

# Topics:

~ QoS

① Market Basket Analysis

② Association Rule Mining

③ Apriori Algo

④ Introduction to Time-series Analysis

✓

✓ → Forecasting

✓ → Resampling

- missing values ✓
- outliers / Anomalies

interactivity

RecSys

# Market Basket (90's)

DMart / walmart

items → Milk, bread, jam, butter --- ↗ few 100's  
of products

$$\mathcal{D} = \{ 1, 2, 3, \dots n \}$$

Transactions:

$T = \{ T_1, T_2, T_3, \dots \}$

$T_1: \{ \underline{1}, \underline{3}, 6, 8 \}$  basket

$T_2: \{ \underline{1}, \underline{3}, 7, 12 \}$

$T_3: \{ \underline{1}, \underline{7}, \underline{3}, 16 \}$

$\vdots$

patterns:

$\{ 1, 3 \}$  itemset

$\{ 1, 3 \}$  bread

Milk

$\{1\} \rightarrow \text{recommend } \{3\}$

RecSys

Task:  $\textcircled{n}$   $\mathcal{D}$ : set of all items  
 $\textcircled{m}$   $\{T\}$ : Set of transactions  
v.large  $T_i \subseteq \mathcal{D}$

$T_i \subseteq \mathcal{D}$   
find item-sels that occur very frequently in  $T_i$

find frequent itemsets

m778



(1)

for each  $T_i$ 

[incomplete]

for each  $T_j \ j \neq i$

$T_i \cap T_j$



(2)

X

K

value  
cat

itemset  
cat

↓

{i,j}  $\neq i,j$ 

{i,j,k}

;

②  
Too slow

all subsets of  $\omega = \{1, 2, \dots, n\}$



$2^n$

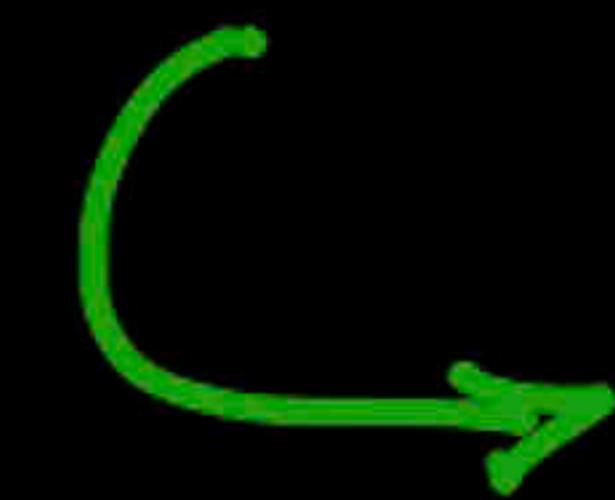
{ ... }

$\{1, 2, 4\} \rightarrow T_1, T_2, \dots, T_M$

$M = \# \text{transactions}$

Time complx:  $O(2^n \times M)$

anecdote:



5PM - 7PM

{ Beer, Diapers } ✓

90's

~~Key Idea:~~

# {2,4} < 100 (let)  
Threshold  
minimum support (c)

# {2,4,6} < 100  
C<sub>Superset of {2,4}</sub>

①  $\rightarrow \{1\}$

$\{1\}^{100}$   $\{2\}^{200} - \{3\}^{150}$

$50 < C = 100$  ✓  
min support

②  $\rightarrow \{1, 2\} \quad \{1, 3\} \quad \{2, 3\}^{60}$

~~$\{1, 2\}$~~   ~~$\{1, 3\}$~~   ~~$\{2, 3\}$~~   
 ~~$\{1, 2, 3\}$~~   $\{2, 3\}$   
 ~~$\{4\}$~~

$$\checkmark O(2^n \times m) \xrightarrow{\sim} \underline{10^b}$$

reduce this  
significantly

$$\tilde{n} = \underline{100}$$

$$2^{100}$$

$$\left. \begin{aligned} 2^{10} &\approx 10^3 \\ 2^{20} &\approx 10^6 \\ 2^{30} &\approx 10^9 \\ 2^{40} &\approx 10^{12} \end{aligned} \right\}$$

not used in e-commerce

↑  
Apxion algorithm (1994-95)

n increases (~1000)  
Very costly

Modifications → FP-growth ↗  
Specialized DS: Tries ✓

frequent itemset mining  
Case  $n = |\mathcal{D}|$  is small



Other applications? (beyond retail)

→ bio-informatics  
 $\{C_1, C_3\}$

$\{\overbrace{AGTC}, \overbrace{ATTG}\}$

→ medicine  
 $\{\underline{M}_1, \underline{M}_2, \underline{M}_3\}$

→ web usage mining }  
 $\{ \tilde{w_1}, \tilde{w_2}, \tilde{w_3} \}$

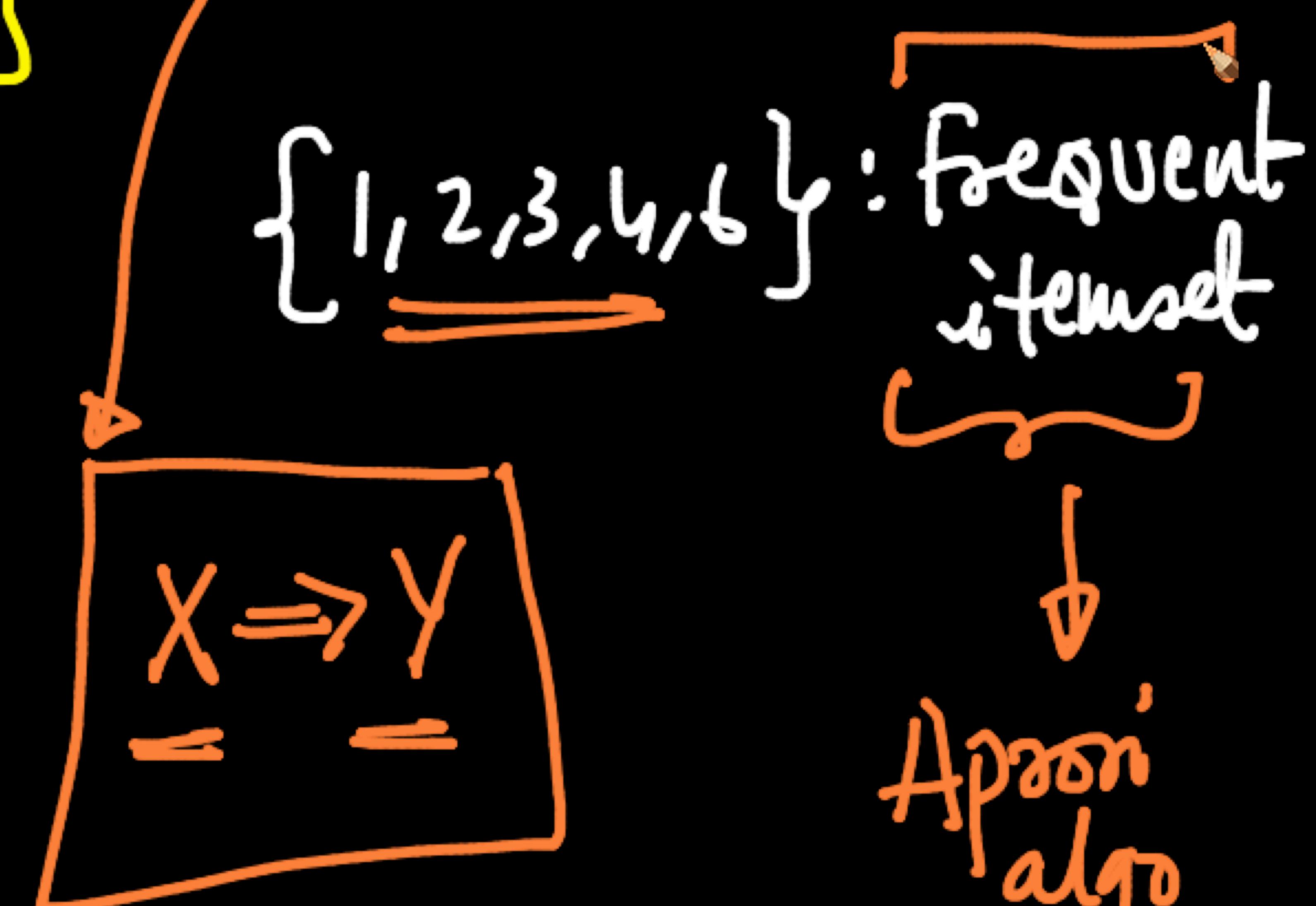
→ Similar words  
 $\{ \tilde{\underline{\underline{w_1}}}, \tilde{w_2}, \tilde{w_3} \}$  You need  
P  
T

# Association Rule Mining

Set of items  
 $\hookrightarrow D = \{1, 2, 3, \dots, n\}$

$X = \{1, 2, 3\}$

let  
 $Y = \{4, 6\}$



e.g.:  $\{ \text{beer} \} \xrightarrow{X} \{ \text{diapers} \}$

$\{ \text{milk, bread} \} \Rightarrow \{ \text{jam, eggs} \}$

Support( $x$ ) =  $\frac{\# \text{Transactions that contain } x}{\text{Total transactions}}$

$$= \frac{100}{10,000} = P(\underline{x} \text{ in } T)$$

$$\text{Confidence} \left( \overbrace{X \Rightarrow Y}^{\text{if } X} \right) = \frac{\# (X \text{ and } Y)}{\# X}$$

$= P(Y|X)$   
(logically)

90% +  
X ∪ Y is a frequent item set

$$\underbrace{\{ \text{lift}(\tilde{x} \Rightarrow \tilde{y}) \}}_{\sim} = \frac{\text{Support}(x \wedge y)}{\text{Supp}(x) \times \text{Supp}(y)}$$

