

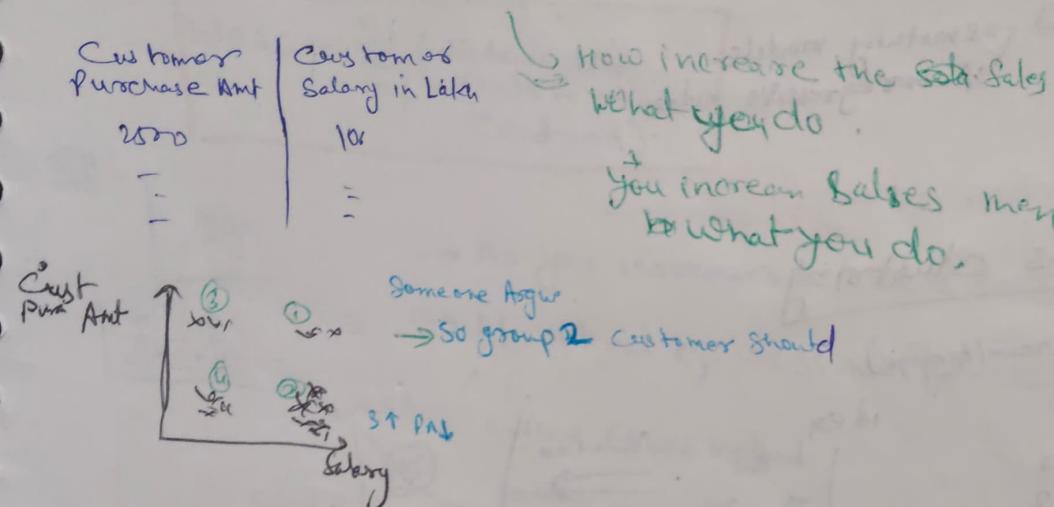
23rd June 2023

UnSupervised learning

↳ Target variable is not given(y)

↳ USL → groups / patterns / segment in Pth

Ex: you work at Zara Store as Data Sci



* Store does not matter what customer come often or not!

→ It depends you and your management to interpret the groups & design business.

→ In higher dimension, you cannot identify / groups pattern normally → So you need USL

⇒ Motivation

3:20 Need to launch a campaign based on income / sales
• Customer Segmentation

→ Different way of money laundering / fraud

→ Group of images / patterns of specific img.

→ Cohort / batch analysis.

Algo

① K-means clustering

② Hierarchical

③ DBSCAN

⇒ Silhouette Score

June-24 → K-Means, DBSCAN, Silhouette Score

① VSL ↘
k-means
Hierarchical
DBSCAN

② Anomaly detection → Isolation forest

→ DBSCAN

→ Local outlier F

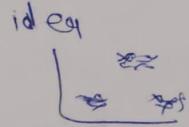
③ Time Series ↘
smoothing models
Autoregressive models

① K-means clustering

K means (mean/Average)

Ex:-

| | f_1 | f_2 | f_3 |
|---|---------------|---------------|----------------|
| 5 | 3 | 2 | |
| 2 | 4 | 3 | |
| 3 | 1 | 5 | |
| | $\frac{6}{3}$ | $\frac{8}{3}$ | $\frac{10}{3}$ |



K-means



what Avg

Avg → Arithmetic Centers of data

Centroid

$$x_{\text{cent}} = \frac{5+3+4}{3} = \frac{12}{3}$$

$$y_{\text{cent}} = \frac{2+5+6}{3} = \frac{13}{3}$$

Ex:-

How to decide the Clusters

→ Steps - ① Initialise centroid (k)

↳ Can be a random dp

↳ Centroid of all dps

↳ A random new dp's

→ ② Points nearer to the centroid will be labelled as the group

→ ③ Recalculate the centroid.

↳ Step 2 → will be repeated until centroid doesn't change.

