# Analysis of Rolling Sales Data - Bronx (04/01/2020 - 03/31/2021)

## Steps

I am going to do the following:

```
1. Import necessary modules
2. Load the prepped data per borough
3. Analyze the data for trends and seasonality
4. Dickey-Fuller Tests and preparing data for ARMA modeling
    - Induce stationarity if needed
5. ARMA model of the data
6. Error analysis of the ARMA model
    - Try to improve ARMA model
7. Comparison with latest data
    -Test data from 04/01/2021 - 04/31/2021
```
8. Observations/Conclusions/Recommendations

## 1. Imports

In [1]:
```python
import pandas as pd
from pandas.plotting import register_matplotlib_converters
import matplotlib.pyplot as plt
import matplotlib as mpl
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
import datetime
from statsmodels.tsa.arima_model import ARMA
from statsmodels.tsa.stattools import adfuller, acf, pacf
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
import statsmodels.api as sm
from statsmodels.tsa.seasonal import seasonal_decompose
import numpy as np
from matplotlib.pylab import rcParams
from sklearn.metrics import mean_squared_error
from math import sqrt
import sklearn
import math

#Supress default INFO logging
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
import logging
logger = logging.getLogger()
logger.setLevel(logging.CRITICAL)
import logging, sys
warnings.simplefilter(action='ignore', category=FutureWarning)
```

## 2. Loading the prepared data

*Observations:*

- Once I loaded the data and sorted it, the SALE DATE values range from 4/1/2020 until 3/31/2021.
- This data was the most recent data when I started working on the project.
- NYC OpenData website updates this data regularly with newer months about every 2-3 months
- The latest data which came out this month gave data up to 4/31/2021, which I can test against the prediction for 30 days

In [2]:
```python
#Loading prepped data
df = pd.read_csv('datasets/rollingsales_brooklyn.xls_prepped_bare.csv')
df.reset_index(drop=True, inplace=True)
df.sort_values('SALE DATE')
```

Out[2]:

| | TAX CLASS AT PRESENT | ZIP CODE | SALE PRICE | SALE DATE |
|---|---|---|---|---|
| **5697** | 2 | 11210.0 | 189000 | 2020-04-01 |
| **5626** | 2 | 11226.0 | 7185567 | 2020-04-01 |
| **5627** | 2 | 11226.0 | 7185567 | 2020-04-01 |
| **5635** | 2 | 11226.0 | 30644330 | 2020-04-01 |
| **5636** | 2 | 11226.0 | 14582474 | 2020-04-01 |
| **...** | ... | ... | ... | ... |
| **4577** | 2 | 11201.0 | 1717500 | 2021-03-31 |
| **10456** | 1 | 11232.0 | 1218500 | 2021-03-31 |
| **5706** | 2 | 11226.0 | 857000 | 2021-03-31 |
| **3039** | 2A | 11221.0 | 3755000 | 2021-03-31 |
| **11296** | 4 | 11249.0 | 75000 | 2021-03-31 |

11624 rows × 4 columns

# 3. Analyzing the data for trends/seasonality

I do the following steps here to help the data work with the modules:

```
1. Convert 'SALE DATE' column to datetime format
2. Create new dataframe with 'SALE DATE' as the index and 'SALE PRICE' a
s the column
3. Since we have multiple sales per day, I will aggregate the data into
 daily data by taking the daily average of sales
4. Check the data for any nulls/NaNs
     -Decide what to do for Nulls/NaNs
5. Use statsmodels to observe the data for trends and seasonality
```

***Observations:***

```
- NaN values came into the data after the data got aggregated.
    - Dropping these rows will result in skewing the data predictions
    - I decided to repalce the NaN values with 0 since no sales were don
e on that day
        -This also preserves the 365 day row length
```

In [4]:
```python
# 1. Convert 'SALE DATE' column to datetime format

df['SALE DATE'] = pd.to_datetime(df['SALE DATE'])
```

In [5]: 
```python
# 2 . Create new dataframe with 'SALE DATE' as the index and 'SALE PRICE' as the

df_price_date = pd.DataFrame(df, columns=['SALE DATE', 'SALE PRICE'])
df_price_date = df_price_date.set_index('SALE DATE')
df_price_date.head()
```

Out[5]:

|  | SALE PRICE |
| --- | --- |
| SALE DATE |  |
| 2020-04-28 | 1300000 |
| 2020-11-30 | 75000 |
| 2020-06-26 | 830000 |
| 2020-07-20 | 1188000 |
| 2021-02-22 | 990000 |

In [6]: 
```python
# 3. Group the sales data by daily average

df_price_date = df_price_date.resample('D').mean()
```

In [8]: 
```python
# 4. We see here number of rows went down 293. Why wasn't it 365 rows to represen
df_price_date.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 365 entries, 2020-04-01 to 2021-03-31
Freq: D
Data columns (total 1 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   SALE PRICE  293 non-null    float64
dtypes: float64(1)
memory usage: 5.7 KB
```

In [9]: 
```python
#Here we see that since we resampled by day, there are NaN values for the days th
df_price_date['SALE PRICE'].isna().sum()
```

Out[9]: 72

In [10]: 
```python
# 4. Instead of dropping the rows, I decided to fill NaN with 0 to reflect no sal
df_price_date['SALE PRICE'].fillna(0, inplace=True)
df_price_date
```

