

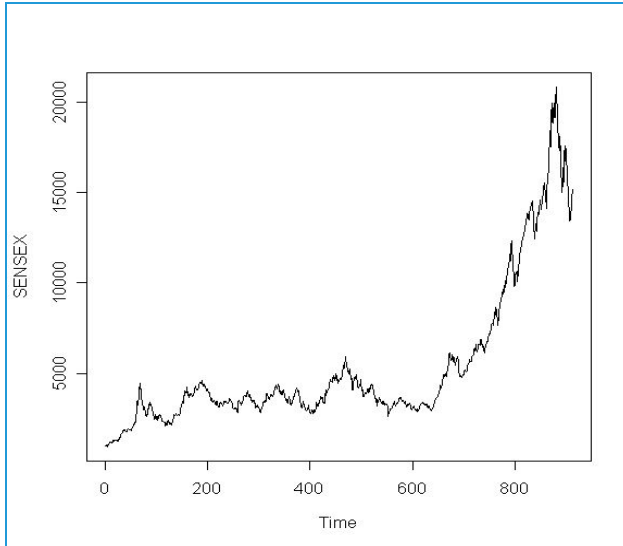
# Introduction to Time Series Analysis

# Contents

1. What is Time Series
2. Components of Time Series
3. Application Areas
4. Time Series Analysis in Python
  - i. Plotting a Time Series
  - ii. Subsetting a Time Series

# What is Time Series ?

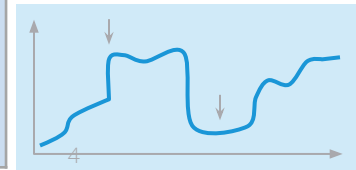
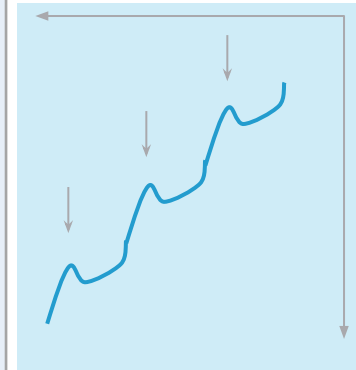
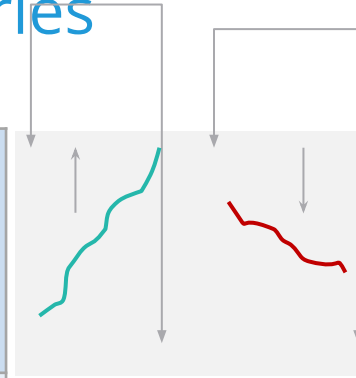
Time Series is a sequence of values observed over time



Types of Time Scale	
<b>Discrete</b>	Value changes after jumping from one time period to other. Example: Dow Jones Index -End of Day Values
<b>Continuous</b>	Value changes within an infinitely short amount of time. Example. Temperature, Dow Jones Index Tracked Real Time

# Components of Time Series

<b>Trend</b>	<p>Long-term increase or decrease in the time series.</p> <p>There may be increase/decrease in short term but overall trend in the long term can be increasing or decreasing.</p>
<b>Seasonality</b>	<p><b>Predictable and recurring trends and patterns</b> over a period of time, normally a year. An example of a seasonal time series is retail data, which sees spikes in sales during holiday seasons like Christmas.</p> <p>Seasonality is reflected only when data is available for more than one year</p>
<b>Cyclic Pattern</b>	<p>Exists when data exhibit rises and falls that are not of fixed period. The duration of these fluctuations is usually of at least 2 years</p>



# Application Areas

<u>Industry</u>	<u>Model/Predict</u>	<u>Based on Information such as:</u>	<u>Purpose</u>
Finance	Price of a Stock	<ul style="list-style-type: none"><li>• Recent price movement of the stock</li></ul>	Forecasting
Economics	Inflation Rates	<ul style="list-style-type: none"><li>• Trend and seasonality in inflation rates</li></ul>	Forecasting
Retail /FMCG	Monthly Sales	<ul style="list-style-type: none"><li>• Location, marketing expenses on TV, print and online media</li></ul>	Predictive and Optimization

# Case Study

## Background

- Annual Sales for a specific company from year 1961 to 2017

## Objective

- To plot a time series object

## Available Information

- Number of cases: 57
- Variables: Year, sales(in 10's GBP)

# Data Snapshot

turnover\_annual data

Variables

Observations on Discrete Time Scale

Year	sales
1961	224786
1962	230034
1963	236562
1964	250960
1965	261615
1966	268316
1967	283589
1968	280160
1969	301422
1970	308018
1971	322025

Columns	Description	Type	Measurement	Possible values
Year	Financial Year	Numeric	-	-
sales	sales(in 10's GBP)	Numeric	In British Pound	Positive values
		1974	364834	
		1975	392503	

# Time Series in Python

```
# Import turnover_annual Data
```

```
import pandas as pd  
salesdata = pd.read_csv('turnover_annual.csv')
```

```
# Creating a Time Series Object
```

```
rng = pd.date_range('01-01-1961', '31-12-2017', freq='Y')  
s = salesdata.sales.values  
salesseries = pd.Series(s, rng)
```

- ❑ **date\_range()** creates pandas date object.
- ❑ When the time series has seasonal components, argument **freq** = can be included. It denotes number of observations per unit of time. Eg. If data is quarterly: **freq = 'Q'**, if data is monthly: **freq = 'M'**.
- ❑ **pd.Series()** combines time series variable object “s” and date object “rng”.  
The new object salesseries will be used for further analysis.



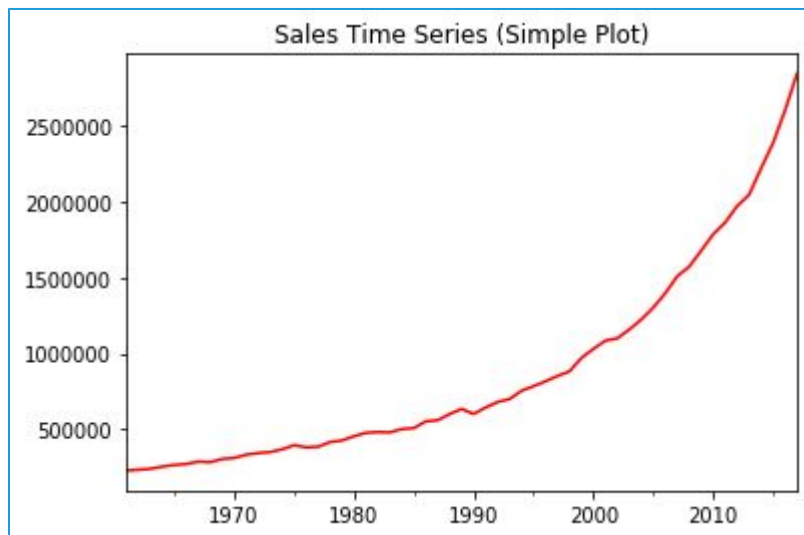
# Plotting Time Series in Python

# Plotting a Time Series Object

```
salesseries.plot(color='red', title ="Sales Time Series (Simple Plot)")
```

# Output

**plot()** generates a simple line chart.



## Interpretation :

- The time-series clearly shows upward trend.

# Subsetting Time Series in Python

- Large volumes of data are required for most real world analytics, time series is no exception.
- Subsetting is an important tool as it facilitates partitioning the data within Python for micro-level specific analysis.

## # Subsetting a Time Series Object

```
salesseries2 = salesseries.loc['1990-12-31':'2016-12-31']
```

**loc[ ]** is a generic function which extracts the subset of the object x observed between the times **specified within the range**.

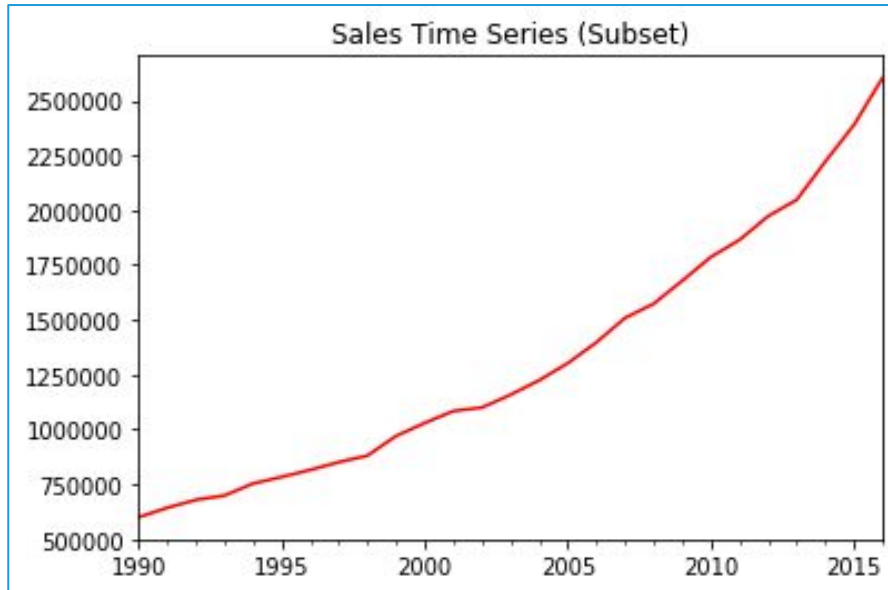


In Pandas we can directly subset the time series object using **loc[ ]** function

# Subsetting Time Series in Python

```
salesseries2.plot(color='red', title = "Sales Time Series (Subset)")
```

# Output



Plot from 1990 to 2016 shows increasing trend

# Quick Recap

## Time Series

- Sequence of values measured over time
- Time scale can be discrete or continuous

## Time Series in Analysis

- Analyze and forecast time series values

## Time Series in Python

- `pd.date_range` created date object.
- `pd.Series` creates time series object by combining date object and time series variable

# Time Series Decomposition

# Contents

1. Components of Time Series
2. Understanding Moving Averages
3. Time Series Decomposition

# Components of Time Series

- As we know, Trend and Seasonality are main components of Time Series.
- If we assume an additive model, we can write

where,

$$Y_t = S_t + T_t + R_t$$

$Y_t$  : Time series value at period t

$S_t$  : Seasonal component at period t

$T_t$  : Trend-cycle component at period t

$R_t$  : Remainder (or irregular or error)  
component at period t

- Alternatively, a multiplicative model would be written as

$$Y_t = S_t * T_t * R_t$$



Additive model is most appropriate if the magnitude of the seasonal fluctuations or the variation around the trend-cycle does not vary with the level of the time series.

# Understanding Moving Averages

- Moving Averages are averages calculated for consecutive data from overlapping subgroups of fixed length
- Moving averages smoothen a time series by filtering out random fluctuations

Day	End of Day Sales		Moving Average of Period 3
1	1500	}	NA
2	2100		NA
3	1750	→	$(1500+2100+1750)/3=1783.33$
4	1900	}	$(2100+1750+1900)/3=1916.67$
5	1650	→	$(1750+1900+1650)/3=1766.67$

The first 2 MA values for length 3 are not calculated

- Period of the moving average depends on type of data
- **Non-seasonal data:** Shorter length (Typically 3 period or 5 period MA is considered)
- **Seasonal data:** Typical period is 12 for monthly data and 4 for quarterly data



# Time Series Decomposition – Simple Method

Decomposition is a statistical method that deconstructs a time series.

Steps to follow :

## Find Trend

Obtain moving averages covering one season – This provides trend component of the time series

## Eliminate Trend

Eliminate trend component from original time series. Calculate  $Y_t - T_t$

## Estimate Seasonality

To estimate the seasonal component for a given time period, simply average the de-trended values for that time period. These seasonal indexes are then adjusted to ensure that they add to zero

The remainder component is calculated by subtracting the estimated seasonal and trend-cycle components

# Time Series Decomposition – Simple Method

Suppose we have monthly time series data, for three years 2014, 2015 and 2016:

## Step 1 □ Calculate Moving Averages



(Consider moving average period of 13 - previous 6 months, next 6 months and current month to calculate moving average of current month)

This gives the trend component  $T_t$

## Step 2 □ Eliminating Trend

Remove  $T_t$  from the original time series  $Y_t$

## Step 3 □ Estimate Seasonal Component

The seasonal index for July is the average of all the de-trended July values in the data i.e. Average of De-trended July 2014, July 2015 and July 2016

# Case Study

## Background

- Monthly Sales Data for 3 Years (2013, 2014, 2015)

## Objective

- To decompose time series into its components and study each component separately.

## Available Information

- Sample size is 36
- Variables: Year, Month, Sales

# Data Snapshot

Sales Data for 3  
Years

Variables

Year	Month	Sales
2013	Jan	123
2013	Feb	142
2013	Mar	164
2013	Apr	173
2013	May	183
2013	Jun	192
2013	Jul	199
2013	Aug	203
2013	Sep	207
2013	Oct	209
2013	Nov	214
2013	Dec	255

Columns	Description	Type	Measurement	Possible values
Year	Year	Numeric	2013, 2014, 2015	3
Month	Month	Character	Jan - Dec	12
Sales	Sales in USD Million	numeric	USD Million	Positive values
	2014	Jul	245	

# Time Series Decomposition in Python

# Simple Decomposition

```
import pandas as pd
salesdata = pd.read_csv("Sales Data for 3 Years.csv")
rng = pd.date_range('01-01-2013', '31-12-2015', freq='M')
s = salesdata.Sales.values
salesseries = pd.Series(s, rng)
```

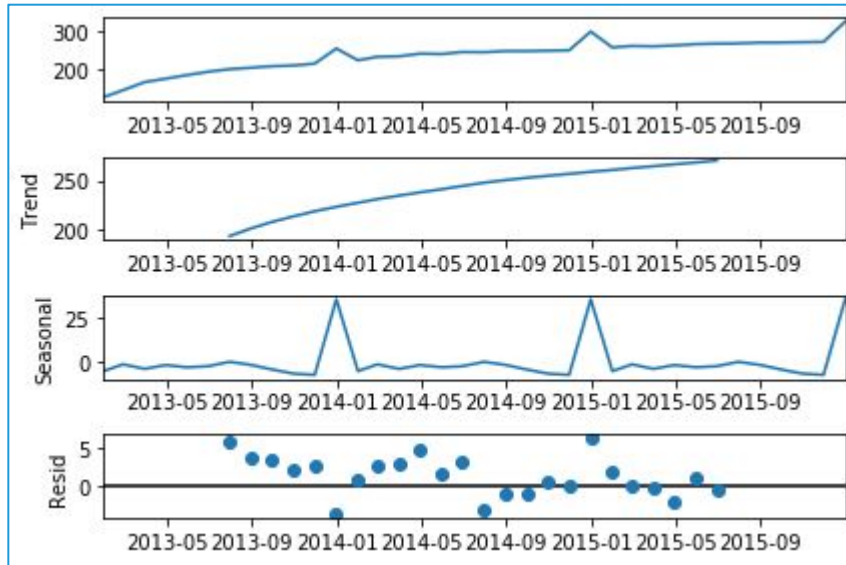
- ❑ **freq = "M"** indicates monthly date data.
- ❑ **pd.Series()** creates time series object using date object as index

```
import statsmodels.api as sm
decomp = sm.tsa.seasonal_decompose(salesseries.interpolate())
decomp.plot()
```

- ❑ **tsa.seasonal\_decompose()** performs a classical seasonal decomposition through moving averages.
- ❑ **plot()** of decompose object gives a 4-level visual representation.

# Time Series Decomposition in Python

# Output



Trend is not estimated for first/last few values  
Seasonal Component repeats from year to year



The `tsa.seasonal_decompose()` function takes time series as additive by default. If you wish to create a multiplicative time series, simply add `model="multiplicative"` in the function.

