

✓ - C.I of metrics - Non Gaussian vs Gaussian - Param vs Non param

Remedial Session: Statistics (3 more topics)

Today

Regular Session:

Linear Algebra &

Coordinate Geometry

Geometric

intuitive
Applicable

Next Sessions
(2 more)

ML concepts

OPS:
≠

each section → chat (Q&A)

general / broad out of the context → Q&A
(30 mins)

Parametric

vs

Non-parametric tests

makes assumptions about the 'distr' of the data

ANOVA: each group's observations are sampled from a Gaussian distr

F-statistic

χ^2 -test

any disb. assumptions about the data

(Q)

 χ^2 -test

→ NO assumptions about underlying data

COMMON
-confusion χ^2 test
statistic

$$\chi^2 \text{ statistic} = \sum_{i=1}^K \frac{(O_i - E_i)^2}{E_i}$$

 χ^2 -distribution

P-value

Why is it non-parametric





✓ t-test:

Parametric → { t-test statis^{tic} ~ t-distr
DOF

2-groups

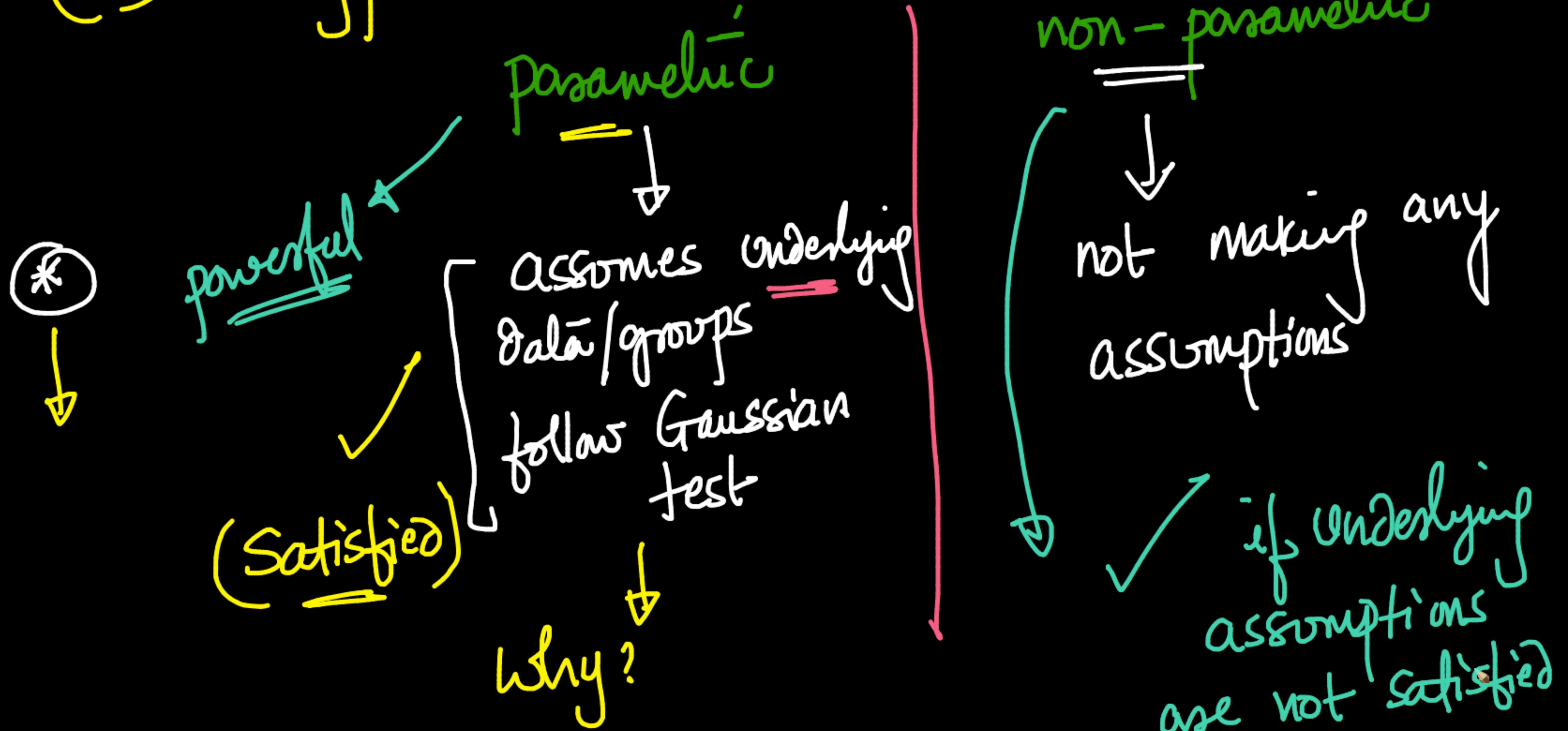
each group's underlying data is gaussian
distr

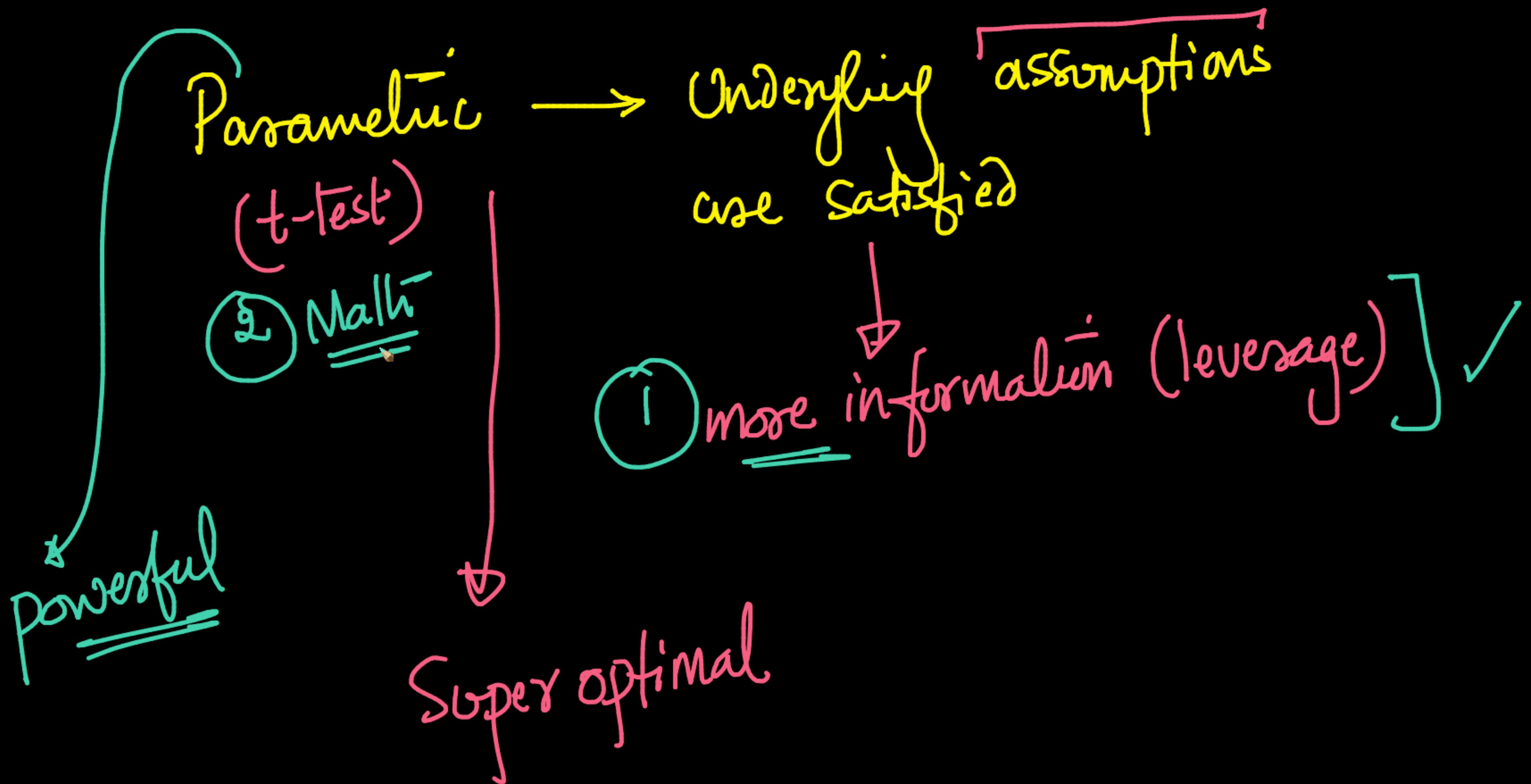
Math & proof

Theory vs Practice

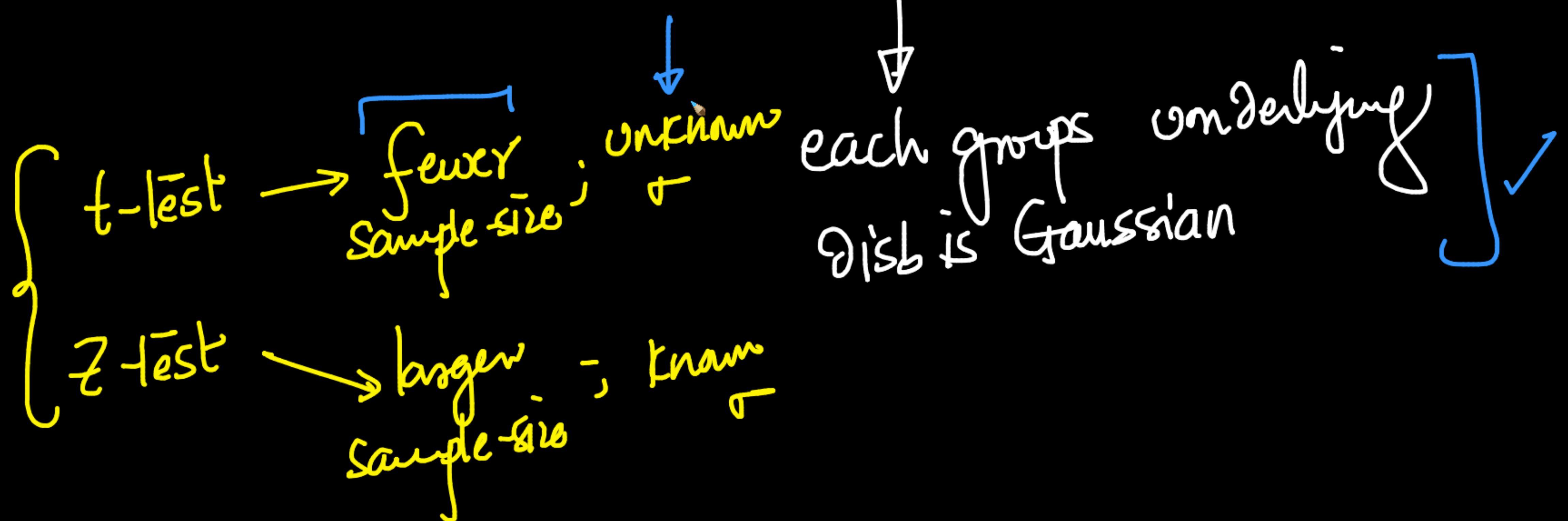
✓ t -statistic $\stackrel{\text{test.}}{\sim}$ t -dist
[Underlying data \sim Gaussian (ideal-case)]

(Q) hyp-test : difference in Means





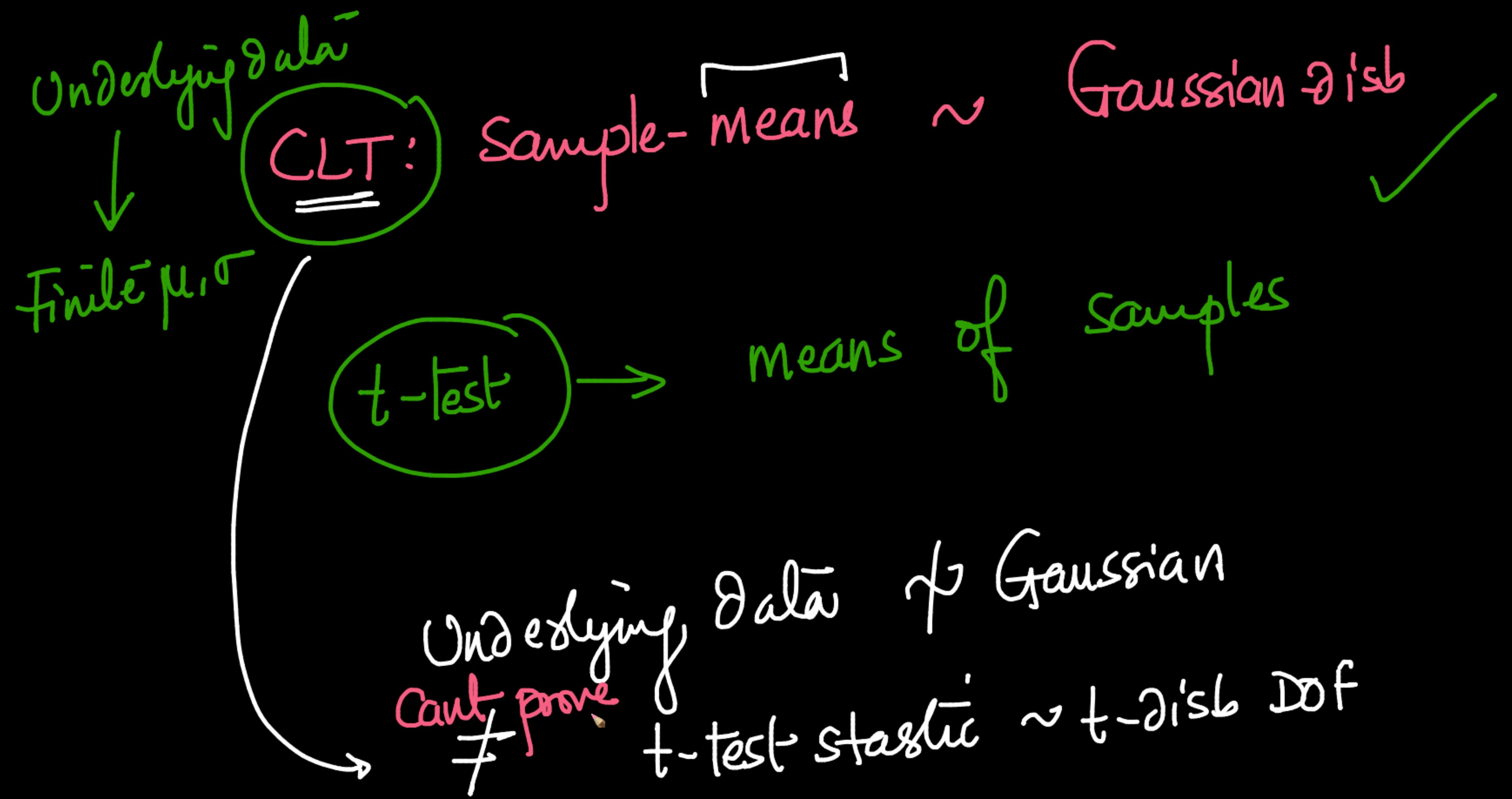
(Q) 2-Sample z-test : parametric or non-param



{ lots of hyp. testing
=

↓ appropriate in the context ; param or non-param

→ 6-step framework
=====



Underlying data

Uniform dist
(finite μ, σ)

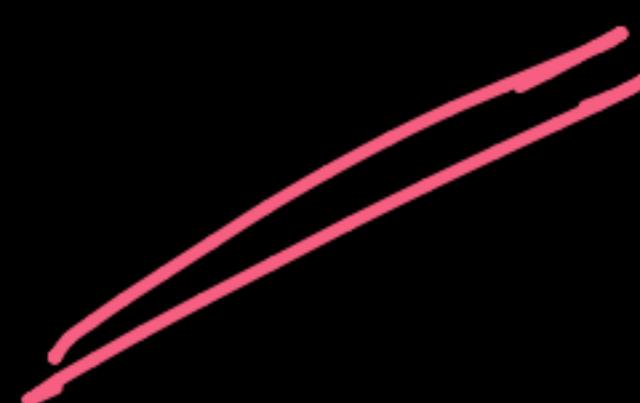
Sample-means \sim Gaussian



NOT converging
non-Gaussian to
Gaussian

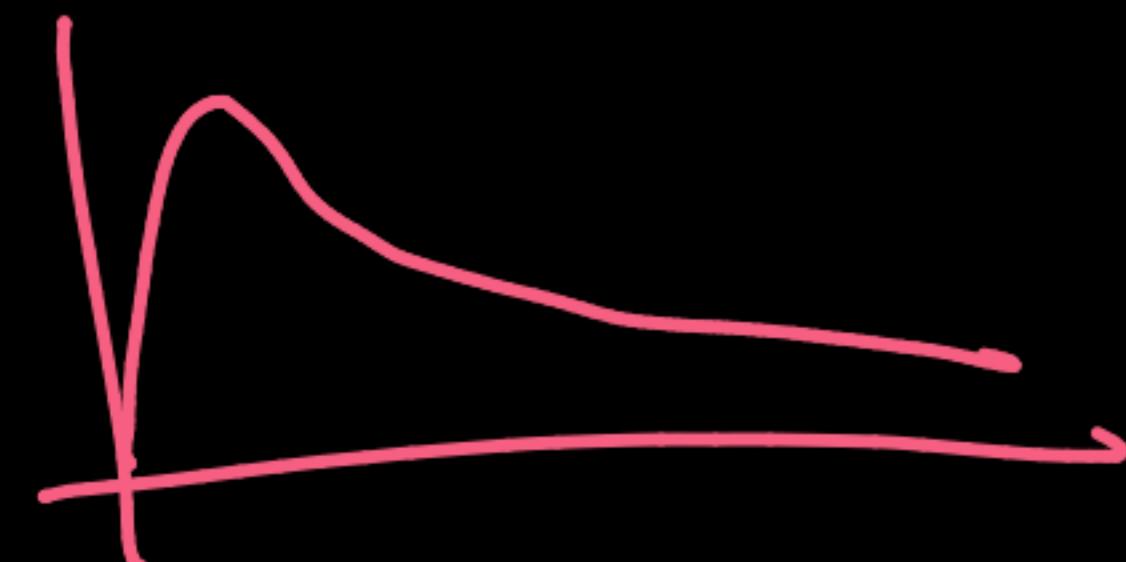


14 / 14



Non-Gaussian \rightarrow Gaussian

$X \rightarrow \{x_1, x_2, \dots\}$ *Observation*



e.g.:

log-normal

$X \sim \text{log-normal}$

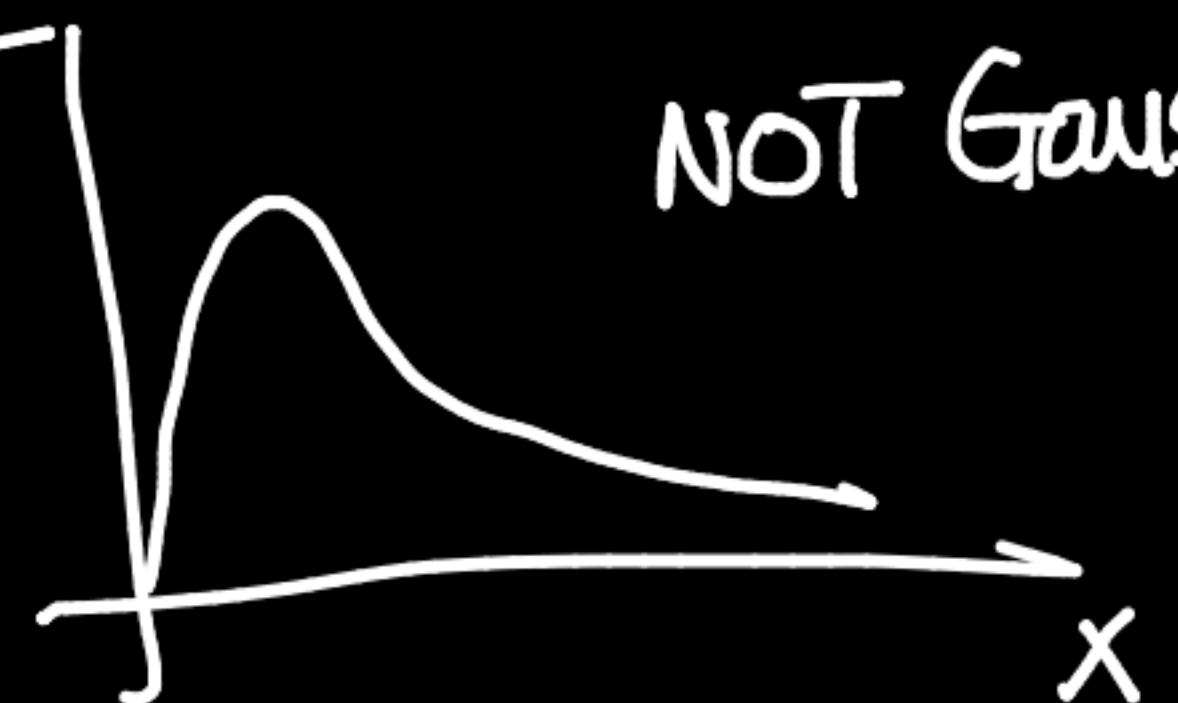
$\log_e(X) = \ln(X) \sim \text{Gaussian-Distr}$

Obs

$$X \Rightarrow \{x_1, x_2, x_3, \dots, x_n\}$$

PDF

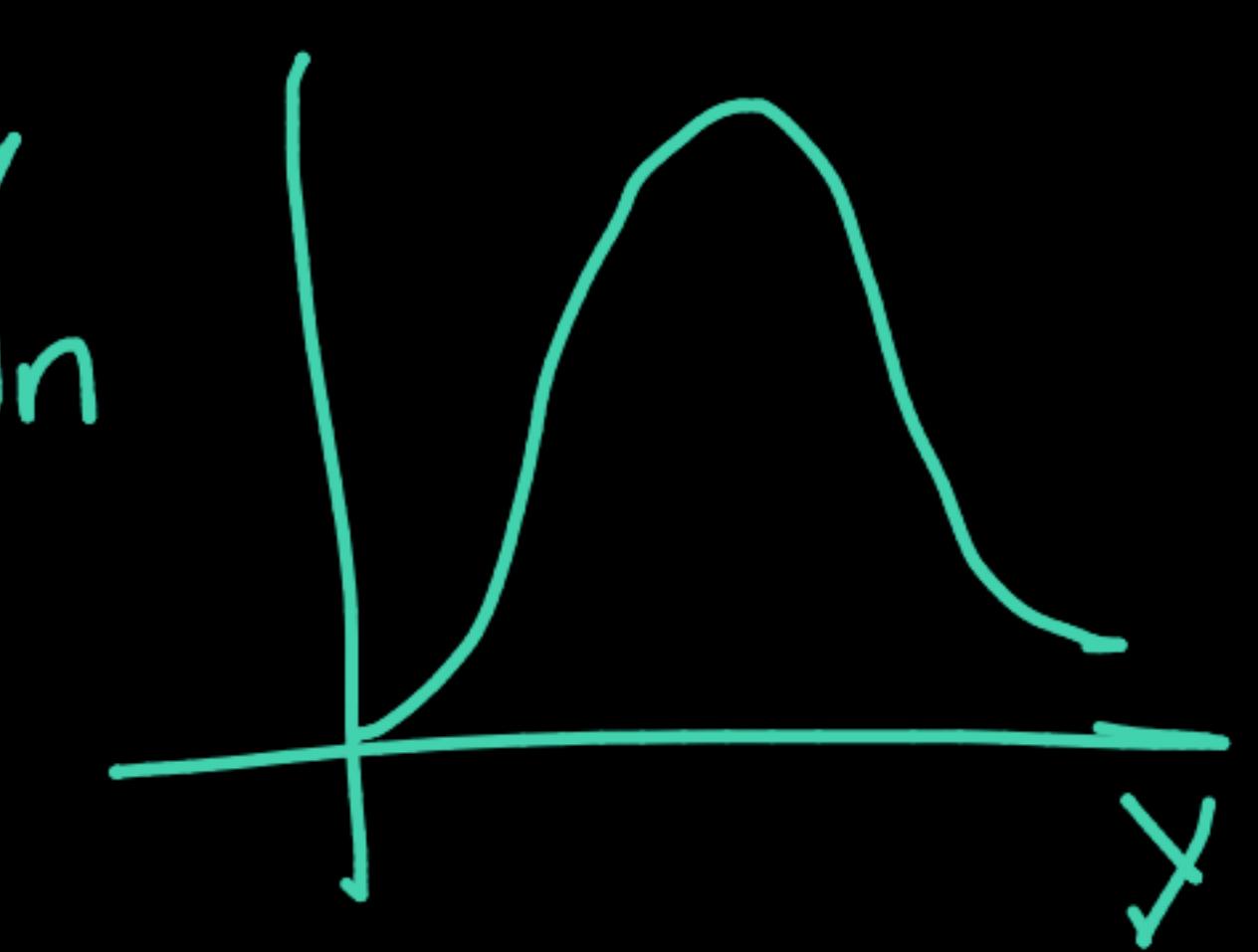
NOT Gaussian



one
transforming

$\ln(x)$: log transform

$$Y = \{ \ln(x_1), \ln(x_2), \ln(x_3), \dots, \ln(x_n) \}$$
$$y_1, y_2, y_3, \dots, y_n$$



why?

analogy: cm → feet
°C → °F

↳ Non-Gaussian
(log-normal)

10, 10

transforming

log_e

Gaussian

t-test

tamping

GX

$$\begin{aligned} 123.6 &\rightarrow 121.6 \\ 250.3 &\rightarrow 213.45 \end{aligned}$$

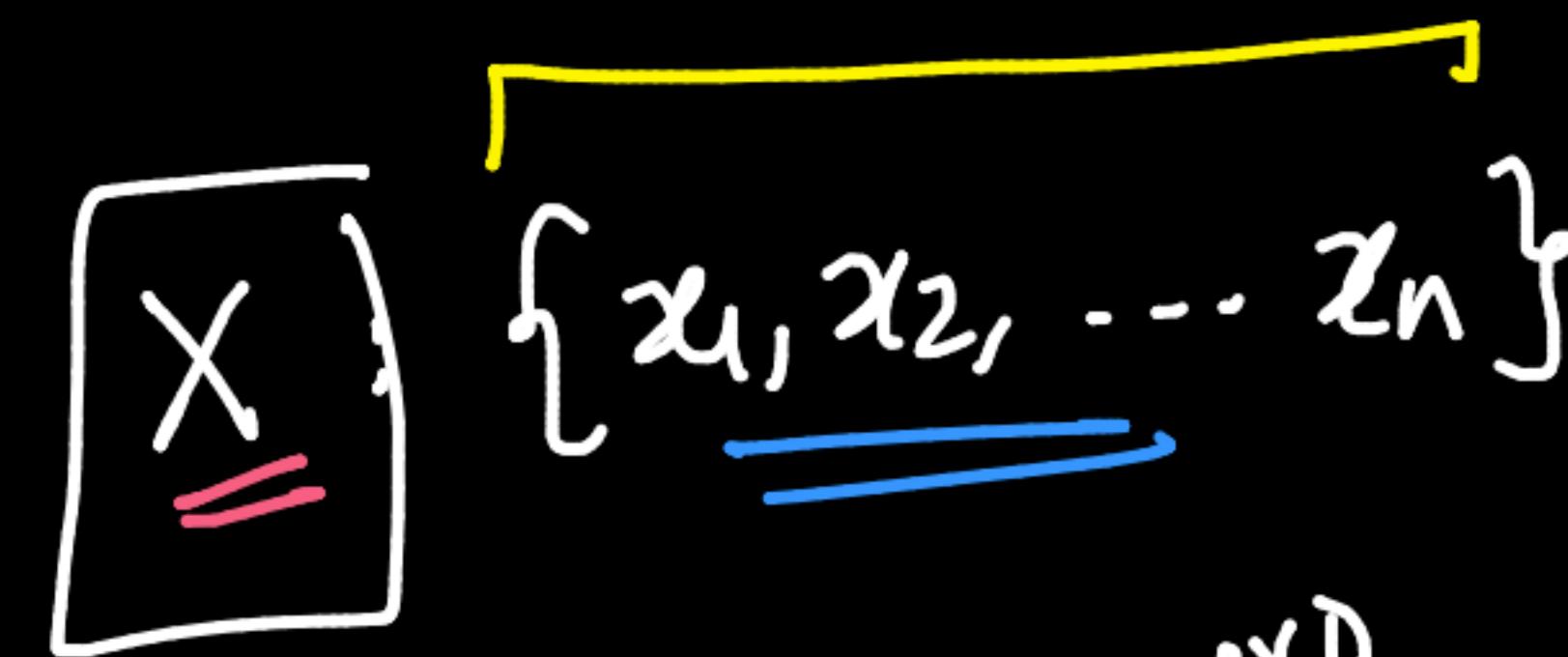
Conclusions

hyp-testing: $H_0: M_1 = M_2$
 $H_a: M_1 \neq M_2$

$$y = \ln(x)$$

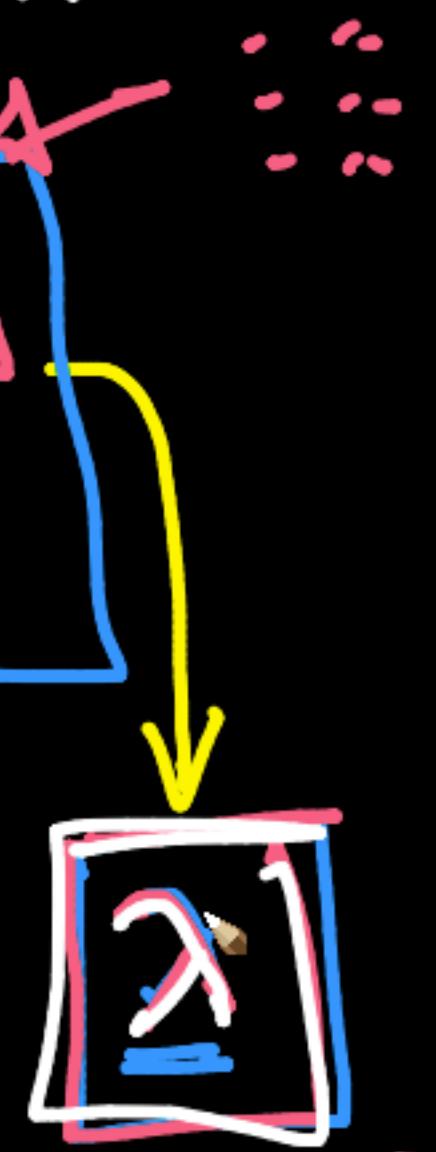
$$x = e^y = \exp(y)$$

✓ { Box-Cox :



non-Gaussian

box-Cox
(scipy)



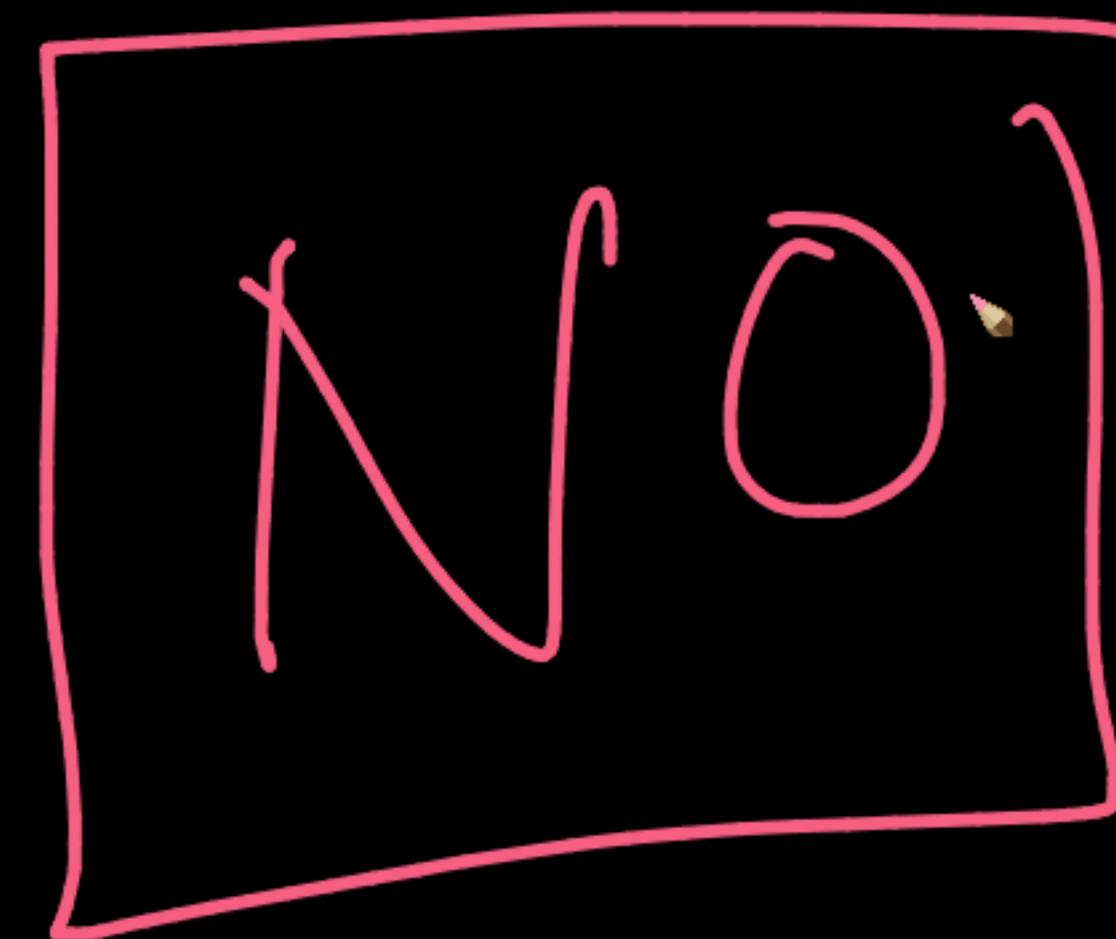
Gaussian

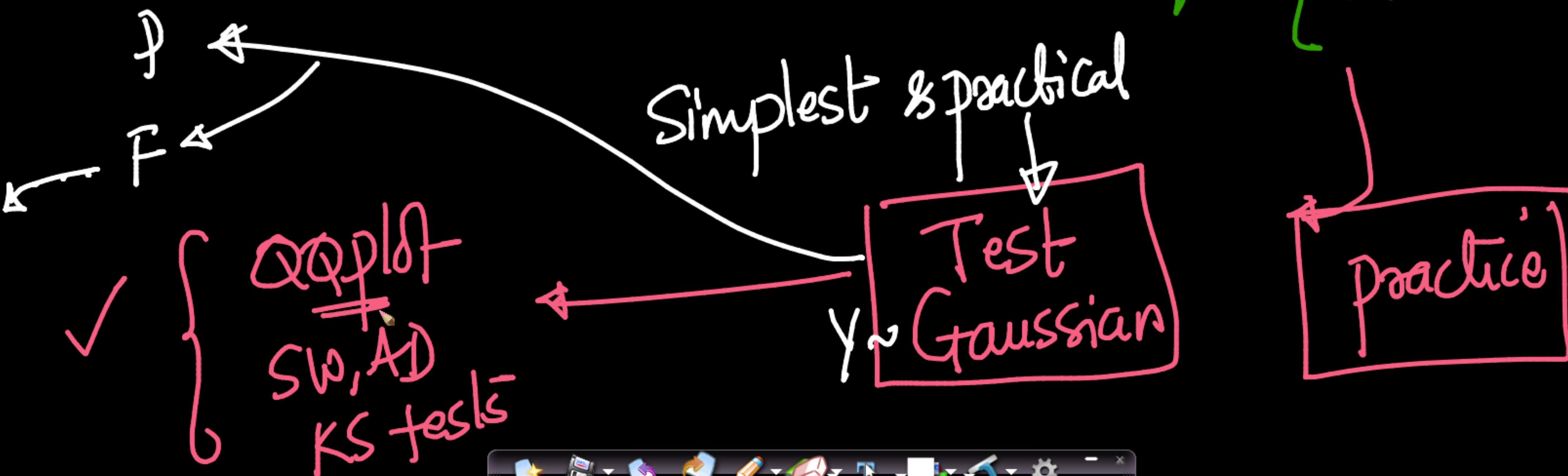
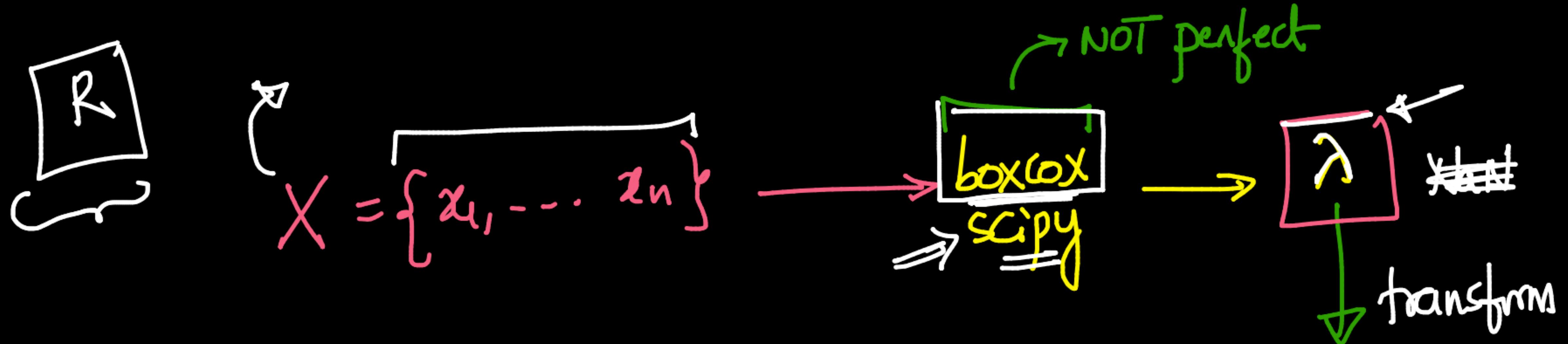
$$y_i = \begin{cases} \frac{x_i - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(x) & \text{if } \lambda = 0 \end{cases}$$