Out[10]:

|  | SALE PRICE |
| --- | --- |
| SALE DATE |  |
| 2020-04-01 | 3.977437e+06 |
| 2020-04-02 | 8.185471e+05 |
| 2020-04-03 | 1.815030e+06 |
| 2020-04-04 | 2.333627e+05 |
| 2020-04-05 | 0.000000e+00 |
| ... | ... |
| 2021-03-27 | 0.000000e+00 |
| 2021-03-28 | 0.000000e+00 |
| 2021-03-29 | 1.002984e+06 |
| 2021-03-30 | 1.058857e+06 |
| 2021-03-31 | 1.126519e+06 |

365 rows × 1 columns

In [11]:
```python
# 5. Checking for trends/seasonality
#Here I check the origional data against its 7-day weekly rolling window to see

df_price_date['roll_avg'] = df_price_date.rolling(window=7).mean()
df_price_date
```

Out[11]:

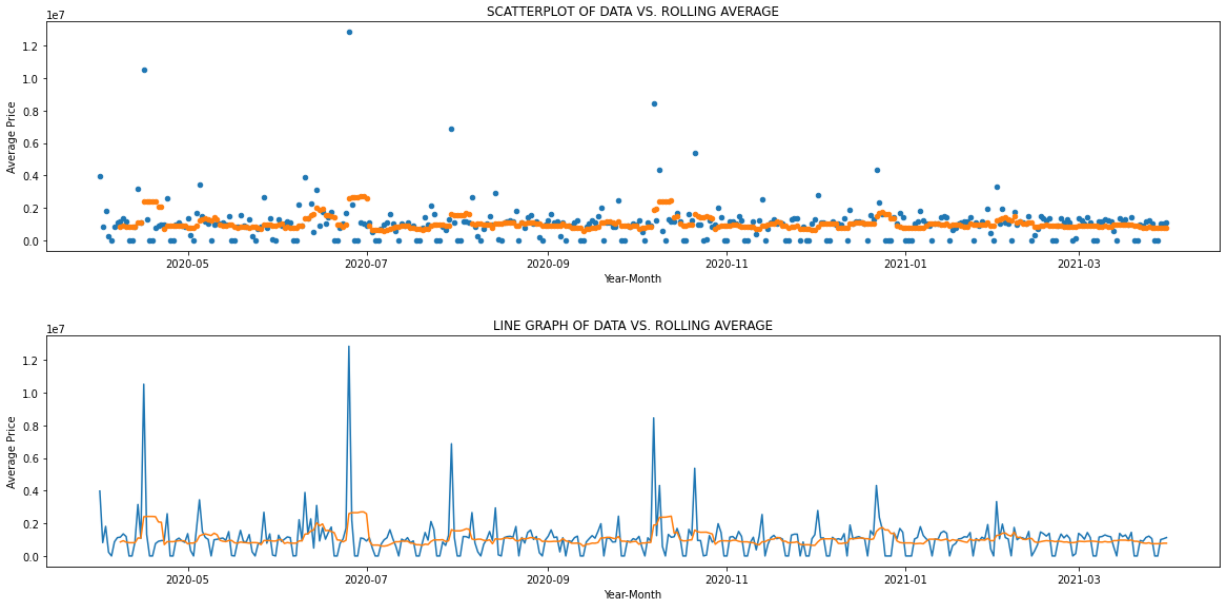| SALE DATE | SALE PRICE | roll_avg |
|---|---|---|
| 2020-04-01 | 3.977437e+06 | NaN |
| 2020-04-02 | 8.185471e+05 | NaN |
| 2020-04-03 | 1.815030e+06 | NaN |
| 2020-04-04 | 2.333627e+05 | NaN |
| 2020-04-05 | 0.000000e+00 | NaN |
| ... | ... | ... |
| 2021-03-27 | 0.000000e+00 | 752500.591285 |
| 2021-03-28 | 0.000000e+00 | 752500.591285 |
| 2021-03-29 | 1.002984e+06 | 760118.473764 |
| 2021-03-30 | 1.058857e+06 | 782883.683764 |
| 2021-03-31 | 1.126519e+06 | 781229.178651 |

365 rows × 2 columns

In [12]:
```python
#Plotting the 7-day rolling average against the origional data

plt.figure(figsize=(20, 4))
plt.title("SCATTERPLOT OF DATA VS. ROLLING AVERAGE")
plt.xlabel("Year-Month")
plt.ylabel("Average Price")

#s=20 to keep dots small in size
plt.scatter(df_price_date.index[:365], df_price_date['SALE PRICE'][:365], s=20)
plt.scatter(df_price_date.index[7:], df_price_date['roll_avg'][7:], s=20);
plt.figure(figsize=(20, 4))

plt.title("LINE GRAPH OF DATA VS. ROLLING AVERAGE")
plt.plot(df_price_date.index[:365], df_price_date['SALE PRICE'][:365])
plt.plot(df_price_date.index[7:], df_price_date['roll_avg'][7:]);
plt.xlabel("Year-Month")
plt.ylabel("Average Price")
```

Out[12]: Text(0, 0.5, 'Average Price')





