

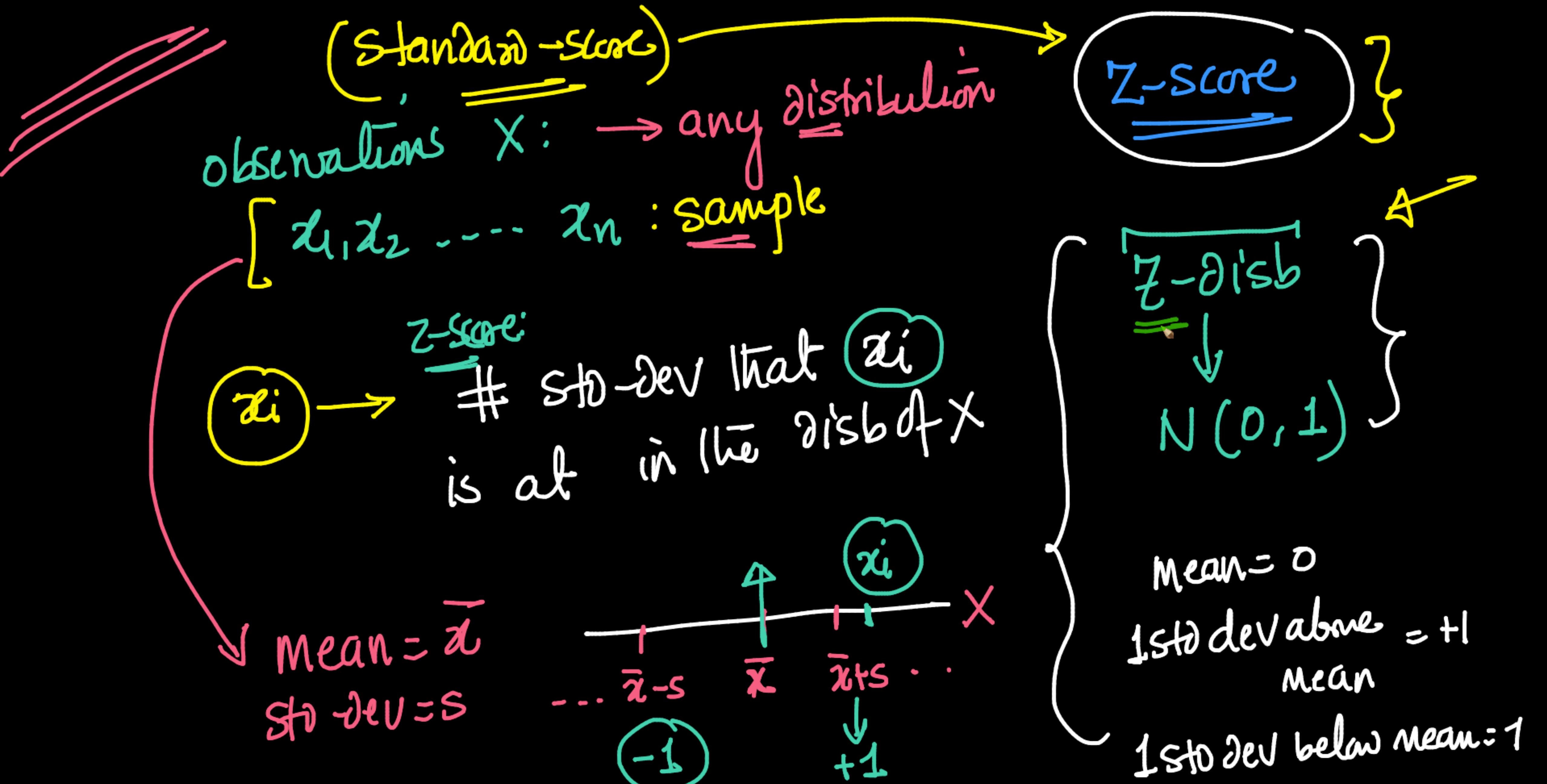
END: ...

Topics:

- Misc: (case-study)
 - { - Z-score
 - Feature Transforms
 - Box-Cox Transform
 - Extended-Bayes Thm
- KS-test ✓
- T-test ✓
- A/B Testing in real-world }

{ Framework
Z-test

- Z-Proportions
= test



✓ Standardization

Z-score (x_i) =

$$\frac{x_i - \bar{x}}{s}$$

any new-term → Google search away

recovery times medicine 1 $\rightarrow \times$

ill patient : 16.5 days

Z-score: 3 -0.5

The diagram illustrates the calculation of a Z-score. A central blue circle contains the number '3'. Two arrows point from this circle to two separate boxes: one containing '-0.5' and another containing 'Z-score'. The word 'Z-score' is written above the '-0.5' box.

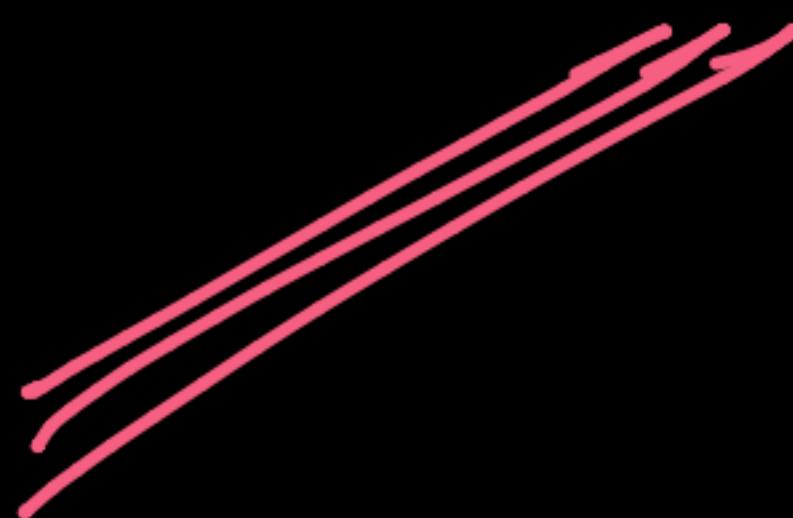


Z-score is -ve :

$$\frac{x_i - \bar{x}}{s} < 0 \Rightarrow \underline{\overline{x_i}} < \bar{x}$$

mean

"Feature" Transforms:



non-gaussian

X:

x_1, x_2, \dots, x_n

x_n : delivery times of Amazon

Stats

$f(\cdot)$

right-skewed ; could be log-normal (let)

$$f(x_i) = y_i$$

{ y_i : y_1, y_2, \dots, y_n } $\stackrel{=}{\sim}$ y_n Gaussian

X ; log-normal

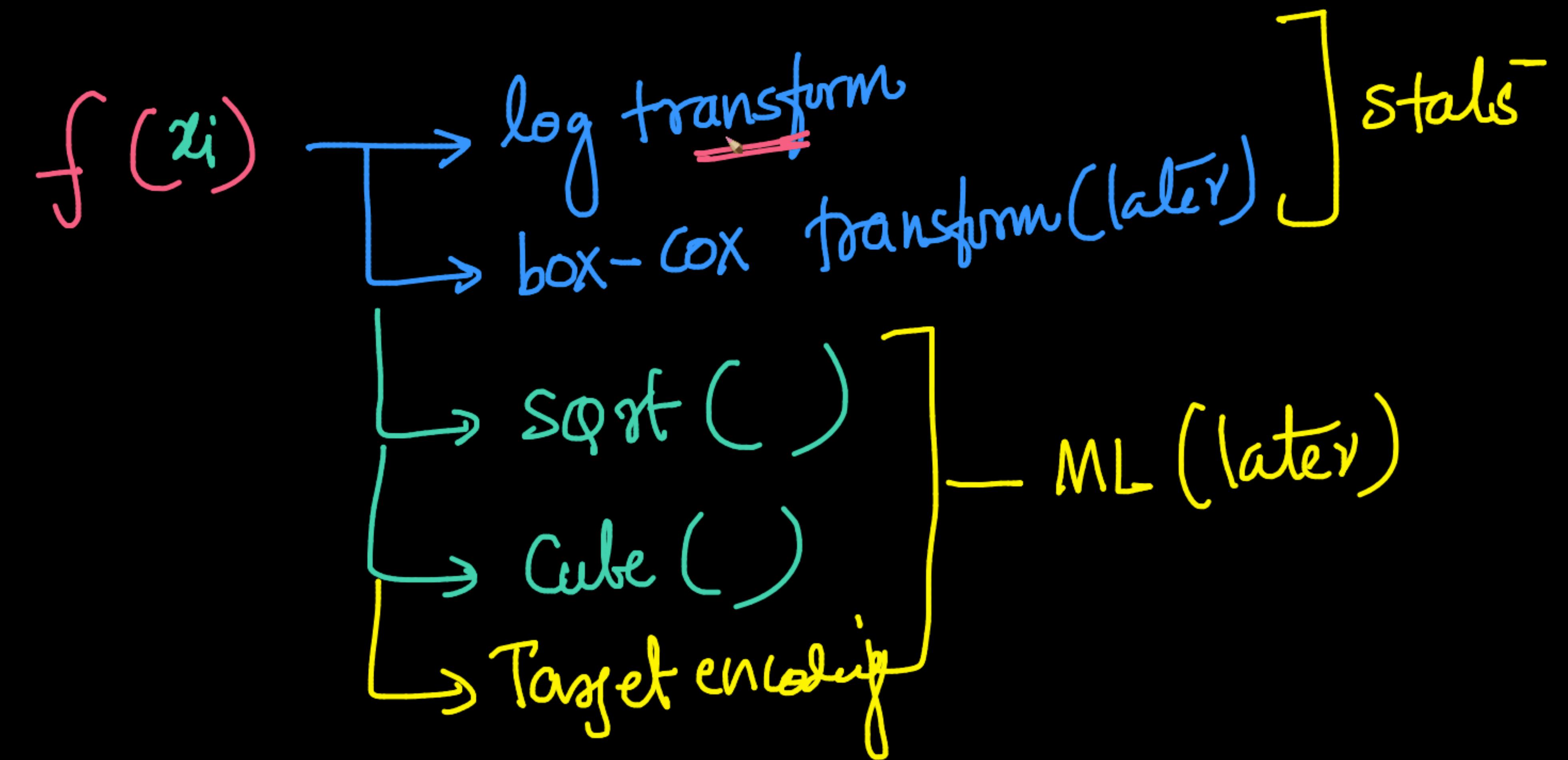
$$\downarrow \quad f(x_i) = y_i$$

$$f(x) = y$$

y : Gaussian

f: log

log-transform



popular transforms

Task:

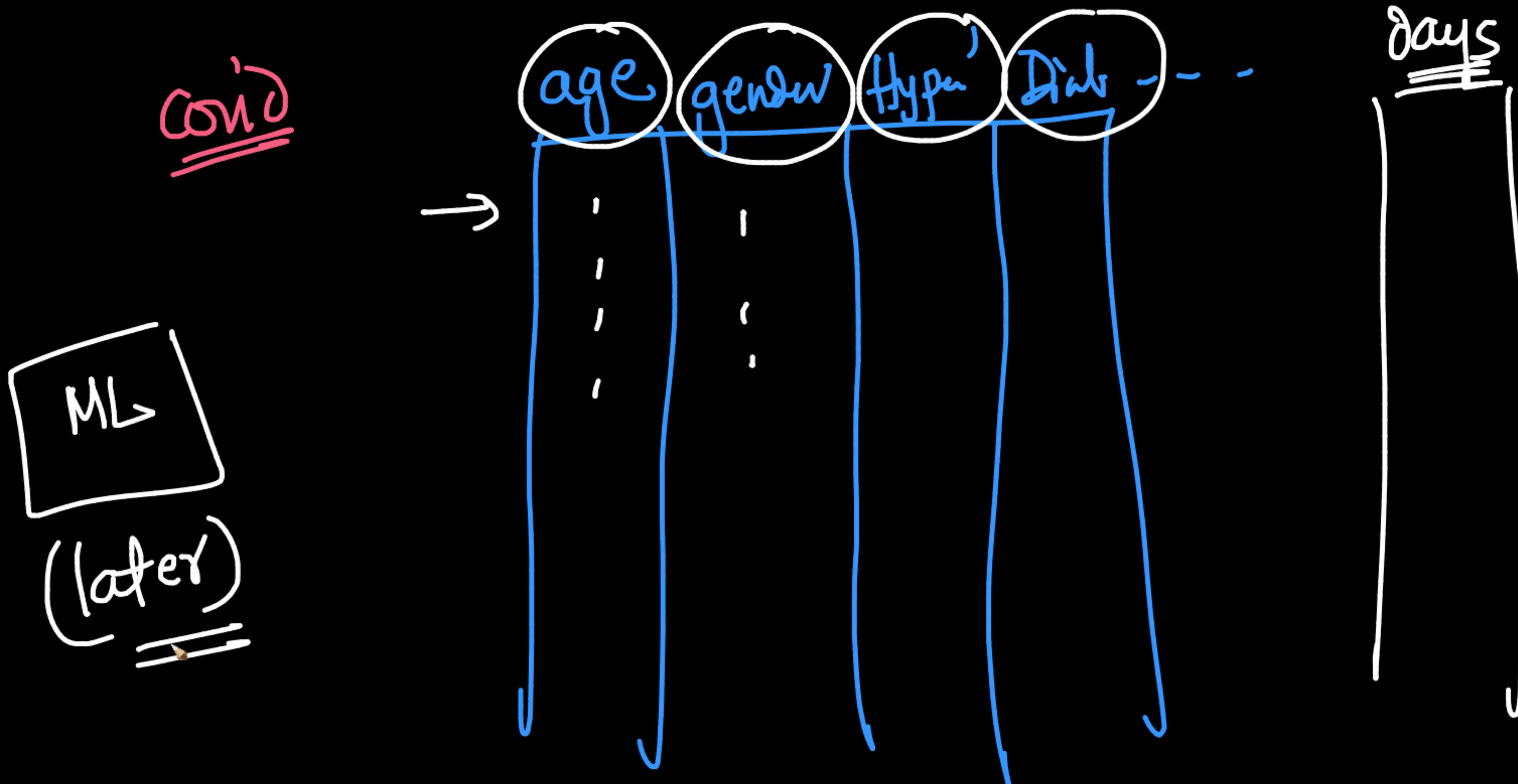
non-Gaussian - Gaussian.

$$X \rightarrow \sqrt{(x_i - \mu)^2} \quad f(x_i) = y_i$$

mean = 0 ; std dev = 1
dist may not change

$$\rightarrow \text{Min-Max} [0, 1] \quad \frac{x_i - \text{Min } x}{\text{Max } x - \text{Min } x} \quad f(x_i) \sqrt{\cdot}$$

[its not guaranteed to convert to Gaussian]



~~SciPy.stats~~

Box-Cox:

- tries its best
- not guaranteed to always work

conditions



generalization of log-transform

$$x_1, x_2, \dots, x_i, \dots, x_n$$

$$y_1, y_2, \dots, y_i, \dots, y_n \sim \text{Gaussian}$$

→ QQ-plot (or) KS-test (later)

n - x_i 's

compute 

→ optimization
technique -
(MLE)

s.t.

$$\hat{y}_i = \begin{cases} \frac{x_i - \bar{x}}{\sigma} & \text{if } \lambda \neq 0 \\ \ln x_i & \text{if } \lambda = 0 \end{cases}$$

Gaussian

log-likelihood

Much
(later)

~~γ_{reg-2}~~

$$\left[\begin{array}{c} x_1 \ x_2 - \bar{x} \ x_n \\ \hline \end{array} \right] \xrightarrow{\downarrow f(x)} \left[\begin{array}{c} y_1 \ y_2 - \bar{y} \ y_n \\ \hline \end{array} \right]$$

~~Edges~~ ✓
~~ANOVA~~ (next classes)

↓
obs to be Gaussian
~~dist~~

▷ ↗

$$y_i = \frac{x_i - 1}{\sigma}$$

$$f_x(x_i) = y_i$$

✓

$$\exp \left\{ \log(\lambda y_i + 1) \right\} = \cancel{\lambda y_i} x_i$$

$$f^{-1}(y_i) = x_i$$

$d_1 \rightarrow x_{11}, x_{12}, \dots, x_{1m} \rightarrow y_{11}, y_{12}, \bar{y}_1, y_{1m}$

$d_2 \rightarrow x_{21}, x_{22}, \dots, x_{2m} \rightarrow y_{21}, y_{22}, \bar{y}_2, y_{2m}$

$\vdots \quad \vdots$

d_{10}

```
# gather some data
data =[ 0.04177737,  0.97977259,  1.19684675,  0.75969411,  0.2772351 ,
        1.20400739,  1.19512711, -1.33315966,  0.47241401,  0.58453053,
        0.21167461,  0.87106215, -0.56663286,  0.3702523 ,  0.72724427,
        0.41126015,  0.33358864,  0.72878097,  0.69929305,  0.72581333,
        1.67334826, -1.54572083, -1.22840893,  0.47103287,  0.895276 ,
        0.16538052, -0.43575904,  1.62784202,  0.98340417,  0.90482144,
       -0.47914975,  0.71812022,  1.14243 , -0.04393411,  1.24946471,
       -0.8699551 ,  1.60196517,  1.00140898,  1.48233878, -0.37088602,
       -0.0954339 ,  1.2969551 ,  0.0457524 , -0.06486335,  0.43257115,
       -0.18945797,  0.46525944,  0.12974487, -0.10501035,  0.94060547,
       -1.57714093,  0.24292938,  0.68759359,  0.24113398,  0.74353881,
        0.0129037 ,  0.47936105, -0.0596165 ,  0.3300311 , -0.19409805,
       -2.15213968, -0.9169724 ,  1.40476752,  0.74067023,  0.36119747,
        1.04507563, -0.54692221,  0.65000261,  0.5359208 ,  0.40091749,
        0.16959609,  0.43828974,  1.69191812, -0.40588725,  0.52772481,
        0.2410331 ,  1.8226663 , -1.36677194,  0.41745297,  0.94050797,
        1.15797033,  0.13883716,  0.9648131 ,  0.71495948,  1.73284151,
        0.9571359 ,  0.38785662,  0.41200029,  1.10201074,  0.41262702]
```

∞ boxcox and z-score.ipynb - Colaboratory | ⚡ scipy.stats.boxcox — SciPy v1.8.0 | W Kolmogorov-Smirnov test - Wikipedia | ∞ KTest_Ttest.ipynb - Colaboratory | +

colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=e3U8mXjRR7uL

+ Code + Text

RAM Disk

[4] 0s

{x}

#Z-score
i=100
xi= x[i] ≈ 10
print(xi)

print ("zscore:", (xi-np.mean(x))/np.std(x))

