



Jacob_s

[Follow](#)

Deloitte- consultant(Data science and data Engineering) <https://www.linkedin.com/in/jacob-stallone-ba4b26149/>

Nov 9, 2017 · 10 min read

Time Series Forecast : A basic introduction using Python.

Time series data is an important source for information and strategy used in various businesses. From a conventional finance industry to education industry, they play a major role in understanding a lot of details on specific factors with respect to time. I recently learnt the importance of Time series data in the telecommunication industry and wanted to brush up on my time series analysis and forecasting information. So I decided to work through a simple example using python and I have explained all the details in this blog.

Time series forecasting is basically the machine learning modeling for Time Series data (years, days, hours...etc.)for predicting future values using Time Series modeling .This helps if your data is serially correlated.

Loading and Handling Time Series in Pandas.

I will be using python in jupyter notebook. Pandas in python has libraries that are specific to handling time series object .You can check out this documentation for more details .

Let's start with the Preliminaries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

To load the data- I have provided the link to my GitHub where the dataset and the code is available. I will be using the AirPassenger dataset from. Are you ready?

```
data = pd.read_csv('AirPassengers.csv')
print data.head()
print '\n Data Types:'
print data.dtypes
```

Looking at the output .

```
1
2 data = pd.read_csv('AirPassengers.csv')
3 print data.head()
4 print '\n Data Types:'
5 print data.dtypes
```

	Month	#Passengers
0	1949-01	112
1	1949-02	118
2	1949-03	132
3	1949-04	129
4	1949-05	121

```
Data Types:
Month          object
#Passengers    int64
dtype: object
```

The data contains a particular month and number of passengers travelling in that month .The data type here is object (month) Let's convert it into a Time series object and use the Month column as our index.

```
from datetime import datetime
con=data['Month']
data['Month']=pd.to_datetime(data['Month'])
data.set_index('Month', inplace=True)
#check datatype of index
data.index
```

```
DatetimeIndex(['1949-01-01', '1949-02-01', '1949-03-01', '1949-04-01',
               '1949-05-01', '1949-06-01', '1949-07-01', '1949-08-01',
               '1949-09-01', '1949-10-01',
               ...,
               '1960-03-01', '1960-04-01', '1960-05-01', '1960-06-01',
               '1960-07-01', '1960-08-01', '1960-09-01', '1960-10-01',
               '1960-11-01', '1960-12-01'],
              dtype='datetime64[ns]', name='Month', length=144, freq=None)
```

You can see that now the data type is 'datetime64[ns]'. Now let's just make it into a series rather than a data frame (this would make it easier for the blog explanation)

```
#convert to time series:
ts = data['#Passengers']
ts.head(10)
```

Let's explore the various properties date-time based index:

```
1 #1. Specific the index as a string constant:
2 ts['1949-01-01']
```

112

```
1 #2. Import the datetime library and use 'datetime' function:
2 from datetime import datetime
3 ts[datetime(1949,1,1)]
```

112

#Get range:

```
1 #1. Specify the entire range:
2 ts['1949-01-01':'1949-05-01']
```

```
Month
1949-01-01    112
1949-02-01    118
1949-03-01    132
1949-04-01    129
1949-05-01    121
Name: #Passengers, dtype: int64
```

```
1 #2. Use ':' if one of the indices is at ends:
2 ts[':1949-05-01']
```

```
Month
1949-01-01    112
1949-02-01    118
1949-03-01    132
1949-04-01    129
1949-05-01    121
Name: #Passengers, dtype: int64
```

```

1 #All rows of 1962:
2 ts['1949']

```

Month	
1949-01-01	112
1949-02-01	118
1949-03-01	132
1949-04-01	129
1949-05-01	121
1949-06-01	135
1949-07-01	148
1949-08-01	148
1949-09-01	136
1949-10-01	119
1949-11-01	104
1949-12-01	118

Name: #Passengers, dtype: int64

STATIONARITY

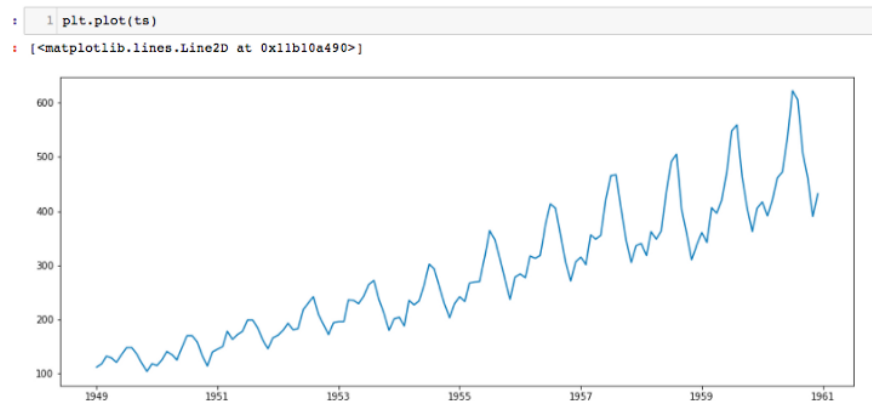
This is a very important concept in Time Series Analysis. In order to apply a time series model, it is important for the Time series to be stationary; in other words all its statistical properties (mean, variance) remain constant over time. This is done basically because if you take a certain behavior over time, it is important that this behavior is same in the future in order for us to forecast the series. There are a lot of statistical theories to explore stationary series than non-stationary series. (Thus we can bring the fight to our home ground!)

In practice we can assume the series to be stationary if it has constant statistical properties over time and these properties can be:

- constant mean
- constant variance
- an auto co-variance that does not depend on time.