Observation

- The spikes in the data where the price goes to the millions or tens of millions is due to buildings being bought.
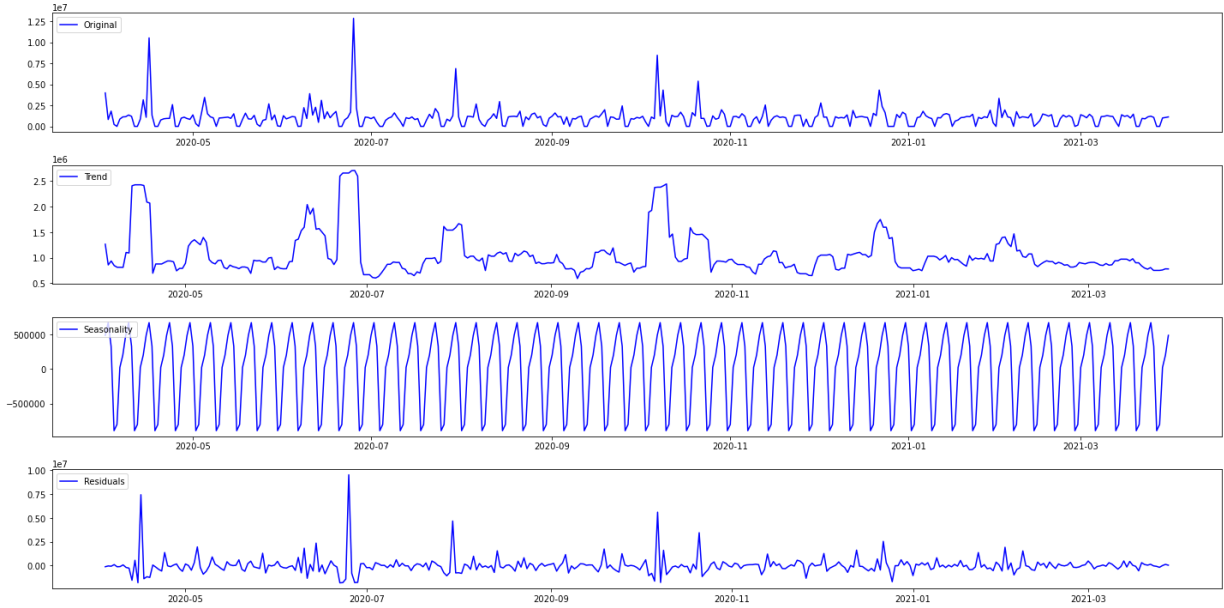- Other than that, the rest are residential properties well under a mill ion in price

In [13]:
```python
# Statsmodels decomposition

# Additive model was chosen here. It would not allow multiplicative with "0" valu
# Period of 7 for weekly lag

decomposition = seasonal_decompose(df_price_date['SALE PRICE'], model='additive'
observed = decomposition.observed
trend = decomposition.trend
seasonal = decomposition.seasonal
residual = decomposition.resid
```

In [14]:
```python
register_matplotlib_converters()
```

In [15]:
```python
plt.figure(figsize=(20,10))
plt.subplot(411)
plt.plot(observed, label='Original', color="blue")
plt.legend(loc='upper left')
plt.subplot(412)
plt.plot(trend, label='Trend', color="blue")
plt.legend(loc='upper left')
plt.subplot(413)
plt.plot(seasonal,label='Seasonality', color="blue")
plt.legend(loc='upper left')
plt.subplot(414)
plt.plot(residual, label='Residuals', color="blue")
plt.legend(loc='upper left')
plt.tight_layout()
```



**Observations:**

- Looks like there may be some seasonality every month

# 4. Dickey-Fuller Tests and preparing data for ARMA modeling

1. First I will run initial Augmented Dickey Fuller (ADF) test to check if the data is already stationary and does not have a unit root.
2. If the data fails the ADF test, I will induce stationarity using the following methods:
   - Differencing
   - Logging the data
   - Rolling mean subtraction

```
In [16]: # Initial test
         dftest = adfuller(df_price_date['SALE PRICE'])
         dfoutput = pd.Series(dftest[0:4], index=['Test Statistic','p-value','#Lags Used'
         for key,value in dftest[4].items():
             dfoutput['Critical Value (%s)'%key] = value
         print(dftest)
         print()
         print(dfoutput)
```

```
(-4.665875408357549, 9.749187516484523e-05, 15, 349, {'1%': -3.449226932880019,
'5%': -2.869857365438656, '10%': -2.571201085130664}, 10619.45547609437)

Test Statistic               -4.665875
p-value                       0.000097
#Lags Used                   15.000000
Number of Observations Used  349.000000
Critical Value (1%)          -3.449227
Critical Value (5%)          -2.869857
Critical Value (10%)         -2.571201
dtype: float64
```

## Augmented Dickey Fuller Test Goals:

Our goal is to induce stationarity and show that the data does not have a unit root.

ADF Test Null Hypothesis: The data has a unit root and is non-stationary.

Requirements for stationarity:

```
    1. If p-value <= 0.05: Reject the null hypothesis (H0), the data does no
    t have a unit root and is stationary.
        - If p-value > 0.05: Fail to reject the null hypothesis (H0), the da
    ta has a unit root and is non-stationary.
    2. If the Test Statistic is lower than the critical values, then reject
     the null hypothesis. Data does not have a unity root and is stationary
```

### Results of ADF Test

#### Test Statistic vs. Critical Values

- Initial test shows Test Statistic of **-4.665875**, this is greater than the critical values for 1% and 5%.
  - *We **REJECT** the null hypothesis! The data does not have a unit root and is stationary*

