



Anamoly Detection Notes

- Darshan R M

"Learn, Share,
and
Collaborate"

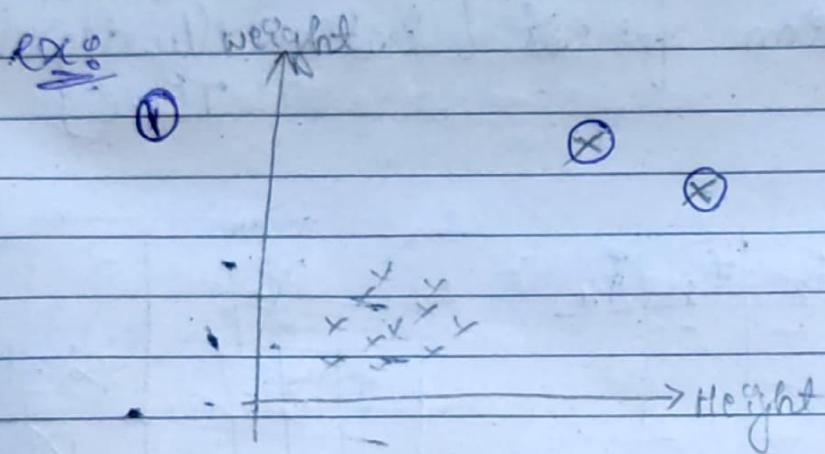
Week - 20

Anomaly detection and

Time series

* Anomaly detection

[To detect outliers] → play a ~~imp~~ role.
↳ used when if outliers is the important parameter for solving a problem.



- ② if person will have a cancer or not it is a very important outlier ↴

In this type of problem statement where outliers play major role we use "Anomaly detection".

- ③ IP address if hacker is hacking from different location we can detect

④

IPL

over runs/over

1 15

2 10

3 12

T 4 100 \rightarrow outlier
as run > 36

16 - 100

* Types of anomaly detection

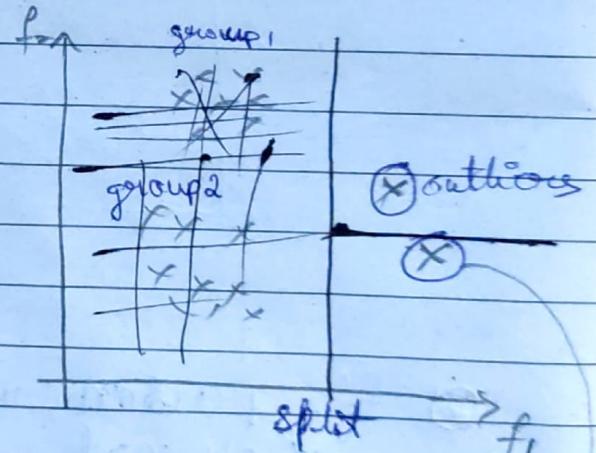
① Isolation forest [internally decision tree]

\Rightarrow consider,

f_1	f_2	f_3	f_4
-	-	-	-
-	-	-	-
-	-	-	-
-	-	-	-

DT called as

~~isolated
tree~~



\Rightarrow Create as DT

for each point.
until we get leaf node.

We make this point completely isolated

~~now~~
quickly
we are
able to

isolate a
data point
quickly.

isolate pt. \rightarrow outliers

* we calculate a "Anomaly score"

↓

so that we get to know
why we are considering
particular point as a
anomaly / isolated point / outlier.
if gross
particular
threshold
then we consider

Maths formula
for Computation
of anomaly score
for a new point

$$s(x, m) = \sigma^{-\frac{E[h(x)]}{c(m)}}$$

where,

$m \Rightarrow$ no of data points

$x \Rightarrow$ Data point

$E[h(x)] \Rightarrow$ Average search depth for x
from the isolated tree.

Now, we are not
considering for

x data $\leftarrow C(m) \Rightarrow$ Average value of $h(x)$ / depth
point in average

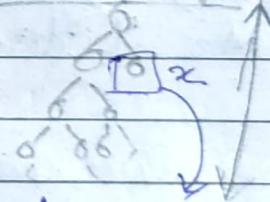
we consider for
all points

what is the

average search
depth of the

x from all
different isolation trees.

Computing only
first 'n' data point



we can consider multiple feature
and split the tree so

we will get multiple different
isolated trees.

* If we have an outlier then

For outlier,

$$E[h(x)] \ll c(m) \Rightarrow [s(x, m) \approx 1]$$

$$\text{as } \frac{E[h(x)]}{c(m)} \approx 0. \rightarrow i \neq j \approx 1$$

Outliers \leftarrow anomaly score

For outlier,

we can keep a

$\boxed{\text{Threshold} \geq 0.5}$ \leftarrow threshold and then categorize as outlier.

+ $E[h(x)] > c(m) \Rightarrow [s(x, m) \approx 0.5]$

\downarrow
Normal points

* Internally, we are creating isolated tree.

* Defference, DT and Isolated tree

\rightarrow we are creating isolated tree for each and every data points.

* Implementation

① Import data.

* Visualize outlier

② Anomaly detection

> from sklearn.ensemble import IsolationForest

> clf = IsolationForest(contamination=0.2)

> clf.fit(df)

> prediction = clf.predict(df)

— —

③ plot

> index = np.where(prediction < 0)

> plt.scatter(df.iloc[:, 0], df.iloc[:, 1])

> plt.scatter(x=index[0], y=index[1],
edgecolors='red')

— —

* Anomaly detection using (DBSCAN.)



Can do non-linear

we detect a ~~cluster~~ outlier here

If those outlier were important ~~for~~ for our problem statement.

then it's a kind of anomaly detection.

* 3 points

No of datapoint \geq min
within ϵ points

- ① Core point
- ② Border point
- ③ Outlier/ Noise

Non-linear clustering

No of datapoint ≤ 1
within ϵ

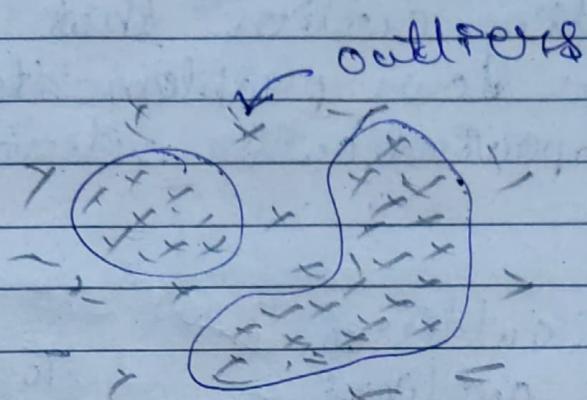
If these outlier are ~~imp~~ for our problem then set is a anomaly detection.

* 2 hyper parameters

- ① minimum points.
- ② Epsilon [ϵ] - radius.

