

Time Series

In time series based problem we are provided with time-based features that means o/p feature or forecasting to be done in time series based problem statement is dependent on time based I/P features.
For Eg; Day-wise Sales of a product.

Non time Series

Eg: Predicting price of house based in non time dependent features such as Size, Location and No. of bedrooms

House Price Prediction

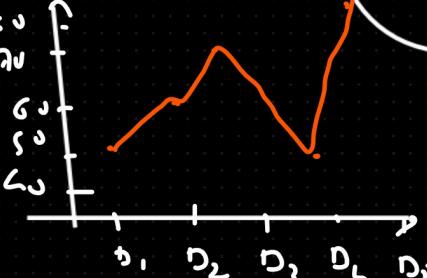


Regression

$$Y = mx + c \rightarrow \text{linear eq}$$

$$Y = m_1x_1 + m_2x_2 + m_3x_3 + \dots + m_nx_n + c$$

~~Time~~



Graphical representation of time series based question can be termed as time based forecasting.

time Series

For Eg; Day-wise Sales of a product.

Time based feature can have any of the below timestamp format

hour
min
sec
Day
month
year

timestamp

Day 1	501c
Day 2	601c
Day 3	701c
Day 4	801c
Day 5	901c

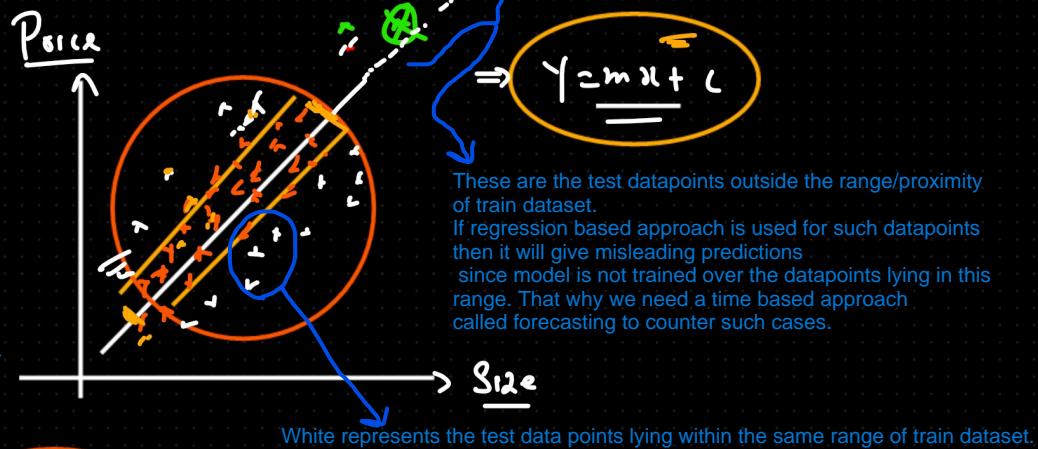


Ques Can we solve this TS with linear regression?

To answer this we need to understand Interpolation and Extrapolation concepts

Interpolation \Rightarrow to find out the value in the range itself.

Best fit line in this case is infinite line. Training dataset captures only a portion of this infinite best fit line. In interpolation Test data points or new datapoints on which prediction is supposed to be made lie within the same range of training datapoint. Whereas, In Extrapolation these data points will lie outside the proximity of the train dataset often denoted by different time frame.



Extrapolation \Rightarrow to find out the value out of Range.



Forecasting

TimeSeries data

extrapolation

Prediction and Forecasting may look similar but are different considering the range of Test dataset. If Test data set within the same time frame as train dataset then Interpolation else Extrapolation.

Since Interpolation consider same timeframe so the features are independent of time. Which in case of Extrapolation is not the case hence features in it is time dependent.

time series problem statement will be extrapolation

- Based on Previous history forecast the future value.

Ques Can we solve this TS with linear regression?

