

Time Series

Time Series → A series with time components.

Area of House	No of rooms	Price of House
-	-	-
-	-	-
-	-	-

You can use
MLR
Regression Problem statement

Area of house	Price of house
1000	60L
950	50L
1100	75L

	↓ (Rearrangement of rows)
1100	75 L
1000	60 L
950	50 L

(times by) month ↓	Sales
Jan	60k
Feb	50k
Mar	40k
-	-
-	-

Can we use MLR?

↓
No because it's
time series
problem

Time Series problem

Mouth	Sales
Jan	50k
Feb	40k
Mar	30k

X rearrangement
is not possible
↓ because

⇒ Time component is involved
⇒ Here order matters.

⇒ because current row depends on former rows
i.e. - March depends on Feb, Feb depends on Jan and so on.

* Interpolation \Rightarrow To find out the value in range itself (Prediction)

Price of house



Area of house

Most of the time (99% of the time),
the test date will come
in the training range

$$\text{training (Area)} = 1000 - 10000$$

$$\text{testing (Area)} = \text{,, (Most of time.)}$$

* If outside training range \Rightarrow wrong prediction

* Why Not Linear Regression for time series?

\rightarrow time component is involved.

\rightarrow Because of Extrapolation, it may lead to wrong Prediction.

\rightarrow LR - Assumes linear relationship but in time series,
the current observation depends on previous obs. \rightarrow which
is not true for non-time series data.

Motivation

① Weather forecasting \Rightarrow Weather, patterns day wise, month wise, seasonal

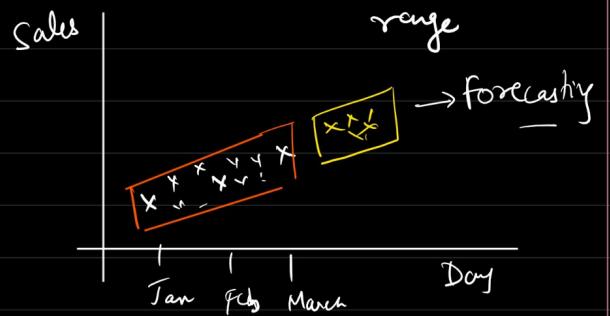
② Medical \rightarrow based on previous medical history, predict future

③ finance) — Sales, Bond price, Stock price

④ E-commerce
Economics — Interest rate, GDP etc.

Extrapolation \rightarrow

To find out the
value out of
the range



* Time series problem
statement will extrapolation
 \rightarrow Based on Previous
History, forecast
the future value.

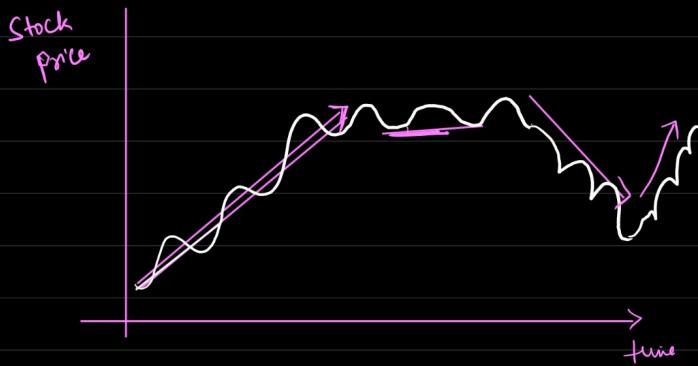
* You can use time series in every domain wherever
you have a time dependent data

Components of Time series

Analysis forecasting time series → Airline passenger.

Time Series Analysis

- ① level
- ② trend
- ③ season / seasonality
- ④ cycle / cyclicity
- ⑤ Noise



① level → the base value of a time series on which other components are added.

② trend — Long term movement or direction in date over a long period of time

① Upward .

② Downward

③ Horizontal/flat



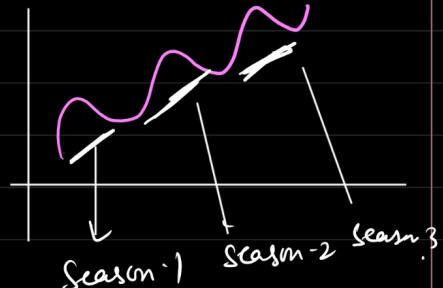
③ season :- Frequent repetition in a date over regular interval such as daily, monthly, annually

Example → Sales of ice-cream increases in summer

→ Sales of electronic gadgets in diwali.

→ traffic in peak hours

→ No of tourist visiting every year.



④ Cycle/cyclicality → The fluctuation in the date over a longer period of time. These periods are not fixed and can vary.

Cycle → Noise = (Anomalous behaviour, Unexpected behaviour)

Cyclicity - Season + Noise
over a long period of time.