-1.59860382
zscore: -2.341243992536608

$x_{100} = -1.598$

$z\text{score}_{100} = -2.34$

qq plot of x vs normal
sm.qqplot(x, stats.norm, fit=True, line="45")

boxcox and z-score.ipynb - Colab

scipy.stats.boxcox — SciPy v1.8.0

Kolmogorov-Smirnov test - Wikipedia

KTest_Ttest.ipynb - Colaboratory

colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=e3U8mXjRR7uL

Update

+ Code + Text

RAM Disk

```
xi= x[i]
print(xi)

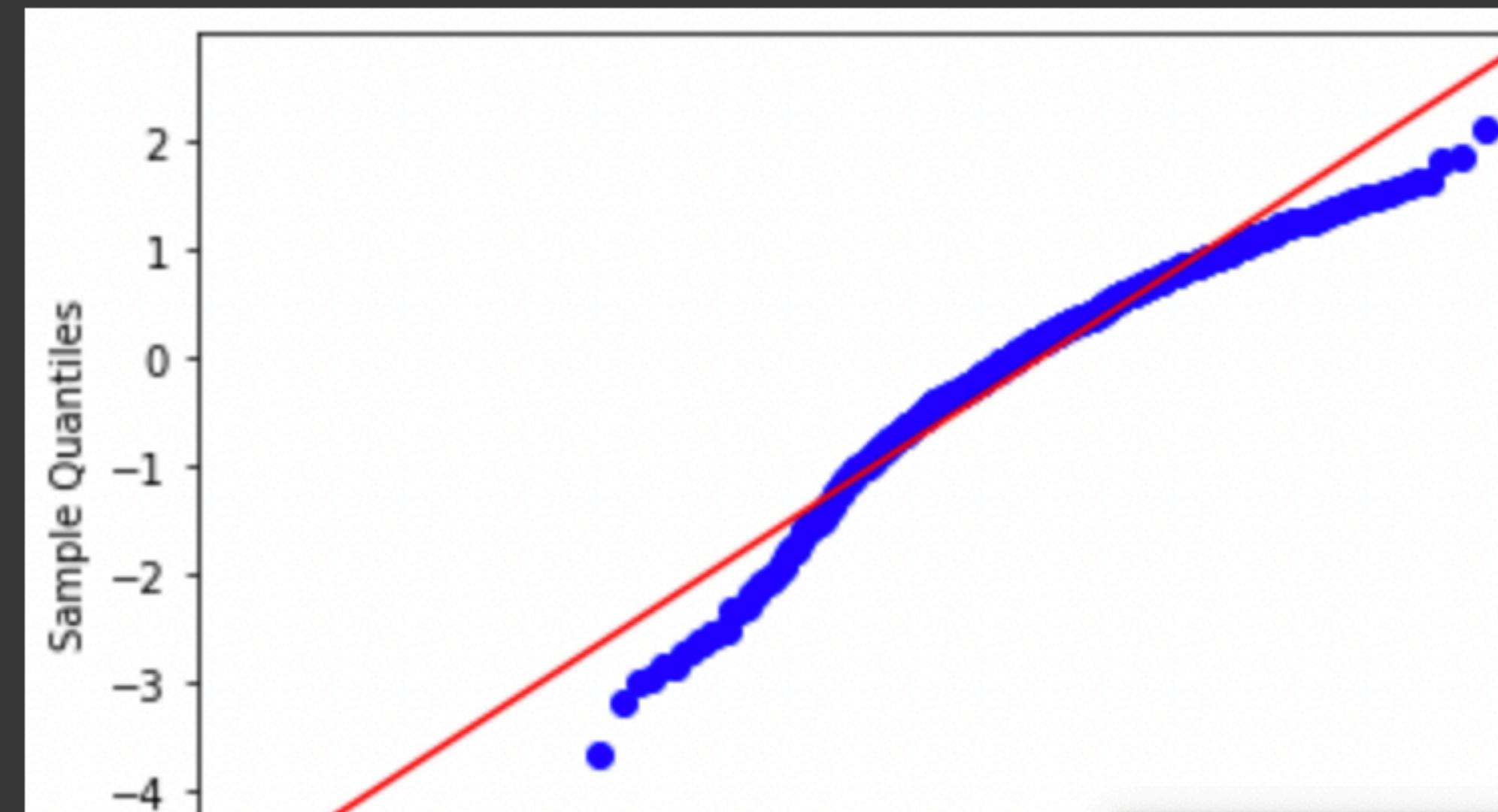
print ("zscore:", (xi-np.mean(x))/np.std(x))
```

-1.59860382

zscore: -2.341243992536608

↑ ↓ ⌂ ⚙ 📈 🗑 ⏺

```
# qq plot of x vs normal
sm.qqplot(x, stats.norm, fit=True, line="45")
```



boxcox and z-score.ipynb - Co

scipy.stats.boxcox — SciPy v1.8

Kolmogorov-Smirnov test - Wi

KTest_Ttest.ipynb - Colaborat

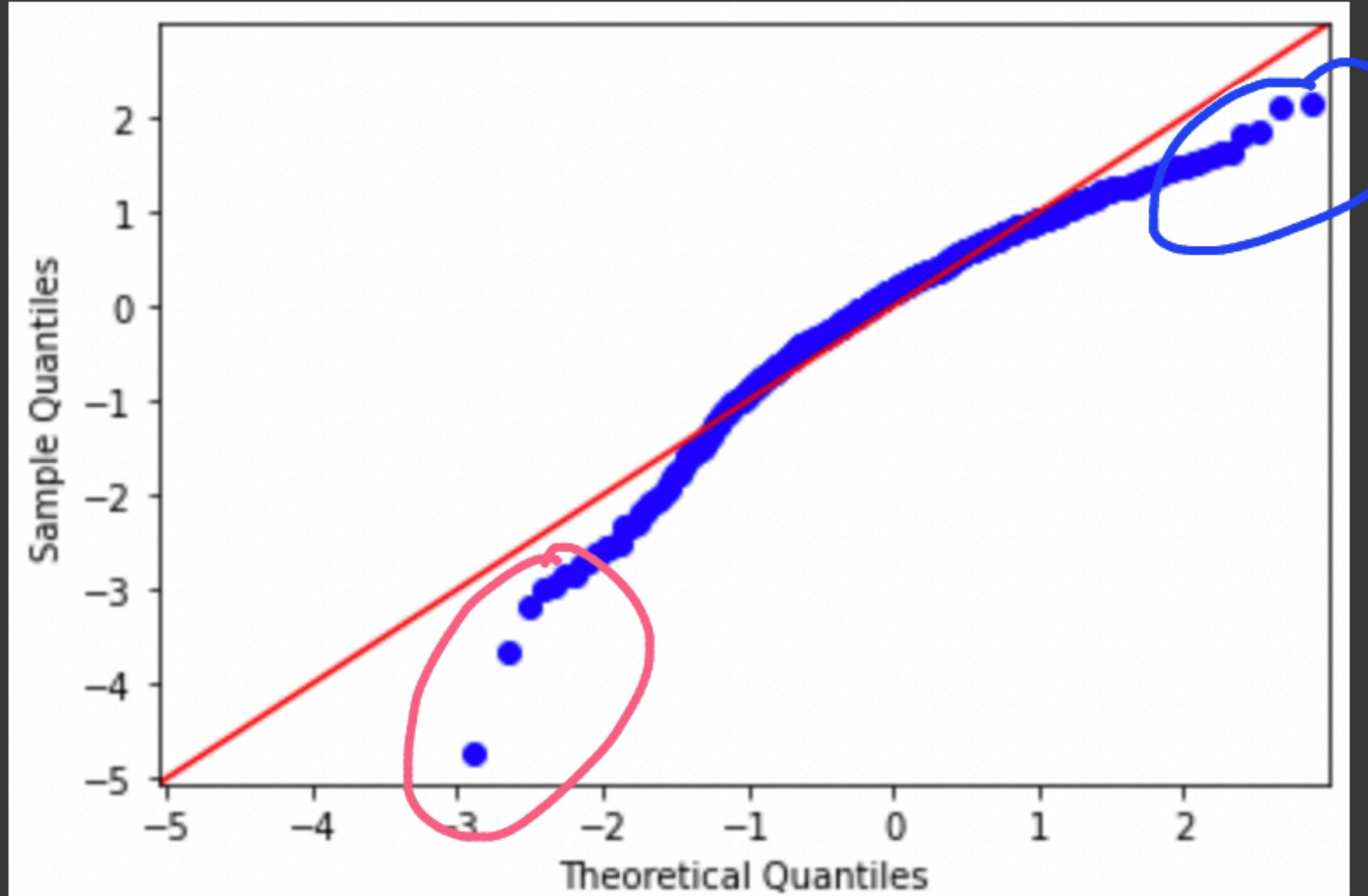
colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=e3U8mXjRR7uL

Update

+ Code + Text

RAM Disk

```
[ ] # qq plot of x vs normal  
sm.qqplot(x, stats.norm, fit=True, line="45")
```



$\times \not\sim \text{Gaussian}$



CO |

boxcox and z-score.ipynb - Co

 `scipy.stats.boxcox` – SciPy v1

Kolmogorov-Smirnov test - Wil

KStest_Ttest.ipynb - Colaborat

 Update 

+ Code + Text

✓ RAM Disk

Theoretical Quantiles

↑ ↓ ⌂ ☰ 🗃 🗑 ⋮

```
[ ] np.min(x)
```

$$\min x = \overbrace{-3.63}^{\sim} \dots$$

-3.63123097

2

```
# box cox transform  
x1 = x + 3.7  
xt, l = stats.boxcox(x1,); # returns x_trasnfomed and lambda  
print("lambda :" + str(l))
```

da box-tot $(x_i + 3 \cdot 7)$

```
# check if xt is gaussian or not using QQ-Plot  
sm.qqplot(xt, stats.norm, fit=True, line="45")
```

λ → lambda : 2.2233087629443724



A diagram illustrating a function f and its inverse f^{-1} . At the top, a pink bracket on the left groups the input x_i and the output y_i . Above x_i , a white arrow points to y_i , with the label $f()$ written in pink above the arrowhead. Below this, a green curved arrow starts at y_i and points back up to x_i , with the label $f^{-1}()$ written in pink below it.

∞ boxcox and z-score.ipynb - Colaboratory | ⚡ scipy.stats.boxcox — SciPy v1.8.0 | W Kolmogorov-Smirnov test - Wilcoxon signed-rank test.ipynb - Colaboratory | +

colab.research.google.com/drive/1HBSL2rle-pokrz50QlhgWGuZiUaxAeUe#scrollTo=Lg17ejlFSxH2

+ Code + Text RAM Disk

Theoretical Quantiles

[] np.min(x)
{x}
-3.63123097

box cox transform
x1 = x + 3.7
xt, l = stats.boxcox(x1,); # returns x_transformed and lambda
print("lambda :" + str(l))

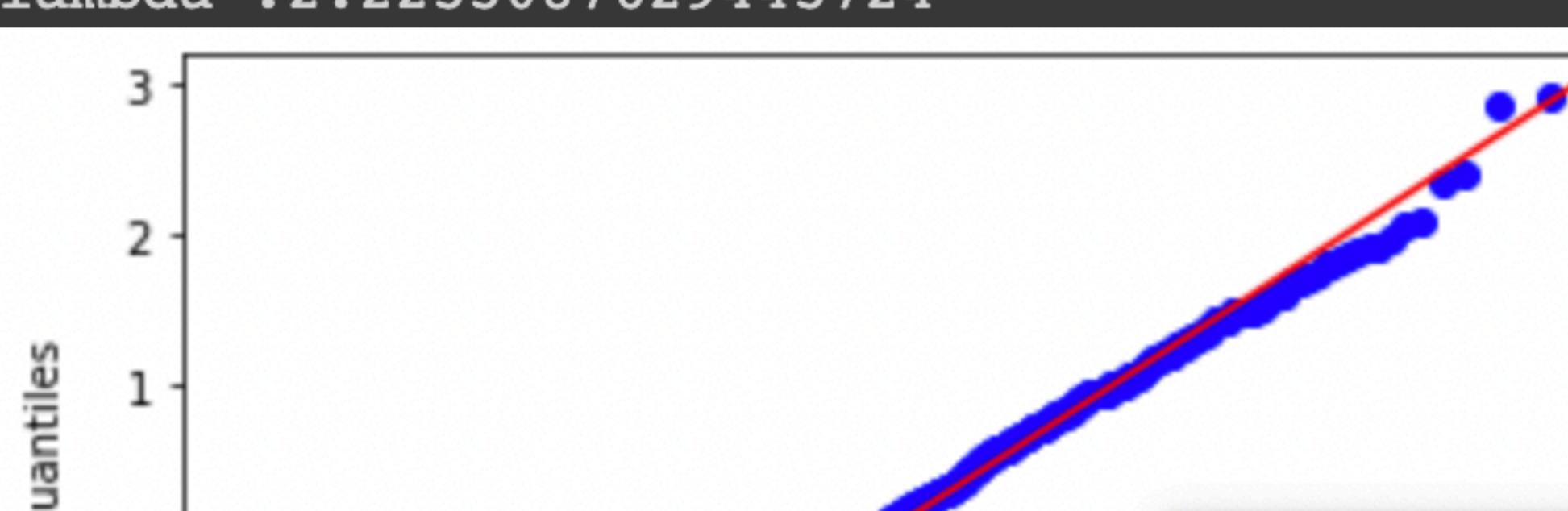
check if xt is gaussian or not using QQ-Plot
sm.qqplot(xt, stats.norm, fit=True, line="45")

lambda : 2.2233087629443724

$x_i + 3.7 \rightarrow h_i = 1 \rightarrow n$

$y_i | h_i; \lambda \rightarrow n \& \lambda \neq 0$

$\lambda = 2.223$



boxcox and z-score.ipynb - Col X

 [scipy.stats.boxcox](#) – SciPy

v1.1 | Kolmogorov-Smirnov test - V

KStest_Ttest.ipynb - Colaborat

Update

+ Code + Text

- ✓ RAM
- Disk

Theoretical Quantile

A set of small, light-gray navigation icons located at the bottom right of the page. From left to right, they include: a double arrow indicating a comparison or split-screen function; a magnifying glass for search; a circular arrow for refresh; a downward-pointing arrow for download; a square icon with a diagonal line for delete; and three vertical dots for more options.

```
[ ] np.min(x)
```

-3.63123097

```
▶ # box cox transform
```

$$x_1 = x + 3.7$$

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

```
# check if xt is gaussian or not using QQ-Plot  
sm.qqplot(xt, stats.norm, fit=True, line="45")
```

```
[1] lambda :2.2233087629443724
```



$$y_i = \frac{x_i - \bar{x}}{\sigma} \quad \text{if } \sigma \neq 0$$

二

4

$$\lambda_i + 3 \cdot 7^{2.223} - 1$$

+ Code + Text

- ✓ RAM
- Disk

Theoretical Quantiles



```
[ ] np.min(x)
```

-3.63123097

```
# box cox transform  
x1 = x + 3.7  
xt, l = stats.boxcox(x1,); # returns x_transformed and lambda  
print("lambda :" + str(l))
```

A hand-drawn illustration featuring a large, irregularly shaped speech bubble. Inside the bubble, the word "box-tox" is written in a bold, sans-serif font. Below "box-tox", there are two parallel horizontal lines. Outside the bottom right corner of the bubble, the phrase "tries its best" is written in a cursive, handwritten style. The entire drawing is done in yellow ink on a black background.

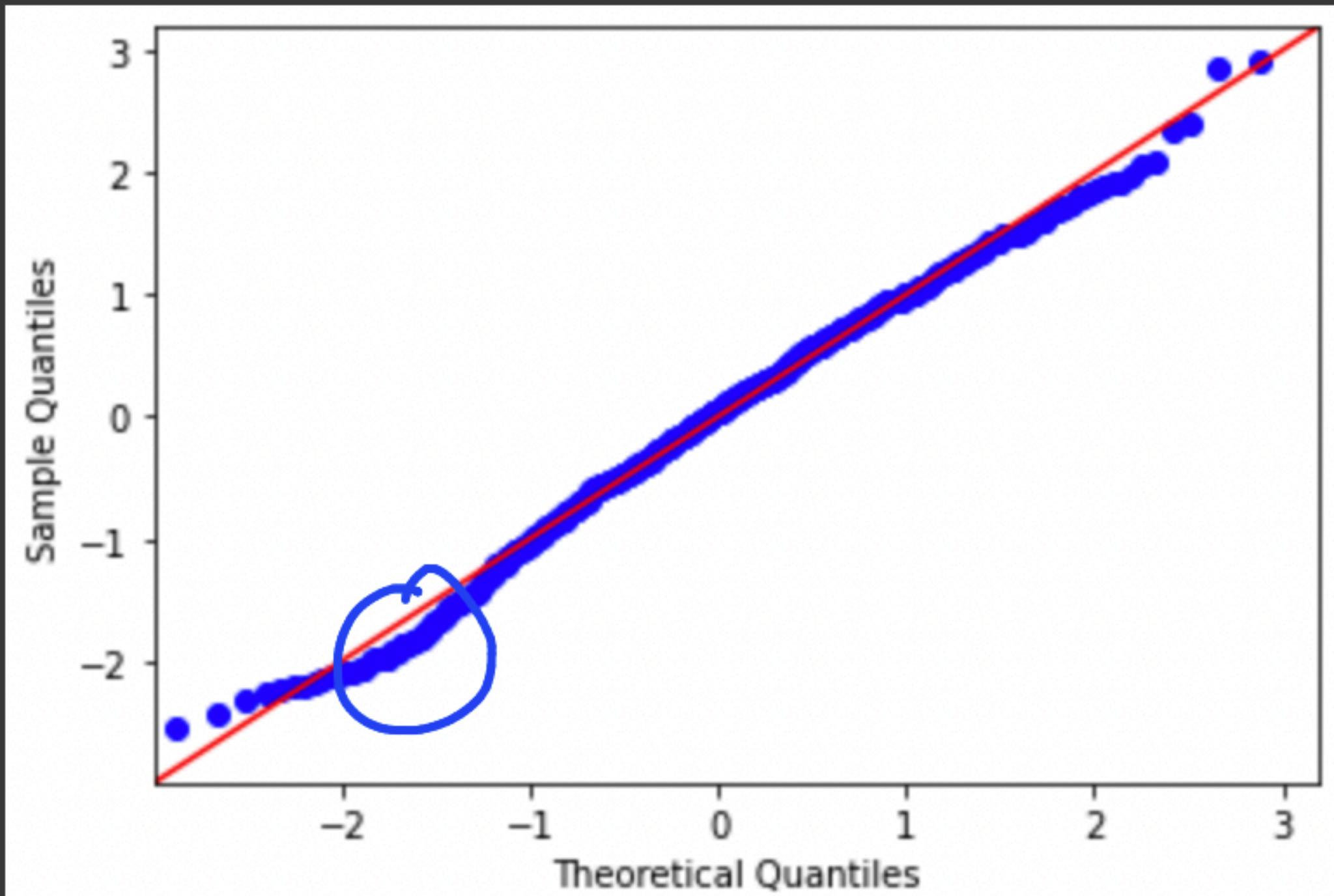
C → lambda : 2.2233087629443724



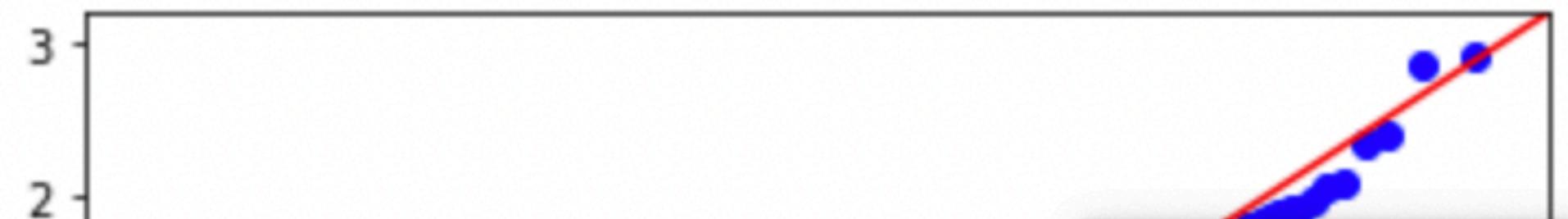
+ Code + Text

```
print("lambda :" + str(l))  
[ ]  
  
# check if xt is gaussian or not using QQ-Plot  
sm.qqplot(xt, stats.norm, fit=True, line="45")
```