* DBSCAN is robust to outliers.

* Outliers



* Implementation

> from sklearn.datasets import make_moons, make_circles.

> from sklearn.cluster import DBSCAN

> $x, y = \text{make_circles}(n_samples=500, \text{factor}=0.3, \text{noise}=0.1)$

> plt.scatter(x[:, 0], x[:, 1])

> dbSCAN = DBSCAN(eps=0.1)

> dbSCAN.fit(x)

> dbSCAN.labels_

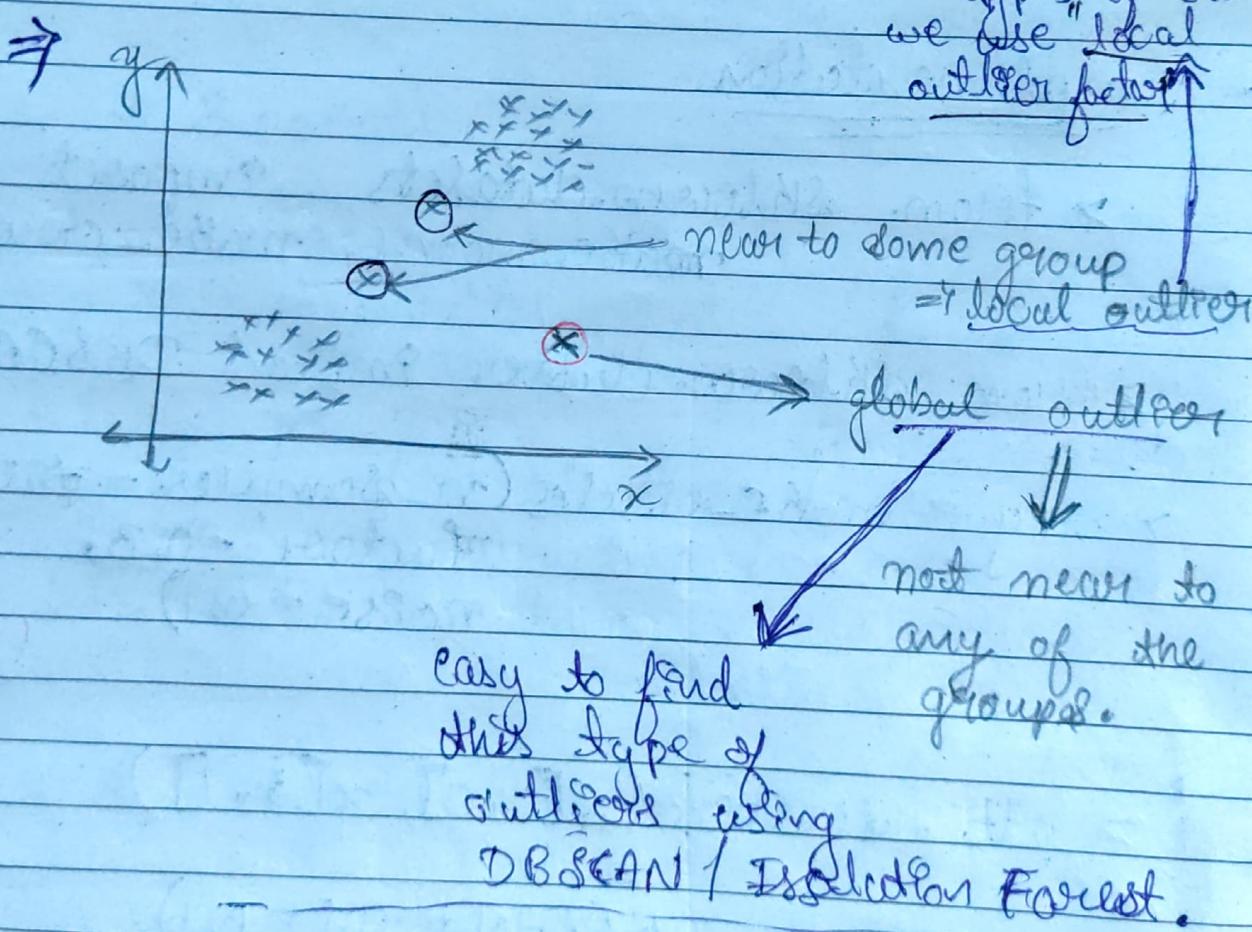
LOF

* "Local outlier factor" Anomaly detection.

→ Finding outliers that are imp for our problem statement
 → Unsupervised outlier detection using LOF.
 ↪ 2 terms:

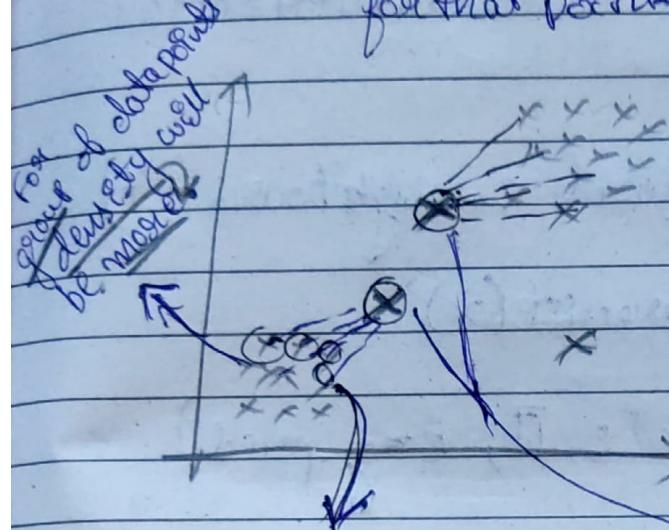
- ① Local outliers. ↪ LOF score
- ② Global outliers.

To find these types of outliers we use "local outlier factor"

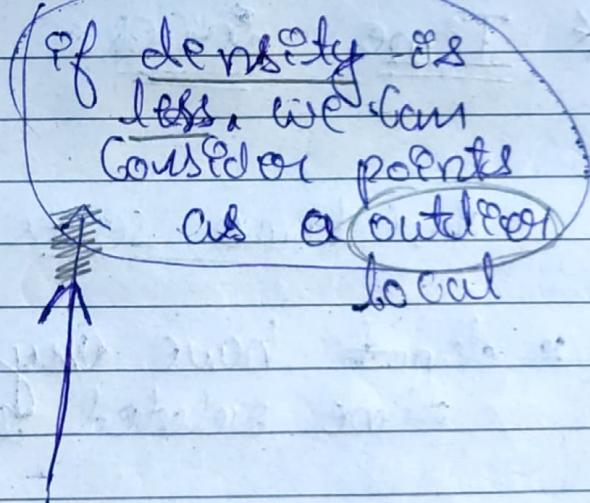


- * Internally it uses K-Nearest neighbour
Local outlier density

* Suppose, To check those local outlier are outliers (or) not. for that point



we consider those points and then we average observe what distance is the density w.r.t neighbour points.