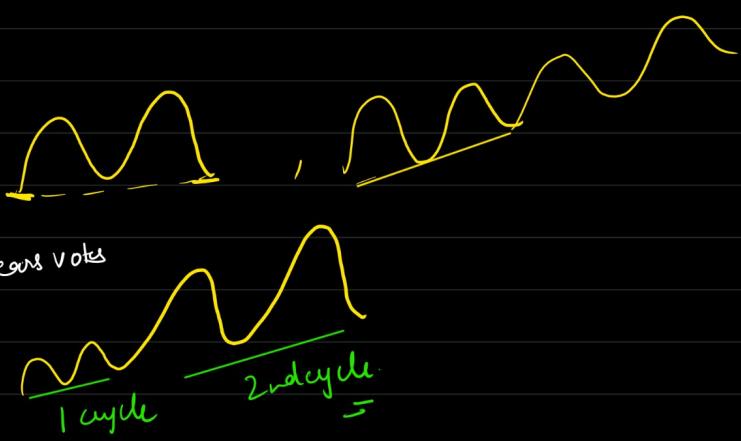
Noise/Anomaly



Ex. Stock price

Ex. Politics - every 5 years votes

Ex. Economics - GDP



⑤ Noise

Anomaly | Outlier \rightarrow Some uncertainty / randomness in time series data because of unexpected reason.



reasons

- \rightarrow News
- \rightarrow Reports
- \rightarrow Pandemic
- \rightarrow Wars
- \rightarrow Elections
- \rightarrow Influencer

+ Dogecoin always goes to move unexpectedly after Elon Musk's tweet.



* Based on Components of a time series, there are two types of time series.

Time Series

Additive

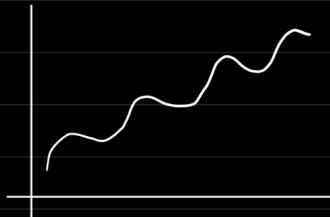
$$\rightarrow y_t = \text{Trend} + \text{Season} + \text{Noise}$$

→ Linear over time

→ Constant variance

→ Increased trend at same difference.

Day 1	100
Day 2	200
Day 3	300
Day 4	400

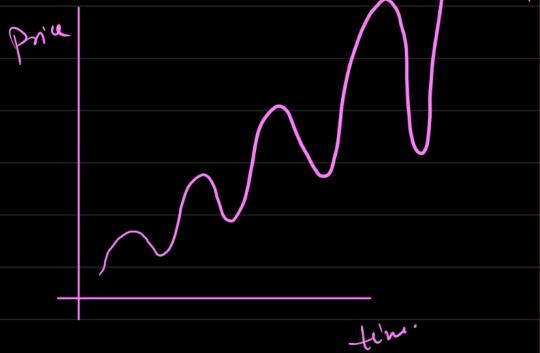


Multiplicative

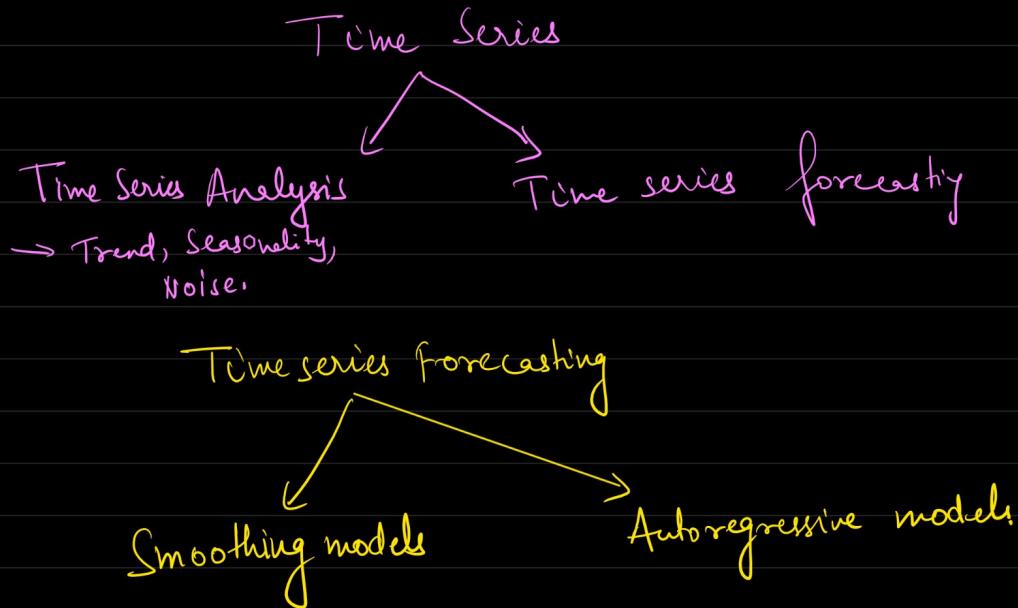
$$\rightarrow y_t = \text{Trend} \times \text{Season} \times \text{Noise}$$