lambda :2.2233087629443724



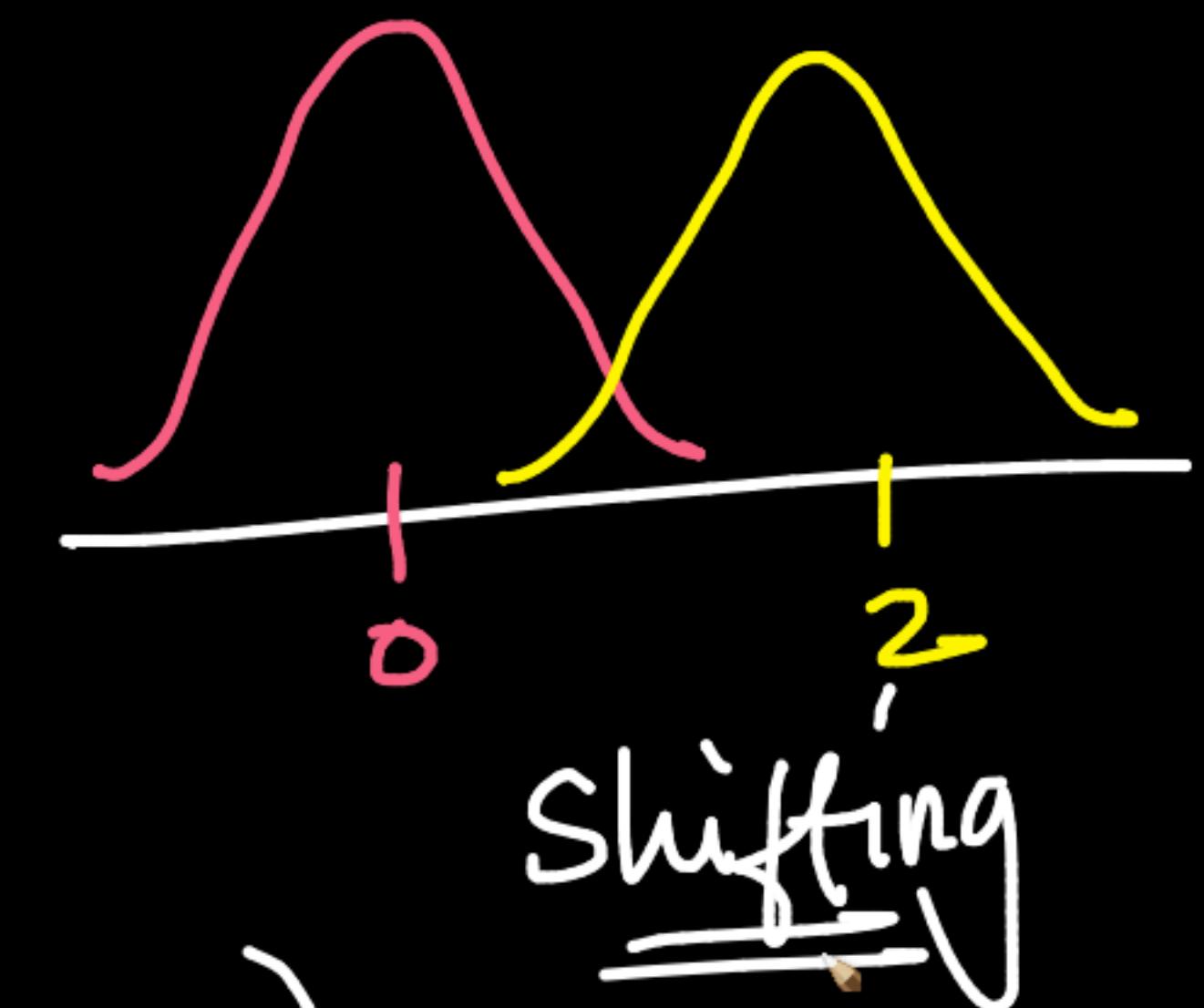
$z_t^i = y_i \sim \text{Gaussian}$



$x_i + \alpha$

$X \sim N(0, 1)$

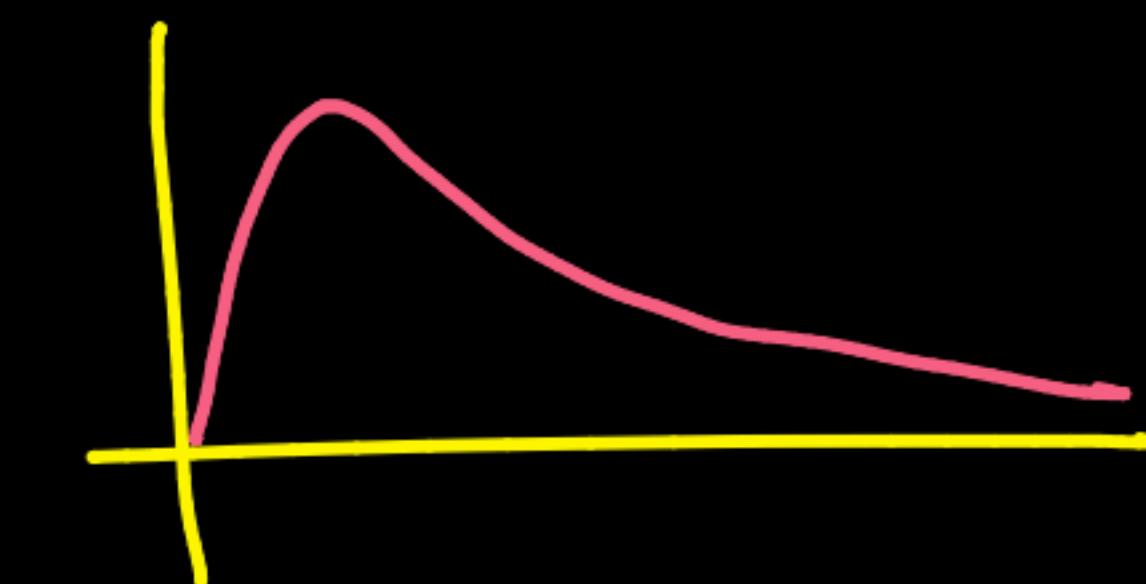
$X_i + \alpha \sim N(2, 1)$

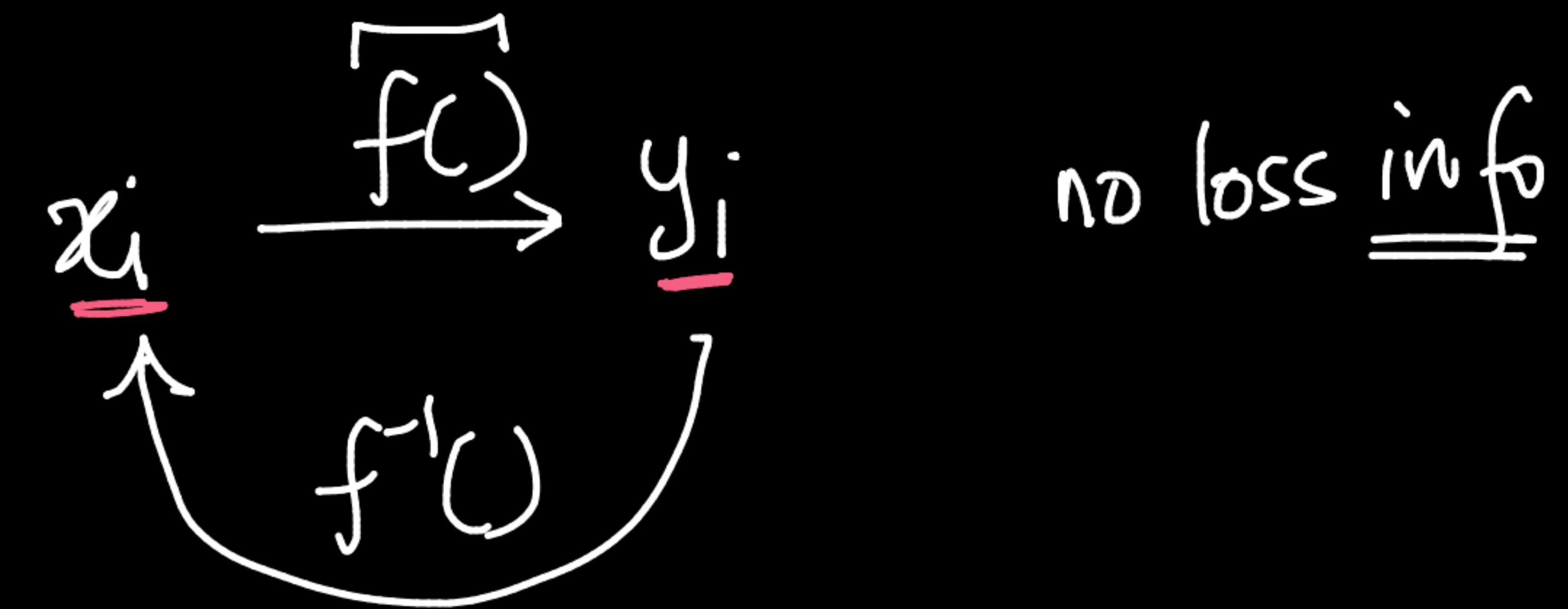


log-normal

↓
log

Gaussian





no loss in info

if f^{-1} does not exist \rightarrow loss of info $\rightarrow \dots$

box \sim Cox
non-gaussian

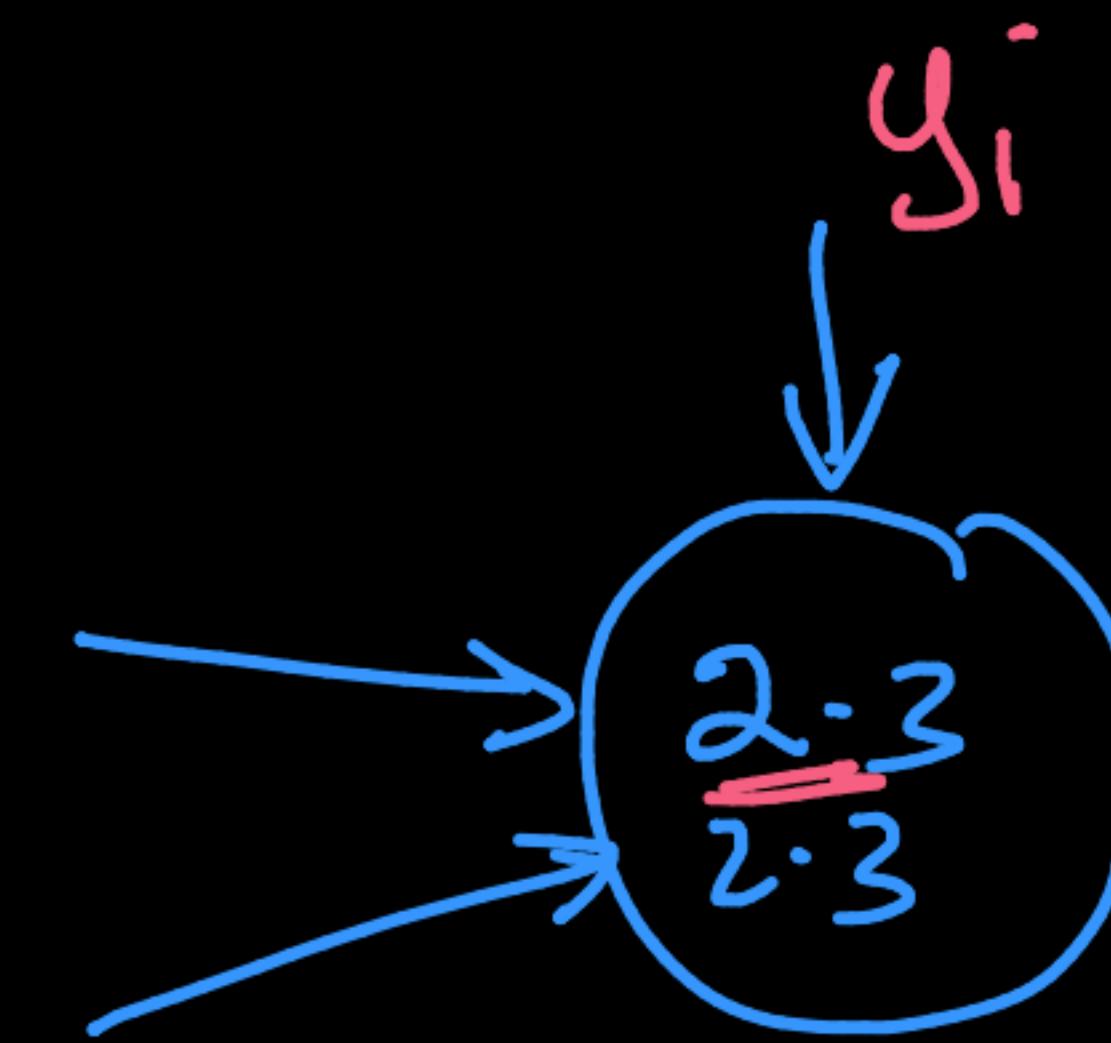
\hookrightarrow Gaussian

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

$|x_i|$: anisotropic

x^2 : non-invertible

$$\left\{ \begin{array}{l} x_1 = \underline{\underline{-2 \cdot 3}} \\ x_2 = \underline{\underline{2 \cdot 3}} \end{array} \right.$$



box-cox $(x_i + 3 \cdot 7)$

Composition

$$\sqrt{\log(x_i + 3 \cdot 7) + x_i^2}$$

Stats → lots of the Math proofs... ✓

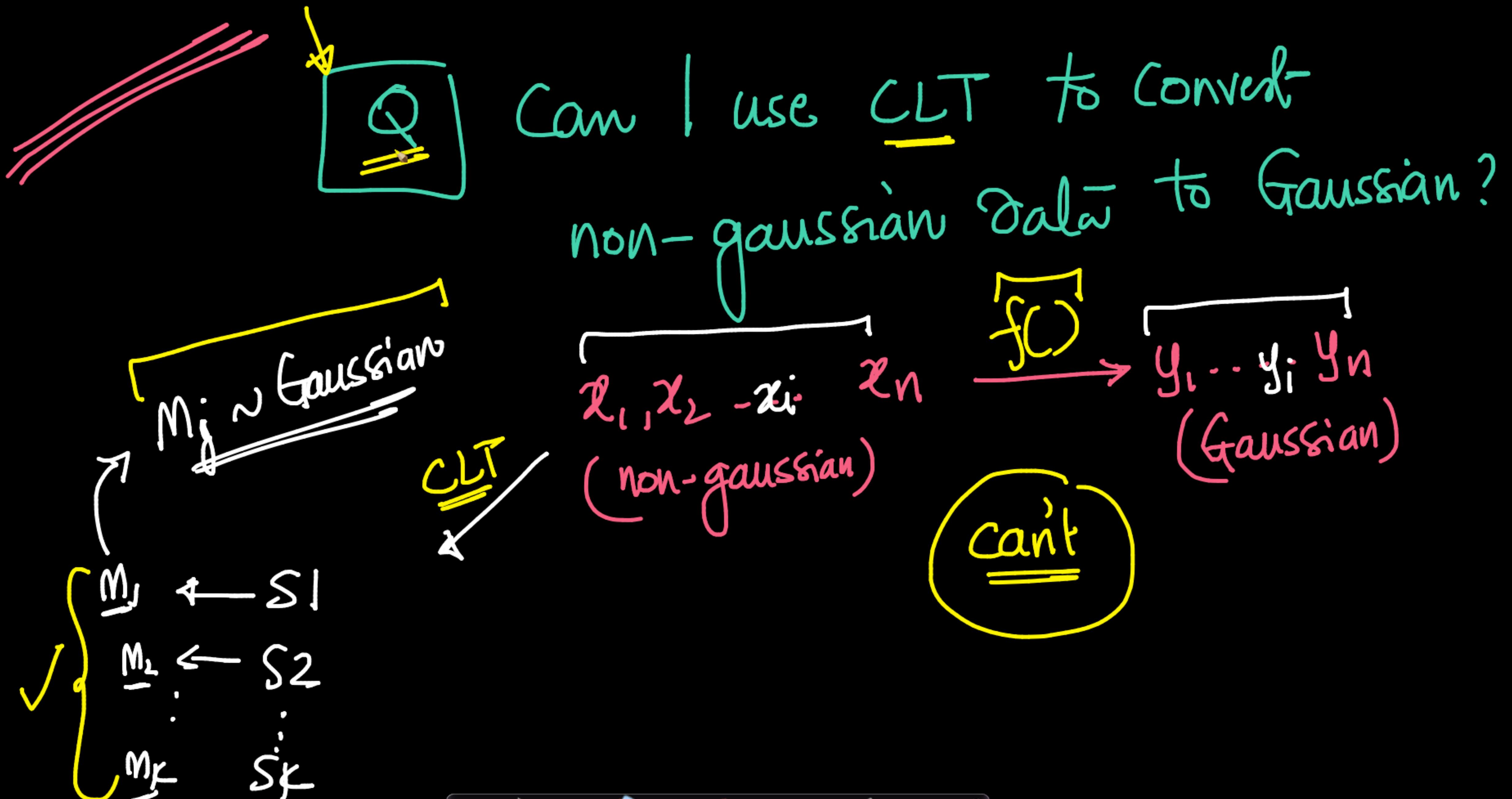
ANOVA

Z-test ~ population μ & σ are finite

Compare means

CLT ~ means are Gaussian

ML: advanced: — GBDT / RF } — don't care
DL



CLT is not a transformation

conditional
prob

$$\underline{P(A|B)} \checkmark$$

$$P(A|\overbrace{B,C}) = P(A|\boxed{B \cap C})$$

$$= \frac{P(A \cap B \cap C) / P(C)}{\cancel{P(B \cap C) / P(C)}}$$

if $P(C) \neq 0$

~~algebra~~

$$= \frac{P(A \cap B | C)}{P(B | C)} =$$

$$= \frac{P(A, B | C)}{P(B | C)}$$

~~Bayes!~~

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$\left. \begin{array}{l} P(B) \neq 0 \\ P(A) \neq 0 \end{array} \right\}$$

~~Extended Bayes!~~

$$\frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

Question

$$P(A | B, C) = \frac{P(A, B, C)}{P(B, C)}$$

$$A, B = A \cap B$$

$$= P(B | A, C) \cdot P(A | C)$$

$$P(B | C) P(C)$$

$$= P(B | A, C) \cdot P(A | C) P(C)$$

$$\cancel{P(B | C) P(C)}$$

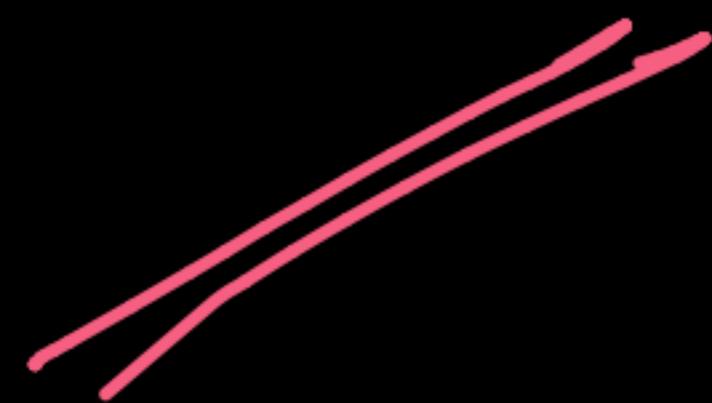
Conditional
Prob & \cap

$$A, B, C \vdash A \cap B \cap C$$

$$\begin{aligned} p(A, B, C) &= p(B, \textcircled{A, C}) \\ &= p(B) \textcircled{A \cap C} p(A \cap C) \end{aligned}$$

✓ ks-test
t-test
A/B testing
z-score

10:22



Z -test:

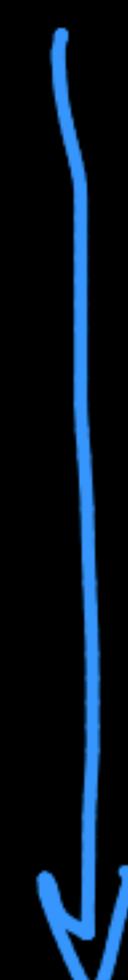
Med 1: x_{11}, \dots, x_{1,n_1} n_1

2-sample

Med 2: x_{21}, \dots, x_{2,n_2} n_2

$$T = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim Z(\sigma_T)$$

1-sample



batch
 $h: h_1, \dots, h_{\text{len}=\underline{n}}$

\bar{h} mean-height

$$\left\{ \begin{array}{l} \mu = 158 \text{ cm} \\ \sigma = 20 \text{ cm} \end{array} \right.$$

$$T = \frac{\bar{h} - \mu}{\sigma / \sqrt{n}} \sim Z(0,1)$$

Sample-std-dev
 $\sim \sigma$ if n
is large

$X: x_1, x_2, \dots, x_n$

$g_s \propto \text{Gaussian}$

Normality
testing: AD



Framework:

→ Data ①

KS-test

Data:

$X: x_1, x_2, \dots, x_n$ [2-sample]

$Y: y_1, y_2, \dots, y_m$

n_{max}
not too
small

Task: disb of X same as disb of Y

\bar{x}, s

$X: x_1, \dots, x_n$

$Y: y_1, \dots, y_m$

Normal($\mu = \bar{x}, \sigma = s$)

$f_{\mathbf{s}} \propto \text{Gaussian}$

State-rvs

②

$H_0: X \sim Y$ have same dist

$H_a: X \neq Y$

④ Sensible

$T_{KS} \rightarrow 0$

Under H_0 :