# Decomposition in Python – Seasonal Component

#Analyzing the decomp object. Each component can be separately viewed

decomp.seasonal

2013-01-31	-5.006944
2013-02-28	-1.181944
2013-03-31	-3.736111
2013-04-30	-1.588194
2013-05-31	-2.831944
2013-06-30	-2.231944
2013-07-31	0.234722
2013-08-31	-1.525694
2013-09-30	-4.094444
2013-10-31	-6.446528
2013-11-30	-7.127778
2013-12-31	35.536806
2014-01-31	-5.006944
2014-02-28	-1.181944
2014-03-31	-3.736111
2014-04-30	-1.588194
2014-05-31	-2.831944
2014-06-30	-2.231944
2014-07-31	0.234722
2014-08-31	-1.525694
2014-09-30	-4.094444
2014-10-31	-6.446528
2014-11-30	-7.127778
2014-12-31	35.536806
2015-01-31	-5.006944
2015-02-28	-1.181944
2015-03-31	-3.736111
2015-04-30	-1.588194
2015-05-31	-2.831944
2015-06-30	-2.231944
2015-07-31	0.234722
2015-08-31	-1.525694
2015-09-30	-4.094444
2015-10-31	-6.446528
2015-11-30	-7.127778

## Interpretation :

- This table shows seasonal component of time series

# Decomposition in Python – Trend Component

decomp.trend

2013-01-31	NaN
2013-02-28	NaN
2013-03-31	NaN
2013-04-30	NaN
2013-05-31	NaN
2013-06-30	NaN
2013-07-31	192.833333
2013-08-31	200.770833
2013-09-30	207.450000
2013-10-31	213.195833
2013-11-30	218.408333
2013-12-31	223.016667
2014-01-31	227.166667
2014-02-28	230.962500
2014-03-31	234.550000
2014-04-30	237.929167
2014-05-31	241.100000
2014-06-30	244.516667
2014-07-31	247.891667
2014-08-31	250.575000
2014-09-30	252.933333
2014-10-31	254.991667
2014-11-30	257.041667
2014-12-31	259.104167
2015-01-31	261.041667
2015-02-28	262.995833
2015-03-31	264.916667
2015-04-30	266.841667
2015-05-31	268.758333
2015-06-30	270.841667
2015-07-31	NaN
2015-08-31	NaN
2015-09-30	NaN

## Interpretation :

- This table shows trend component of time series

?

Why are NAs getting generated?

This is because trend not estimated for first/last few values. Consequently, the same will be reflected in random component as well.



# Decomposition in Python – Random Component

decomp.resid

2013-01-31	NaN
2013-02-28	NaN
2013-03-31	NaN
2013-04-30	NaN
2013-05-31	NaN
2013-06-30	NaN
2013-07-31	5.931944
2013-08-31	3.754861
2013-09-30	3.644444
2013-10-31	2.250694
2013-11-30	2.719444
2013-12-31	-3.553472
2014-01-31	0.840278
2014-02-28	2.719444
2014-03-31	2.986111
2014-04-30	4.759028
2014-05-31	1.731944
2014-06-30	3.315278
2014-07-31	-3.126389
2014-08-31	-0.949306
2014-09-30	-0.838889
2014-10-31	0.554861
2014-11-30	0.086111
2014-12-31	6.359028
2015-01-31	1.965278
2015-02-28	0.086111
2015-03-31	-0.180556
2015-04-30	-1.953472
2015-05-31	1.073611
2015-06-30	-0.509722
2015-07-31	NaN
2015-08-31	NaN
2015-09-30	NaN

## Interpretation :

- This table shows random component of time series

?

Why are NAs getting generated?

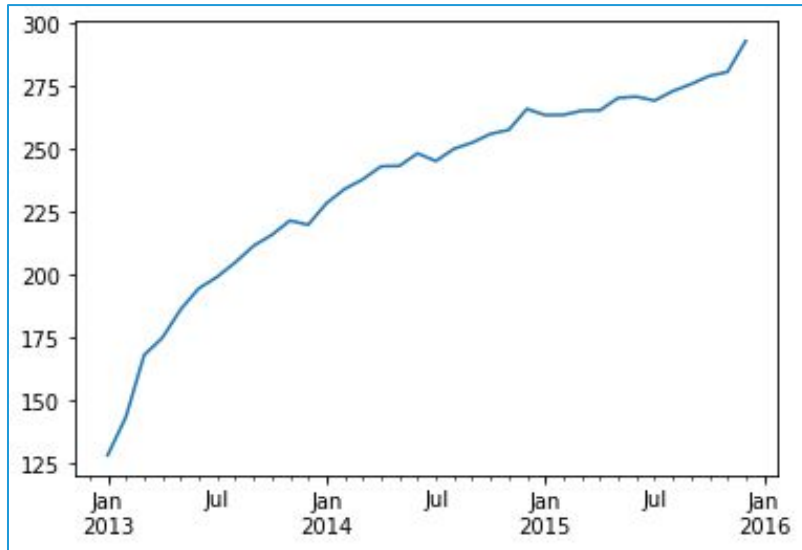
This is because trend not estimated for first/last few values. Consequently, the same will be reflected in random component as well.

# Seasonally Adjusted Time Series

# Doing Seasonal Adjustment

```
salesadj = salesseries - decomp.seasonal  
salesadj.plot()
```

# Output



## Interpretation :

- This plot shows seasonally adjusted time series

# Quick Recap

In this session, we learnt about **time series decomposition** and **exponential smoothing**:

## What is Decomposition

- A time series is made up of multiple components such as seasonality, trend, randomness
- Sometimes, studying these components separately provides a more comprehensive insight about the series

## Decomposition in Python

- **`sm.tsa.seasonal_decompose()`** carries out simple seasonal decomposition

# Time Series Analysis

## Stationarity of Time Series - I

# Contents

1. What is Stationarity ?
2. **Stationary Time Series – White Noise Process**
3. **Non-Stationary Time Series – Random Walk**
4. Importance of Stationary Time Series
5. Identifying Stationary Time Series
6. Concept of Autocorrelation & Correlograms

# What is Stationarity of Time Series ?

Time series process is called **Stationary** if statistical properties of the process remain unchanged over

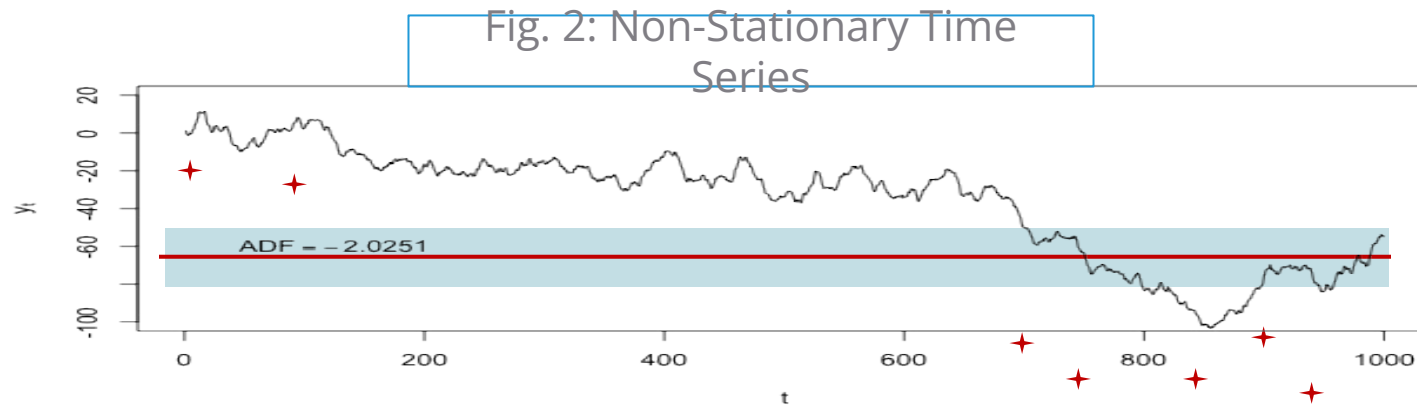
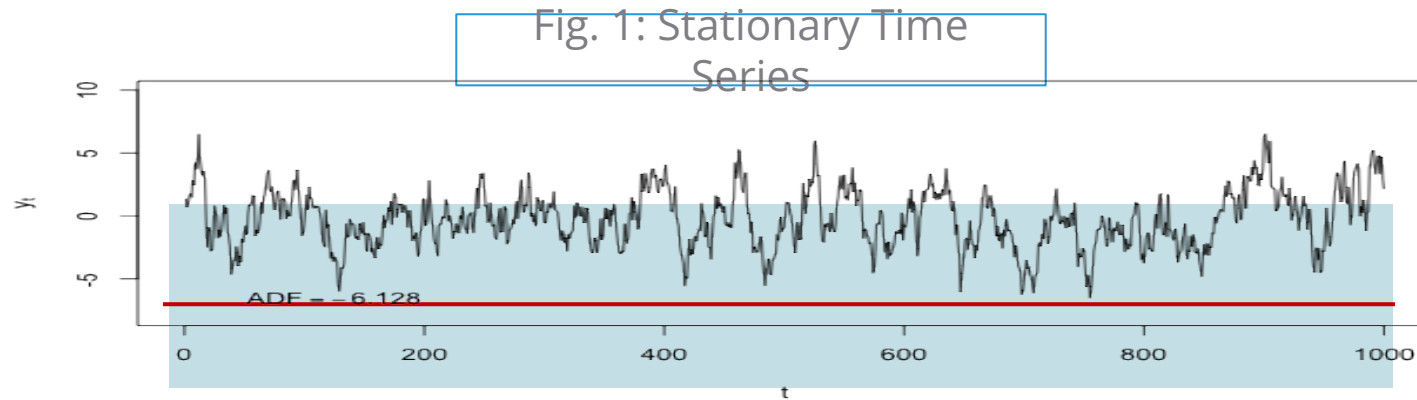
✓  
If  $Y_t$  is a **stationary** time series where  $t=1,2,3,\dots$   
then,

$$E(Y_t) = \mu_t = \mu \text{ (constant)}$$

$$\text{Var}(Y_t) = \sigma_t^2 = \sigma^2 \text{ (constant)}$$

$\text{cov}(Y_t, Y_{t+s})$  depends only on **s** (lag), and is independent of **t** (time)

# Stationary vs. Non-Stationary Time Series

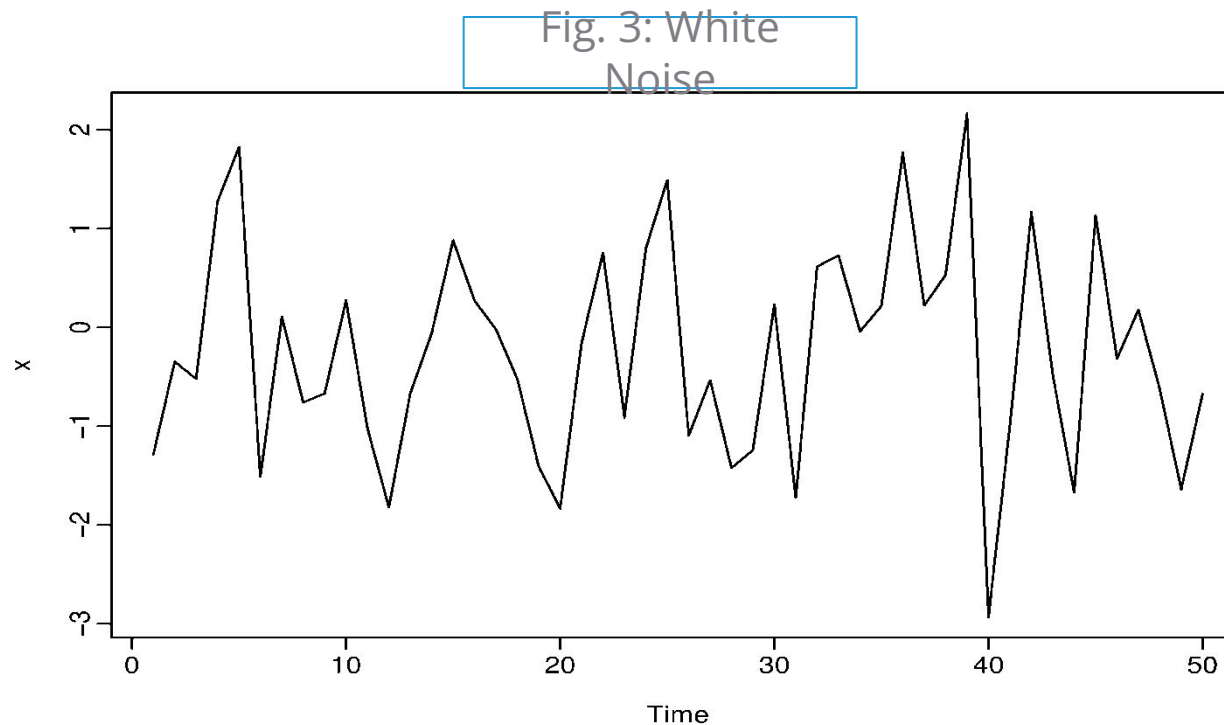


## Interpretation :

- A stationary time series has a constant long term mean and variance.
- The first diagram shows a stationary time series whereas the second shows a non-stationary series.

# Stationary Time Series - White Noise Process

- **White noise** is the simplest example of stationary time series.
- White Noise time series has **zero mean**, **constant variance** and **zero covariance** with lagged time series.





# Non Stationary Time Series - Random Walk

- Random Walk is the simplest case of Non-stationary time series. It is of the form

$$Y_t = Y_{t-1} + U_t$$
$$t=1,2,3,\dots$$

- We assume that  $U_t$  is a random series with,
  - Constant mean  $\mu$
  - Constant variance  $\sigma^2$
  - Serially uncorrelated

Value of Y at time t is equal to its value at time (t-1) plus a random shock

# Why Random Walk is Non Stationary

Let  $Y_t = 0$  at time  $t = 0$

$$Y_1 = U_1$$

$$Y_2 = Y_1 + U_2 = U_1 + U_2$$

$$Y_3 = Y_2 + U_3 = U_1 + U_2 + U_3$$

$$\text{So, } Y_t = \sum U_t \quad t = 1, 2, 3 \dots$$

$Y_t = \sum U_t$	
$E(Y_t) = E(\sum U_t) = t \cdot \mu$ i.e. Mean is not constant	$\text{Var}(Y_t) = \text{Var}(\sum U_t) = t \cdot \sigma^2$ i.e. Variance is not constant

Therefore, Random Walk is a non-stationary time series

# Importance of Stationary Time Series

- **Calibration** (estimation of model parameters using historical data) is an important concept in the **forecasting** of time series values.
- In the calibration of time series models we need a stationary time series.
- With a non stationary time series we get into **spurious** regression which badly affects forecasting.

# How to Make a Non Stationary Time Series Stationary?

## Two Methods for Making Time Series Stationary

### Differencing

$$Y_t = Y_{t-1} + U_t ; t=1,2,3, \dots$$

$U_t$  is a random series with **Constant mean**  $\mu$ , **Constant variance**  $\sigma^2$ , and is **serially uncorrelated** i.e ( $U_t$  is stationary).

Hence,  $Y_t$  is differenced:

$$Y_t - Y_{t-1} = \Delta Y = U_t$$

Differencing can be well applied in case of stochastic time series

### De-trending

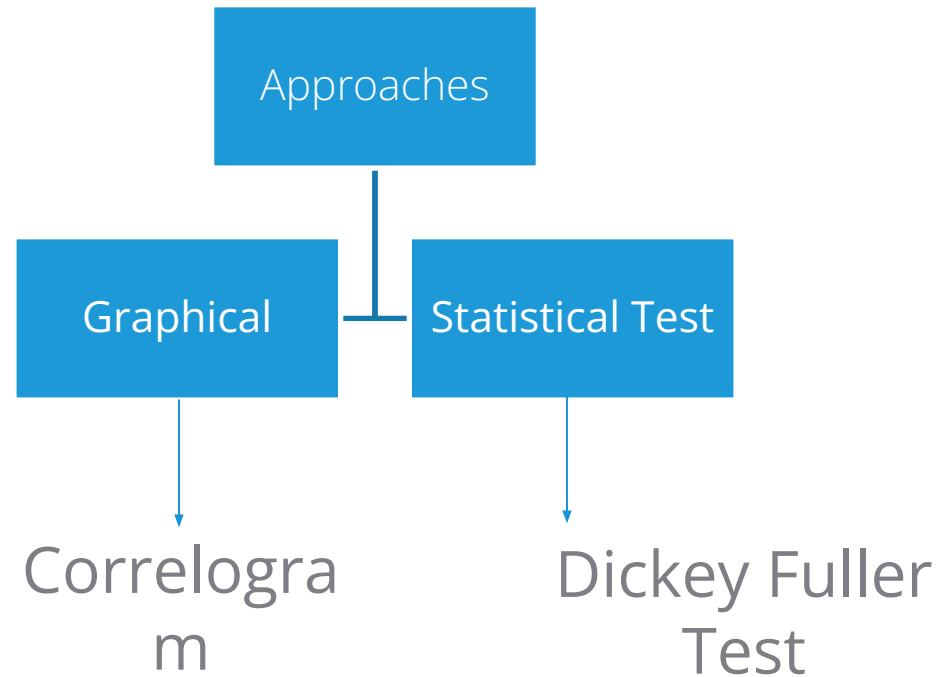
$$Y_t = \beta_1 + \beta_2 t + U_t ; t=1,2,3, \dots$$

$U_t$  is a stationary with **zero mean** and **constant variance**  $\sigma^2$ . When Trend element  $(\beta_1 + \beta_2 t)$  is subtracted, the result is a stationary process :

$$Y_t - (\beta_1 + \beta_2 t) = U_t$$

De-trending is useful when trend is deterministic

# Identifying Stationary Time Series & Concept of Autocorrelation

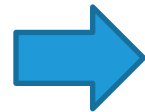


# Autocorrelation of Lag 1

Data

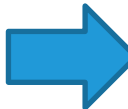
Year	Sales(m)
Jan	110
Feb	113
Mar	121
Apr	123
May	126
Jun	120
Jul	119
Aug	127
Sep	129
Oct	131
Nov	130
Dec	132