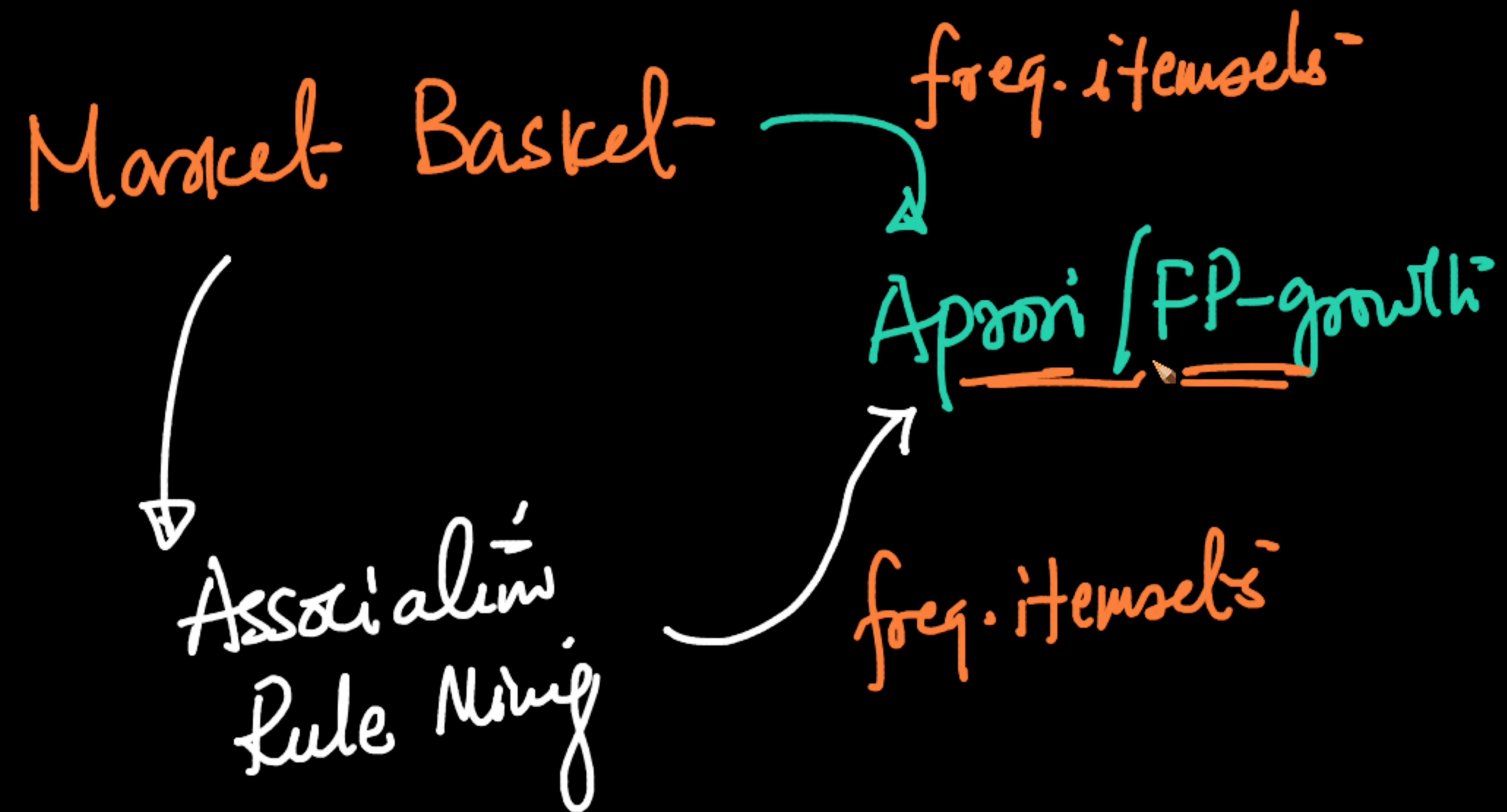
$$= \frac{P(\underline{x} \text{ and } \underline{y})}{P(x) P(y)} \quad \text{intuitively}$$

{  $x \wedge y$  and up

$P(X \text{ and } Y) = P(X) P(Y) \text{ if } X \text{ & } Y \text{ are indep}$

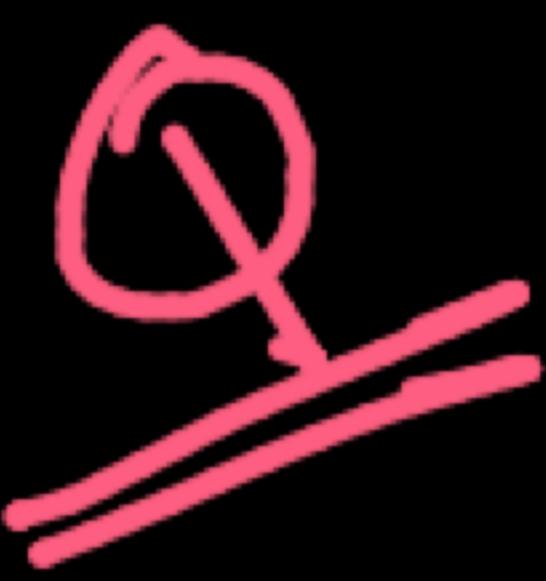
$\text{lif} (x \Rightarrow y) = 1 \text{ if } X \text{ & } Y \text{ are indep}$

7.1



{milk, bread} → frequent itemset  
X Y  
(Apaoni)

[  
Supp(x)  
Supp(y)  
Conf( $x \Rightarrow y$ ) ✓  
]  
lift( $x \Rightarrow y$ )



View → purchased

[MF]

e-commerce!

✓ ✓ ✓  
l<sub>1</sub>, l<sub>2</sub>, l<sub>3</sub> ... l<sub>10</sub>  
V V V  
B

Naive Bayes

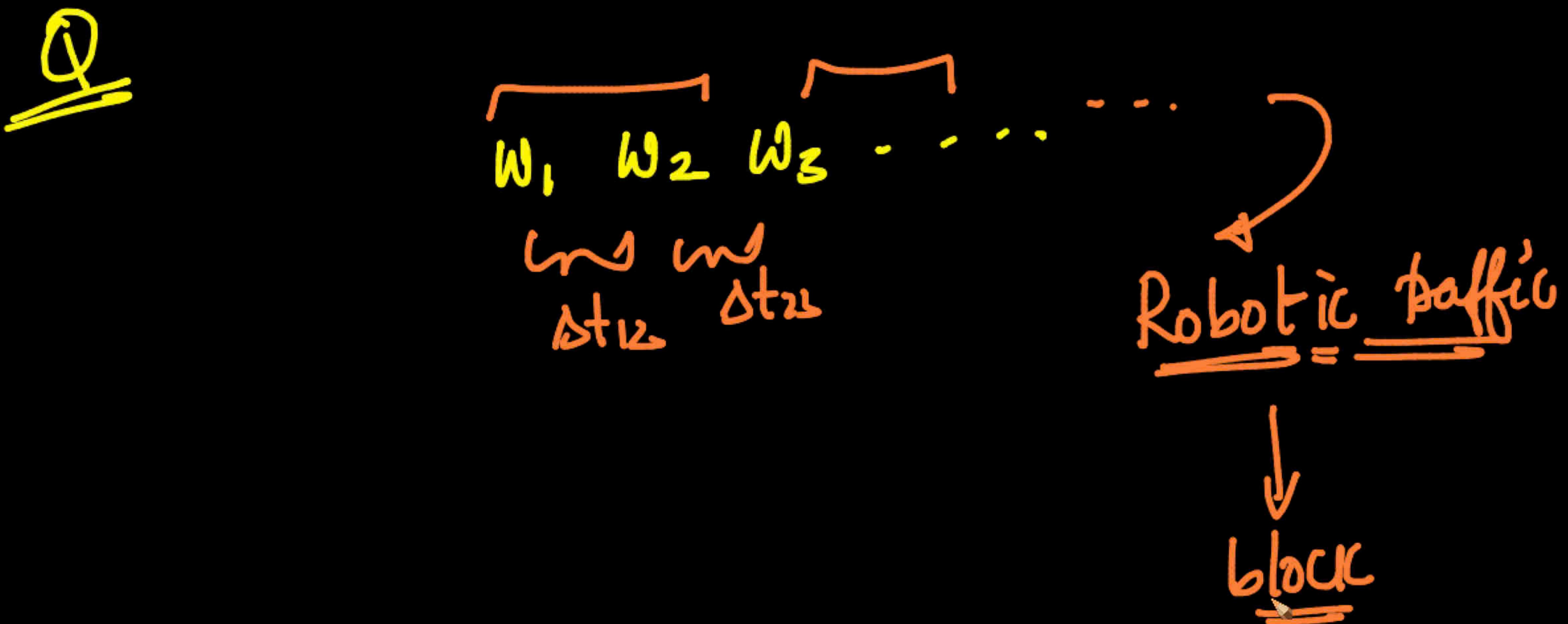
$$P(\underline{l_{10}} | l_1, l_2, \dots)$$

