

Analysis of 77084 zip code

Imports and loading csv

```
In [39]: #Imports
import pandas as pd
import numpy as np
from pandas.plotting import register_matplotlib_converters
import matplotlib.pyplot as plt
from matplotlib.pylab import rcParams
register_matplotlib_converters()

from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error

from scipy import stats
from random import gauss as gs
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

#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)
```

```
In [40]: df = pd.read_csv('Data Files/df_zillow_77084_prepped_fbprophet.csv')
```

```
In [41]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 265 entries, 0 to 264
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  -
0    ds      265 non-null      object
1    y        265 non-null      float64
dtypes: float64(1), object(1)
memory usage: 4.3+ KB
```

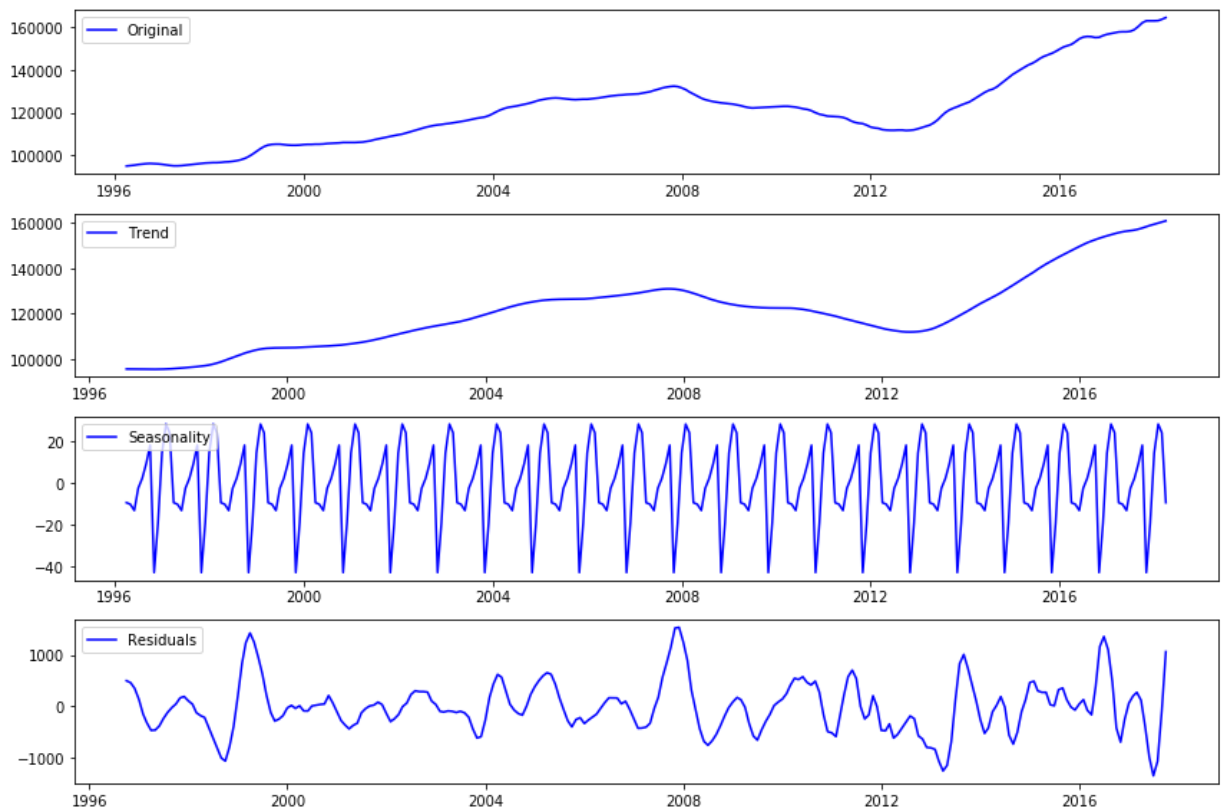
Decomposition and plots

```
In [42]: df.index = pd.to_datetime(df['ds'])  
df= df.drop(columns='ds')
```

```
In [43]: decomposition = seasonal_decompose(df.y)  
observed = decomposition.observed  
trend = decomposition.trend  
seasonal = decomposition.seasonal  
residual = decomposition.resid
```

```
In [44]: register_matplotlib_converters()
```

```
In [45]: plt.figure(figsize=(12,8))  
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()
```



I want to see if the data correlates with earlier data of

What to see in the data generated with earlier data on itself

1) Get rolling average with window of 4

- Couldn't see much with window of 1-3

2) Plot data against itself with rolling avg to see visual of the graph.

```
In [46]: df.rolling(window=2).mean().head()
```