#### P-Value Analysis

- Our current p-value is **0.000097** which is REALLY close to zero.
  - *This means: p-value <= 0.05:*
  - *We **REJECT** the null hypothesis! The data does not have a unit root and is stationary*
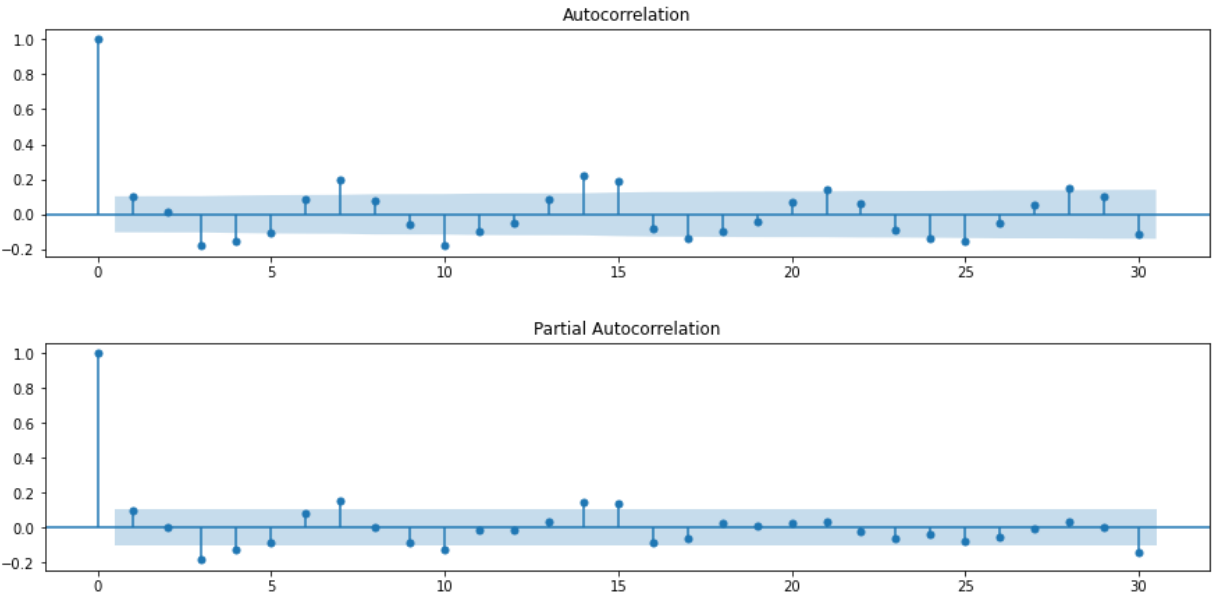
# 5. ARMA MODELING

Because ADF test shows data was stationary and does not have a unit root, we can proceed with ARMA model setup.

ACF and PACF will be used to determine the parameters.

In [17]:
```python
# ACF AND PACF

rcParams['figure.figsize'] = 15, 3
plot_acf(df_price_date['SALE PRICE'], lags=30, alpha=0.05);

rcParams['figure.figsize'] = 15, 3
plot_pacf(df_price_date['SALE PRICE'], lags=30, alpha=0.05);
```



In [19]:
```python
# Instantiate & fit model with statsmodels
#p = num lags - ACF
p = 17

# q = lagged forecast errors - PACF
q = 17

#d = number of differences - will compare differenced data RMSE with this model
# d=


# Fitting ARMA model and summary
ar = ARMA(df_price_date['SALE PRICE'],(p,q)).fit()
ar.summary()
```

Out[19]:

ARMA Model Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | SALE PRICE | **No. Observations:** | 365 |
| **Model:** | ARMA(17, 17) | **Log Likelihood** | -5588.777 |
| **Method:** | css-mle | **S.D. of innovations** | 1047550.052 |
| **Date:** | Sun, 20 Jun 2021 | **AIC** | 11249.554 |
| **Time:** | 15:26:45 | **BIC** | 11389.951 |
| **Sample:** | 04-01-2020 | **HQIC** | 11305.350 |
| | - 03-31-2021 | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **const** | 1.07e+06 | 6.32e+04 | 16.933 | 0.000 | 9.46e+05 | 1.19e+06 |
| **ar.L1.SALE PRICE** | -0.2998 | 0.266 | -1.125 | 0.260 | -0.822 | 0.222 |
| **ar.L2.SALE PRICE** | 0.0612 | 0.132 | 0.465 | 0.642 | -0.197 | 0.319 |

In [20]:
```python
#plot of ARMA model
plt.figure(figsize=(20,10))
fig, ax = plt.subplots()
# ax = df_price_date['SALE_PRICE_LOGGED'].plot(ax=ax, title='FORECAST')
ax = df_price_date['SALE PRICE'].plot(ax=ax, title='FORECAST',figsize=(15,5))
fig = ar.plot_predict(365, 395, dynamic=True, ax=ax, plot_insample=True)

handles, labels = ax.get_legend_handles_labels()
labels = ['SALE PRICE', 'FORECAST and 95% confidence interval']
ax.legend(handles, labels)

plt.show()
```