CDF are close to each other

Under H_a : T_{KS} : large val



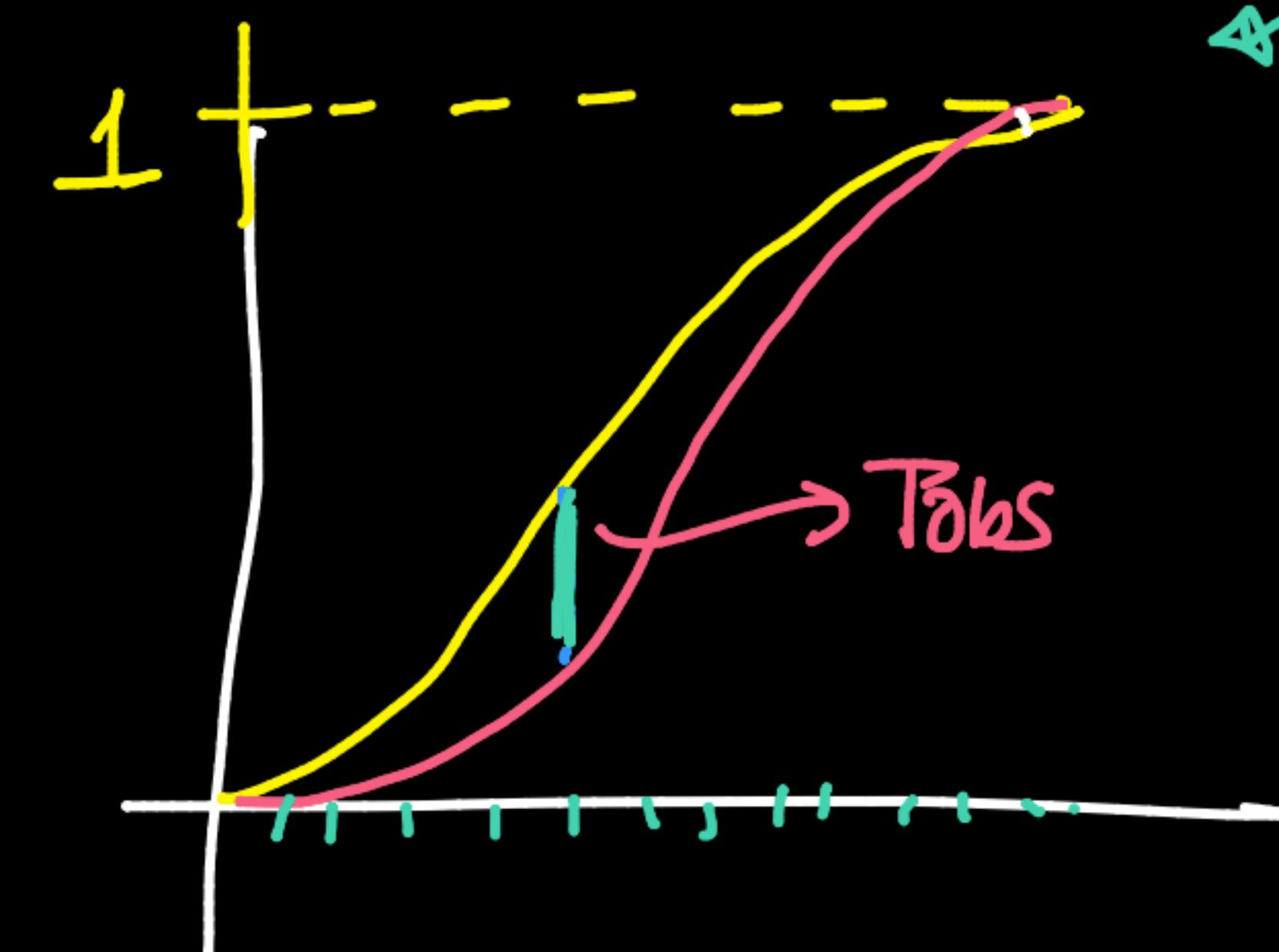
$T = \sup_{x \in \mathbb{R}} |F_n(x) - F_m(x)|$

Maximal

Supremum

gap b/w CDFs

\rightarrow QQ-plot
 \rightarrow compare CDFs

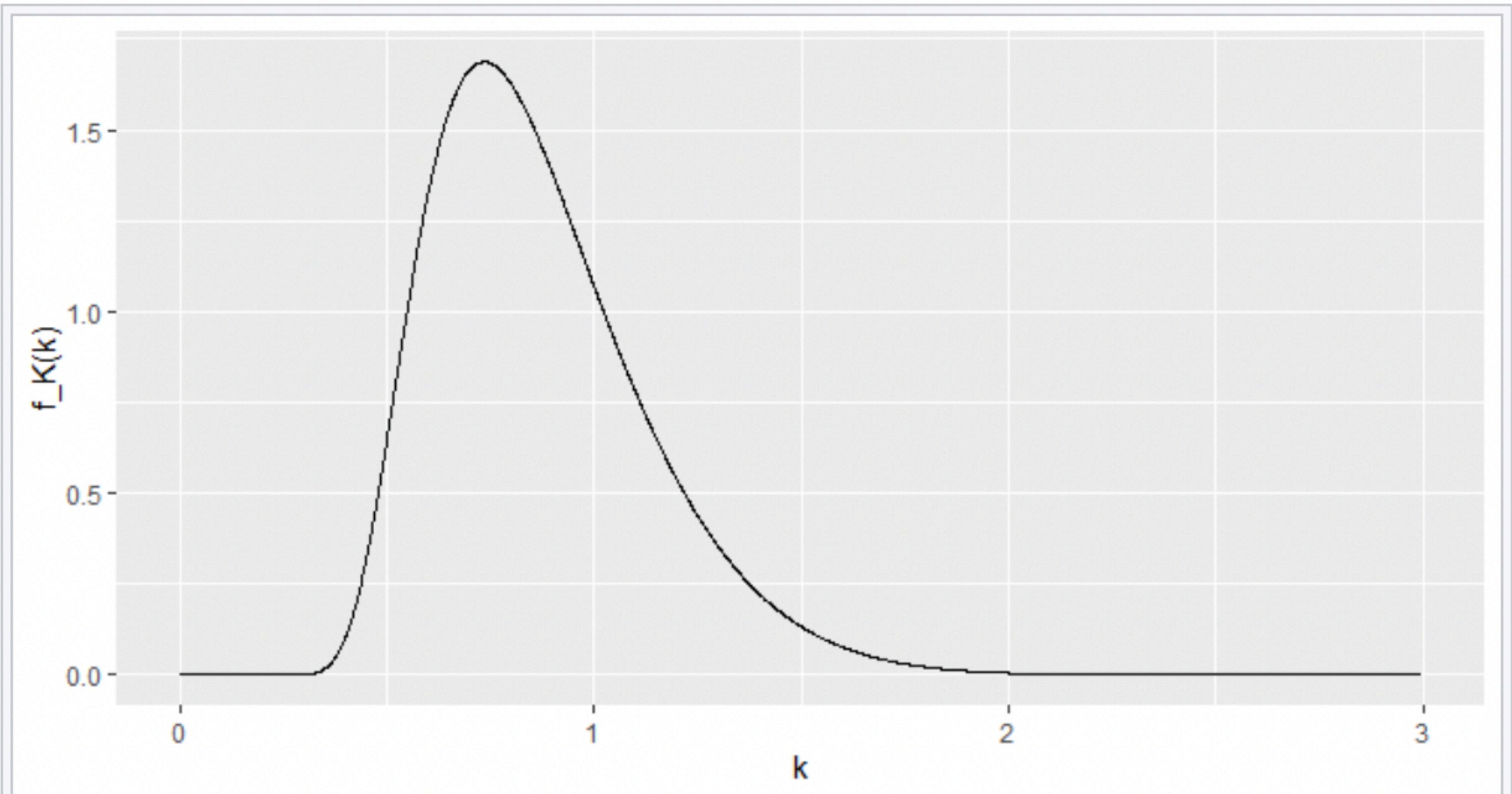


T_{KS} under H_0

$\hookrightarrow K_S$

Kolmogorov distribution [edit]

The Kolmogorov distribution is the distribution of the random variable



$$K = \sup_{t \in [0,1]} |B(t)|$$

where $B(t)$ is the Brownian bridge. The cumulative distribution function of K is given by^[3]

such as the [Anderson–Darling test statistic](#)) to properly reject the null hypothesis.

Physics

Kolmogorov distribution [edit]

The Kolmogorov distribution is the distribution of the random variable

KS-test

Brownian Motion

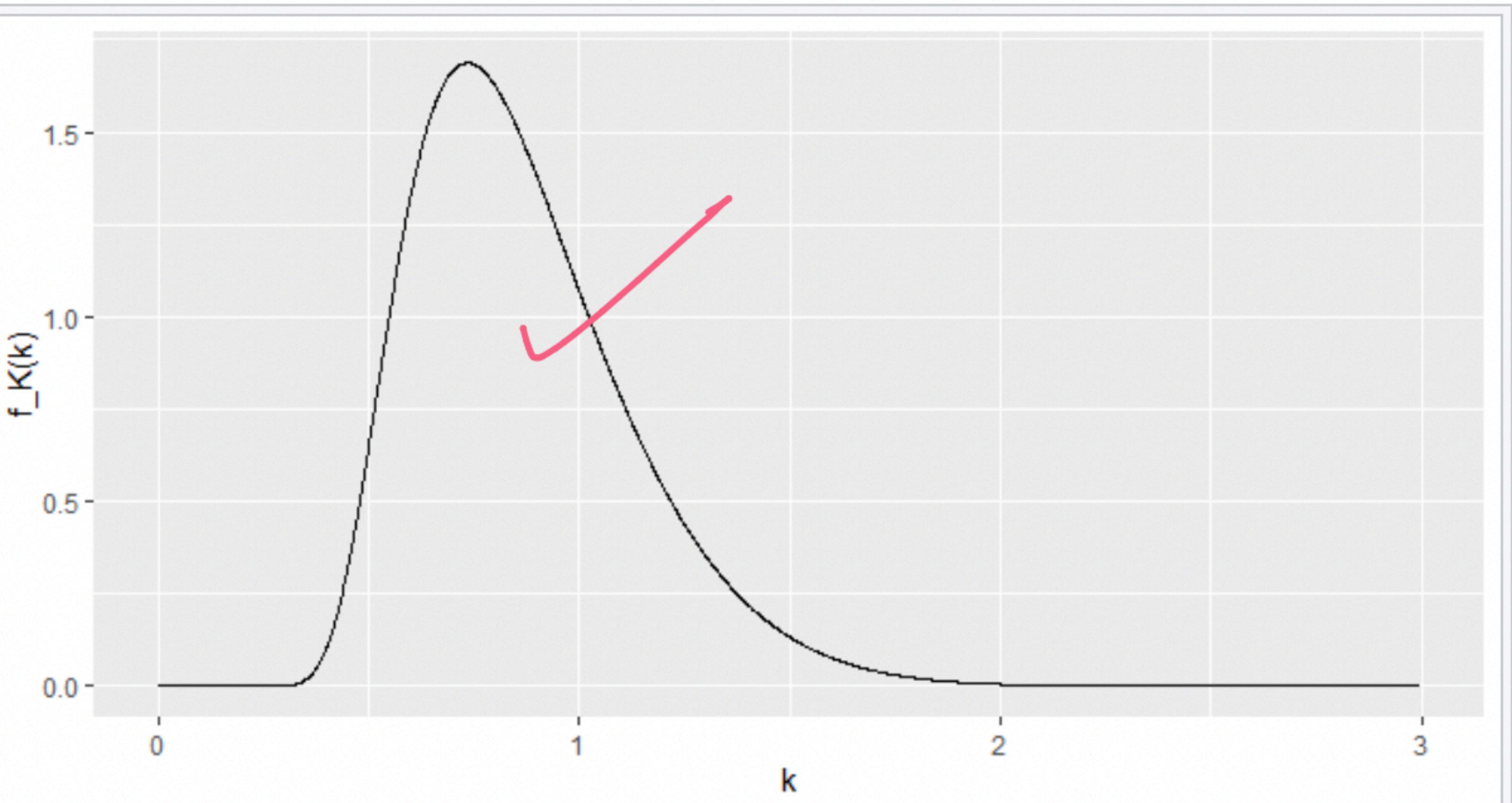
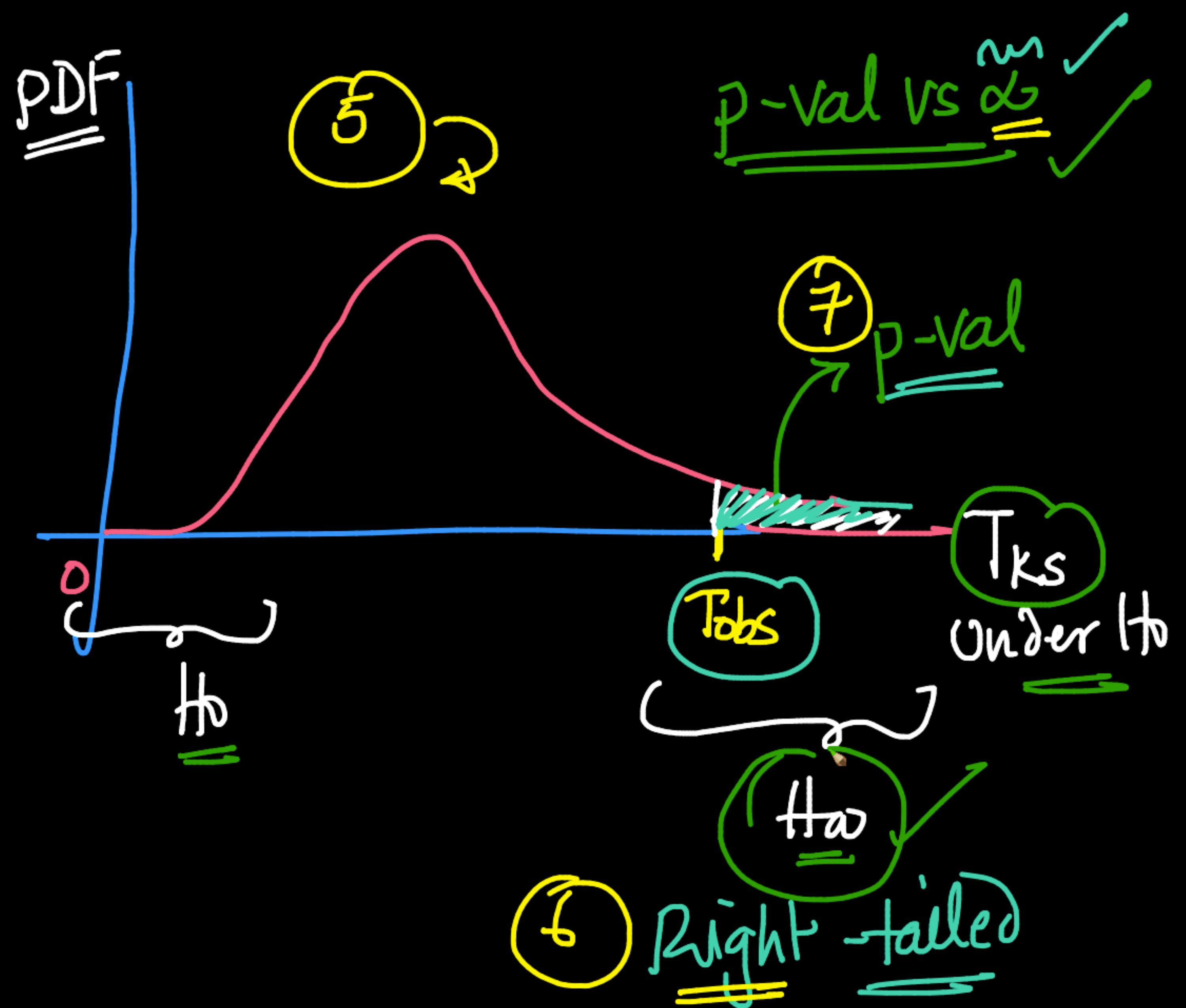


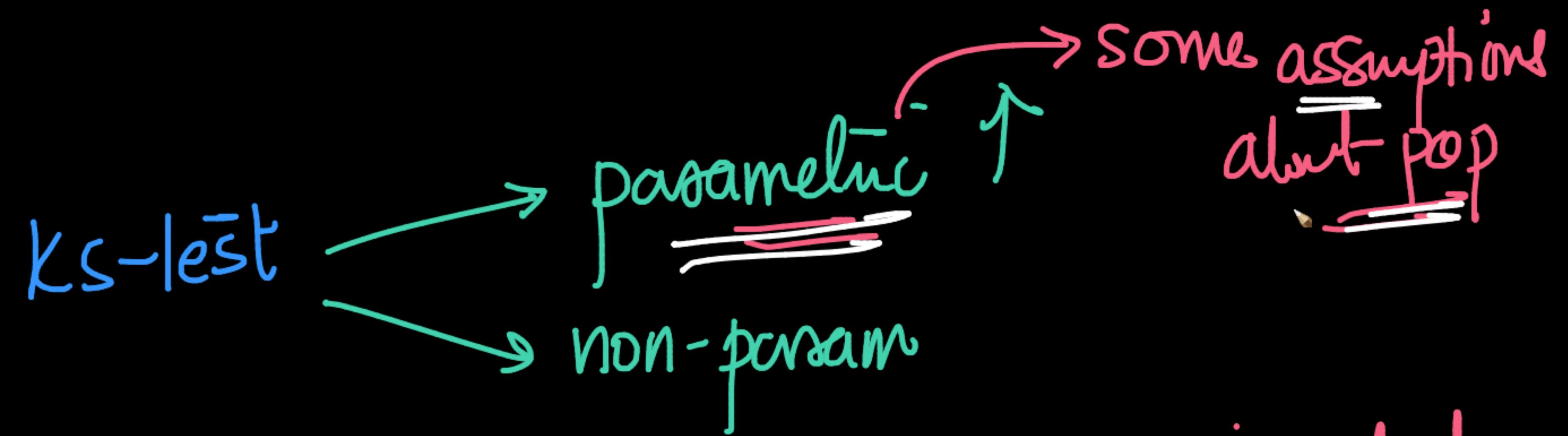
Illustration of the Kolmogorov distribution's PDF.

$$K = \sup_{t \in [0,1]} |B_t|$$



Jobs =



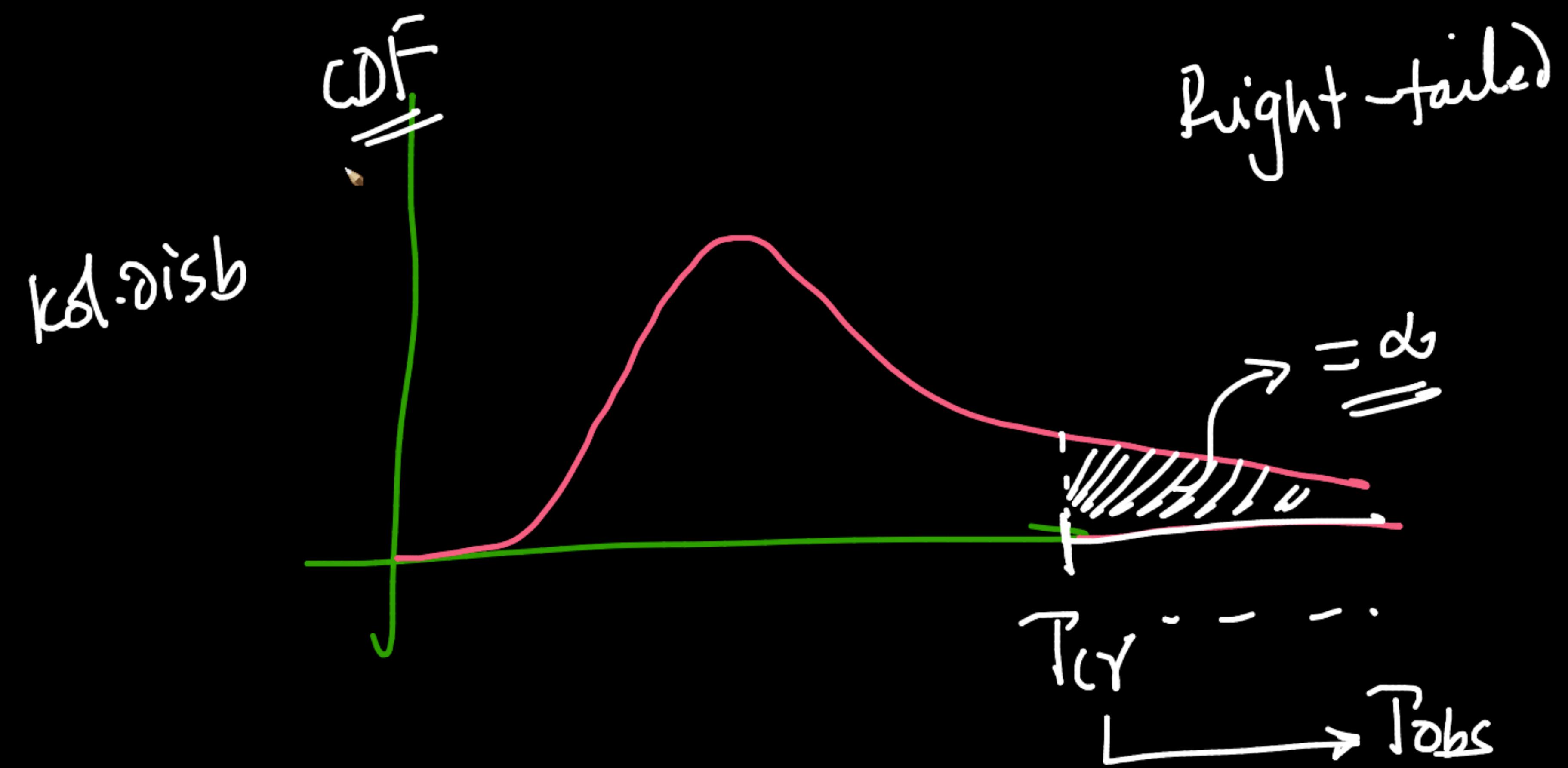


$T_{KS} \sim \text{kol-disb}$

{ any assumption about
the pop. disb of
X or Y }

Z-test :- pop finite $\mu \& \sigma \rightarrow$ param

$$T_z \sim z(0,1) \times$$

+ Code + Text

- ✓ RAM
- Disk

▼ KS-Test

z-test:

Med 1 : } Compose - 2'ish
Med 2 : }

```
[1] from scipy import stats  
import numpy as np  
import matplotlib.pyplot as pl
```

```
# recovery times of patients who took medicine-1  
r1 = [8.82420842, 7.47774471, 7.55712098, 7.98131439, 6.82771606,  
      7.48566433, 9.15385732, 5.84040502, 8.26124313, 8.4728876 ,  
      6.82582186, 7.00490974, 8.43423058, 6.72099932, 6.97495982,  
      5.93748053, 5.40707847, 6.16385557, 6.71421056, 4.42396183,  
      6.87285228, 8.00313581, 6.69035041, 7.83622942, 8.70984957,  
      5.56284584, 9.08093437, 4.98165193, 7.67769408, 6.04738478,  
      7.64921582, 7.31051639, 6.74463303, 7.27356973, 8.16787232,  
      6.90990965, 7.06439167, 6.62921957, 6.08283539, 6.2458137 ,  
      8.65173634, 5.76080646, 6.20573219, 8.91561004, 6.22560201,  
      5.67542104, 6.97412435, 8.31354697, 8.14172701, 8.26099345,  
      7.87612791, 6.24835109, 9.95324783, 6.59504627, 6.17365145,  
      6.05676895, 7.23030223, 7.71311809, 7.37163804, 5.69798738,
```

boxcox and z-score.ipynb - Colab Notebook | [scipy.stats.boxcox — SciPy v1.8.0 documentation](#) | [Anderson-Darling test - Wikipedia](#) | KStest_Ttest.ipynb - Colaboratory | [Normality test - Wikipedia](#)

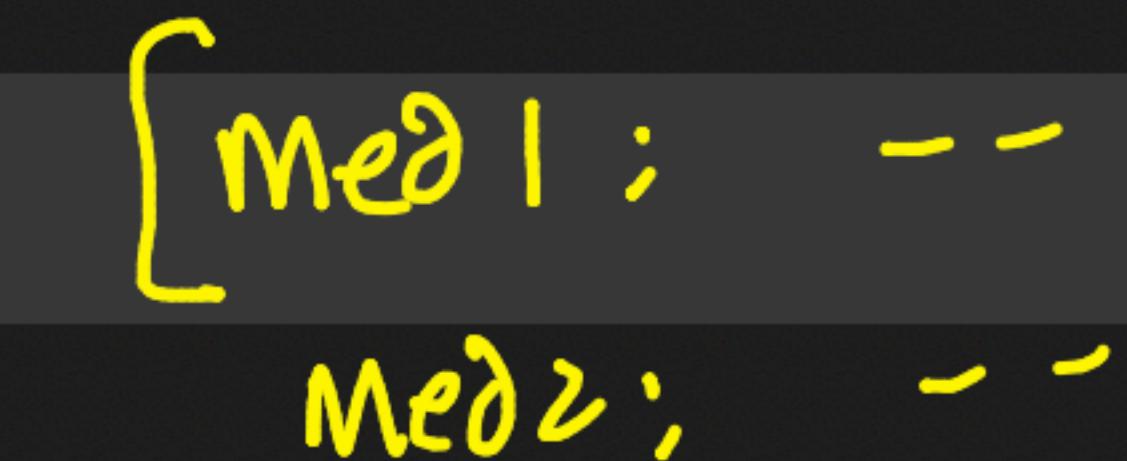
colab.research.google.com/drive/1IWuS8AxULBBold2GIfBKUeT1RSj3fFAJ#scrollTo=wVXdQQi1WSWT

+ Code + Text RAM Disk  Update

[6] n1
0s
100

[7] n2
0s
120

[8] #2-sample KS Test
stats.ks_2samp(d1, d2)



Ks_2sampResult(statistic=0.3233333333333333, pvalue=1.516338798324135e-05)

p-value = 0.00001516 < 0.01 = alpha

=> Reject H0 (r1 and r2 same distribution)

=> Accept Ha (the distributions of r1 and r2 are different)

[9] plt.grid()
a = plt.hist(d1, bins=100, cumulative=True, label='CDF', density=True, histtype='step')

boxcox and z-score.ipynb - Colab Notebook | scipy.stats.boxcox — SciPy v1.8.0 | Anderson-Darling test - Wikipedia | KStest_Ttest.ipynb - Colaboratory | Normality test - Wikipedia

colab.research.google.com/drive/1IWuS8AxULBBold2GIfBKUeT1RSj3fFAJ#scrollTo=wVXdQQi1WSWT

+ Code + Text

RAM Disk

[6] n1
0s
100

[7] n2
0s
120

#2-sample KS Test
stats.ks_2samp(d1, d2)

Ks_2sampResult(statistic=0.3233333333333333, pvalue=1.516338798324135e-05)

p-value = 0.00001516 < 0.01 = alpha
=> Reject H0 (r1 and r2 same distribution) ✓
=> Accept Ha (the distributions of r1 and r2 are different)

plt.grid()
a = plt.hist(d1, bins=100, cumulative=True, label='CDF', density=True, histtype='step')

0.01
 $\alpha = 0.01$

$\alpha = 0.1$,
0.00

T_{KS}
1.5x10⁻⁵

p-val < 0.00]

RAM Disk

57 / 57

boxcox and z-score.ipynb - Colab Notebook | [scipy.stats.boxcox — SciPy v1.8.0 documentation](#) | [Anderson–Darling test - Wikipedia](#) | KStest_Ttest.ipynb - Colaboratory | [Normality test - Wikipedia](#)

colab.research.google.com/drive/1IWuS8AxULBBold2GIfBKUeT1RSj3fFAJ#scrollTo=nYpehMuXNvD0

+ Code + Text

RAM Disk

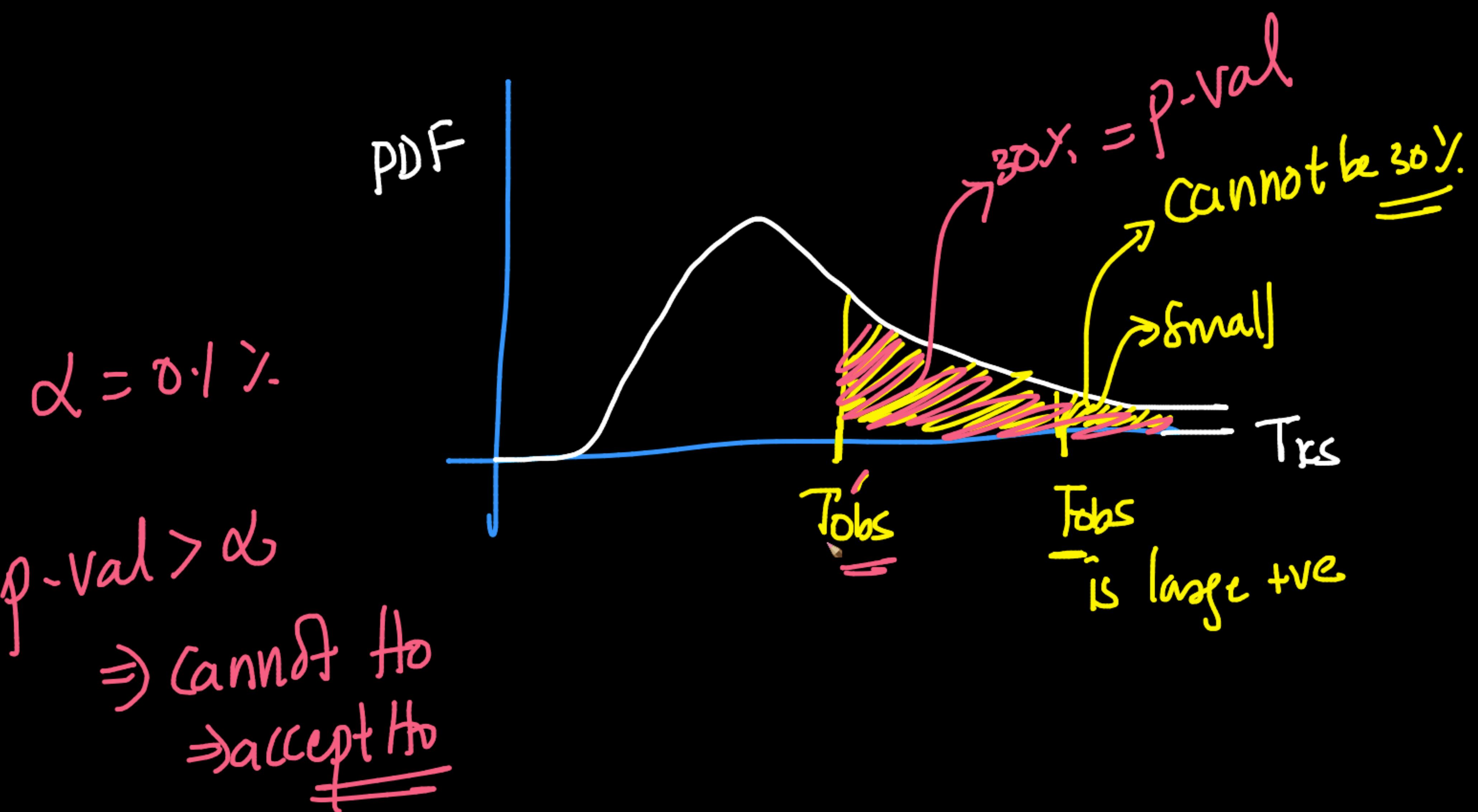
Q {x} Double-click (or enter) to edit

0s

plt.grid()
a = plt.hist(d1, bins=100, cumulative=True, label='CDF', density=True, histtype='step')
b = plt.hist(d2, bins=100, cumulative=True, label='CDF', density=True, histtype='step')
plt.show()

↑ ↓ ↻ ⚙️ 📈 🗑️ ⋮

Double-click (or enter) to edit



boxcox and z-score.ipynb - Co | scipy.stats.boxcox — SciPy v1.8 | Anderson-Darling test - Wikipedia | KTest_Ttest.ipynb - Colaboratory | Kolmogorov-Smirnov test - Wikipedia +

en.wikipedia.org/wiki/Kolmogorov-Smirnov_test#Kolmogorov_distribution

Kolmogorov distribution [edit]

The Kolmogorov distribution is the distribution of the random variable

Normal(0, 1)

(let)

$\checkmark x_1 \dots x_{10}$

$\checkmark y_1 \dots y_{10}$

Samples

Illustration of the Kolmogorov distribution's PDF.

$$K = \sup_{t \in [0,1]} |B(t)|$$

where $B(t)$ is the Brownian bridge. The cumulative distribution function of K is given by^[3]

Kolmogorov distribution [edit]