1 Because of extrapolation. we can not do.

2 Because of outliers we can not use.

Time Series

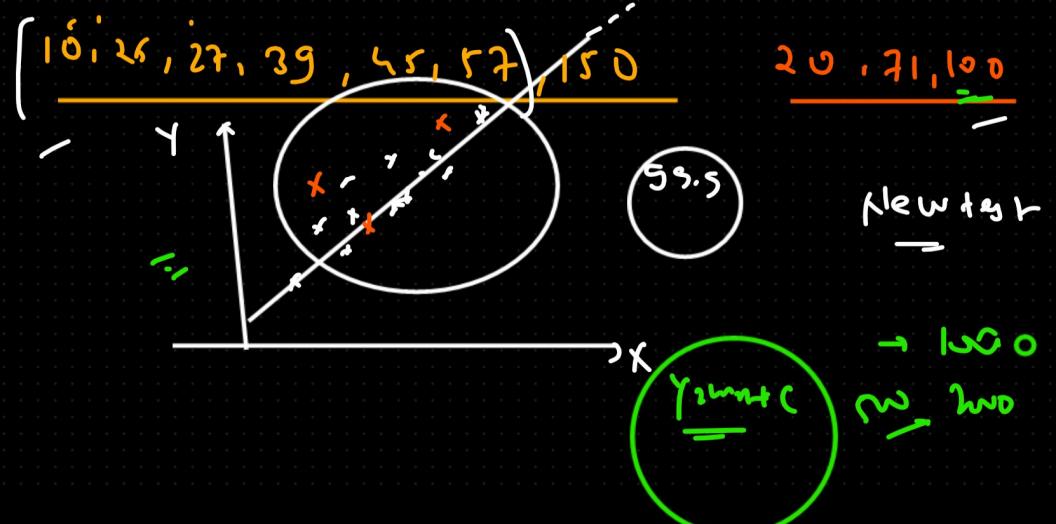
2 TS Data is more complex

3 Linear reg. Assume that ^{there} should linear relationship but in time you won't get linear relationship

5 In Non-TS Data there is no such effect of time but ^{time} in TS - Data there is effect of time meaning current TS is dependent on prev. time stamp.

$$\text{Data} \Rightarrow [10, 20, 25, 27, 39, 45, 57, 71, 100, 150]$$

$$[10, 25, 27, 39, 45, 57, 150] \quad 20, 71, 100$$



Example = ① economics forecasting \Rightarrow GDP, interest rate, inflation rate

= ② finance \rightarrow Sales, Bond Price

= ③ weather forecasting \Rightarrow Pattern of Diff season, Prediction of weather

④ Medical \Rightarrow Based on prev. history of patient need to predict future condition.