→ Non-linear

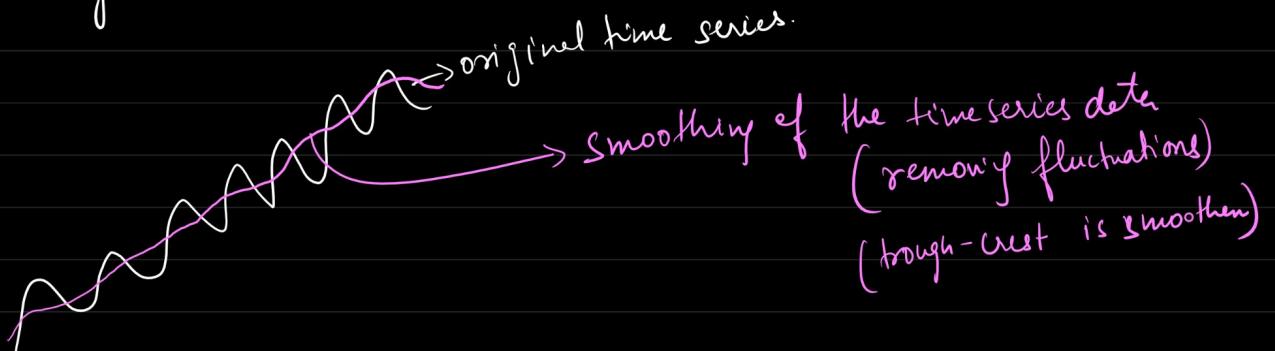
→ Non-constant variance.



Time series forecasting



* Smoothing models.



- ① Simple moving avg (SMA)
- ② Cumulative moving avg (CMA)
- ③ Exponential Weighted Moving Average (EWMA)

① Simple moving Average

$$\text{Avg} = \frac{\text{Sum of all nos}}{\text{No of values.}}$$

$$2, 3, 4, 5 \Rightarrow \frac{2+3+4+5}{4} = \frac{14}{4} = 3.5$$

Moving avg.



move over the
time axis in
a specific
window.

→ window size
→ avg value.

		SM A
① Window = 3	Jan	50
② Avg.	Feb	65
	Mar	70
	April	85
	May	90
	Jun	100
	July	110
		62
		74
		82
		-
		-

$\rightarrow \text{1st avg.} = \frac{50 + 65 + 70}{3} \Rightarrow 62$
 $\rightarrow \text{2nd avg.} = \frac{65 + 70 + 85}{3} \Rightarrow 74$
 $\rightarrow \text{3rd avg.} = \frac{70 + 85 + 90}{3} = 82$

1) window 2



→ Due to simple moving avg, the time series will smoothen.
after smoothing.

Why Smoothing

→ To remove all the effects from date

→ You want to see trend of date

→ reduce effect of outlier

→ pattern recognition from the data

→ visualisation



② Cumulative Moving Avg.

→ Find out the avg of all the data points up to the given time stamp.

		CMA
Jan	10	10
Feb	12	$\frac{10+12}{2}$
March	15	$\frac{10+12+15}{3}$
Apr	14	$\frac{(10+12+15+14)}{4}$
May	16	$()/5$
June	17	$()/6$
July	18	$()/7$

When?

→ for long time period.

→ give you exponential trend



③ EMA or EWMA

→ We give more weightage / importance / priority to the recent data point / timestamp.

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

V_t = EMA at time t

$\beta \Rightarrow 0 < \beta < 1 \Rightarrow$ generally its 0.9.

D₁ → less weight

D₂

D₃

D₄

D₅ → recent obs / more weight to be given

V_{t-1} = EMA at previous timestamp

Θ_t → Date at current timestamp t.

	<u>E_{MA}</u>
V_0	or 25 ($V_0=0$ or $V_0=25$)
V_1	1.3
V_2	2.87
V_3	-
V_4	-

$$\begin{aligned}
 V_1 &= \beta \times V_0 + (1-\beta) \Theta \\
 &= 0.9 \times 0 + (1-0.9) \times 13 \\
 V_1 &= 0 + 0.1 \times 13 \\
 &= 1.3
 \end{aligned}$$

$$\begin{aligned}
 V_2 &= \beta \times V_1 + (1-\beta) \Theta \\
 &= 0.9 \times 1.3 + 0.1 \times 17 \\
 &= 1.17 + 1.7 \\
 V_2 &= 2.87
 \end{aligned}$$

SMA or CMA

disadvantage

→ It gives equal priority to all the time period values.