`<Figure size 1440x720 with 0 Axes>`



# 6. Error analysis of ARMA model

```
In [21]: predictions = list(ar.predict(276, 365))
         test = list(df_price_date['SALE PRICE'][275:365])

         print("\033[1m" + '\033[4m'+ 'Length of Predictions' + "\033[0m", ': ', len(pred:
         print("\033[1m" + '\033[4m'+ 'Length of Test data' + "\033[0m", ': ', len(test))

         #RMSE
         mse = sklearn.metrics.mean_squared_error(test, predictions)
         rmse = math.sqrt(mse)
         print("\033[1m" + '\033[4m'+ 'RMSE' + "\033[0m", ': ', rmse)

         #standard error
         stderr = ar.bse.const
         print("\033[1m" + '\033[4m'+ 'Standard Error' + "\033[0m", ': ', stderr)

         #plot of all
         plt.figure(figsize=(15,5))
         plt.plot(predictions, label='PREDICTIONS (90 days)', color='blue')
         plt.plot(test, label='TEST (90 days)', color='green')

         x=[0,90]
         y=[rmse,rmse]
         plt.plot(x,y, label=rmse, color='red')

         x=[0,90]
         y=[stderr,stderr]

         plt.plot(x,y, label=stderr, color='orange')

         plt.legend(loc='best')
         plt.xlabel("The last 90 days of data")
         plt.ylabel("Average Price (in hundreds of thousands)")
```
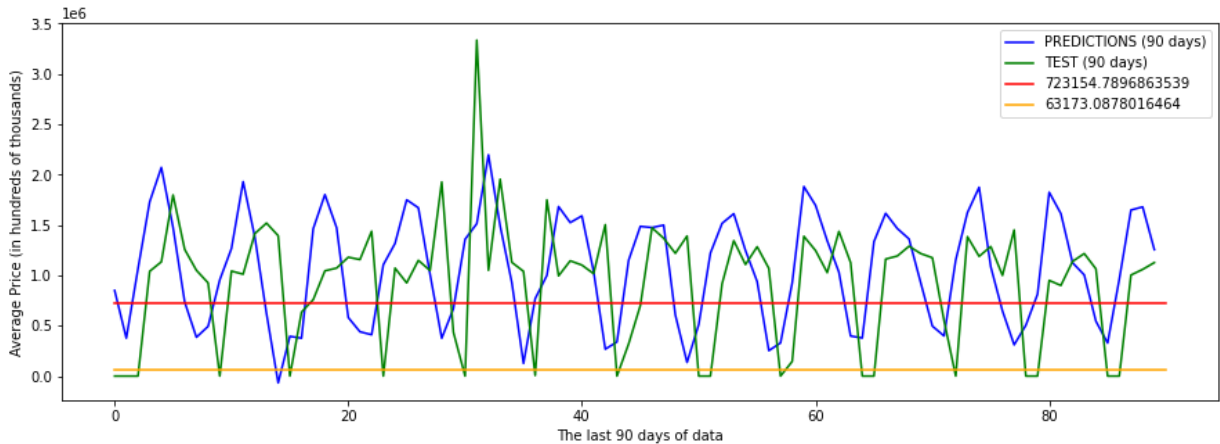
**Length of Predictions** :   90
**Length of Test data** :   90
**RMSE** :   723154.7896863539
**Standard Error** :   63173.0878016464

Out[21]: Text(0, 0.5, 'Average Price (in hundreds of thousands)')



***Observation:***

RMSE is not too high and not too low compared to the data. Does not indicate a bad fit nor a good fit

- RMSE is 723154.8
- Standard error is 63173.1

# 6a. Testing parameters to improve ARMA model

- I will try p of 9 per ACF
- I will try q of 9 per PACF
- I will try d = 7 to difference weekly

```
In [25]: # Instantiate & fit model with statsmodels
         #p = num lags - ACF
         p = 9

         # q = lagged forecast errors - PACF
         q = 9

         #d = number of differences
         d = 7


         # Fitting ARMA model and summary
         ar1 = ARMA(df_price_date['SALE PRICE'],(p,d,q)).fit()
         ar1.summary()
```

Out[25]:

ARMA Model Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | SALE PRICE | **No. Observations:** | 365 |
| **Model:** | ARMA(9, 7) | **Log Likelihood** | -5591.866 |
| **Method:** | css-mle | **S.D. of innovations** | 1071076.702 |
| **Date:** | Sun, 20 Jun 2021 | **AIC** | 11219.731 |
| **Time:** | 15:47:15 | **BIC** | 11289.929 |
| **Sample:** | 04-01-2020 | **HQIC** | 11247.629 |
| | - 03-31-2021 | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **const** | 1.07e+06 | 312.772 | 3420.002 | 0.000 | 1.07e+06 | 1.07e+06 |
| **ar.L1.SALE PRICE** | -0.1020 | 2.44e-05 | -4171.827 | 0.000 | -0.102 | -0.102 |
| **ar.L2.SALE PRICE** | 0.0734 | 6.49e-05 | 1131.002 | 0.000 | 0.073 | 0.074 |
| **ar.L3.SALE PRICE** | -0.0284 | 4.52e-05 | -628.503 | 0.000 | -0.028 | -0.028 |
| **ar.L4.SALE PRICE** | -0.0871 | 5.2e-05 | -1674.051 | 0.000 | -0.087 | -0.087 |
| **ar.L5.SALE PRICE** | -0.0650 | nan | nan | nan | nan | nan |
| **ar.L6.SALE PRICE** | -0.0123 | 4.83e-05 | -253.864 | 0.000 | -0.012 | -0.012 |
| **ar.L7.SALE PRICE** | 0.9356 | nan | nan | nan | nan | nan |
| **ar.L8.SALE PRICE** | 0.0120 | 7.6e-05 | 157.563 | 0.000 | 0.012 | 0.012 |
| **ar.L9.SALE PRICE** | -0.1141 | 6.51e-05 | -1753.442 | 0.000 | -0.114 | -0.114 |
| **ma.L1.SALE PRICE** | 0.0939 | 0.013 | 6.994 | 0.000 | 0.068 | 0.120 |
| **ma.L2.SALE PRICE** | 0.0164 | 0.013 | 1.222 | 0.222 | -0.010 | 0.043 |
| **ma.L3.SALE PRICE** | -0.0069 | 0.016 | -0.445 | 0.657 | -0.037 | 0.024 |
| **ma.L4.SALE PRICE** | 0.0802 | 0.015 | 5.405 | 0.000 | 0.051 | 0.109 |
| **ma.L5.SALE PRICE** | 0.0592 | 0.013 | 4.520 | 0.000 | 0.034 | 0.085 |
| **ma.L6.SALE PRICE** | -0.0196 | 0.015 | -1.345 | 0.178 | -0.048 | 0.009 |
| **ma.L7.SALE PRICE** | -0.9639 | 0.008 | -125.340 | 0.000 | -0.979 | -0.949 |