↑

Time Series \Rightarrow time dependent data \Rightarrow any domain

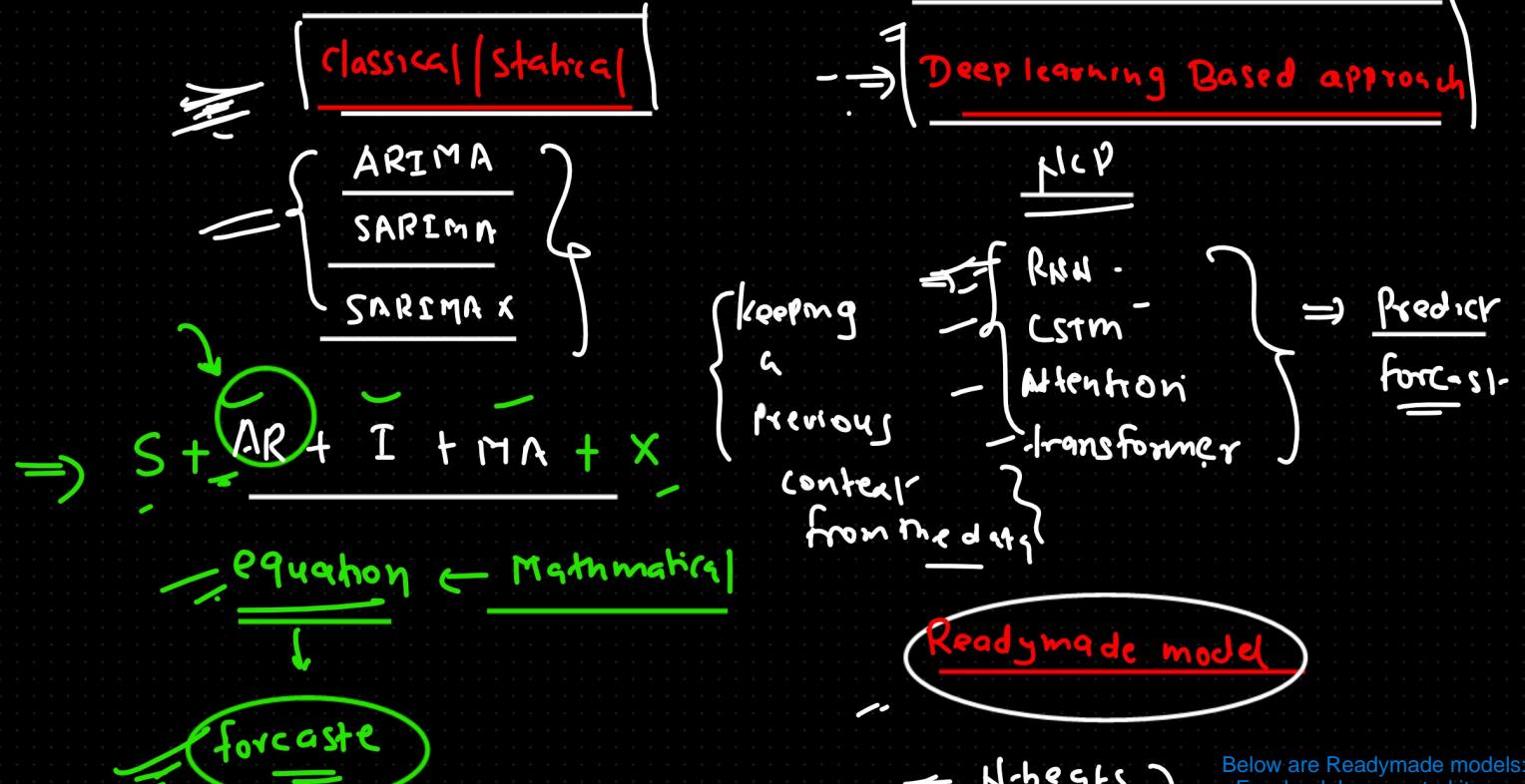
Dq+q → Time dependent Data

EDA → —

Preprocessing → —

1 model → —

Model evaluation → —



Lets, now discuss some of the common terminologies to be used while working with a time series based problem statement:

- Trend
- Season
- Noise
- Cycle

Below are Readymade models:

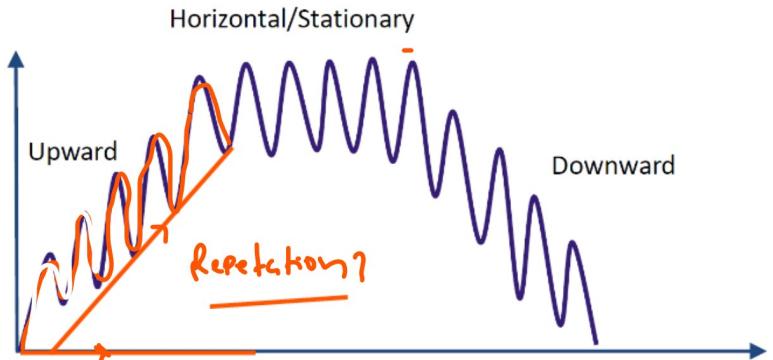
- Facebook has created its own time series model called Prophet.
- Similarly, Google has also created it's own time series based model called Temporal Fusion transformer.
- N-bets can also be used for time series based problem statement.

Google these to know more.

Timeseries =

- ① trend
- ② Season
- ③ Noise
- ④ Cycle

• Trends



Trend = ① Upward trend

Trend tells us that how the data is behaving with the time. It has 3 types:
Upward, Downward and horizontal or stationary.

- ② Downward trend
- ③ flat (horizontal)

Season = frequent (repetition) (Daily, monthly, yearly)

Seasonality is a characteristic of a time series in which the data experiences regular and predictable changes that recur every calendar year.

- ① Sales of the ice-cream in summer
- ② traffic situation at 6 pm in my area
- ③ houses in goa at the end of year (Dec, Jan)

Noise =

Some uncertainty or some randomness in my data. bcoz of unpredictable reason.
our data because

UNPREDICTABLE REASONS

Eg: Ronaldo removed Pepsi bottle during interview that affected stock price of Pepsi company.

Eg: Stock market crashes at time of war, pandemic, publishing report that defames the company.