Out[46]:

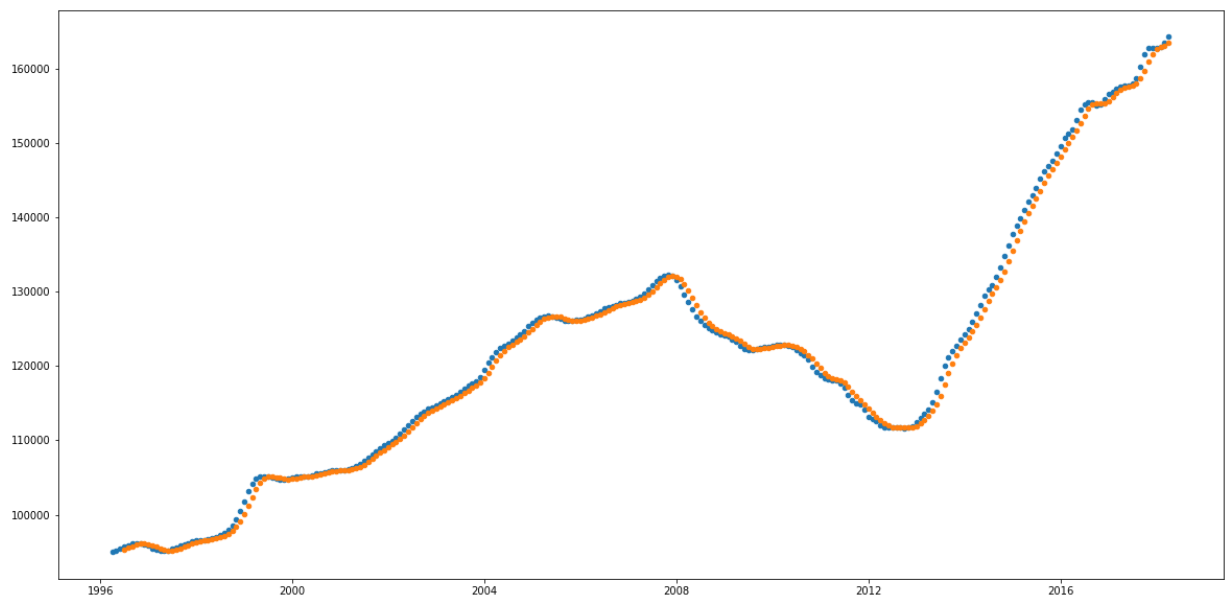
	y
ds	
1996-04-01	NaN
1996-05-01	95100.0
1996-06-01	95300.0
1996-07-01	95550.0
1996-08-01	95800.0

```
In [47]: df['roll_avg'] = df.rolling(window=4).mean()  
df.corr()
```

Out[47]:

	y	roll_avg
y	1.00000	0.99905
roll_avg	0.99905	1.00000

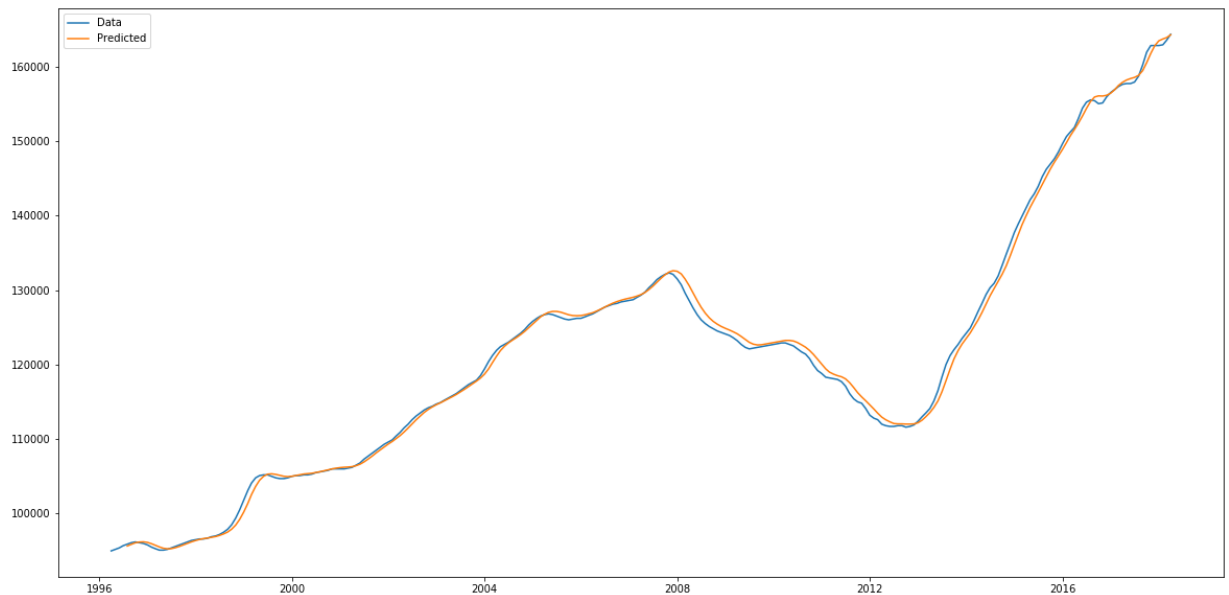
```
In [48]: plt.figure(figsize=(20, 10))  
plt.scatter(df.index[:265], df['y'][:265], s=20)  
plt.scatter(df.index[1:265], df['roll_avg'][1:265], s=20);
```



```
In [49]: lr = LinearRegression()
lr.fit(df[['roll_avg']][4:], df['y'][4:])
```

```
Out[49]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [50]: plt.figure(figsize=(20, 10))
plt.plot(df.index[:265], df['y'][:265], label='Data')
plt.plot(df.index[4:265], lr.predict(df[['roll_avg']][4:265]),
         label='Predicted')
plt.legend();
```