* If the average distance is more then density is less.

means points are far.

* Here, it compares local density of samples with local density of its neighbors.

Can identify samples that have a lower density than their neighbors.

outliers

density than their neighbors.

* Implementation

> from sklearn.neighbors import LocalOutlierFactor

> lof = LocalOutlierFactor(n_neighbors=7)

> ypred = lof.fit_predict(x)

> plt.scatter(x[:, 0], x[:, 1], c=ypred)

* Time Series

Non time series

* doesn't have any time related column.

e.g. House price prediction

size . location . bedroom | price

No time representing column

regression problem.

Time series

* At one column will be time

e.g.

Sales data

Day 1	50K
Day 2	60K
Day 3	70K
Day 4	70K
Day 5	60K

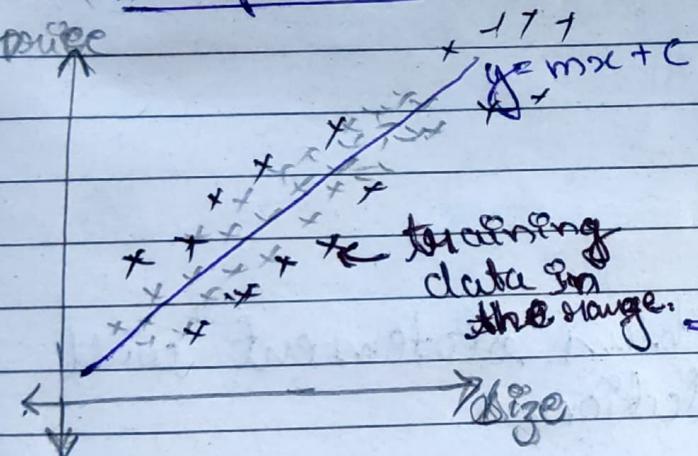
time stamp

* Time stamp can be hour, min, sec, day, month

Q Can we solve time-series problem with linear regression.

⇒ Yes, but not recommended.
because of noise

* Interpolation ⇒ To find out the value
within the range



→ Supervised learning [regression]
[test data will be
out of range.]

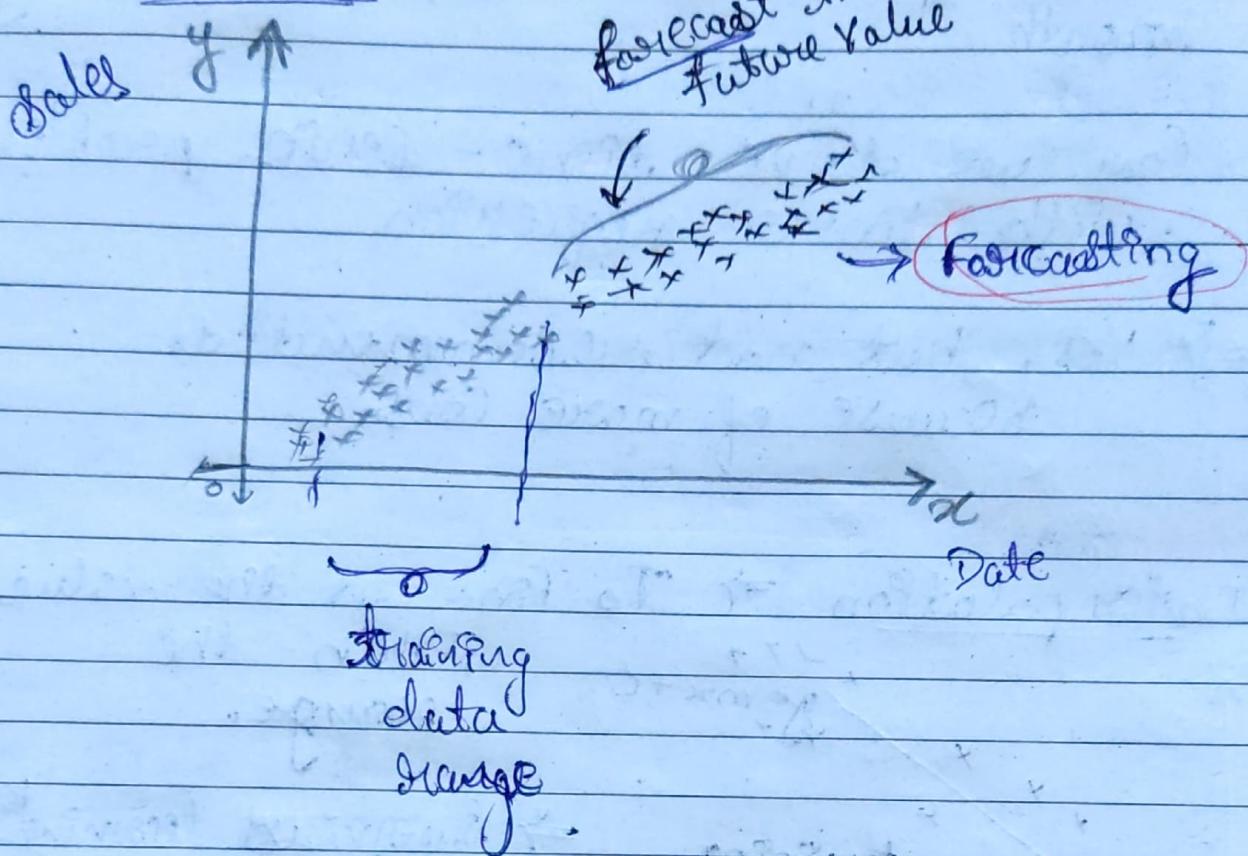
* Extrapolation ⇒ To find out the value
out ~~within~~ the range.

if training data
range is 0-100

if test data
comes like → 60, 70, 40, 20, 200, 1000,

↓
lead to wrong
prediction.

Extrapolation



* Time series problem statement will be a extrapolation.

⇒ Based on previous history forecast the future values

* If we use linear regression for the time series it may lead to the wrong prediction.

* ~~Limit of regression on Time series data~~

① Because of extrapolation.

② Because of outlier.

③ There should be a linear relationship.

but in time series ~~data~~ data must be
well a linear.

⇒ Data will be complicated.

④ ~~Non-time~~

Non-time series data → ~~no~~ effect of previous time

~~but~~

time series data will have → previous time
effect of

⇒ Non time series data

size	location	price
1000	Delhi	55L
2000	Mumbai	1cr.

Independent of previous data

⇒ Time series data

Time.

dependency

Day1	sales
Day2	50K 60K

dependent on previous data.