(Any iteration Avg will be centroid)



again 2nd centroid



again 3

→ How to calculate distance

① Euclidean $\Rightarrow \sqrt{(x_2-x_1)^2 + (y_2-y_1)^2}$

② Manhattan Distance $\Rightarrow |x_2-x_1| + |y_2-y_1|$

→ How to decide K

↳ Use WCSS (Within Cluster Sum of Square distance)

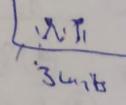
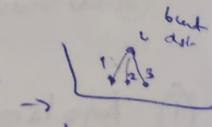
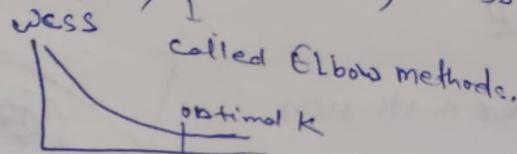
$$WCSS = \sum_{i=1}^n (\text{distance b/w Points to nearest Centroid})^2$$

e.g.:



→ As you increases k, WCSS decreases.

→ At optimal k, WCSS will not change



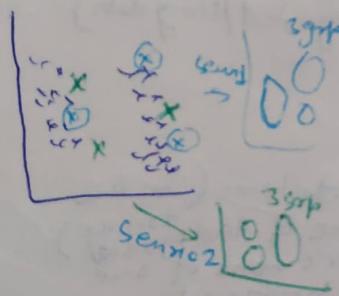
Soft class label not changing → Centroid not Change.

If all the dp are classified to the nearest centroid class.

• How k is decided \rightarrow Elbow Methods

obtuseness team also

→ Random initialization k-step!



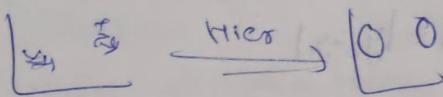
* No of Clusters depends on the initialization centroid.

where you see

-② Hierarchical clustering

Hierarchical Clustering says you give me data & I will give you clusters

1.29



→ but in HC we will not say K

• So HC says buddy I will give you complete picture then you decide K or not

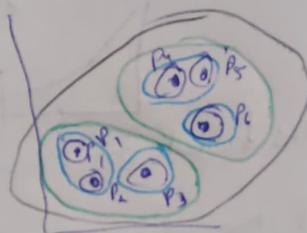
types

- ① Agglomerative clustering (Combining)
- ② divisive clustering

① Agglomerative (Combining)

① Steps:

- ② Each point is a cluster in its own
- ③ find the nearest point and create a new cluster
- ④ keep on doing step 2 until we get a single group



→ So How you are combining?

use ① Euclidean distance

② Manhattan distance

③ Cosine Similarity → To calculate distance, b/w two vectors
 $\cos \theta = \frac{a \cdot b}{\|a\| \|b\|}$ (Categorical/String data)

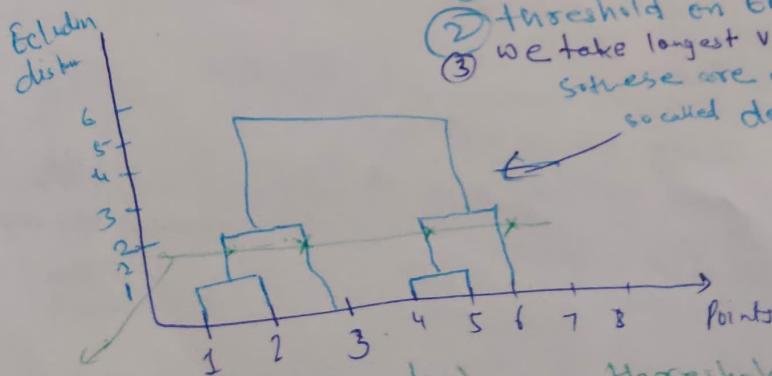
→ Can we visualize

→ Here How do you decide no. of K

① business know

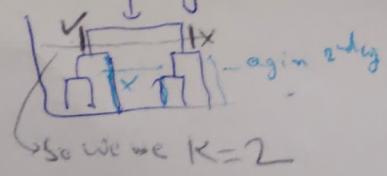
② threshold on Euclidean distance (can be tricky)

③ we take longest vertical line of dendrogram
These are dots where none of the so-called dendrogram horizontal line is passing.

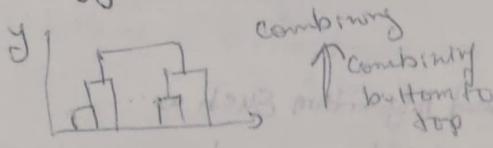


take this Euclidean distance threshold → So 4 groups

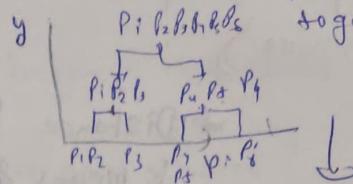
* but second Man take another Euclidean distance



→ Agglomerative



divisive



it says something
Put all data points
together.