Pandemic \Rightarrow corona

War \Rightarrow R|U

Report \Rightarrow Hindenburg report



Cycle

time Series behavior over the long Period of time

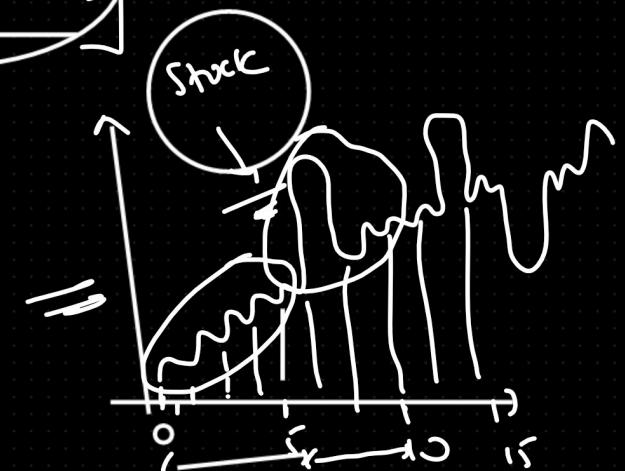
Cycle \Rightarrow seasonality + Noise (fluctuation)

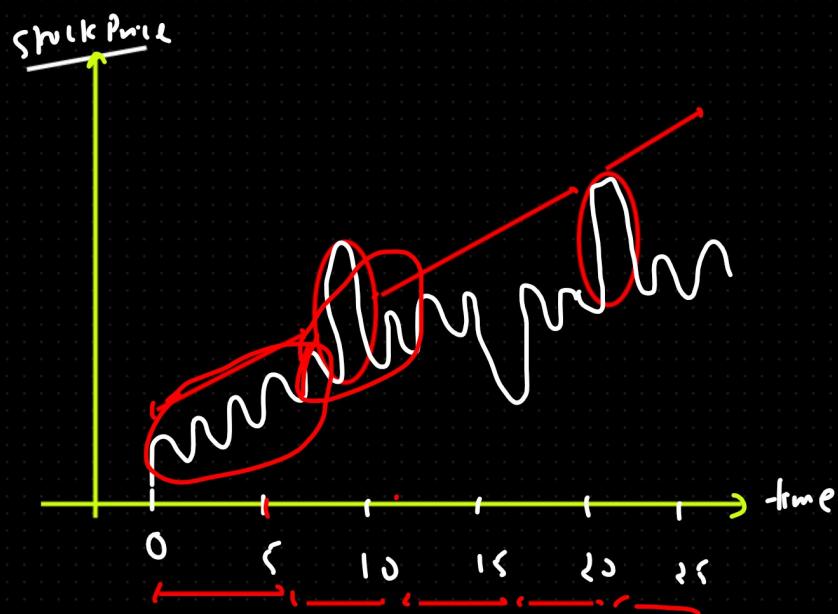
Cycle = Season + Noise
When GDP report is published after every 1 year (Season) then based on the sentiment of the report, there is sudden fluctuation (Noise) on the price of stocks listed on stock market.

Election \Rightarrow 5-year \Rightarrow cycle

Census Counting \Rightarrow 10 years

Economics \Rightarrow GDP





Two components of timeseries:

- Multiplicative
- Additive

Component 1 - \Rightarrow Multiplicative \Rightarrow

Applied when there is

- 1 Non-linear Pattern
- 2 Non-constant variation

$$y_t = T \times S \times N$$

50K
60K
70K
80K
90K
100K

$$y_t = T + S + N$$

Applied when there is

- 1 linear Pattern
- 2 If it is having constant variation

Moving Average

- ① Simple moving average (SMA)
- ② Cumulative moving average (CMA)
- ③ EMA or EWMA

Why we calculate MA?

SMOOTHING OF TIME SERIES DATA

SMA =>

Simple moving average



of

$$\left[\frac{10, 12, 15, 13, 11}{\downarrow} \right]$$

$$\frac{10 + 12 + 15 + 13 + 11}{5}$$

$$= 12.5 \text{ APPROX}$$

Moving

move over the time axis
in given window size



Moving + average

window size

calculating average

Move over the time axis in the given window size followed by calculating the average of the corresponding time dependent output feature for the respective windows.