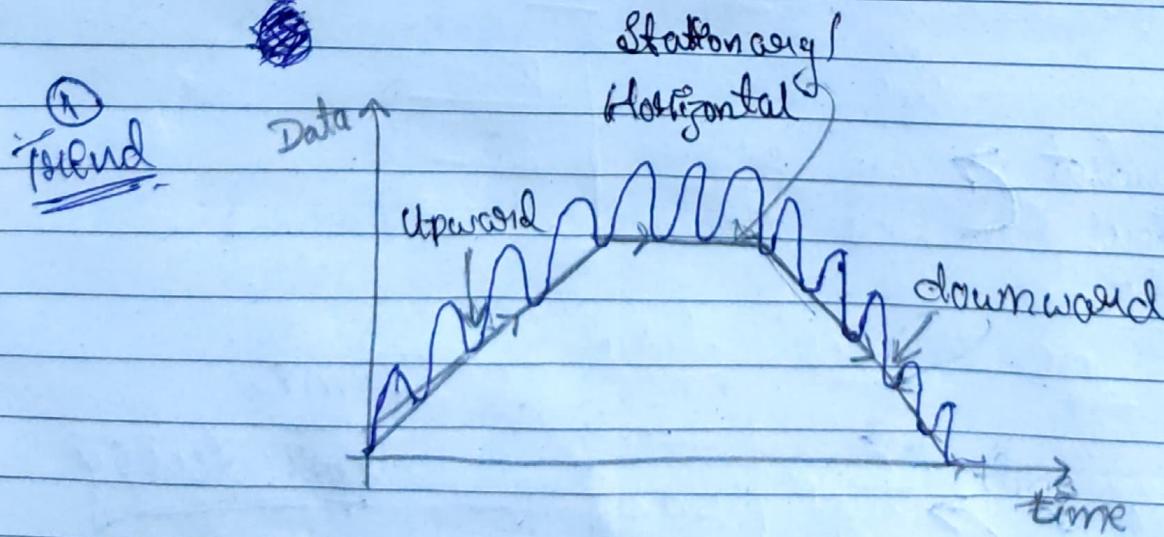
* Example of time series

- ① Economics forecasting. \rightarrow GDP, interest, ...
- ② Finance \rightarrow Sales, bond price
- ③ Weather Forecasting \rightarrow weather, pattern
- ④ Medical. \rightarrow based on previous hist for future conclusion.

* Time series. $\xrightarrow{\text{can be applied to}}$ time dependencies is present.

* Time series components

- ① Trend.
- ② Season.
- ③ Cycle.
- ④ Noise.



Trend

- ① Upward.
- ② Horizontal/ Stationary.
- ③ Downward.

③ Season

based on daily, yearly, monthly, hourly.

→ it is a frequent repetition

e.g. Sales of Scream is Summer,

② Traffic at 5pm in my area.

③ Tourist in Goa at the end of year. [Dec/Jan]

④ Cycle

→ (cyclic pattern)

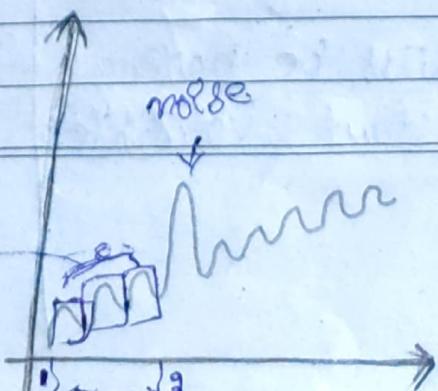
→ time series behavior over the long time.

Cycle \Rightarrow season + noise
[fluctuation]

Ex

stock price

$BIDP = ?(\text{season})$

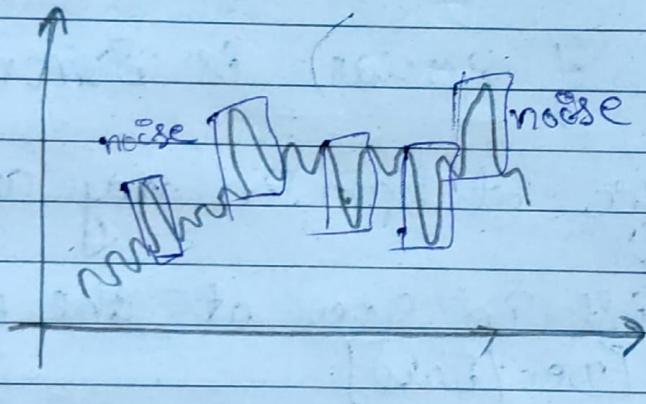


④ Noise

↳ Some uncertainty / randomness in data because of unpredictable reason.

Cri^o pandemic, war, news, -

Corona, palestine/ Adani, Spain



* $y_t = T + S + N$ → multiplicative TS

$y_t = T + S + N$ → additive TS

Additive
Multiplicative

Multiplicative

* It will be linear over time.

* Non-linear

* It will be having constant variable.

* Non-constant Variable

* Moving average

Types

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

a) why we calculate Moving average?

↳ Smoothing of time series data.

① Simple moving average:

Average =

$$\Rightarrow \text{data} = [10, 12, 15, 13, 11]$$

$$\Rightarrow \text{average} = \frac{\text{sum of all values}}{\text{no of values}}$$
$$= \frac{10 + 12 + 15 + 13 + 11}{5}$$

$$\text{Average} = 13.2$$

Moving → move over the time axis
in specific window time.

- { ① window size.
② average calculating.

Consider,

time Sales

D1	10	}	window size 3
D2	12		
D3	15		
D4	13		
D5	14		
D6	16		
D7	17		

$$\text{1st avg} = \frac{D1 + D2 + D3}{3} = \frac{10 + 12 + 15}{3} = \underline{\underline{12.33}}$$

2nd avg

$$= \underline{\underline{13.33}}$$

$$= \underline{\underline{14}}$$

$$= \underline{\underline{14.33}}$$

$$= \underline{\underline{15.66}}$$

3rd avg

Calculating due to smoothing of data.

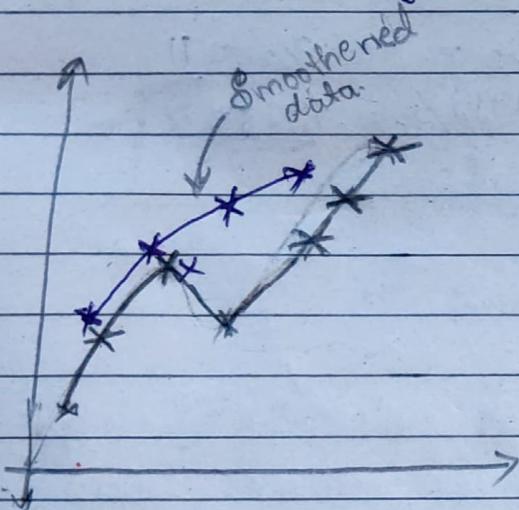


To remove all effect from the data.