First Order  
Lag



Year	Y <sub>t</sub>	Y <sub>t-1</sub>
Jan	110	
Feb	113	110
Mar	121	113
Apr	123	121
May	126	123
Jun	120	126
Jul	119	120
Aug	127	119
Sep	129	127
Oct	131	129
Nov	130	131
Dec	132	130

Autocorrelation of Lag 1,  
where  $\bar{Y}$  = Average of sales


$$\sum_{t=2}^n \frac{(Y_t - \bar{Y})(Y_{t-1} - \bar{Y})}{\sum_{t=1}^n (Y_t - \bar{Y})^2}$$

- Autocorrelation is a correlation between a time series ( $Y_t$ ) and another time series representing lagged values of the same time series ( $Y_{t-k}$ )

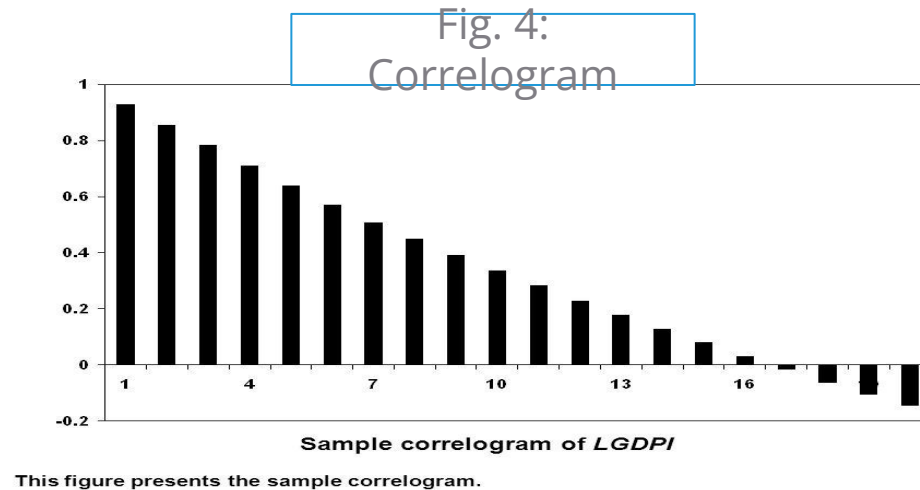
# Autocorrelation Function & Correlogram

- Plot of the sample autocorrelation function against lag is called Correlogram.
- ACF (autocorrelation function) is a general expression for lag k autocorrelation.
- Correlogram is mainly used in checking stationarity of a series.

$$r_k \text{ (Sample autocorrelation of lag } k) = \sum_{t=k+1}^n \frac{(Y_t - \bar{Y})(Y_{t-k} - \bar{Y})}{\sum_{t=1}^n (Y_t - \bar{Y})^2}$$

# Correlogram for Checking Stationarity

- For a non-stationary time series showing trend, slow decay pattern will be observed.
- If a time series is characterized by seasonal fluctuations, then the correlogram would also exhibit oscillations at the same frequency.





# Partial Autocorrelation Function

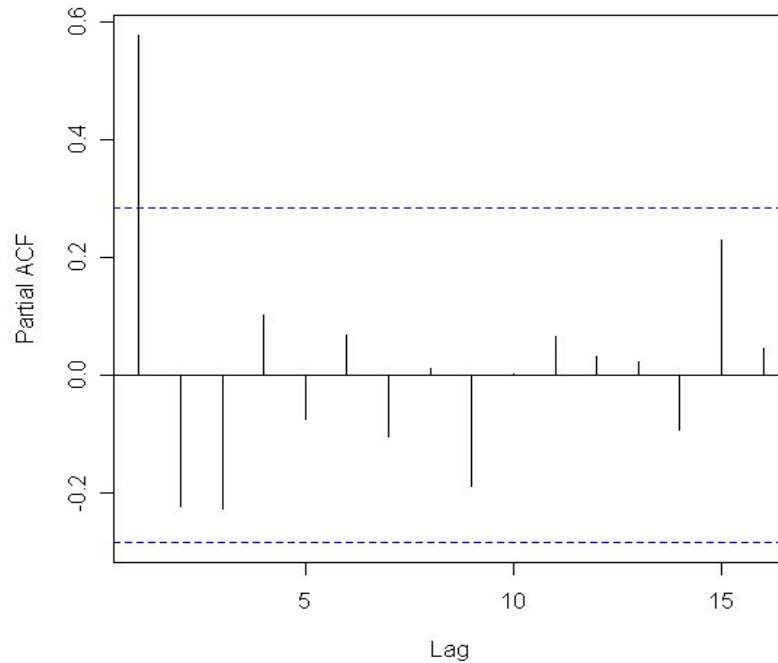
- Partial autocorrelation (PACF) is the

Autocorrelation between  $y_t$  and  $y_{t-h}$  after removing any linear dependence on  $y_1, y_2, \dots, y_{t-h+1}$

- PACF gives partial autocorrelations at various lags.
- PACF is mainly used to identify order of moving average present in the process.

# Partial Correlogram

Fig. 5: Partial ACF Plot



- Y-Axis = Partial Autocorrelation Function
- X-Axis = Lag
- Here we see that, after lag 1, the PACF drops dramatically and all PACFs after lag 1 are statistically insignificant.

# Quick Recap

## What is Stationarity

- Time series is stationary if the statistical properties of the process remain unchanged over time
- Time series which has zero mean, constant variance and zero covariance with lagged time series is known as **“White Noise”**

## Random Walk

- Is the simplest form of non-stationary time series where  $Y_t - Y_{t-1} = U_t$  ( $U_t$  is a random error term)

## Making a Non-Stationary Time Series Stationary

- Two approaches: **Differencing** (when time series is stochastic) and **De-trending** (when trend is deterministic)

## Identifying Stationarity of Time Series

- **Graphical method:** ACF and PACF Correlograms
- **Analytical method:** Dickey Fuller Test

# Time Series Modeling

## Seasonal ARIMA Model

# Contents

1. Case Study
2. Creating and Plotting Time Series in Python
3. Checking Stationarity in Python
4. Dickey Fuller (DF) Test

# Case Study

## Background

- Annual Sales for a specific company from year 1961 to 2017

## Objective

- To assess stationarity of time series

## Available Information

- Number of cases: 57
- Variables: Year, sales(in 10's GBP)

# Data Snapshot

turnover\_annual data

Variables

Observations on Discrete Time Scale

Year	sales
1961	224786
1962	230034
1963	236562
1964	250960
1965	261615
1966	268316
1967	283589
1968	280160
1969	301422
1970	308018
1971	320035

Columns	Description	Type	Measurement	Possible values
Year	Financial Year	Numeric	-	-
sales	sales(in 10's GBP)	Numeric	In British Pound	Positive values
		1974	364834	
		1975	392503	

# Creating and Plotting Time Series in Python

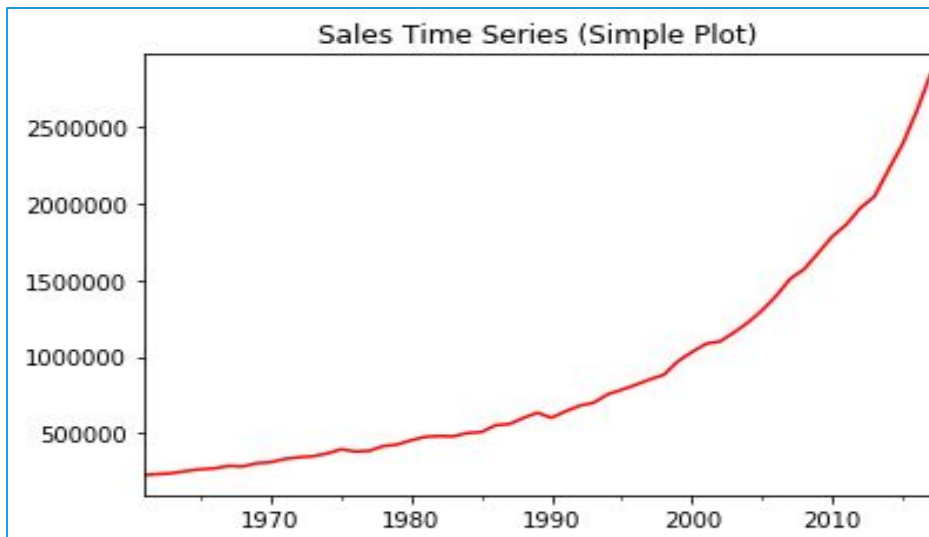
#Importing turnover\_annual data

```
import pandas as pd
salesdata=pd.read_csv('turnover_annual.csv')
```

#Creating and Plotting a Time Series Object

```
rng = pd.date_range('01-01-1961','31-12-2017',freq='Y')
s = salesdata.sales.values
salesseries = pd.Series(s, rng)
```

```
salesseries.plot(color='red', title ="Sales Time Series  
(Simple Plot)")
```



- **date\_range()**  
creates pandas date object.
- **freq='Y'**  
indicates yealy data
- **pd.Series()**  
creates time series object
- Plot function gives line chart

## Interpretation :

- The time-series clearly shows a positive trend.



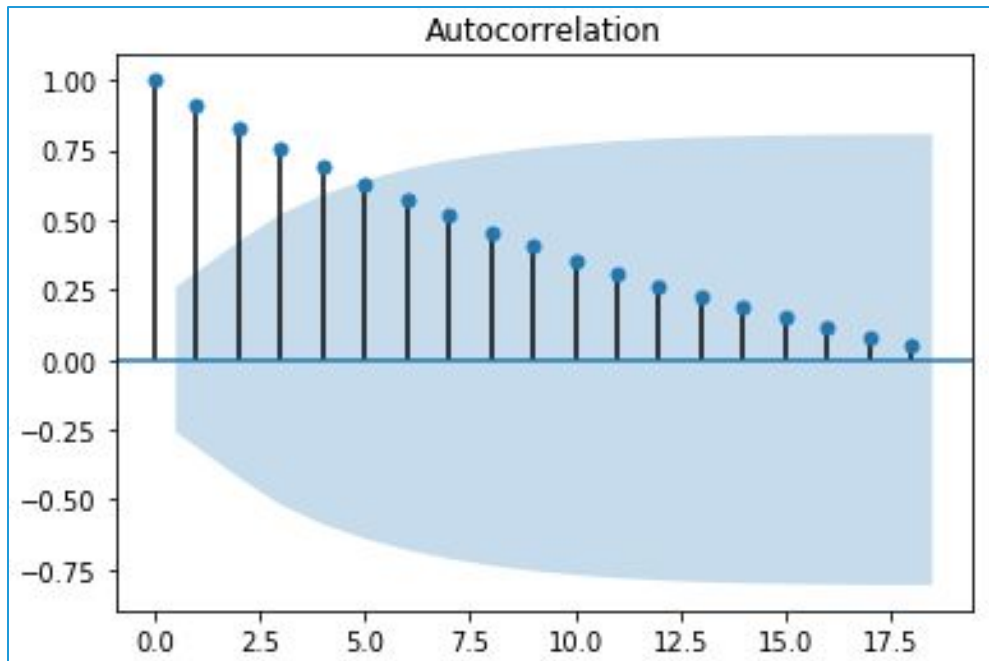
# Checking Stationarity – Correlogram

# ACF Plot

```
import matplotlib.pyplot as plt
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
plot_acf(salesseries)
```

□ **plot\_acf()** returns an ACF (Auto Correlation Function) plot.

# Output



## Interpretation :

□ We can observe that there is a very slow decay which is a sign of Non-stationarity.

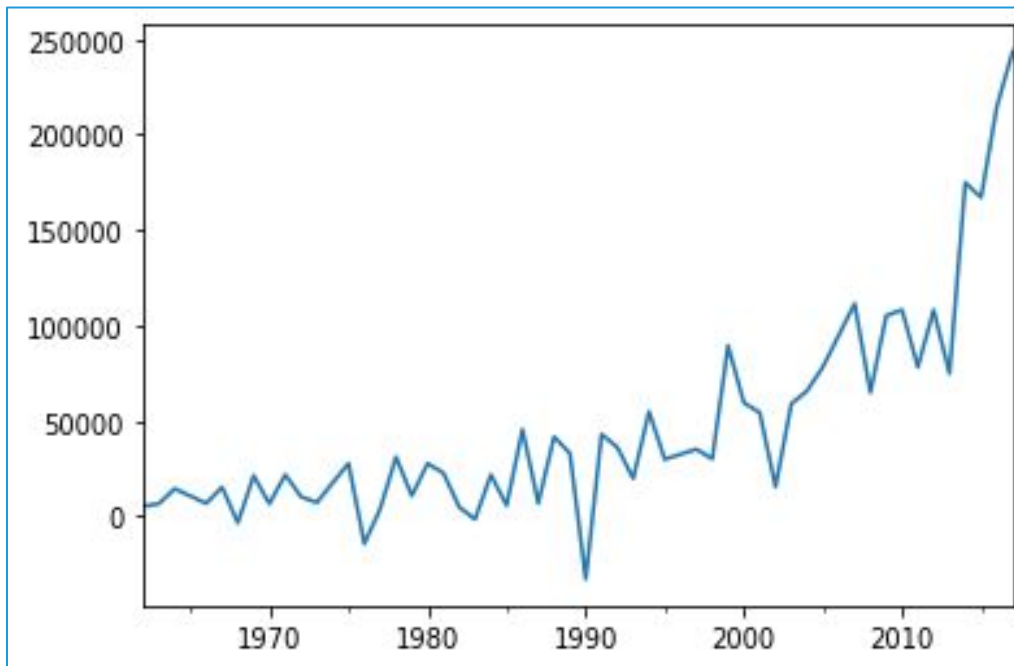
# Plot of 1<sup>st</sup> Order Differenced Time Series

# Creating and Plotting a Difference Series

```
from statsmodels.tsa.statespace.tools import diff
salesdiff = diff(salesseries)
salesdiff.plot()
```

- ❑ diff() gives 1<sup>st</sup> order differences
- ❑ plot function gives line chart for differenced series

# Output



## Interpretation :

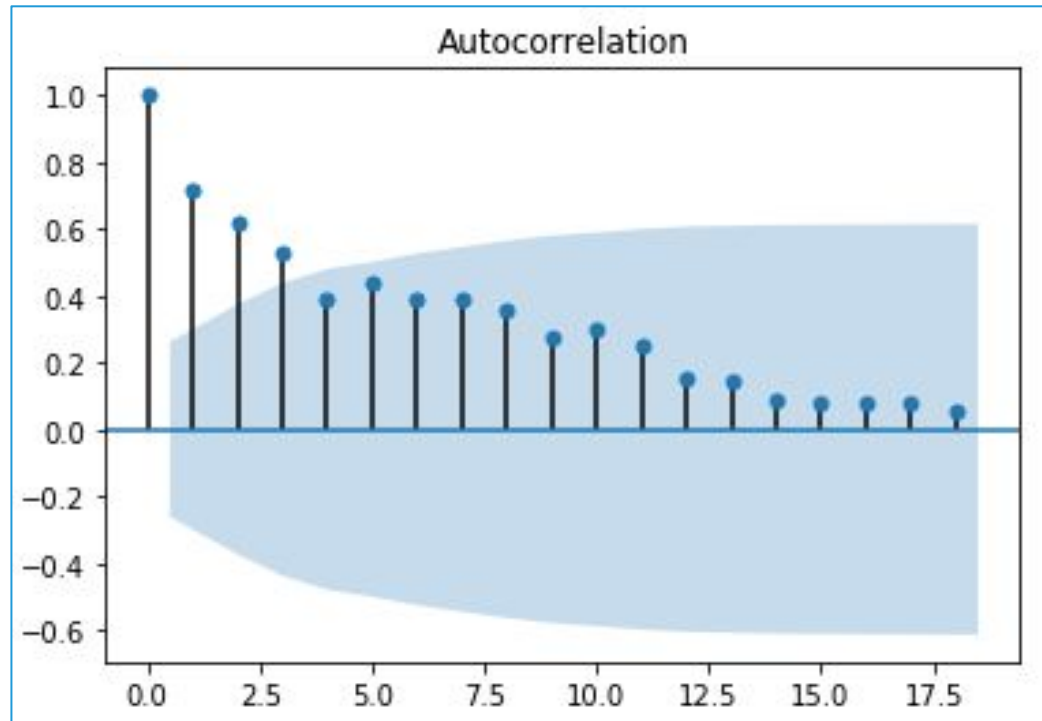
- ❑ Even after first order differencing, the series looks non-stationary.