These details can be easily retrieved using stat commands in python.

The best way to understand you stationarity in a Time Series is by eyeballing the plot:



It's clear from the plot that there is an overall increase in the trend, with some seasonality in it.

I have written a function for it as I will be using it quite often in this Time series explanation. But before we get to that, let me explain all the concepts in the function.

Plotting Rolling Statistics : The function will plot the moving mean or moving Standard Deviation. This is still visual method

NOTE: moving mean and moving standard deviation—At any instant 't', we take the mean/std of the last year which in this case is 12 months)

Dickey-fuller Test : This is one of the statistical tests for checking stationarity. First we consider the null hypothesis: the time series is non-stationary. The result from the test will contain the test statistic and critical value for different confidence levels. The idea is to have Test statistics less than critical value, in this case we can reject the null hypothesis and say that this Time series is indeed stationary (the force is strong with this one !!)

More details for Dickey fuller Test.

Function details:

- mean
- Standard deviation (instead of variance)
- Plot original series

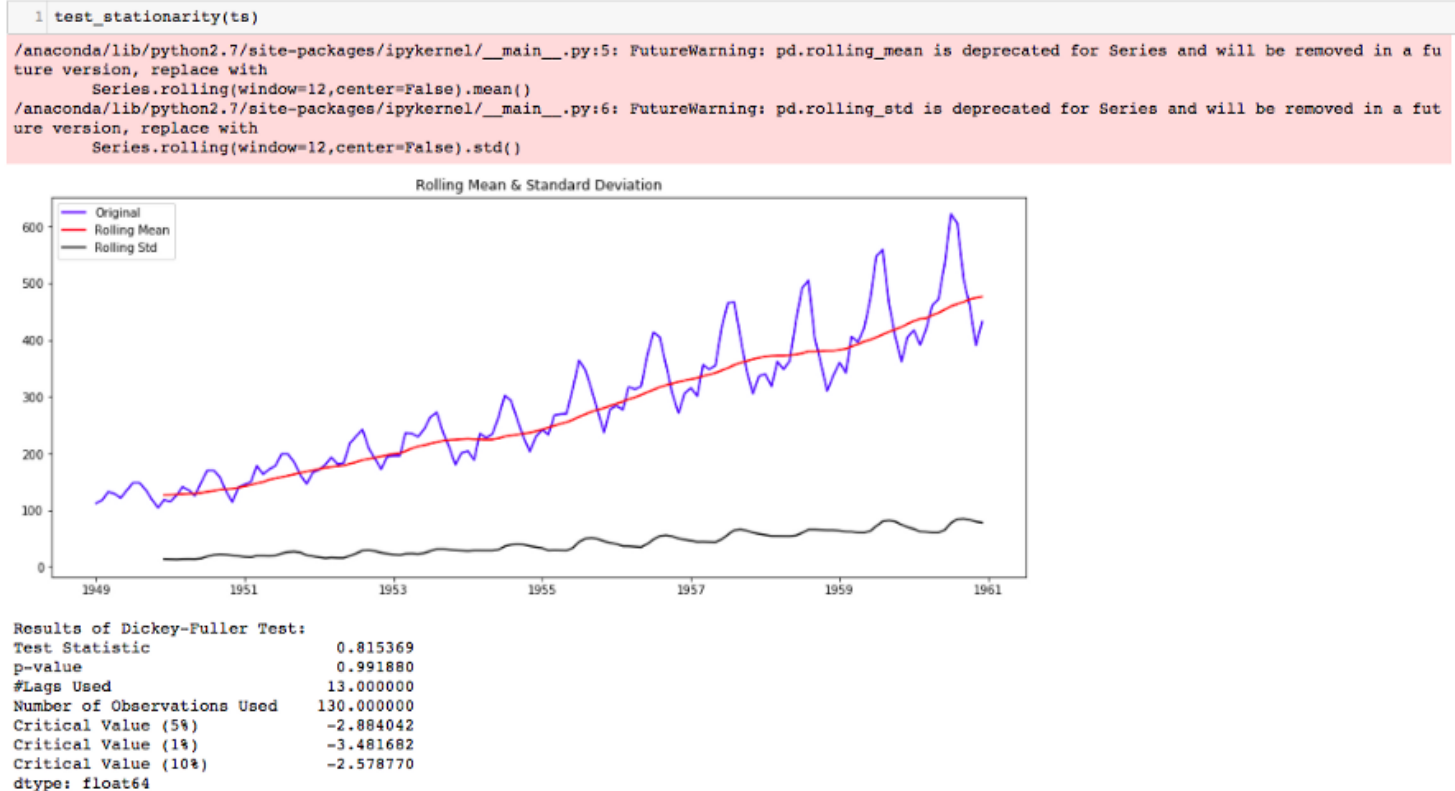
- Plot mean
- Plot std
- Plot Dickey-Fuller test

```
from statsmodels.tsa.stattools import adfuller
def test_stationarity(timeseries):

    #Determining rolling statistics
    rolmean = pd.rolling_mean(timeseries, window=12)
    rolstd = pd.rolling_std(timeseries, window=12)

    #Plot rolling statistics:
    plt.plot(timeseries, color='blue',label='Original')
    plt.plot(rolmean, color='red', label='Rolling Mean')
    plt.plot(rolstd, color='black', label = 'Rolling Std')
    plt.legend(loc='best')
    plt.title('Rolling Mean & Standard Deviation')
    plt.show()
    #Perform Dickey-Fuller test:
    print 'Results of Dickey-Fuller Test:'
    dftest = adfuller(timeseries, autolag='AIC')
    dfoutput = pd.Series(dftest[0:4], index=['Test
Statistic','p-value','#Lags Used','Number of Observations
Used'])
    for key,value in dftest[4].items():
        dfoutput['Critical Value (%s)'%key] = value
    print dfoutput
```

Now let's parse our time series data into this function:



This is not stationary because :

- mean is increasing even though the std is small.
- Test stat is > critical value.
- Note: the signed values are compared and the absolute values.