We can perform smoothing with the help of moving average.

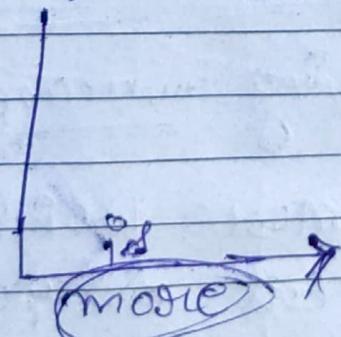
Advantage of smoothing

- ① Pattern recognition from the data.
- ② You can analyse the trend of the data.
- ③ You can reduce the effect of outlier.
- ④ Enhancing the visualization of data.



* Window size affect the smoothing.

If window size is less → then it will be more specific to point.
(not good) Smoothing



→ Loose smoothing

② Cumulative moving average [CMA]

→ Find out average of all the data point upto given time stamp.

Consider

D ₁	10	D ₁ = 10
D ₂	12	D ₁ + D ₂ /2 = 10 + 12/2 = <u>11</u>
D ₃	15	10 + 12 + 15 / 3 = <u>12.33</u>
D ₄	19	
D ₅	16	
D ₆	17	
D ₇	11	

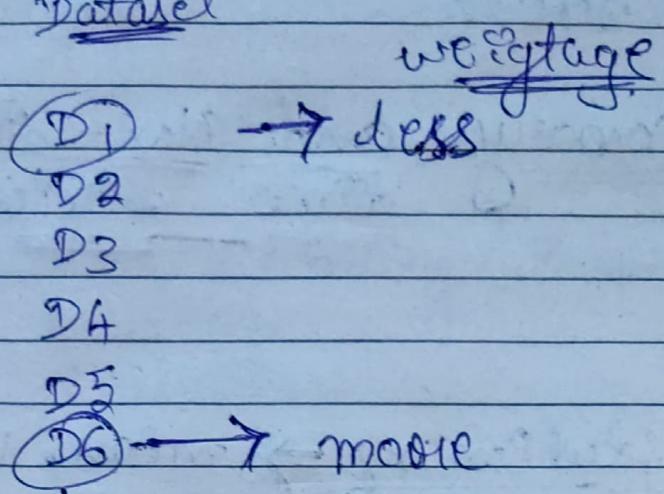
- * Used for long time period.
- * give exponential trend.

* Here, also we do smoothing of data

③ Exponential moving average [EMA]:

→ we give more weight to recent datapoint / timestamp.

Consider, Dataset



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

V_t = EMA at time t ,

amount ← β → $0 < \beta < 1 \Rightarrow [0.9]$ be value may of smoothing

V_{t-1} ⇒ EMA at previous time stamp.

$\theta_t \Rightarrow$ Data at time stamp t .

Consider,

$V_0 \rightarrow D_1$	25
$V_1 \rightarrow D_2$	13.
$V_2 \rightarrow D_3$	17
D_4	31
D_5	43

$$V_0 = 0 / V_0 = 25$$

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

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

$$= 0 + 0.1 \times 13 = \underline{\underline{1.3}}$$

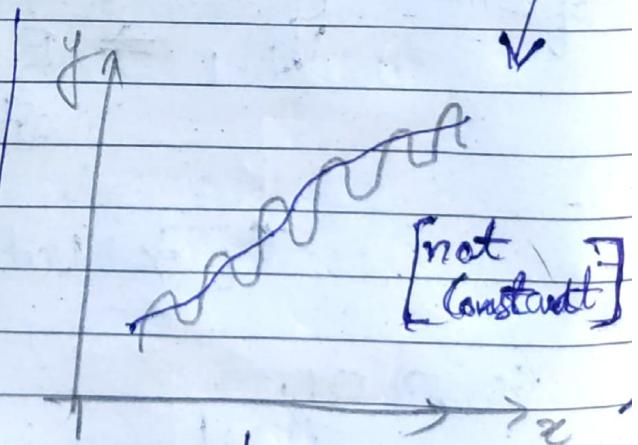
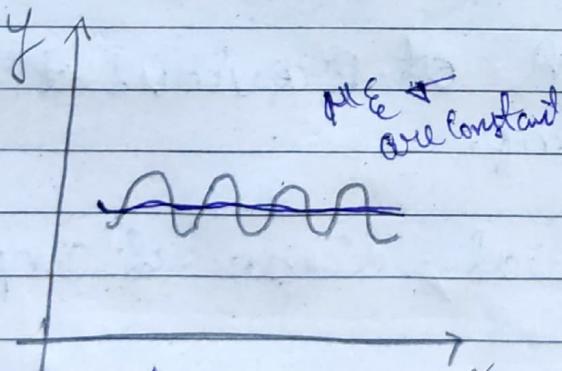
* ~~Time series~~

* Stationary and Non-Stationary Time Series

* Non-stationary \rightarrow mean, Variance
time series will not be constant.

* stationary time series \rightarrow mean, Variance and Constant.

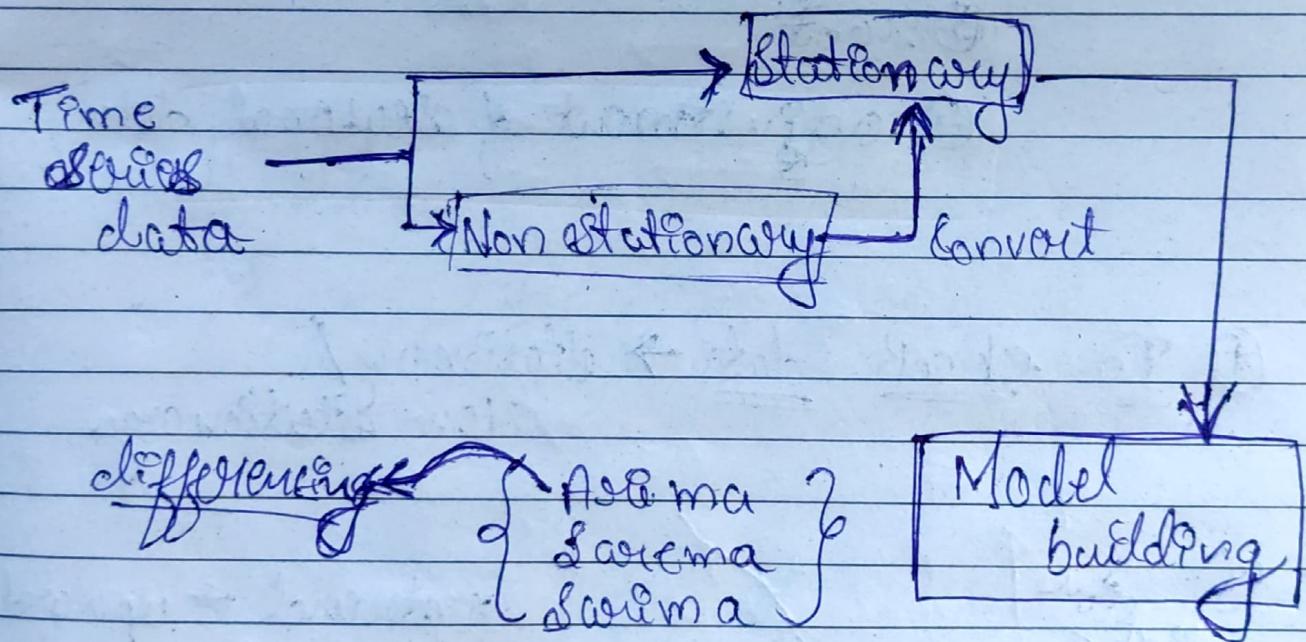
* Constant \rightarrow values are not varying about time axis.