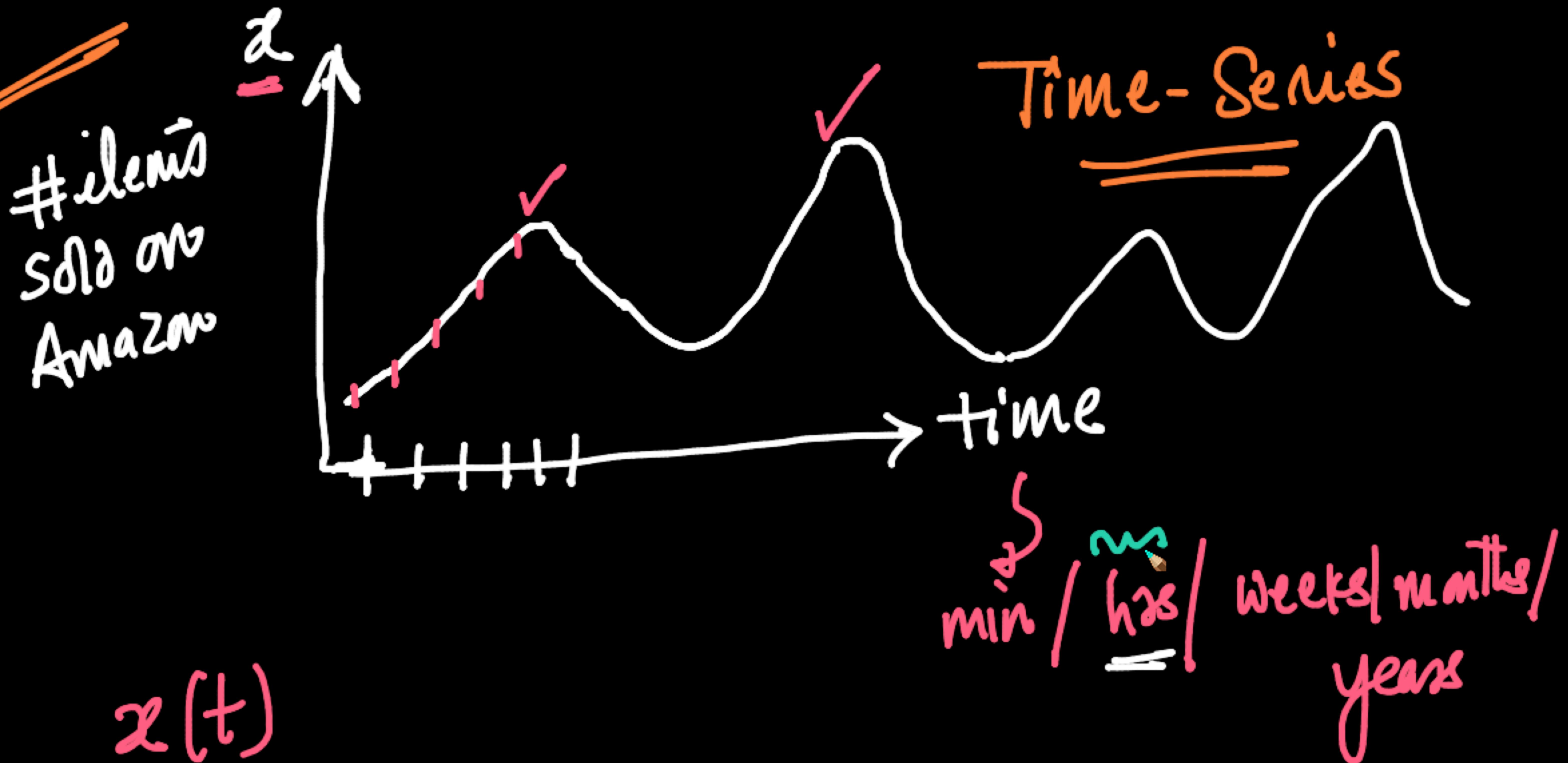


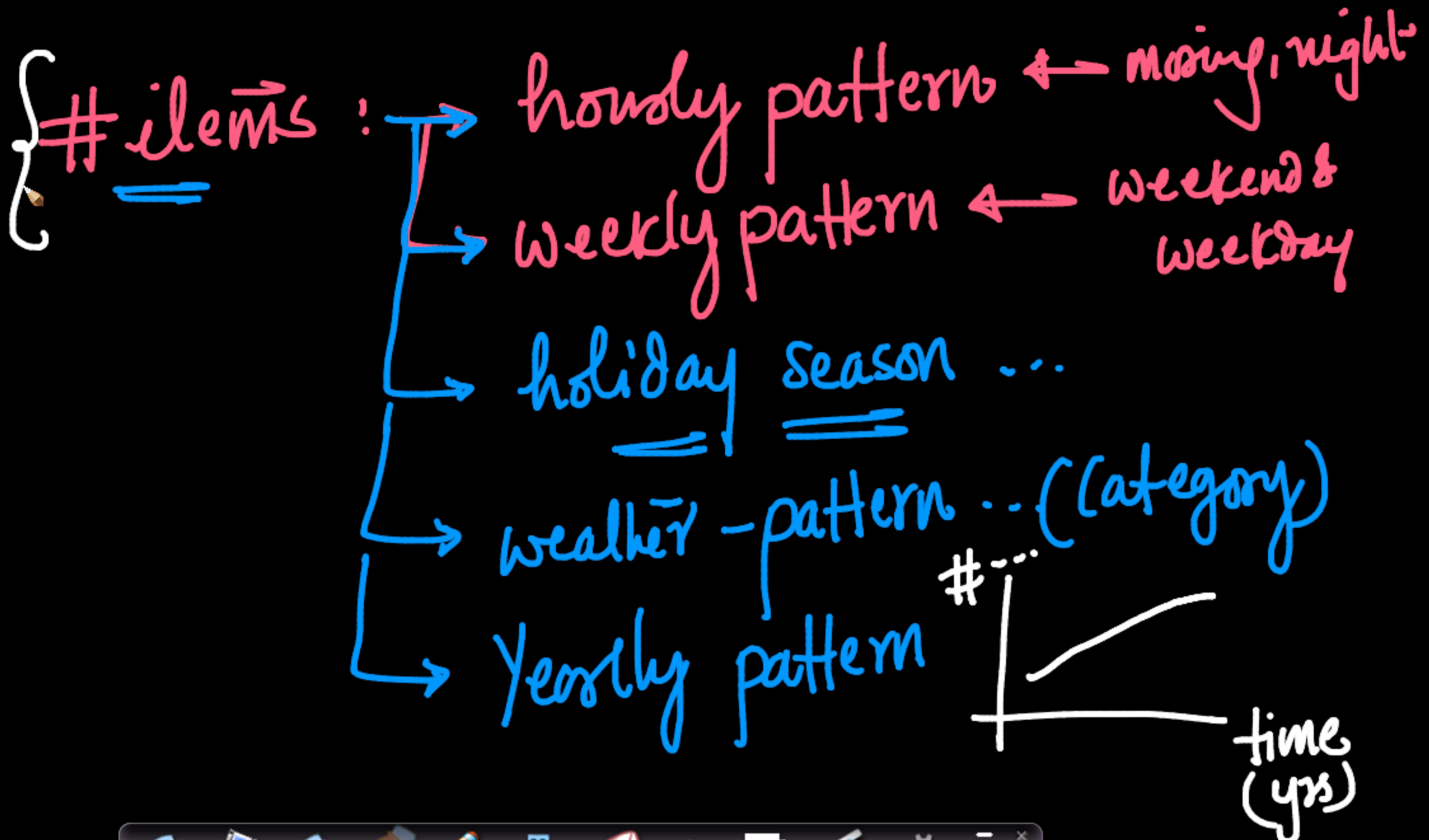
i<sub>1</sub>, i<sub>2</sub>, i<sub>3</sub>, i<sub>4</sub>.  
IV headphones  
last k-items

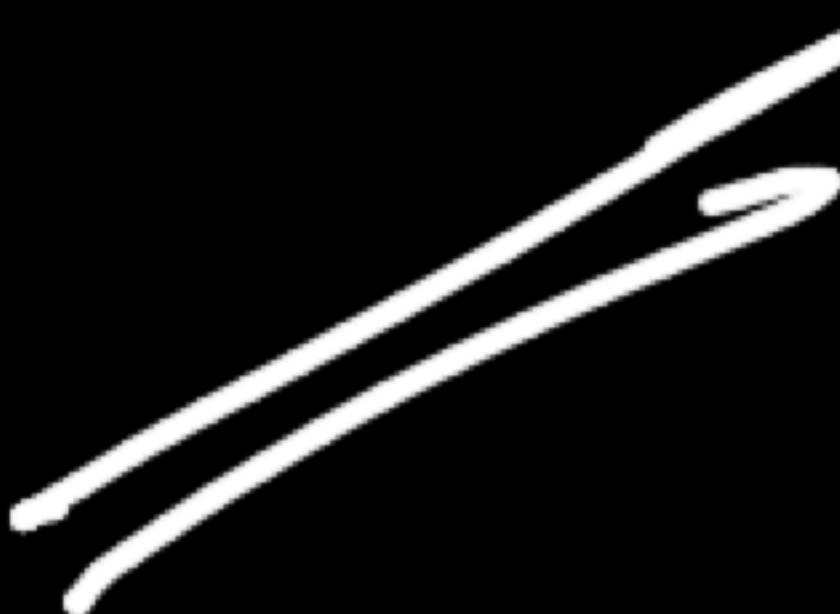
→ all similar items  
(MF or item)  
DL











# Forecasting

Task:

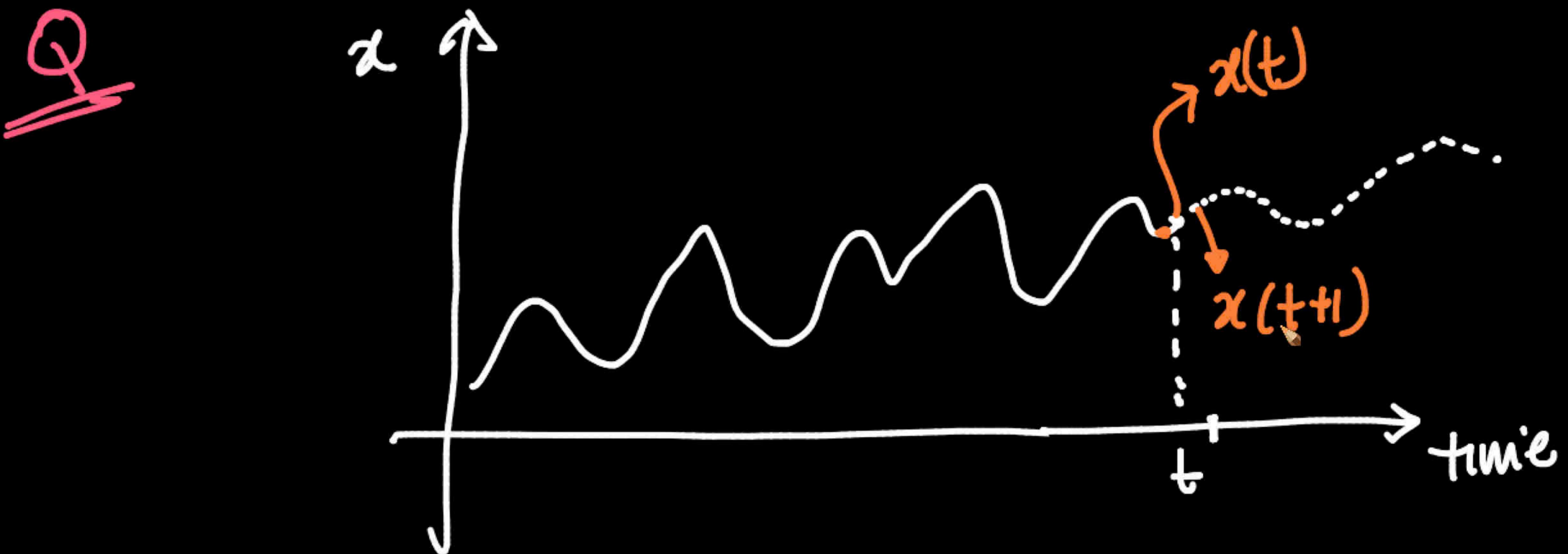
Given

$x_1 \ x_2 \ \dots \ x_{t-1} \ x_t$

historical  
data

Predict

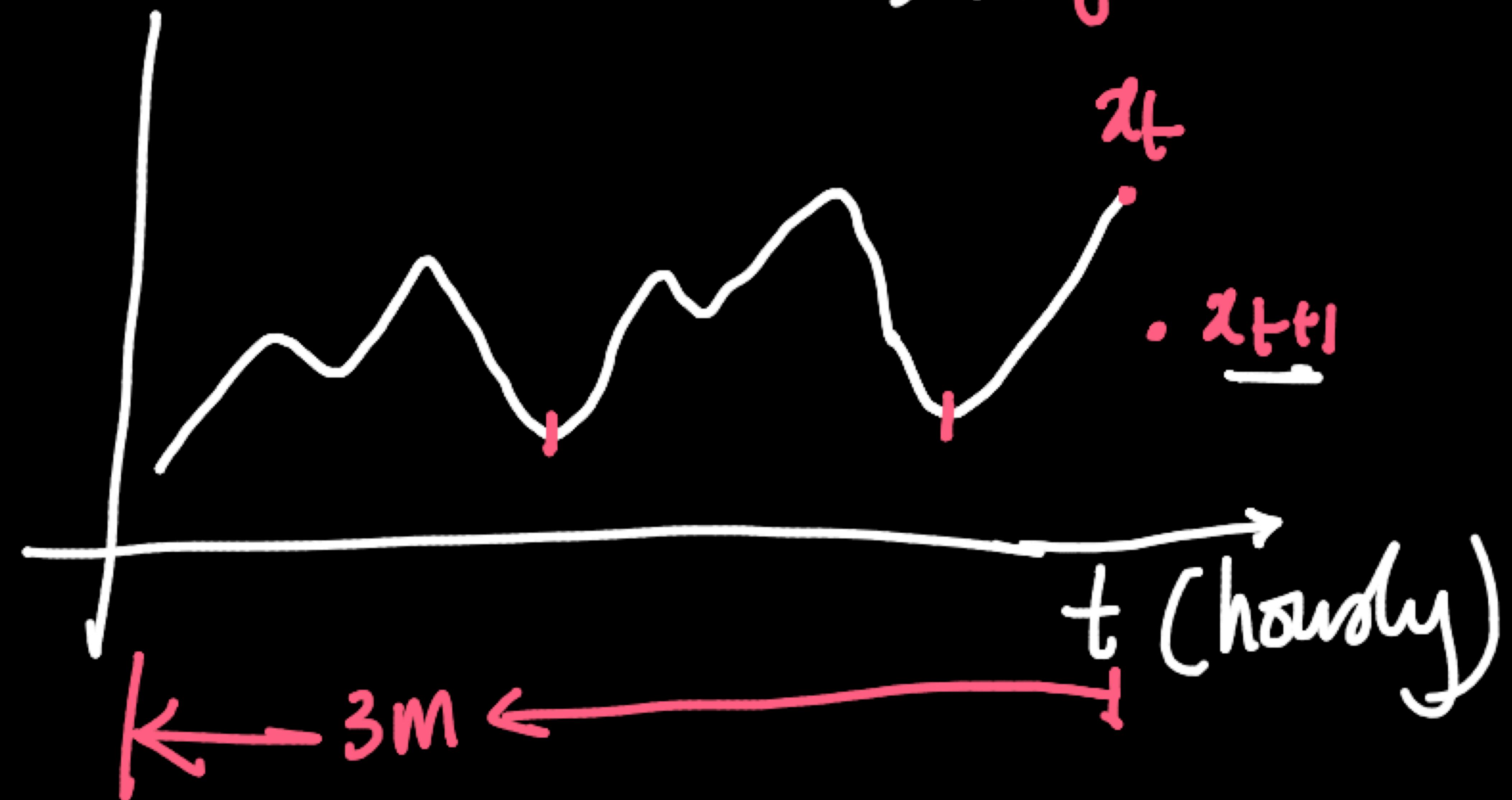
$x_{t+1} \ x_{t+2} \ \dots$



i

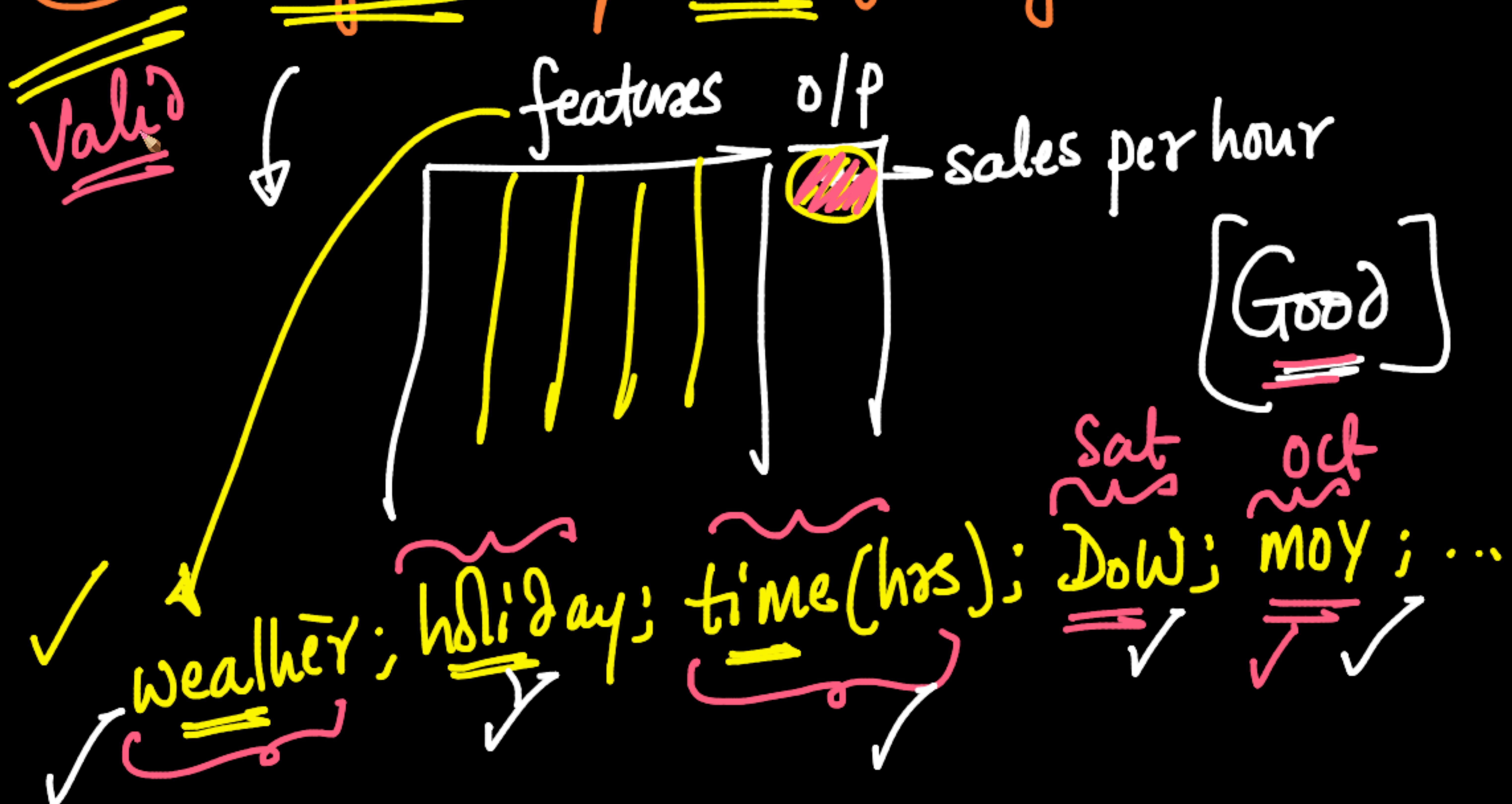
average of last ~~3 months~~  
~~daily~~

may not work



2

# regression / curve fitting



flaw:

may not capture trend over years

sln:

feature: ✓

2012 → 0

2013 → 1

2014 → 2

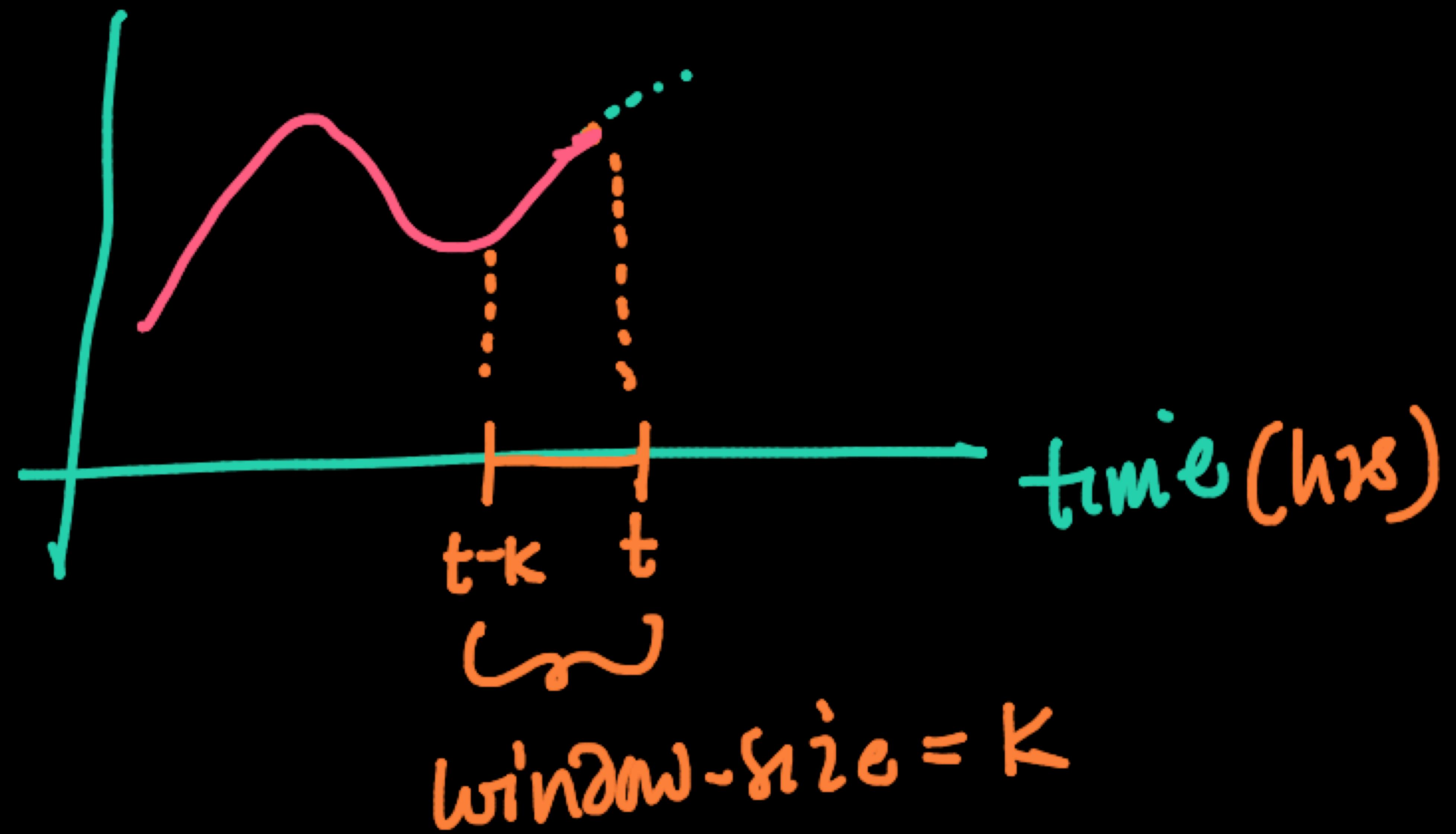
⋮ ↗ Year of operation



2012

(3)

Weighted moving average  $\rightarrow$  classical TS



$$\hat{x}_{t+1} = \frac{\alpha_0 \hat{x}_t + \alpha_1 \hat{x}_{t-1} + \alpha_2 \hat{x}_{t-2} + \dots + \alpha_k \hat{x}_{t-k}}{\alpha_0 + \alpha_1 + \alpha_2 + \dots + \alpha_k}$$