→ Here keep dividing
until every individual dp's
becomes closer in its way

- In Industry Most time use Agglomerative because → In divisive get all data in 1 go it will be memory consume.
But in Agglo → combine based

→ Kmeans Vs Hierarchical

#DBSCAN (Density Based spatial clustering of application with noise)

- Why come

eg :- 

we need

* Distance based Algorithm Such as
K Means & hierarchical can not
be used here with Non linear Clusters

So what we do

- \rightarrow Core datapoint
 - \rightarrow border datapoint
 - \rightarrow outliers
- * min no of dp - 4 }
* radius (ϵ) - 1 } \rightarrow hyperparameters

→ Core \rightarrow the dp have atleast minimum no of dp (4)
in its radius ($\epsilon = 1$ unit)



② Border \rightarrow No of dp within its boundary its boundary (ϵ) will be less than minimum dp,



So how many dp (3) but required dp is 4 So this is Border dp.

③ Outliers \rightarrow No dp in radius (ϵ)

\rightarrow

\rightarrow Why DBSCAN? use

- find Patterns in Non-linear data
- also find outliers

\rightarrow DBSCAN with Pedia

Silhouette Score / Coefficient :-

(-1 to 1) $\rightarrow \frac{1+3}{4}$

More near to 1 better the cluster would be.

WCSS
↳ within cluster distance

OCSS
↳ outside of the distance of b/w clusters

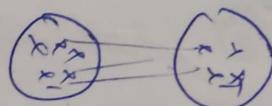
Q When we say a given cluster is a good cluster?

Point 1 : Within cluster dp's are compact (tightly packed)

↳ So WCSS \rightarrow Should be minimum ~~inter cluster~~
 \rightarrow inter cluster distance should be minimum



Point 2 : Outside cluster Sum of Square should be as maximum as possible (OCSS) from nearest cluster.
(Inter cluster distance)



SC

* Separate cluster / good cluster

WCSS $\rightarrow (a(i)) \rightarrow \min^m$

OCSS $\rightarrow (b(i)) \rightarrow \max^m$

-1 \rightarrow

0 \rightarrow

1 \rightarrow

$$\text{Silhouette Score} = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

(S_i)

i is any dp

Lect 30 June Anomaly Detection & Time Series

Anomaly detection (VSL)

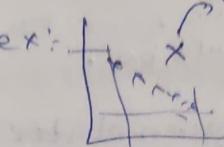
↓

Abnormality

e.g. ① fraud transaction

② gmail spam or not

③ fraud IP

ex: 
cannot be detected using traditional methods like boxplot, violinplot.

⇒ Three methods

① Isolation forest Anomaly detection

② DBSCAN

③ Local outlier factor Anomaly detection

⇒ USE CASE

↳ Specifically Solving an. outliers / Anomaly Problem

↳ DBSCAN → e.g.: find only the transaction which are abnormal / outliers.

① Isolation forest (VSL)

→ multiple Isolation trees

→ Isolation is like a decision tree.

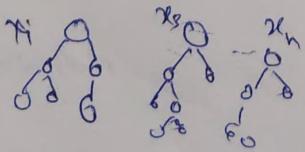
$x_1 | x_2 | x_3 | x_4$

Construction of Isolation tree :-

↳ Select a feature randomly x .

↳ Randomly choose a split value with range of x_1

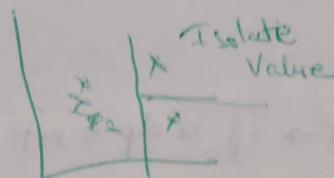
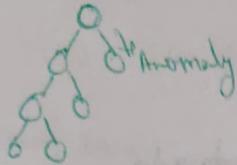
↳ Repeat the process recursively to build the tree.



(each of these isolation tree will be built to leaf node / pure node)

⇒ How you will get Outliers?

Ans Anomaly due to their distinctiveness tend to end up in leaf node at shortest path.