upward trend

Machine learning

- ① Data ingestion.
- ② EDA.
- ③ Feature engineering / Preprocessing.
- ④ Model building.
- ⑤ Model evaluation.



* For checking ~~my own~~ Time series data stationary / Non-stationary time series

- ① Visualization.
- ② Stats based test.

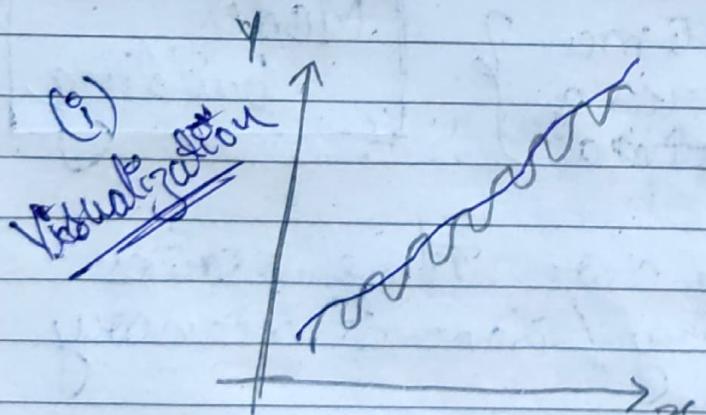
* If Non-stationary Time series data is received then do convert it into ~~Time~~ stationary Time series data.



- ① Differencing.
- ② Log
- ③ Root

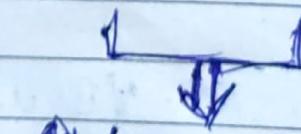
- ④ Adjustment of seasonal data.

① To check data \rightarrow stationary / Non-stationary.

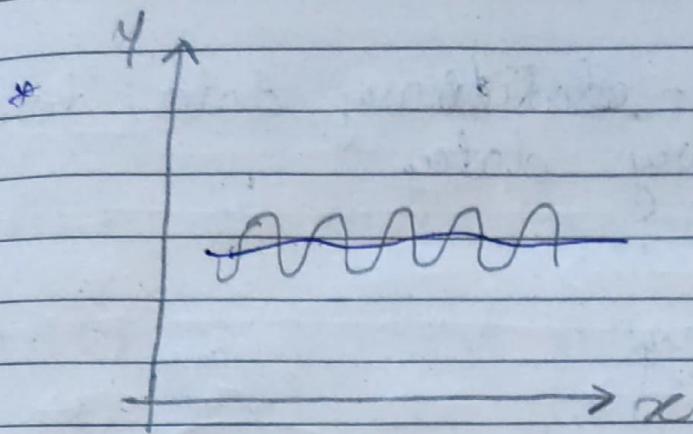


+ trend \Rightarrow upward

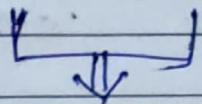
+ moving average \Rightarrow increase over time



Non-stationary



* Trend \Rightarrow flat
+
* moving average \Rightarrow constant over time



Stationary

(ii) Stats test

\rightarrow we use ADF test.

[Augmented Dickey Fuller Test]

① Stats test.

② P-value \rightarrow observing p-value we can tell

③ Critical value

TS data is stationary/not

* $H_0 \Rightarrow$ data is Non-stationary.

$H_1 \Rightarrow$ data is stationary.

Stationary $\leftarrow [P \leq 0.05\right] \rightarrow$ reject (H_0) null hypothesis

Non-stationary $\leftarrow [P > 0.05\right] \rightarrow$ fail to reject null hypothesis

② Convert Non-stationary data to stationary data.

① Differencing

Consider,

		if Non-stationary	
		First differencing	Second differencing
D ₁	5	N/A	
D ₂	10	5	
D ₃	6	-4	
D ₄	8	2	
D ₅	15	7	
D ₆	7	8	

↓
test TS data

↓
if Stationary/
Non-stationary

* difference the data until we
get a stationary data.

* ACF, PACF, Auto regression

ACF \rightarrow Auto Correlation Function

PACF \rightarrow Partial Auto Correlation Function

In ACF \Rightarrow AUTO + CORrelation

Correlation within
some feature for
particular time
interval.

relationship b/w
2 features

: Types / methods

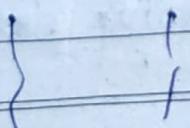
$$r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{x} \cdot \sqrt{y}}$$

- ① Pearson.
- ② Spearman rank.
- ③ Kendall.

$$r_{xy} = \frac{\text{Cov}[R(x), R(y)]}{\sqrt{R(x)} \cdot \sqrt{R(y)}}$$

~~PACF~~ measure the correlation between
time series and lag value.

$$\begin{aligned} &\text{Corr}(Y_t, Y_{t-1}) \\ &\text{Corr}(Y_t, Y_{t-2}) \end{aligned}$$

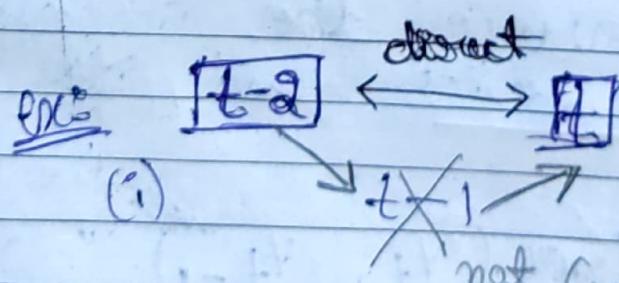


	<u>dimensions</u>	<u>1st lag</u>	<u>2nd lag</u>
D1	10	N/A	
D2	25	\leftrightarrow ^{relation} D ₁ \rightarrow 10	N/A
D3	14	\leftrightarrow D ₂ \rightarrow 25	\leftrightarrow D ₁ \rightarrow 10
D4	16	\leftrightarrow D ₃ \rightarrow 14	\leftrightarrow D ₂ \rightarrow 25
D5	20	\leftrightarrow D ₄ \rightarrow 16	\leftrightarrow D ₃ \rightarrow 14
D6	32	\leftrightarrow D ₅ \rightarrow 20	\rightarrow D ₄ \rightarrow 16