Weighted avg

$$= \frac{\sum_{i=0}^k \alpha_i x_{t-i}}{\sum_{i=0}^k \alpha_i}$$

logically ✓  $\underline{\alpha_0} > \alpha_1 > \alpha_2 \dots > \underline{\alpha_k}$

$$x_{t-k} \dots x_t \rightarrow \hat{x}_{t+1}$$

$$x_{t-k+1} \rightarrow \hat{x}_{t+1}$$
$$\vdots$$
$$\hat{x}_{t+2}$$

Type A regression with manual weights

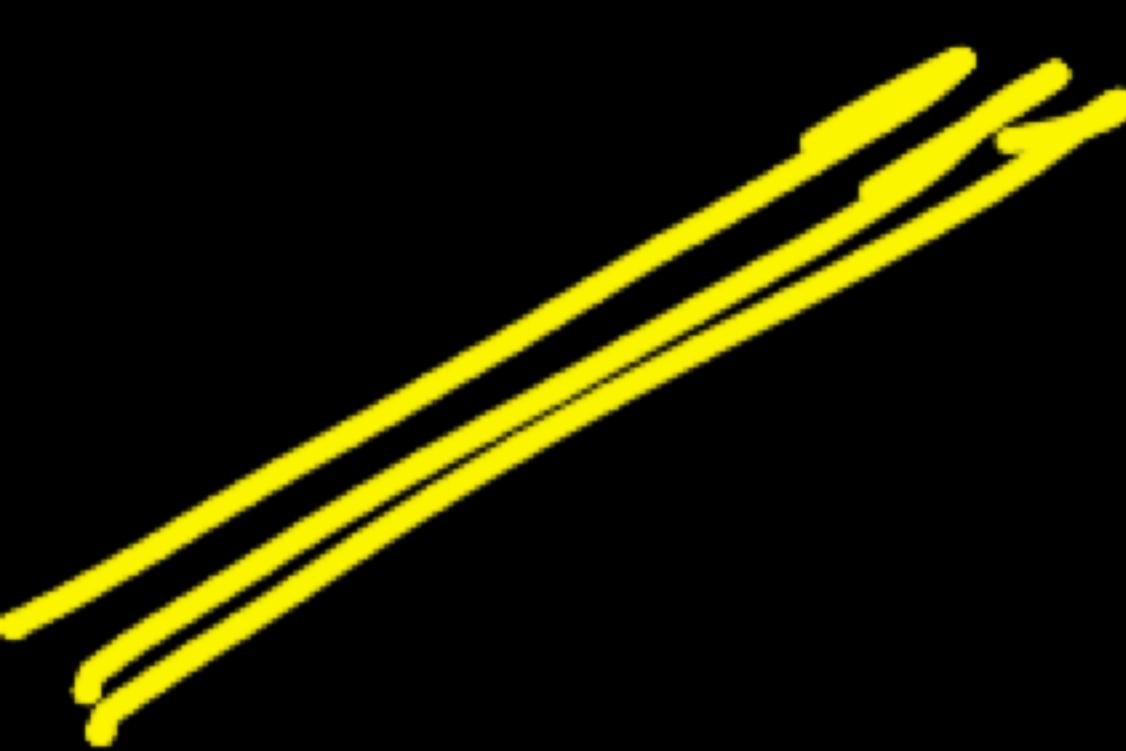
$$x_{t+1} = \frac{\sum_{i=0}^k \alpha_i^* x_{t-i}}{\sum_{i=0}^k \alpha_i^*}$$

wMA with Computed weights



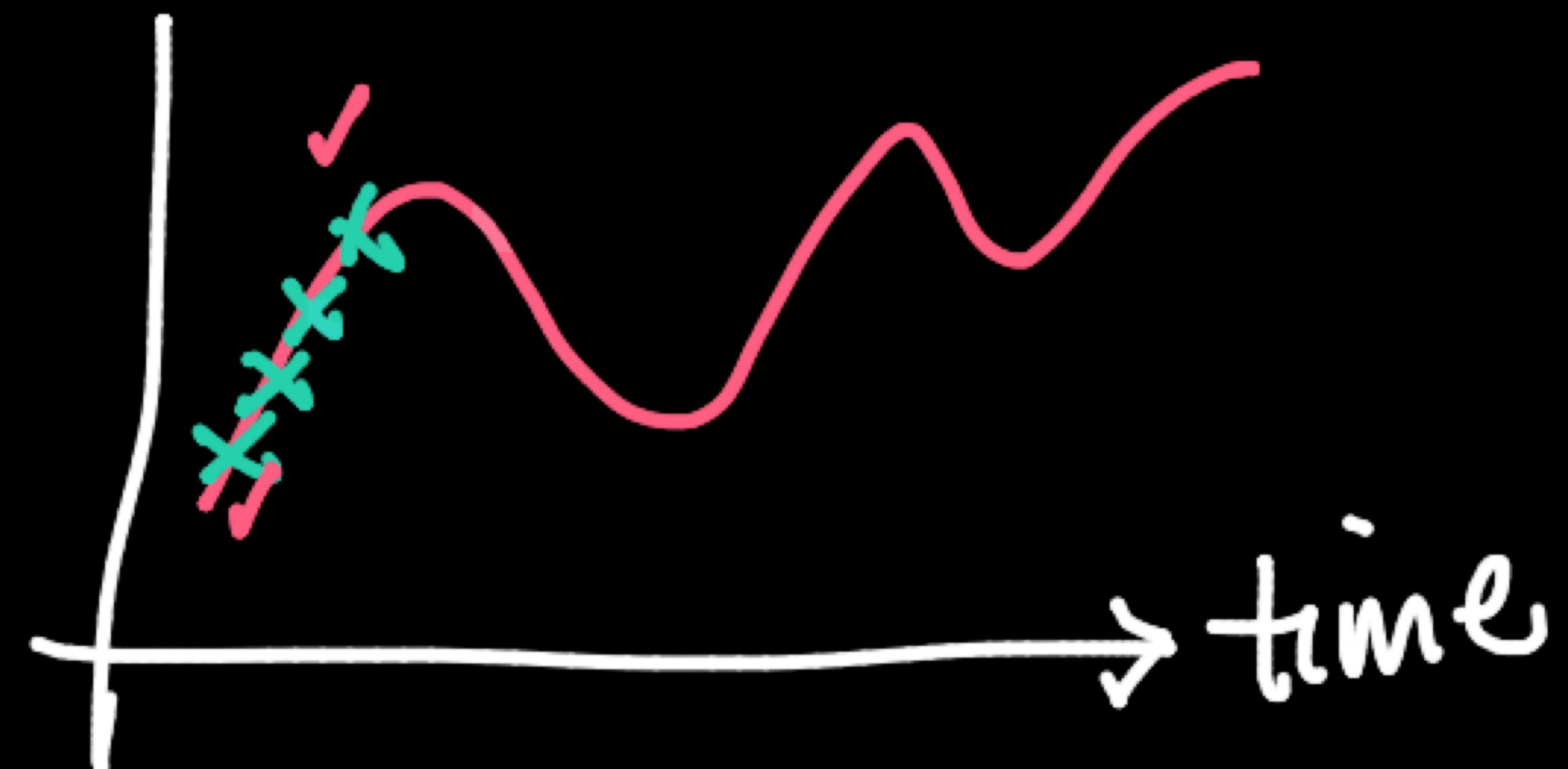
$d_i$ : weights

(Linear-reg)