1

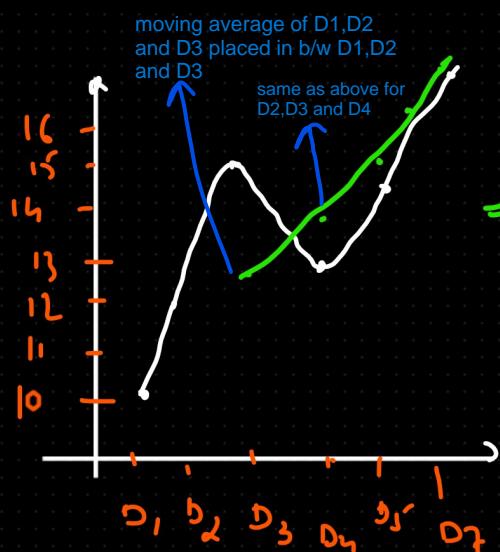
Moving average

Lets understand the simple moving average using the below example:

Window
Size = 3

representing 1st average upto D₃. Above since there out not in the window size so the respective average value is replaced by NAN

D ₁	10	NAN
D ₂	12	NAN
-	-	-
D ₃	15	→ 12.33
D ₄	13	→ 13.33
D ₅	14	→ 14
D ₆	16	→ 14.33
D ₇	17	→ 15.66



In above plot white line represents general time series plot whereas, green line represents the smoothen time series curve obtained using SMA.

$$1^{\text{st}} \text{ avg} = \frac{D_1 + D_2 + D_3}{3} = \frac{10 + 12 + 15}{3} = 12.33$$

$$2^{\text{nd}} \text{ avg} = \frac{D_2 + D_3 + D_4}{3} = \frac{12 + 15 + 13}{3} = 13.33$$

$$3^{\text{rd}} \text{ avg} = \frac{D_3 + D_4 + D_5}{3} = \frac{15 + 13 + 14}{3} = 14$$

$$4^{\text{th}} = \frac{D_4 + D_5 + D_6}{3} = \frac{13 + 14 + 16}{3} = 14.33$$

$$5^{\text{th}} = \frac{D_5 + D_6 + D_7}{3} = \frac{14 + 16 + 17}{3} = 15.66$$

Why we do
Smoothing?
TIME

- 1 to remove all the effect from the data
- 2 Pattern recognition from the data
(to summarize the overall trend of the data)
- 3 You can easily analysis trend of Data
- 4 You can reduce the effect of outlier
- 5 You can easily visualizing the data.

CMA \Rightarrow Find out the avg of all the Data Points up to the given time stamp

$$D_1 = 10 \Rightarrow D_1 = 10$$

$$D_2 = 12 \Rightarrow \frac{D_1 + D_2}{2} = \frac{10 + 12}{2} =$$

$$D_3 = 15 \Rightarrow \frac{D_1 + D_2 + D_3}{3}$$

$$D_4 = 14 \Rightarrow \frac{D_1 + D_2 + D_3 + D_4}{4}$$

$$D_5 = 16 \Rightarrow \frac{D_1 + \dots + D_5}{5}$$

$$D_6 = 17 \Rightarrow \frac{D_1 + D_2 + \dots + D_6}{6}$$

$$D_7 = 18 \Rightarrow \frac{D_1 + D_2 + D_3 + \dots + D_7}{7}$$

Ans

① It will give you the exponential trend.

② We use CMA for the long time period.

based time series problem statement.

For nth datapoint CMA will be this =

$$\frac{D_1 + D_2 + D_3 + \dots + D_n}{n}$$

Now with the values of respective cumulative averages we can plot the smoothen curve of time series based plot.

EMA OR EWMA \Rightarrow in EMA we give more weightage to the recent Data Point

(Exponential moving average)

Or give More weightage to the recent time stamp.

Please note that SMA and CMA are the non weighted moving averages. Whereas, EMA is the weighted moving average technique.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