Upon brief visual look, there might be some correlation. We will set up for our model by using the Dickey-Fuller test and ACF (Auto-correlation) and PACF (Partial-autocorrelation)

Checking for Stationarity

```
In [51]: dfctest = adfuller(df.y)
dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lags Used']
for key,value in dfctest[4].items():
    dfoutput['Critical Value (%s)'%key] = value
print(dfctest)
print()
print(dfoutput)
```

```
(-1.0526382056244998, 0.7335852181016486, 13, 251, {'1%': -3.4566744514553016,
'5%': -2.8731248767783426, '10%': -2.5729436702592023}, 3301.1131948148727)
```

```
Test Statistic          -1.052638
p-value                 0.733585
#Lags Used              13.000000
Number of Observations Used 251.000000
Critical Value (1%)      -3.456674
Critical Value (5%)      -2.873125
Critical Value (10%)     -2.572944
dtype: float64
```

Dickey Fuller Test

- We see that test statistic value is -1.052638
- We see that the critical values are LESS than the test statistic. (-3.45, -2.87, -2.57)
- From just the baseline data, the test statistic I have is MORE than the critical value.
- We accept the null that the time series is not stationary!

P-Value analysis

1. If p-value > 0.05: Fail to reject the null hypothesis (H0), the data has a unit root and is non-stationary.
 - Our current p-value is 0.733585
 - This means: p-value > 0.05: Fail to reject the null hypothesis (H0), the data has a unit root and is non-stationary.
2. If p-value <= 0.05: Reject the null hypothesis (H0), the data does not have a unit root and is stationary.
 - Our goal is to make the data stationary

