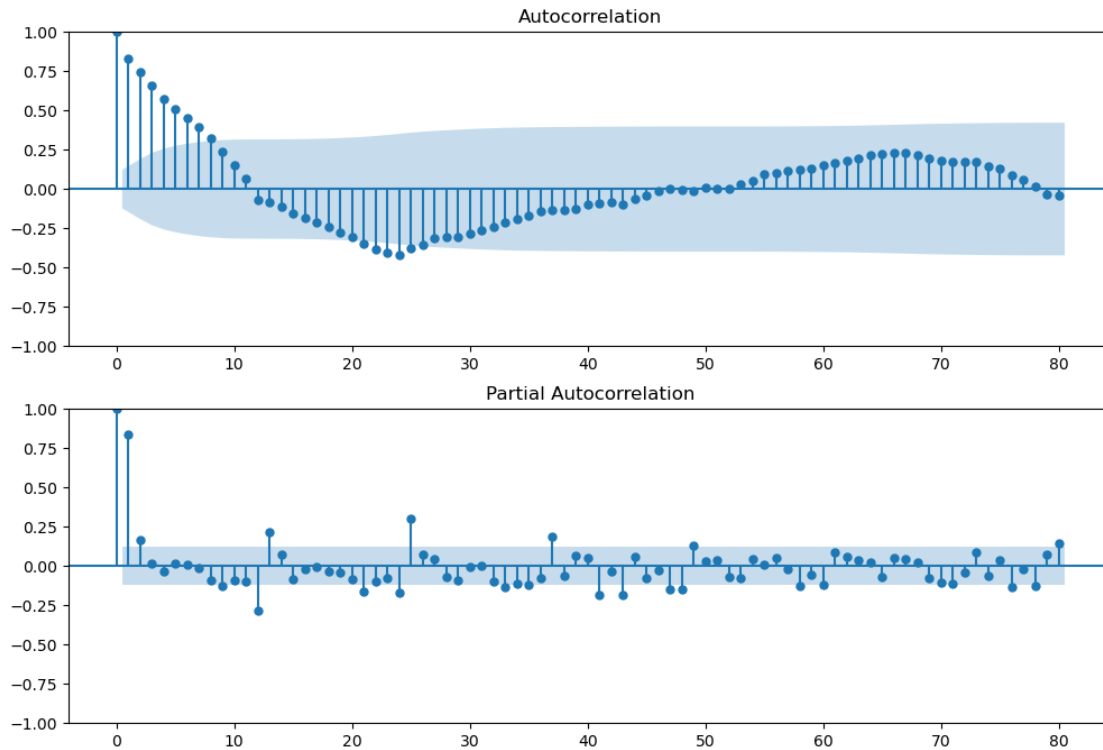
Seasonal difference plot - type2

### 1.3.3 Auto Regressive Model

```
[ ]: fig = plt.figure(figsize=(12,8))
     ax1 = fig.add_subplot(211)
     fig = sm.graphics.tsa.plot_acf(df_by_month['seasonal_first_difference'].iloc[13:
       ↪],lags=80,ax=ax1)
     ax2 = fig.add_subplot(212)
     fig = sm.graphics.tsa.plot_pacf(df_by_month['seasonal_first_difference'].
       ↪iloc[13:],lags=80,ax=ax2)
```

/Users/xuyanchong/opt/anaconda3/lib/python3.9/site-
packages/statsmodels/graphics/tsaplots.py:348: FutureWarning: The default method
'yw' can produce PACF values outside of the [-1,1] interval. After 0.13, the
default will change tounadjusted Yule-Walker ('ywm'). You can use this method
now by setting method='ywm'.
  warnings.warn(