MAKING THE TIME SERIES STATIONARY

There are two major factors that make a time series non-stationary.
They are:

- Trend: non-constant mean
- Seasonality: Variation at specific time-frames

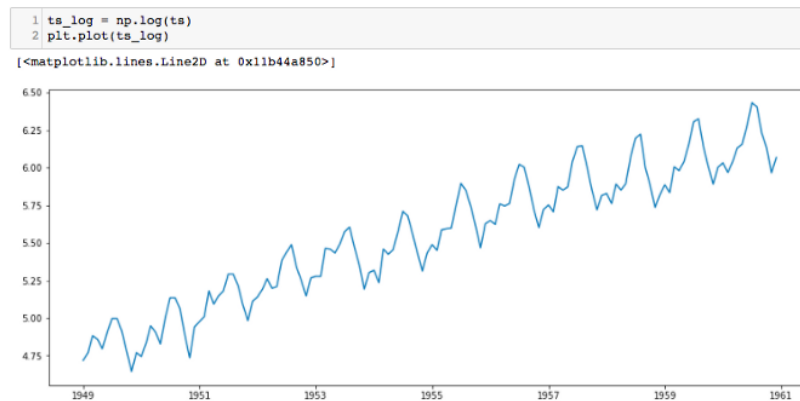
The basic idea is to model the trend and seasonality in this series, so we can remove it and make the series stationary. Then we can go ahead and apply statistical forecasting to the stationary series. And finally we

can convert the forecasted values into original by applying the trend and seasonality constraints back to those that we previously separated.

Let's start by working on the trend piece.

Trend

The first step is to reduce the trend using transformation, as we can see here that there is a strong positive trend. These transformation can be log, sq-rt, cube root etc . Basically it penalizes larger values more than the smaller. In this case we will use the logarithmic transformation.



There is some noise in realizing the forward trend here. There are some methods to model these trends and then remove them from the series. Some of the common ones are:

- Smoothing: using rolling/moving average
- Aggression: by taking the mean for a certain time period (year/month)

I will be using Smoothing here.

Smoothing:

In smoothing we usually take the past few instances (rolling estimates) We will discuss two methods under smoothing- Moving average and Exponentially weighted moving average.

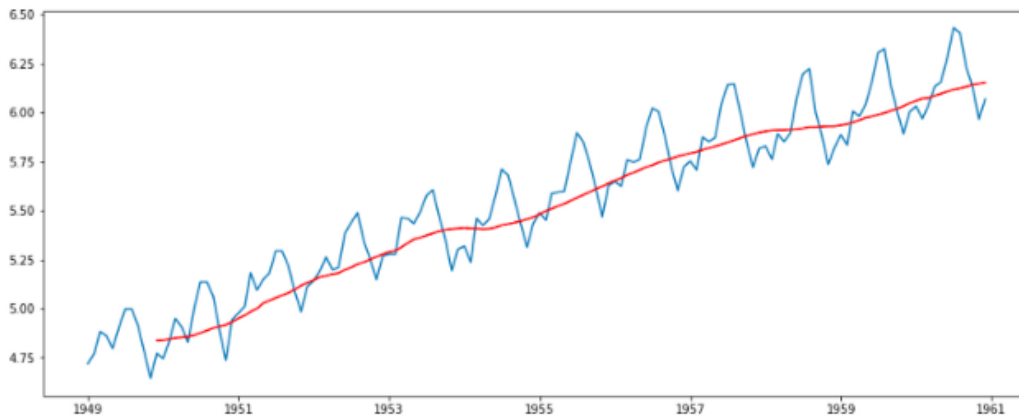
Moving average -

First take x consecutive values and this depends on the frequency if it is 1 year we take 12 values. Lucky for us that Pandas has a function for rolling estimate (“alright alright alright” -Matthew McConaughey!)

```
1 moving_avg = pd.rolling_mean(ts_log,12)
2 plt.plot(ts_log)
3 plt.plot(moving_avg, color='red')
```

```
/anaconda/lib/python2.7/site-packages/ipykernel/__main__.py:1: FutureWarning: pd.rolling_mean is deprecated for Series and will be removed in a future version, replace with
    Series.rolling(window=12,center=False).mean()
if __name__ == '__main__':
```

```
[<matplotlib.lines.Line2D at 0x11a52c5d0>]
```



Now subtract the rolling mean from the original series.

```
1 ts_log_moving_avg_diff = ts_log - moving_avg
2 ts_log_moving_avg_diff.head(12)
```

```
Month
1949-01-01      NaN
1949-02-01      NaN
1949-03-01      NaN
1949-04-01      NaN
1949-05-01      NaN
1949-06-01      NaN
1949-07-01      NaN
1949-08-01      NaN
1949-09-01      NaN
1949-10-01      NaN
1949-11-01      NaN
1949-12-01  -0.065494
Name: #Passengers, dtype: float64
```