→ Since outliers are different from normal datapoints during construction of Isolated trees it will be isolated as a separated leaf node.

↓
How much confidence ??
↓
use

* Anomaly Score: $S(x, m) = \frac{E(h(x))}{C(m)}$
for one dp →

M = total no of datapoints.

x = datapoint for which you want to check anomaly score.

$E(h(x))$ - Average ~~Search~~ Search depth of x in all isolated tree.

$C(m)$ → Average depth of all datapoints in all isolation tree.

$$E(h(x)) \ll C(m)$$

→ Since $(C(m))$ is avg depth of all dp's so $(C(m))$ be far greater than $E(h(x))$

$S(x, m) = \frac{E(h(x))}{C(m)}$

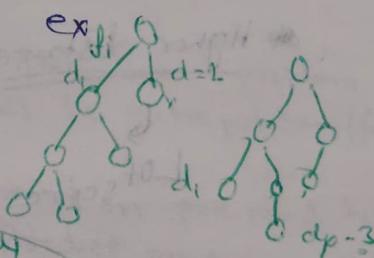
→ very small
→ very large
→ very very small

So
2 - very small no

$$\text{if } S(x, m) \geq 0.5 \Rightarrow \text{outlier}$$

threshold
 $< 0.5 \rightarrow$ Normal dp

$$\frac{1}{2} \approx 0.5$$



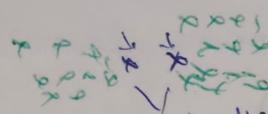
② DBSCAN

↳ Studied in USL

→ Also detects the outliers → if a dp is not core/boundary dp
↳ that's an outlier.

problems if the outlier dp are near then

③ Local Outlier factor



Local outlier (near to any cluster)

↳ global outlier (far from all the clusters)

↳ easy to detect.

How to detect local outliers

↳ Local Outlier factor

→ Calculates distance of a dp with its k-nearest neighbours
+ dp's which are far will have more distance from its neighbours & vice-versa.

↳ Higher distance \Rightarrow density less \Rightarrow Outlier.

$$\text{LoF score} \approx \frac{\text{Interest dp distance}}{\text{All dp's distance of neighbourhood}}$$

$$\text{LoF} \rightarrow 0 \rightarrow \infty$$

- $\sim 1 \sim \text{Normal}$
- $> 1 \rightarrow \text{outlier}$
- $< 1 \rightarrow \text{Normal}$

Time Series → A series with time component

Previously solve
↳ x_1
Area of house | x_2 No. of Room | x_3 Price of house

III

MLR (Multi Linear Regression)

→ So

Month Sales

| | |
|-----|-----|
| Jan | 60K |
| Feb | 10K |
| Mar | 50K |
| Apr | 80K |
| ... | ... |

Time Stamp
↳ Month
day
Years
Minute
Hours
Seconds

↳ Time component is involved

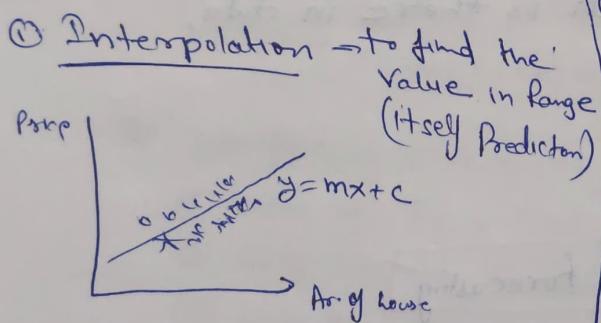
| | | | |
|-----|-----|-----------|---------|
| Jan | 15k | Mar | 60k |
| Feb | 20k | # Jan 15k | Feb 20k |
| Mar | 60k | | |

Here order matters
because the current time period value depends on the previous time period.

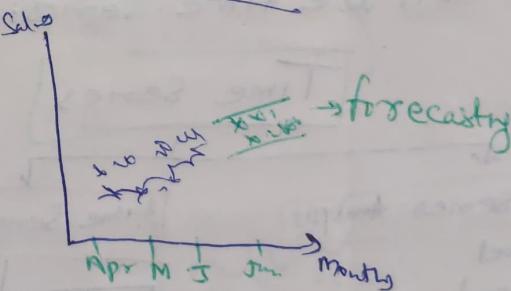
model for time

→ So Can we use Regression Series data?

