Introduction

Different Types Of Dataset

1. <u>Cross-sectional data</u>: Cross sectional data can obtained when there is multiple observations is taken from multiple individuals at same point time. In case of cross sectional data Time does not play any important role in analysis. Analysis of cross sectional data starts with visualization of basic statistical properties i.e. central tendency, dispersion, skewness, kurtosis.

- 2. <u>Time series data</u>: Timeseries data can obtained when there is multiple observations is taken form same source at different points of time. Time series data can be describe by trend, seasonality, stationarity, autocorrelation, and so on.
- 3. <u>Panel data</u>: Panel data is collection of multiple entities over multiple points in time. Panel data also known as longitudinal data.

Here we have Nifty50 data from Apr-2010 to Mar-2018

https://www.nseindia.com/products/content/equities/indices/historical_index_data.htm

This is time series data with the frequency of business days.

First we need to explore the data to check the presence of null values

```
import pandas as pd
Nifty_data=pd.read_csv("E:/summer/NIFTY50.csv",parse_dates=['Date'],index_col=
['Date'])
Nifty_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 1987 entries, 2010-04-01 to 2018-03-28
Data columns (total 6 columns):
Open
                    1987 non-null float64
High
                     1987 non-null float64
                     1987 non-null float64
Low
                     1987 non-null float64
Close
Shares Traded
                     1987 non-null int64
Turnover (Rs. Cr)
                    1987 non-null float64
dtypes: float64(5), int64(1)
memory usage: 108.7 KB
```

Let's understand the arguments one by one:

parse dates: This specifies the column which contains the date-time information. As we say above, the column name is 'date'.

index col: This argument tells pandas to use the 'date' column as index.

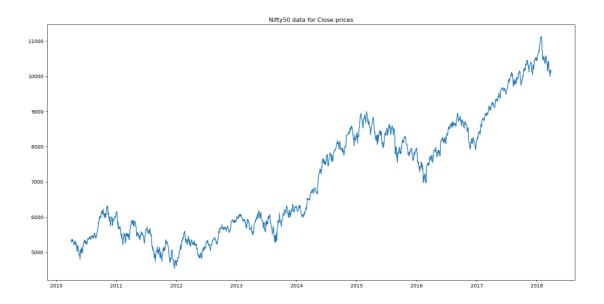
In this dataset we have DateTime series as index and 6 columns, 1987 entries for each. 5 of them are float and one is integer type also we don't have any Null values.

Internal structure of Timeseries

A time series is combination of trend, seasonal, cyclical, and irregular components

• **General trend**: A general trend can be identify by it's upward or downward movement in a long run. General trend can be seen with plotting the data, here we are using matplotlib library to plot Close prices.

```
import matplotlib.pyplot as plt
plt.plot(Nifty_data.Close)
plt.title('Nifty50 data for Close prices')
plt.show()
```



The graph is showing movement in upward direction, this is clear sign of presence of trend component.

General trend might not capable of being noticed during short run of time series.

- <u>Seasonality</u>: If in time series data, there are patterns that repeat over known periods of time. For example the consumption of ice-cream during summer is more than winter and hence an ice-cream dealer's sales would be higher in summer months. Mostly, presence of seasonality can be reveals by exploratory data analysis.
- <u>Cyclical movements</u>: If there are movements observe after every few units of time, but they are not as frequent as seasonal components, are known as Cyclic components. cyclic components do not have fixed periods of variations.
- <u>Unexpected variations</u>: These are sudden changes occurring in a time series which are unlikely to be repeated. They are components of a time series which cannot be explained by trends, seasonal or cyclic

movements. These variations are sometimes called residual or random components. These variations, though accidental in nature, can cause a continual change in the trends, seasonal and cyclical oscillations during the forthcoming period.

The objective of time series analysis is to decompose timeseries into it's elements and develop mathematical models to predict future values.