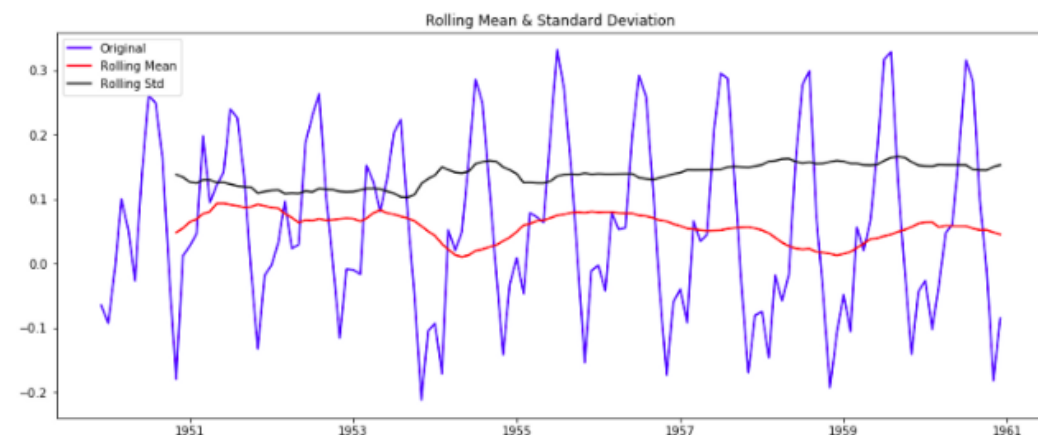
```
1 ts_log_moving_avg_diff.dropna(inplace=True)
2 ts_log_moving_avg_diff.head()
```

```
Month
1949-12-01  -0.065494
1950-01-01  -0.093449
1950-02-01  -0.007566
1950-03-01   0.099416
1950-04-01   0.052142
Name: #Passengers, dtype: float64
```

The reason there are null values is because we take the average of first 12 so 11 values are null. We can also see that in the visual representation. Thus it is dropped for further analysis. Now let's parse it to the function to check for stationarity.

```
1 test_stationarity(ts_log_moving_avg_diff)
```

```
/anaconda/lib/python2.7/site-packages/ipykernel/_main_.py:5: FutureWarning: pd.rolling_mean is deprecated for Series and will be removed in a future version, replace with
    Series.rolling(window=12,center=False).mean()
/anaconda/lib/python2.7/site-packages/ipykernel/_main_.py:6: FutureWarning: pd.rolling_std is deprecated for Series and will be removed in a future version, replace with
    Series.rolling(window=12,center=False).std()
```



```
Results of Dickey-Fuller Test:
Test Statistic      -3.162908
p-value             0.022235
#Lags Used          13.000000
Number of Observations Used  119.000000
Critical Value (5%)  -2.886151
Critical Value (1%)  -3.486535
Critical Value (10%) -2.579896
dtype: float64
```

We notice two things:

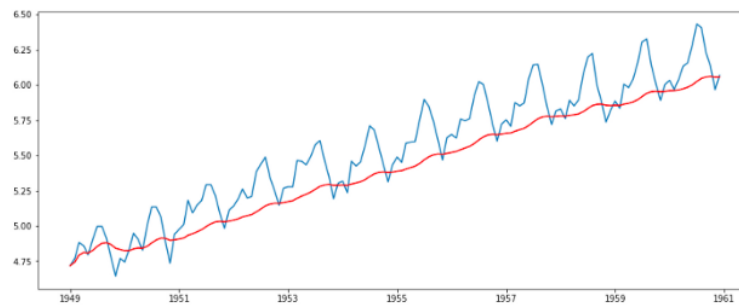
- The rolling values are varying slightly but there is no specific trend.
- The test statistics is smaller than the 5 % critical values. That tells us that we are 95% confident that this series is stationary.

In this example we can easily take a time period (12 months for a year), but there are situations where the time period range is more complex like stock price etc. So we use the exponentially weighted moving average (there are other weighted moving averages but for starters, let's use this). The previous values are assigned with a decay factor. Pandas again comes to the rescue with some awesome functions for it, like:

```

1 expwighted_avg = pd.ewm(ts_log, halflife=12)
2 plt.plot(ts_log)
3 plt.plot(expwighted_avg, color='red')
4
/anaconda/lib/python2.7/site-packages/ipykernel/_main_.py:1: FutureWarning: pd.ewm_mean is deprecated for Series and
d will be removed in a future version, replace with
Series.ewm(halflife=12, ignore_na=False, min_periods=0, adjust=True).mean()
if __name__ == '__main__':
[<matplotlib.lines.Line2D at 0x121d34790>]