# Correlogram for 1<sup>st</sup> Order Differenced Time Series

# ACF Plot

```
plot_acf(salesdiff)
```

# Output



## Interpretation :

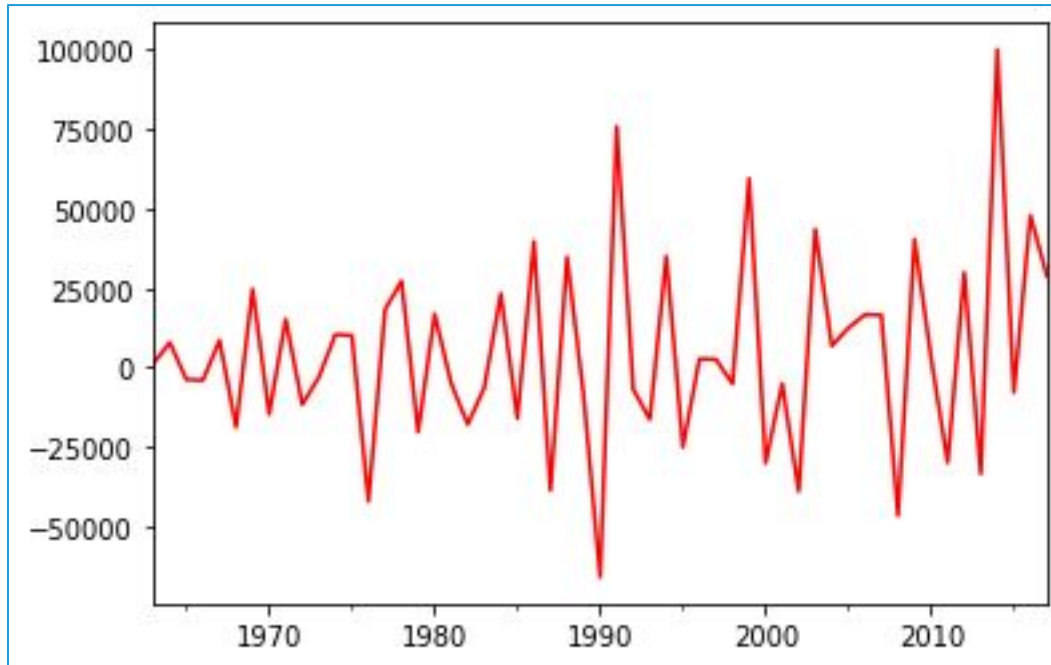
- ACF plot shows slow decay
- Stationarity is not achieved with first difference.

# Plot of 2nd Order Differenced Time Series

#Creating and Plotting 2<sup>nd</sup> Difference Series

```
salesdiff2 = diff(salesdiff)  
salesdiff2.plot(color='red')
```

# Output



## Interpretation :

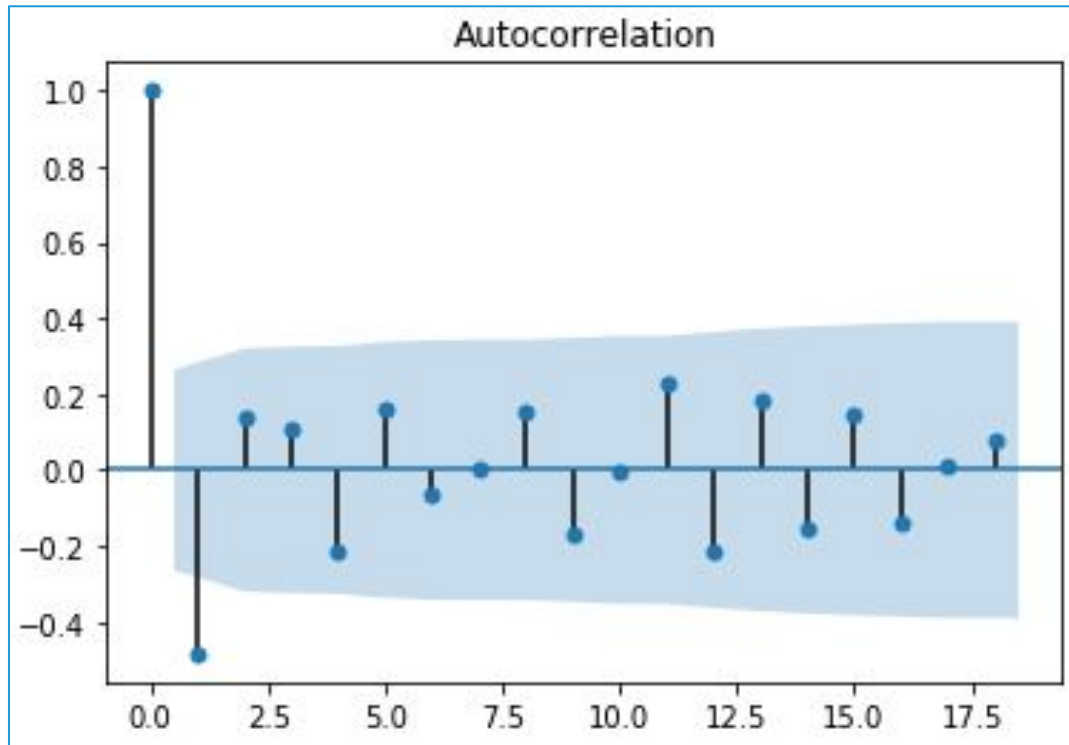
- After 2nd order differencing, the series looks **stationary**.

# Correlogram for 2<sup>nd</sup> Order Differenced Time Series

# ACF Plot

```
plot_acf(salesdiff2)
```

# Output



## Interpretation :

- Stationarity is achieved with 2<sup>nd</sup> order difference.

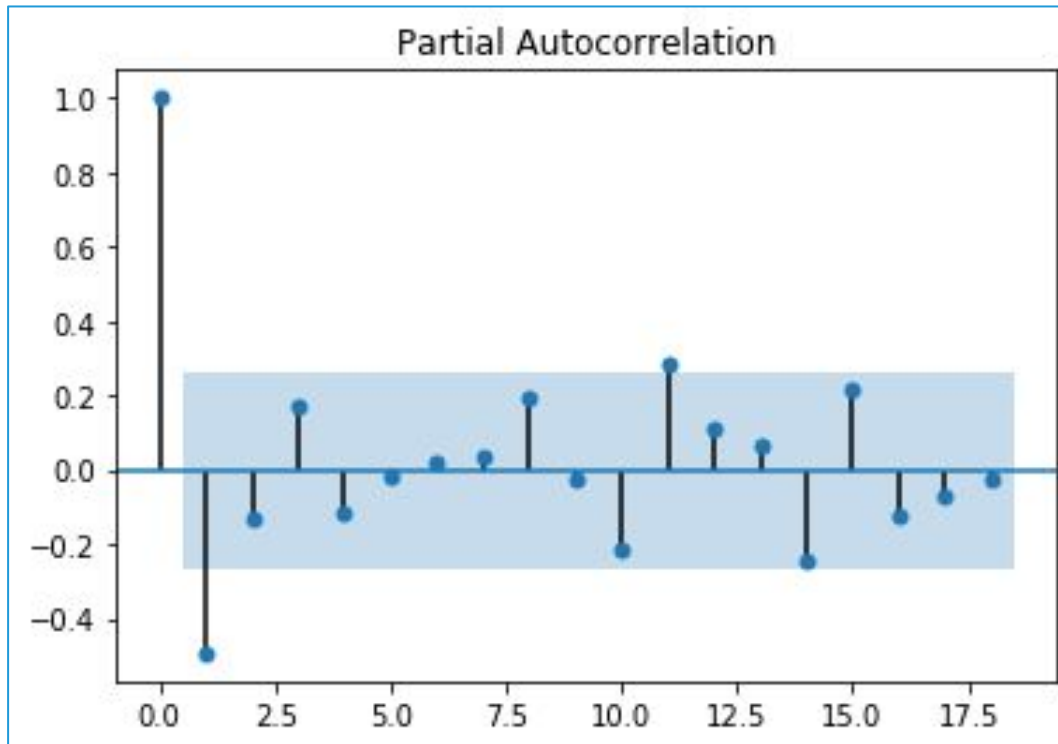
# Partial Autocorrelations for 2<sup>nd</sup> Order Differenced Time Series

# PACF Plot

```
plot_pacf(salesdiff2)
```

# Output

- **plot\_pacf()** returns an PACF (Partial Auto Correlation Function) plot.



## Interpretation :

- Stationarity is achieved with 2<sup>nd</sup> order difference.

# Analytical Method – Dickey Fuller (DF) Test

- A linear stochastic process has a unit root if 1 is the root of the process's characteristic equation. Such a process is non-stationary.
- Dickey and Fuller pioneered idea of testing for unit roots for stationarity checking.

Consider  $X_t$  ( $t=1,2,3,\dots$ ) is a time series of the form

$$X_t = \rho X_{t-1} + U_t \dots\dots\dots (1)$$

If  $\rho=1$  then  $X_t$  becomes a random walk

- We assume that  $U_t \sim \text{IID} (0, \sigma^2)$ , i.e  $U_t$  is a white noise
- Therefore, we are interested in testing for  $\rho=1$



IID means independent and identically distributed

# Dickey Fuller (DF) Unit Root Test

Objective	To test the <b>null hypothesis</b> that <b>time series is not stationary</b>
-----------	--

Null Hypothesis  $H_0: \rho=1$

Alternate Hypothesis  $H_1: \rho < 1$

$$(X_t - X_{t-1}) = \Delta X_t = (\rho - 1) X_{t-1} + U_t \dots \dots \dots \text{from (1)}$$

$$H_0: \rho^* = 0, H_1: \rho^* < 0, \rho^* = (\rho - 1)$$

Test Statistic	$(\rho^* / SE(\rho^*))$ Test statistic follows DF distribution under null
Decision Criteria	Reject the null hypothesis <b>tcal</b> < DF table value



# Dickey Fuller Test

```
# Install "arch"
```

```
pip install arch
```

```
# Import "ADF" from library "arch"
```

```
from arch.unitroot import ADF
```

```
adf = ADF(salesseries,lags=0,trend='nc')  
adf.summary()
```

- ❑ **ADF()** performs a Dickey Fuller unit root test on time series data.
- ❑ **lags=** allows to mention the number of lags to use in the ADF regression. We have used zero.
- ❑ **trend='nc'** specifies no trend and constant in regression

```
# Output
```

```
Augmented Dickey-Fuller Results  
=====
```

Test Statistic	19.275
P-value	1.000
Lags	0

```
-----
```

```
Trend: No Trend  
Critical Values: -2.61 (1%), -1.95 (5%), -1.61 (10%)  
Null Hypothesis: The process contains a unit root.  
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- ❑ Time series is non-stationary as value of test statistic is greater than 5% critical value.

# Dickey Fuller Test

# Checking stationarity for series with difference of order 2

```
adf = ADF(salesdiff2,lags=0,trend='nc')  
adf.summary()
```

# Output

```
Augmented Dickey-Fuller Results  
=====
```

Test Statistic	-11.908
P-value	0.000
Lags	0

```
-----  
  
Trend: No Trend  
Critical Values: -2.61 (1%), -1.95 (5%), -1.61 (10%)  
Null Hypothesis: The process contains a unit root.  
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- Time series is stationary as value of test statistic is less than 5% critical value.

# Quick Recap

## Correlograms

- **plot\_acf()** & **plot\_pacf()** function in Python generate Correlograms

## Differencing a Time Series

- Simple numeric function **diff()** can be used to difference a series

## Dickey Fuller Test

- **ADF()** function from the package **arch** performs a Dickey Fuller test
- The output gives test statistic and critical values for the test statistic

# Time Series Modeling

## ARIMA Model

# Contents

1. Box-Jenkins (ARIMA) Models
2. Five Step Iterative Procedure
  - i. Stationarity Checking
    - Differencing, Correlograms, and Dickey Fuller Test in Python
  - ii. Model Identification
  - iii. Parameter Estimation
    - Simple and Automated Model Estimation in Python
    - Running ARIMA in Python
  - iv. Diagnostic Checking
    - Residual plot in Python
  - v. Forecasting
    - Predictions on ARIMA Model in Python

# Box-Jenkins (ARIMA) Models

- ARIMA (Auto Regressive Integrated Moving Average) models are Regression models that use lagged values of the dependent variable and/or random disturbance term as explanatory variables.
- ARIMA models rely heavily on the autocorrelation pattern in the data.
- ARIMA models can also be developed in the presence of seasonality in the time series.
- ARIMA models thus essentially ignore domain theory (by ignoring “traditional” explanatory variables)

# When to Use ARIMA Models

Little or nothing is known about the dependent variable being forecasted

The independent variables known to be important cannot be forecasted effectively

Objective is to obtain short term forecasts

# Basic ARIMA Models

1. Autoregressive model of order  $p$  (AR( $p$ )):

$$y_t = \delta + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t$$

Where  $y_t$  depends on its  $p$  previous values

2. Moving Average model of order  $q$  (MA( $q$ ))

$$y_t = \delta + \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_q \varepsilon_{t-q}$$

$y_t$  depends on  $q$  previous random error terms

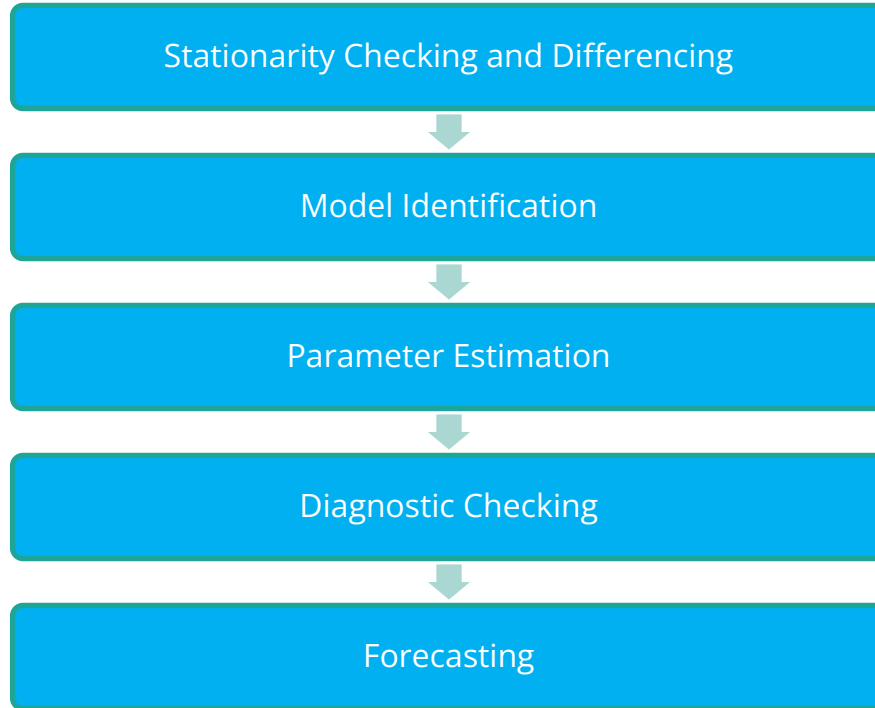
3. Autoregressive-Moving Average model of order  $p$  and  $q$  (ARMA( $p, q$ ))

$$y_t = \delta + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} \\ + \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_q \varepsilon_{t-q}$$

$y_t$  depends on its  $p$  previous values and  $q$  previous random error terms



# Five-Step Iterative Procedure




## Step 1: Stationarity Checking

# Differencing

- Differencing continues until stationarity is achieved

$$\Delta y_t = y_t - y_{t-1}$$

$$\Delta^2 y_t = \Delta(\Delta y_t) = \Delta(y_t - y_{t-1}) = y_t - 2y_{t-1} + y_{t-2}$$

- The differenced series has  $n-1$  values after taking the first-difference,  $n-2$  values after taking the second difference, and so on
- The number of times that the original series must be differenced in order to achieve stationarity is called the order of integration, denoted by  $d$  

?

How many times should a series be differenced?

In practice, it is not required to go beyond second difference.