Stationary time series

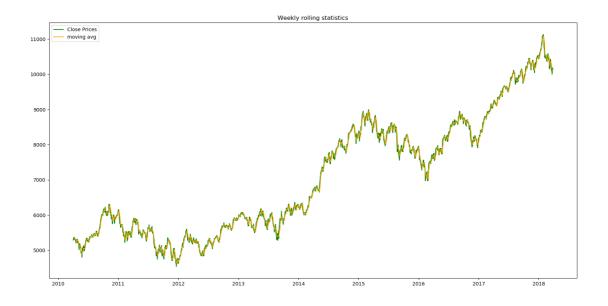
When a time series free from general trend and seasonality, it become stationary. Statistical properties like mean, variance, autocorrelation etc. of an stationary time series remains constant over time. It is important to gain stationarity before forecasting because most statistical forecasting methods are applicable on stationary time series.

we can check stationarity using the following:

- Plotting rolling statistics: we can plot moving mean and moving variance and see if these terms are
 varying with time. let's plot the moving average for business week days, business month days,

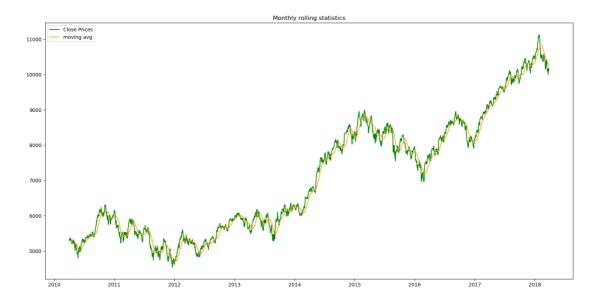
 Quarterly and yearly. we have .rolling() method in pyhon to calculate rolling statistics.
- weekly

```
plt.plot(Nifty_data.Close,label='Close Prices',color='green')
plt.plot(Nifty_data['Close'].rolling(window=5).mean(),label='moving
avg',color='orange')
plt.legend()
plt.title('Weekly rolling statistics')
plt.show()
```



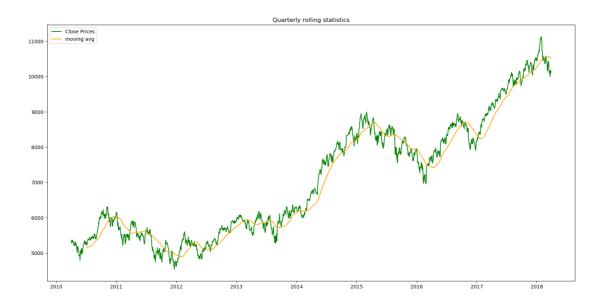
Monthly

```
plt.plot(Nifty_data.Close,label='Close Prices',color='green')
plt.plot(Nifty_data['Close'].rolling(window=21).mean(),label='moving
avg',color='orange')
plt.legend()
plt.title('Monthly rolling statistics')
plt.show()
```



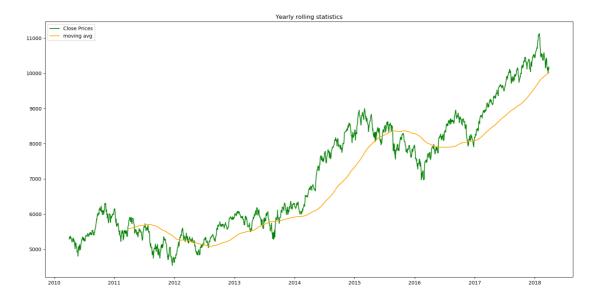
Quarterly

```
plt.plot(Nifty_data.Close,label='Close Prices',color='green')
plt.plot(Nifty_data['Close'].rolling(window=63).mean(),label='moving
avg',color='orange')
plt.legend()
plt.title('Quarterly rolling statistics')
plt.show()
```



yearly

```
plt.plot(Nifty_data.Close,label='Close Prices',color='green')
plt.plot(Nifty_data['Close'].rolling(window=252).mean(),label='moving
avg',color='orange')
plt.legend()
plt.title('Yearly rolling statistics')
plt.show()
```



As we can see moving average is changing with time so, Time series is not stationary.