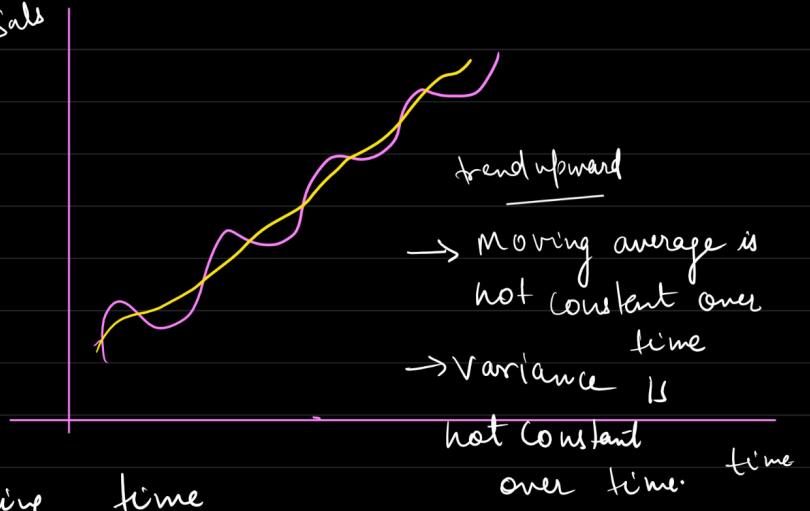
→ In time series current observation is highly influenced by recent observations.

Jan 10
Feb 15
Mar 16

∴ You need to give priority to recent observations.

Stationary and Non stationary time series

Month	Sales
Jan	40
Feb	50
Mar	60
...	...
8	80
9	90

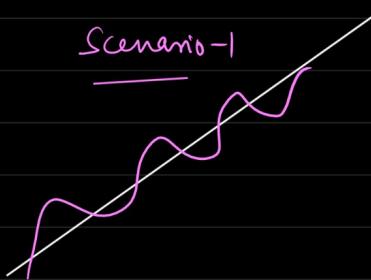


- * You can build a smoothing time series on any time series data.

Time series forecasting



- * To build autoregressive models, statistical properties of a time series like mean / variance should be constant (not change over time).



Scenario-2

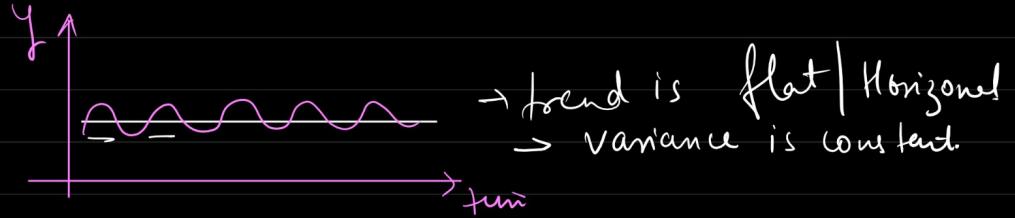


→ It will be easy to build a Autoregressive forecasting model in Scenario-2

Non stationary → Mean - variance will not be constant

Stationary → Mean - variance will be constant

Constant → Over time, value is not varying (change)



ML → Data ingestion → Analysis → Preprocessing / transformation → Model building

Time series → Time series is

Stationary (mean | variance constant) $\xrightarrow{\text{Yes}}$ Model

Non-stationary

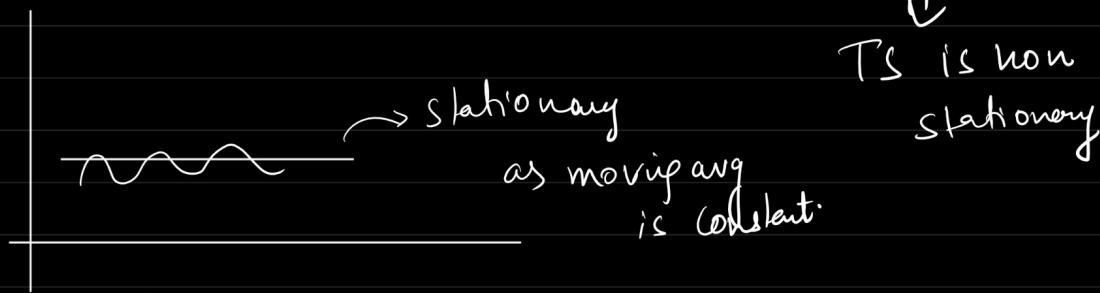
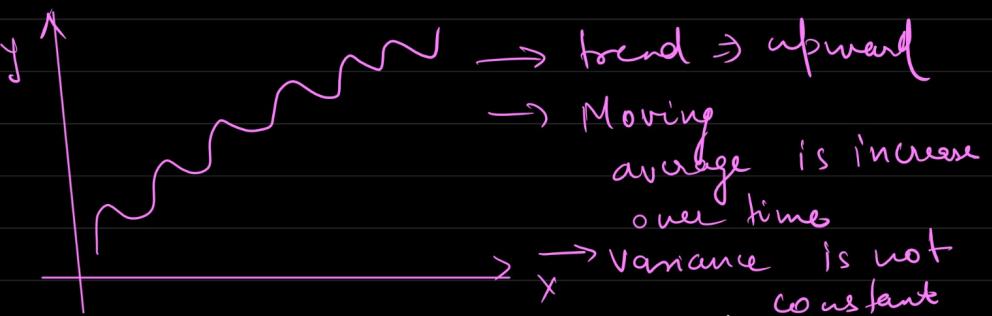
Check

if nonstationary → convert to stationary time series

(ARIMA ,
 SARIMA)