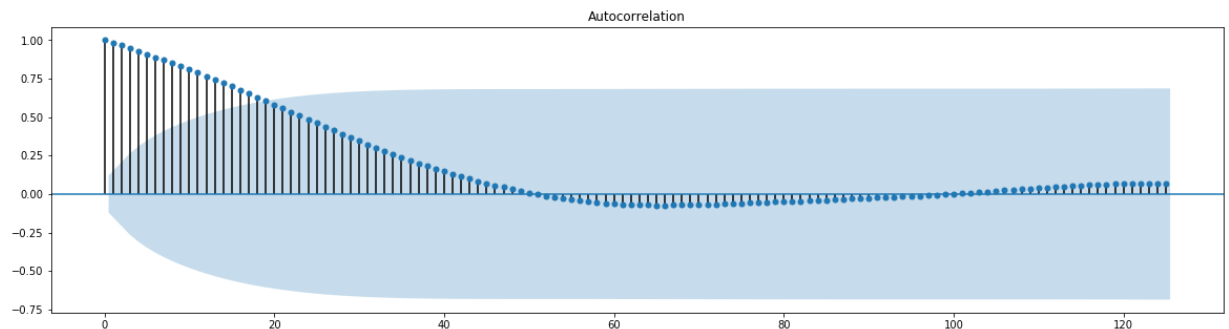
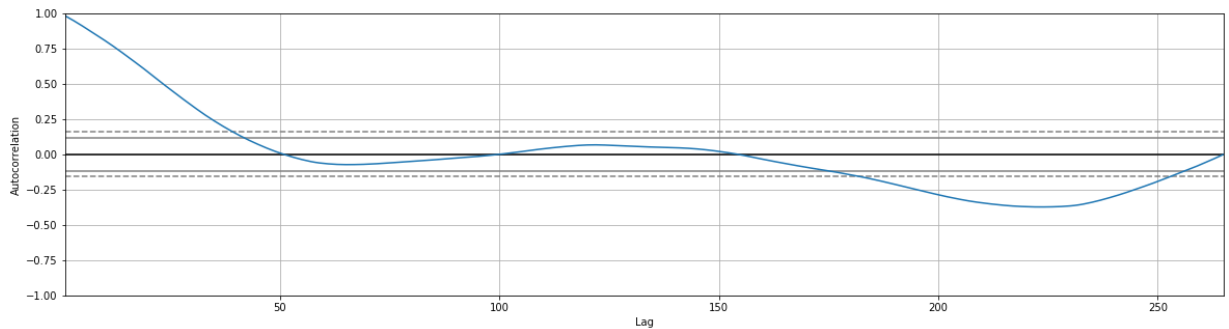
Auto-Correlation and Partial Auto-Correlation Check

```
In [52]: acf(df['y'], nlags=20, fft=False)
```

```
Out[52]: array([1.          , 0.98320081, 0.96581333, 0.94786212, 0.92920807,
0.9098244 , 0.88974082, 0.86943976, 0.84933628, 0.82935796,
0.80910934, 0.78824935, 0.76673984, 0.74466391, 0.72219496,
0.69940295, 0.67631589, 0.65306276, 0.6298105 , 0.60623204,
0.58207846])
```

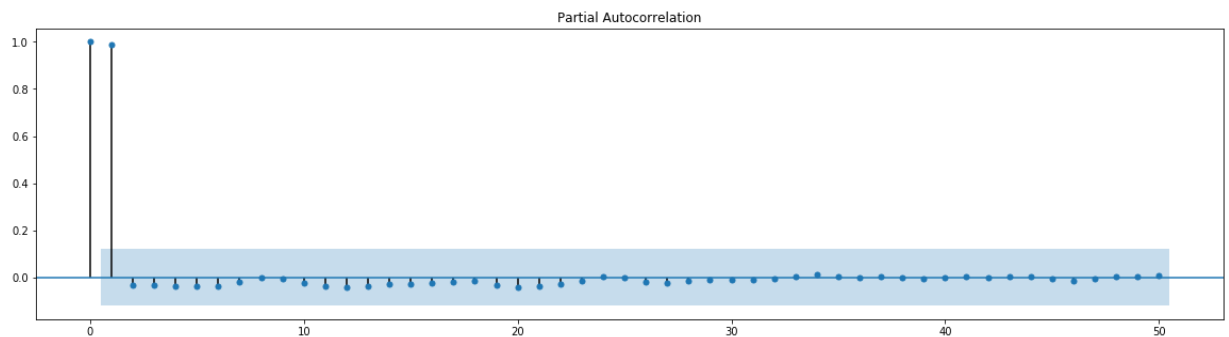
```
In [53]: #ACF using plotting
plt.figure(figsize=(20, 5))
pd.plotting.autocorrelation_plot(df['y']);

#Statsmodels ACF
rcParams['figure.figsize'] = 20, 5
plot_acf(df['y'], lags=125, alpha=0.05);
```



PACF

```
In [54]: pacf(df['y'], nlags=20)
rcParams['figure.figsize'] = 20, 5
plot_pacf(df['y'], lags=50, alpha=0.05);
```



Observations of ACF and PACF

We see the following:

- We know that the ACF describes the autocorrelation between an observation and another observation at a prior time step that includes direct and indirect dependence information.

- After about 18 lags, the line goes into our confidence interval (light blue area).
- This can be due to seasonality of every 18 months in our data.
- We know that the PACF only describes the direct relationship between an observation and its lag.
 - PACF cuts off after lags = 2
 - This means there are no correlations for lags beyond 2

**** Granted the data is not stationary, we will have to transform the data to make it stationary and satisfy the Dicky-Fuller test****