V_t = EMA at time t

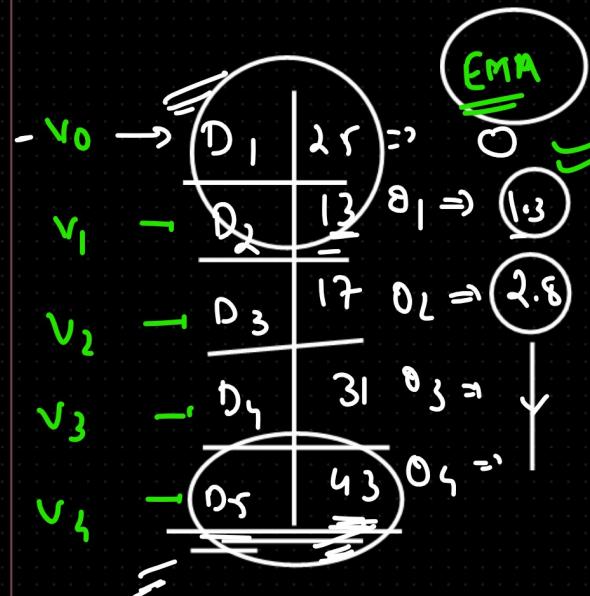
$$\Rightarrow \beta = 0 < \beta < 1 \Rightarrow 0.9$$

Beta is the parameter whose value range b/w 0 and 1. In most of the cases it is taken as 0.9

V_{t-1} = EMA at previous time stamp

θ_t = Data value at t time stamp

* V_0 will be 0 or 25 based on the research paper documenting this approach



$$V_0 = 0$$

$$V_0 = 25$$

$$V_1 = \beta \times V_0 + (1-\beta) \theta_1$$

$$= 0.9 \times 0 + (1-0.9) \times 13$$

$$= 0 + 0.1 \times 13 = 1.3$$

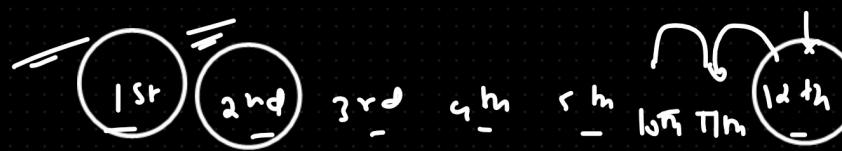
$$V_2 = \beta \times V_1 + (1-\beta) \theta_2$$

$$= 0.9 \times 1.3 + (1-0.9) \times 17$$

$$= 1.17 + (0.1 \times 17)$$

$$= 1.17 + 1.7$$

$$= 2.8$$



Now with the values of V_0, V_1, V_2, \dots so on we can plot the smoothen curve of time series based plot.

EMA

P_1 15k
20k
25k
30k
35k

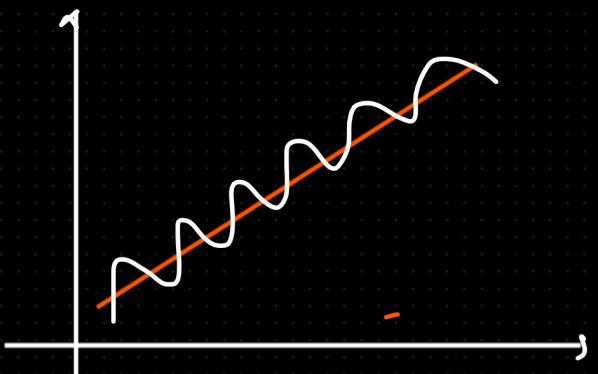
1st 2nd

10m 11m 12m

DL

Gradient
Chamrile

Stationary + Non-stationary

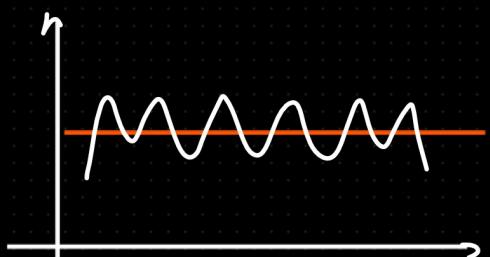


trend = upward

MA (SMA)
window + average

\Rightarrow my moving is not constant over the time.

\Rightarrow Variance is not constant over the time.