generalization
of log-normal

optimization

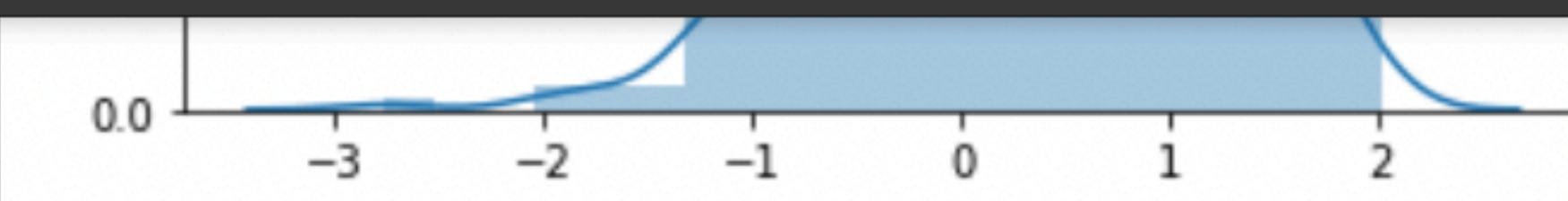
any data \rightarrow Gaussian [No]





+ Code + Text

✓ [5]
0s

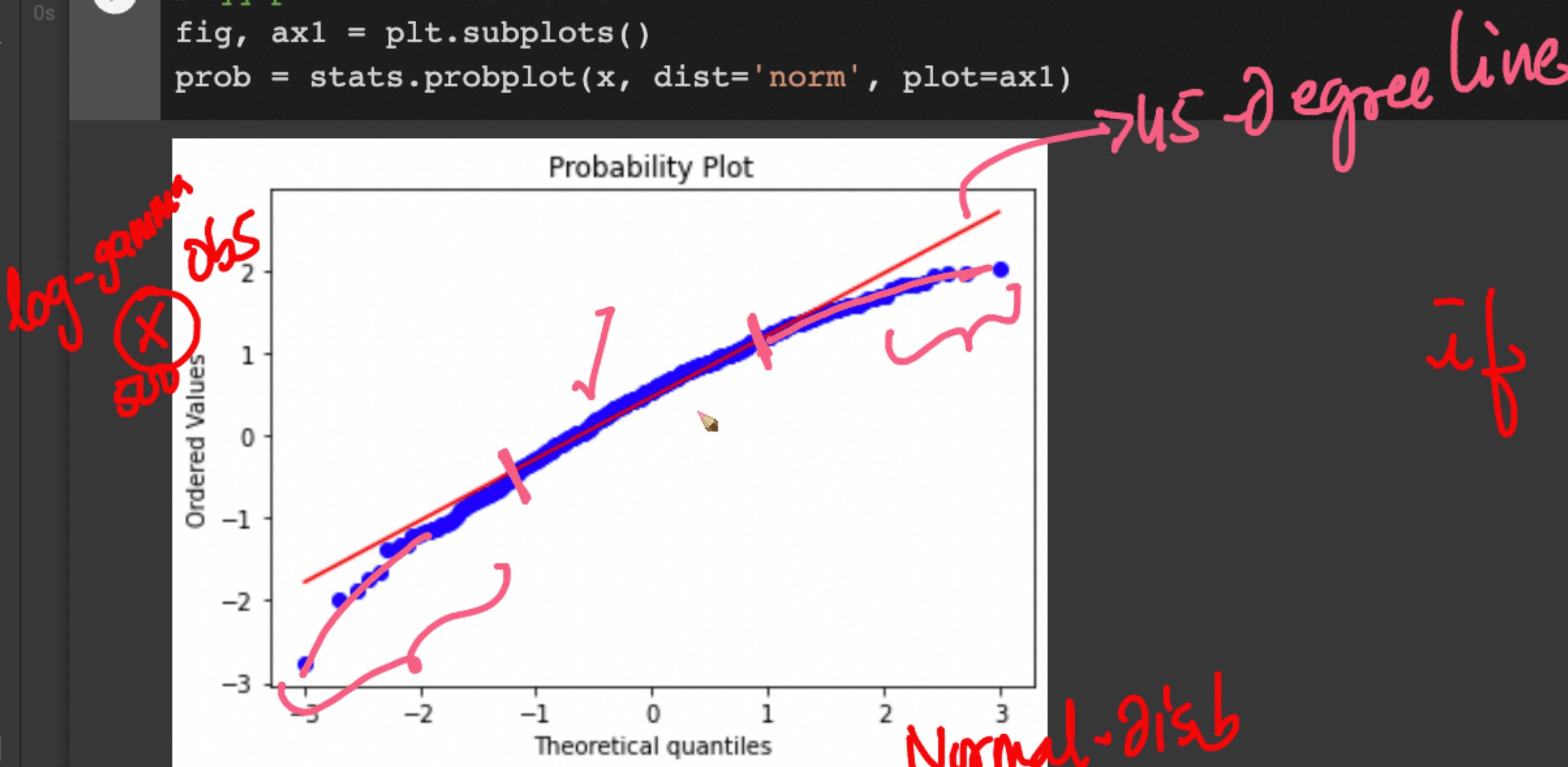


- ✓ RAM Disk

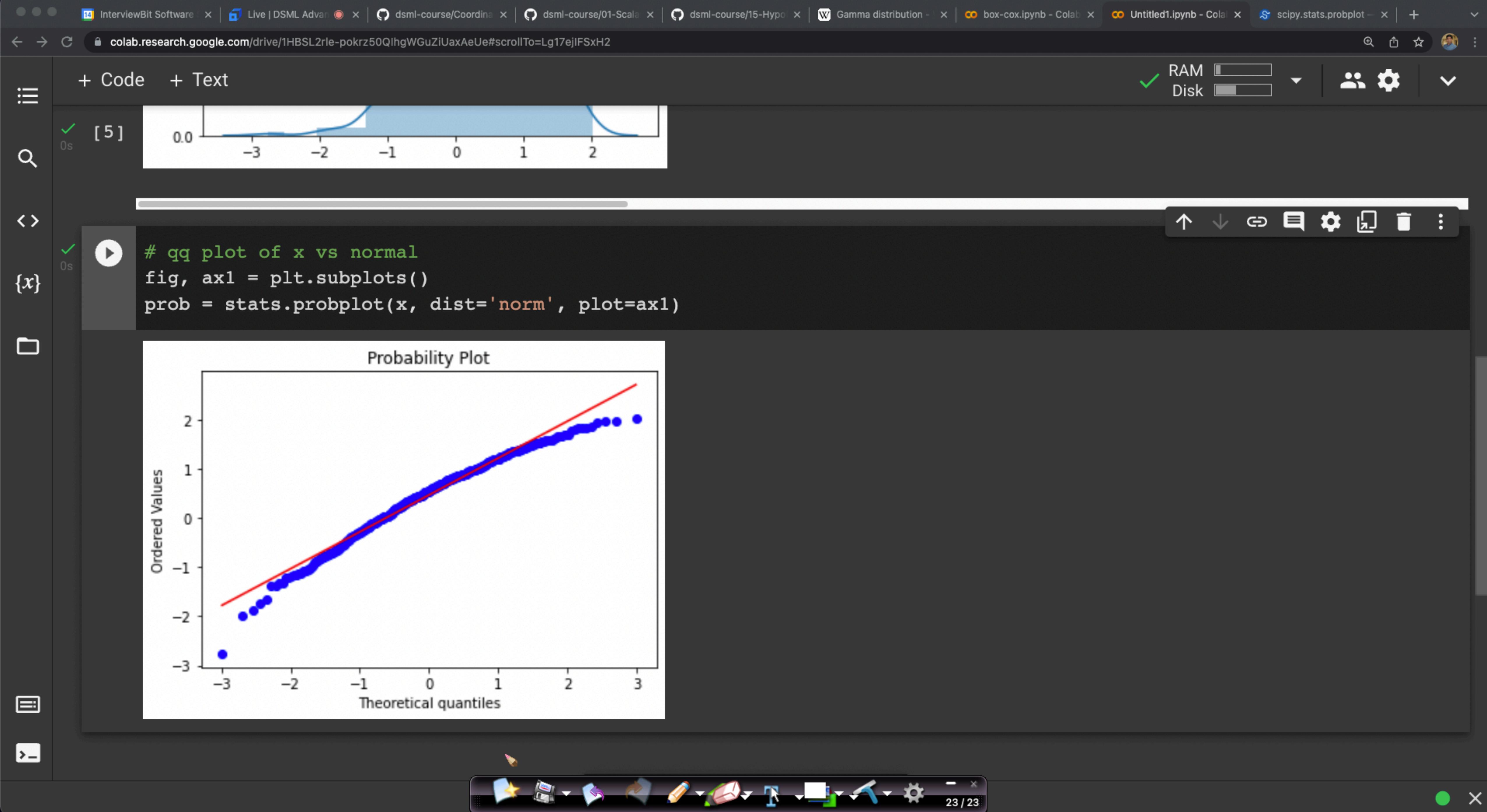
1

```
✓ 0s # qq plot of x vs normal
fig, ax1 = plt.subplots()
prob = stats.probplot(x, dist='norm', plot=ax1)
```

A horizontal bar containing several icons: an upward arrow, a downward arrow, a link icon, a message icon, a gear icon, a refresh/circular arrow icon, a trash bin icon, and three vertical dots.



if $X \sim \text{Normal}$



InterviewBit Software x | Live | DSML Advan x | dsml-course/Coordina x | dsml-course/01-Scala x | dsml-course/15-Hypo x | Gamma distribution - x | box-cox.ipynb - Colab x | Untitled1.ipynb - Colab x | scipy.stats.boxcox - x

colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=D6-0DIS2T_Pb

RAM Disk

+ Code + Text

Theoretical quantiles

abs →

```
# box cox transform  
x1 = x+5  
xt, l = stats.boxcox(x,); # returns x_trasnformed and lambda  
print("lambda :" + str(l))
```

```
# check if xt is gaussian or not using QQ-Plot  
fig, ax2 = plt.subplots()  
prob = stats.probplot(xt, dist='norm', plot=ax2)
```

$x_1^2, x_2^{-2}, \dots, x_n$

```
-----  
ValueError Traceback (most recent call last)  
<ipython-input-10-1c3c724c108c> in <module>()  
      1 # box cox transform  
----> 2 xt, l = stats.boxcox(x); # returns x_trasnformed and lambda  
      3 print("lambda :" + str(l))  
      4  
      5 # check if xt is gaussian or not using QQ-Plot
```

+5

-5

$x_1 + 5, x_2 + 5, \dots, x_n + 5$

```
/usr/local/lib/python3.7/dist-packages/scipy/stats/morestats.py in boxcox(x, lmbda, alpha)  
    1043  
    1044      if any(x <= 0):  
-> 1045          raise ValueError("Data must be positive.")  
    1046  
    1047      if lmbda is not None: # single transformation
```



InterviewBit Software x | Live | DSML Advan x | dsml-course/Coordina x | dsml-course/01-Scala x | dsml-course/15-Hypo x | Gamma distribution - x | box-cox.ipynb - Colab x | Untitled1.ipynb - Colab x | scipy.stats.boxcox - x | +

colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=D6-0DIS2T_Pb

+ Code + Text

✓ RAM Disk

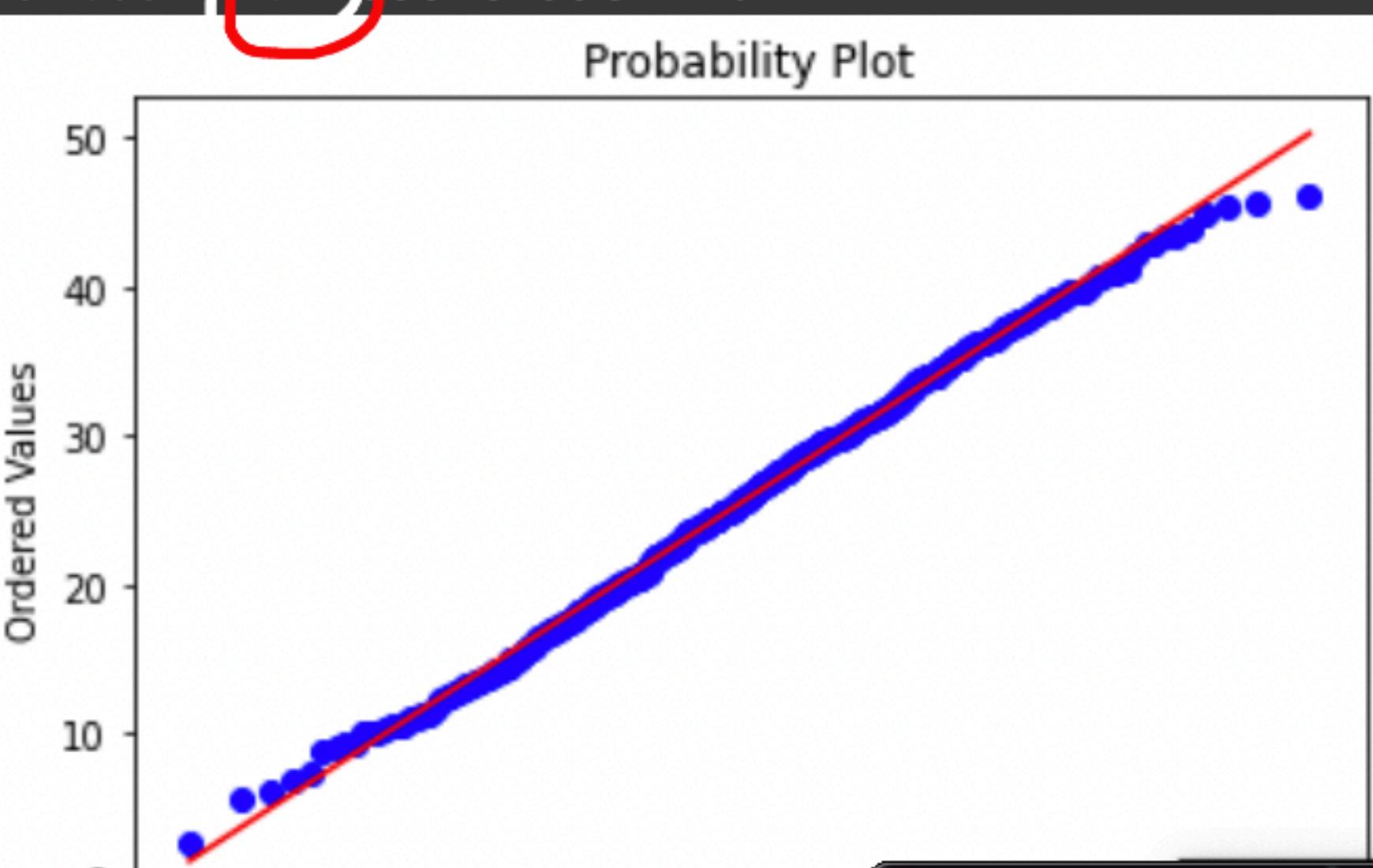


[]



```
# box cox transform  
x1 = x+5  
xt, l = stats.boxcox(x1,); # returns x_trasnformed and lambda  
print("lambda :" + str(l))  
  
# check if xt is gaussian or not using QQ-Plot  
fig, ax2 = plt.subplots()  
prob = stats.probplot(xt, dist='norm', plot=ax2)
```

