

## Stationarity.

A time series is stationary when its statistical behavior doesn't change over time. Concretely that usually means the mean, variance, and auto-covariance stay same across time. If this properties drift, change volatility or have seasonal patterns, the series is non-stationary. Stationarity matters because many models assume that the relationship they learn is stable if the underlying distribution shifts. Models can fail or be biased.

### types of Stationarity:

- ⇒ Strict: the joint distribution of any set of time points is invariant under time shifts.
- ⇒ Weak Stationarity: mean is constant, variance is finite and constant, and co-covariance depend only on lag not absolute time.
- ⇒ Raw stock prices are typically non-stationary

### Estimate trend ( $T_t$ )

$$\hat{T}_t = \frac{1}{m} \sum_{j=K}^K y_{t+j}$$

where m is the window size often equal to

### Estimate Seasonality ( $s_t$ )

\* Remove the trend to isolate the seasonal pattern

Additive  $y_t - T_t$  • Multiplicative  $y_t / T_t$

Group by season (month, quarter) and average to estimate seasonal

### Estimate Residuals ( $r_t$ ):

$$r_t = y_t - T_t - s_t \quad R_t = \frac{y_t}{T_t \times s_t}$$

Descriptive: ex EDA

Predictive: what likely to happen

Prescriptive: what we need to do achieve the Predictive.

Forecasting: It about time

And Predictive may (or) may not about time.

Data types in Time Series:

1. Cross Sectional data: Data is collected at single point in time on one (or) more variables. Here data not sequential and usually, data points are independent of one another. Regression, random forest, neural network methods have been applied widely on these data.

2. Time Series data: univariate (or) multivariate data is observed across time in a sequential manner at pre determined and equal-spaced time interval ex yearly monthly ... ordering among data points is important and cannot be destroyed.

3. Combination of cross sectional & time series data:

This is complex data where information on same variables are collected over various points of time.

## Time Series:

A collection of observations that has been observed at regular time intervals for a certain variables over given duration is called Time Series.

\* All observations are dependent: In time series data, each observation is expected to depend on the past data observations.

\* Missing data must be imputed: Because all the data points are sequential in time series. If any data point is missing it must be imputed before the actual analysis process commences, otherwise proper ordering is not preserved.

\* Two different types of intervals cannot be mixed.

Time Series data is observed on the same variables over a given period of time with fixed and regular time intervals. Though data can be collected at various intervals such as yearly, weekly, hourly etc any specific time interval, The time-interval must remain same throughout out the entire range.

\* Objective:

Time Series forecasting is applied to extract information from historical series and is used to predict future behaviors of the same series based on past pattern.

\* Approaches used for Time Series:

1. Decomposition: This method is based on extraction of individual components of time series.

2. Regression: This method is based on regression on past observations;

Decomposition: what part of the past projects to the future and what part does not.

The part that can be forecast which is systematic and not an error. that part which is not systematic (or) random error we are going to ignore.

$$y = a + bx + (\varepsilon_0) \rightarrow \text{residual}$$

In forecast we ignore the residual  
Because the residual is applied to current data set  
not the future data set.

- ① Part that is predictable
- ② Part that is not predictable. (Ignore)

\* what we are going to project (i) Predict the future  
there is nothing to do with the past data.

\* If we are trying to get any information about the  
any situation so we need get information from similar  
situations already happened.

\* If nothing like this exist then ?  
Here we make hypothesis

we guess this is what is usually happen, if not this  
happen we keep the safety margins around it to  
~~protect~~ predict which is confidence intervals suppose to  
do.

what is happening long term period and short term  
period.

$$y = f(x)$$

$$y_{t+1} = f(y_t, y_{t-1}, y_{t-2}) \rightarrow \text{test situation}$$

future as a function of the past.

$$y_t = f(y_{t-1}, y_{t-2}, y_{t-3}, \dots) \rightarrow \text{train situation.}$$

Supervised algo.  $\rightarrow$  time Series forecast algo

Past observations as inputs to the supervised algo. (neural network).

we create the Regression on the present and past to forecast the future

$$y_t = \text{trend} + \text{Seasonality} + \text{Residual}$$

$$y_t = \text{trend} \times \text{Seasonality} \times \text{Residual}$$

trend  $\Rightarrow t$ , long term upward/downward movement

$\Rightarrow$  rising stock price over years.

Seasonality  $S_t \Rightarrow$  Regular pattern repeating over time  
 $\Rightarrow$  monthly energy usage

cyclical  $C_t \Rightarrow$  fluctuations due to business/economic cycles  
 Boom - bust cycles.

Residual/irregular  $\Rightarrow$  Random, unexplained variation  
 Noise, errors, sudden shocks.