## Autocorrelation

## Partial Autocorrelation

### 1.3.4 Implementing Seasonal Arima Model
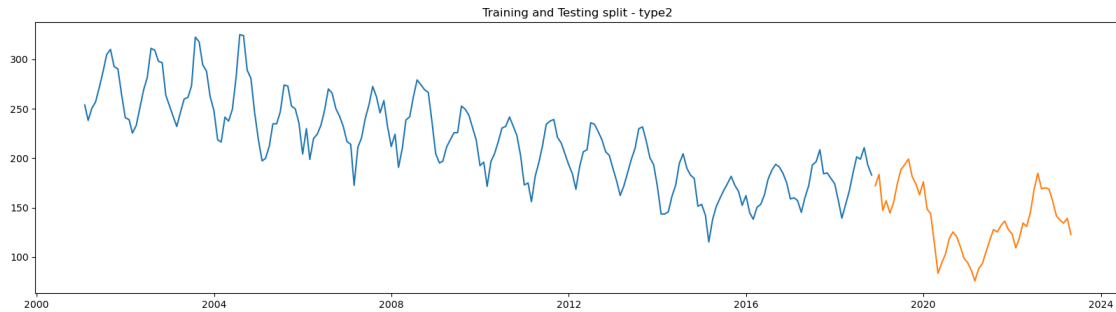
```
[ ]: adf_test=ADFTest(alpha=0.05)
     adf_test.should_diff(df_by_month[target])
```

```
[ ]: (0.01, False)
```

```
[ ]: start=int(df_by_month.shape[0]*0.8)
     train=df_by_month[:start]
     test=df_by_month[start:]
     plt.figure(figsize=fig_size)
     plt.plot(train[target])
     plt.plot(test[target])
     plt.title('Training and Testing split - '+ crime)
     plt.show()
```

Training and Testing split - type2

```
model=auto_arima(train[target],start_p=0,d=1,start_q=0,
        max_p=10,max_d=10,max_q=10, start_P=0,
        D=1, start_Q=0, max_P=10,max_D=10,
        max_Q=10, m=12, seasonal=True,
        error_action='warn',trace=True,
        supress_warnings=True,stepwise=True,
        random_state=20,n_fits=50)
```

```
Performing stepwise search to minimize aic
 ARIMA(0,1,0)(0,1,0)[12]             : AIC=1589.575, Time=0.06 sec
 ARIMA(1,1,0)(1,1,0)[12]             : AIC=1538.297, Time=0.71 sec
 ARIMA(0,1,1)(0,1,1)[12]             : AIC=1493.290, Time=4.14 sec
 ARIMA(0,1,1)(0,1,0)[12]             : AIC=1554.976, Time=0.12 sec
 ARIMA(0,1,1)(1,1,1)[12]             : AIC=1493.261, Time=3.89 sec
 ARIMA(0,1,1)(1,1,0)[12]             : AIC=1528.746, Time=0.56 sec
 ARIMA(0,1,1)(2,1,1)[12]             : AIC=1492.553, Time=31.03 sec
 ARIMA(0,1,1)(2,1,0)[12]             : AIC=1508.505, Time=5.70 sec
 ARIMA(0,1,1)(3,1,1)[12]             : AIC=1493.474, Time=33.00 sec
 ARIMA(0,1,1)(2,1,2)[12]             : AIC=1493.588, Time=50.52 sec
 ARIMA(0,1,1)(1,1,2)[12]             : AIC=1494.060, Time=28.50 sec
 ARIMA(0,1,1)(3,1,0)[12]             : AIC=1495.957, Time=8.14 sec
 ARIMA(0,1,1)(3,1,2)[12]             : AIC=1495.349, Time=31.81 sec
 ARIMA(0,1,0)(2,1,1)[12]             : AIC=1524.448, Time=3.48 sec
 ARIMA(1,1,1)(2,1,1)[12]             : AIC=1490.620, Time=5.36 sec
 ARIMA(1,1,1)(1,1,1)[12]             : AIC=1491.348, Time=0.90 sec
 ARIMA(1,1,1)(2,1,0)[12]             : AIC=1507.600, Time=1.53 sec
 ARIMA(1,1,1)(3,1,1)[12]             : AIC=1491.074, Time=6.20 sec
 ARIMA(1,1,1)(2,1,2)[12]             : AIC=1491.691, Time=2.64 sec
 ARIMA(1,1,1)(1,1,0)[12]             : AIC=1527.738, Time=0.13 sec
 ARIMA(1,1,1)(1,1,2)[12]             : AIC=1492.186, Time=1.65 sec
 ARIMA(1,1,1)(3,1,0)[12]             : AIC=1493.405, Time=0.97 sec
 ARIMA(1,1,1)(3,1,2)[12]             : AIC=1493.010, Time=3.59 sec
 ARIMA(1,1,0)(2,1,1)[12]             : AIC=1503.019, Time=0.85 sec
 ARIMA(2,1,1)(2,1,1)[12]             : AIC=1492.105, Time=1.88 sec
 ARIMA(1,1,2)(2,1,1)[12]             : AIC=1492.040, Time=2.49 sec
 ARIMA(0,1,2)(2,1,1)[12]             : AIC=1491.564, Time=1.11 sec
```