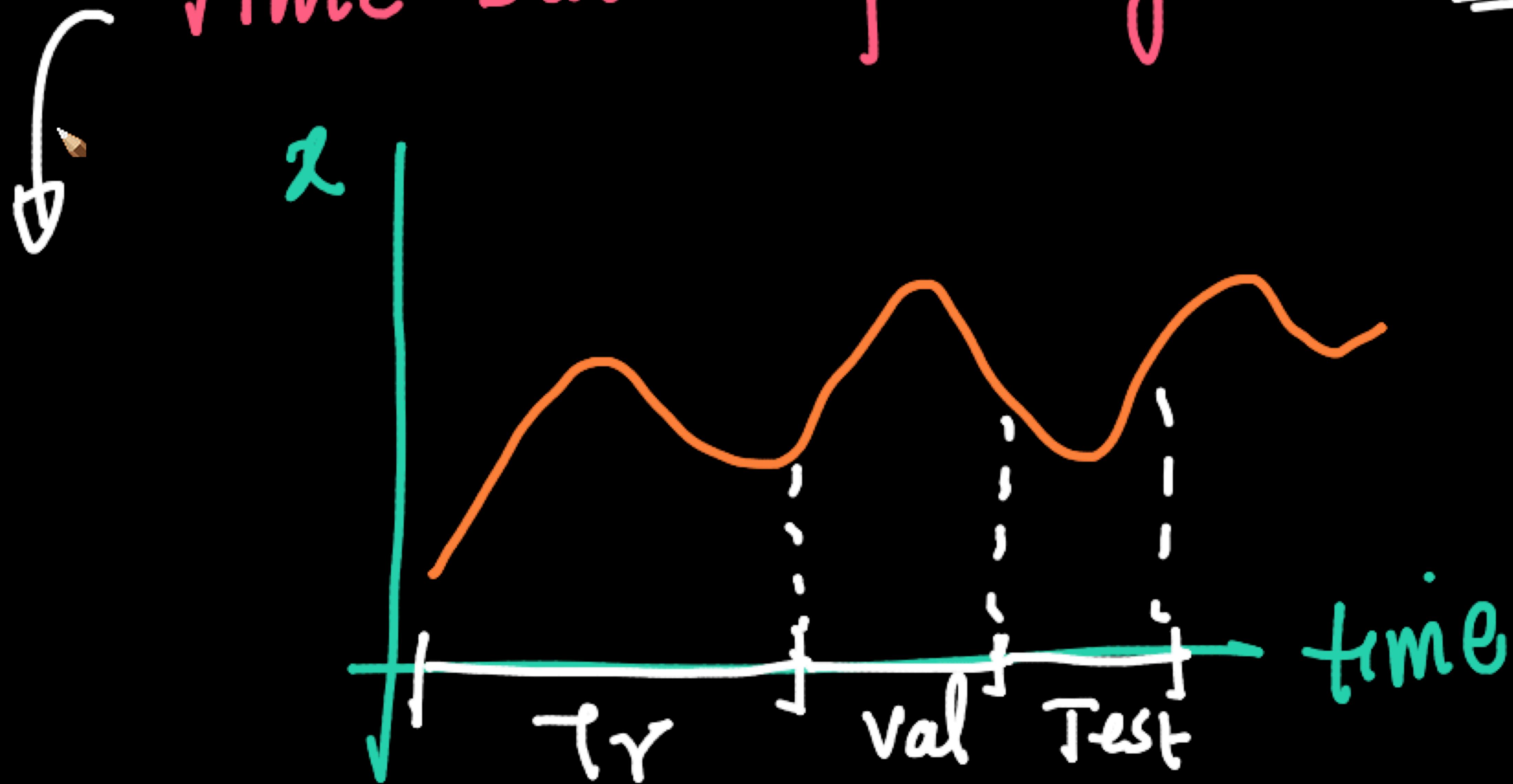


Split your data  
{ Train - Val - Test}

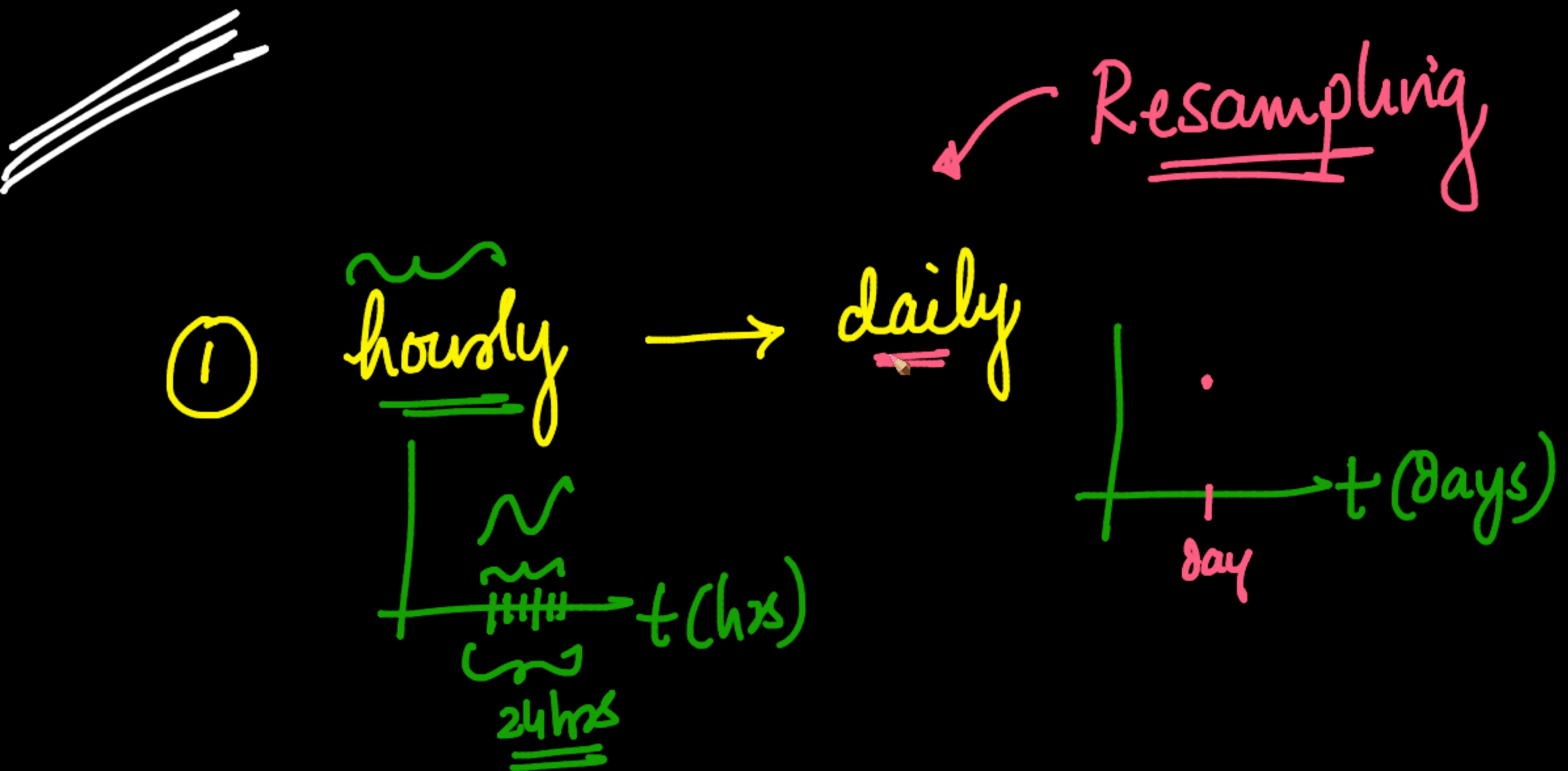
random-splitting  
random-splitting X



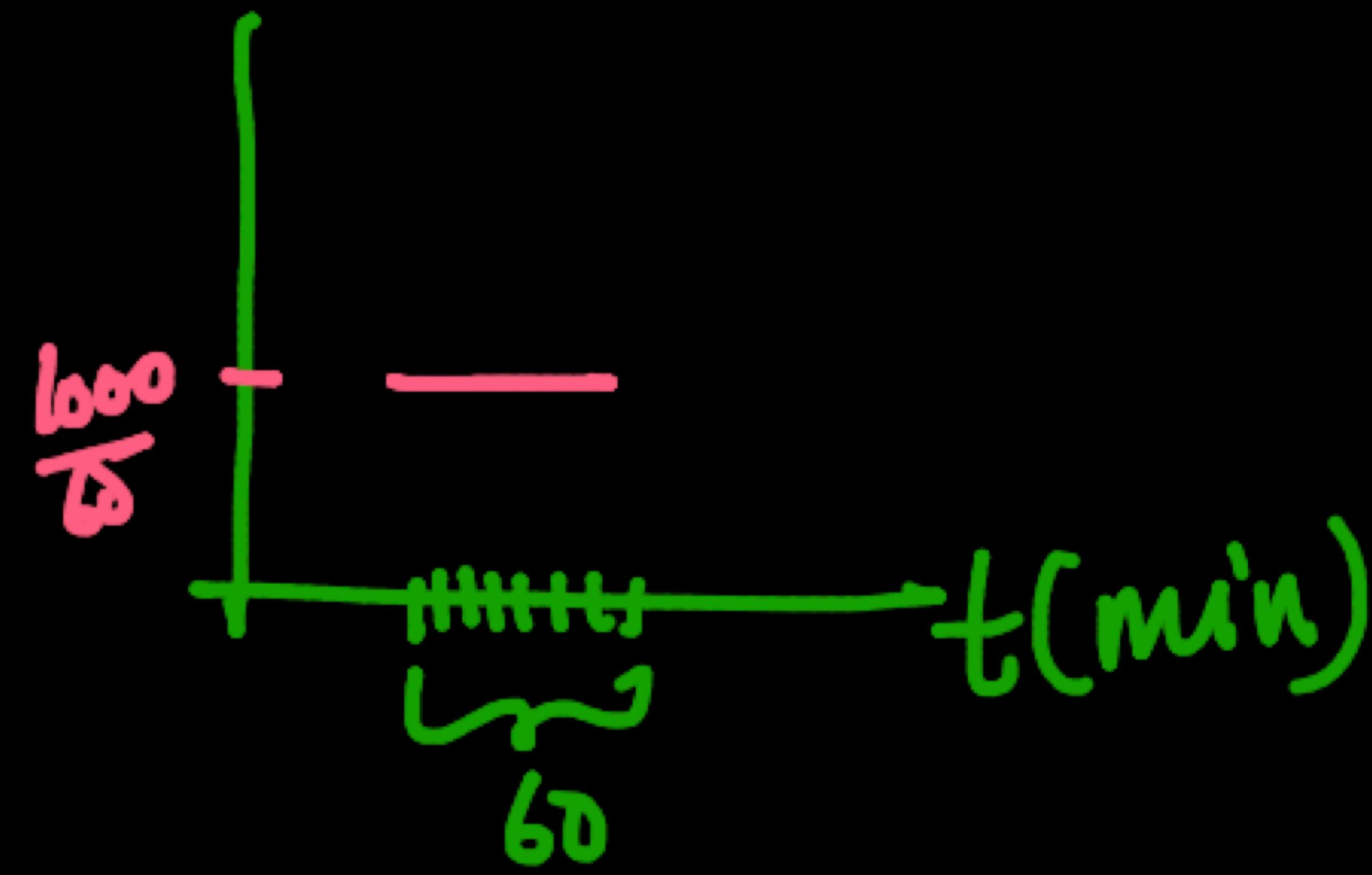
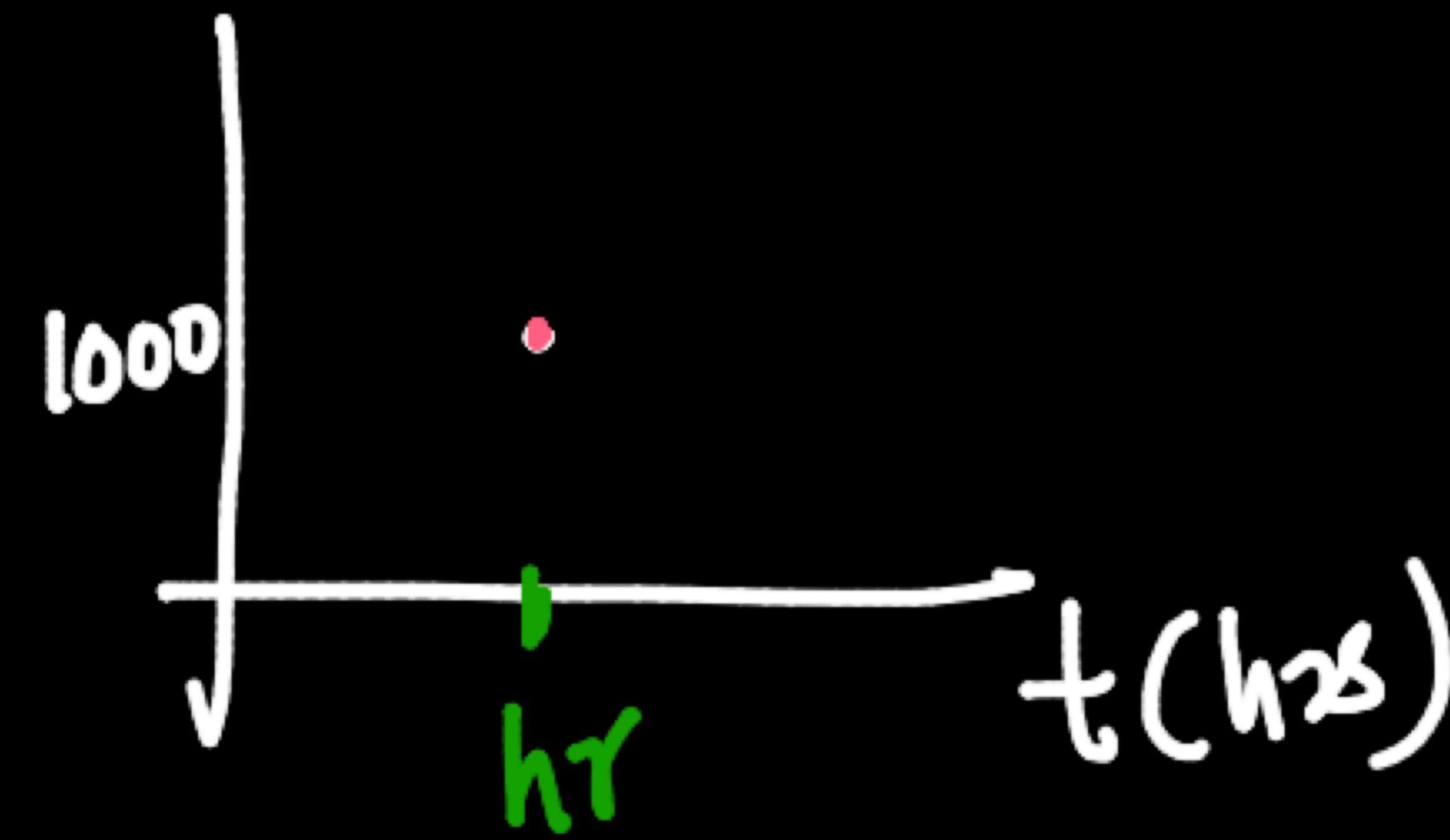
# Time-based Splitting