lambda 2.4208828199512202



✓
box cox $(x_i + 5)$

$\lambda = 2.4 \neq 0$

$$\left\{ x_t = \frac{(x_i + 5)^{2.42} - 1}{2.42} \right\}$$

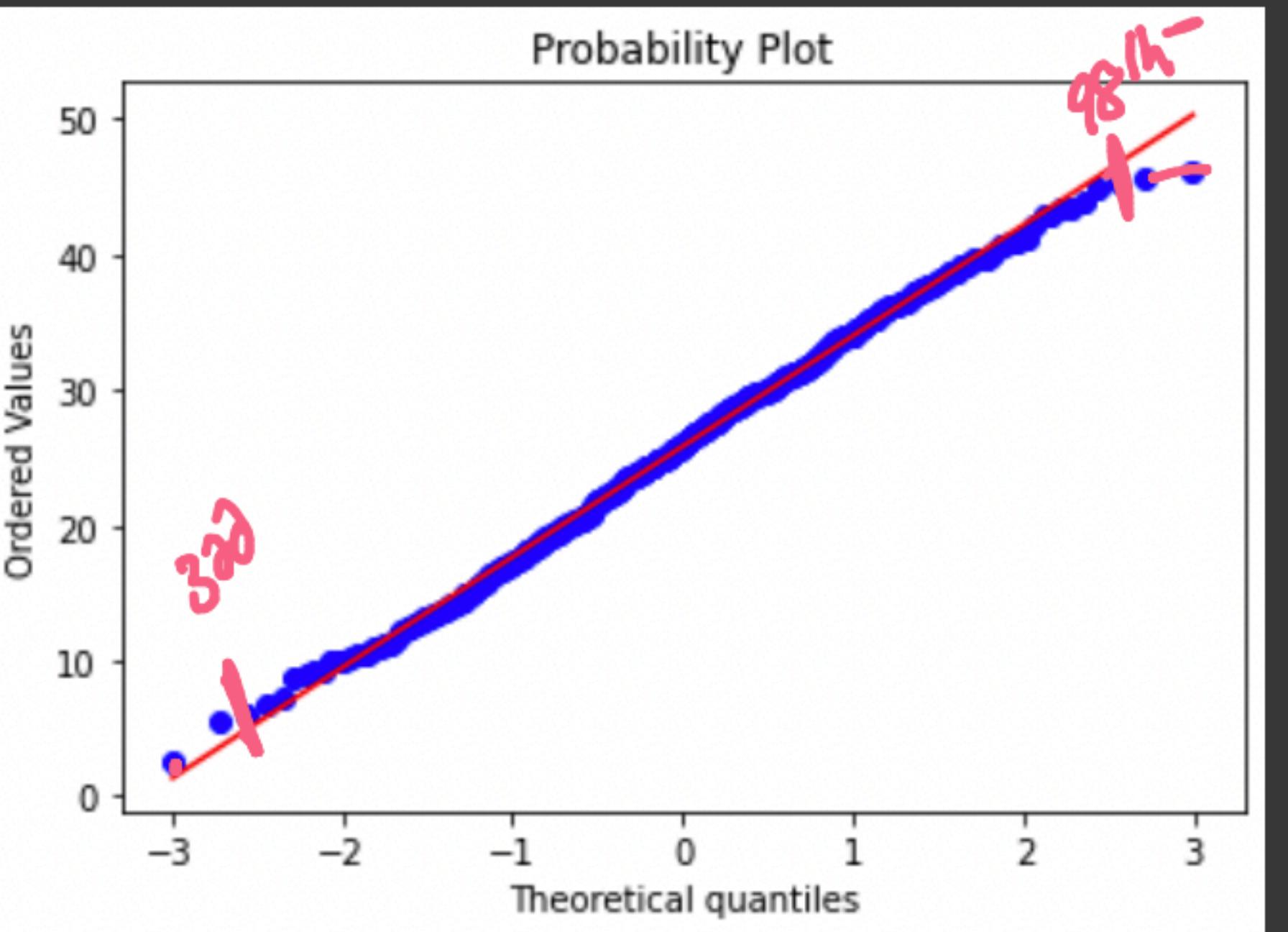
+ Code + Text

✓ RAM Disk

```
xt, l = stats.boxcox(x1,); # returns x_transformed and lambda
print("lambda :" + str(l))

# check if xt is gaussian or not using QQ-Plot
fig, ax2 = plt.subplots()
prob = stats.probplot(xt, dist='norm', plot=ax2)
```

lambda :2.4268828199512202



better
~ Gaussian
approx

Confidence Interval using bootstrapping

~~per item sale~~
~~price~~

e.g.: $S \rightarrow$

$$\{S_1, S_2, S_3, \dots, S_{100}\}$$

$n=100$

→ randomly sampling
Millions

not gaussian disb

only ~~data~~

(Q) find the

95% C.I. of the

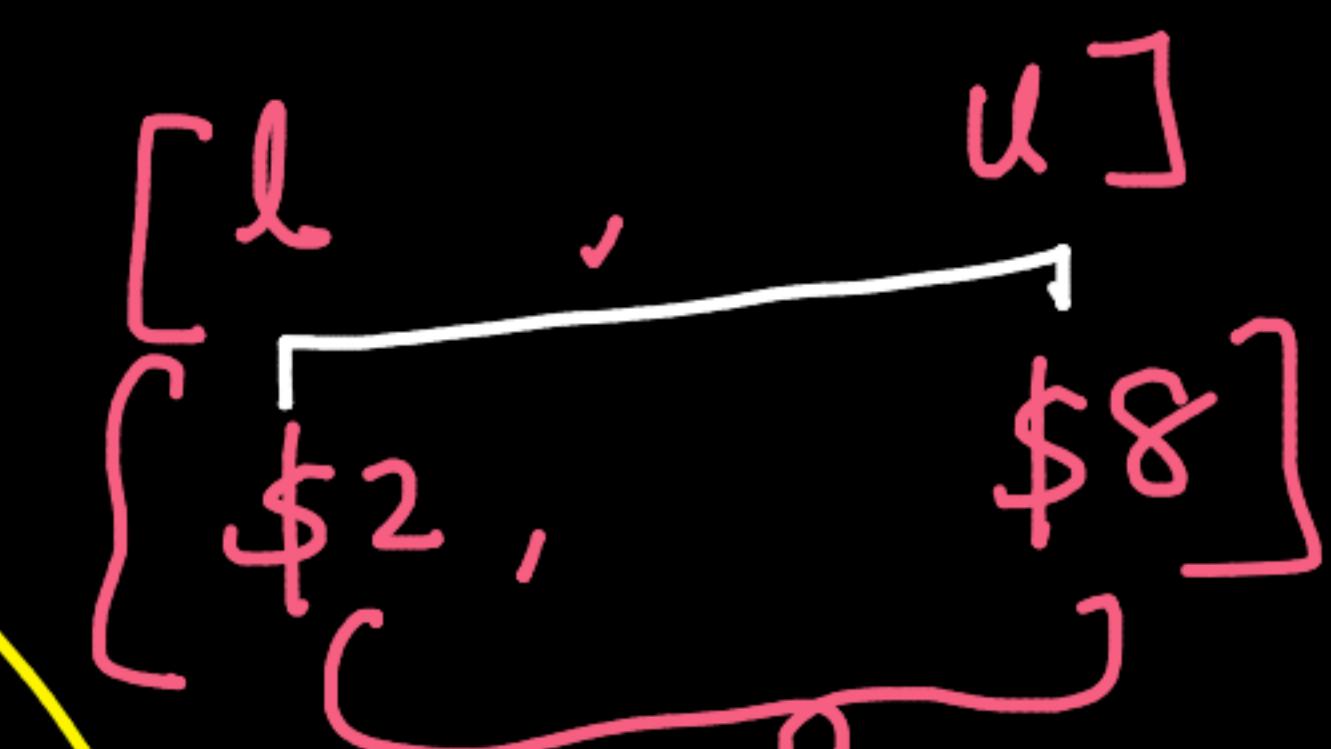
median

Sale

price @ Amazon

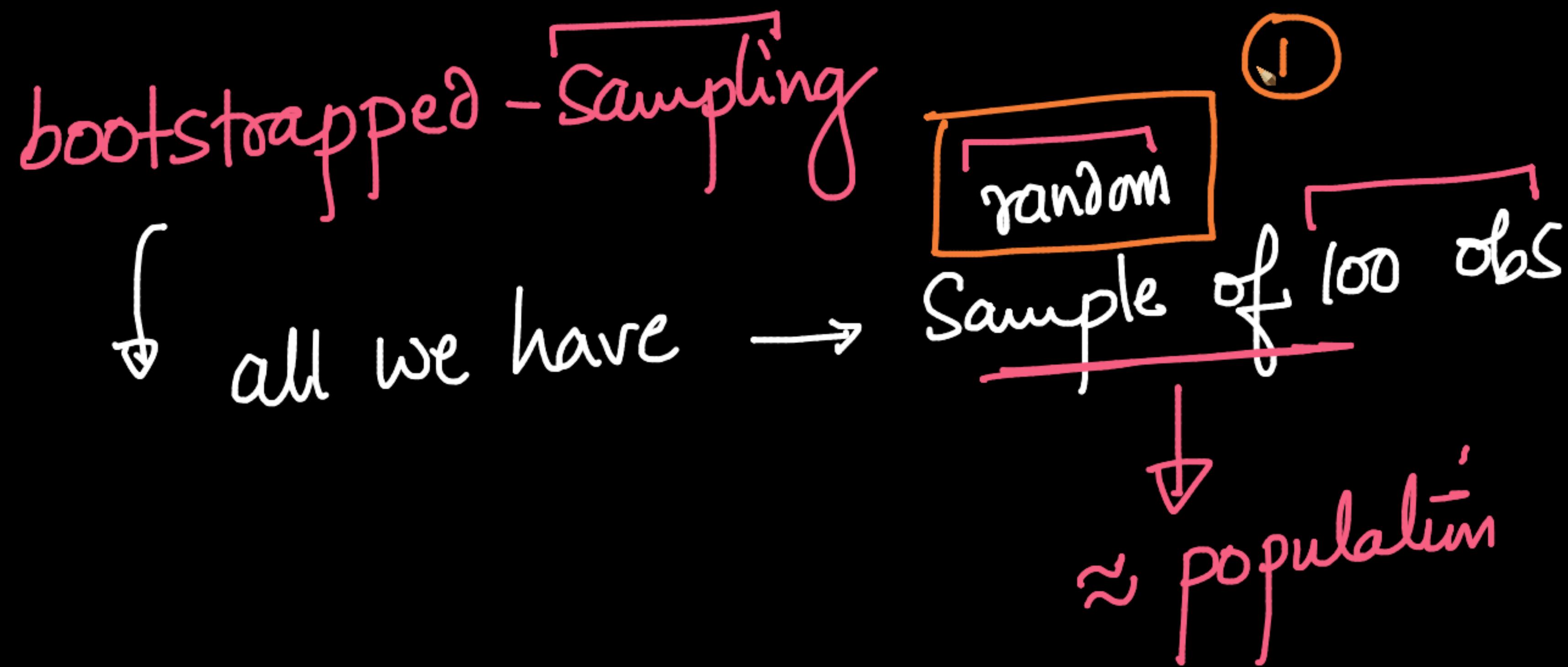
(population)

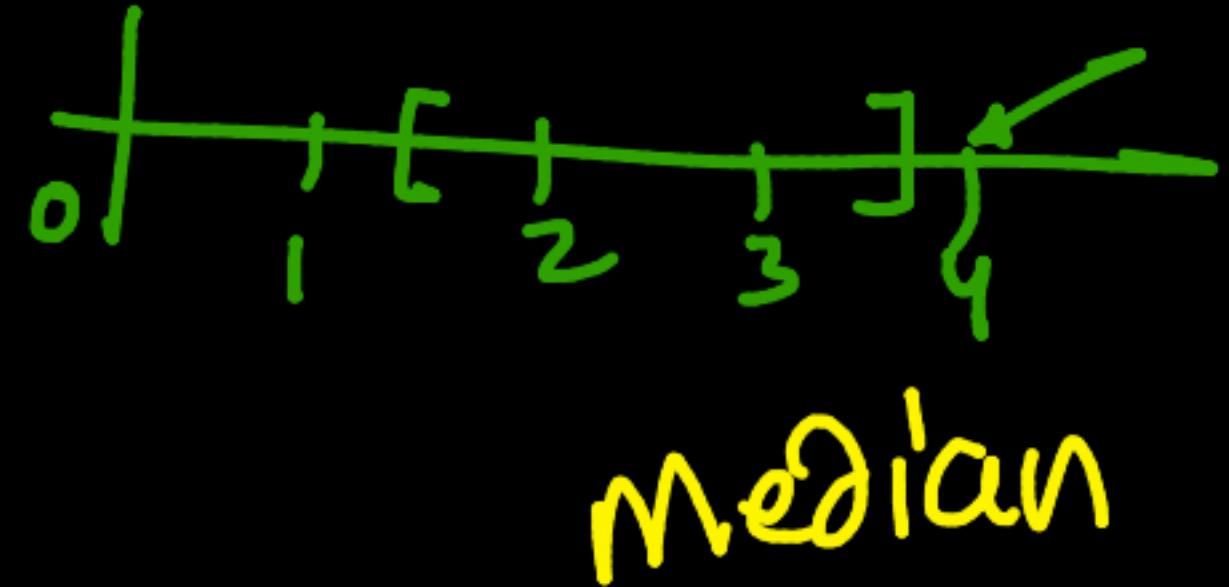
→ bootstrapped samples



Given: $\mu \approx \text{Median}$

95% C.I. of the median of the population





point estimate
\$ 4 = median (sample)

DO
Not
population

estimate population-median

✓ $[\$2, \$6]$ 95% C.I.

left
 $(\$4.5, \$4.5)$
pop-median

drugs for COVID

population median

median recovery time

|S|

1000

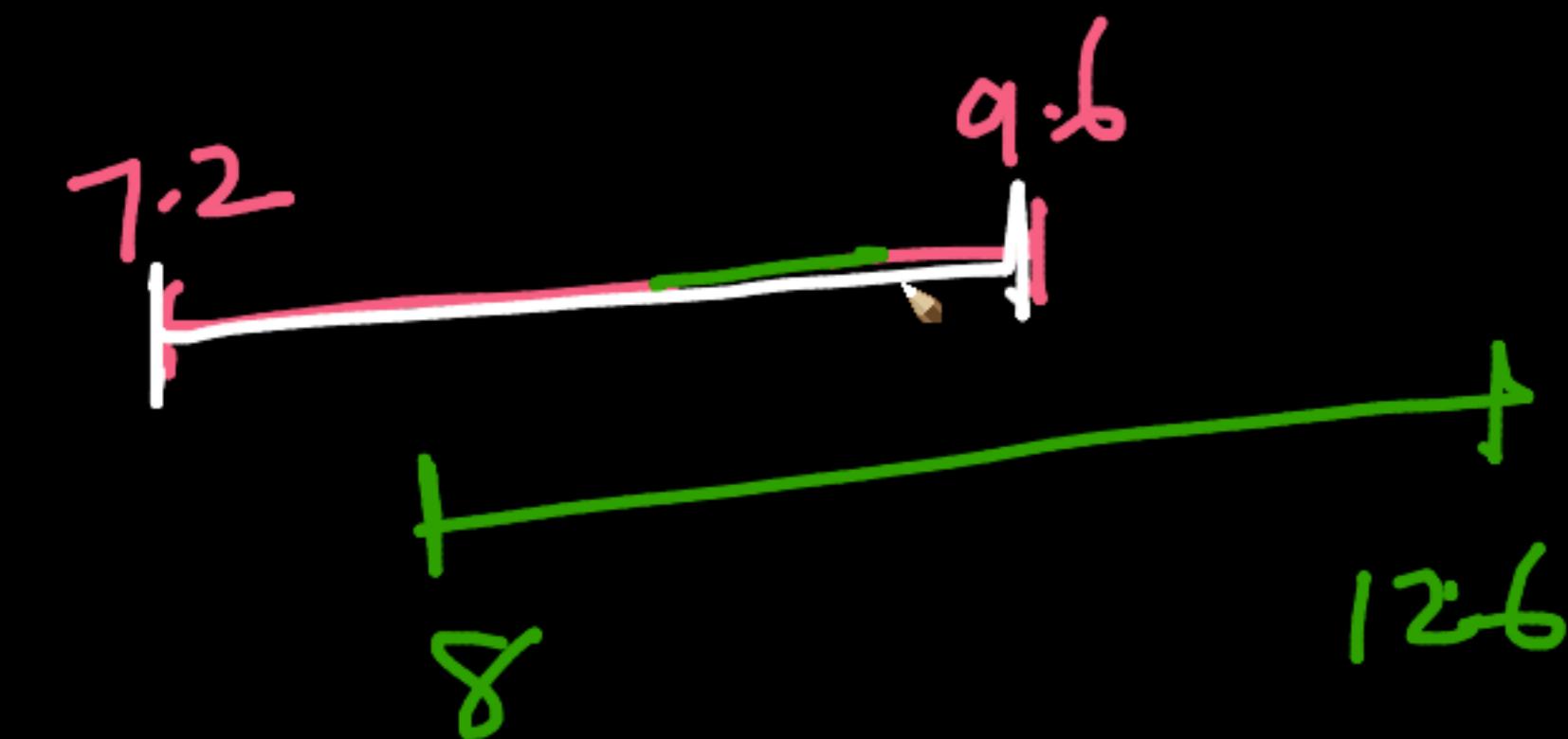
2000

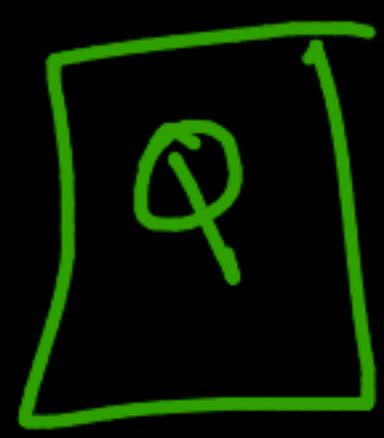
$$d_1 \rightarrow [7.2, 9.6]$$

$$d_2 \rightarrow [8, 12.6]$$

$\checkmark d_3$

d_4





Sample-size is fixed

95% C.I
[2, 3.2]

99% C.I
[1.2, 4.6]

wider or narrower

$[-\infty, \infty]$

medians = []