7

```
ARIMA(2,1,0)(2,1,1)[12]                  : AIC=1497.351, Time=1.12 sec
ARIMA(2,1,2)(2,1,1)[12]                  : AIC=1494.499, Time=1.98 sec
ARIMA(1,1,1)(2,1,1)[12] intercept        : AIC=1492.326, Time=2.29 sec

Best model:  ARIMA(1,1,1)(2,1,1)[12]
Total fit time: 236.360 seconds
```

[ ]: `model.summary()`

[ ]: 
```
<class 'statsmodels.iolib.summary.Summary'>
"""
                                SARIMAX Results
================================================================================
==========
Dep. Variable:                            y   No. Observations:
214
Model:             SARIMAX(1, 1, 1)x(2, 1, 1, 12)   Log Likelihood
-739.310
Date:                          Sun, 23 Apr 2023   AIC
1490.620
Time:                                  01:35:32   BIC
1510.440
Sample:                                01-31-2001   HQIC
1498.640
                                     - 10-31-2018
Covariance Type:                            opg
================================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
--------------------------------------------------------------------------------
ar.L1          0.2987      0.141      2.119      0.034       0.022       0.575
ma.L1         -0.7356      0.098     -7.524      0.000      -0.927      -0.544
ar.S.L12       0.0505      0.126      0.401      0.689      -0.196       0.297
ar.S.L24      -0.1595      0.073     -2.186      0.029      -0.303      -0.016
ma.S.L12      -0.6788      0.104     -6.518      0.000      -0.883      -0.475
sigma2        87.5092      7.234     12.098      0.000      73.332     101.687
================================================================================
===
Ljung-Box (L1) (Q):                     0.06   Jarque-Bera (JB):
12.46
Prob(Q):                                0.81   Prob(JB):
0.00
Heteroskedasticity (H):                 0.75   Skew:
0.27
Prob(H) (two-sided):                    0.25   Kurtosis:
4.09
================================================================================
===
```
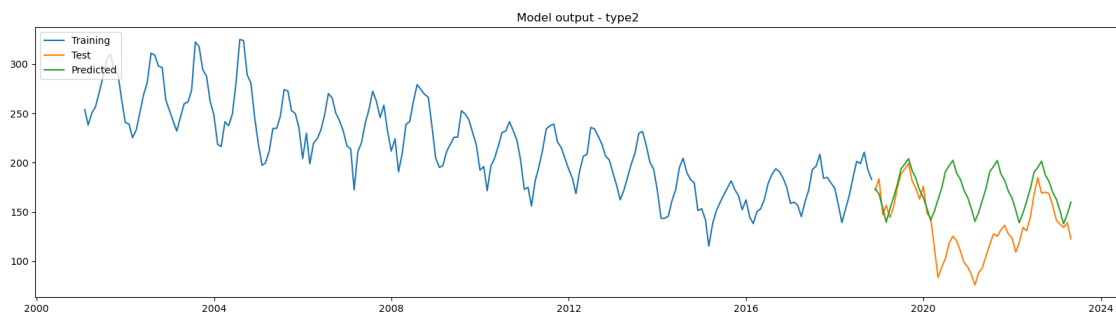
```
Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-
step).
"""
```

```python
prediction = pd.DataFrame(model.predict(n_periods = train.shape[0]),index=test.
 ↪index)
prediction.columns = ['predicted_crime']
plt.figure(figsize=fig_size)
plt.plot(train[target],label="Training")
plt.plot(test[target],label="Test")
plt.plot(prediction,label="Predicted")
plt.legend(loc = 'upper left')
plt.savefig('../output/%s_%s_pred.jpg' % (city,crime))
plt.title('Model output - '+crime)
plt.show()
```



```python
np.sqrt(np.square(np.subtract(test[target].values,prediction['predicted_crime'].
 ↪values)).mean())
```

```
45.26006197362022
```