```

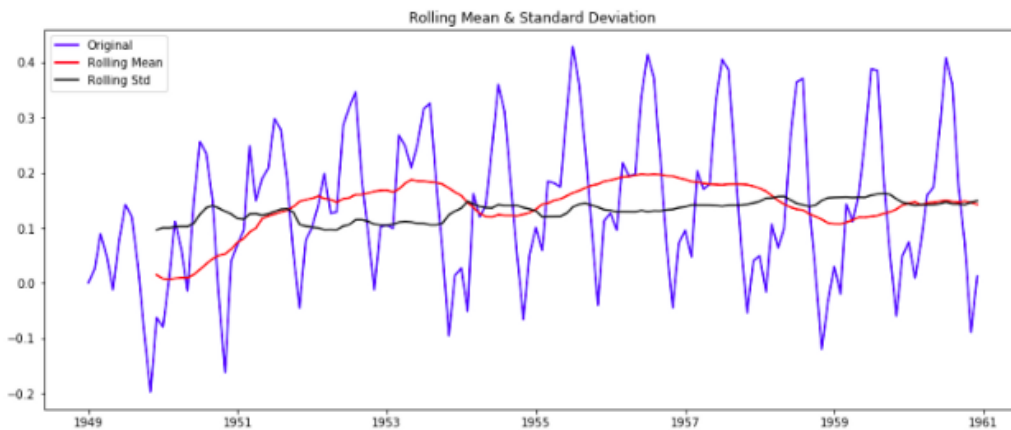


the parameter (halflife) is assumed 12, but it really depends on the domain. Let's check stationarity now:

```

1 ts_log_ewma_diff = ts_log - expwighted_avg
2 test_stationarity(ts_log_ewma_diff)
/anaconda/lib/python2.7/site-packages/ipykernel/_main_.py:5: FutureWarning: pd.rolling_mean is deprecated for Series and will be removed in a fu
ture version, replace with
Series.rolling(window=12, center=False).mean()
/anaconda/lib/python2.7/site-packages/ipykernel/_main_.py:6: FutureWarning: pd.rolling_std is deprecated for Series and will be removed in a fut
ure version, replace with
Series.rolling(window=12, center=False).std()

```



```

Results of Dickey-Fuller Test:
Test Statistic      -3.601262
p-value             0.005737
#Lags Used          13.000000
Number of Observations Used  130.000000
Critical Value (5%)  -2.884042
Critical Value (1%)  -3.481682
Critical Value (10%) -2.578770
dtype: float64

```

It is stationary because:

- Rolling values have less variations in mean and standard deviation in magnitude.
- the test statistic is smaller than 1% of the critical value. So we can say we are almost 99% confident that this is stationary.

Seasonality (along with Trend)

Previously we saw just trend part of the time series, now we will see both trend and seasonality. Most Time series have trends along with seasonality. There are two common methods to remove trend and seasonality, they are:

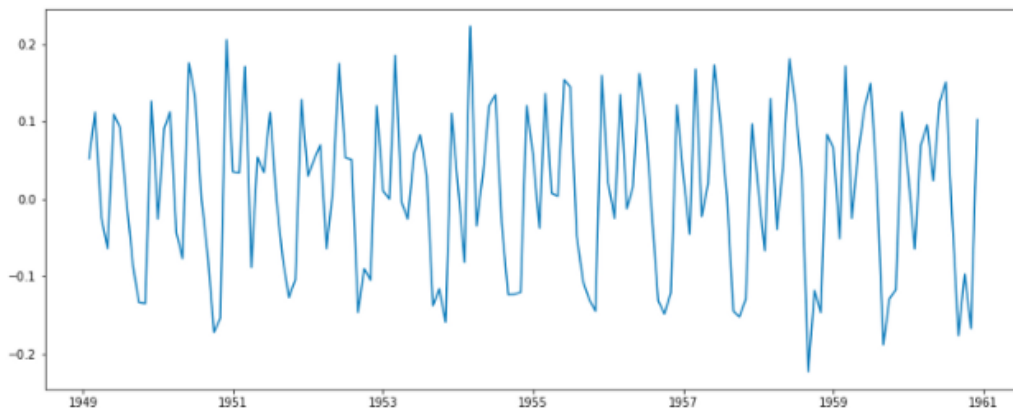
- Differencing: by taking difference using time lag
- Decomposition: model both trend and seasonality, then remove them

Differencing:

Here we first take the difference of the value at a particular time with that of the previous time. Now let's do it in Pandas.

```
1 #Take first difference:
2 ts_log_diff = ts_log - ts_log.shift()
3 plt.plot(ts_log_diff)
```

[<matplotlib.lines.Line2D at 0x11a74da50>]



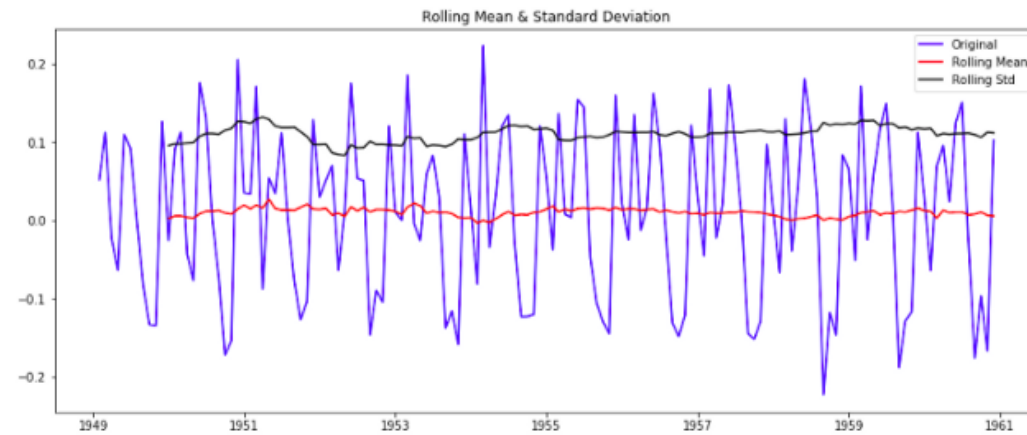
Looks ok to me but let's parse it using our stationary testing function

```

1 ts_log_diff.dropna(inplace=True)
2 test_stationarity(ts_log_diff)

/anaconda/lib/python2.7/site-packages/ipykernel/__main__.py:5: FutureWarning: pd.rolling_mean is deprecated for Series and will be removed in a future version, replace with
    Series.rolling(window=12,center=False).mean()
/anaconda/lib/python2.7/site-packages/ipykernel/__main__.py:6: FutureWarning: pd.rolling_std is deprecated for Series and will be removed in a future version, replace with
    Series.rolling(window=12,center=False).std()

```



```

Results of Dickey-Fuller Test:
Test Statistic      -2.717131
p-value              0.071121
#Lags Used           14.000000
Number of Observations Used  128.000000
Critical Value (5%)    -2.884398
Critical Value (1%)    -3.482501
Critical Value (10%)   -2.578960
dtype: float64

```

It is stationary because:

- the mean and std variations have small variations with time.
- test statistic is less than 10% of the critical values, so we can be 90 % confident that this is stationary.