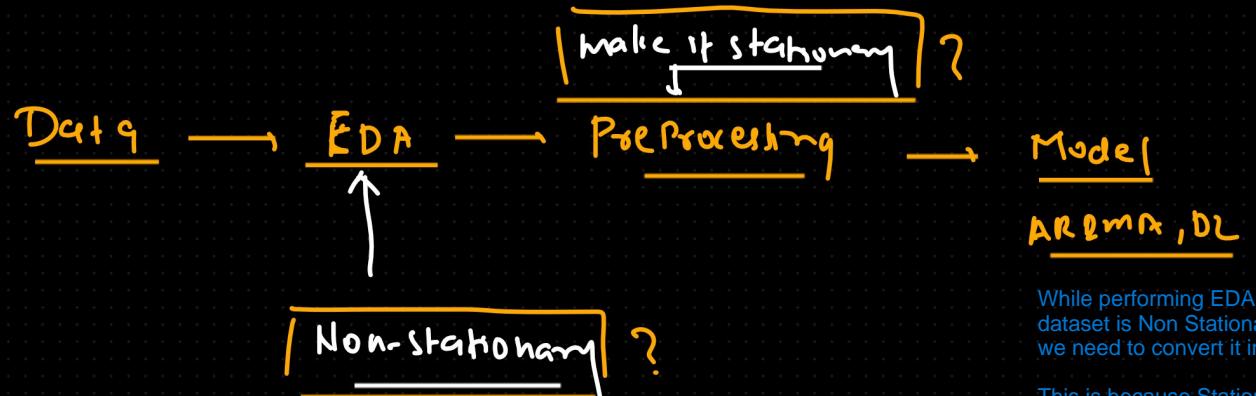


trend is horizontal

Non-stationary = Mean and Var will not be const

Stationary \Rightarrow Mean and Var will be const

{the meaning is there is no change in
Mean, var over the time axis?}



ARIMA, DL

While performing EDA if it is found that the time series based dataset is Non Stationary then during the Pre Processing stage we need to convert it into Stationary time series based dataset.

This is because Stationary dataset will give better results over Non Stationary dataset.

- A nqly srs \Rightarrow
- 1 Visualization
 - 2 Stats based test

Ways to identify whether time series is Stationary/Non stationary

- Using Visualization techniques
- Using Stats based test

Preprocess \Rightarrow (Non-st to st) \Rightarrow

- 1 Differencing
- 2 log
- 3 root
- 4 Adjustment of seasonal Data

These are the different approaches that can be used for converting the Non-Stationary time series based dataset into Stationary time series based dataset.

ACF and PACF and ARIMA

ACF \rightarrow auto corr function

PACF \Rightarrow Partial auto corr fn

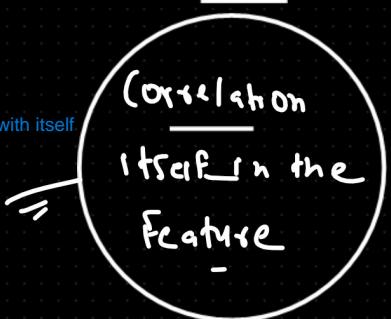
ACF \Rightarrow

Auto

+

Correlation

Auto means correlation of time dependent feature with itself



It is a relationship b/w two feature



{ Pearson
Spearman rank
Kendall}

ACF measure the corr b/w time series and its lag value

↓

1st lag

and lag

3rd lag

D ₁	10	NA	NA
D ₂	25	D ₁ - 10	NA
D ₃	14	D ₂ - 25	D ₁
D ₄	16	D ₃ - 14	D ₂
D ₅	15	D ₄ - 16	D ₃
D ₆	32	D ₅ - 15	D ₄

In 1st lag we are basically comparing or measuring the correlation b/w current and previous time series dataset. Since, we are comparing with the same column but in current to previous manner so that's why it is called auto corr fun^.

y_t ↓ y_{t-1} x_1 x_2 y_{t-2} y_{t-3}

representing 1st lag, 2nd lag, 3rd lag....
so on with $y_t, y_{t-1}, y_{t-2} \dots$ so on.

$$\left\{ \begin{array}{cc} D_2 & D_1 \\ D_3 & D_2 \\ D_4 & D_3 \\ D_5 & D_4 \\ D_6 & D_5 \end{array} \right\} = \boxed{\begin{array}{|c|c|} \hline 25 & 10 \\ \hline 14 & 25 \\ \hline 16 & 14 \\ \hline 20 & 16 \\ \hline 32 & 20 \\ \hline \end{array}}$$