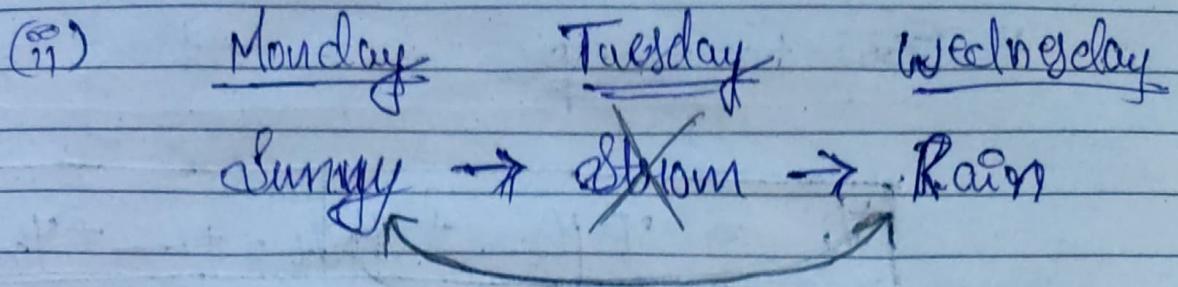
Auto Correlation

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

* PACF [Partial Auto Correlation Function]



not consider
remove t-1 from
finding correlation



if we want to find relationship b/w monday and wednesday
 \Rightarrow In PACF we ignore / do not consider - Tuesday
 bc we assume that there will be no storm b/w monday and wednesday.

* PACF \Rightarrow we remove the intermediate effect.

* Auto regression

Auto + Regression

Regression & itself within the Variable.

$x \leftarrow y$ \downarrow relationship
 simple linear re

$$h_y(x) = \theta_0 + \theta_1 x,$$

$x \rightarrow$ independent
 $y \rightarrow$ dependent

* ~~Fit~~ formula

$$Y_t = \psi_1 Y_{t-1} + \psi_2 Y_{t-2} + \dots + \psi_n Y_{t-n} + C + \epsilon$$

where,

$Y_t \Rightarrow$ Value at current time stamp

$\psi \Rightarrow$ Coefficient term

$C \Rightarrow$ Constant

$\epsilon \Rightarrow$ error

epsilon

* To know until which term we have to consider we know in ARIMA:

* ARIMA
p d q

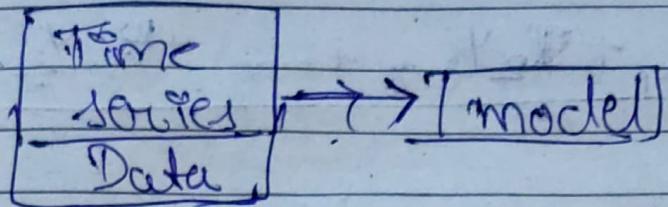
→ ~~Fit~~ (Auto regression) (integrated)
moving average. parts

Machine
learning

[Data] → [Model]

linear regression
SVM

KNN



ARIMA

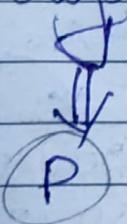
SARIMA, SARIMAX

[Deep learning] DL: RNN; Attention, Transformer.

ARIMA

AR

[Auto regression]



Value $[0, 1, 2, \dots, n]$

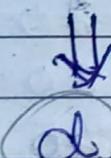
lag value

\Rightarrow PACF

[Corr. logram graph]

I

[Integration]



$[0, 1, 2, \dots, n]$

lag value

\Rightarrow differencing

[stationarity]

& now

stationary

MA

[Moving average model]



we analyze moving average

& calculate

error from current timestamp to previous timestamp

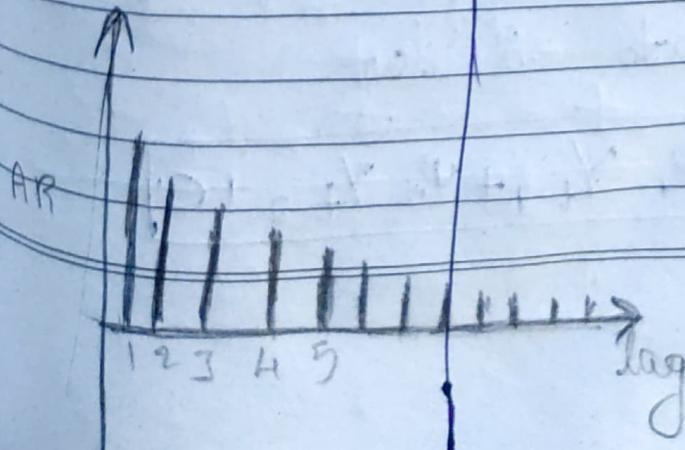


$[0, 1, 2, \dots, n]$

lag value

\Rightarrow ACF

[Corr. logram]



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



* AR(1) \Rightarrow upto 1 lag

Current value \downarrow Previous \uparrow

$$Y_t = \psi_1 Y_{t-1} + c$$

lag = 1
AR = 1
P = 1

Consider:

TS Data y_t

D ₁	10
D ₂	20
D ₃	30
D ₄	40
D ₅	50

y_{t-1}

lag 1 TS

N/A

D₁ = 10

D₂ = 20

D₃ = 30

D₄ = 40

Previous Cova

20	10
30	20
40	30
50	40

Future
Forecasting

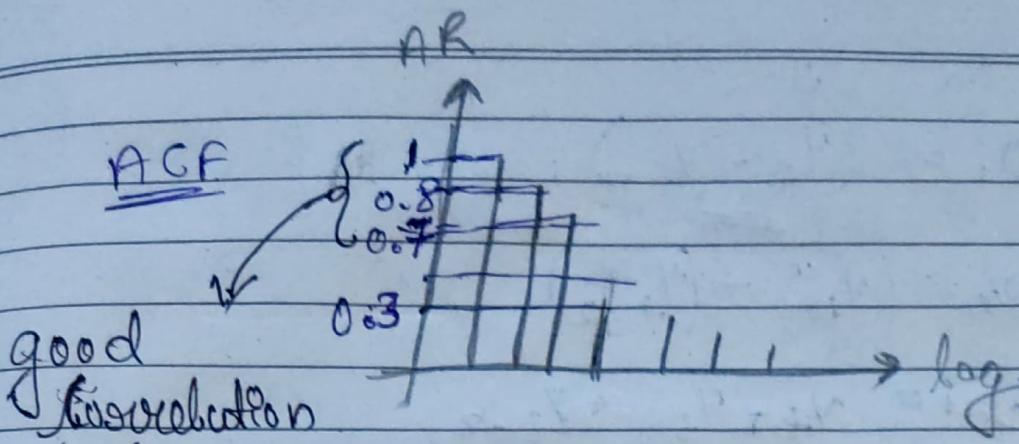
22	45
22	2

$$Y_t = \psi_1 Y_{t-1} + c$$

y_t y_{t-1}

* AR(2) \Rightarrow upto 2nd lag

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



↳ so we take $AB(3)$

because we don't

Consider intermediate

Intermediate effect also"

Here, usually
we don't take
(ACE)

D_1	y_1	$(y_t - y_{t-2})$
D_2	20	
D_3	30	40-20
D_4	40	20
D_5	50	60-40
D_6	60	20

→ we don't
consider the
intermediate
effect.