→ There is 2 concept



② Extrapolation



* Most of time (99%) the test data will come in training range

e.g.: training (Area) = 1000 - 10000
testing (Area) = ...

* if outside training range ⇒ Wrong prediction

* Time Series problem will be extrapolation (forecasting)

* Based on previous history predict next 6 months.

① Why not Linear Regression for time series?

↳ Time Component is involved;

↳ Because of Extrapolation, Prediction may be wrong.

→ LR → assume linear relationship but in time series the current observation depends on previous observation → which is not true for non-time series data.

Motivation

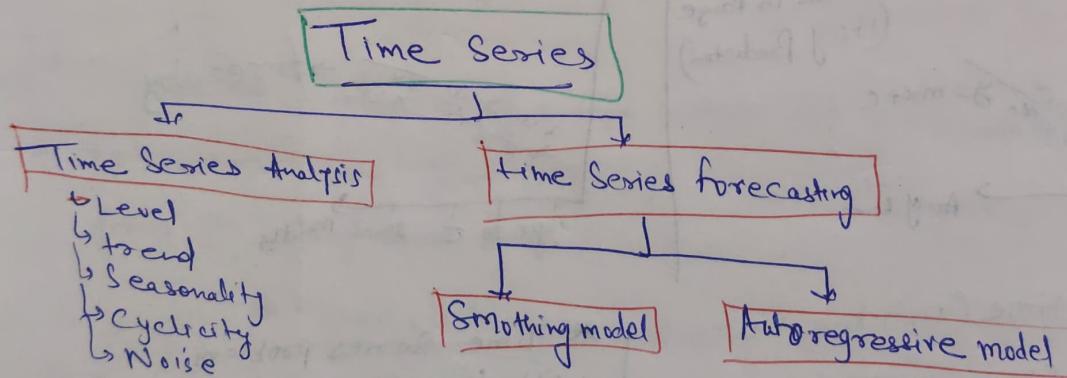
① Weather forecasting → Weather pattern day wise / month wise.

② Medical

③ Stock Mkt

④ E-commerce / finance → Sales, bond price.

Whenever time component is there in data,
then can we time series!



⇒ Components of Time Series

① Level

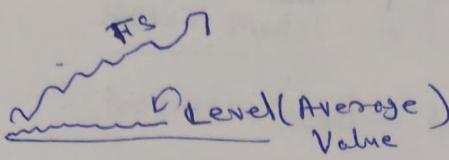
② trend

③ Season / Seasonality

④ Cycle / Cyclicity

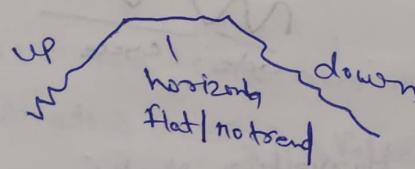
⑤ Noise

① Level :- The base value of a time series on which other components are added.



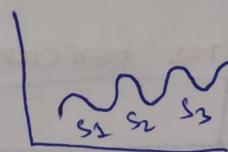
② Trend :- Long term moment or direction in data over a long period of time.

- (a) Upward
- (b) downward
- (c) Horizontal

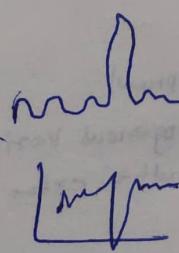


③ Season :- frequent repetition in a data over regular interval (daily, monthly, annually)

- ex:-
- ① traffic in peak hours
 - ② Sales of TV in diwali
 - ③ No of tourist in peak season
 - ④ No of Ice-cream sales.



④ Noise/Anomaly :- Some uncertainty/randomness in time series data because of unexpected reason.