2. <u>Augmented Dickey Fuller test</u>: This is one of the unit root tests for checking stationarity. In this this test we check for null hypothiesis, where H_0 stats that time series is stationary whereas H_A stats that

time series is not Stationary. In python we have stattools.adfuller function to check the stationarity of time series.

```
from statsmodels.tsa.stattools import adfuller
stnry_test=adfuller(Nifty_data['Close'],autolag='AIC')
stnry_rslt = pd.Series(stnry_test[0:4], index=['Test Statistic','p-value','#Lags
Used','Number of Observations Used'])
for key,value in stnry_test[4].items():
    stnry_rslt['Critical Value (%s)'%key] = value
print(stnry_rslt)
if(stnry_test[1]>0.05):
    print("Time Series is not stationary")
else:
    print("Time series is stationary")
```

Test Statistic	-0.371803
p-value	0.914701
#Lags Used	1.000000
Number of Observations Used	1985.000000
Critical Value (1%)	-3.433649
Critical Value (5%)	-2.862997
Critical Value (10%)	-2.567546
dtype: float64	
Time Series is not stationary	

^{**}Note - ** if p-value is greater than 0.05 reject null hypothesis.

Here, p-value is **0.914701** which is greater than **0.05**. Hence it is confirmed that our time series is not stationary.

To get stationary time series we need to remove trend and seasonality.

Methods to detrending data:

A time series can be detrended using following methods -

- Differencing
- Regression
- using functions
- 1. <u>Method to detrening timeseries using Differencing</u>: Differencing is the process of taking difference between successive occurrence of time series $\Delta x_t = x_t x_{t-1}$.

Where, Δx_t is stationary time series.

x_t is original time series.

 x_{t-1} is time series with lag 1.

In python we have .shift() method to create a series with lag.

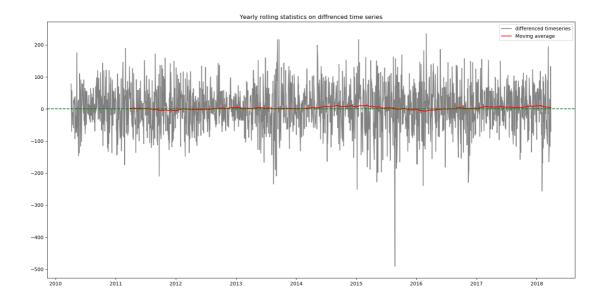
```
diff=Nifty_data['Close']-Nifty_data['Close'].shift(1)
```

There will be null values because of lag. It is important to remove null values otherwise adfuller test function will show an error i.e. "SVD did not converge"

```
diff.dropna(inplace=True)
```

lets plot rolling statistics

```
plt.plot(diff,label='differenced timeseries',color='grey')
plt.plot(diff.rolling(window=252).mean(),label='Moving average',color='red')
plt.title('weekly rolling statistics on diffrenced time series')
plt.axhline(y=0,color='green')
plt.legend()
plt.show()
```



Now, the moving average is constant.

We can confirm the stationarity by apply Augmented DickeyFullerTest on diff time series -

```
stnry_test_diff=adfuller(diff,autolag='AIC')
stnry_rslt_diff = pd.Series(stnry_test_diff[0:4], index=['Test Statistic','p-
value','#Lags Used','Number of Observations Used'])
for key,value in stnry_test[4].items():
    stnry_rslt_diff['Critical Value (%s)'%key] = value
print(stnry_rslt_diff)
if(stnry_test_diff[1]>0.05):
```

```
print("Time Series is not stationary")
else:
   print("Time series is stationary")
```

Test Statistic -41.047846 p-value 0.000000 #Lags Used 0.000000 Number of Observations Used 1985.000000 Critical Value (1%) -3.433649 Critical Value (5%) -2.862997 Critical Value (10%) -2.567546 dtype: float64 Time series is stationary

p-value is 0.000000 which is less than 0.05.

So, the time series is now stationary.

The method of differencing discussed above is first order differencing. After applying first order differencing It is possible that the time series may not become stationary then we have to perform second order differencing. To perform second second order differencing, take difference another time.

```
x'_{t} - x'_{t-1} = (x_{t} - x_{t-1}) - (x_{t-1} - x_{t-2}) = x_{t} - 2x_{t-1} - x_{t-2}
```

This was the first method of making a time series stationary.

2. <u>Method to detrening timeseries using Regression</u>: Regression is another way to detrending timeseries data. Objective of choosing regression over differencing is to get trend line. Trend line can later be used as prediction of long run movement of time series.

Regression analysis is a form of predictive modeling technique which looks into the relationship between a target and predictor variables. Regression analysis is used for forecasting, time series modeling.