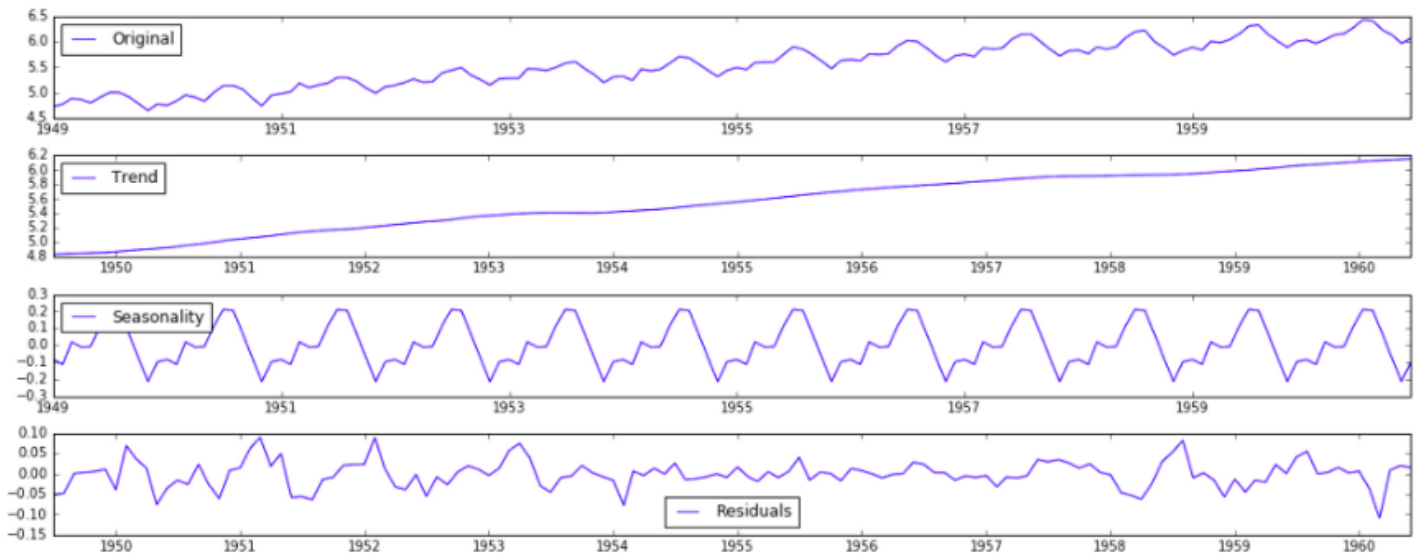
Decomposing:

Here we model both the trend and the seasonality, then the remaining part of the time series is returned. Guess what? Yup, we have some awesome function for it. Let's check it out:

```

1 from statsmodels.tsa.seasonal import seasonal_decompose
2 decomposition = seasonal_decompose(ts_log)
3
4 trend = decomposition.trend
5 seasonal = decomposition.seasonal
6 residual = decomposition.resid
7
8 plt.subplot(411)
9 plt.plot(ts_log, label='Original')
10 plt.legend(loc='best')
11 plt.subplot(412)
12 plt.plot(trend, label='Trend')
13 plt.legend(loc='best')
14 plt.subplot(413)
15 plt.plot(seasonal, label='Seasonality')
16 plt.legend(loc='best')
17 plt.subplot(414)
18 plt.plot(residual, label='Residuals')
19 plt.legend(loc='best')
20 plt.tight_layout()

```

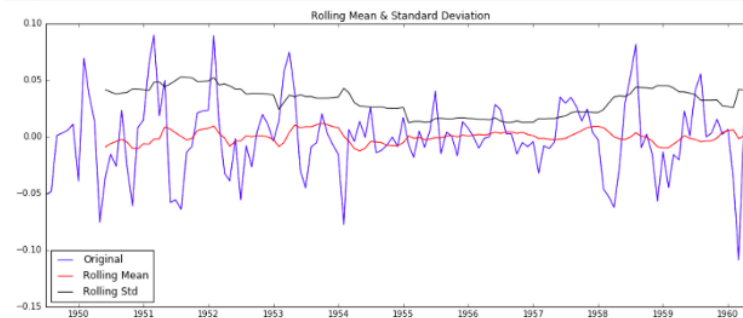


Remove the trend and seasonality from the Time series and now we can use the residual values. Let's check stationarity.

```

1 ts_log_decompose = residual
2 ts_log_decompose.dropna(inplace=True)
3 test_stationarity(ts_log_decompose)

```



```

Results of Dickey-Fuller Test:
Test Statistic      -6.332387e+00
p-value             2.885059e-08
#Lags Used          9.000000e+00
Number of Observations Used  1.220000e+02
Critical Value (5%)  -2.88538e+00
Critical Value (1%)  -3.485122e+00
Critical Value (10%) -2.579569e+00
dtype: float64

```

This is stationary because:

- test statistic is lower than 1% critical values.
- the mean and std variations have small variations with time.

Forecasting a Time Series

Now that we have made the Time series stationary, let's make models on the time series using differencing because it is easy to add the error, trend and seasonality back into predicted values.