\Rightarrow Corr itself in the
feature

|Autocorr with lag 1|

\Rightarrow Corr with 1st lag

$\rightarrow \text{Corr}(y_t, y_{t-1})$

\Rightarrow Corr with 2nd lag

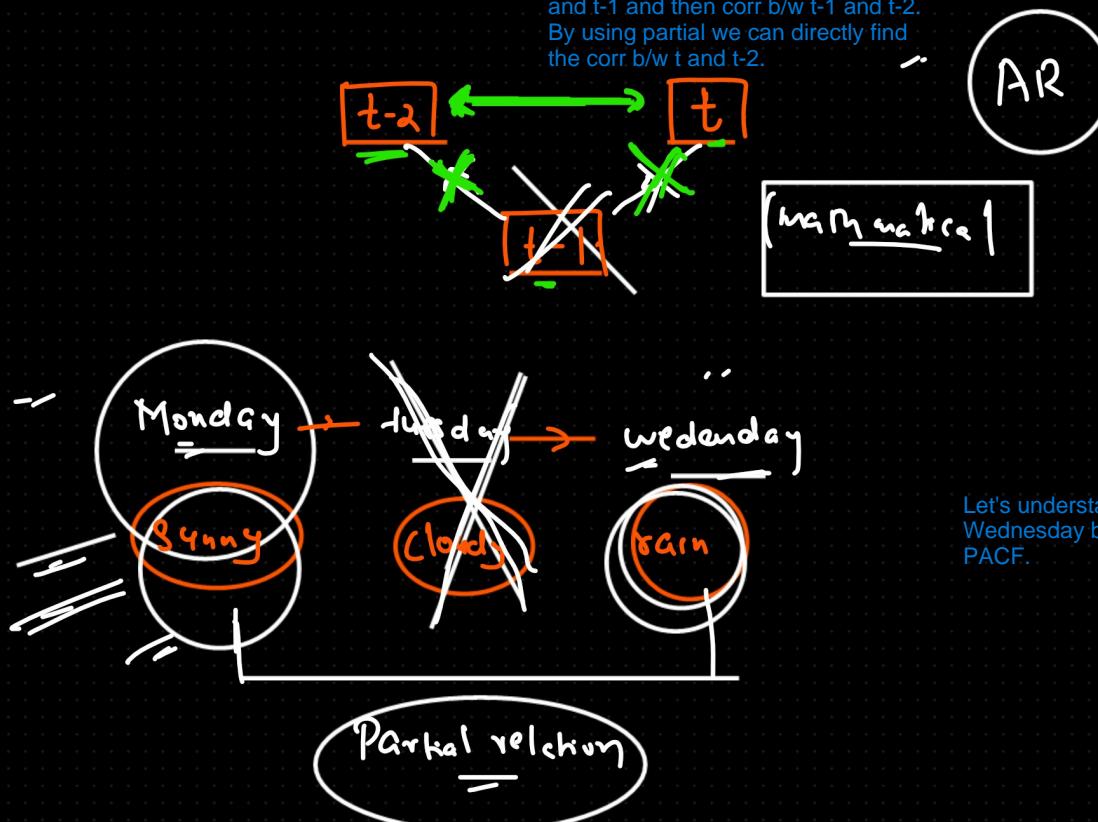
$\text{Corr}(y_t, y_{t-2})$

\Rightarrow Corr of 3rd lag

\Rightarrow Corr(y_t, y_{t-3})

PACF \Rightarrow Partial auto corr

In auto we are first finding corr b/w t and t-1 and then corr b/w t-1 and t-2.
By using partial we can directly find the corr b/w t and t-2.



Let's understand with this example. If we attempt to predict forecast of Wednesday based on Monday by skipping Tuesday then this will lie under PACF.

Difference b/w auto correlation and correlation:

Correlation is correlation b/w 2 or more different features whereas, auto correlation is the correlation of a feature with itself.

Auto regression



Auto

Regression

regression itself
in the var



X → independent feature
Y → dependent feature

$$y = mx + c$$

m = slope
c = intercept

$$y = m_1x_1 + m_2x_2 + m_3x_3 + \dots + m_nx_n + c$$

$$\phi \left\{ \begin{array}{l} y_t = \text{Value at the current time stamp} \\ \psi = \text{coeff term} \\ c = \text{constant} \\ \epsilon = \text{error} \end{array} \right\}$$

Auto regression

AR eq. with lag 1

$$y_t = \psi y_{t-1} + c$$

AR eq. with lag 2

$$y_t = \psi_1 y_{t-1} + \psi_2 y_{t-2} + c$$

AR eq. with lag n

$$y_t = \psi_1 y_{t-1} + \psi_2 y_{t-2} + \dots + \psi_n y_{t-n} + c$$