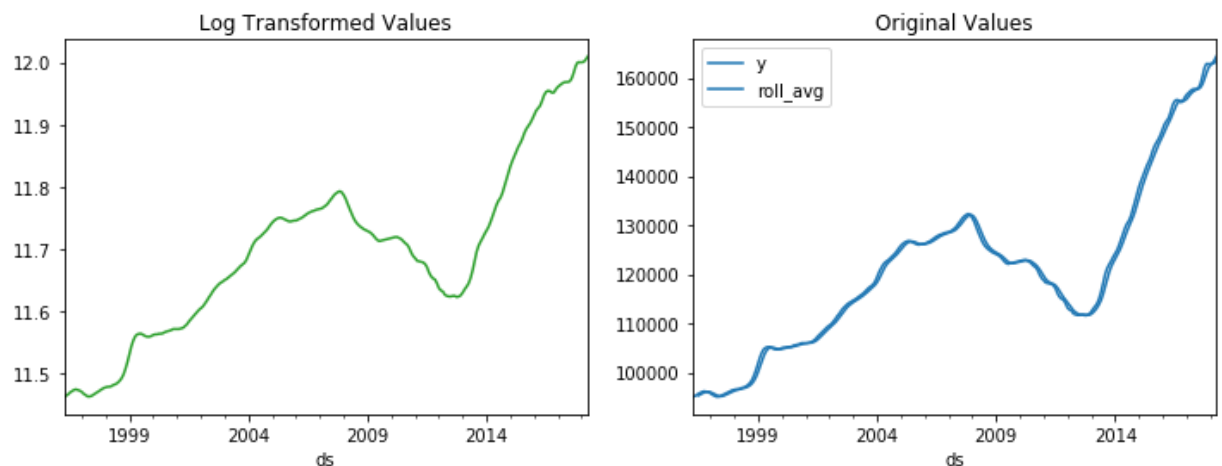
De-trending and transforming the data

1. I will try the following
 - Log transform
 - Subtract rolling mean
 - Run Dickey-Fuller test with each transform to see if I can reject/accept the null hypothesis
 - Null-Hypothesis for Dickey-Fuller test is: The null-hypothesis for the test is that the time series is not stationary. So if the test statistic is less than the critical value, we reject the null hypothesis and say that the series is stationary.

Log-transform on data and testing for stationarity

```
In [55]: logged_df = df['y'].apply(lambda x : np.log(x))
```

```
In [56]: ax1 = plt.subplot(121)
logged_df.plot(figsize=(12,4) ,color="tab:green", title="Log Transformed Values")
ax2 = plt.subplot(122)
df.plot(color="tab:blue", title="Original Values", ax=ax2);
```



```
In [57]: dfctest = adfuller(logged_df)
dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lags Used']
for key,value in dfctest[4].items():
    dfoutput['Critical Value (%s)'%key] = value
print(dfctest)
print()
print(dfoutput)
```

```
(-1.2051234223665641, 0.6713878707920429, 13, 251, {'1%': -3.4566744514553016,
'5%': -2.8731248767783426, '10%': -2.5729436702592023}, -2519.397286409148)
```

```
Test Statistic          -1.205123
p-value                  0.671388
#Lags Used               13.000000
Number of Observations Used 251.000000
Critical Value (1%)      -3.456674
Critical Value (5%)      -2.873125
Critical Value (10%)     -2.572944
dtype: float64
```

Observations after log-transform

1. Test Statistic is still larger than Critical Values. We accept the null-hypothesis that the time series is not stationary!
 - Test Statistic -1.205123
 - Critical Value (1%) -3.456360
 - Critical Value (5%) -2.872987
 - Critical Value (10%) -2.572870
2. P value is 0.671388
 - This means: p-value > 0.05: Fail to reject the null hypothesis (H0), the data has a unit root and is non-stationary.

Subtracting Rolling Mean from logged data and a better window size

```
In [58]: #Try breakdown with data minus rollmean. It looks like there is seasonality but
# Window of 11

logged_df_roll_mean = logged_df.rolling(window=11).mean()
logged_df_minus_roll_mean1 = logged_df - logged_df_roll_mean
logged_df_minus_roll_mean1.dropna(inplace=True)
```



```
In [59]: logged_df_minus_roll_mean1.head()
```

```
Out[59]: ds
1997-02-01    -0.002274
1997-03-01    -0.004657
1997-04-01    -0.006663
1997-05-01    -0.006376
1997-06-01    -0.004849
Name: y, dtype: float64
```

```
In [60]: dfctest = adfuller(logged_df_minus_roll_mean1)
# Extract and display test results in a user friendly manner
dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lags Used']
for key,value in dfctest[4].items():
    dfoutput['Critical Value (%s)'%key] = value
print(dfctest)
print()
print(dfoutput)
```

```
(-2.939616555350866, 0.04093128873141314, 13, 241, {'1%': -3.4577787098622674,
'5%': -2.873608704758507, '10%': -2.573201765981991}, -2458.8607786167463)
```

Test Statistic	-2.939617
p-value	0.040931
#Lags Used	13.000000
Number of Observations Used	241.000000
Critical Value (1%)	-3.457779
Critical Value (5%)	-2.873609
Critical Value (10%)	-2.573202
dtype:	float64

--- Observations from Dickey Fuller Test ---

We are getting close.