# Case Study

## Background

- Annual Sales for a specific company from year 1961 to 2017

## Objective

- To develop time series model and forecast sales for next 3 years

## Available Information

- Number of cases: 57
- Variables: Year, sales(in 10's GBP)

# Data Snapshot

## turnover\_annual data

### Variables

Observations on Discrete Time Scale

Year	sales
1961	224786
1962	230034
1963	236562
1964	250960
1965	261615
1966	268316
1967	283589
1968	280160
1969	301422
1970	308018
1971	322025

Columns	Description	Type	Measurement	Possible values
Year	Financial Year	Numeric	-	-
sales	sales(in 10's GBP)	Numeric	In British Pound	Positive values
		1974	364834	
		1975	392503	

# Creating and Plotting Time Series in Python

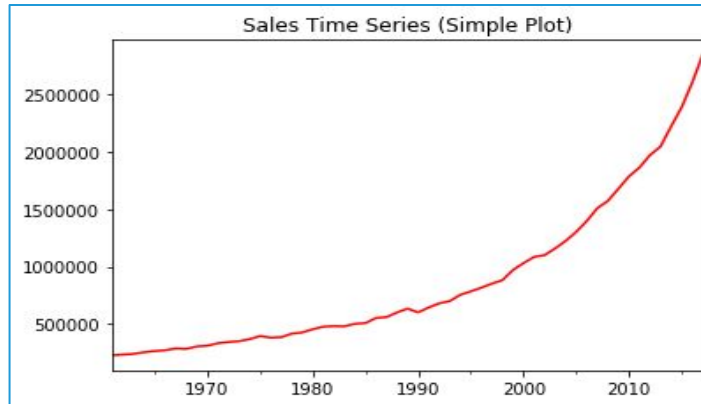
```
#Importing turnover_annual data
```

```
import pandas as pd  
salesdata=pd.read_csv('turnover_annual.csv')
```

```
#Creating and Plotting a Time Series Object
```

```
rng = pd.date_range('01-01-1961','31-12-2017',freq='Y')  
s = salesdata.sales.values  
salesseries = pd.Series(s, rng)
```

```
salesseries.plot(color='red', title ="Sales Time Series  
(Simple Plot)")
```



- ❑ **date\_range()** creates pandas date object.
- ❑ **freq='Y'** indicates yearly data
- ❑ **pd.Series()** creates time series object
- ❑ Plot function gives line chart

## Interpretation :

- ❑ The time-series clearly shows a positive trend.

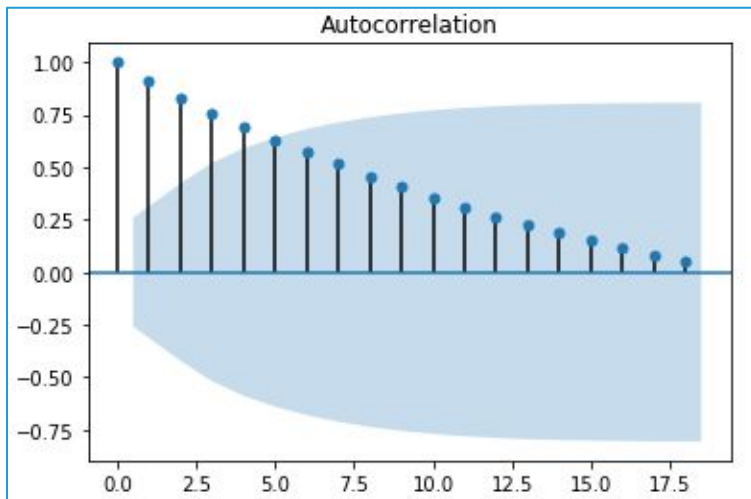
# Checking Stationarity – Correlogram

# ACF Plot

```
import matplotlib.pyplot as plt
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
plot_acf(salesseries)
```

- **plot\_acf()** returns an ACF (Auto Correlation Function) plot.

# Output



## Interpretation :

- We can observe that there is a very slow decay which is a sign of Non-stationarity.

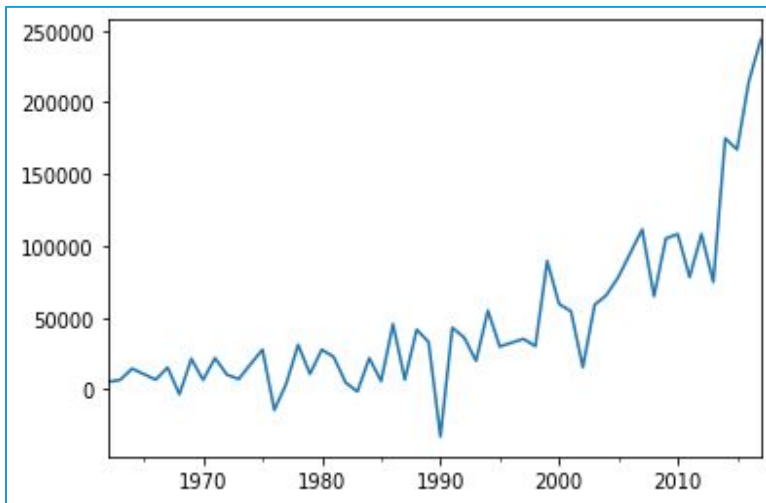
# Plot of 1<sup>st</sup> Order Differenced Time Series

# Creating and Plotting a Difference Series

```
from statsmodels.tsa.statespace.tools import diff  
salesdiff = diff(salesseries)  
salesdiff.plot()
```

- `diff()` gives 1<sup>st</sup> order differences
- `plot` function gives line chart for differenced series

# Output



## Interpretation :

- Even after first order differencing, the series looks non-stationary.

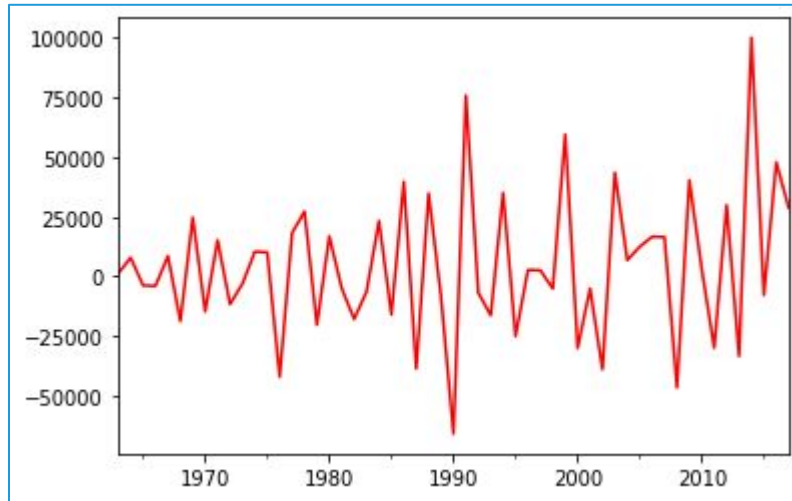


# Plot of 2nd Order Differenced Time Series

#Creating and Plotting 2<sup>nd</sup> Difference Series

```
salesdiff2 = diff(salesdiff)  
salesdiff2.plot(color='red')
```

# Output



## Interpretation :

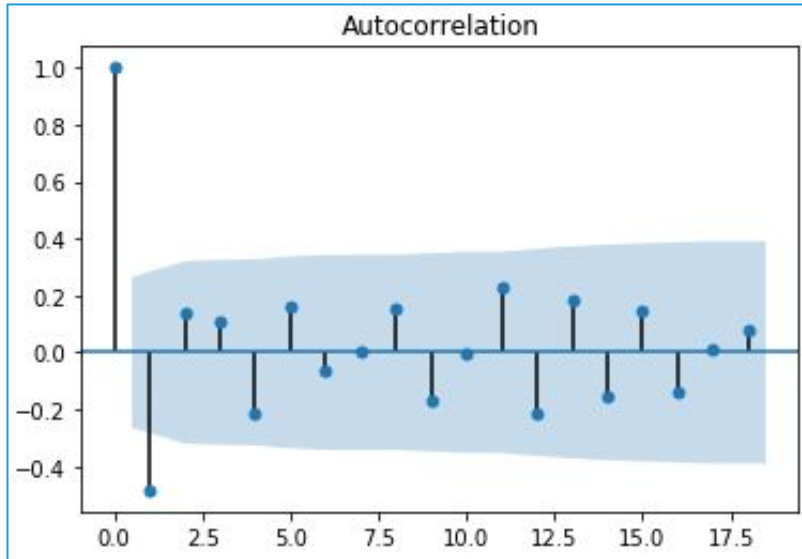
- After 2nd order differencing, the series looks **stationary**.

# Correlogram for 2<sup>nd</sup> Order Differenced Time Series

# ACF Plot

```
plot_acf(salesdiff2)
```

# Output



## Interpretation :

- Stationarity is achieved with 2<sup>nd</sup> order difference.

# Dickey Fuller Test

```
# Install "arch"
```

```
pip install arch
```

```
# Import "ADF" from library "arch"
```

```
from arch.unitroot import ADF
```

```
adf = ADF(salesseries,lags=0,trend='nc')  
adf.summary()
```

- ❑ **ADF()** performs a Dickey Fuller unit root test on time series data.
- ❑ **lags=** allows to mention the number of lags to use in the ADF regression. We have used zero.
- ❑ **trend='nc'** specifies no trend and constant in regression

# Output

```
Augmented Dickey-Fuller Results  
=====
```

Test Statistic	19.275
P-value	1.000
Lags	0

```
-----  
  
Trend: No Trend  
Critical Values: -2.61 (1%), -1.95 (5%), -1.61 (10%)  
Null Hypothesis: The process contains a unit root.  
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- ❑ Time series is non-stationary as value of test statistic is greater than 5% critical value.

# Dickey Fuller Test

```
# Checking stationarity for series with difference of order 2
```

```
adf = ADF(salesdiff2,lags=0,trend='nc')  
adf.summary()
```

```
# Output
```

```
Augmented Dickey-Fuller Results  
=====
```

Test Statistic	-11.908
P-value	0.000
Lags	0

```
-----  
  
Trend: No Trend  
Critical Values: -2.61 (1%), -1.95 (5%), -1.61 (10%)  
Null Hypothesis: The process contains a unit root.  
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- Time series is stationary as value of test statistic is less than 5% critical value.

## Step 2: Model Identification

# Model Identification

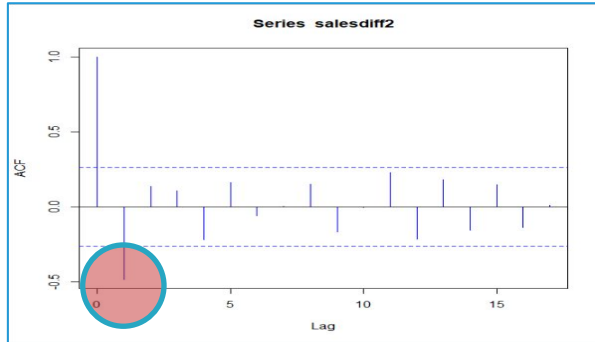
- When the data are confirmed stationary, proceed to tentative identification of models through **visual inspection of correlogram and partial correlogram**

Model	AC	PAC
	Dies down	Cuts off after lag $p$
	Cuts off after lag $q$	Dies down
	Dies down	Dies down

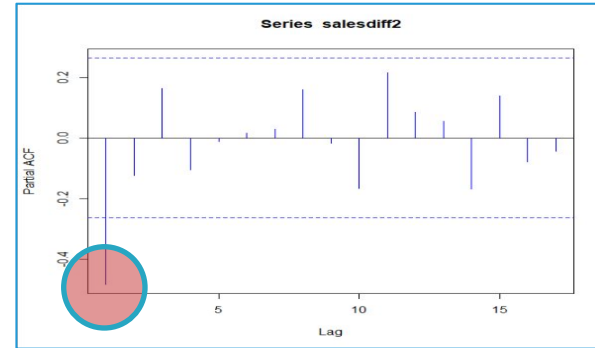
# Model Identification

- ARIMA model is expressed as  $\text{arima}(p,d,q)$  where
  - $p$  = no. of autoregressive terms
  - $d$  = order of differencing
  - $q$  = no. of moving average terms

ACF Plot



PACF Plot



- ACF and PACF correlograms will help in determining the MA and AR values respectively.

Indicative Model :  
 $\text{arima}(2,2,2)$

## Step 3: Parameter Estimation



# Parameter Estimation in Python

# Simple Estimation

```
from statsmodels.tsa.arima_model import ARIMA  
model = ARIMA(salesseries, order=(2, 2, 2)).fit(trend='nc')
```

- ❑ **ARIMA()** fits a model to a univariate time series.
- ❑ **order=** argument gives the model (p,d,q) order.

```
model.params  
model.aic
```

**params()** and **aic()** return the model coefficients and AIC value.

# Output

```
ar.L1.D2.None    -1.235413  
ar.L2.D2.None    -0.670320  
ma.L1.D2.None     0.785072  
ma.L2.D2.None     0.330062  
dtype: float64
```

# Output

```
1285.9836066562698
```

## Interpretation :

- ❑ Smaller the AIC value, better is the model. We need to try out various combinations of AR and MA terms to arrive at final model.

# Automatic Estimation of Model Parameters

# Automatic Model Identification and Parameter Estimation

```
import pmdarima as pm
model = pm.auto_arima(salesseries,max_p=2, max_q=2, d=2,
                      seasonal=False, trace=True)
```

- ❑ **auto\_arima()** generates the best order arima model. The function conducts a search over possible model within the order constraints provided.
- ❑ **trace=** True returns the list of all models considered.
- ❑ **max\_p and max\_q** gives maximum values of p and q respectively.
- ❑ **seasonal=** allows you to specify whether to fit a seasonal ARIMA or not.

# Automatic Estimation of Model Parameters

## # Output

```
Fit ARIMA: order=(2, 2, 2) seasonal_order=(0, 0, 0, 0); AIC=1294.586, BIC=1306.630, Fit
time=0.056 seconds
Fit ARIMA: order=(0, 2, 0) seasonal_order=(0, 0, 0, 0); AIC=1295.387, BIC=1299.402, Fit
time=0.006 seconds
Fit ARIMA: order=(1, 2, 0) seasonal_order=(0, 0, 0, 0); AIC=1291.568, BIC=1297.590, Fit
time=0.019 seconds
Fit ARIMA: order=(0, 2, 1) seasonal_order=(0, 0, 0, 0); AIC=1292.152, BIC=1298.174, Fit
time=0.018 seconds
Fit ARIMA: order=(0, 2, 0) seasonal_order=(0, 0, 0, 0); AIC=1294.532, BIC=1296.540, Fit
time=0.006 seconds
Fit ARIMA: order=(2, 2, 0) seasonal_order=(0, 0, 0, 0); AIC=1291.358, BIC=1299.387, Fit
time=0.022 seconds
Fit ARIMA: order=(2, 2, 1) seasonal_order=(0, 0, 0, 0); AIC=1292.717, BIC=1302.754, Fit
time=0.044 seconds
Fit ARIMA: order=(1, 2, 1) seasonal_order=(0, 0, 0, 0); AIC=1293.092, BIC=1301.122, Fit
time=0.029 seconds
Total fit time: 0.203 seconds
```

```
ARIMA(order=(2, 2, 0))
```

Lowest AIC

### Interpretation :

- Model with the lowest AIC value is selected as the best model.

# ARIMA Model Using BEST Order

```
# Run arima() for cross checking parameters based on model suggested  
# by auto.arima
```

```
from statsmodels.tsa.arima_model import ARIMA  
model = ARIMA(salesseries, order=(2, 2, 0)).fit(trend='nc')
```

```
model.params  
model.aic
```

```
# Output
```

```
ar.L1.D2.None    -0.506320  
ar.L2.D2.None    -0.101797  
dtype: float64
```

```
# Output
```

```
1285.114088844849
```

## Step 4: Diagnostic Checking

# Residual Analysis

If an ARMA( $p, q$ ) model is an adequate representation of the data generating process then the residuals should be 'White Noise'

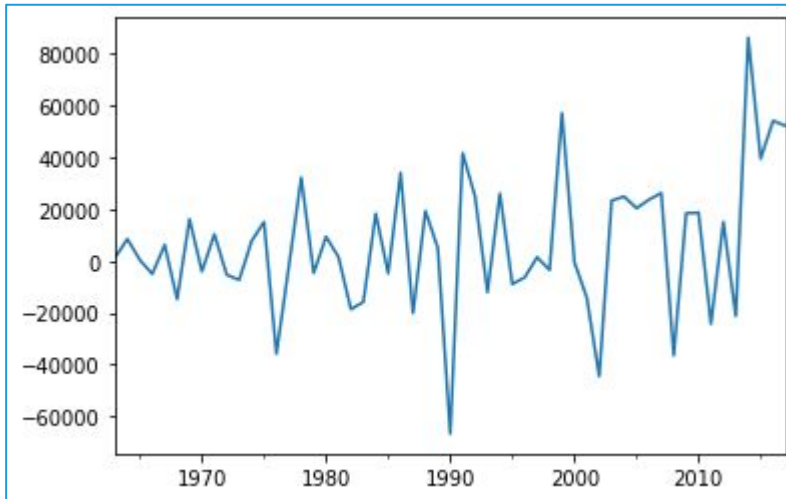
- White Noise time series has zero mean, constant variance and zero covariance with lagged time series.
- Residual plot is commonly used method for checking if the residuals are white noise process.