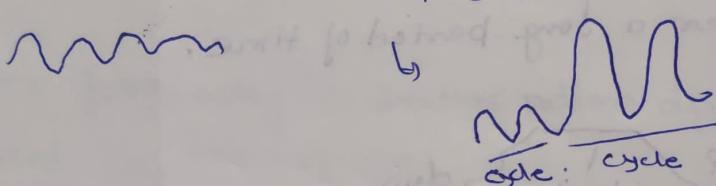


- Reason → News
- Reports
- Pandemic
- War
- Election
- Influences

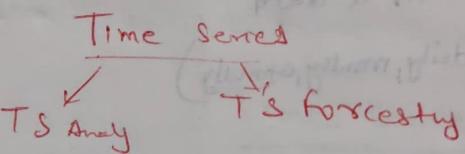
⑤ Cycle → The fluctuation in the data over a longer period of time. These periods are not fixed and can vary.

$$\text{Cyclicity} = \text{Season} + \text{Noise}$$

Q) Pharmaceutical industry like Pire



⑥ Factors making company → Harvesting Month.



Time Series forecasting

Smoothing Model

- ① Simple moving Avg (SMA) (Window Avg)
- ② Cumulative Moving Avg (CMA)
- ③ Exponential Weighted Moving Avg (EWMA)

Autoregressive Model

- SAR
- SMA (moving Avg)
- PARIMA
- SARIMA
- SARIMAX (optional)

Stationary
Over the time, value is not very change

Exogenous Variable
(Outlier, extra effect)

Smoothing

Smoothing Models :- Smoothing of TS \rightarrow Removing fluctuations

i) Naïve Model \rightarrow Last observed value is forecast

" Average Model \rightarrow Average Value is the prediction

(1) Simple Moving Avg (SMA)

$$\text{Avg} = \frac{\text{Sum of All nos}}{\text{no of Values}}$$

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

MA \rightarrow moves over the time axis in a specific window (window size, Avg Value)

eg:

window 3

Av

| | Jan | S0 | SMA |
|-----|-----|----|-----|
| f | 65 | | N1 |
| m | 70 | | N2 |
| A | 85 | | 74 |
| May | 90 | | 82 |
| J | 100 | | 80 |
| J | 110 | | 90 |
| | | | 95 |

$$\frac{50+65+70}{3} = 62$$

$$2^{\text{nd}} \text{ window} = \frac{65+71+85}{3} = 74$$

(2) Why Smoothing

\rightarrow To check trend of data

\rightarrow To remove all the effect from date

\rightarrow reduce the effect of outliers

\rightarrow Visualization

② Cumulative MA :-

→ find all the Avg. if all the datapoint upto given timestamp.

→ for long period of time

→ Exponential trend

Eg: Jan 10

CMA

$$10$$

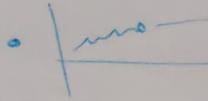
$$\frac{10+12}{2} \Rightarrow$$

Mar 15

$$\frac{10+12+15}{3}$$

Apr 14

$$\frac{10+12+15+14}{4}$$



SMA, CMA, Naive, Avg → It gives equal priority to all the value.

but ↓

→ In ~~Time Series~~ Current observation is highly influenced by last few observations

∴ You need to give priority to recent observations.

③ EMA or EWMA (Exponentially Weighted MA)

→ We give more weightage / importance / priority to the recent datapoint / time stamp

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

V_t = EMA at time t

$\beta = 0 < \beta < 1 \rightarrow$ generally it's 0.9 (Weight)

V_{t-1} = EMA at previous t_s .

θ_t → Data at current stamp.

eg:-

| V_{t-1} | Month | Value | End |
|-----------|----------------|-------|------|
| v_0 | D ₁ | 25 | 0 |
| v_1 | D ₂ | 13 | 1.3 |
| v_2 | D ₃ | 17 | 2.87 |
| v_3 | D ₄ | 31 | |
| v_4 | D ₅ | 43 | |

$$\beta = 0.9 \text{ (To reduce the previous time stamp effect)}$$

$$\frac{1}{1-\beta}$$

$$V_1 = \beta \cdot V_0 + (1-\beta) \theta$$

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

$$= 0.1 \times 13 \Rightarrow 1.3$$

$$V_2 = \beta \cdot V_1 + (1-\beta) \theta$$

$$= 0.9 \times 1.3 + 0.1 \times 17$$

$$= 2.87$$

In

Time series 2 way to calculate

Additive

$$Y_t = \text{Trend} + \text{Season} + \text{Noise}$$

→ Linear over time

→ Constant Variance

→ Increased trend at same difference

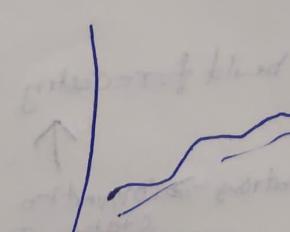
Multiplication

$$Y_t = \text{Trend} \times \text{Seasonality} \times \text{Noise}$$

→ Non linear

→ Non Constant Variance

| | |
|----------|-----|
| Eg Day 1 | 100 |
| Day 2 | 200 |
| Day 3 | 300 |
| Day 4 | 400 |



So TS → forecasting →
Smoothing Autoregression

→ You can build a Smoothing time series model on any time series data.