- 2.94 Test statistic which is within the range of the critical values, lower than 5% and 10%, I can attempt ARMA
- p value is ~0.04
- I can reject null hypothesis since $p < 0.05$

Differencing the data and re-running Dickey Fuller

```
In [61]: logged_df_diff = logged_df.diff(periods=1)
```

```
In [62]: logged_df_diff_roll_mean = logged_df_diff.rolling(window=11).mean()
logged_df_diff_roll_mean1 = logged_df_diff - logged_df_diff_roll_mean
logged_df_diff_roll_mean1.dropna(inplace=True)
```

```
In [63]: logged_df_diff_roll_mean1.head()
```

```
Out[63]: ds
1997-03-01    -0.002383
1997-04-01    -0.002005
1997-05-01     0.000286
1997-06-01     0.001527
1997-07-01     0.002574
Name: y, dtype: float64
```

```
In [64]: dfctest = adfuller(logged_df_diff_roll_mean1)
# Extract and display test results in a user friendly manner
dfoutput = pd.Series(dfctest[0:4], index=['Test Statistic', 'p-value', '#Lags Used',
for key,value in dfctest[4].items():
    dfoutput['Critical Value (%)'%key] = value
print(dfctest)
print()
print(dfoutput)
```

```
(-4.428505432438004, 0.00026418725013481774, 16, 237, {'1%': -3.4582467982399105, '5%': -2.8738137461081323, '10%': -2.5733111490323846}, -2445.992120196823)
```

Test Statistic	-4.428505
p-value	0.000264
#Lags Used	16.000000
Number of Observations Used	237.000000
Critical Value (1%)	-3.458247
Critical Value (5%)	-2.873814
Critical Value (10%)	-2.573311

dtype: float64

--- Observations of Dickey-Fuller Test ---

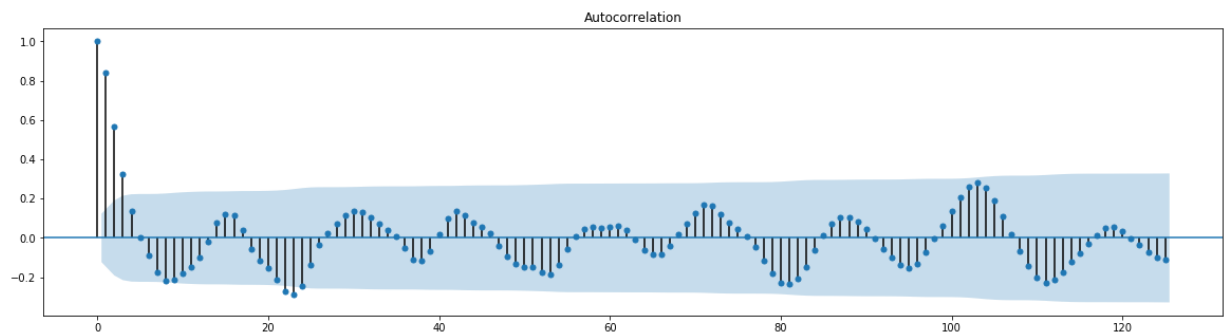
- We see Test Statistic is less than the Critical values, this satisfies the stationarity assumption. We can reject the null and say series is stationary.

- Test Statistic	-4.428505
- Critical Value (1%)	-3.458247
- Critical Value (5%)	-2.873814
- Critical Value (10%)	-2.573311

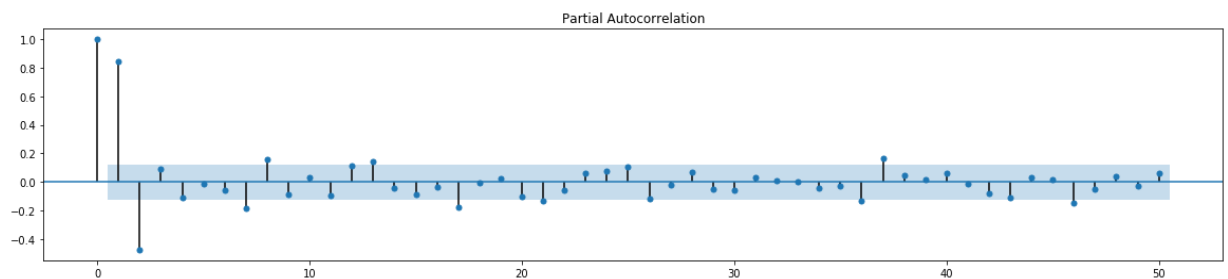
- We see that p-value = 0.000264. Since $p \leq 0.05$, I can reject the null hypothesis (H_0 = series is non-stationary). The data does not have a unit root and is stationary.