We will use statistical modelling method called ARIMA to forecast the data where there are dependencies in the values.

Auto Regressive Integrated Moving Average (ARIMA)—It is like a linear regression equation where the predictors depend on parameters (p,d,q) of the ARIMA model.

Let me explain these dependent parameters:

- p : This is the number of AR (Auto-Regressive) terms. Example—if p is 3 the predictor for $y(t)$ will be $y(t-1), y(t-2), y(t-3)$.
- q : This is the number of MA (Moving-Average) terms. Example—if p is 3 the predictor for $y(t)$ will be $y(t-1), y(t-2), y(t-3)$.
- d : This is the number of differences or the number of non-seasonal differences.

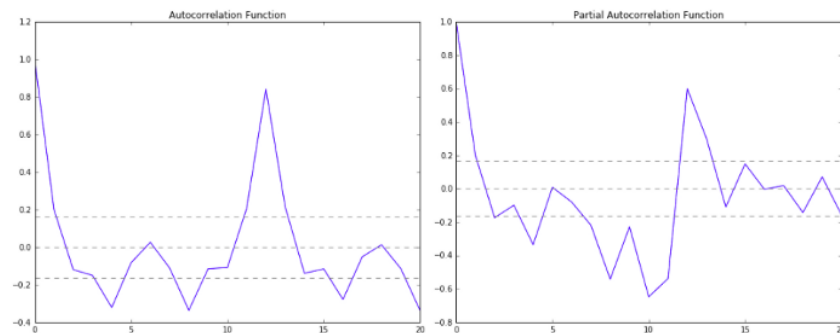
Now let's check out on how we can figure out what value of p and q to use. We use two popular plotting techniques; they are:

- Autocorrelation Function (ACF): It just measures the correlation between two consecutive (lagged version). example at lag 4, ACF will compare series at time instance $t_1 \dots t_2$ with series at instance $t_1-4 \dots t_2-4$
- Partial Autocorrelation Function (PACF): is used to measure the degree of association between $y(t)$ and $y(t-p)$.

```

1 from statsmodels.tsa.arima_model import ARIMA
2 #ACF and PACF plots:
3 from statsmodels.tsa.stattools import acf, pacf
4
5 lag_acf = acf(ts_log_diff, nlags=20)
6 lag_pacf = pacf(ts_log_diff, nlags=20, method='ols')
7
8 #Plot ACF:
9 plt.subplot(121)
10 plt.plot(lag_acf)
11 plt.axhline(y=0, linestyle='--', color='gray')
12 plt.axhline(y=-1.96/np.sqrt(len(ts_log_diff)), linestyle='--', color='gray')
13 plt.axhline(y=1.96/np.sqrt(len(ts_log_diff)), linestyle='--', color='gray')
14 plt.title('Autocorrelation Function')
15
16 #Plot PACF:
17 plt.subplot(122)
18 plt.plot(lag_pacf)
19 plt.axhline(y=0, linestyle='--', color='gray')
20 plt.axhline(y=-1.96/np.sqrt(len(ts_log_diff)), linestyle='--', color='gray')
21 plt.axhline(y=1.96/np.sqrt(len(ts_log_diff)), linestyle='--', color='gray')
22 plt.title('Partial Autocorrelation Function')
23 plt.tight_layout()

```



The two dotted lines on either sides of 0 are the confidence intervals. These can be used to determine the 'p' and 'q' values as:

- p: The first time where the PACF crosses the upper confidence interval, here its close to 2. hence $p = 2$.
- q: The first time where the ACF crosses the upper confidence interval, here its close to 2. hence $p = 2$.

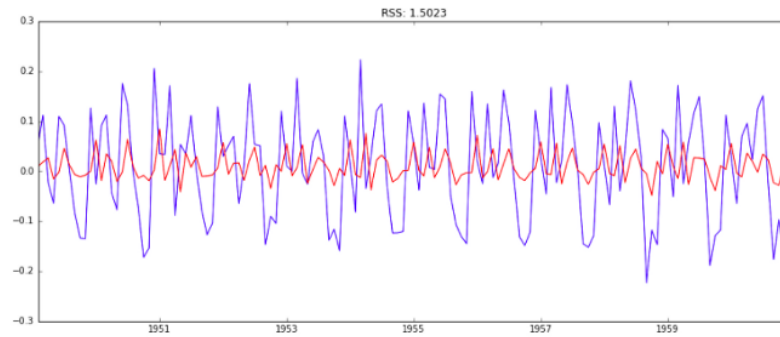
Now, using this make 3 different ARIMA models considering individual as well as combined effects. I will also print the RSS for each. Please note that here RSS is for the values of residuals and not actual series.

AR Model


```

1 #MA model:
2 model = ARIMA(ts_log, order=(2, 1, 0))
3 results_AR = model.fit(disp=-1)
4 plt.plot(ts_log_diff)
5 plt.plot(results_AR.fittedvalues, color='red')
6 plt.title('RSS: %.4f' % sum((results_AR.fittedvalues-ts_log_diff)**2))
<matplotlib.text.Text at 0x1103cae50>