* To check if a time series data is stationary or Non-Stationary

① Visualization



② Statistical based test

① ADF (Augmented dickey fuller test)

H_0 : TS data is non stationary
 H_A : TS data is stationary

$P\text{value} \leq 0.05 \rightarrow \text{reject the } H_0$

Conclusion: TS data is stationary

② KPSS test

(Kwiatkowska -
Philips - Schmidt -)
Shin test

$H_0 \rightarrow \text{my data is stationary.}$

* How to convert non-stationary time series data to stationary time series data.

① Differencing

② log transformation → take the log of value

③ root → take the sq root of value

④ seasonal adjustment

AR \int MA

Integrated stands

$$\text{Differencing} = Y_t - Y_{t-1}$$

(current) (previous)

Month	Price	1st difference	2nd order differencing	3rd order differencing	for differencing
Jan	5	NA	NA		
Feb	10	$10 - 5 = 5$	NA		
Mar	6	$6 - 10 = -4$	$-4 - 5 = 9$		
Apr	8	$8 - 6 = 2$	$2 - (-4) = 6$		
May	15	7	$7 - 2 = 5$		
June	2	-8	-15		

check stationary

After differencing check if TS is stationary

stats test (ADF)
visualisation

→ if stationary build the model
else do differencing

again and keep doing so
until you get a stationary TS.

ACF, PACF and Auto regression

ACF \rightarrow Auto Correlation function

PACF \rightarrow Partial Auto Correlation function.

ACF \rightarrow Auto + Correlation
 ↓
Correlation ↓
 Relationship b/w two feature
 itself in the feature.

* ACF measures the correlation between time series & its lag value

Date	y
Jan	10
Feb	20
Mar	30
Apr	40

Month	y_t	1st lag	2nd lag	3rd lag.	↓ Time series data
Jan	10	NA	NA	NA	
Feb	25	10 (Jan)	NA	NA	
Mar	35	25 (Feb)	10 (Jan)	NA	
Apr	42	35 (Mar)	25 (Feb)	10 (Jan)	
May	50	42	35 (Mar)	25 (Feb)	
June	55	50	42	35	
July	62	55	50	42	

Y_t Y_{t-1} Y_{t-2} Y_{t-3}

Auto Correlation

$$\text{Corr}(Y_t, Y_{t-1})$$

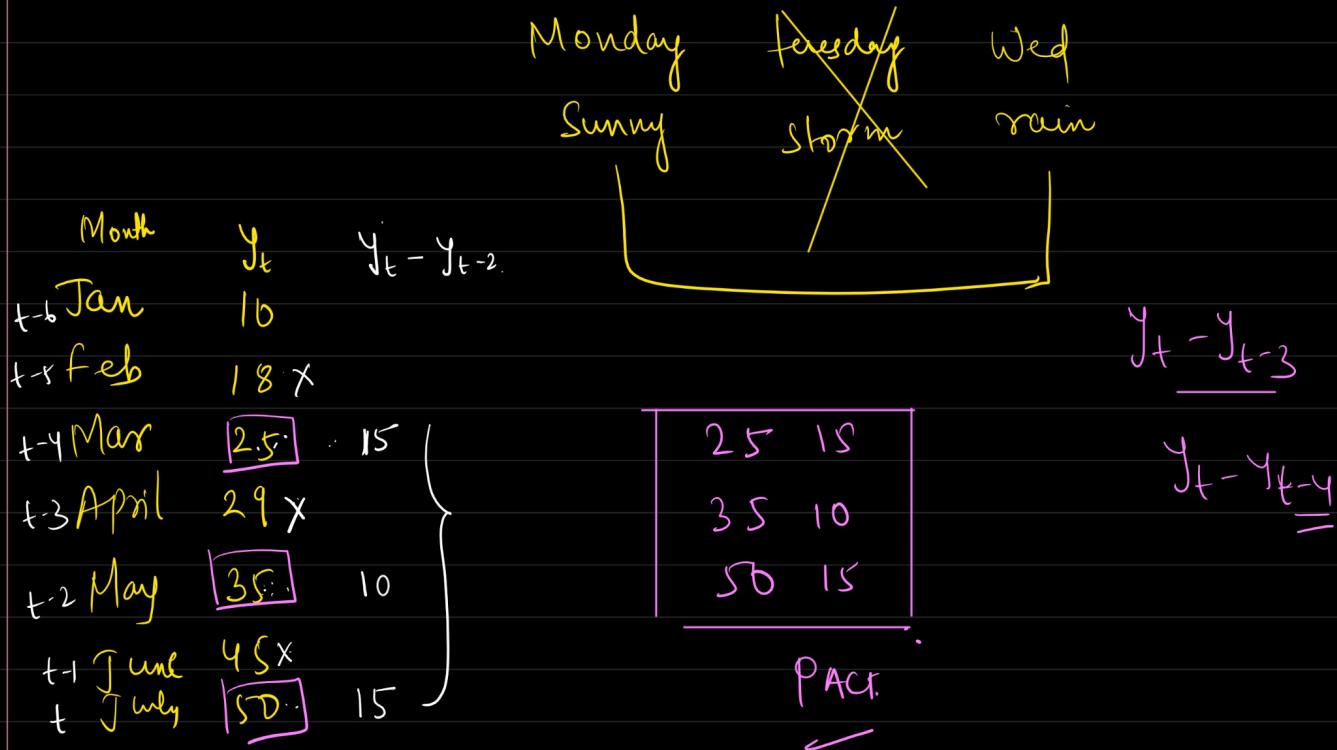
$\left\{ \begin{array}{l} \text{Pearson, } \\ \text{Spearman, } \\ \text{rank, } \\ \text{Kendall} \end{array} \right\}$