The Kolmogorov distribution is the distribution of the random variable

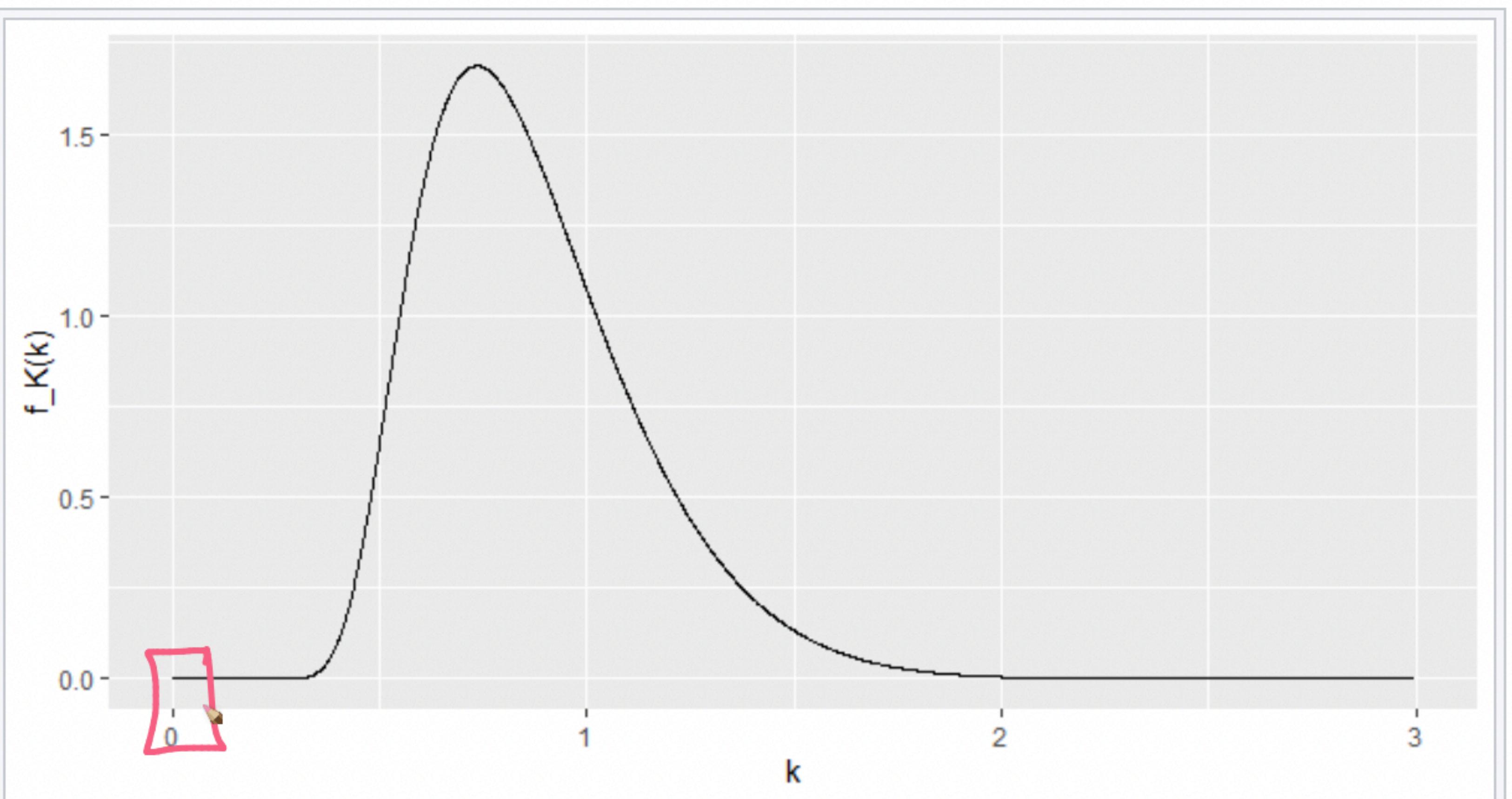


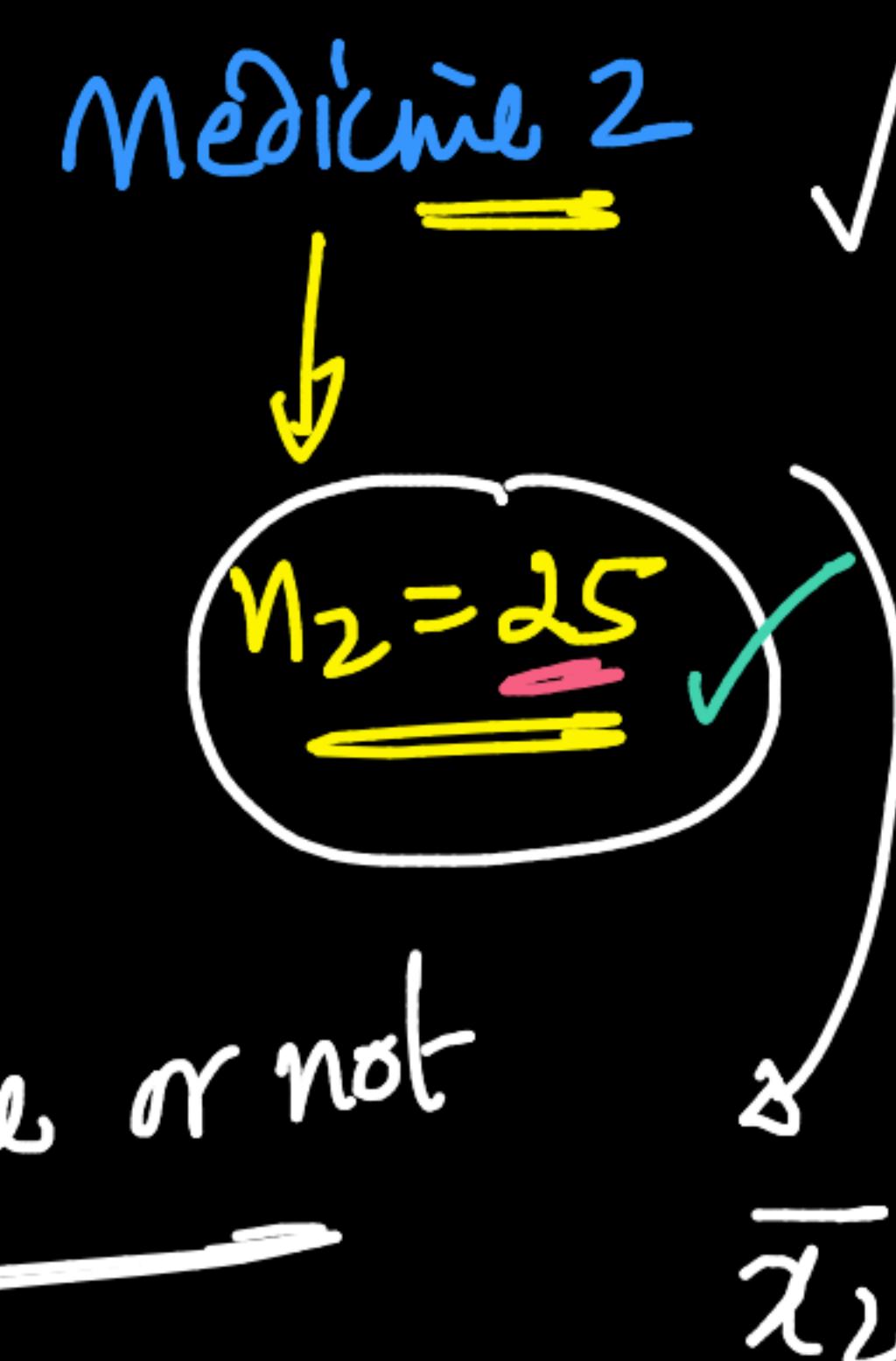
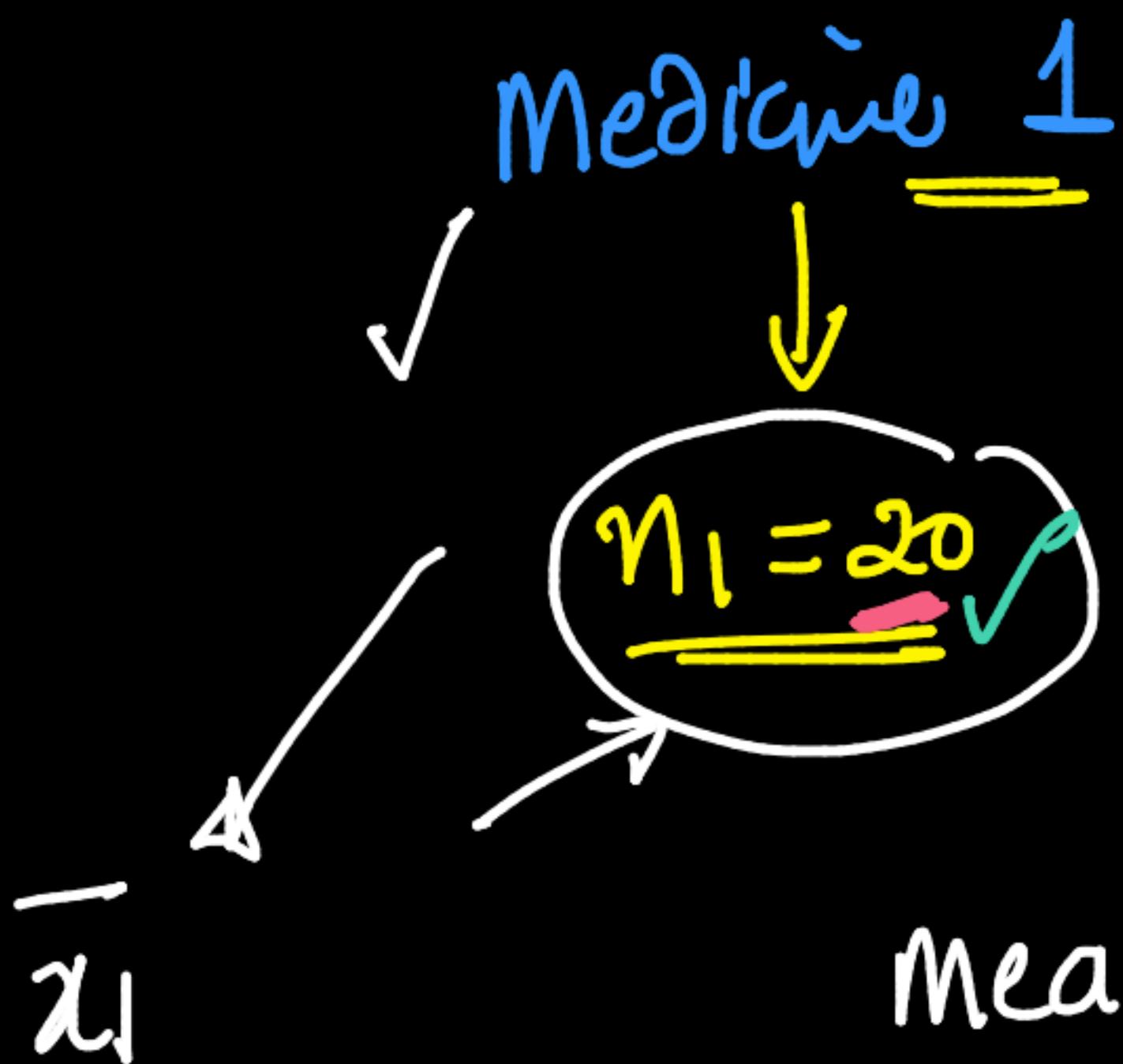
Illustration of the Kolmogorov distribution's PDF.

$$K = \sup_{t \in [0,1]} |B(t)|$$

where $B(t)$ is the Brownian bridge. The cumulative distribution function of K is given by^[3]



2-Sample



T-test ✓

sample sizes
are small

estimate σ_1 & σ_2

(1)
can't use

Z-test ✓

Framework:

$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 \neq \mu_2$$

real-world: n_1 & n_2 are small
↳ extreme medical expls. } \rightarrow surveys }
→ high-cost ← physics

Assumptions:

T-test

- ① sample-means
 \bar{x}_1 & \bar{x}_2 → Gaussian dist (post)
CLT → population μ & σ are finite
- ② obs are random and indep - .
- ③ σ_1 & σ_2 are not known/estimatable because n_1 & n_2 are small
 $n_1, n_2 \leq 30$

$T \rightarrow$ dis
 $T(t)$

$$T_t =$$

$$\bar{x}_1 - \bar{x}_2$$

$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

Sensible

s_1^2 & s_2^2 are sample
 $s_1 \rightarrow \text{dev}$

$$T_z = \bar{z} - \bar{x}_2$$

$\sqrt{\frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}}$

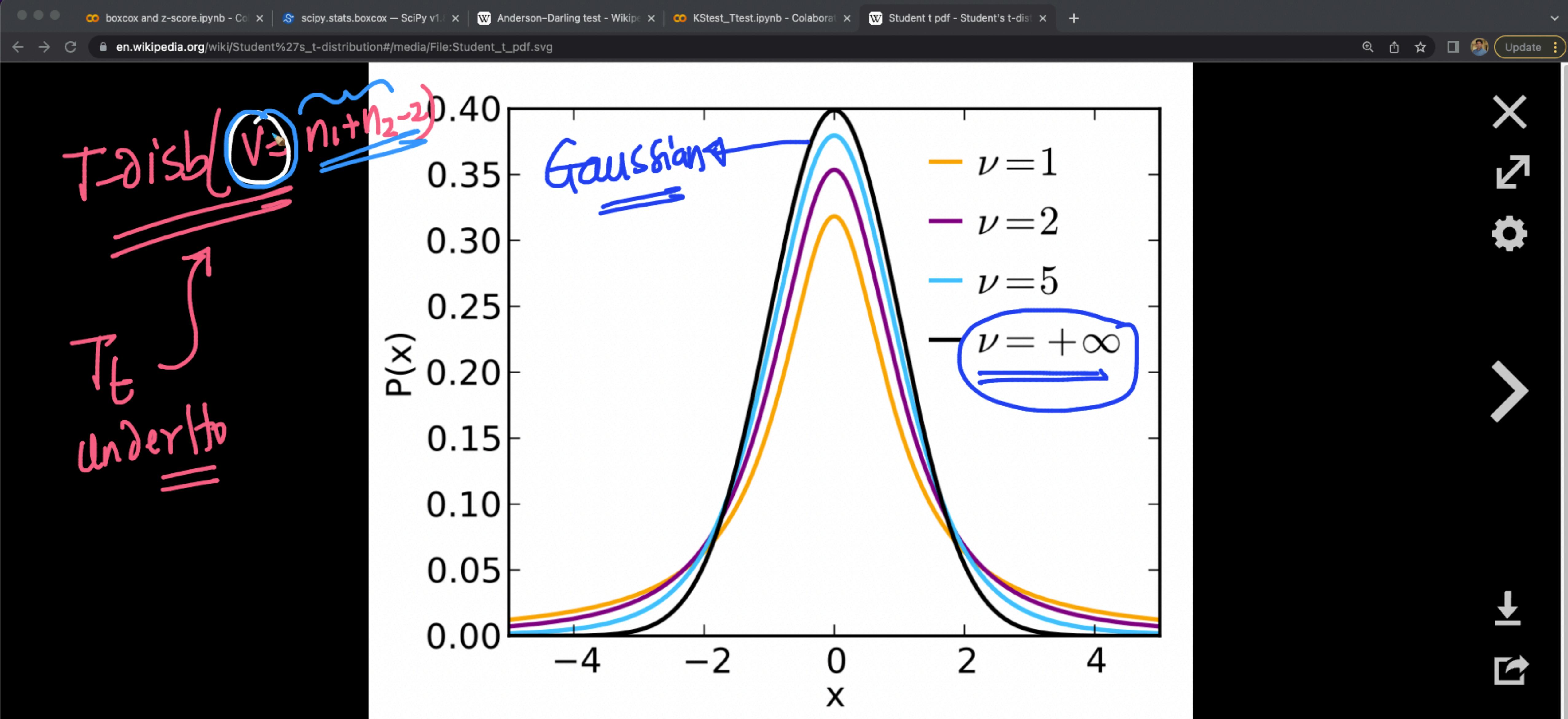
pop. std - dev \rightarrow given ✓
 \rightarrow estimate
hem wll

✓

$$T_t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$\sim T \text{ distb}(v) \xrightarrow{\text{DOF}}$

$n_1 + n_2 - 2$
(pasam)



Plot of the density function for several members of the Student t family.

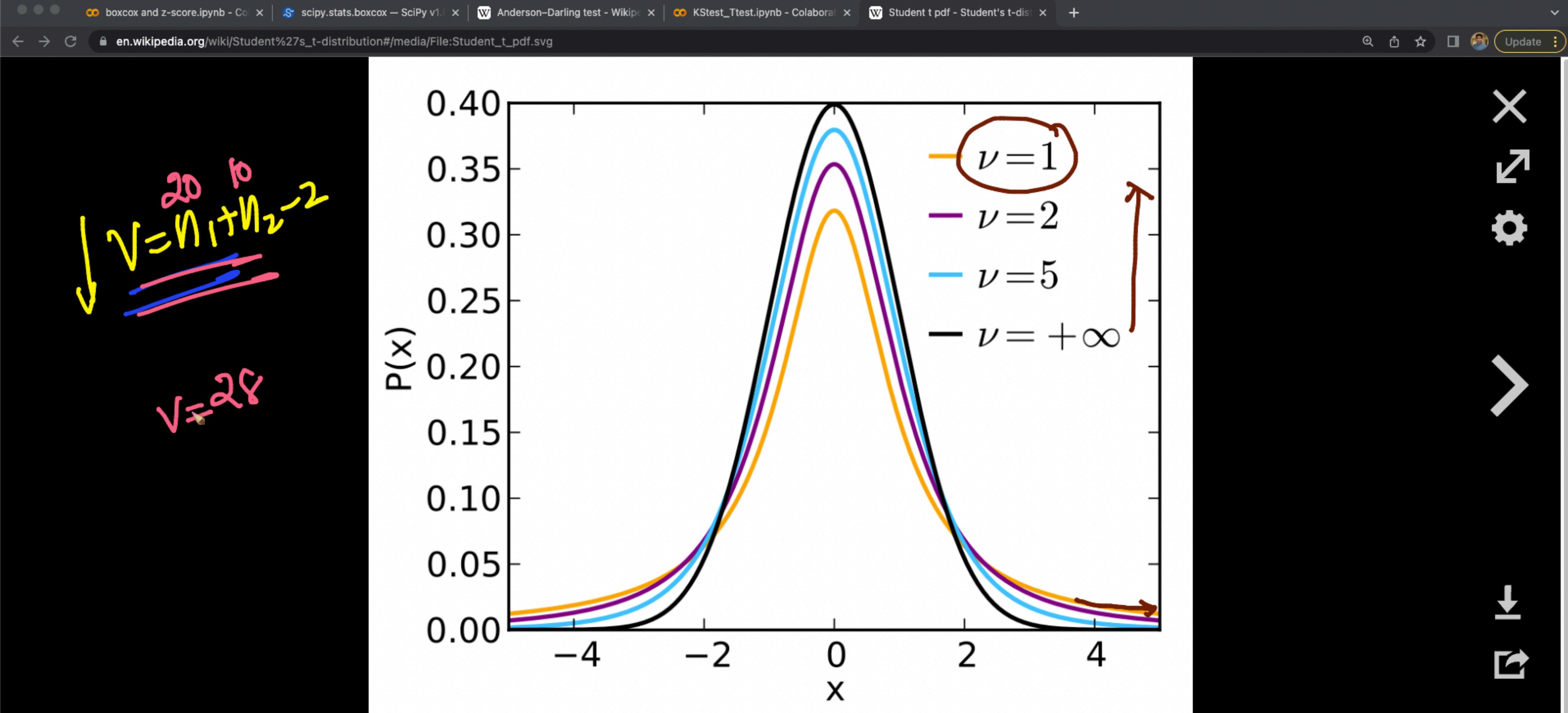
More details

✓ $\sqrt{V = (n_1 + n_2 - 2)}$ ↑ ⇒ larger-sample sizes ↓

$T_t \rightarrow$ Gaussian

$$\begin{aligned} S_1^2 &\approx \sigma_1^2 \\ S_2^2 &\approx \sigma_2^2 \end{aligned}$$

↓ $T_t \approx T_z \rightarrow$ Gaussian



More details

In other projects

Wikimedia Commons

Languages 

العربية

Deutsch

Español

Français

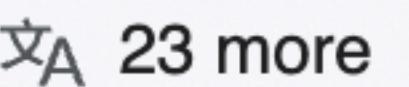
Jawa

日本語

Português

Русский

中文

 23 more

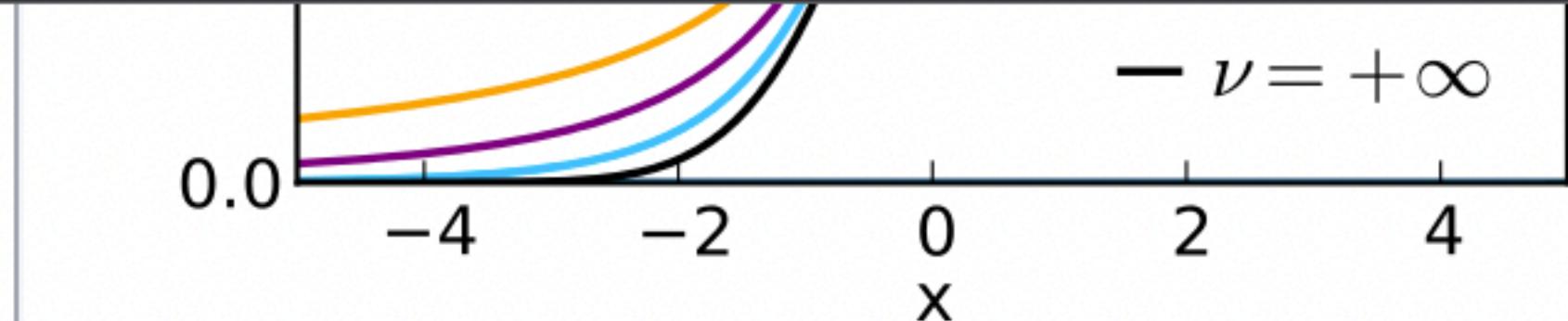
 Edit links

en.wikipedia.org/wiki/Student%27s_t-distribution

distribution. However, the t -distribution has heavier tails, meaning that it is more prone to producing values that fall far from its mean. This makes it useful for understanding the statistical behavior of certain types of ratios of random quantities, in which variation in the denominator is amplified and may produce outlying values when the denominator of the ratio falls close to zero. The Student's t -distribution is a special case of the [generalized hyperbolic distribution](#).

Contents [hide]

- 1 History and etymology
- 2 How Student's distribution arises from sampling
- 3 Definition
 - 3.1 Probability density function
 - 3.2 Cumulative distribution function
 - 3.3 Special cases
- 4 How the t -distribution arises
 - 4.1 Sampling distribution
 - 4.2 Bayesian inference
- 5 Characterization
 - 5.1 As the distribution of a test statistic
 - 5.1.1 Definition
 - 5.2 As a maximum likelihood estimator



$\nu = +\infty$

Parameters	$\nu > 0$ degrees of freedom (real)
Support	$x \in (-\infty, \infty)$
PDF	$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$
CDF	$\frac{1}{2} + x\Gamma\left(\frac{\nu+1}{2}\right) \times \frac{{}_2F_1\left(\frac{1}{2}, \frac{\nu+1}{2}; \frac{3}{2}; -\frac{x^2}{\nu}\right)}{\sqrt{\pi\nu}\Gamma\left(\frac{\nu}{2}\right)}$
where ${}_2F_1$ is the hypergeometric function	
Mean	0 for $\nu > 1$, otherwise undefined
Median	0
Mode	0
Variance	$\frac{\nu}{\nu-2}$ for $\nu > 2$, ∞ for $1 < \nu \leq 2$, otherwise undefined
Skewness	0 for $\nu > 3$, otherwise undefined
Ex. kurtosis	$\frac{6}{\nu-4}$ for $\nu > 4$, ∞ for $2 < \nu \leq 4$, otherwise undefined
Entropy	$\frac{\nu+1}{2} \left[{}_2F_1\left(\frac{1+\nu}{2}, \frac{1+\nu}{2}; \frac{3}{2}; -\frac{1}{\nu}\right) - {}_2F_1\left(\frac{\nu}{2}, \frac{\nu}{2}; \frac{3}{2}; -\frac{1}{\nu}\right) \right]$

Under H_0 : $T_t \sim T\text{-dist} \left(V = n_1 + n_2 - 2 \right)$

Sensible as this is similar to T_2

2-sided

Observed Data

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

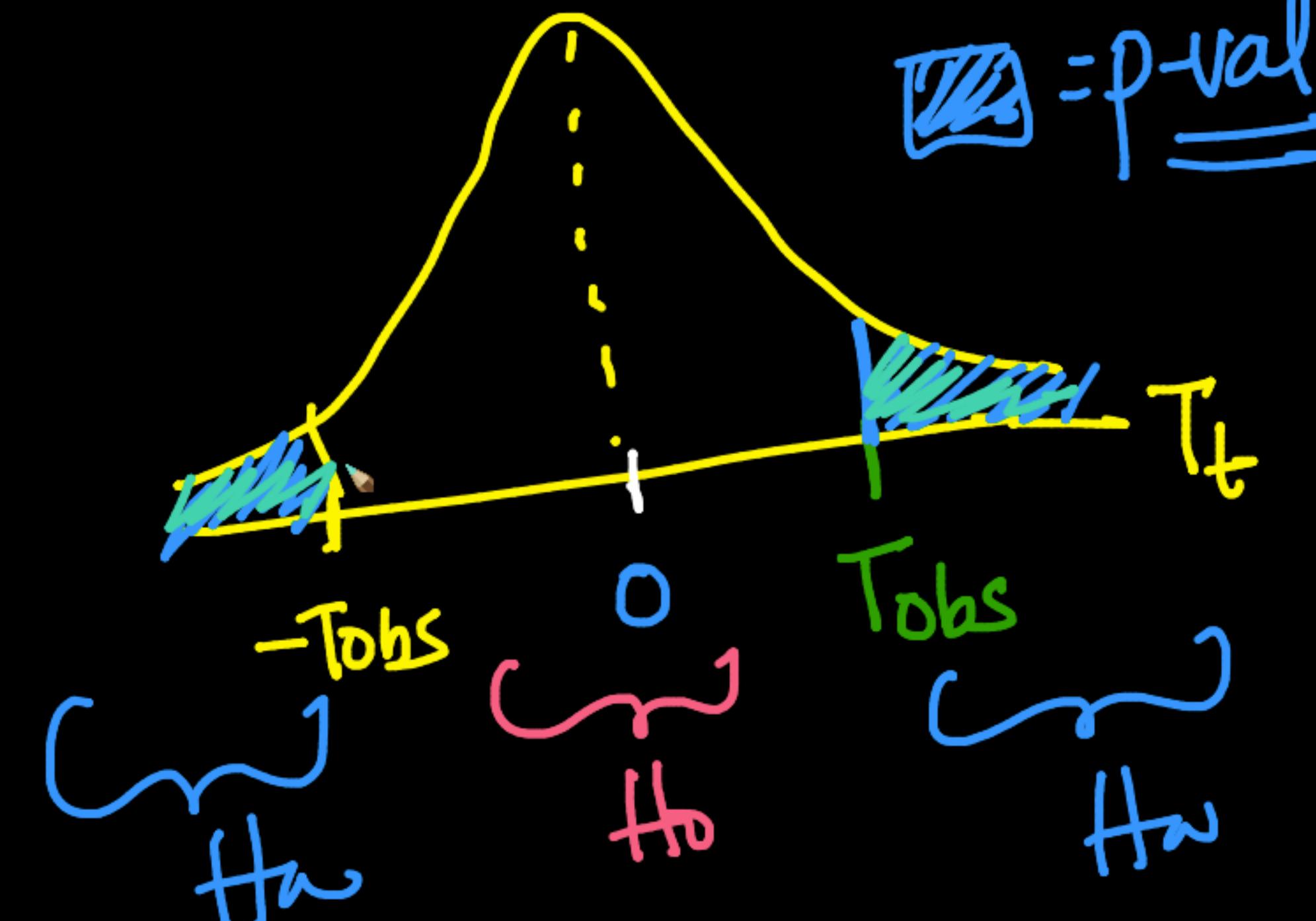
Under H_0 :

$$T_t \rightarrow 0$$

Under H_a :