* \Rightarrow Integration \Rightarrow difference

Given $D=0$ $D=\beta$

$y_{t-4} D_1$

$y_{t-3} D_2$

$y_{t-2} D_3$

$y_{t-1} D_4$

$y_t D_5$

$D_2 - D_1$

$D_3 - D_2$

$D_4 - D_3$

$D_5 - D_4$

$$D_0 = Y_t$$

$$D_1 = \underline{(Y_t - Y_{t-1})}$$

$$D_2 = (Y_{t-1} - Y_{t-2})$$

$$D_3 = Y_t - \underline{Y_{t-1}} + Y_{t-2}$$

Moving average

↳ model build with the errors.

$$Y_t = E_{t-1} \Psi + c$$

Coefficient

Coefficient term

Constant

Moving
average
model

$$Y_t = E_{t-1} \Psi_{t-1} + E_{t-2} \Psi_{t-2} + \dots + E_{t-n} \Psi_{t-n} + c$$

general
 Σe^n

built on
error term

ARIMA ↗

$$Y_t = (Y_{t-1} \Psi_{t-1} + Y_{t-2} \Psi_{t-2} + \dots + Y_{t-n} \Psi_{t-n} + c)$$

$$Y_t = (Y_{t-1} \Psi_{t-1} + Y_{t-2} \Psi_{t-2} + \dots + Y_{t-n} \Psi_{t-n} + c) + \\ (Y_{t-1} - Y_t - Y_{t-2} - \dots - Y_{t-n}) + \\ (E_1 \Psi_{t-1} + E_2 \Psi_{t-2} + \dots + E_n \Psi_{t-n} + c)$$

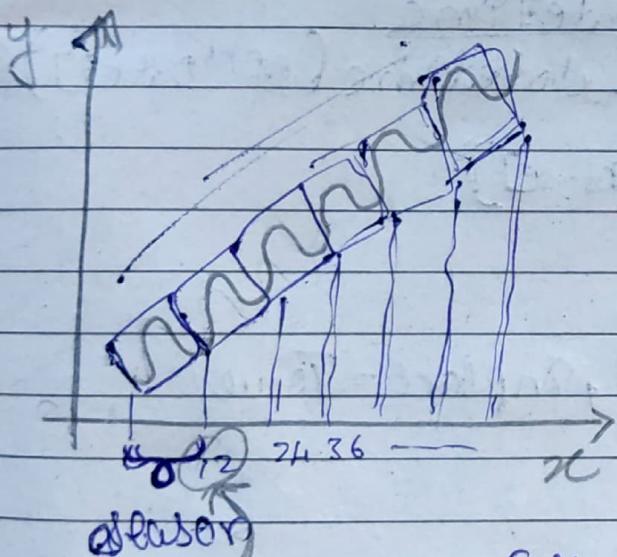
this same eqn is
used to forecast
the future values.

~~if~~ $AB(1), P=1$
 $I(2), d=2 \Rightarrow Y_t = (\psi_{t-1} \psi_{t-2} + c) +$
 $MA(1), q=1 \Rightarrow Y_t - Y_{t-1} - Y_{t-2} +$
 $(\epsilon_{t-1} \psi_{t-1} + c)$

* Variants of ARIMA:

SARIMA

↳ seasonal



$p, d, q \neq \text{full}$

$(P, D, Q) \neq \text{seasonal}$

$(p, d, q)(P, D, Q)_S$ 12 month

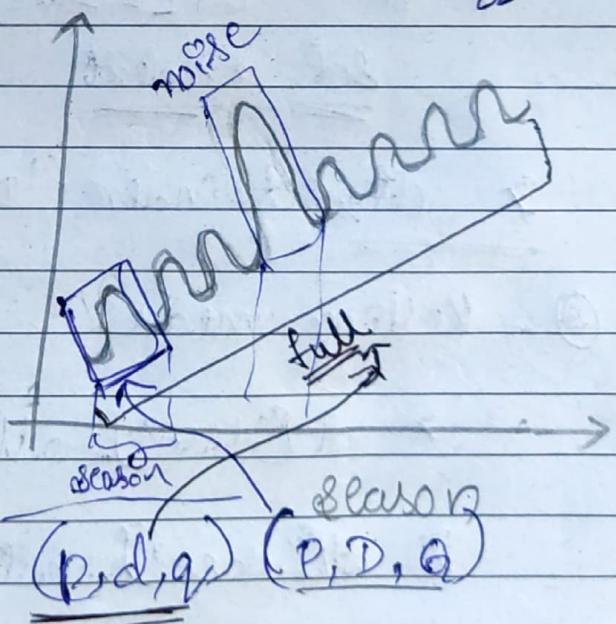
$(1, 1, 2)(1, 0, 2)_{12}$

$\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$

AR I MA / ARIMA

SARI MX

↳ exogenous
[outlier]
[extra effect]



$\Rightarrow (0, 1, 2)(1, 1, 1)$

* Implementation

steps

- ① Data ingestion.
- ② EDA.
- ③ processing of the data.
- ④ Model building.
- ⑤ Model evaluation.

① \Rightarrow Convert object to Datetime
 $\Rightarrow df['Date'] = pd.to_datetime(df['Date'])$

$\Rightarrow df = df[['Date', 'Close']]$

set index

$\Rightarrow df.set_index('Date', inplace=True)$

② \Rightarrow Rolling mean

$\Rightarrow df['Close'].rolling(12)$

$\Rightarrow df['close'].rolling(window=12).mean()$

③ To convert stationary to Non-stationary

$\Rightarrow dfc = df['close']$

$\Rightarrow dfe = dfc.diff()$

$\Rightarrow dfe.dropna()$

5

> import statsmodels as sm
> from ~~statsmodel~~ statsmodel.tsa.arima.model
import ARIMA
> from sklearn.metrics import
mean_squared_error.

> model = ARIMA(history, order=(1,1,1))
> model.fit()

> model.summary()

> model.forecast()