$$\text{Corr}(Y_t, Y_{t-2})$$

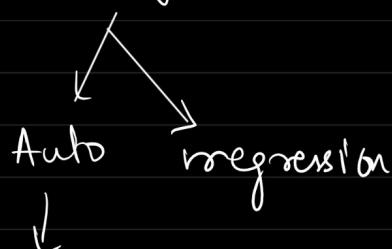
$$\text{Corr}(Y_t, Y_{t-3})$$

Partial Autocorrelation Function (PACF)

$t-2 \leftarrow \cancel{t-1} \leftarrow t$
 remove the intermediate effect



* Auto regression



regression with itself

$$Y_t = \psi Y_{t-1} + C$$

$$Y_t = \psi_1 Y_{t-1} + \psi_2 Y_{t-2} + C$$

$$Y_t = \psi_1 Y_{t-1} + \psi_2 Y_{t-2} + \psi_3 Y_{t-3} + C$$

Y_t - value at current timestamp
 ψ - coeff

$$\text{MLR} \Rightarrow Y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 - \dots - \theta_n x_n$$

$$\theta_0, \theta_1, \theta_2, \dots = \text{coeff}$$

C - constant.
 $\varepsilon \rightarrow \text{Error}$.

$$\begin{array}{c} \text{regression} \\ | \\ X \quad Y \\ | \\ \text{IV} \quad \text{DV} \end{array}$$

$$Y = mX + C$$

slope m intercept C

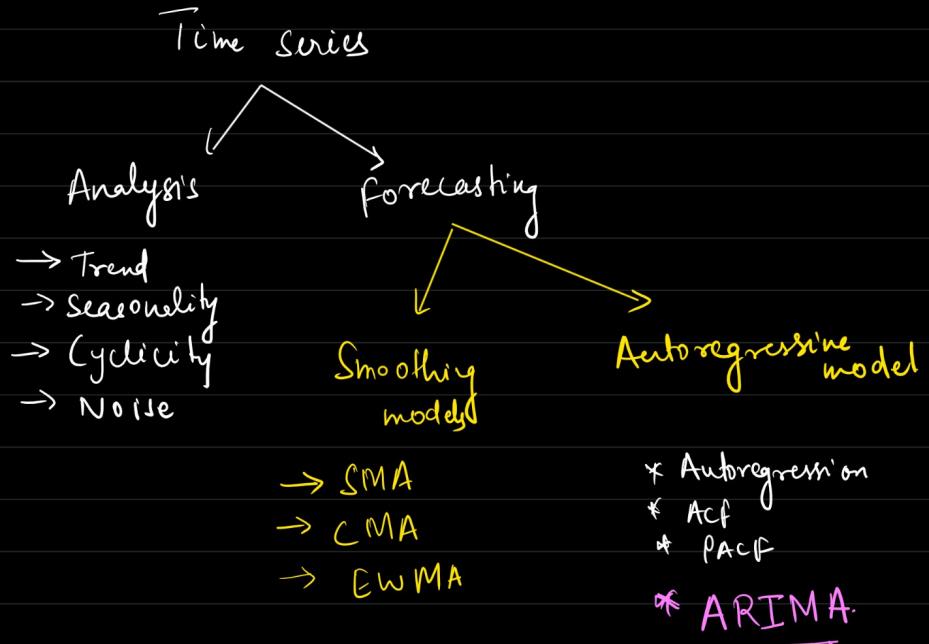
$$y_t = \psi_1 y_{t-1} + \psi_2 y_{t-2} + \psi_3 y_{t-3} + \dots + \psi_n y_{t-n} + c$$

How to decide if you have to go till
 $t-1$, or $t-2$ or $t-3$??