# Residual Plot In Python

```
resi = model.resid  
resi.plot()
```

**resid()** calculates residual values.

# Output



**Interpretation :**

- Errors follow white noise process.

## Step 5: Forecasting



# Forecasting

```
# Forecast for next 3years
```

```
model.forecast(steps=3)
```

**forecast()** function here gives predicted values for 3 years

```
# Output
```

```
(array([3072357.58495709, 3303466.18670477, 3533054.4598124 ]),  
array([27092.85937931, 48699.96239955, 75757.91735836]),  
array([[3019256.55633543, 3125458.61357874],  
       [3208016.0143532 , 3398916.35905634],  
       [3384571.67024626, 3681537.24937854]]))
```

## Interpretation :

- **forecast()** returns three arrays:
- array of three forecasts
- array of these standard error of the forecasts
- array of the confidence interval for the forecast

# Quick Recap

Stationarity Checking	<ul style="list-style-type: none"><li>• Plot correlogram using <code>plot_acf()</code> and ADF() for Dickey-Fuller Test</li></ul>
Model Identification	<ul style="list-style-type: none"><li>• Tentative identification of models through visual inspection of correlogram and partial correlogram</li></ul>
Parameter Estimation	<ul style="list-style-type: none"><li>• <b>auto_arima()</b> is recommended for obtaining best ARIMA model</li><li>• It uses AIC as the model selection criteria</li></ul>
Diagnostic Checking	<ul style="list-style-type: none"><li>• <b>Residual plot</b> for checking whether errors follow white noise process</li></ul>
Forecasting	<ul style="list-style-type: none"><li>• Use <b>forecast()</b> to generate forecasts</li></ul>

# Time Series Modeling

## Seasonal ARIMA Model

# Contents

1. Seasonal Box-Jenkins (ARIMA) Models
2. Five Step Iterative Procedure
  - i. Stationarity Checking and Seasonal Differencing
    - Differencing, Correlograms, and Dickey Fuller Test in Python
  - ii. Model Identification
  - iii. Parameter Estimation
    - Simple and Automated Model Estimation in Python
    - Running ARIMA in Python
  - iv. Diagnostic Checking
    - Residual plot in Python
  - v. Forecasting
    - Predictions on ARIMA Model in Python

# Seasonal Box-Jenkins (ARIMA) Models

- ARIMA (Auto Regressive Integrated Moving Average) models are Regression models that use lagged values of the dependent variable and/or random disturbance term as explanatory variables.
- Seasonal ARIMA (Often abbreviated as SARIMA) Model is formed by including seasonal terms in the ARIMA model.
- Several real world time series have a seasonal component. Some examples are: Sales of woolen clothes, demand for fertilizers, electricity consumption, etc.

# Seasonal Box-Jenkins (ARIMA) Models

- The **seasonal ARIMA model** incorporates both non-seasonal and seasonal factors in a multiplicative model.
- Shorthand notation for the model is,

$$\text{ARIMA } (p, d, q) \times (P, D, Q)_S,$$

with,

$p$  = non-seasonal AR order,

$d$  = non-seasonal differencing,

$q$  = non-seasonal MA order,

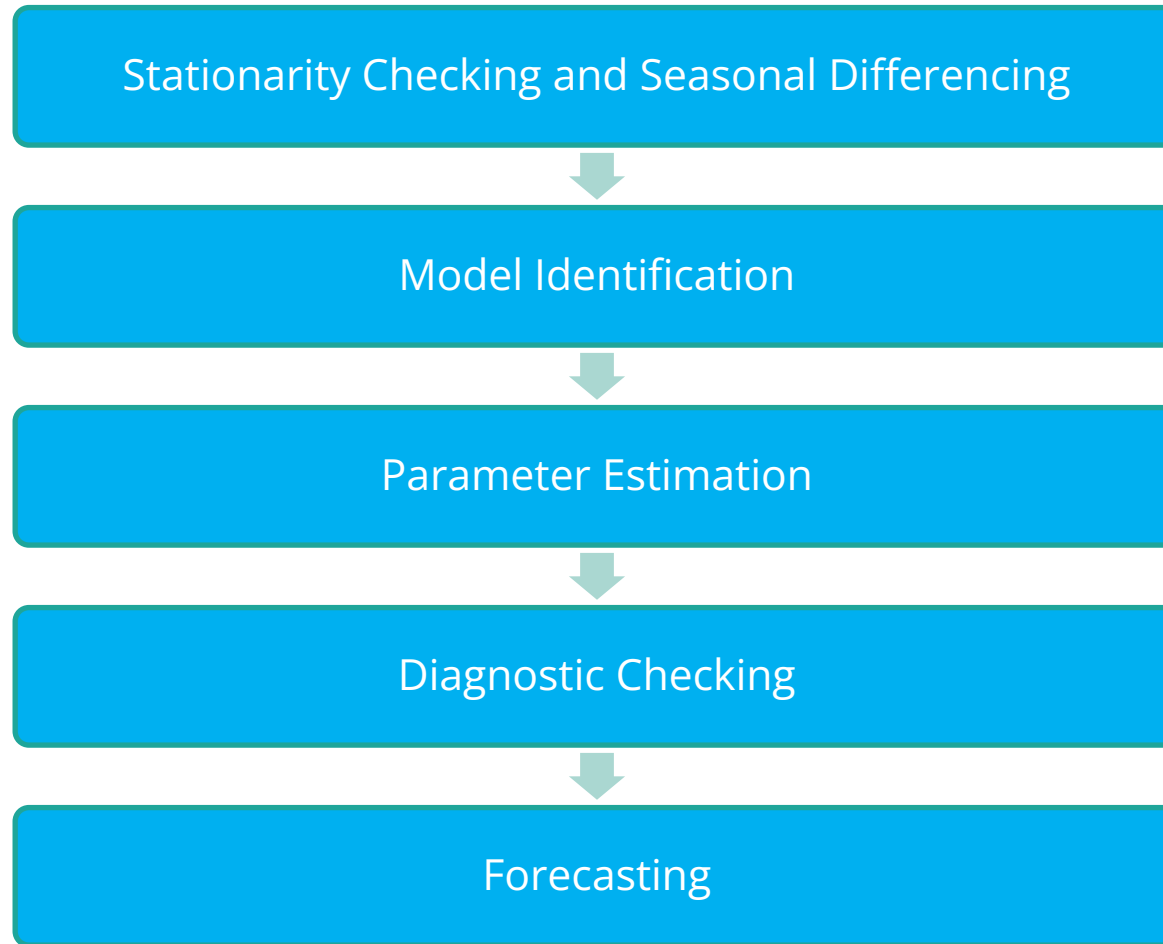
$P$  = seasonal AR order,

$D$  = seasonal differencing,

$Q$  = seasonal MA order, and

$S$  = time span of repeating seasonal pattern.

# Five-Step Iterative Procedure



# Step 1: Stationarity Checking



# Assessing Stationarity of Time Series

- Stationarity of a time series can be assessed using:



- If a time series is non-stationary then it can be converted via

Differencing

De-trending



Concept of Stationarity is explained previously in 'Stationarity in Time Series' ppt in detail.

# Seasonal Differencing

- Seasonal differencing is denoted as ,

$$\Delta_s y_t = y_t - y_{t-s}$$

Where,

$s$  denotes **frequency of season**

$s = 12$  if data is monthly;  $s = 4$  if data is quarterly and so on

- First and seasonal span differencing for monthly data is,

$$\Delta_1 \Delta_s y_t = \Delta_1 (y_t - y_{t-s}) = y_t - y_{t-1} - y_{t-s} + y_{t-s-1}$$

# Case Study

## Background

- Sales Data for 3 Years (2013, 2014, 2015)

## Objective

- To fit a Seasonal ARIMA Model and forecast next 3 Months sales.

## Available Information

- Sample size is 36
- Variables: Year, Month, Sales

# Data Snapshot

## Sales Data for 3 Years

Variables

Monthly Observations

Year	Month	Sales
2013	Jan	123
2013	Feb	142
2013	Mar	164
2013	Apr	173
2013	May	183
2013	Jun	192
2013	Jul	199
2013	Aug	203
2013	Sep	207
2013	Oct	209
2013	Nov	214
2013	Dec	255

Columns	Description	Type	Measurement	Possible values
Year	Year	numeric	2013, 2014, 2015	3
Month	Month	character	Jan - Dec	12
Sales	Sales in USD Million	numeric	USD Million	Positive values
	2014	Jul	245	

# Plotting a Time Series in Python

# Importing the Data

```
import pandas as pd  
salesdata = pd.read_csv('Sales Data for 3 Years.csv')
```

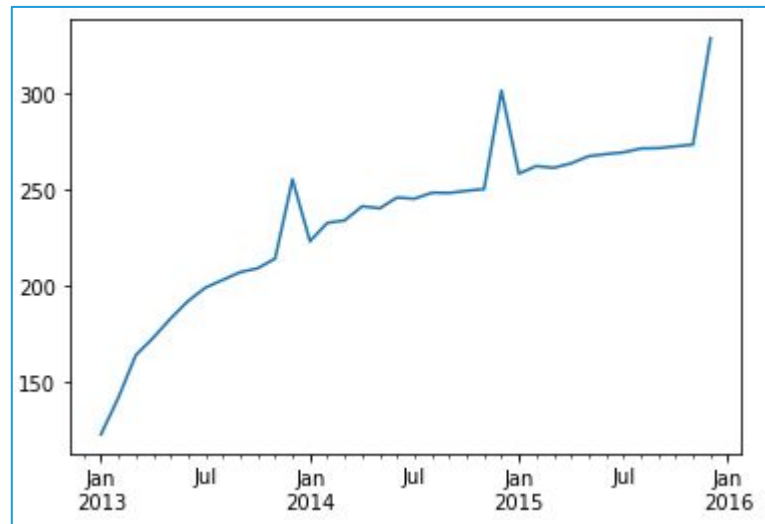
#Creating and Plotting a Time Series Object

```
rng = pd.date_range('01-01-2013', '31-12-2015', freq='M')
```

```
s = salesdata.Sales.values  
salesseries = pd.Series(s, rng)
```

```
salesseries.plot()
```

# Output



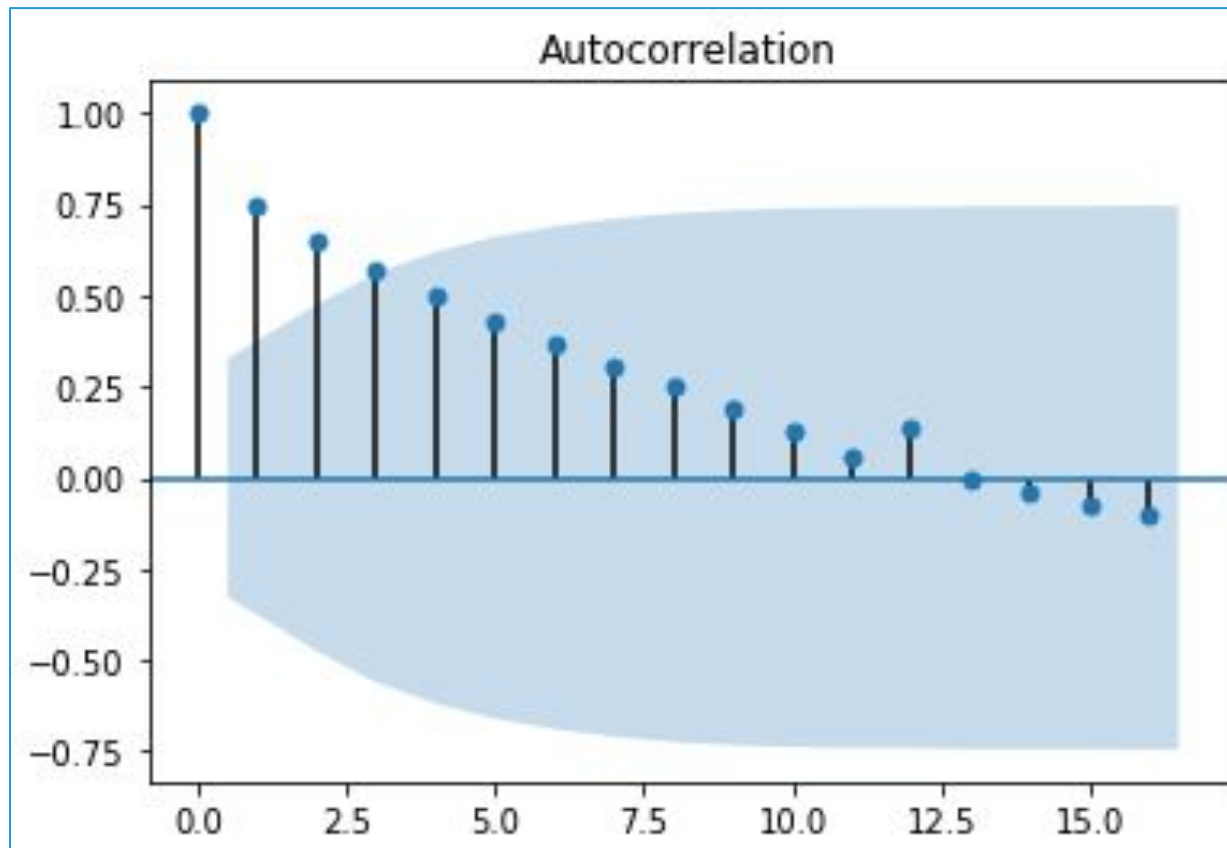
## Interpretation :

- The time series shows periodic peaks, indicative of seasonality.

# Correlogram

```
from statsmodels.graphics.tsaplots import plot_acf  
plot_acf(salesseries)
```

# Output



## Interpretation :

- ACF plot shows a slow decay indicating non-stationarity.

# Dickey Fuller Test

# Dickey Fuller Test

```
from arch.unitroot import ADF

adf = ADF(salesseries,lags=0,trend='nc')
adf.summary()
```

# Output

```
Augmented Dickey-Fuller Results
=====
Test Statistic          1.621
P-value                 0.975
Lags                    0
-----

Trend: No Trend
Critical Values: -2.63 (1%), -1.95 (5%), -1.61 (10%)
Null Hypothesis: The process contains a unit root.
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- Time series is non-stationary. Value of test statistic is greater than 5% critical value.

# Dickey Fuller Test – Differenced Series

```
# Dickey Fuller Test for Difference Series
```

```
from statsmodels.tsa.statespace.tools import diff
salesdiff = diff(salesseries)
(ADF(salesdiff,lags=0,trend='nc')).summary()
```

```
# Output
```

```
Augmented Dickey-Fuller Results
=====
Test Statistic              -6.891
P-value                     0.000
Lags                        0
-----

Trend: No Trend
Critical Values: -2.63 (1%), -1.95 (5%), -1.61 (10%)
Null Hypothesis: The process contains a unit root.
Alternative Hypothesis: The process is weakly stationary.
```

## Interpretation :

- Time series is stationary.  
Value of test statistic is less than 5% critical value.



## Step 2: Model Identification

# Model Identification

- When the data are confirmed stationary, proceed to tentative identification of models through visual inspection of correlogram and partial correlogram

Model	AC	PAC
	Dies down	Cuts off after lag $p$
	Cuts off after lag $q$	Dies down
	Dies down	Dies down

# Model Identification

- Seasonal ARIMA model is expressed as  $\text{arima}(p,d,q) (P,D,Q)$  where
  - $p$  = no. of autoregressive terms
  - $d$  = order of differencing
  - $q$  = no. of moving average terms
  - $(P,D,Q)$  are seasonal equivalents of autoregressive, difference and moving average terms

Fig. 3: ACF  
Plot

Series salesdiff

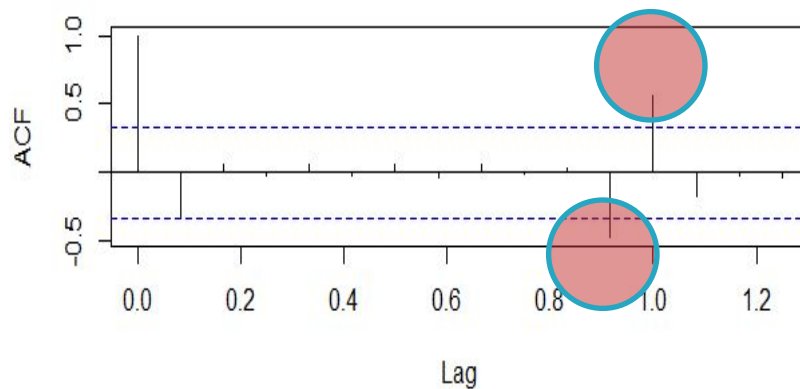
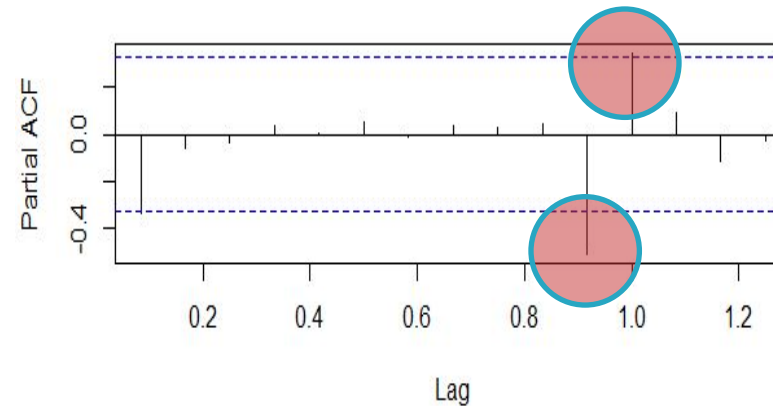


Fig. 4: PACF  
Plot

Series salesdiff



Indicative Model :  
 $\text{arima}(2,1,2)(2,1,2)$

## Step 3: Parameter Estimation

# Parameter Estimation

- There are two ways in which parameters of arima models can be estimated

1. Ordinary Least Squares

2. Maximum Likelihood Method – when the model involves MA component

- Given  $n$  observations  $y_1, y_2, \dots, y_n$ , the likelihood function  $L$  is defined as - the probability of obtaining the data actually observed
- The maximum likelihood estimators (MLE) are those values of the parameters for which the data actually observed are most likely, that is, the values that maximize the likelihood function  $L$ .

# Parameter Estimation in Python

# Automatic Model Identification and Parameter Estimation