### a) additive model:

Used when the magnitude of fluctuations is constant over time.

e.g.: Temperature data, Sales without exponential growth

$$y_t = T_t + S_t + R_t$$

### b) Multiplicative Model.

Used when seasonal fluctuations increase (a) decrease proportionally with the level of the series.

$$y_t = T_t \times S_t \times R_t$$

e.g.: Product sales where peaks growth

**Trend:** when the series increased (a) decreased over the entire length of time. the price of a share may increase (a) decrease linearly over a period of time.

**Seasonality:** when a series is observed with more frequency than a year, the series are subjected to rhythmic fluctuations which are stable and repeatable in each year.

Normal dataset It's called the ETS  
for the time series data it's called as the data decomposition.

where it gives the full picture of the trend & seasonality and noise.

Methods: STL, X-11/X-13 ARIMA, Wavelet, EMD, VMD, Fourier

## Stationary:

→ No Long term trend & no Seasonality

→ Constant mean and constant variance

How to transform the data to - stationary.

1. Logarithmic transform (impacts variance)
2. Differencing (impacts mean)

### ADF test: (Augmented Dickey-Fuller) test

Determines if a time Series is stationary or not.

Null hypothesis  $H_0$  = non stationary

Alternative hypothesis  $H_a$  = stationary.

Significance level  $\alpha = 0.05$

If P-value  $< \alpha = 0.05$  suggests rejecting the null hypothesis, meaning the series is ~~not~~ stationary.

### KPSS Test: (Kwiatkowski-Phillips-Schmidt-Shin) test

Determines if time Series is Stationary.

Null hypothesis  $H_0$  = stationary,  $H_1$  = Non stationary.

$\alpha = 0.05$

P-value  $< 0.05$  suggests rejecting the null hypothesis meaning the series is non-stationary.

## ACF: Autocorrelation function

Displays the correlation between the time series and its own past values (at lags).

X-axis represents lag (time difference)

Y-axis represents correlation coefficient (-1 to 1)

### How to interpret:

Gradual Decay: indicates non-stationary (trend is present in the data).

Sharp cutoff: indicates stationarity and potential use of an autoregressive (AR) model.

Repeated peaks: suggests seasonality in the data

Eg: correlation repeats every month

## PACF: Partial auto correlation function:

Measures the correlation between the series and its lag while controlling for the influence of intermediate lags.

A slowly decaying PACF plot may indicate the presence of the seasonality in the data.

If the PACF does not show a clear cut-off it might suggest that an AR model may not be a good fit or data may benefit from additional pre-processing or transformation.

Covariance always measures the relationship between two variables.

for same variable: Then we need to do auto-covariance which measures how the same variable relates to itself at different time lags.

data points: [100, 102, 101, 103, 105, 107, 106, 108, 110, 111]

Autocovariance of lag k for points  $x_1, x_2, x_3, \dots, x_m$

$$r_k = \frac{1}{m-k} \sum_{t=k+1}^m (x_t - \bar{x})(x_{t-k} - \bar{x})$$

$m$  = total no. of data points

$\bar{x}$  = mean of the series.

$x_t$  = current value

$x_{t-k}$  = value  $k$  steps before.

$$\text{Lag}_1 \quad r_1 = \frac{1}{9} \sum_{t=2}^{10} (x_t - 105.3)(x_{t-1} - 105.3)$$

$$\text{Lag}_1 = 9.36 \quad \text{Lag}_2 = 6.35 \quad \text{Lag}_3 = 2.20$$

Lag<sub>1</sub> is high because the each day's price is strongly related to the previous day's.

Further days (2, 3, ...) become weaker because the connection fades older day's influence the present less.

we can capture the direct and indirect co-variance of the previous days.

ACF

$$P_k = \frac{\text{cov}(x_t, x_{t-k})}{\text{Var}(x_t)}$$

Here ACF at lag  $k$  tells us how much does today's value  $x_t$  move together with the value  $k$  days ago  $x_{t-k}$

It tells the overall linear relationship it includes both

Direct effect of  $x_{t-k}$  on  $x_t$   
indirect effect of all the values between.

$x_t$  is strongly related to  $x_{t-1}$

$x_{t-1}$  is strongly related to  $x_{t-2}$

Then even if  $x_t$  never talks directly to  $x_{t-2}$  they will still show correlation because of the indirect path

$$x_t \leftarrow x_{t-1} \leftarrow x_{t-2}$$

PACF: Removing the indirect paths.

PACF measures the direct correlation between  $x_t$  and  $x_{t-k}$  after removing the effects of all intermediate lags

Mathematically it's the correlation between the residuals of two regressions:

Regress  $x_t$  on  $x_{t-1}, x_{t-2}, \dots, x_{t-(k-1)}$

Regress  $x_{t-k}$  on the same set of intermediate terms.