wMA: hyperparameter

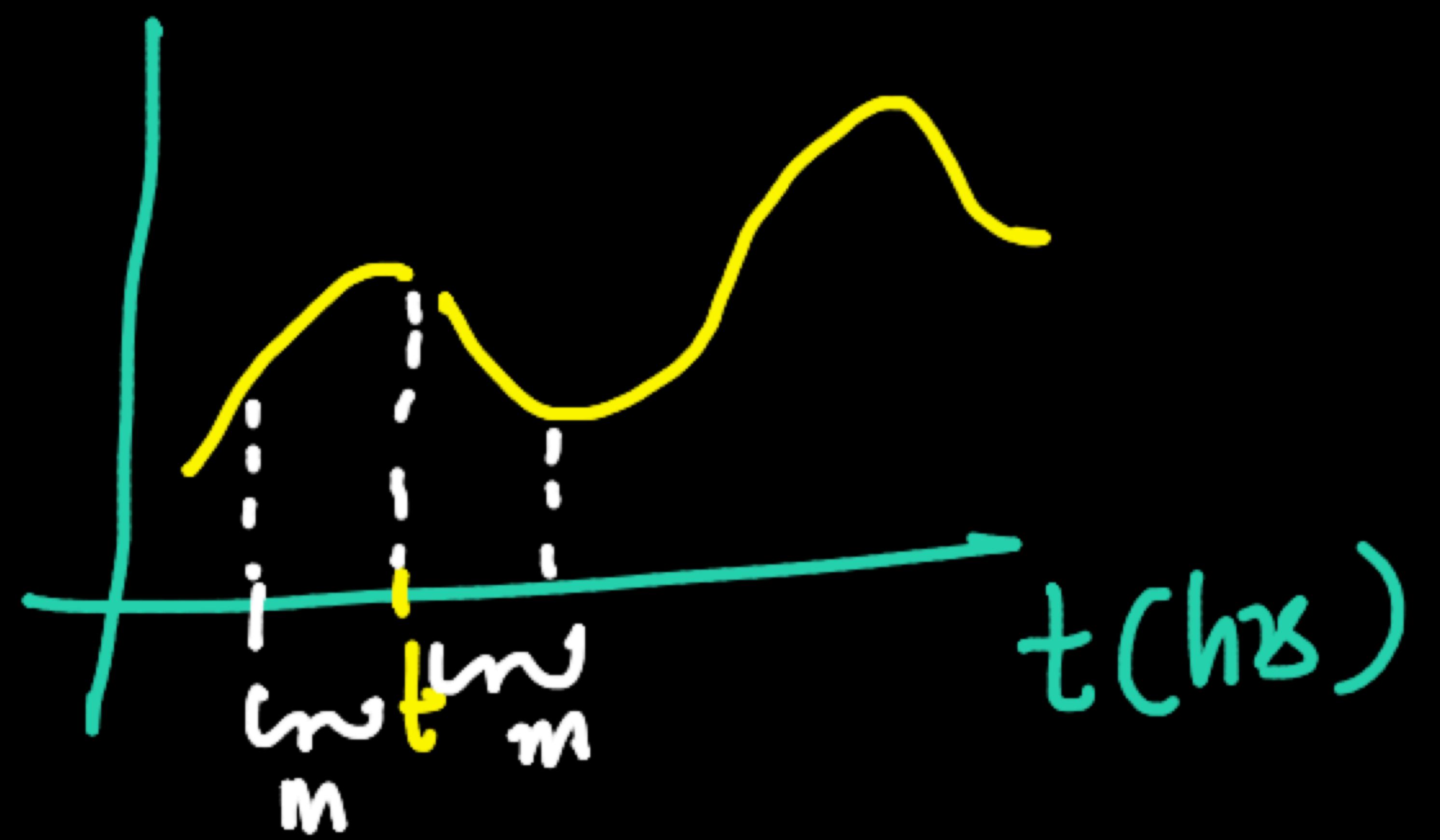


② hourly → minutely



assume uniform spread  
across each min

Missing - data  
Impulse



→ Cannot drop

$$\rightarrow x_t = x_{t-1} \quad ①$$

②

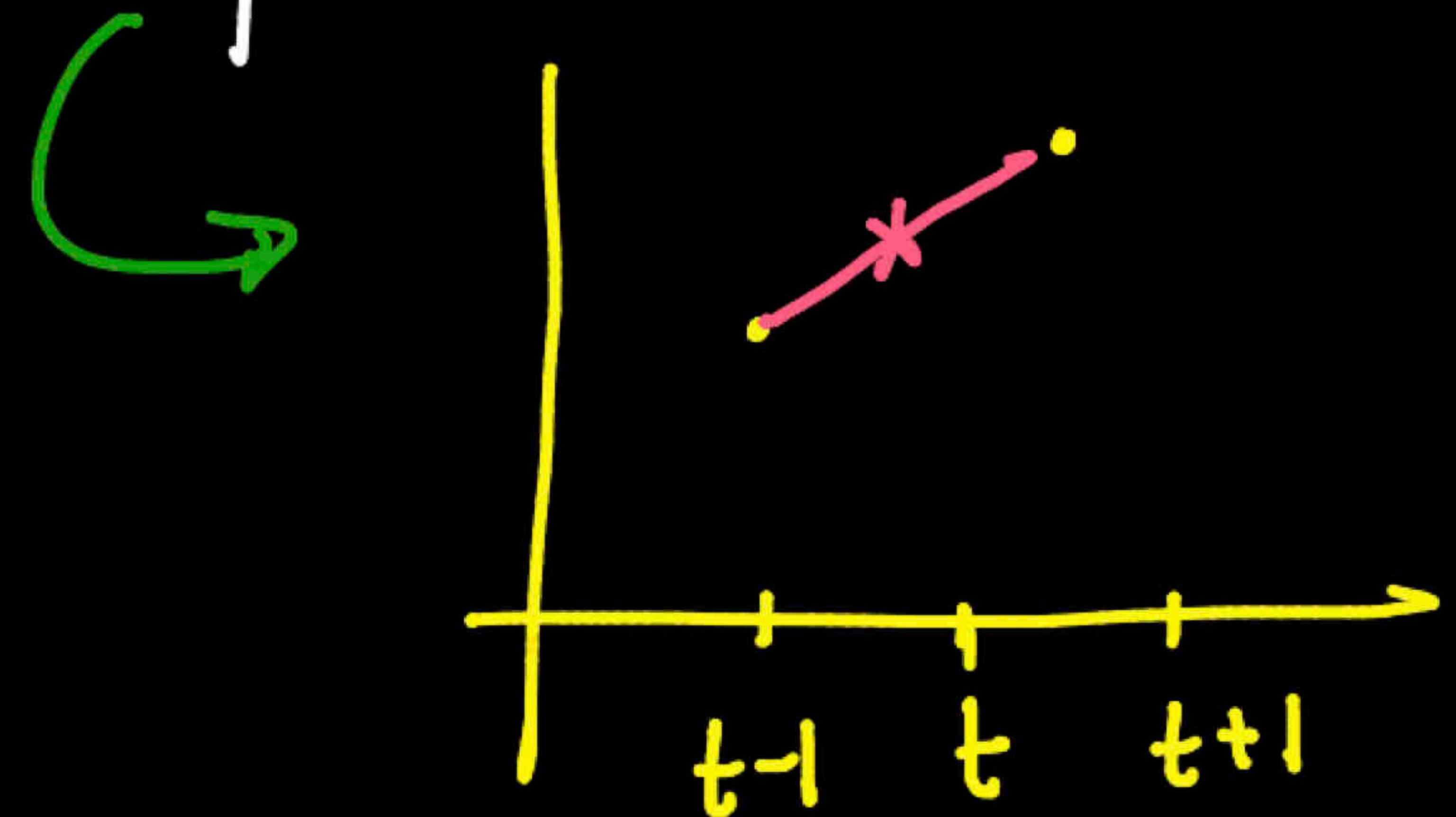
$$\checkmark x_t = \frac{x_{t-1} + x_{t+1}}{2}$$

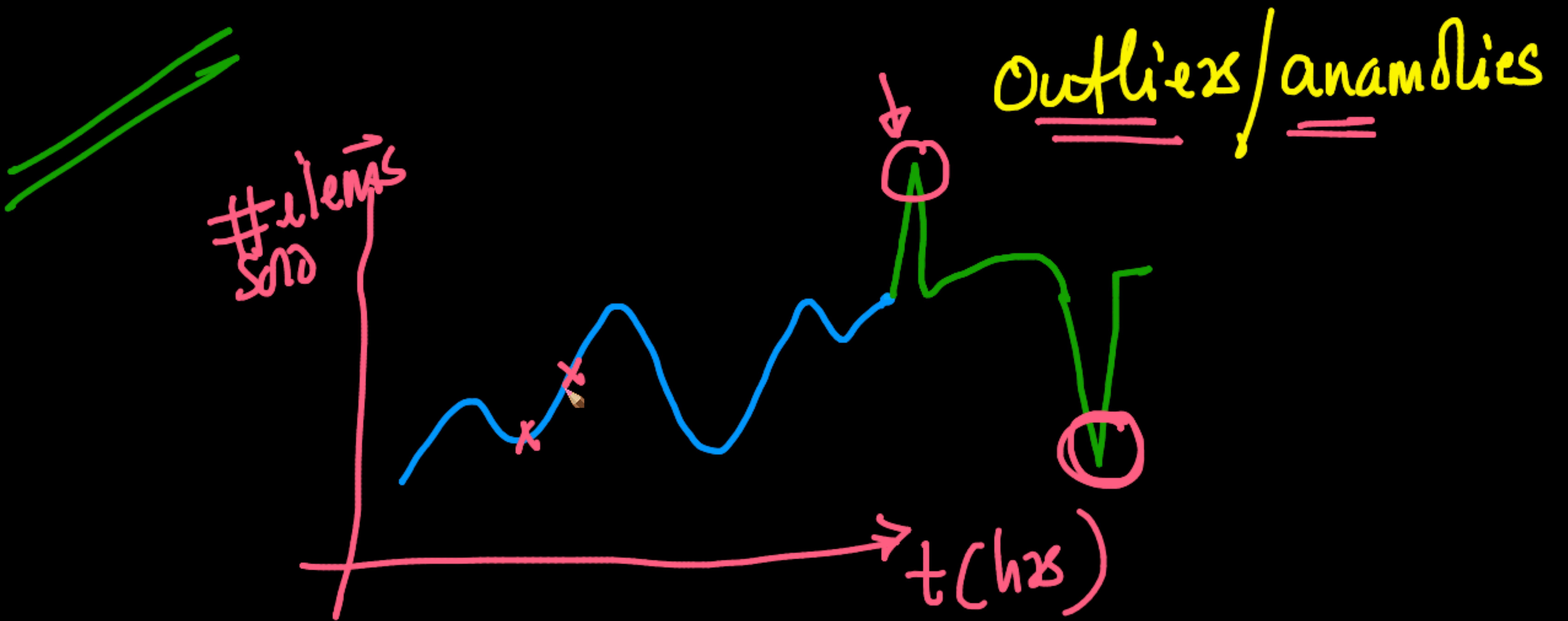
✓ ③ Centered MA

$$x_t = \frac{1}{2^M} \sum_{i=t-M}^{t+M} x_i \cdot 1$$

generalization  
of looking back  
& forward

interpolation:  $x_t = \frac{x_{t-1} + x_{t+1}}{2}$   $\approx$  center MA 4 width = 1







$$\frac{x_t - x_{t-1}}{x_{t-1}} = \Delta_t^{\text{rel}}$$

relative change

$$\Delta_1 \ \Delta_2 \dots \ \Delta_t \rightarrow \overbrace{\quad}^{\text{IQR}}$$

(2)