$S = \{S_1, S_2, \dots, S_n\}$

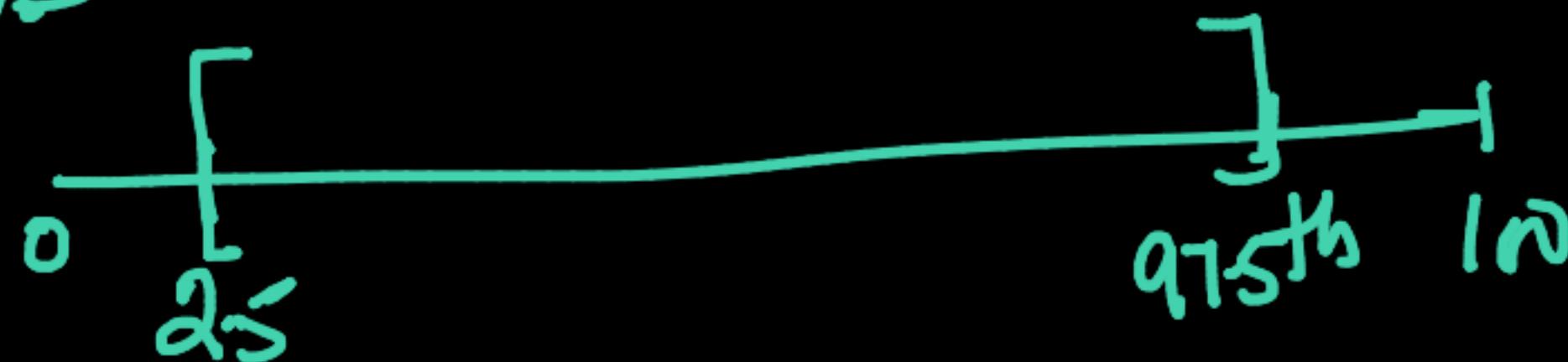
for i = 1 to 1000 # resampling

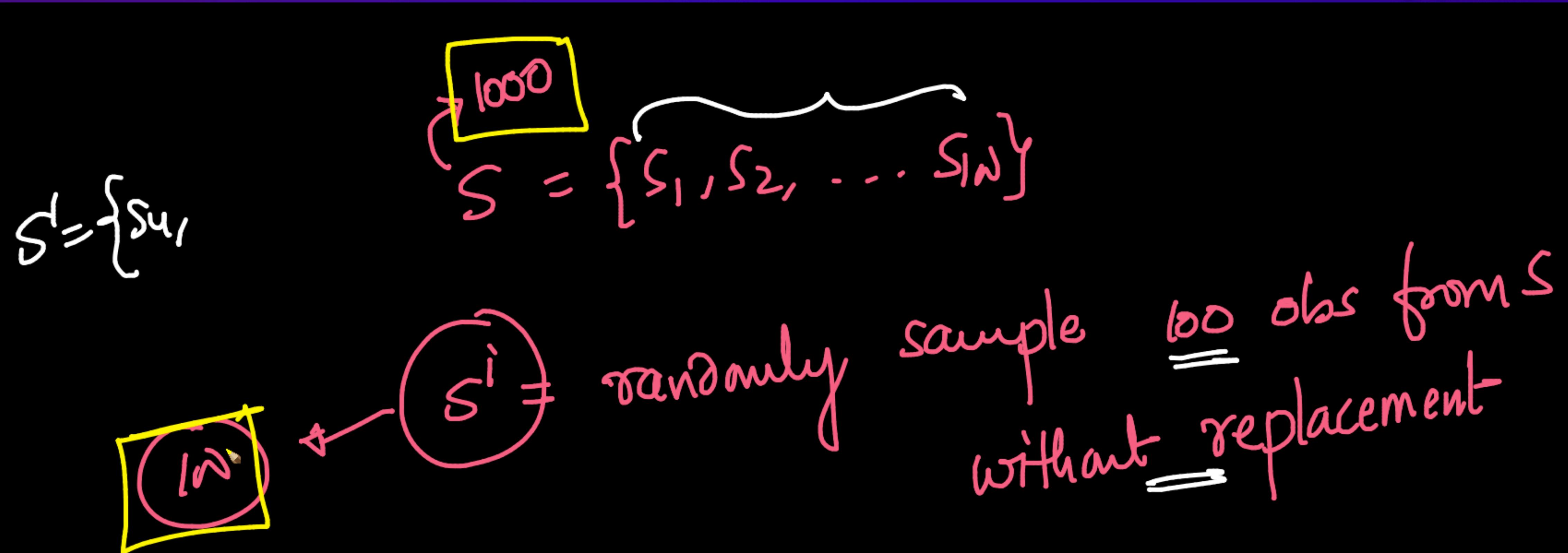
1. $S_i \xrightarrow{i}$ Sample 100 prices from S with replacement
2. $\underline{\text{medians}}[i] = \text{median}(S_i)$

Q Why not without replacement -

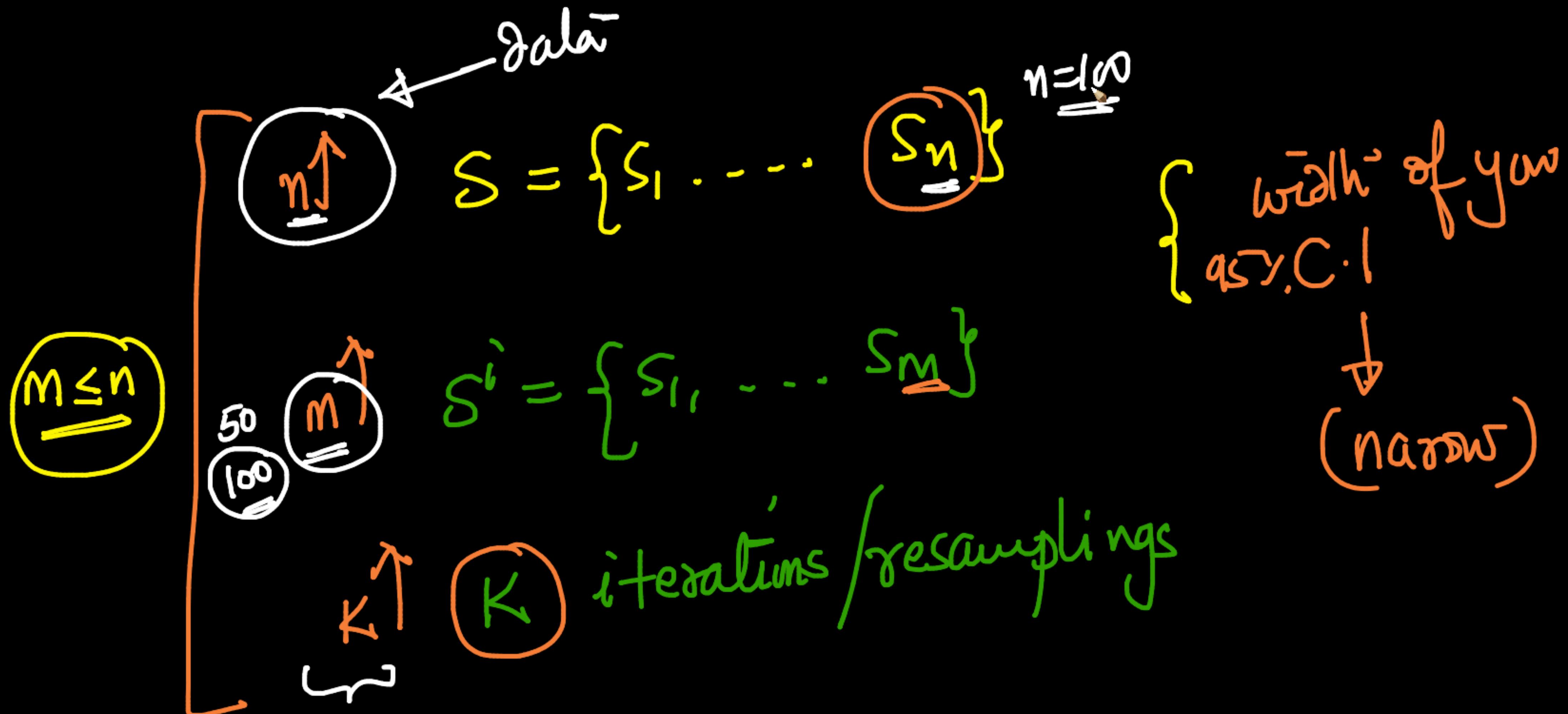
Sort medians

→





$$S^1 = S^2 = S^3 = S^4 \dots$$





q.s.c. of Median

[2.2 , p.6]

4

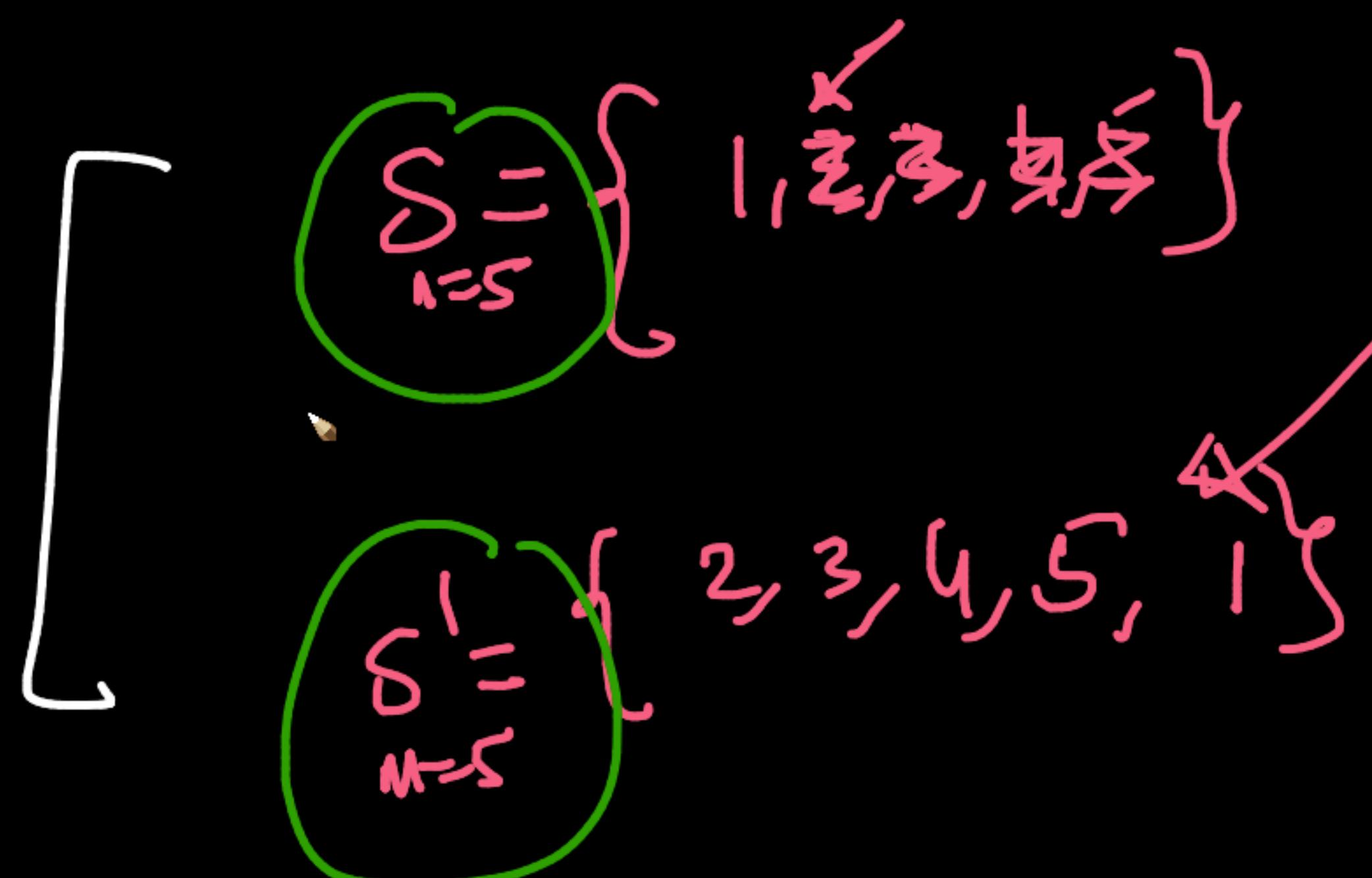
population - median
=



$$S^1 = S^2 = S^3 \dots$$

if Sampling without replacement

n=M

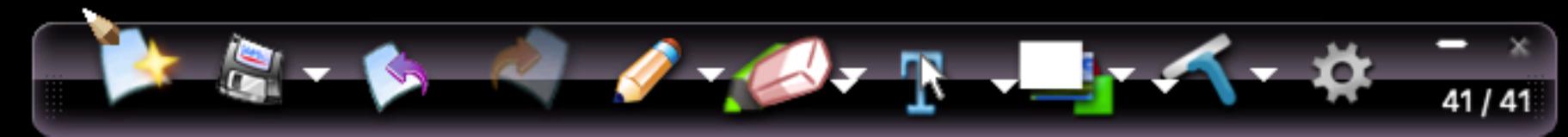


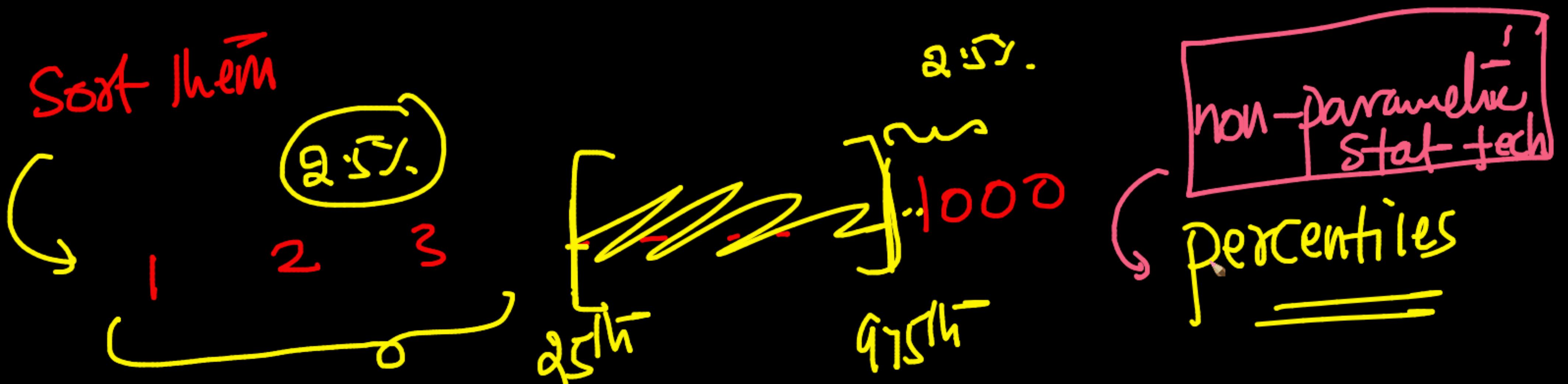
with repl:

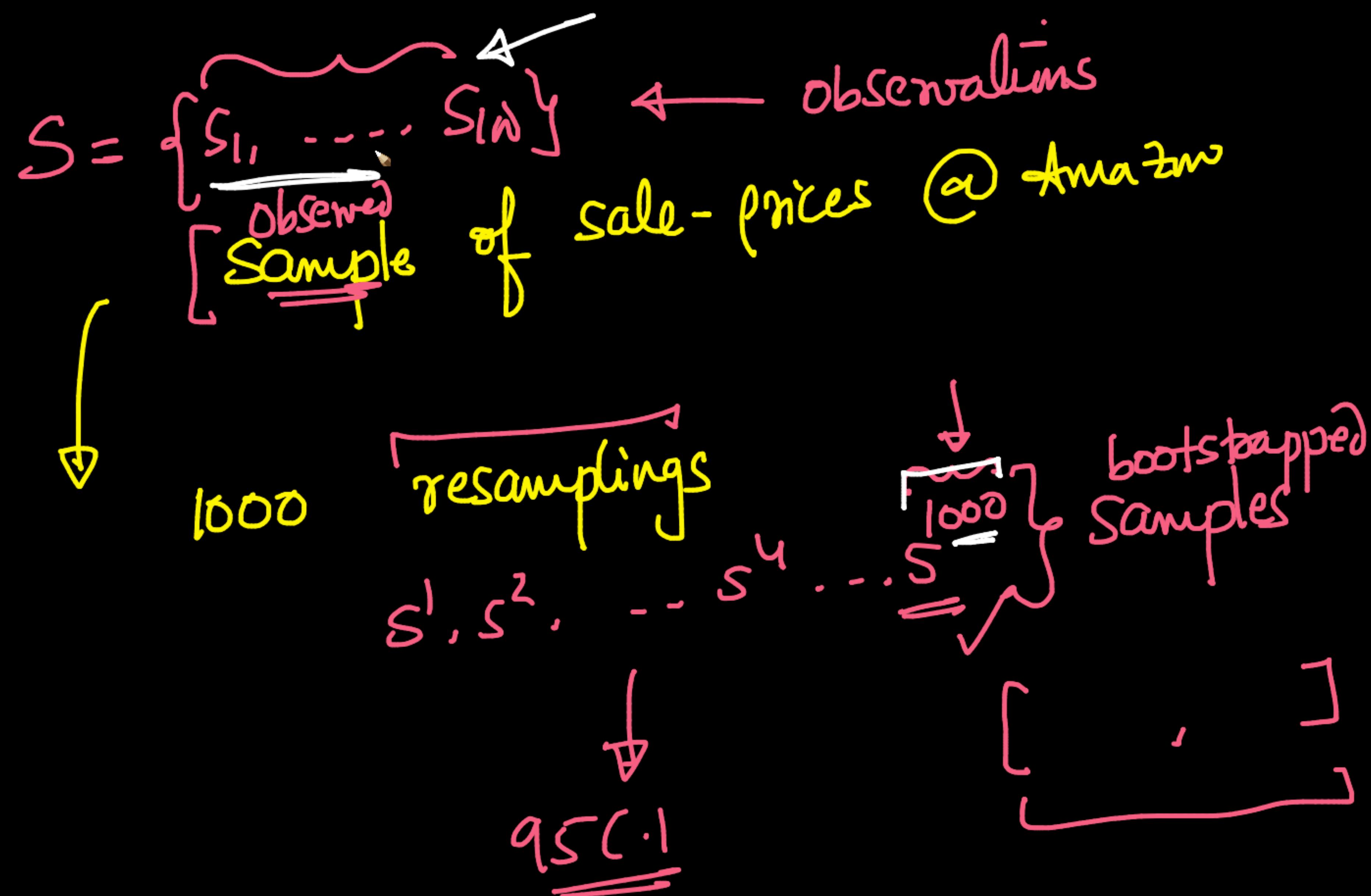
$$S = \{1, 2, 3, 4, 5\}$$

$$S' = \{3, 4, 3, 5, 1\}$$

✓ bootstrapping → Means
Medians
Percentiles







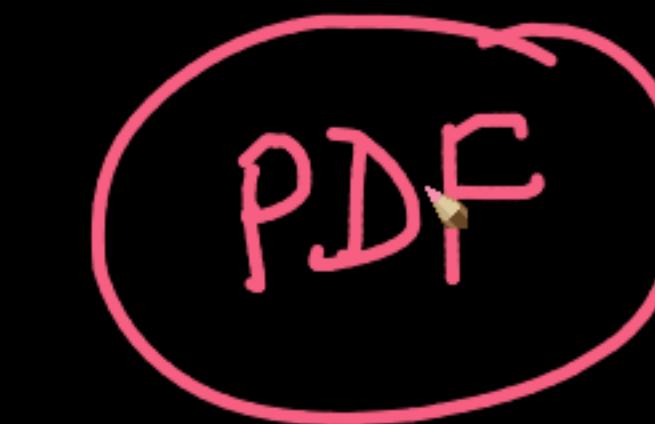
{ CLT → means

q6.23)