* To build Auto-regressive model, statistical prediction properties of a time series like mean/Variance should be constant (Not change over time)

e.g.: Scenario 1



Scenario 2



- It will be easy to build a forecasting model in sc-2
- less variance in data.

* Non Stationarity → Mean-Variance will not be constant.

* Stationarity → Over the time, value is not varying (changing)

* So → To build a Auto-regressive model make the time series stationary.

In ML → Data Ingestion → Analysis → Preprocessing → Model building

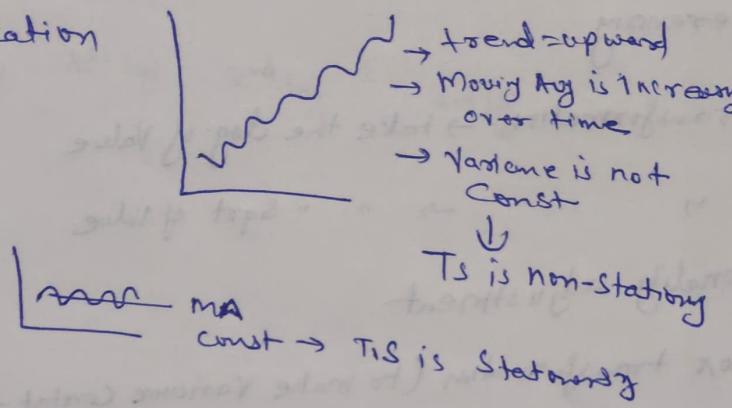
→ Time Ser →

→ TSA → TS is stationary → build forecasting Model
if non stat → convert to stationary TS

⇒ To check if a TS is stationary or not

b.

① Visualization



② Statistical test

④ ADF (Augmented dickey fuller test)

H_0 = TS data is Non Stationary

H_A = TS " " Stationary

p Value $\leq 0.05 \rightarrow$ reject H_0

Conclude → TS data is Stationary

⑤ KPSS test (kutkouesk - philips - schmidt - schin test)

H_0 : TS data is Stationary

Q) How to Convert non-st TS to st TS

- ① Differencing
- ② log transformation → take the log of Value
- ③ root " → " " Sqrt of Value
- ④ Seasonality adjustment
- ⑤ Box-Cox transformation (to make Variance Constant \rightarrow +ve)
- ⑥ Yeo-Johnson (All TS data)

\Rightarrow ① Differencing \rightarrow Difference ($y_t - y_{t-1}$)

current \downarrow previous

| | Month | Price | 1st diff | 2nd order | 3rd |
|---|-------|-----------|-------------|-----------|-----|
| J | 5 | NA | NA | NA | NA |
| F | 10 | $10-5=5$ | NA | NA | NA |
| M | 6 | $6-10=-4$ | $-4-(-5)=1$ | NA | NA |
| A | 8 | $8-6=2$ | $2-(-4)=6$ | NA | NA |
| M | 15 | 7 | NA | NA | NA |
| J | 7 | -8 | NA | NA | NA |

Check st \rightarrow 3rd order

After diff, check if TS is Stationary

Statistical test (ADP) \swarrow \searrow
 Visculation (using MA)

* if Stationary \rightarrow build the forecasting model

else \rightarrow Again do differencing

* ACF \rightarrow Auto Correlation function

* PACF \rightarrow Partial Autocorrelation function

* Autoregression & ACF \rightarrow Auto + Correlation

\downarrow
Correlation Hself in feature

\hookrightarrow Relationship b/w 2 Variable

* ACF measures the correlation b/w time series
& its lag value.

* Month y_t 1st lag 2nd lag 3rd lag