↓
ACF | PACF

ARIMA

ARIMA



ML → Date → LR, SVR, DTR and so on.

TS data → Forecasting models. (Smoothing models already we have seen)

- AR
- MA
- ARIMA
- SARIMA
- SARIMAX

A R

(Autoregression)



P

(0, 1, 2, ..., n)

lag value

PACF

Partial Auto-correlation Function
(Correlogram)

I

(Integrated)



d

(0, 1, 2, 3, ..., n)

lag value

Differencing
(Stationary)

M A.

(Moving average)



(0, 1, 2, 3, ..., n)

lag value

ACF

(Auto-correlation function)

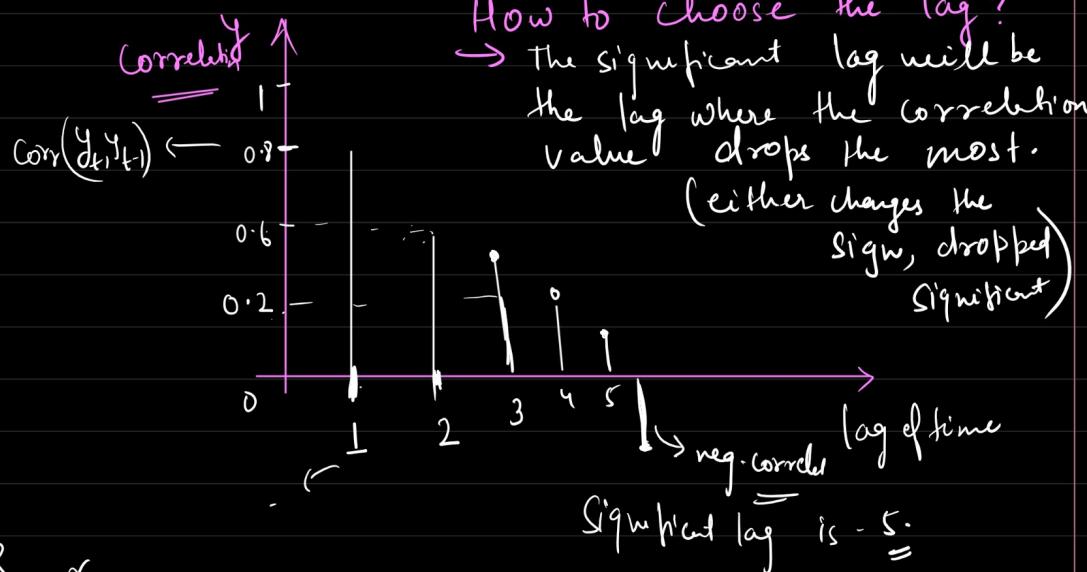
(Correlogram)

ACF (for MA)

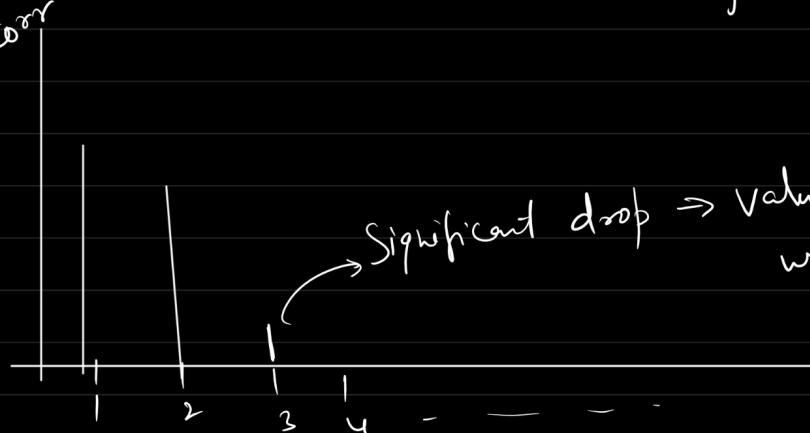
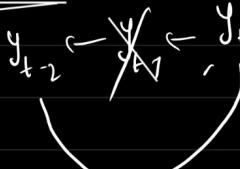
$$\text{Corr}(y_t, y_{t-1}) \checkmark$$

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

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



PACF \rightarrow p \rightarrow AR



* Autoregression

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

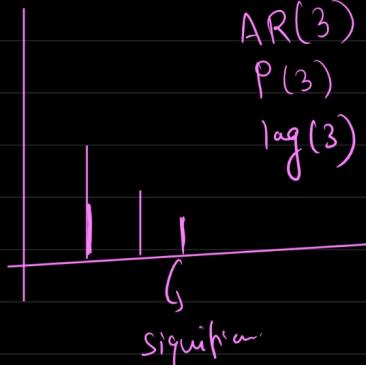
Till how much lag we have to model ??