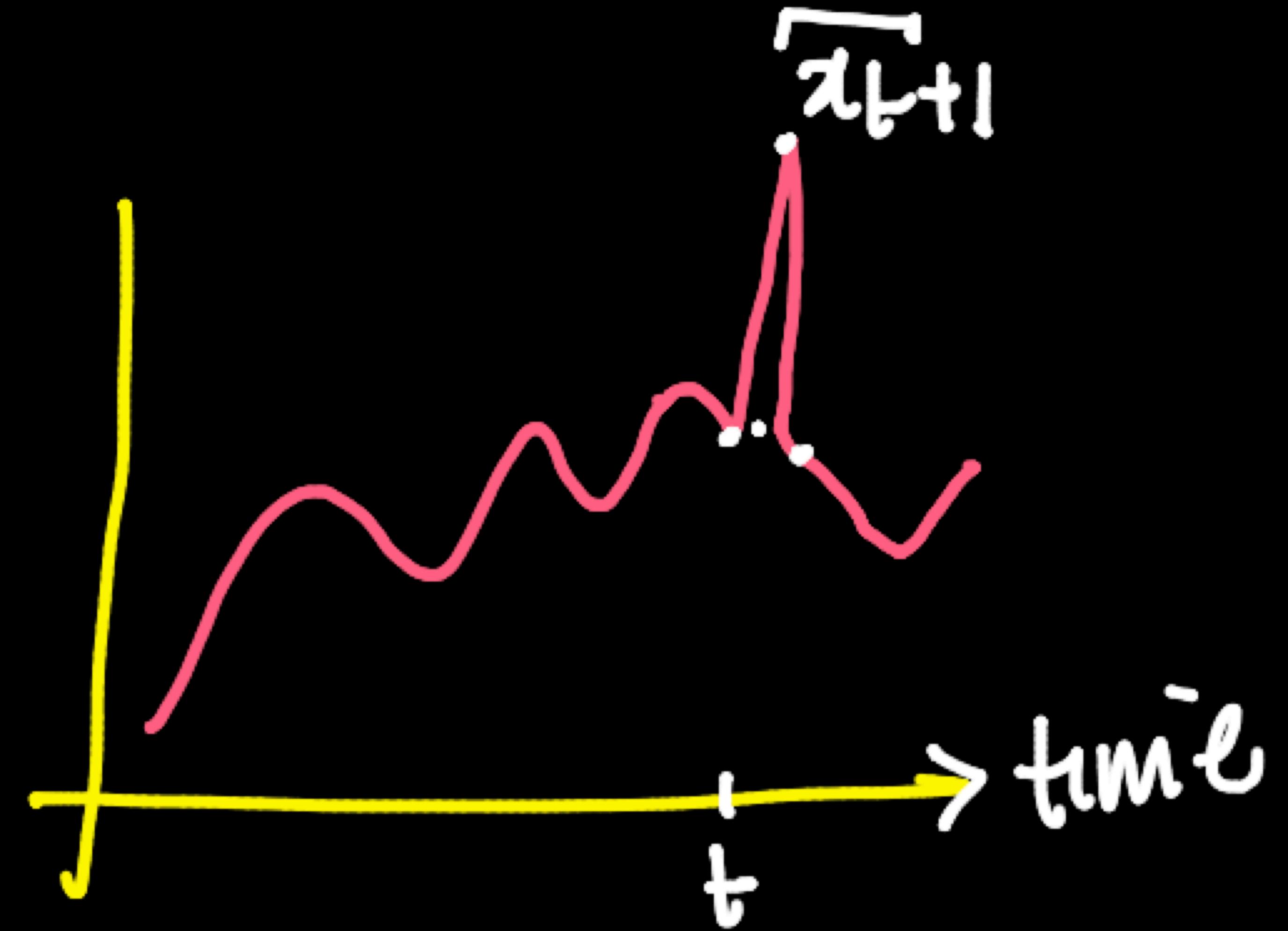
iForest /one class SUM /DBSCAN

↳ Tabular data



Time Series  
Sequence / series

(3)



assumption

- historical data has no outliers
- pattern -

$\hat{x}_{t+1}$  : regression ✓  
~~wMA~~ ...  
DL model

✓ 
$$\frac{\hat{x}_{t+1} - \hat{x}_{t+1}}{\hat{x}_{t+1}} = \delta_t$$
 is high/low  
↳ 200%

Classical ML

Next-class

→ Code

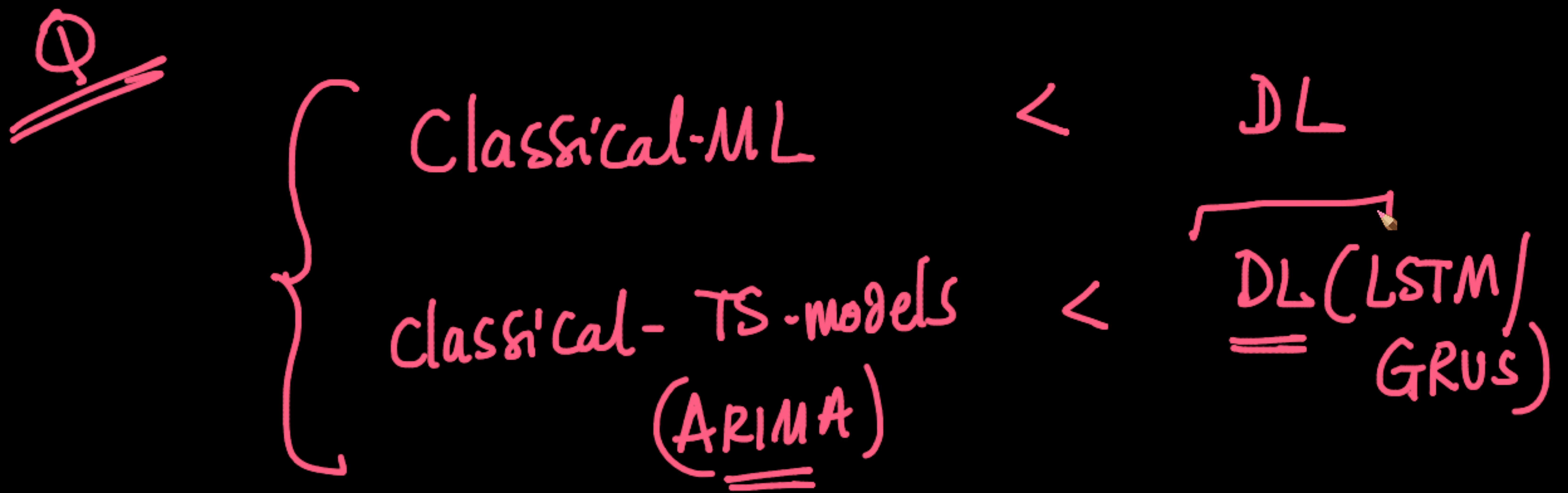
— Trend, Seasonality

⋮

— ARIMA, SARIMA

...

time  
series



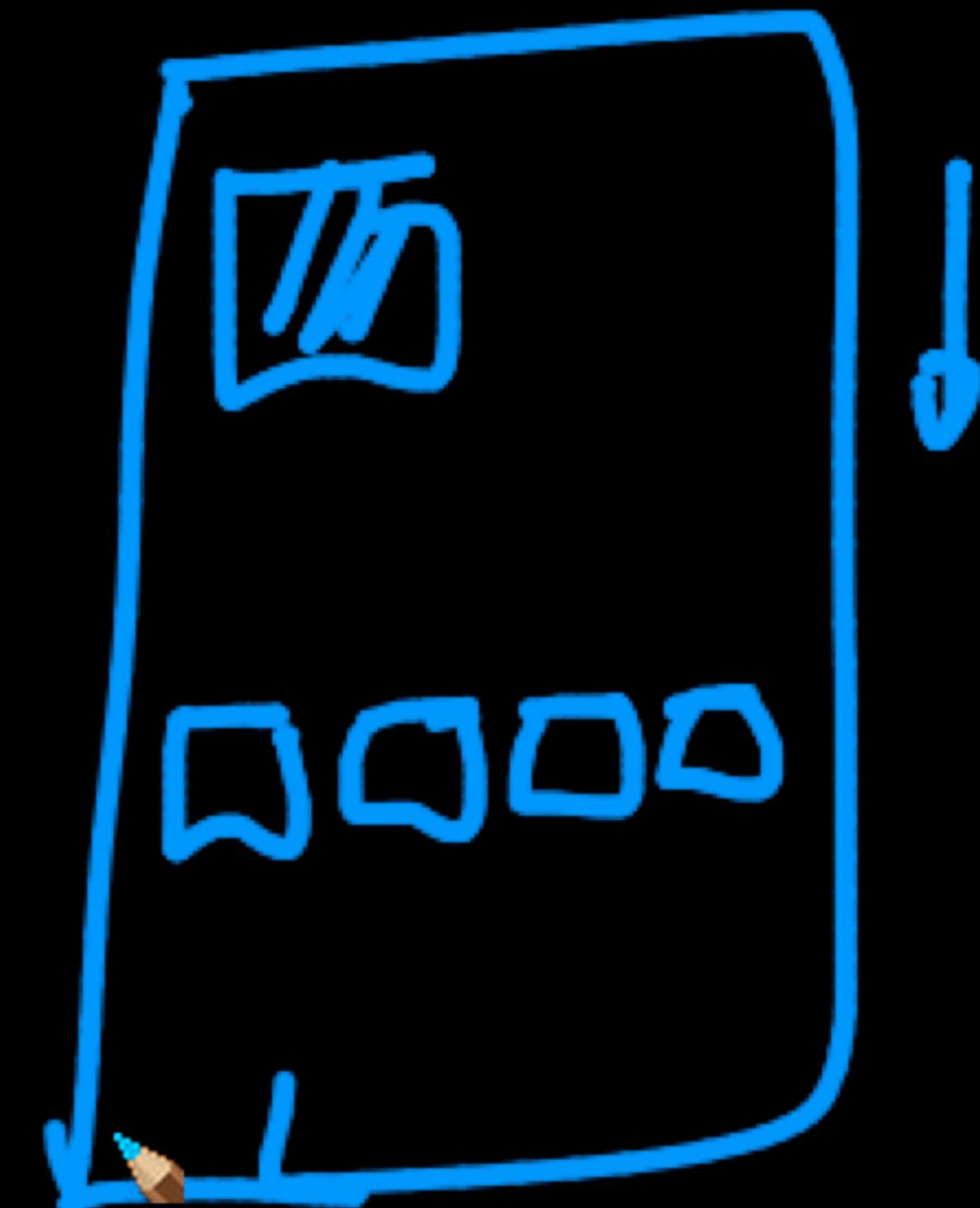


RecSys →

RMSE; ~~log-loss~~; ...

↳ click-rate

→ click-through-rate





Gmail Images



Google

Search Google or type a URL



Colaboratory



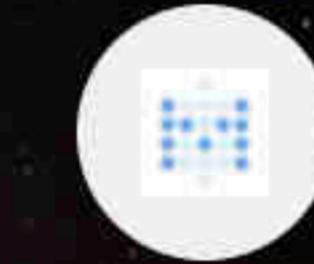
YouTube



My Drive



InterviewBit S...



Learning



GitHub



Reduce the fil...



Scaler Academ...



Amazon.in



Add shortcut