Here we fit a trend line to the training data, in such a manner that the distance between data points and trend line is minimum.

lets perform linear regression, but first we have to seperate training data and testing data

```
# separating training and testing data
train = Nifty_data['Close'].iloc[:1750]
test = Nifty_data['Close'].iloc[1751:]
# building linear model
import numpy as np
from sklearn import linear_model
lm = linear_model.LinearRegression(normalize=True, fit_intercept=True)
```

Details of parameters:

1. **Normalize**: This parameter is ignored when **fit_intercept** is set to False. If True, the regressors X will be normalized before regression. By default this is set to False.

2. **fit_intercept**: whether to calculate the intercept for this model. If set to False, no intercept will be used in calculations. By default this is set to True.

For more informatin about parameters you can see LinearRegression documentation.

http://scikit-

 $learn. org/stable/modules/generated/sklearn. linear_model. LinearRegression. html \#sklearn. linear_model. LinearRegression + tml \#sklearn. linearRegression + tml \#sklearn. linearRegression + tml \#sklearn. linearRegression + tml \#sklearn. linearRegression + tml Regression + tml R$

Now we have to fit linear_model to Nifty_data timeseries

```
lm.fit(np.arange(np.array(len(train.index))).reshape((-1,1)), train)
```

Now lets check the accuracy of fitted model

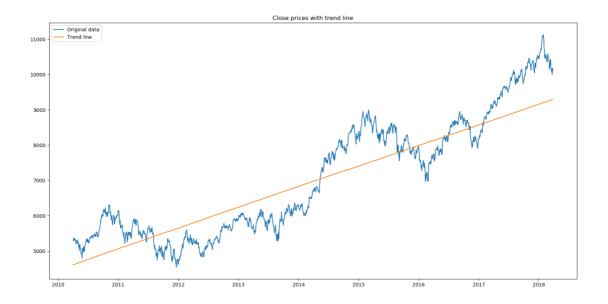
```
lm.score(np.arange(np.array(len(train.index))).reshape((-1,1)), train)
```

```
In [110]: lm.score(np.arange(np.array(len(train.index))).reshape((-1,1)), train)
Out[110]: 0.7783615453922795
```

Now, we are going to plot close prices with trend line

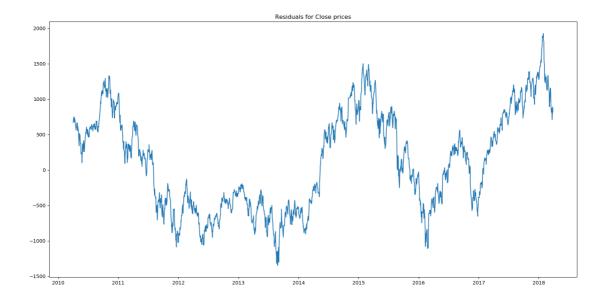
```
plt.plot(Nifty_data['Close'],label='Original data')
plt.plot(pd.Series(lm.predict(np.arange(np.array(len(Nifty_data.index))).reshape((
-1,1))),index=Nifty_data.index),label='Trend line')
plt.title('Close prices with trend line')
plt.legend()
plt.show()
```

Output will be -



Now we are going to calculate residuals which will be used later in forecasting

```
Residuals=Nifty_data['Close']-
lm.predict(np.arange(np.array(len(Nifty_data.index))).reshape((-1,1)))
plt.plot(pd.Series(Residuals,index=Nifty_data.index),label=Residuals)
plt.title('Residuals for Close prices')
plt.show()
```



3. <u>Method to detrending timeseries using Functions</u>: There are many libraries in python which are specially made to perform Statistical operations, one of them is <u>statsmodels</u>.

```
from statsmodel import seasonal
decompose = seasonal.seasonal_decompose(Nifty_data['Close'],freq=252)
```

seasonal_decompose() method returns,

- Observed data (i.e. Close prices from Nifty_data dataset)
- Trend components
- Seasonal components
- Residauls

We can access Trend, seasonal components and Residuals by decompose.trend, decompose.seasonal, decompose.resid respectively.

Note: Because we use 252 as frequency so we will have 252 null values in trend component and in residuals

```
decompose.trend.dropna(inplace=True)
decompose.resid.dropna(inplace=True)
```

Now, we are going to plot trend component.

```
decompose.trend,plot()
```