```

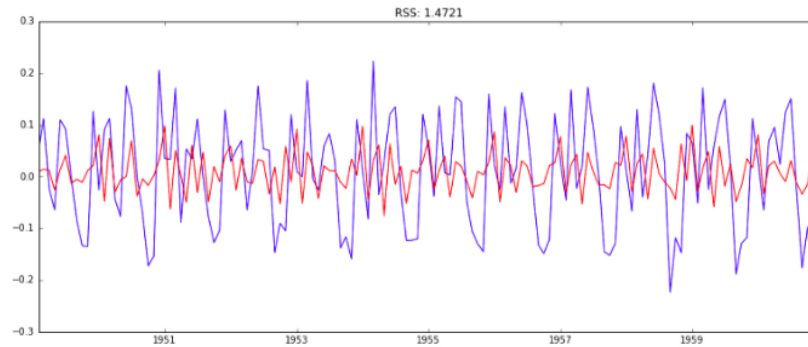


MA Model

```

1 model = ARIMA(ts_log, order=(0, 1, 2))
2 results_MA = model.fit(disp=-1)
3 plt.plot(ts_log_diff)
4 plt.plot(results_MA.fittedvalues, color='red')
5 plt.title('RSS: %.4f' % sum((results_MA.fittedvalues-ts_log_diff)**2))
<matplotlib.text.Text at 0x1106b3e50>

```

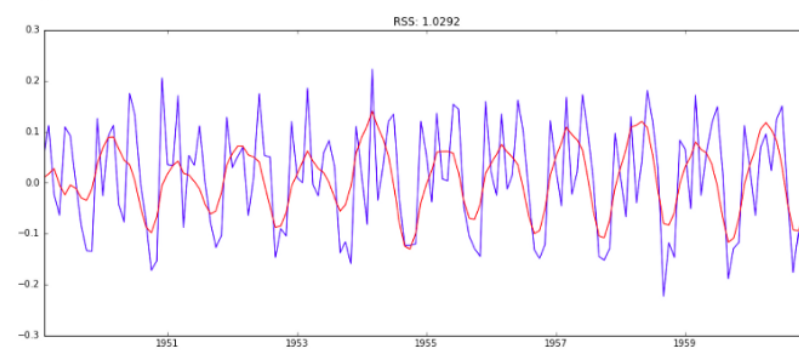


ARIMA MODEL

```

1 model = ARIMA(ts_log, order=(2, 1, 2))
2 results_ARIMA = model.fit(disp=-1)
3 plt.plot(ts_log_diff)
4 plt.plot(results_ARIMA.fittedvalues, color='red')
5 plt.title('RSS: %.4f'% sum((results_ARIMA.fittedvalues-ts_log_diff)**2))

```



RSS values:

- AR=1.5023
- MA=1.472
- ARIMA = 1.0292

ARIMA has the best RSS values.

FINAL STEP: BRINGING THIS BACK TO THE ORIGINAL SCALE

Steps involved:

- First get the predicted values and store it as series. You will notice the first month is missing because we took a lag of 1 (shift).
- Now convert differencing to log scale: find the cumulative sum and add it to a new series with a base value (here the first-month value of the log series).

```

1 predictions_ARIMA_diff = pd.Series(results_ARIMA.fittedvalues, copy=True)
2 print predictions_ARIMA_diff.head()

```

Month
1949-02-01 0.009580
1949-03-01 0.017491
1949-04-01 0.027670
1949-05-01 -0.004521
1949-06-01 -0.023889
dtype: float64

```

1 predictions_ARIMA_diff_cumsum = predictions_ARIMA_diff.cumsum()
2 print predictions_ARIMA_diff_cumsum.head()

```

Month
1949-02-01 0.009580
1949-03-01 0.027071
1949-04-01 0.054742
1949-05-01 0.050221
1949-06-01 0.026331
dtype: float64

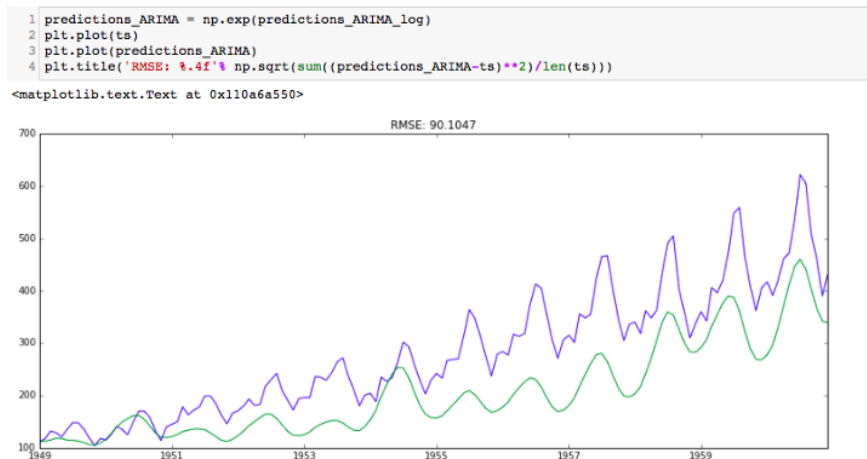
```

1 predictions_ARIMA_log = pd.Series(ts_log.ix[0], index=ts_log.index)
2 predictions_ARIMA_log = predictions_ARIMA_log.add(predictions_ARIMA_diff_cumsum, fill_value=0)
3 predictions_ARIMA_log.head()

```

Month
1949-01-01 4.718499
1949-02-01 4.728079
1949-03-01 4.745570
1949-04-01 4.773241
1949-05-01 4.768720
dtype: float64

- Next -take the exponent of the series from above (anti-log) which will be the predicted value—the time series forecast model.
- Now plot the predicted values with the original.
- Find the RMSE



The result can be further refined to get a better model. The scope of the blog was to quickly introduce Time Series Forecasting. Hope you guys enjoyed the blog, there a lot more details with respect Time series analysis and forecasting. This is a good place to understanding the theory behind the practical techniques .