Day	Value	y_{t-1}	y_{t-2}	y_{t-3}
1	10			
2	20	↓		
3	30	X		
4	40			
5	50	20		
6	60			
7	60	20		

$$(y_t, y_{t-1}), (y_t, y_{t-2}), (y_t, y_{t-3})$$

PACF \downarrow \Rightarrow Significant lag $\rightarrow p$

AR(3)
P(3)
lag(3)



$$y_t = \psi_1 y_{t-1} + \psi_2 y_{t-2} + \psi_3 y_{t-3} + c$$

$$\text{for } 2 \text{ lags} \quad \underline{\underline{y_t = \psi_1 y_{t-1} + \psi_2 y_{t-2} + c}}$$

$(P=2)$

* We decide p (PACF) of AR, by seeing PACF plot.

and q (ACF) of MA by seeing ACF plot.

ACF		lag 1	lag 2
10	10	ψ_1	ψ_2
20	10		ψ_2
30	20		10
40	30		20
50	40		30

* Integrated - Differencing

$$\begin{array}{c}
 \stackrel{D_0}{=} \\
 \stackrel{D_1}{\circlearrowleft} d_1 \\
 \stackrel{D_2}{\circlearrowleft} d_1 \\
 \stackrel{D_3}{\circlearrowleft} d_1 \\
 \stackrel{D_4}{\circlearrowleft} d_1 \\
 \stackrel{D_5}{\circlearrowleft} d_1
 \end{array}
 \quad y_t \quad y_{t-1}$$

$$(y_t - y_{t-1}) \rightarrow D=1$$

$$(y_t - y_{t-2}) = D=2.$$

* Moving Error :-

(No same as Smoothing model average)

→ It models the error.

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

D1	10.	8
D2	15	?
D3	20	
D4	25	
D5	30	

$\varepsilon \rightarrow \text{error}$

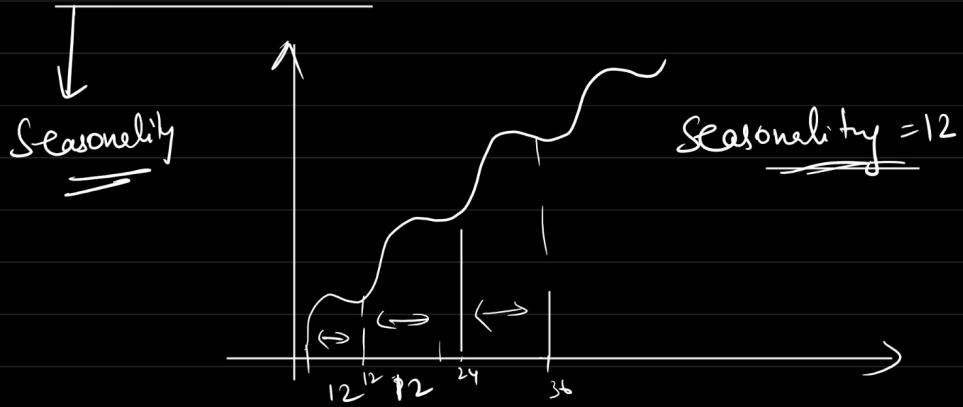
$$\begin{aligned}
 y_t = & \varepsilon_{t-1} \psi_{t-1} + \varepsilon_{t-2} \psi_{t-2} + \varepsilon_{t-3} \psi_{t-3} \\
 & \cdots \varepsilon_{t-n} \psi_{t-n} + c
 \end{aligned}$$

$$ARIMA = \frac{Y_{t-1} \psi_{t-1} + Y_{t-2} \psi_{t-2} + \dots + Y_{t-n} \psi_{t-n} + C}{AR(p)} + \frac{(Y_{t-1} - Y_{t-2} - \dots - Y_{t-n})}{\text{Integrated } I_d} + \frac{\varepsilon_{t-1} \psi_{t-1} + \varepsilon_{t-2} \psi_{t-2} + \dots + \varepsilon_{t-n} \psi_{t-n} + C}{MA(q)}$$

$AR(1)$
 $I(2)$
 $MA(1)$

$$Y_t = (Y_{t-1} \psi_{t-1} + C) + (Y_t - Y_{t-1} - Y_{t-2}) + \underline{\varepsilon_{t-1} \psi_{t-1} + C}$$

SARIMA.



$$(P, d, q_s) (P, D, Q)_s$$

$$\Rightarrow (1, 1, 2) (1, 0, 2)_{12}$$

→ Same P, d, Q for Seasonal Component.

SARIMAX

Seasonal

Exogenous Variable (Outlier, extra effect)
Noise



$$(P, d, q)(P, D, Q)_s$$

$$(0, 1, 2) (1, 0, 2)_{12} \cdot y_t \text{ (some extra variable)}$$