Prob(a person who weighs 96.23 kgs)

$\rightarrow 96.23000000\ldots$

Continuous r.v



prob (median = $\tilde{3.62}$) $\rightarrow 0$

↑
point-estimate
↓

InterviewBit Software x Live | DSML Advan x dsml-course/Coordina x dsml-course/01-Scala x dsml-course/15-Hypo x Gamma distribution - x box-cox.ipynb - Colab x Untitled1.ipynb - Colab x scipy.stats.boxcox - x +

scaler.com/meetings/i/dsml-advanced-linear-algebra-1-remedial/live

GEOmRTT

DSML Advanced : Linear Algebra 1- Remedial | Lecture

95% C.I of Median

[2.2 8.6]

You are sharing your screen now

Stop Sharing

Srikanth Varma Chekuri (You) (Screen)

02:03:12

8:00

Srikanth Varma Chekuri (You)

Chat

Notify me about Nothing

Pin a message +

0

Prateek Gupta To: Everyone 11:00 pm approx zero

Udit Manav To: Everyone 11:00 pm ~0

abhilash singh To: Everyone 11:01 pm close to 0

Rahul Shivani To: Everyone 11:01 pm 0

Prateek Gupta To: Everyone 11:01 pm zero

Rahul Shivani To: Everyone 11:01 pm ~0

Prateek Gupta To: Everyone 11:02 pm area under a point is zero

Rohit Sinha To: Everyone 11:02 pm let's say if we want to approx it to some decimal values or very small interval ? how do we approximate it ?

11:03 pm

Start Doubt Session

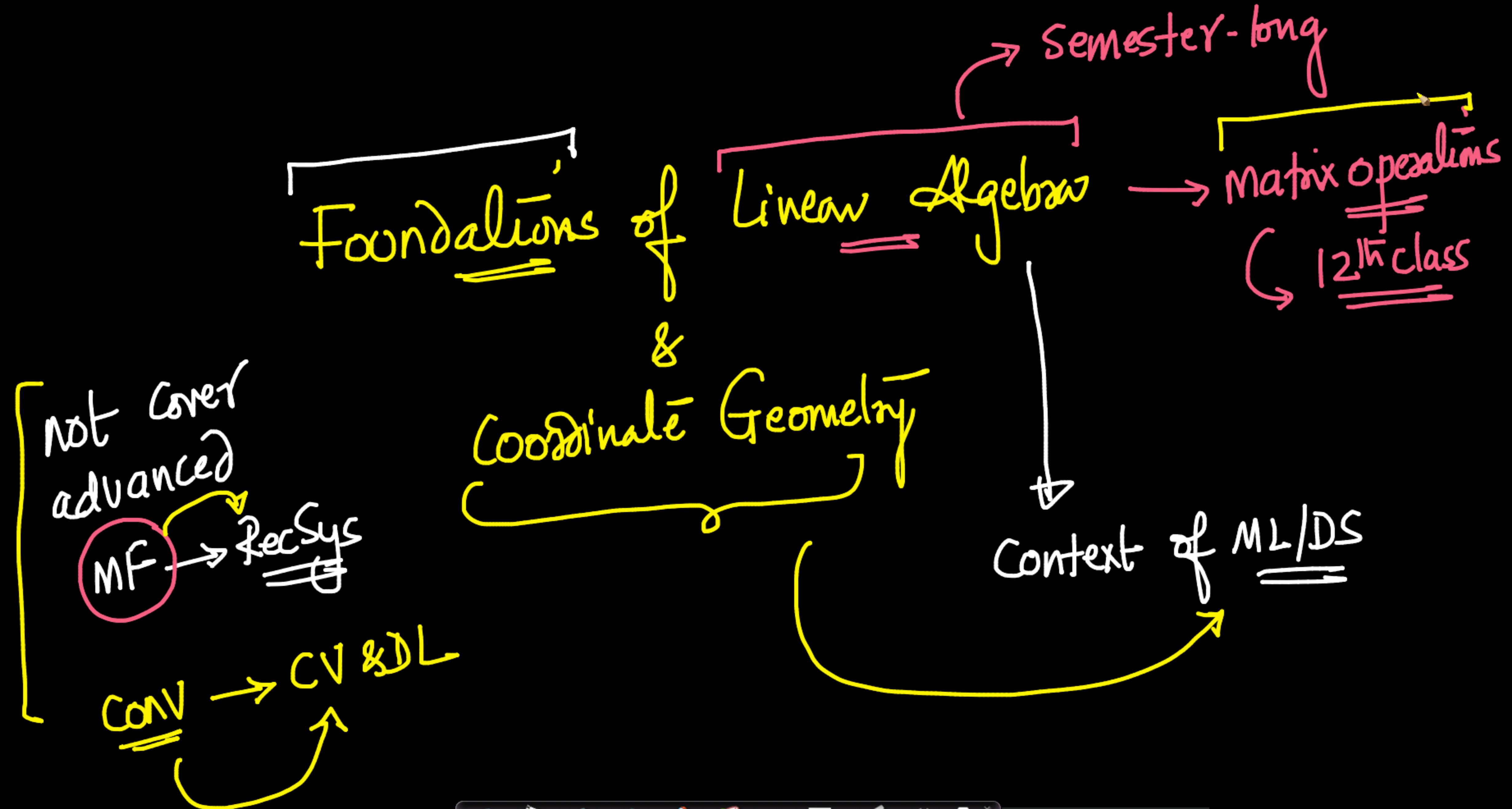
To: Everyone Enable/Disable Chat

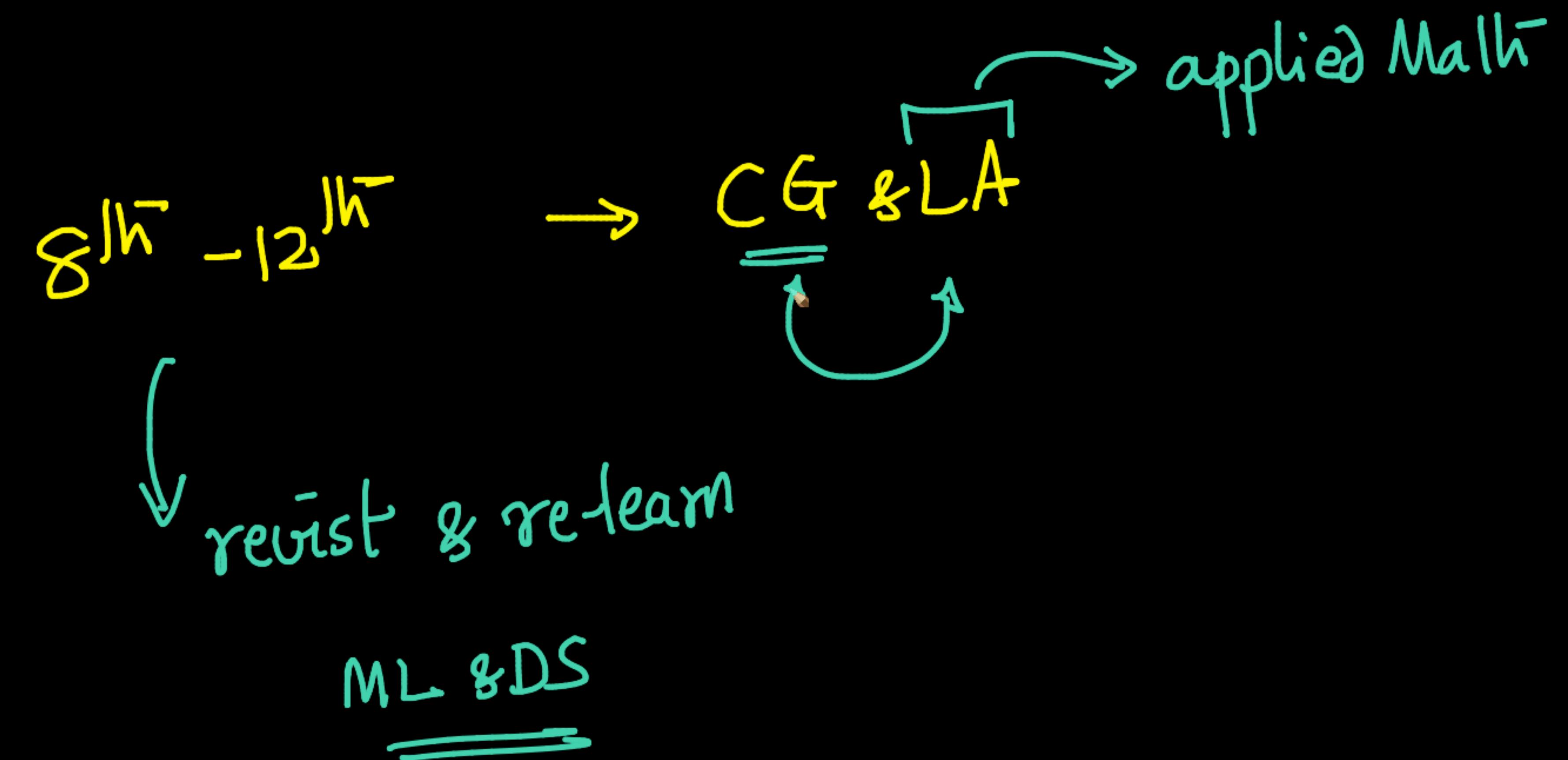
Type message

47 / 47

Break:

23:07

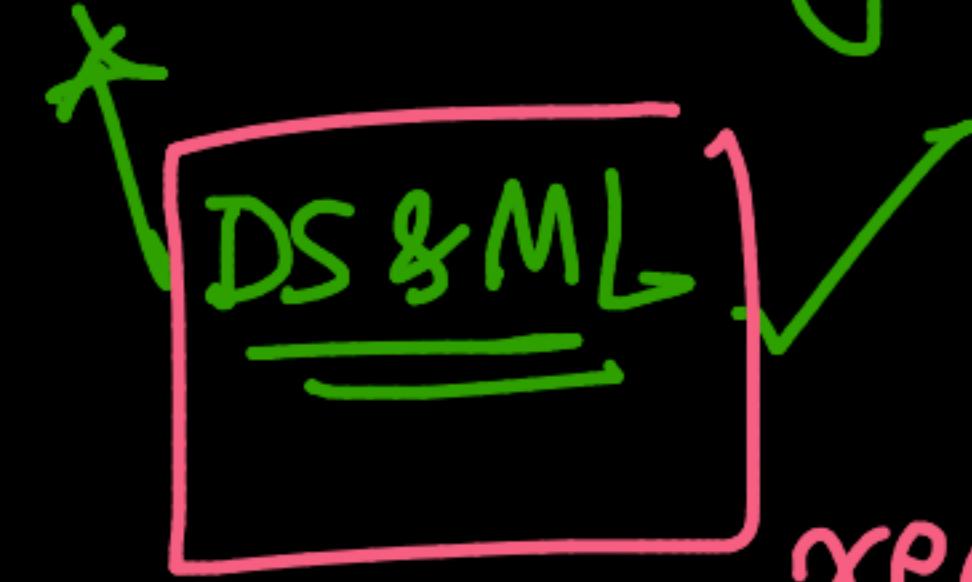




Lines, plane → high-school & college

vectors; matrices; dot-product → Physics X

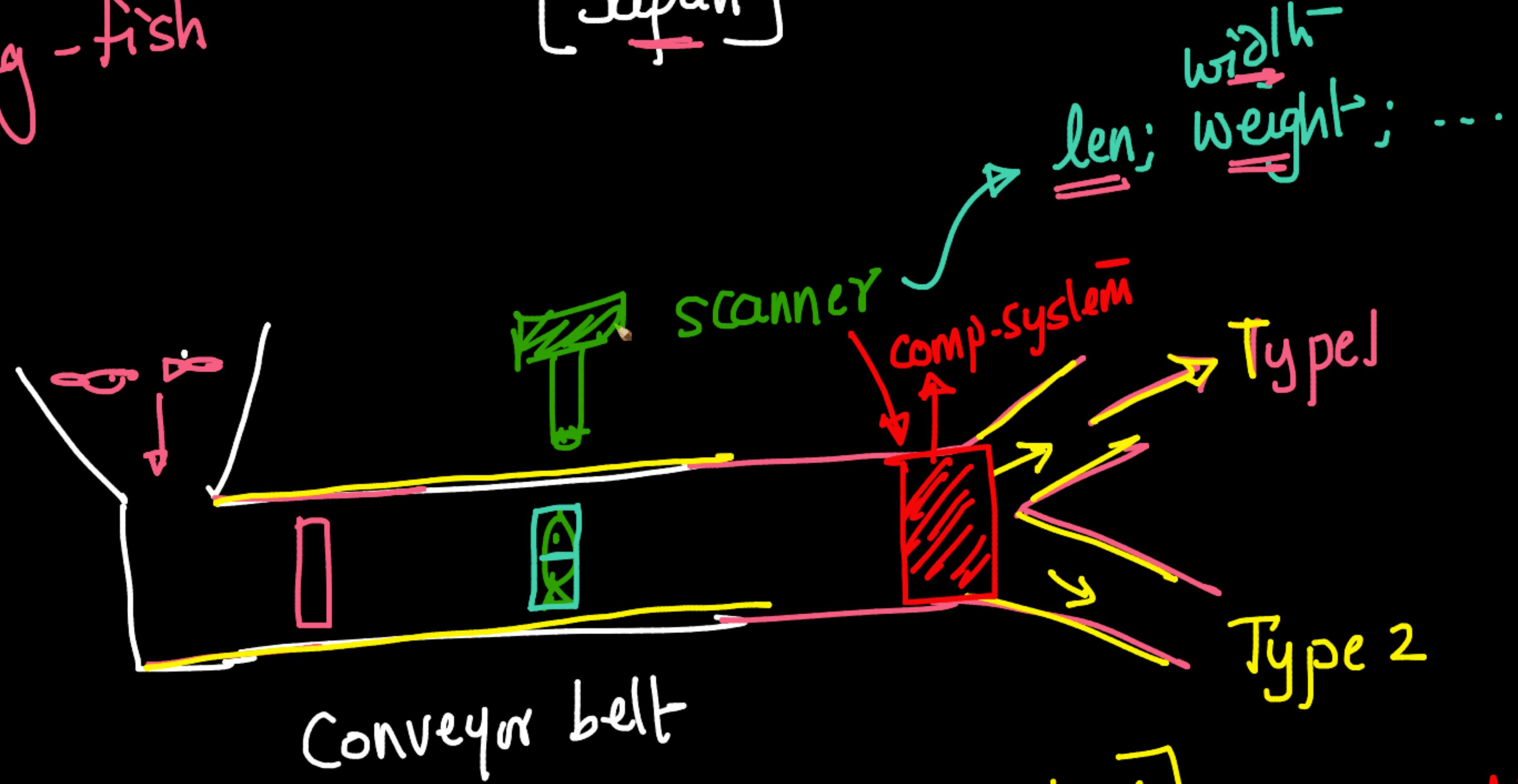
circle; ellipses; spheres



real-world

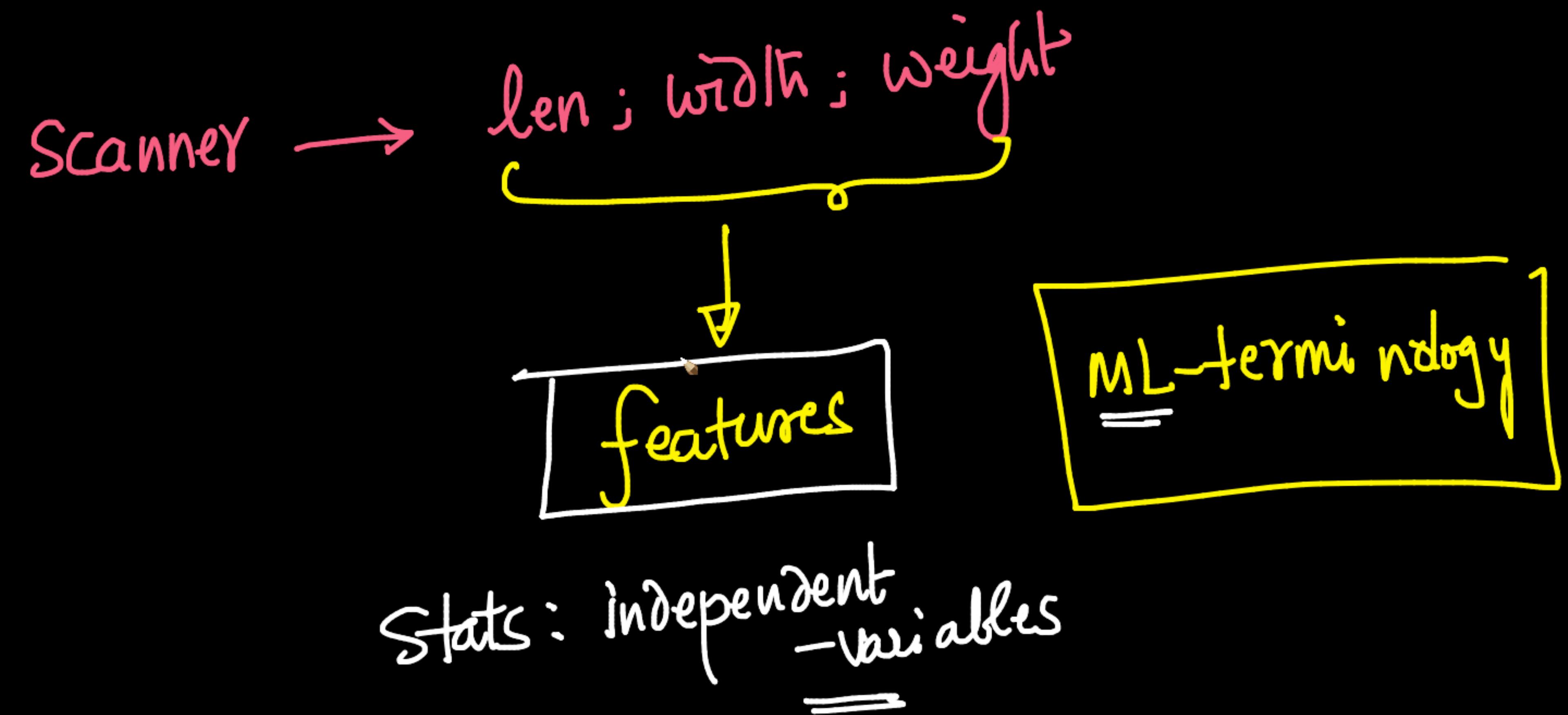
Sorting-fish

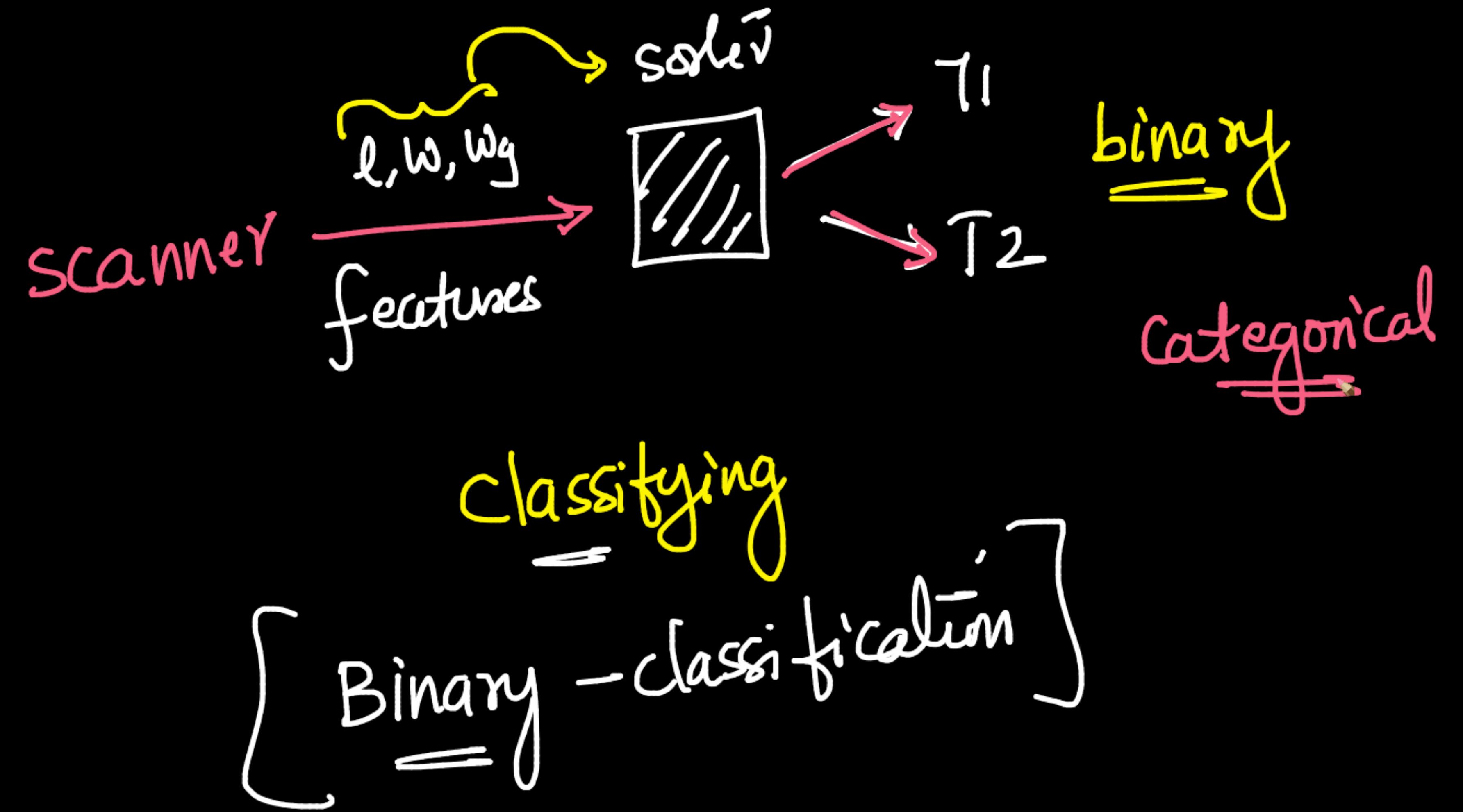
[Japan]

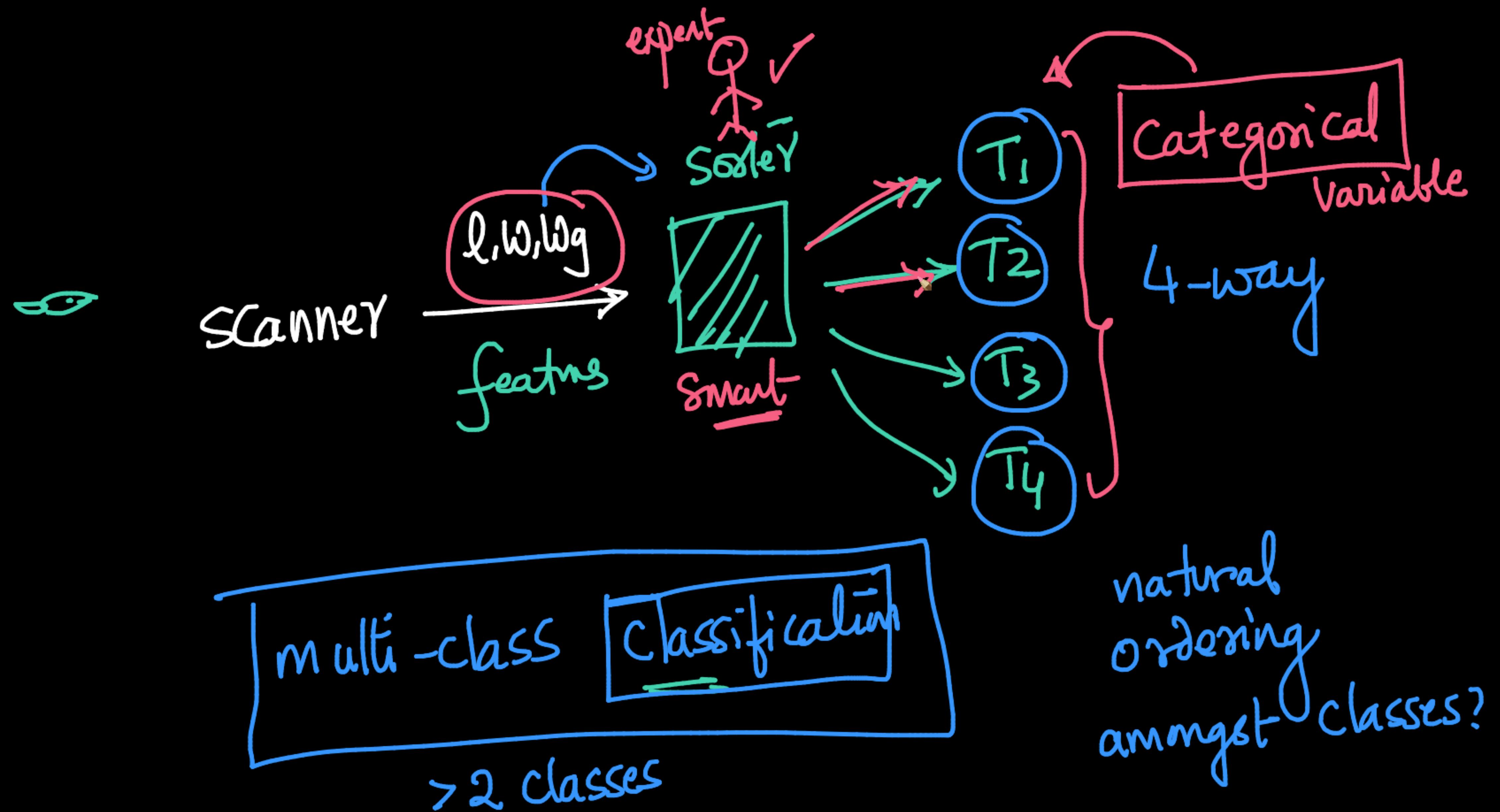


[credits: Duda - hart]

automated
sorted



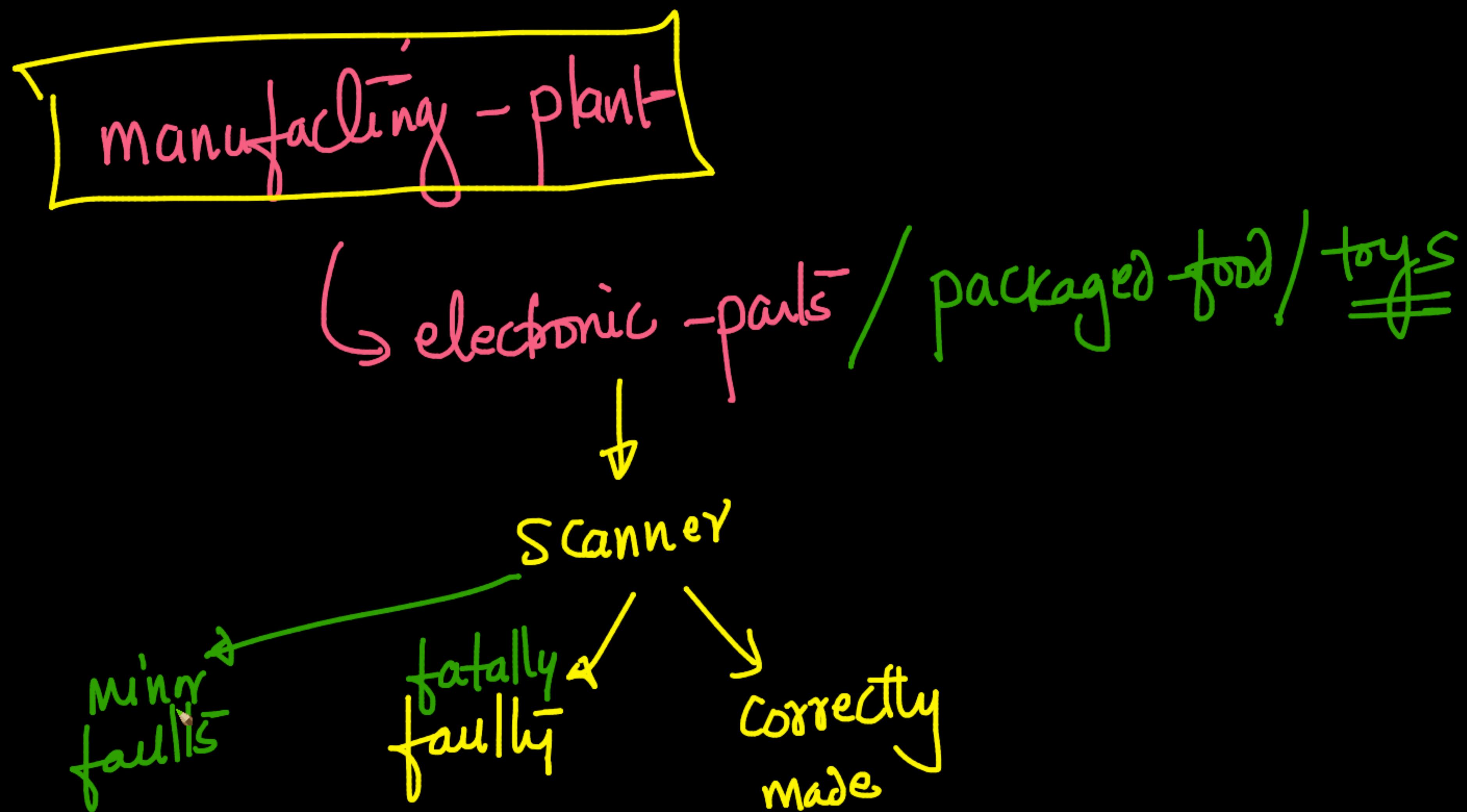




Predict age of fish using l, w, w_g

↓
numerical





Sorting fish → Linear Alg & Co-ordinate

(Q)

predict Marks = {1, 2, 3, 4 - - 100} → discrete
; features

regression or

multi-class classification

practical →

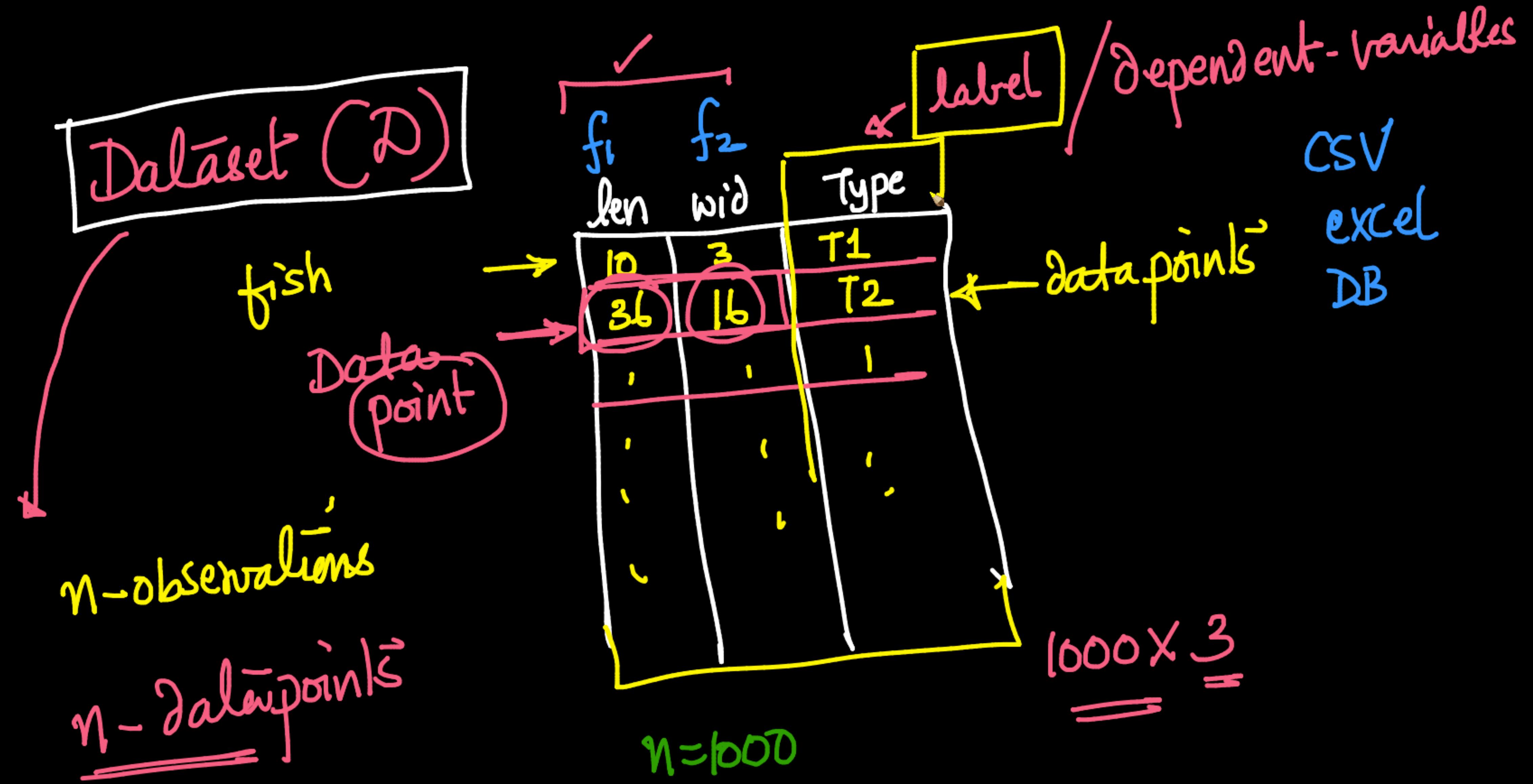
natural ordering

no ordering amongst the classes

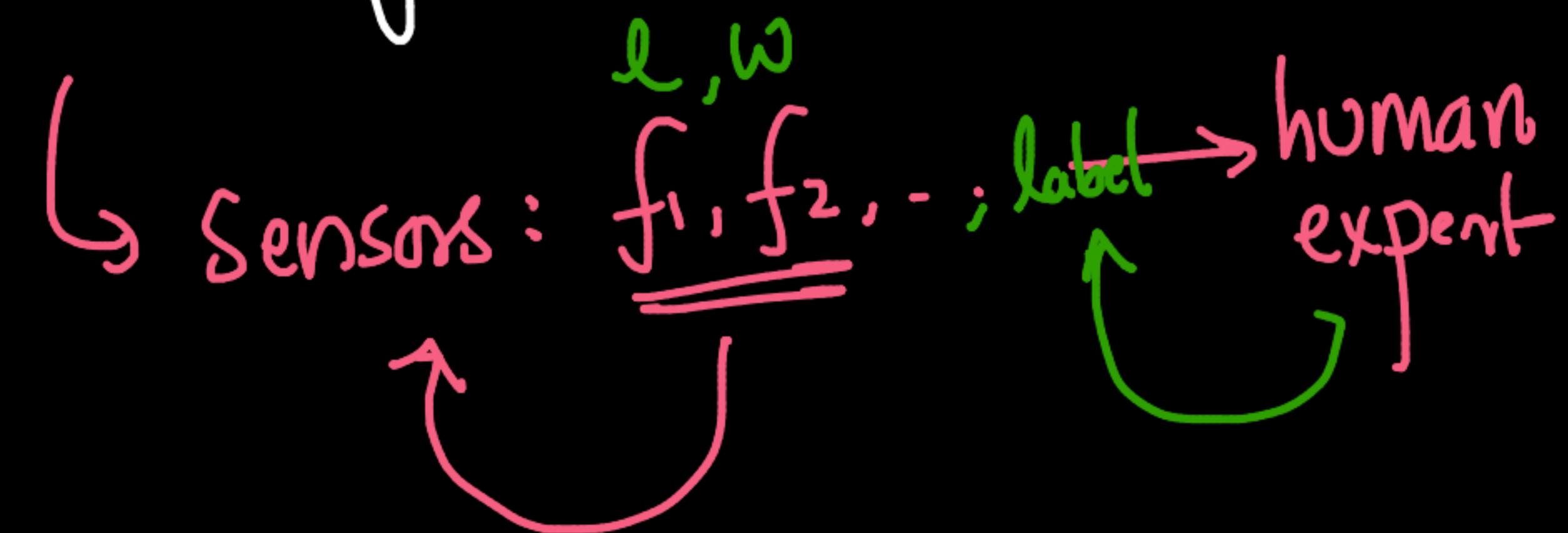
round(3.2)