ACF and PACF

```
In [65]: #Statsmodels ACF
rcParams['figure.figsize'] = 20, 5
plot_acf(logged_df_diff_roll_mean1, lags=125, alpha=0.05);
```



```
In [66]: #PACF plot
rcParams['figure.figsize'] = 20, 4
plot_pacf(logged_df_diff_roll_mean1, lags=50, alpha=0.05);
```



--- Observations of ACF and PACF ---

1. After about 4 lags, the line goes into our confidence interval (light blue area).
 - This can be due to seasonality of every 4 months in our data.
2. PACF trails off after 2-3 lags.
 - Also slight slight sinusoidal behavior but nothing crazy
 - This means there are no high correlations for lags beyond 2-3
3. Based on above information and that the data is stationary, we can use the p and q values for the ARMA model
 - $p = 4$ (per ACF)
 - $q = 2$ (per PACF)

ARMA Modeling

```
In [67]: # Instantiate & fit model with statsmodels
#p = num lags - ACF
p = 4

# q = lagged forecast errors - PACF
q = 2

# Fitting ARMA model and summary
ar = ARMA(logged_df_minus_roll_mean1, (p, q)).fit()
ar.summary()
```

Out[67]: ARMA Model Results

Dep. Variable:	y	No. Observations:	255
Model:	ARMA(4, 2)	Log Likelihood	1323.888
Method:	csm-mle	S.D. of innovations	0.001
Date:	Thu, 29 Apr 2021	AIC	-2631.777
Time:	14:10:07	BIC	-2603.447
Sample:	02-01-1997	HQIC	-2620.381
	- 04-01-2018		

	coef	std err	z	P> z	[0.025	0.975]
const	0.0105	0.005	1.996	0.047	0.000	0.021
ar.L1.y	0.8239	0.119	6.901	0.000	0.590	1.058
ar.L2.y	0.7635	0.148	5.157	0.000	0.473	1.054
ar.L3.y	-0.6511	0.160	-4.058	0.000	-0.966	-0.337
ar.L4.y	0.0179	0.116	0.154	0.877	-0.210	0.245
ma.L1.y	1.4562	0.103	14.088	0.000	1.254	1.659
ma.L2.y	0.5610	0.104	5.370	0.000	0.356	0.766

Roots

	Real	Imaginary	Modulus	Frequency
AR.1	-1.1264	+0.0000j	1.1264	0.5000
AR.2	1.1722	+0.0000j	1.1722	0.0000
AR.3	1.2047	+0.0000j	1.2047	0.0000
AR.4	35.0787	+0.0000j	35.0787	0.0000
MA.1	-1.2979	-0.3131j	1.3351	-0.4623
MA.2	-1.2979	+0.3131j	1.3351	0.4623

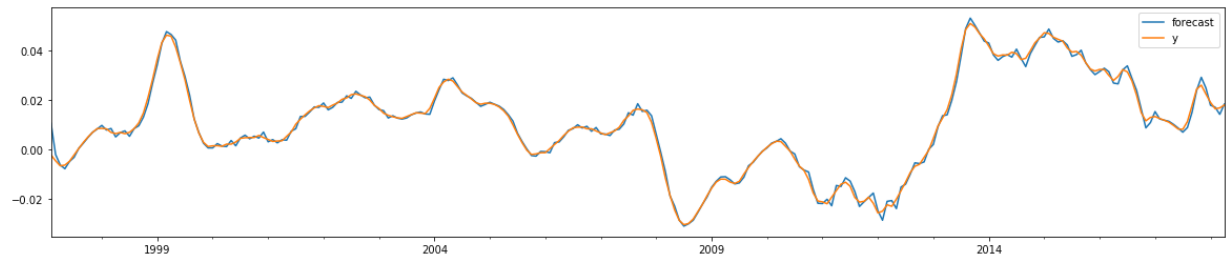
```
In [68]: r2_score(logged_df_minus_roll_mean1, ar.predict())
```

Out[68]: 0.992967791313284

- This means that 99.2 percent of the variation in the y data is due to variation in the x data
- This might indicate overfitting, but we chose our params from a stationary time series ACF and PACF.