Roots

| | Real | Imaginary | Modulus | Frequency |
|---|---|---|---|---|
| **AR.1** | -0.9023 | -0.4356j | 1.0020 | -0.4284 |
| **AR.2** | -0.9023 | +0.4356j | 1.0020 | 0.4284 |
| **AR.3** | -0.2404 | -0.9746j | 1.0039 | -0.2885 |
| **AR.4** | -0.2404 | +0.9746j | 1.0039 | 0.2885 |
| **AR.5** | 0.6222 | -0.7828j | 1.0000 | -0.1431 |
| **AR.6** | 0.6222 | +0.7828j | 1.0000 | 0.1431 |
| **AR.7** | 1.0659 | -0.0000j | 1.0659 | -0.0000 |
| **AR.8** | -2.8105 | -0.0000j | 2.8105 | -0.5000 |

| | | | | |
|---|---|---|---|---|
| **AR.9** | 2.8907 | -0.0000j | 2.8907 | -0.0000 |
| **MA.1** | -0.9008 | -0.4342j | 1.0000 | -0.4285 |
| **MA.2** | -0.9008 | +0.4342j | 1.0000 | 0.4285 |
| **MA.3** | -0.2506 | -0.9681j | 1.0000 | -0.2903 |
| **MA.4** | -0.2506 | +0.9681j | 1.0000 | 0.2903 |
| **MA.5** | 0.6225 | -0.7826j | 1.0000 | -0.1431 |
| **MA.6** | 0.6225 | +0.7826j | 1.0000 | 0.1431 |
| **MA.7** | 1.0375 | -0.0000j | 1.0375 | -0.0000 |

In [26]:
```python
#plot of ARMA model
plt.figure(figsize=(20,10))
fig, ax = plt.subplots()
ax = df_price_date['SALE PRICE'].plot(ax=ax, title='FORECAST',figsize=(15,5))
fig = ar1.plot_predict(365, 395, dynamic=True, ax=ax, plot_insample=True)

handles, labels = ax.get_legend_handles_labels()
labels = ['SALE PRICE', 'FORECAST and 95% confidence interval']
ax.legend(handles, labels)

plt.show()
```

<Figure size 1440x720 with 0 Axes>



# 6a - Error Analysis of new model

In [27]:
```python
predictions = list(ar1.predict(276, 365))
test = list(df_price_date['SALE PRICE'][275:365])

print("\033[1m" + '\033[4m'+ 'Length of Predictions' + "\033[0m", ': ', len(pred:
print("\033[1m" + '\033[4m'+ 'Length of Test data' + "\033[0m", ': ', len(test))

#RMSE
mse = sklearn.metrics.mean_squared_error(test, predictions)
rmse = math.sqrt(mse)
print("\033[1m" + '\033[4m'+ 'RMSE' + "\033[0m", ': ', rmse)

#standard error
stderr = ar1.bse.const
print("\033[1m" + '\033[4m'+ 'Standard Error' + "\033[0m", ': ', stderr)

#plot of all
plt.figure(figsize=(15,5))
plt.plot(predictions, label='PREDICTIONS (90 days)', color='blue')
plt.plot(test, label='TEST (90 days)', color='green')

x=[0,90]
y=[rmse,rmse]
plt.plot(x,y, label=rmse, color='red')

x=[0,90]
y=[stderr,stderr]
plt.plot(x,y, label=stderr, color='orange')

plt.legend(loc='best')
plt.xlabel("The last 90 days of data")
plt.ylabel("Average Price (in hundreds of thousands)")
```
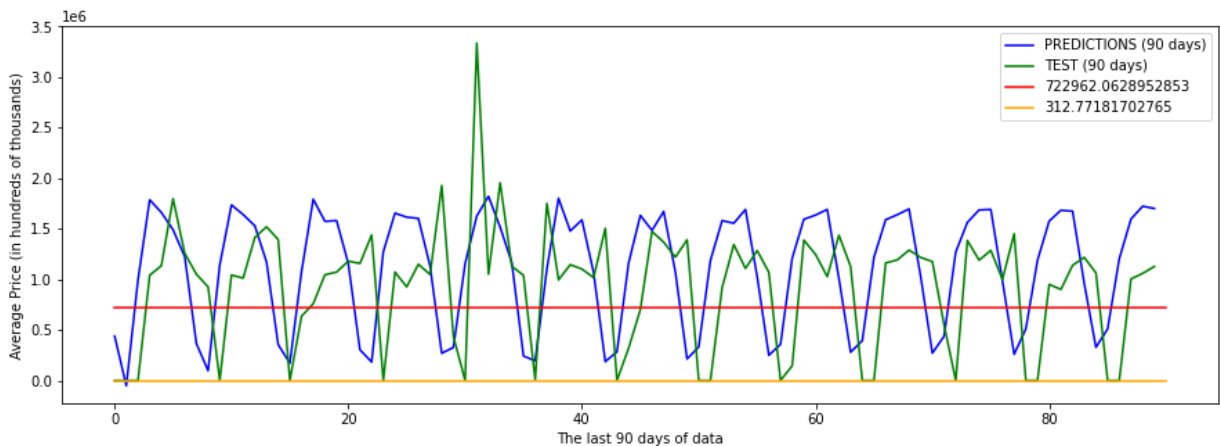
**Length of Predictions** :  90
**Length of Test data** :  90
**RMSE** :  722962.0628952853
**Standard Error** :  312.77181702765