```
import pmdarima as pm
model = pm.auto_arima(salesseries,
                      max_p=2, max_q=2,
                      max_P=2,max_Q=2,
                      d=1,m=12,
                      seasonal=True,
                      D=1, suppress_warnings=True,
                      trace=True)
```

model

- ❑ **auto\_arima()** generates the best order arima model. The function conducts a search over possible model within the order constraints provided.
- ❑ Seasonal model requires **max\_D,max\_P** and **max\_Q** arguments as well.
- ❑ **trace= True** returns the list of all models considered.

- ❑ **D=** gives order of seasonal differencing

# Automatic Model Identification

## # Output

```
Fit ARIMA: order=(2, 1, 2) seasonal_order=(1, 1, 1, 12); AIC=157.695, BIC=166.779, Fit
time=1.214 seconds
Fit ARIMA: order=(0, 1, 0) seasonal_order=(0, 1, 0, 12); AIC=152.639, BIC=154.910, Fit
time=0.075 seconds
Fit ARIMA: order=(1, 1, 0) seasonal_order=(1, 1, 0, 12); AIC=156.618, BIC=161.160, Fit
time=0.140 seconds
Fit ARIMA: order=(0, 1, 1) seasonal_order=(0, 1, 1, 12); AIC=156.623, BIC=161.165, Fit
time=0.118 seconds
Fit ARIMA: order=(0, 1, 0) seasonal_order=(0, 1, 0, 12); AIC=156.110, BIC=157.245, Fit
time=0.017 seconds
Fit ARIMA: order=(0, 1, 0) seasonal_order=(1, 1, 0, 12); AIC=154.629, BIC=158.035, Fit
time=0.089 seconds
Fit ARIMA: order=(0, 1, 0) seasonal_order=(0, 1, 1, 12); AIC=154.629, BIC=158.035, Fit
time=0.056 seconds
Fit ARIMA: order=(0, 1, 0) seasonal_order=(1, 1, 1, 12); AIC=156.629, BIC=161.171, Fit
time=0.127 seconds
Fit ARIMA: order=(1, 1, 0) seasonal_order=(0, 1, 0, 12); AIC=154.633, BIC=158.040, Fit
time=0.068 seconds
Fit ARIMA: order=(0, 1, 1) seasonal_order=(0, 1, 0, 12); AIC=154.636, BIC=158.042, Fit
time=0.067 seconds
Fit ARIMA: order=(1, 1, 1) seasonal_order=(0, 1, 0, 12); AIC=155.485, BIC=160.027, Fit
time=0.115 seconds
Total fit time: 2.118 seconds
```

```
ARIMA(order=(0, 1, 0), seasonal_order=(0, 1, 0, 12))
```

## Interpretation :

- Model with the lowest AIC value is selected as the best model.



Do note that the model identified by visually analysing correlograms was different. This shows it is not safe to rely just on visual review.

# ARIMA Model in Python

# Obtaining Coefficient

```
salesseries = pd.to_numeric(salesseries.astype(float))

from statsmodels.tsa.statespace.sarimax import SARIMAX
salesmodel = SARIMAX(salesseries, order=(0,1,0),
seasonal_order=(0,1,0,12)).fit(trend='nc')

salesmodel.params
salesmodel.aic
```

# Output

```
sigma2    47.585895
dtype: float64
```

# Output

```
156.10960236428923
```



# Model Selection Criteria

- Akaike Information Criterion (AIC)

$$\text{AIC} = -2 \ln(L) + 2k$$

- Schwartz Bayesian Criterion

(SBC, also called Bayesian Information Criterion - BIC)

$$\text{SBC} = -2 \ln(L) + k \ln(n)$$

where  $L$  = Likelihood function

$k$  = Number of parameters to be estimated

$n$  = Number of observations

Ideally, the AIC and SBC should be as small as possible

## Step 4: Diagnostic Checking

# Residual Analysis

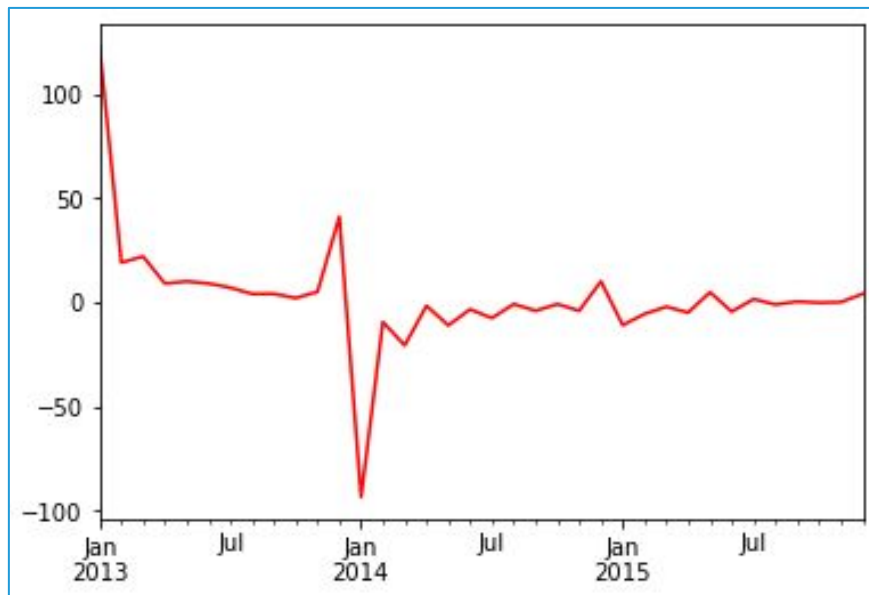
If an ARMA(p,q) model is an adequate representation of the data generating process then the residuals should be 'White Noise'

- White Noise time series has **zero mean, constant variance and zero covariance with lagged time series.**
- Residual plot is used for checking if the residuals are white noise process.

# Residual Plot In Python

```
resi = salesmodel.resid  
resi.plot(color="red")
```

# Output



## Step 5: Forecasting

# Forecasting

# Forecast for next 3years

```
salesmodel.forecast(steps=3)
```

# Output

```
2016-01-31    285.0  
2016-02-29    288.9  
2016-03-31    288.0  
Freq: M, dtype: float64
```

**forecast()** function is used to yield forecasts for a time series

# Quick Recap

## Stationarity Checking

- Plot correlogram using **plot\_acf()** and validate stationarity using **ADF()**

## Model Identification

- Tentative identification of models through visual inspection of correlogram and partial correlogram

## Parameter Estimation

- **auto\_arima()** is recommended for obtaining best ARIMA model & **SARIMAX()** for fitting the best model
- It uses AIC as the model selection criteria

## Diagnostic Checking

- **Residual plot** for checking whether errors follow white noise process

## Forecasting

- Use **forecast()** to generate forecasts

# Time Series Analysis –Exponential Smoothing Methods for Forecasting

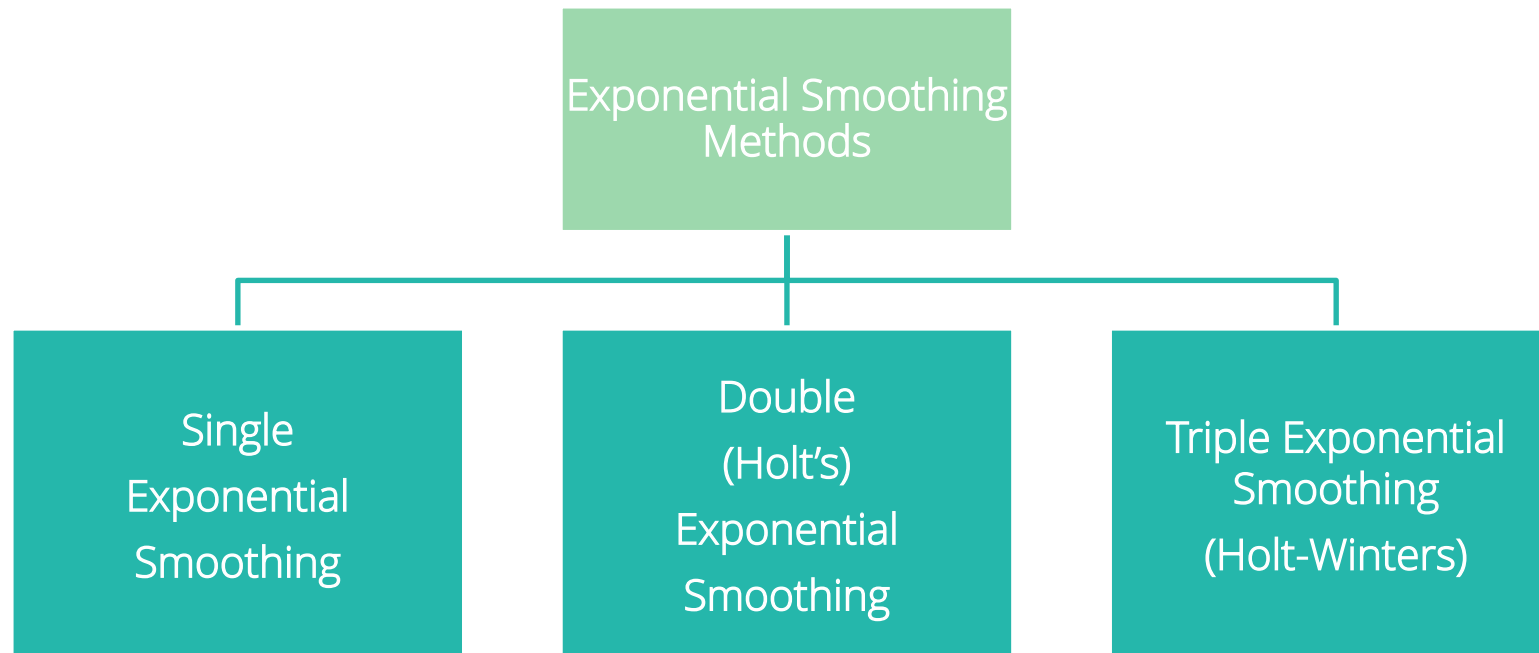


# Contents

1. Forecasting Using Smoothing Methods
2. Exponential Smoothing in Python
  - i. Single Exponential Smoothing
  - ii. Double Exponential Smoothing
  - iii. Triple Exponential Smoothing

# Forecasting Using Smoothing Methods

- Random, unexplained variation in a time series can have an undesirable impact on forecasts
- **Smoothing** can **cancel or reduce** such impacts
- Smoothing can either be Simple (using Moving Averages) or Exponential



# Single Exponential Smoothing Model

Mathematical Model :

$$F_{t+1} = \alpha Y_t + (1 - \alpha) F_t$$

Where,

$F_{t+1}$  : Forecast value for period  $t + 1$

$F_t$  : Forecast value for period  $t$

$Y_t$  : Actual value for period  $t$

$\alpha$  : Alpha (Smoothing constant)

# Single Exponential Smoothing Model

Assume  $\alpha=0.8$

t	yt	Ft	
1	23	-	
2	24	23	
3	26	23.80	$=0.8*24+0.2*23$
4	23.5	25.56	$=0.8*26+0.2*23.8$
5	27	23.91	
6	26.1	26.38	
7	28	26.16	
8	27	27.63	
9	29	27.13	
10	29.3	28.63	
11	28.2	29.17	
12	27	28.39	
		27.28	



1<sup>st</sup> future value will always be the previous value & then for rest future values exponential smoothing mathematical formula is applied.



# Single Exponential Smoothing Model - Smoothing Constant $\alpha$

## Values of $\alpha$

**close to one** ➡ have less of a smoothing effect and give greater weight to recent changes in the data

**closer to zero** ➡ have a greater smoothing effect and are less responsive to recent changes

- There is no formally correct procedure for choosing  $\alpha$ . Sometimes the statistician's judgment is used to choose an appropriate factor.
- Alternatively,  $\alpha$  can be decided based on statistical measure such as Root Mean Squared Error.

# Get an Edge!

Why the Name “Exponential”?

- This method gives weights to past observation in exponentially decreasing manner.

$$\begin{aligned}F_{t+1} &= \alpha y_t + \alpha(1-\alpha) y_{t-1} + \alpha(1-\alpha)^2 y_{t-2} + \alpha(1-\alpha)^3 y_{t-3} \dots \\&= \alpha y_t + (1-\alpha) [\alpha y_{t-1} + \alpha(1-\alpha) y_{t-2} + \alpha(1-\alpha)^2 y_{t-3} \dots] \\&= \alpha y_t + (1-\alpha) F_t\end{aligned}$$

- Larger alpha gives more weight to recent values.



# Case Study

## Background

- Sales Data for 3 Years (2013, 2014, 2015)

## Objective

- To apply Decomposition & Exponential Smoothing to Time Series data using different methods.

## Available Information

- Sample size is 36
- Variables: Year, Month, Sales



# Data Snapshot

## Sales Data for 3 Years

Variables

Observations

Year	Month	Sales
2013	Jan	123
2013	Feb	142
2013	Mar	164
2013	Apr	173
2013	May	183
2013	Jun	192
2013	Jul	199
2013	Aug	203
2013	Sep	207
2013	Oct	209
2013	Nov	214
2013	Dec	255

Columns	Description	Type	Measurement	Possible values
Year	Year	factor	2013, 2014, 2015	3
Month	Month	factor	Jan - Dec	12
Sales	Sales in USD Million	numeric	USD Million	Positive values



# Simple Exponential Smoothing in Python

# Import data

```
import pandas as pd
salesdata = pd.read_csv("Sales Data for 3 Years.csv")
rng = pd.date_range('2013', '2016', freq='M')
s = salesdata.Sales.values
salesseries = pd.Series(s, rng)
```

- ❑ **freq** = tells Python the frequency of time period in the data, 'M' for monthly data.
- ❑ **pd.Series()** converts a column from a data frame to a simple time series object.

# Simple Exponential Smoothing in Python

#Single Exponential Smoothing

```
from statsmodels.tsa.holtwinters import SimpleExpSmoothing
model = SimpleExpSmoothing(salesseries)
fit1 = model.fit()

fit1.predict()
fit1.summary()
```

**SimpleExpSmoothing()** from  
holtwinters in statsmodels  
undertakes exponential smoothing.

# Simple Exponential Smoothing in Python

# Output

```
2016-01-31    313.35045  
Freq: M, dtype: float64
```

Predicted  
value

SimpleExpSmoothing Model Results

<b>Dep. Variable:</b>	None	<b>No. Observations:</b>	36
<b>Model:</b>	SimpleExpSmoothing	<b>SSE</b>	10897.844
<b>Optimized:</b>	True	<b>AIC</b>	209.661
<b>Trend:</b>	None	<b>BIC</b>	212.828
<b>Seasonal:</b>	None	<b>AICC</b>	210.951
<b>Seasonal Periods:</b>	None	<b>Date:</b>	Tue, 17 Oct 2023
<b>Box-Cox:</b>	False	<b>Time:</b>	14:05:33
<b>Box-Cox Coeff.:</b>	None		

	coeff	code	optimized
<b>smoothing_level</b>	0.7351707	alpha	True
<b>initial_level</b>	129.80567	1.0	True

alpha

## Interpretation :

It returns predicted future value & value of alpha.



# Double (Holt) and Triple(Holt-Winters) Exponential Smoothing Methods

Double exponential smoothing has two equations

First equation is similar to single exponential smoothing method

Second equation updates trend using constant beta.

Double exponential smoothing method is used when there is a trend in the time series.

Triple exponential smoothing has three equations

First 2 equations are similar to double exponential smoothing method

Third equation updates seasonal component using constant gamma.

Triple exponential smoothing method is used when there is trend + seasonality in the time series.

# Double Exponential Smoothing Model

Mathematical Model :



Where,

$F_{t+1}$  : Forecast value for period  $t + 1$

$F_t$  : Forecast value for period  $t$

$T_t$  : Trend component for period  $t$

$T_{t+1}$  : Trend component for period  $t + 1$

$Y_t$  : Actual value for period  $t$

$\alpha$  : Alpha (Smoothing constant)

$\beta$  : Beta (Second smoothing constant)

# Double Exponential Smoothing in Python