-Future work: investigate more tweaks to the model

```
In [69]: #plot of ARMA model  
fig = ar.plot_predict()
```



Change the params, maybe it will affect r^2

In [70]: *# Try p = 4 and q = 3*

```
# Instantiate & fit model with statsmodels
#p = num lags - ACF
p = 4

# q = lagged forecast errors - PACF
q = 3

# Fitting ARMA model and summary
ar = ARMA(logged_df_minus_roll_mean1,(p, q)).fit()
ar.summary()
```

Out[70]:

ARMA Model Results

Dep. Variable:	y	No. Observations:	255
Model:	ARMA(4, 3)	Log Likelihood	1324.427
Method:	csmle	S.D. of innovations	0.001
Date:	Thu, 29 Apr 2021	AIC	-2630.854
Time:	14:10:08	BIC	-2598.983
Sample:	02-01-1997	HQIC	-2618.034
	- 04-01-2018		

	coef	std err	z	P> z	[0.025	0.975]
const	0.0105	0.006	1.759	0.080	-0.001	0.022
ar.L1.y	1.5524	0.262	5.921	0.000	1.039	2.066
ar.L2.y	0.1763	0.266	0.662	0.508	-0.345	0.698
ar.L3.y	-1.2503	0.242	-5.161	0.000	-1.725	-0.775
ar.L4.y	0.5103	0.161	3.173	0.002	0.195	0.826
ma.L1.y	0.7359	0.267	2.755	0.006	0.212	1.259
ma.L2.y	-0.4861	0.382	-1.272	0.205	-1.235	0.263
ma.L3.y	-0.3900	0.167	-2.341	0.020	-0.717	-0.063

Roots

	Real	Imaginary	Modulus	Frequency
AR.1	-1.1134	-0.0000j	1.1134	-0.5000
AR.2	1.0689	-0.0000j	1.0689	-0.0000
AR.3	1.2473	-0.3013j	1.2832	-0.0377
AR.4	1.2473	+0.3013j	1.2832	0.0377
MA.1	1.4025	-0.0000j	1.4025	-0.0000
MA.2	-1.3244	-0.2721j	1.3521	-0.4677
MA.3	-1.3244	+0.2721j	1.3521	0.4677

```
In [71]: r2_score(logged_df_minus_roll_mean1, ar.predict())
```

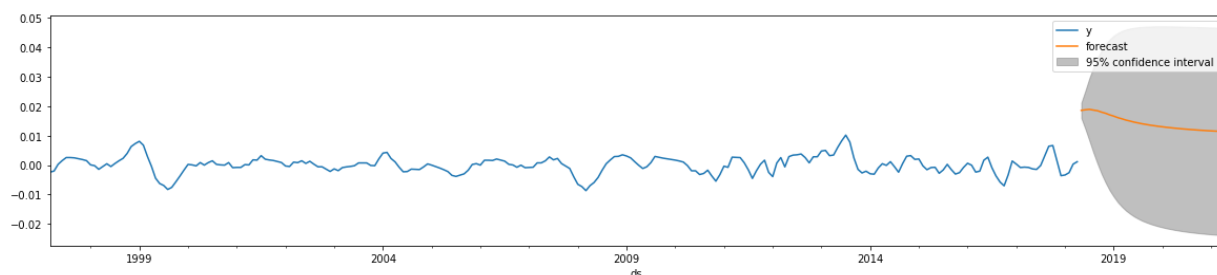
```
Out[71]: 0.992992404006822
```

Forecasting

```
In [81]: logged_df_diff_roll_mean1
```

```
Out[81]: ds
1997-03-01    -0.002383
1997-04-01    -0.002005
1997-05-01     0.000286
1997-06-01     0.001527
1997-07-01     0.002574
...
2017-12-01    -0.003588
2018-01-01    -0.003356
2018-02-01    -0.002566
2018-03-01     0.000335
2018-04-01     0.001154
Name: y, Length: 254, dtype: float64
```

```
In [82]: #plot of ARMA model
fig, ax = plt.subplots()
ax = logged_df_diff_roll_mean1.plot(ax=ax)
fig = ar.plot_predict('2018-05-01', '2021-04-01', dynamic=True, ax=ax, plot_insample=False)
plt.show()
```



We see prices are predicted to be higher. Could be a good investment zip code.

```
In [73]: #Future work, try SARIMAX prediction
# Need to install modules properly for SARIMAX to work.
```