$T_t \rightarrow$ large \pm -ve values

~~Symmetric~~
~~T-Dist~~ ($V = n_1 + n_2 - 2$)



2-sided

 p-val vs α

Code → || Youtube -dalā



Z-test
T-test

KS-test

{ 1-ad → ----
2-ad → ~~~

- Z-proportions
- Chi-test
- ANOVA

boxcox and z-score.ipynb - Colaboratory | scipy.stats.boxcox — SciPy v1.8.0.dev0+11.gf3d333e | Anderson-Darling test - Wikipedia | KStest_Ttest.ipynb - Colaboratory | Normal distribution - Wikipedia +

en.wikipedia.org/wiki/Normal_distribution

Permanent link
Page information
Cite this page
Wikidata item
Print/export
Download as PDF
Printable version

In other projects
Wikimedia Commons

Languages 

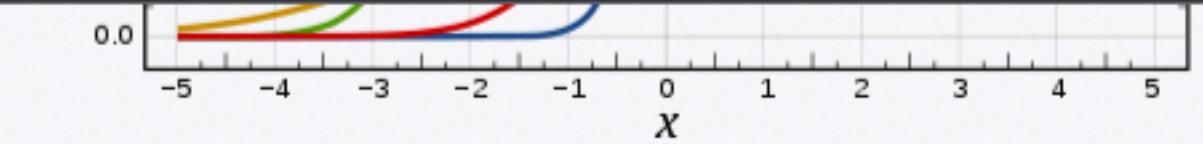
العربية
Español
Français
हिन्दी
मराठी
Português
தமிழ்
اردو
中文
文 A 60 more

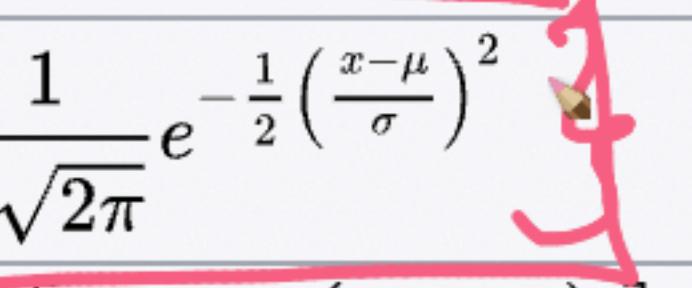
is partly due to the [central limit theorem](#). It states that, under some conditions, the average of many samples (observations) of a random variable with finite mean and variance is itself a random variable—whose distribution [converges](#) to a normal distribution as the number of samples increases. Therefore, physical quantities that are expected to be the sum of many independent processes, such as [measurement errors](#), often have distributions that are nearly normal.^[4]

Moreover, Gaussian distributions have some unique properties that are valuable in analytic studies. For instance, any linear combination of a fixed collection of normal deviates is a normal deviate. Many results and methods, such as [propagation of uncertainty](#) and [least squares](#) parameter fitting, can be derived analytically in explicit form when the relevant variables are normally distributed.

A normal distribution is sometimes informally called a **bell curve**.^[5] However, many other distributions are bell-shaped (such as the [Cauchy](#), [Student's t](#), and [logistic distributions](#)).

The univariate probability distribution is generalized for vectors in the [multivariate normal distribution](#) and for matrices in the [matrix normal distribution](#).

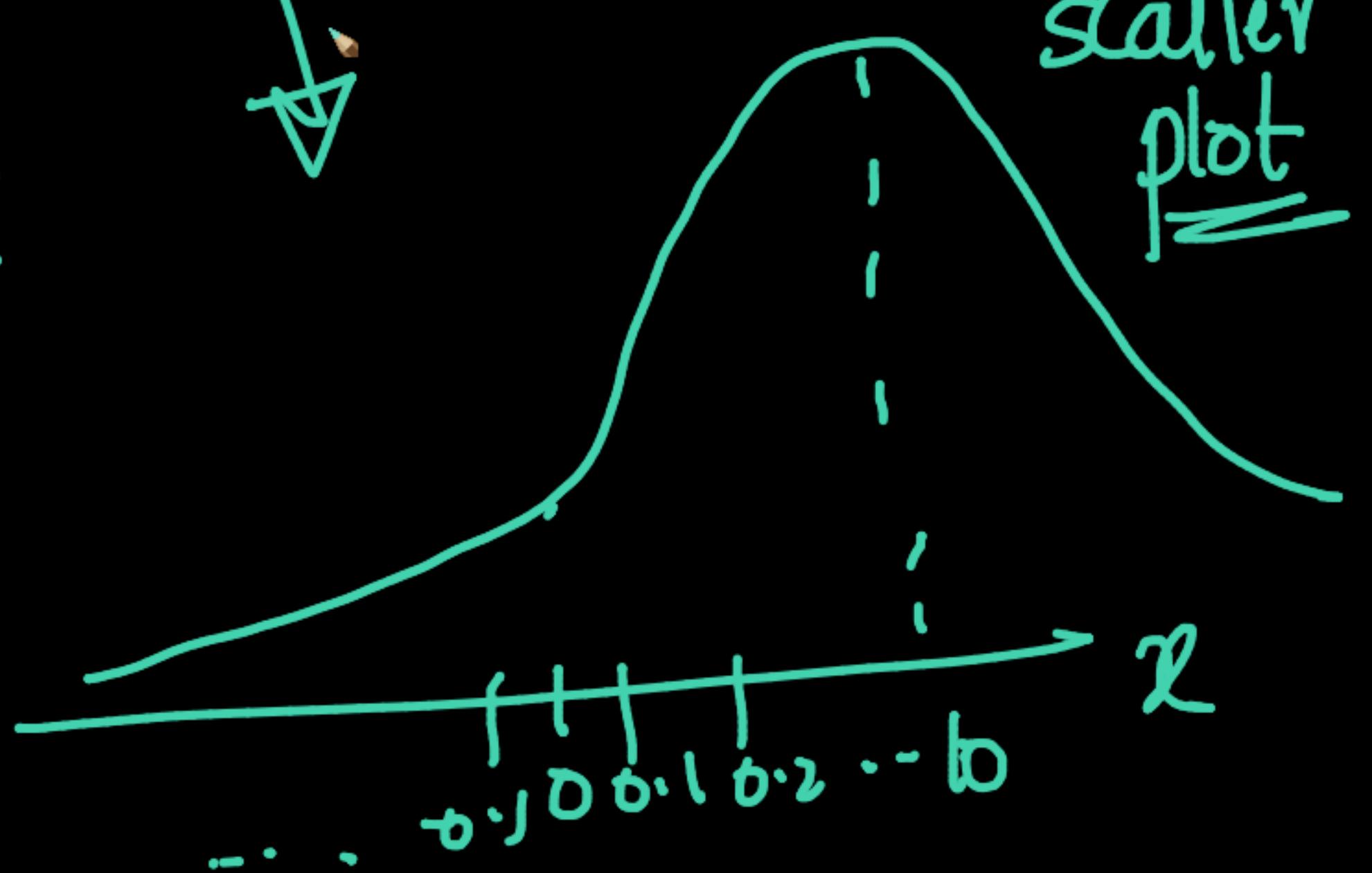
A graph showing the probability density function of a standard normal distribution. The x-axis is labeled 'x' and ranges from -5 to 5. The y-axis represents probability density. Three curves are shown: a red curve peaking at x=0, a green curve shifted to the right, and a blue curve shifted further to the right. All three curves are symmetric and approach zero as |x| increases.

Notation	$\mathcal{N}(\mu, \sigma^2)$
Parameters	$\mu \in \mathbb{R}$ = mean (location) $\sigma^2 \in \mathbb{R}_{>0}$ = variance (squared scale)
Support	$x \in \mathbb{R}$
PDF	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$ 
CDF	$\frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right]$
Quantile	$\mu + \sigma\sqrt{2} \operatorname{erf}^{-1}(2p - 1)$
Mean	μ
Median	μ
Mode	μ
Variance	σ^2
MAD	$\sigma\sqrt{2/\pi}$
Skewness	0
Ex. kurtosis	0
Entropy	$\frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2}$
MGF	$\exp(\mu t + \sigma^2 t^2/2)$