```
#Double Exponential Smoothing
from statsmodels.tsa.holtwinters import ExponentialSmoothing
model = ExponentialSmoothing(salesseries, trend='add',seasonal = None)
fit2 = model.fit()
print(fit2.predict())
fit2.summary()
```

# Output

```
2016-01-31      295.970308
Freq: M, dtype: float64
```

ExponentialSmoothing Model Results

Dep. Variable:		None	No. Observations:	
Model:	ExponentialSmoothing		SSE	8649.636
Optimized:	True		AIC	205.343
Trend:	Additive		BIC	211.677
Seasonal:	None		AICC	208.240
Seasonal Periods:	None		Date:	Tue, 17 Oct 2023
Box-Cox:	False		Time:	14:15:55
Box-Cox Coeff.:	None			

	coeff	code	optimized
smoothing_level	0.3039444	alpha	True
smoothing_trend	0.3039444	beta	True
initial_level	127.50907	l.0	True
initial_trend	11.565240	b.0	True

Predicted  
value

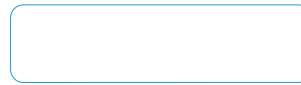
**Interpretation :**  
It returns  
predicted  
future value,  
value of  
alpha and  
beta.

alpha, b

# Triple Exponential Smoothing Model

Mathematical Model :

$$F_{t+1} = \alpha \frac{Y_t}{S_{t+1-k}} + (1 - \alpha) F_t - T_t$$



where,

$S_{t+1-k}$  : Seasonal smoothing value for period  $t + 1$

$F_t$  : Forecast value for period  $t$

$F_{t+1}$  : Forecast value for period  $t + 1$

$F_t$  : Forecast value for period  $t$

$T_t$  : Trend component for period  $t$

$T_{t+1}$  : Trend component for period  $t + 1$

$Y_t$  : Actual value for period  $t$

$\alpha$  : Alpha (Smoothing constant)       $\beta$  : Beta (Second smoothing constant)  
: Gamma (Third smoothing constant)

# Triple Exponential Smoothing in Python

#Triple Exponential Smoothing

```
from statsmodels.tsa.holtwinters import ExponentialSmoothing
model = ExponentialSmoothing(salesseries,seasonal_periods=12,
trend='add', seasonal='add')
fit3 = model.fit()
print(fit3.predict())
fit3.summary()
```

# Output

```
2016-01-31    293.447983
Freq: M, dtype: float64
```

Predicted  
value

**Interpretation :**  
It returns predicted  
future value &  
value of alpha, beta  
and gamma.

ExponentialSmoothing Model Results

Dep. Variable:	None	No. Observations:	
Model:	ExponentialSmoothing	SSE	539.668
Optimized:	True	AIC	129.468
Trend:	Additive	BIC	154.804
Seasonal:	Additive	AICC	169.703
Seasonal Periods:	12	Date:	Tue, 17 Oct 2023
Box-Cox:	False	Time:	14:06:06
Box-Cox Coeff.:	None		

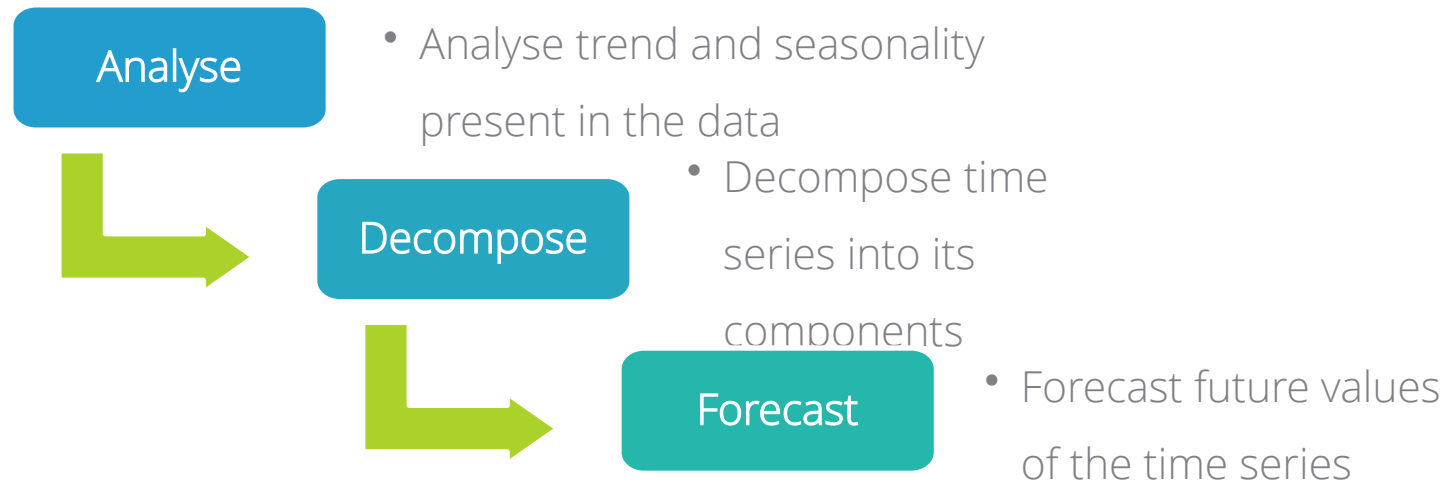
	coeff	code	optimized
smoothing_level	1.0000000	alpha	True
smoothing_trend	0.3703030	beta	True
smoothing_seasonal	1.9375e-08	gamma	True
initial_level	121.18636	l.0	True
initial_trend	13.392973	b.0	True

alpha, beta,  
gamma



# Get an Edge!

Always approach time series analysis in a systematic manner



- For vector time series, investigate connections between two or more time series with the aim of using values of some of the processes to predict those of the others. (Eg. Pairs trading in stock market)

# Quick Recap

In this session, we learnt about **exponential smoothing**:

## Smoothing

- Smoothing gives weights to past observations, in order to give more significance to seasonality and trend components of a time series

## Smoothing in Python

- From **statsmodels.tsa.holtwinters** :
- Use **SimpleExpSmoothing()** to carry out simple exponential smoothing
- Use **ExponentialSmoothing()** to carry out double and triple exponential smoothing

# The Concept of Co-integration

# Contents

1. Concept of Co-integration
2. What is Pairs Trading
3. Co-integration and Pairs Trading
4. Engle-Granger Test
5. EG Test in Python

# Concept of Co-integration

- The stock prices are believed to be **non-stationary** in nature.
- If there exists a linear combination of 2 stock prices , such that errors of the long term relationship are stationary (mean reverting) in nature, then the 2 stocks are said to be co-integrated.
- **Engle - Granger test** is used to identify the co-integrated pairs in a particular sector.

# Co-Integration

## Examples of Co-integrated Series

Income and Consumption

Money, National Incomes, Price Rates, Interest Rates

Price of a Commodity and Taxes Levied on that Commodity, Inflation Rate

Temperatures, Electricity Consumption

Prices of Two Stocks

# Co-Integration-Formal Definition

- “Order of integration” tells you the minimum number of differences needed to get a stationary series.
- A series of successive differences,  $d$ , can transform the time series into one with stationarity. The differences are denoted by  $I(d)$ , where  $d$  is the order of integration.
- Suppose  $Y_t$  and  $X_t$  are two time series integrated of order  $d$ , then
  - Any linear combination of such two series will also be integrated of order  $d$  (denoted as  $I(d)$ ). This is called ‘Integration’
- However,

If there exists a vector  $\beta$  such that  $u_t = Y_t - \beta X_t$  is of a lower order of integration ( $I(d - b)$ , where  $b > 0$ ) then  $Y_t$  and  $X_t$  are defined as Cointegrated of order  $(d, b)$

- The idea of cointegration was introduced by Engle and Granger in 1987



Usually, the order of integration is either  $I(0)$  or  $I(1)$ ; It's rare to see values for  $d$  that are 2 or more.



**DATA SCIENCE**  
INSTITUTE

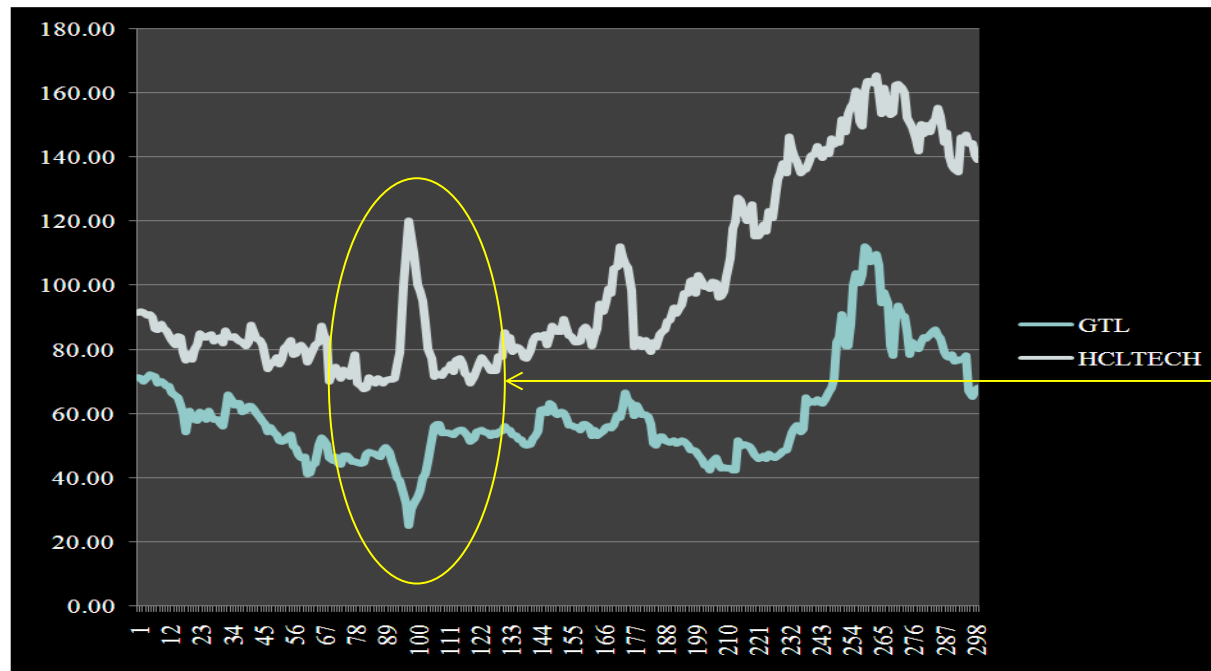
# What is Pairs Trading

- Pairs trading finds its roots in the area of securities trading
- A quest to generate returns, irrespective of market behavior, has led to the evolution of this strategy
- Certain securities, often competitors in the same sector, are associated in their day-to-day price movements. When the association breaks down, Traders take advantage of short term over pricing or under pricing of securities & invest accordingly, betting that the "spread" between the two would eventually converge
- Pairs Trading strategy works best in volatile market conditions



# Simple Example of Pairs Trading

Fig. 1: Pairs Trading



This is the right opportunity for Pairs Trade

TRADER WILL GAIN PROFIT IF HE SELLS HCLTECH AND BUYS GTL

# Co-Integration and Pairs Trading

- The **stock prices** are believed to be **non-stationary** in nature.
- If there exists a linear combination of 2 stock prices , such that errors of the long term relationship are stationary (mean reverting) in nature, then the 2 stocks are said to be co-integrated.
- Pairs trading is in fact an extended application of the theory of co-integration.
- **Engle -Granger test** is used to identify the co-integrated pairs in a particular sector.

# Engle Granger Test for Stock Prices

Let  $Y_t$  : TCS stock price at time  $t$

Let  $X_t$  : IBM stock price at time  $t$

Both the stock prices form non-stationary time series over time  $t$ .

If there exists a linear combination, such that the errors form stationary time series then TCS and IBM stocks are said to be co-integrated.

$$Y_t = \alpha + \beta X_t + \varepsilon_t$$

$\alpha$  = model intercept

$\beta$  = regression coefficient

$\varepsilon_t$  = error term

# Engle Granger Test for Stock Prices

Regression coefficient  $\beta$  is estimated using ordinary least square (OLS) method

If  $\beta < 0$  then the pair is discarded

Selection of regression equation is done by comparing  $\beta$

$$Y_t = \alpha_1 + \beta_1 X_t + \varepsilon_t \dots\dots\dots 1$$

$$X_t = \alpha_2 + \beta_2 Y_t + \varepsilon_t \dots\dots\dots 2$$

The equation with higher  $\beta$  is selected

In case of TCS and IBM stocks  $\beta_1 > \beta_2$  hence equation 1 was selected

# Engle Granger Test

Engle Granger test is used to check stationarity of error time series.

Consider,

$$\Delta \varepsilon_t = \gamma^* \varepsilon_{t-1} + W_t$$

$$\Delta \varepsilon_t = \varepsilon_t - \varepsilon_{t-1}$$

$\varepsilon_{t-1}$  = errors with lag 1

Objective	To test the null hypothesis that time series is not stationary
-----------	--

Null Hypothesis H0: Two stocks are not co-integrated  
Alternate Hypothesis H1: Two stocks are co-integrated

Test Statistic	$(\gamma^*/SE(\gamma^*))$ Test statistic follows DF distribution under null
Decision Criteria	Reject the null hypothesis if $tcal < c(p)$



# Engle Granger Test in Python

# Installing and Loading Required Libraries

```
import yfinance as yf
import pandas as pd
```

- **yfinance** is a library that allows import of data directly from yahoo finance.

#EG Test

```
# Define the stock symbols for Wipro and Infosys
tickers = ["WIPRO.NS", "INFY.NS"]
```

```
start_date = "2022-01-01"
end_date = "2022-12-31"
```

```
# Download the historical data
data = yf.download(tickers, start=start_date, end=end_date)
adj_close_data = data['Adj Close']
adj_close_data.columns = ['INFY', 'WIPRO']
```

```
infy = adj_close_data.INFY
wipro = adj_close_data.WIPRO
```

- **tickers** is a set symbols used to extract data.
- **start\_date, end\_date** is the start date and end date of the date range

# Engle Granger Test in Python

```
# Import "coint" from "statsmodels"  
from statsmodels.tsa.stattools import coint  
coint(infy,wipro)
```

- **coint()** is a function that runs the engel garner test and returns t-statistic, p-value, Critical values for the test statistic at the 1 %, 5 %, and 10 % levels based on regression curve.

# Output

```
(-2.154422825452808,  
0.44797584761041354,  
array([-3.94132922, -3.36097908, -3.06166543]))
```

## Interpretation :

- As  $t_{cal} > c(p)$  at 5%, we accept null Hypothesis.
- Stocks INFY & WIPRO are not cointegrated.



What if we want to get data from other sources?

There are different packages using database APIs for data from other indices, eg. `insepy` for India's National Stock Exchange (NSE) prices



**DATA SCIENCE**  
INSTITUTE

# Get an Edge!

There are a number of alternative approaches for testing co-integration, apart from the Engle Granger test, such as

- **CRDW (Cointegrating Regression Durbin Watson) test** ✉ Simple regression of one variable on the other, and the standard Durbin-Watson test on the residuals.
- **Error Correction Test** ✉ The tendency of cointegrated variables to revert to common stochastic trends is expressed in terms of error-correction
- **Johansen's Multivariate VAR Approach** ✉ Examines the number of independent linear combinations ( $k$ ) for an  $m$  time series variables set that yields a stationary process





# Quick Recap

In this session, we learnt about **co-integration in time series**:

## Co-integration

- When the time series is integrated and linear combination of variables in that series is also integrated, but having an order lower than the whole series, then the variables are said to be co-integrated.

## Pairs Trading

- When the association between two cointegrated stocks breaks down, Traders take advantage of short term over pricing or under pricing of securities & invest accordingly, betting that the "spread" between the two would eventually converge.

## EG Test in Python

- **coint()** is used to perform Engle Granger test in Python, to check if a pair of time series is co-integrated.