→ $\boxed{3.2}$

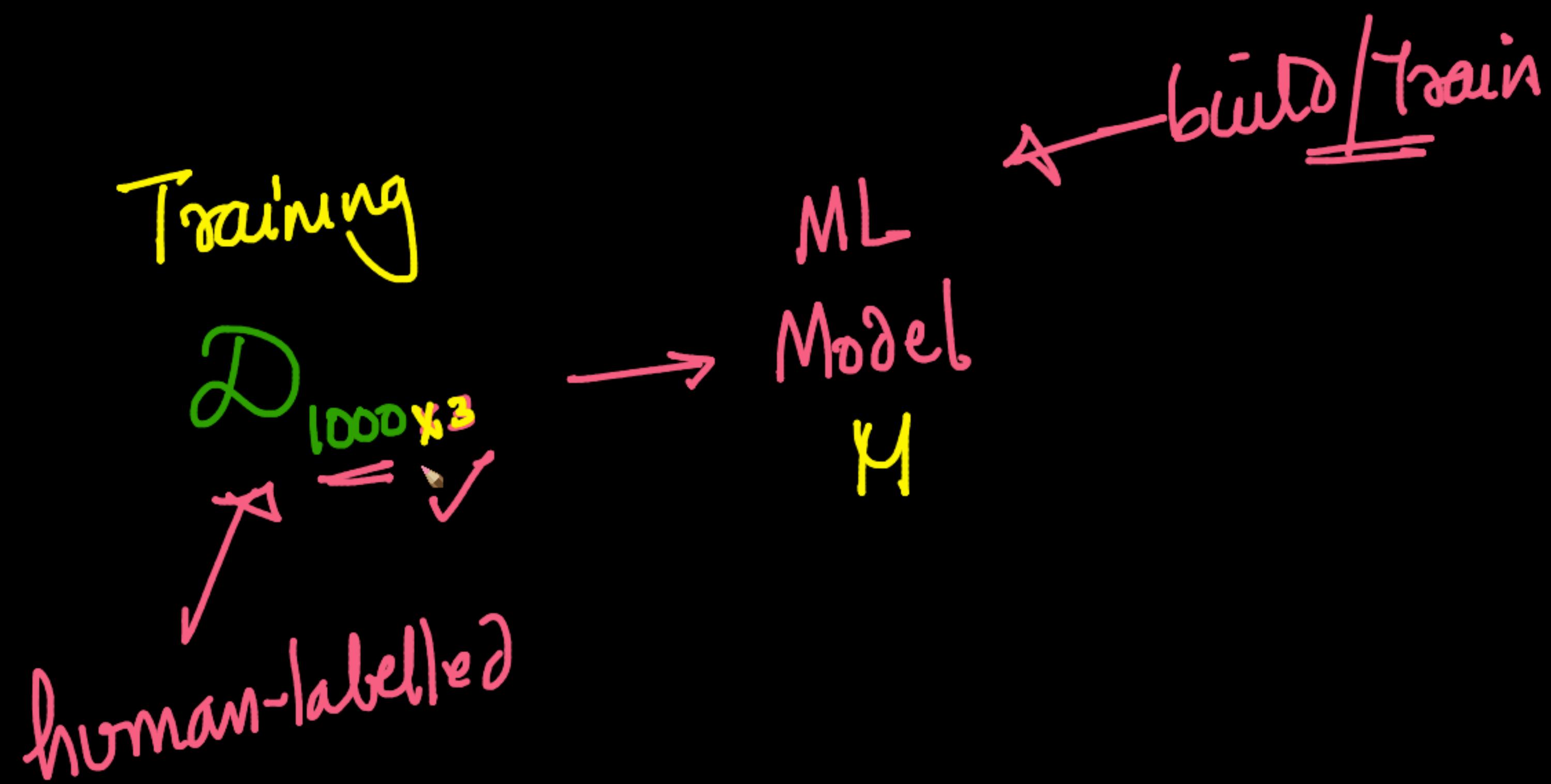
3

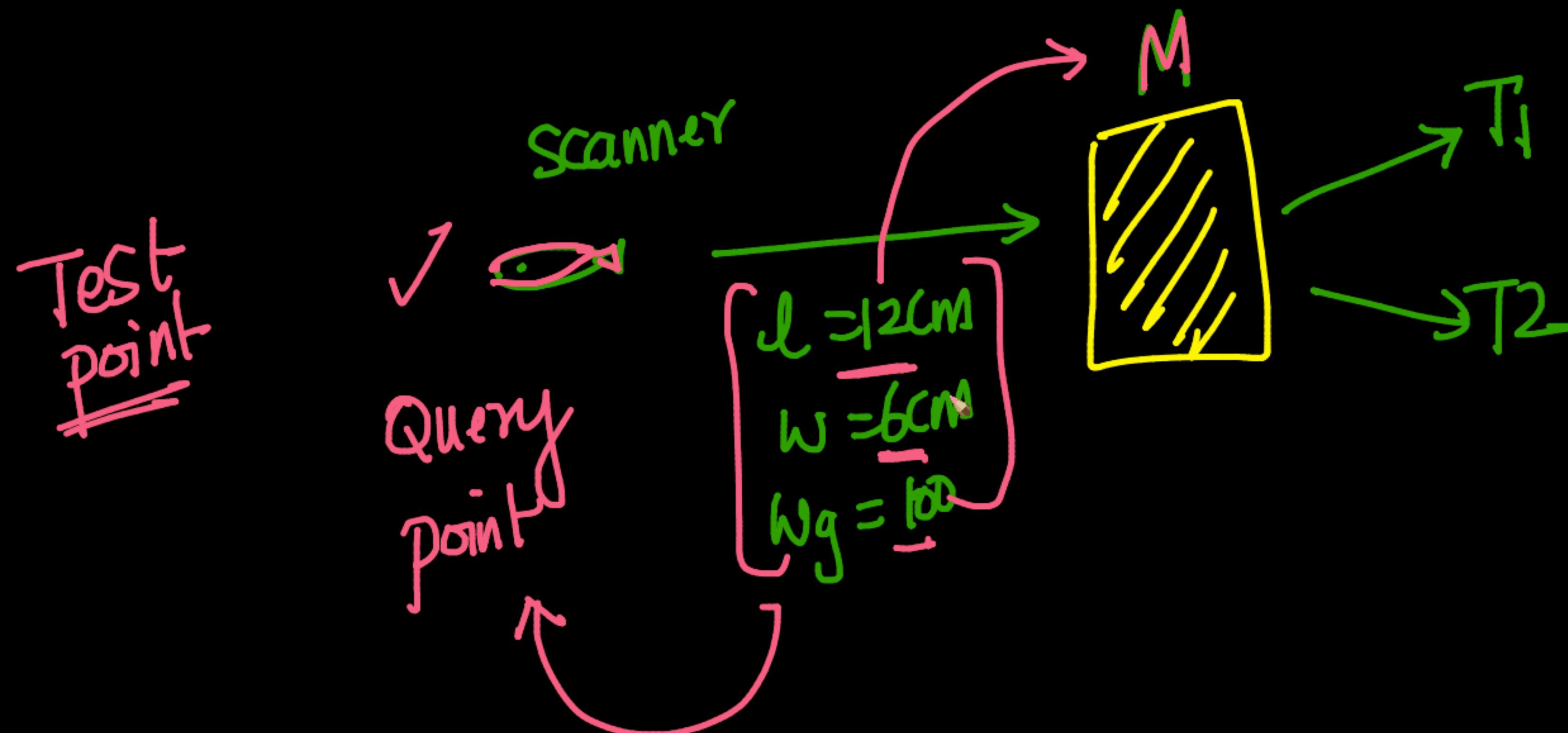


(Q) How do we get the data?



human-labelled
(annotated-data)



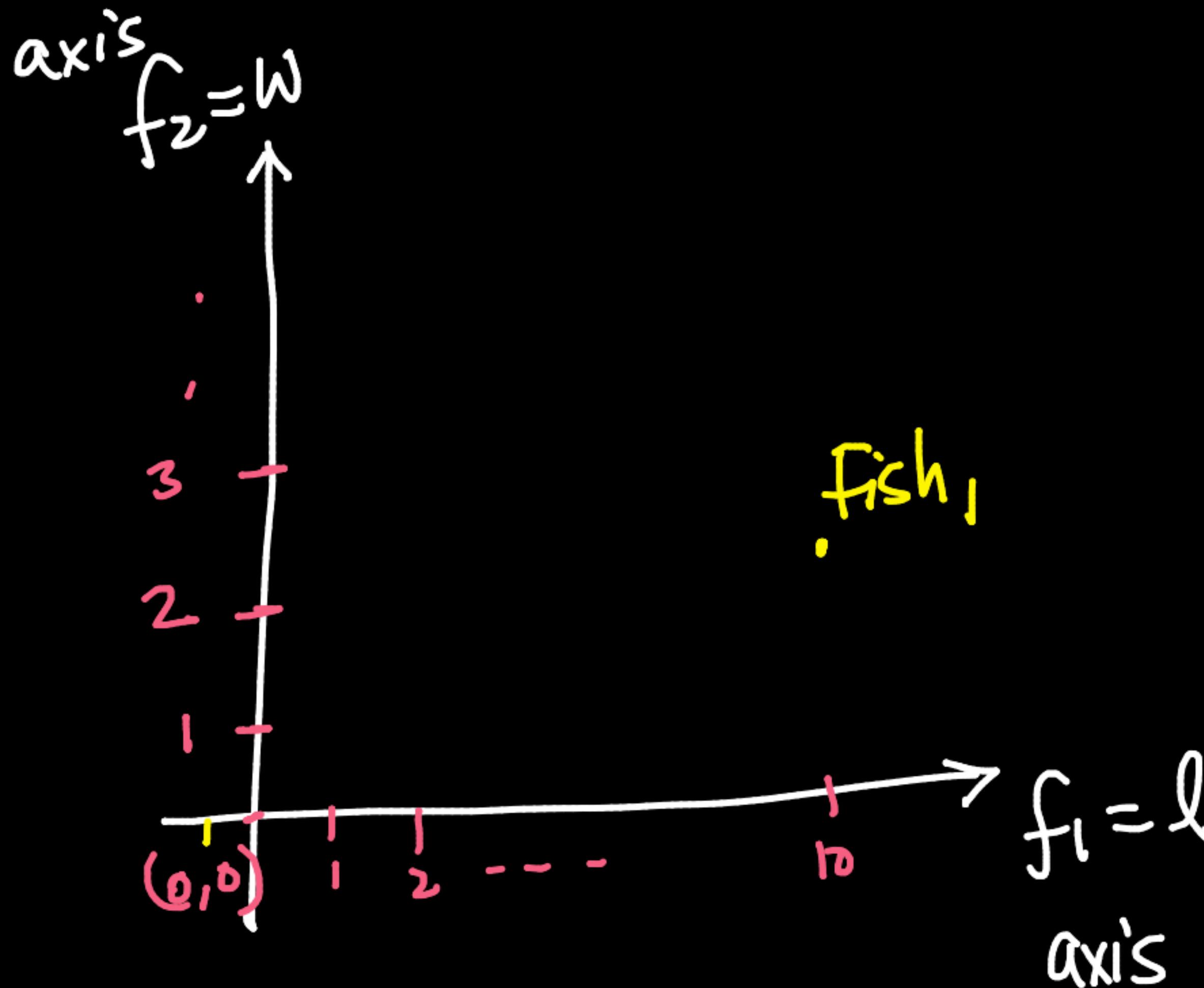


2D -
geometrī

$$\underline{\text{Fish}}_1 = \begin{bmatrix} l \\ w \end{bmatrix}$$

data point

data vector



InterviewBit Software x Live | DSML Advan x dsml-course/Coordina x dsml-course/01-Scala x dsml-course/15-Hypo x Gamma distribution - x box-cox.ipynb - Colab x Untitled1.ipynb - Colab x scipy.stats.boxcox - x +

scaler.com/meetings/i/dsml-advanced-linear-algebra-1-remedial/live

DSML Advanced : Linear Algebra 1- Remedial | Lecture

GEO MRTT

CLT

asy. c. l m mean

percentile-based C.I → Mean

=

You are sharing your screen now

Stop Sharing

Sample Means \sim Gaussian

Ankit Gupta

Srikanth Varma Chekuri (Host, You)

Rishav Kumar

Adarsh Vinayak

Aditya Yadav

Amit Srivastava

Arijit Bhowmick

Atin Gupta

Avijit Das

Gowtham

Hemant Parashar

MANISH PATEL

Mohammad Rashid

Muthu kamalan .M

Navroop

Nikunj Prajapati

Raghu Vamsi Yaram

Rahul Maramreddy

Start Doubt Session

02:59:05

65 / 65

ing your screen. Stop sharing Hide

People

Search

Chat

Questions

30 People

InterviewBit Software x Live | DSML Advan x dsml-course/Coordina x dsml-course/01-Scala x dsml-course/15-Hypo x Gamma distribution - x box-cox.ipynb - Colab x Untitled1.ipynb - Colab x scipy.stats.boxcox - x +

scaler.com/meetings/i/dsml-advanced-linear-algebra-1-remedial/live

DSML Advanced : Linear Algebra 1- Remedial | Lecture

large
K-resamplings → Median

CLT ≈ percentile based ✓

You are sharing your screen now

Stop Sharing

Sample-Size = $\tilde{n} = \underline{\underline{100}}$

size of bootstrap sample = $M = \underline{\underline{1MM}}$

Sameer Udgirkar has raised hand Unmute

03:01:41

Ankit Gupta

Srikanth Varma Chekuri (Host, You)

Aditya Yadav

Amit Srivastava

Arijit Bhowmick

Atin Gupta

Avijit Das

Gowtham

MANISH PATEL

Mohammad Rashid

Muthu kamalan .M

Navroop

Nikunj Prajapati

Raghu Vamsi Yaram

Rahul Maramreddy

Rahul Shivani

Start Doubt Session

66 / 67

Stop sharing Hide

People

Search

Chat

Questions

28 People

27
People

Search

SV Srikanth Varma Chekuri (Host, You)

AG Ankit Gupta

RK Rishav Kumar

SD Sameer Udgirkar

AV Adarsh Vinayak

AY Aditya Yadav

AS Amit Srivastava

AB Arijit Bhowmick

AG Atin Gupta

AD Avijit Das

G Gowtham

MP MANISH PATEL

MR Mohammad Rashid

MK Muthu kamalan .M

N Navroop

RY Raghu Vamsi Yaram

RM Rahul Maramreddy

RS Rahul Shivani

Start Doubt Session

DSML Advanced : Linear Algebra 1- Remedial | Lecture

✓ { left-skewed
Not-normal } → 2-sample t-test

right-skewed

You are sharing your screen now

Stop Sharing

H_a ①

t-statistic ②

{ behavior ③
dist of t-statistic }



Srikanth Varma Chekuri (You)

Srikanth Varma Chekuri (You) (Screen)

InterviewBit Software x Live | DSML Advan x dsml-course/Coordina x dsml-course/01-Scala x dsml-course/15-Hypo x Gamma distribution - x box-cox.ipynb - Colab x Untitled1.ipynb - Colab x scipy.stats.boxcox - x +

scaler.com/meetings/i/dsml-advanced-linear-algebra-1-remedial/live

DSML Advanced : Linear Algebra 1- Remedial | Lecture

left-skewed data → boxcox → ~Normal

You are sharing your screen now

Stop Sharing

✓ { Ha: t-statistic
t-distr

People

Search

SV Srikanth Varma Chekuri (Host, You)

AG Ankit Gupta

RK Rishav Kumar

Sameer Udgirkar

AV Adarsh Vinayak

Aditya Yadav

AS Amit Srivastava

Arijit Bhowmick

Atin Gupta

Avijit Das

G Gowtham

MANISH PATEL

Muthu kamalan .M

Navroop

Raghu Vamsi Yaram

Rahul Maramreddy

Ranita Bhattacharya

Rishabh Singh

Start Doubt Session

03:05:25

68 / 69

Stop sharing Hide

scaler.com/meetings/i/dsml-advanced-linear-algebra-1-remedial/live

You have left the meeting

We get frequent requests for your notes!

Notes written by you helps in understanding the topic better. You can upload the notes in two simple steps mentioned below

- 1 Scan the QR code with your iPad
Scanner should be present in the top menu on your iPad
- 2 Upload Notes on the generated link
All notes uploaded will be visible in the saved version of this session

OR

Drag and drop files or [click here to upload](#)

Files Uploaded from your computer appear here

69 / 70