Out[27]:  Text(0, 0.5, 'Average Price (in hundreds of thousands)')



## Observation:

- Here RMSE is lower than original model. We will stick with new model.

# 7. Comparing predictions with fresh data from June 2021 dataset (4/1/2021 - 4/31/2021)

## Here I do the following:

1. Load data with only specific columns to borough
   - Sale price
   - Sale data
   - Borough

2. Clean the data to get rid of issues when plotting/calculating errors

- This dataset was in .csv format, different from the origional rolling dataset
- I had to filter the data and change columns from strings to int
- Change 'SALE DATE' to datetime
- Resample the data to match origional rolling data
  - aggregate by day

3. Plot the new data versus the predicted data and calculate RMSE

In [28]:
```python
#Loading the data and reset the index

excel_df = pd.read_csv('NYC_Citywide_Rolling_Calendar_Sales.csv', usecols=['BOROI
excel_df = excel_df[excel_df['BOROUGH']=='BROOKLYN']
excel_df.reset_index(drop=True, inplace=True)
```

In [29]:
```python
#Fixes to the data

excel_df['SALE PRICE'] = excel_df['SALE PRICE'].str.replace(',','')
excel_df['SALE PRICE'] = excel_df['SALE PRICE'].astype(int)
excel_df['SALE DATE'] = pd.to_datetime(excel_df['SALE DATE'])
```

In [30]:
```python
#Create new dataframe and aggregate to days like I did with origional rolling dai

excel_price_date = pd.DataFrame(excel_df, columns=['SALE DATE', 'SALE PRICE'])
excel_price_date = excel_price_date.set_index('SALE DATE')

#aggregate by day
excel_price_date = excel_price_date.resample('D').mean()
```

In [31]:
```python
# Again, if I drop NaN here, it will change the dates which will affect the ploti
# I decide to fillna(0) similar to origional rolling data

excel_price_date = excel_price_date.fillna(0)
```

In [32]:
```python
excel_price_date.head()
```

Out[32]:

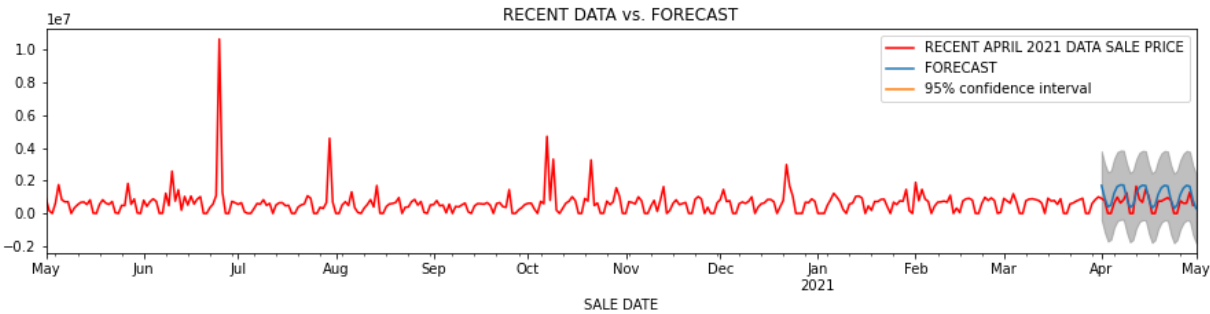| SALE DATE | SALE PRICE |
| --- | --- |
| 2020-05-01 | 9.110564e+05 |
| 2020-05-02 | 1.650000e+05 |
| 2020-05-03 | 0.000000e+00 |
| 2020-05-04 | 6.289992e+05 |
| 2020-05-05 | 1.752086e+06 |

In [35]:
```python
# Plotting the data versus the ar.plot_predict values

fig, ax = plt.subplots()

ax = excel_price_date['SALE PRICE'].plot(title='RECENT DATA vs. FORECAST', color
fig = ar1.plot_predict(365, 395, dynamic=True, ax = ax, plot_insample=True)

handles, labels = ax.get_legend_handles_labels()
labels = ['RECENT APRIL 2021 DATA SALE PRICE', 'FORECAST','95% confidence interva
ax.legend(handles, labels)

plt.show()
```
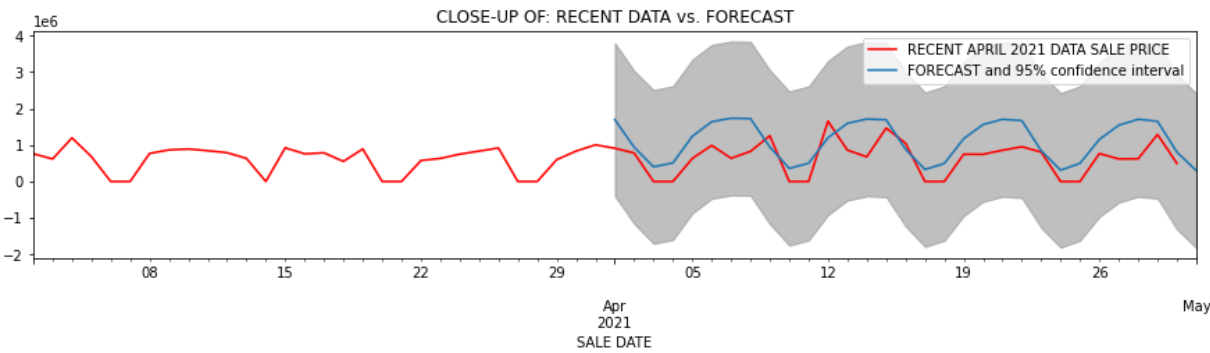


In [36]:
```python
# Plotting the data versus the ar.plot_predict values
#Here I do a close up

fig, ax = plt.subplots()

ax = excel_price_date['SALE PRICE'][305:365].plot(title='CLOSE-UP OF: RECENT DATA
fig = ar1.plot_predict(365, 395, dynamic=True, ax = ax, plot_insample=True)

handles, labels = ax.get_legend_handles_labels()
labels = ['RECENT APRIL 2021 DATA SALE PRICE', 'FORECAST and 95% confidence inter
ax.legend(handles, labels)

plt.show()
```



## Observation

- We see that the model looks like it fits well versus the test data of 4/1/2021 until 4/31/2021

In [39]:
```python
#RMSE, Standard error

# last 30 days of data
predictions = list(ar1.predict(365, 394))
test = list(excel_price_date['SALE PRICE'][335:365])

print("\033[1m" + '\033[4m'+ 'Length of Predictions' + "\033[0m", ': ', len(pred:
print("\033[1m" + '\033[4m'+ 'Length of Test data' + "\033[0m", ': ', len(test))

#RMSE
mse = sklearn.metrics.mean_squared_error(test, predictions)
rmse = math.sqrt(mse)
print("\033[1m" + '\033[4m'+ 'RMSE' + "\033[0m", ': ', rmse)

#standard error
stderr = ar1.bse.const
print("\033[1m" + '\033[4m'+ 'Standard Error' + "\033[0m", ': ', stderr)

#plot of all
plt.figure(figsize=(15,5))
plt.plot(predictions, label='PREDICTIONS (90 days)', color='blue')
plt.plot(test, label='TEST (90 days)', color='green')

x=[0,30]
y=[rmse,rmse]
plt.plot(x,y, label=rmse, color='red')

x=[0,30]
y=[stderr,stderr]
plt.plot(x,y, label=stderr, color='orange')

plt.legend(loc='best')
plt.xlabel("The last 30 days of data")
plt.ylabel("Average Price (in hundreds of thousands)")
```
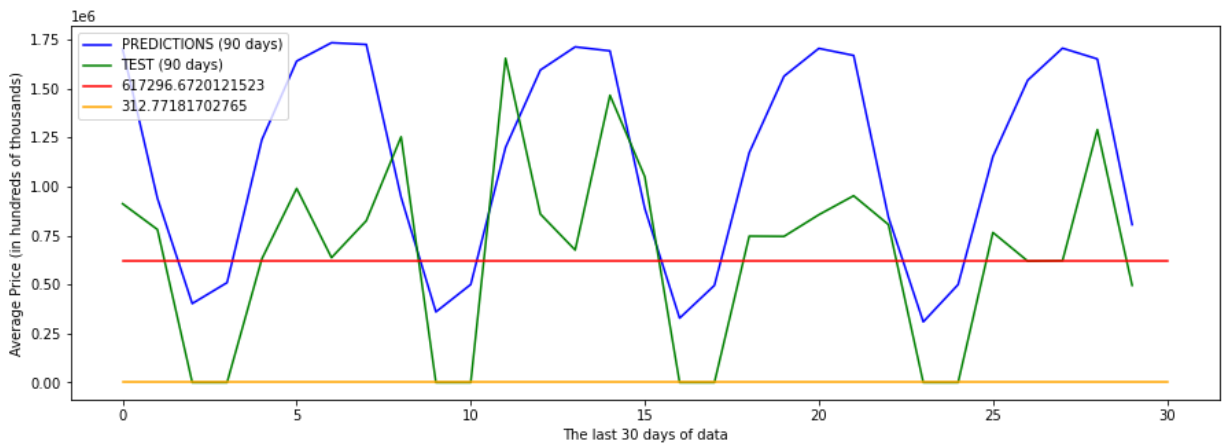
**Length of Predictions** :  30
**Length of Test data** :  30
**RMSE** :  617296.6720121523
**Standard Error** :  312.77181702765

Out[39]: Text(0, 0.5, 'Average Price (in hundreds of thousands)')



## Observation

1. RMSE is lower here when comparing the new month data with predicted values

# 8. Observations/Conclusions/Recommendations

1. The point of this analysis was to see if the borough was good to invest in
2. Based on the model:
   - We can enter to buy or exit to sell based on when the market will do well
3. The borough sales look predictable
   - There is predicable fluctuation in Brooklyn
4. We can look at the top 10 building permit heavy locations further

1. The point of this analysis was to see if the borough was good to invest in
2. Based on the model:
   - We can enter to buy or exit to sell based on when the market will do well
3. The borough sales look predictable
   - There is predicable fluctuation in Brooklyn
4. We can look at the top 10 building permit heavy locations further