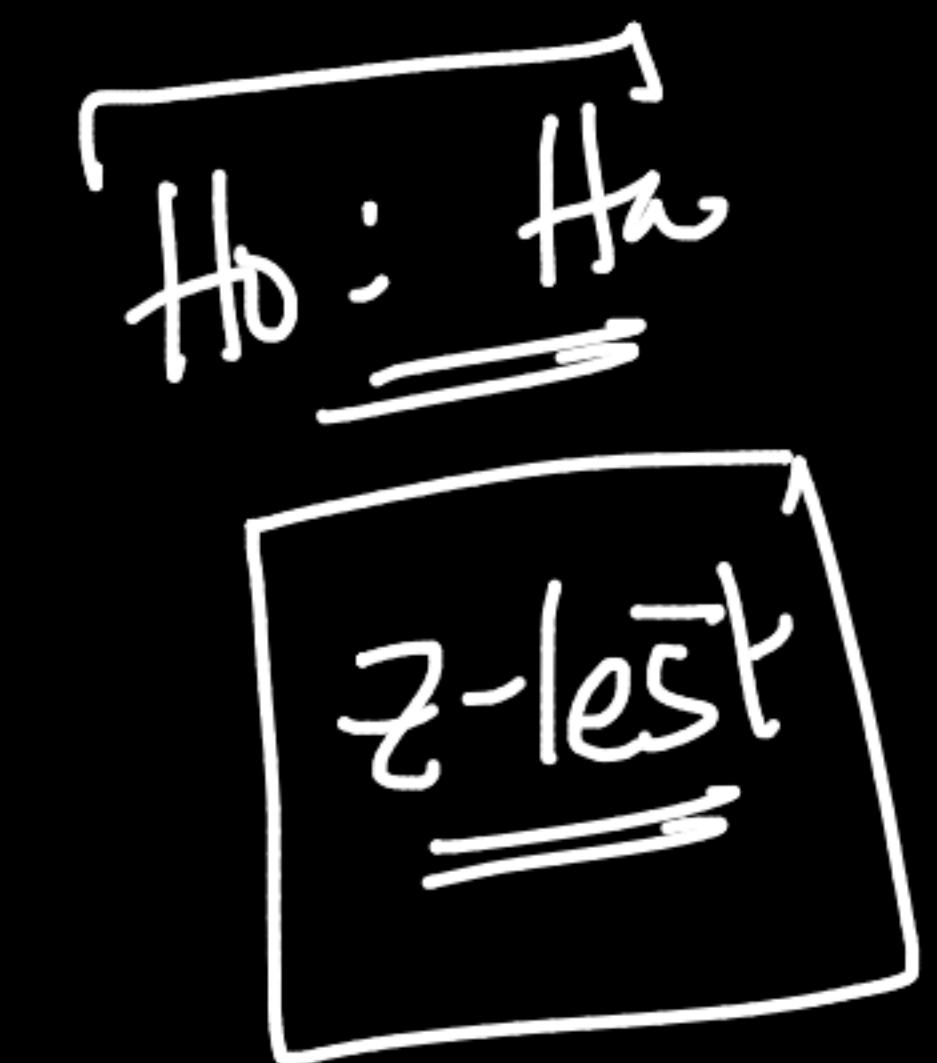
Pdf

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}$$

$$\underline{\mu=10}; \quad \underline{\sigma=1}$$



x_1, x_2, \dots, x_{100} : numerical
↓
CDF \rightarrow PDF

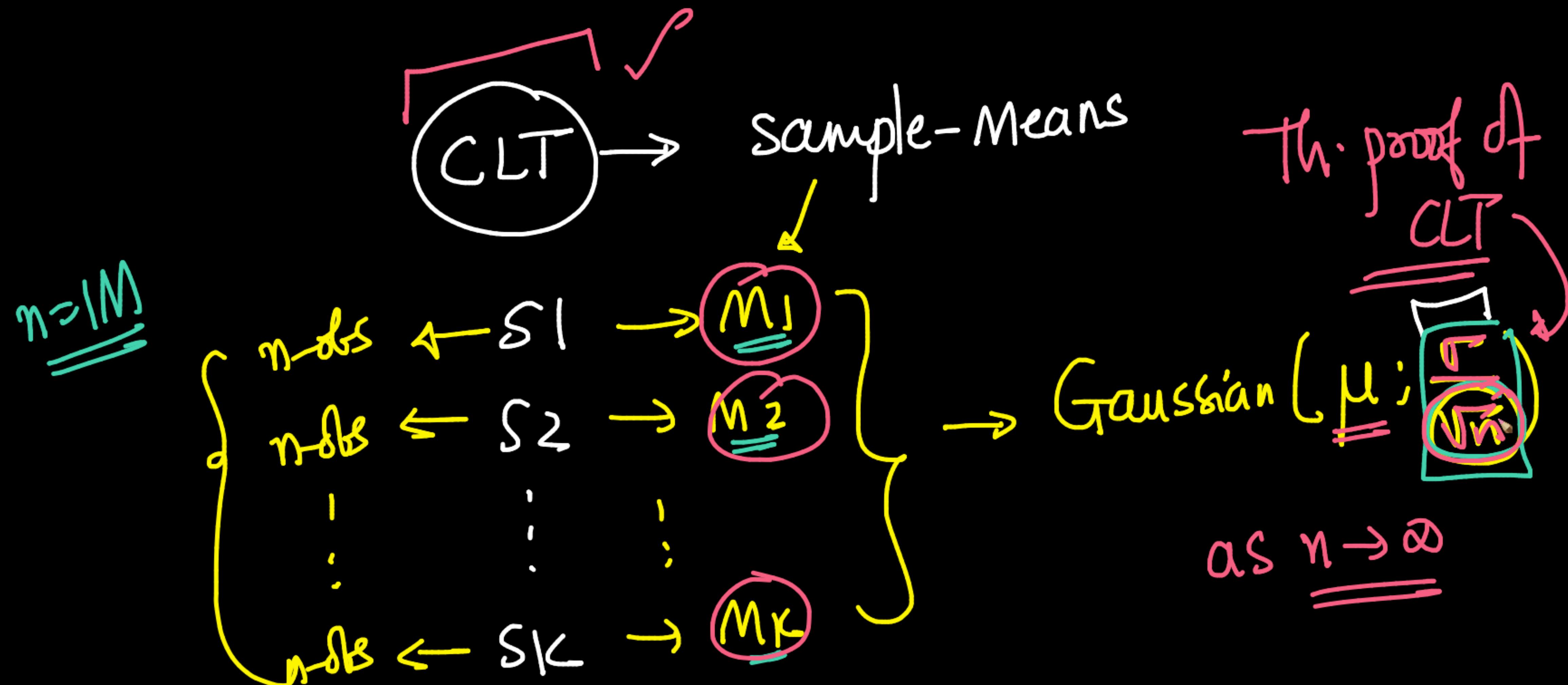


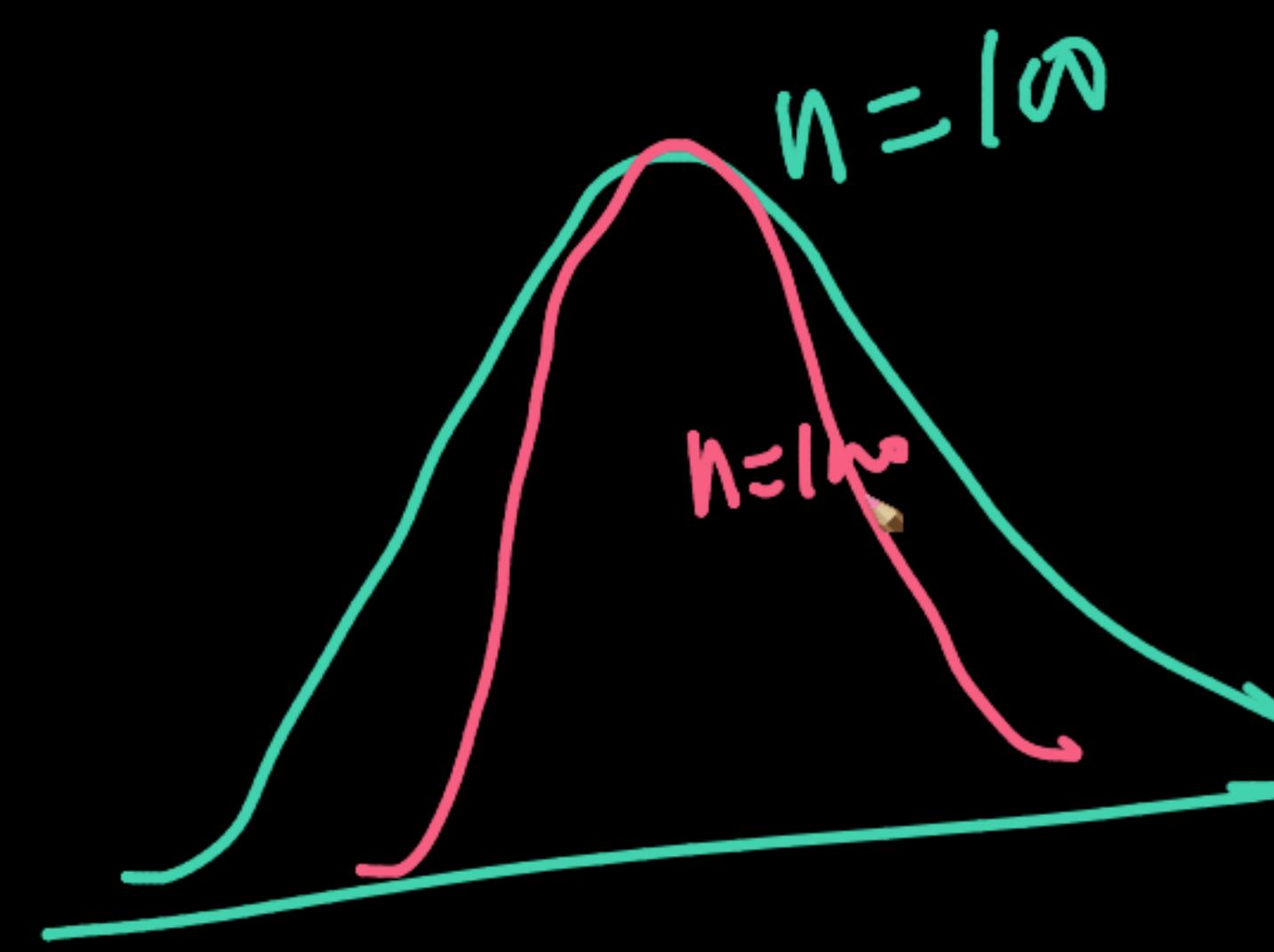
tem vs placebo

p-val < 1% \Rightarrow reject H₀

T_z

{ 95%.
95%. placebo



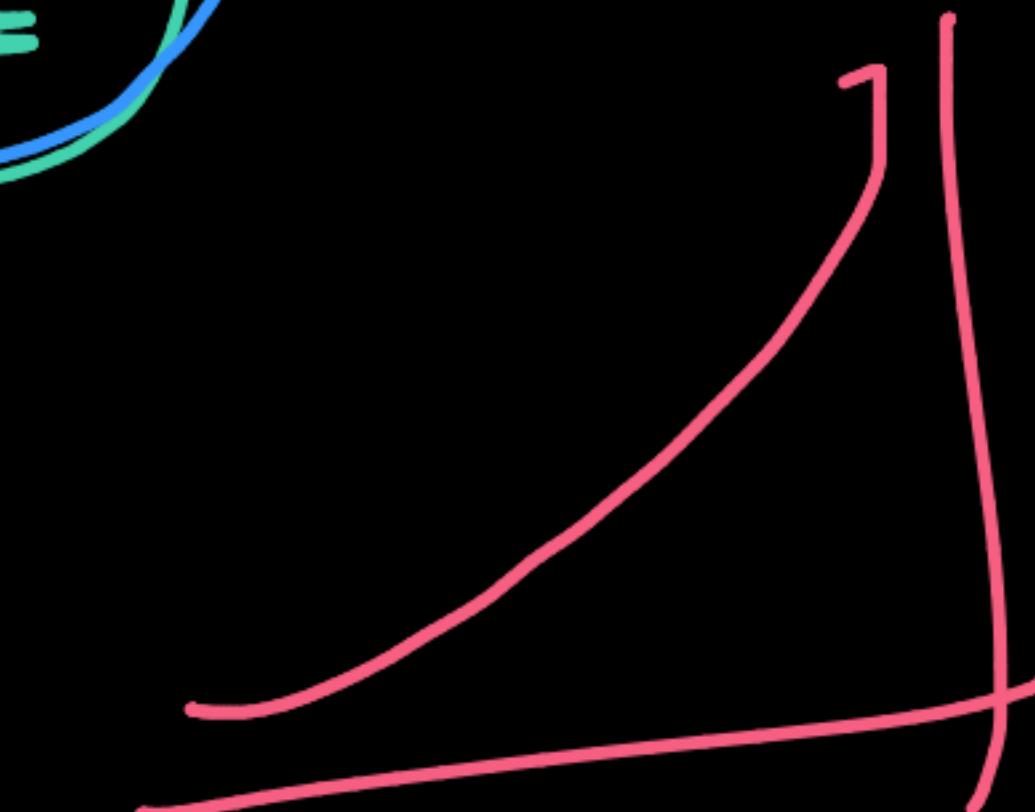
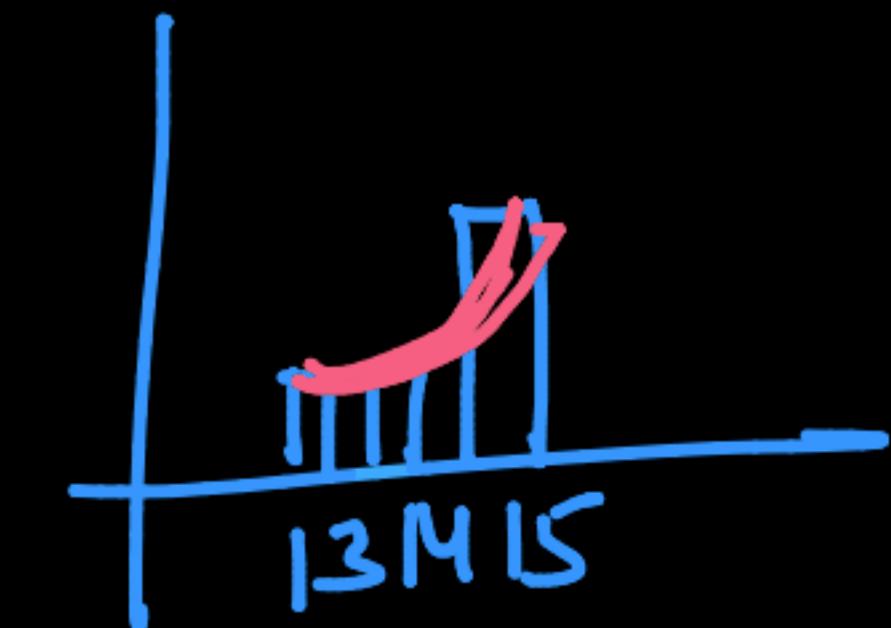


boxcox

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

log-normal x_i 's:

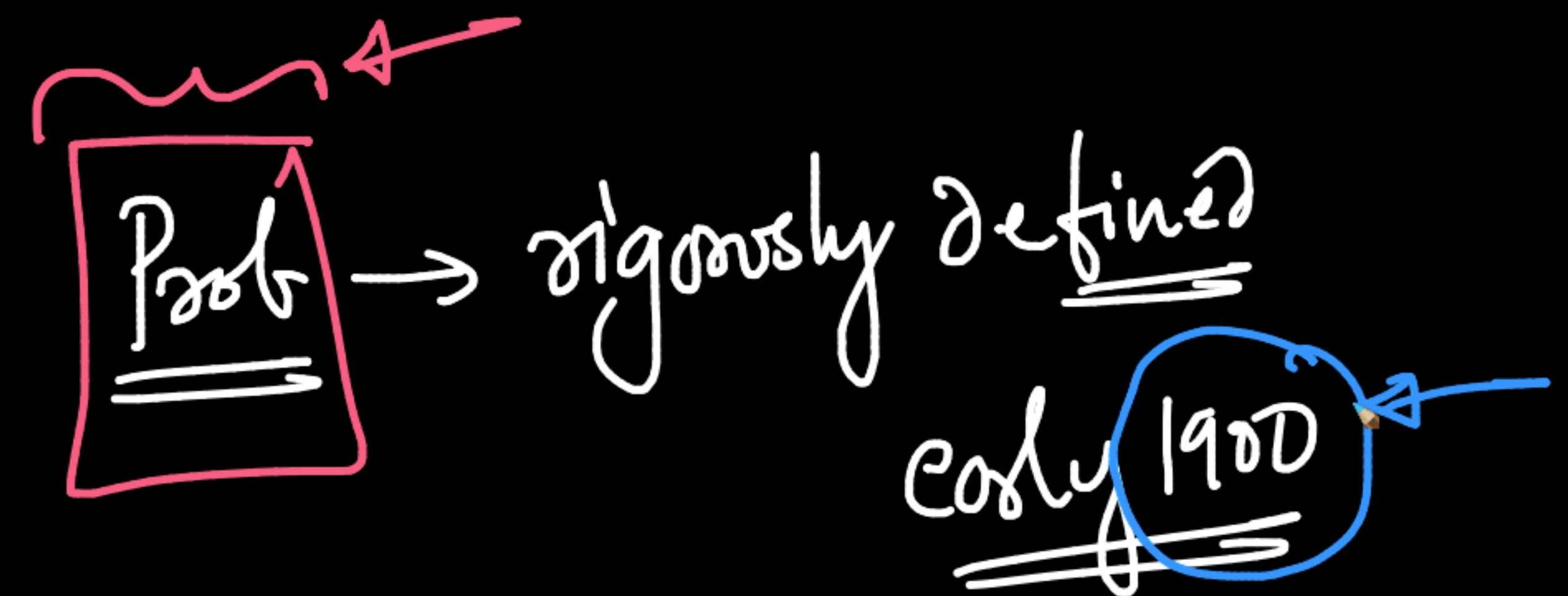
13, 14, 15, 15
 ↓ ↓ ↓ ↓


 y_i 's:KS-test Q-Q plot

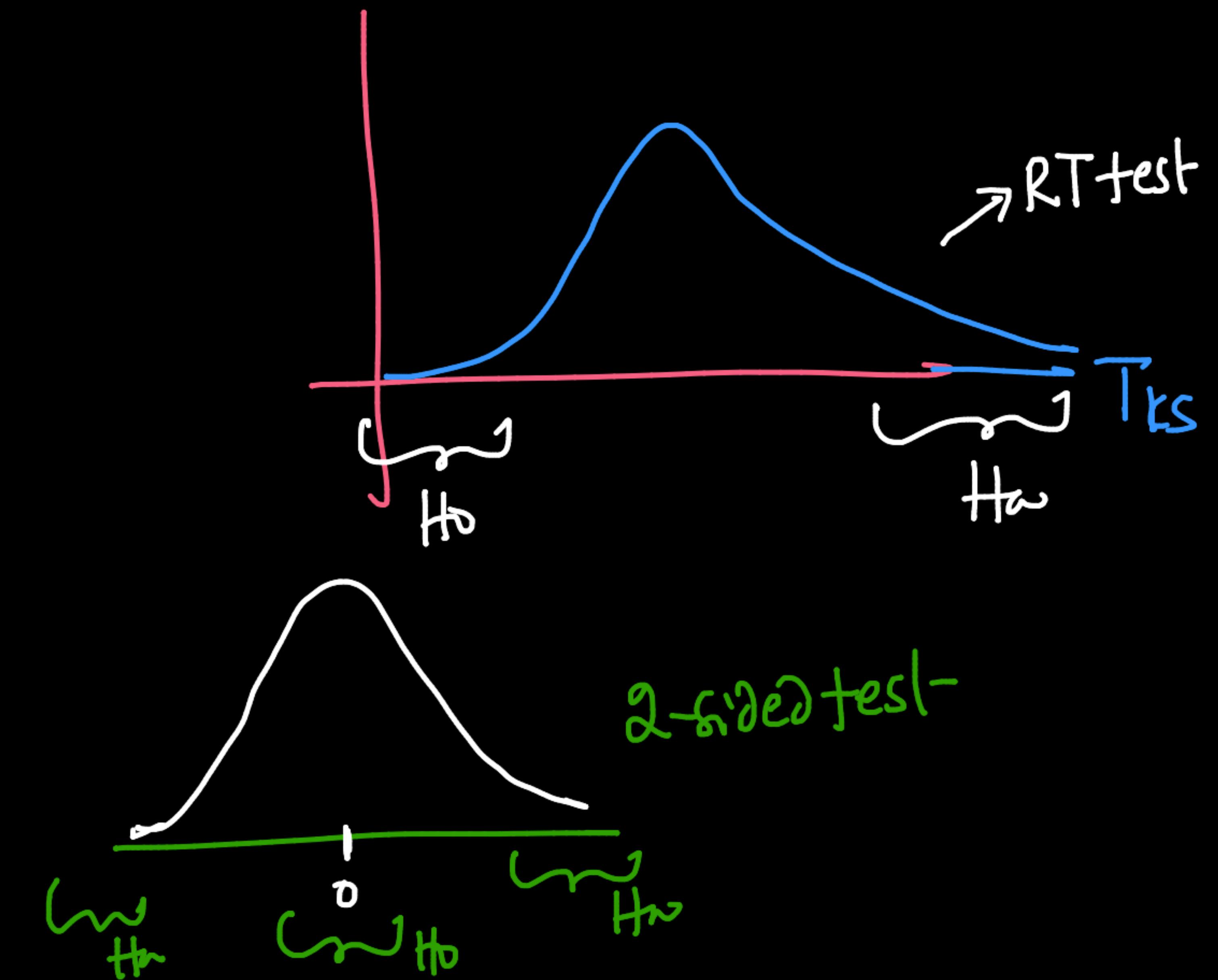
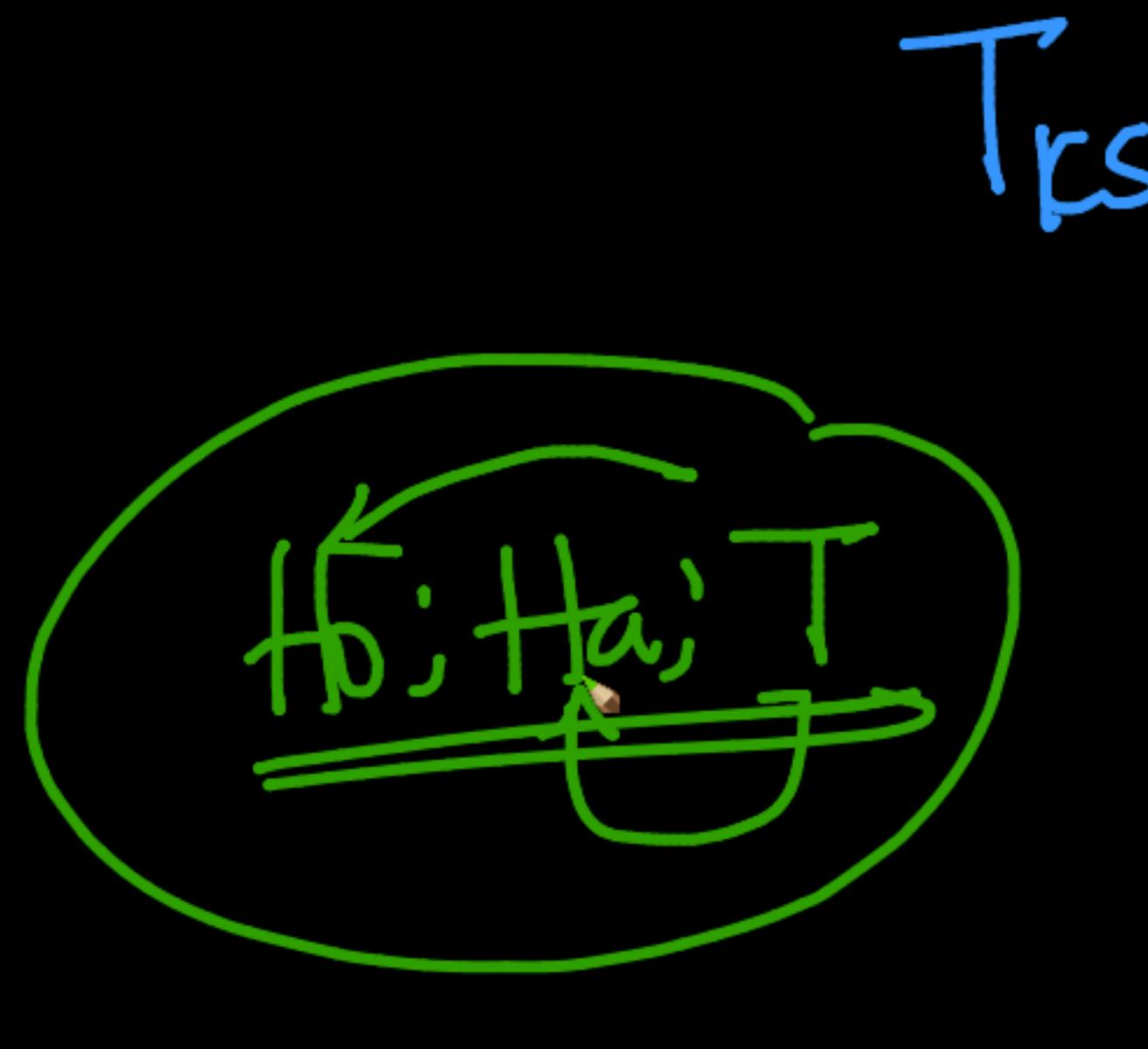


pdf

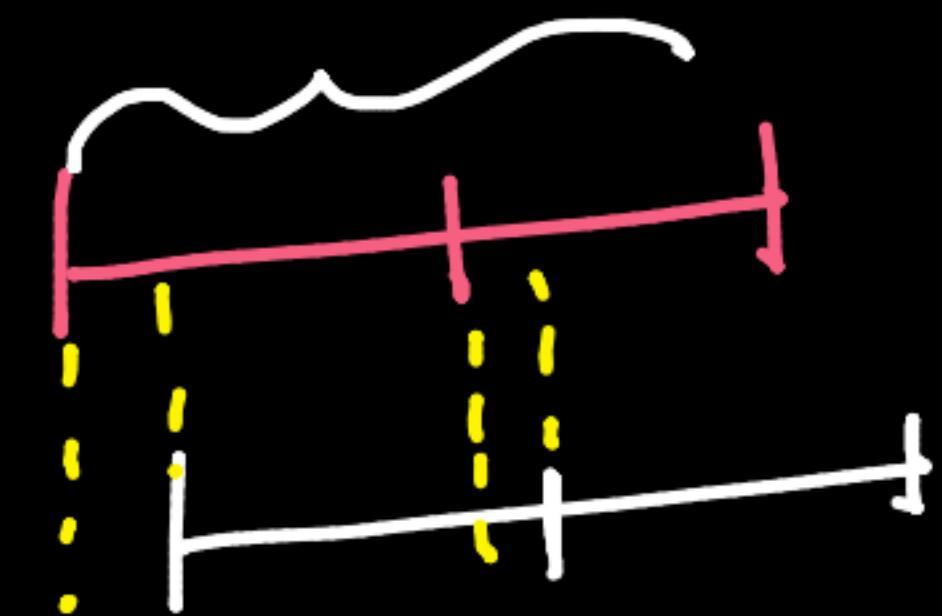
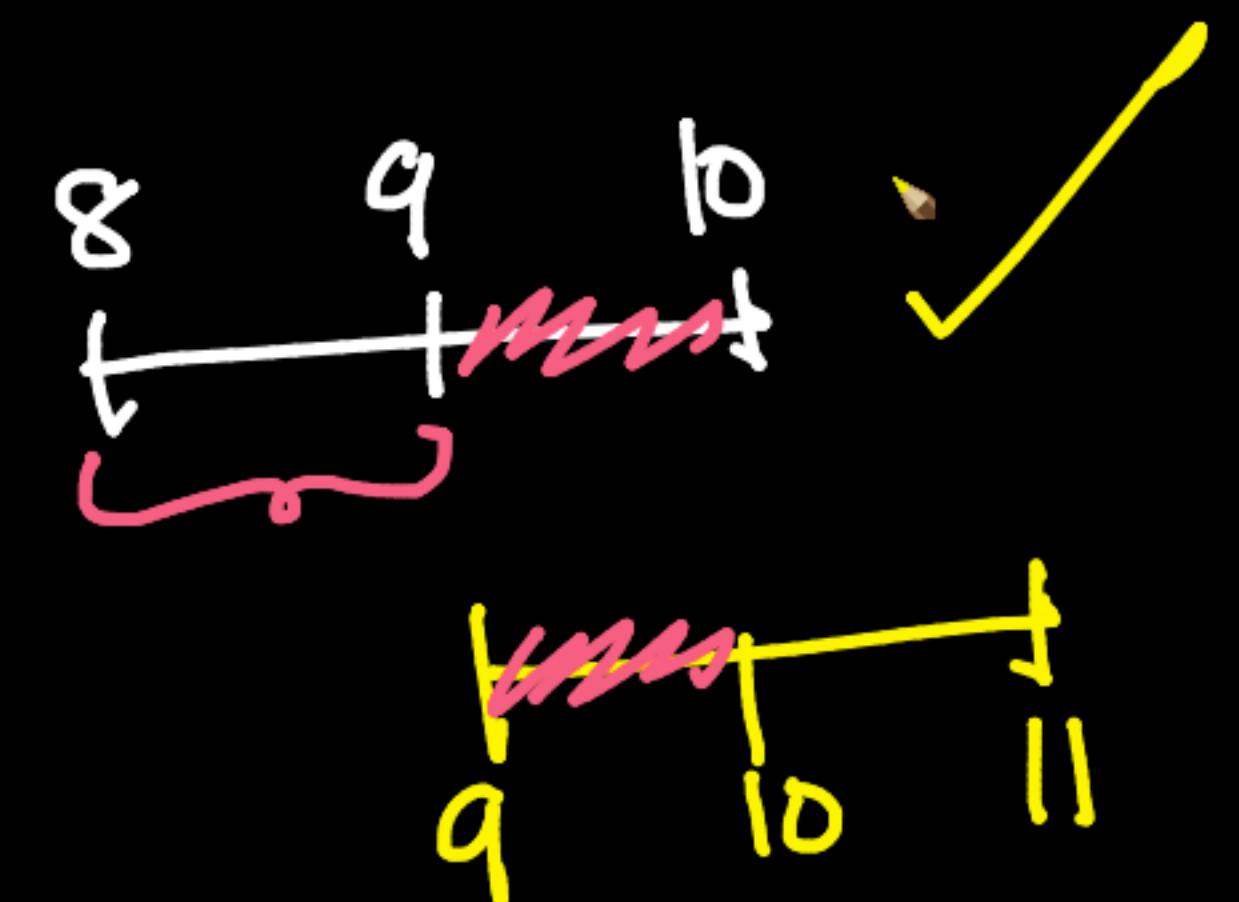
looks similar to smoothened histograms



cough 1900



✓ Med 1: Mean rec-time
Med 2: Mean dec-time



→ no conclusion

boxcox and z-score.ipynb - Colaboratory | scipy.stats.boxcox — SciPy v1.8.0.dev0+11.gf3f3d9d | Anderson–Darling test - Wikipedia | KStest_Ttest.ipynb - Colaboratory | Power transform - Wikipedia +

en.wikipedia.org/wiki/Power_transform#Yeo–Johnson_transformation

not be a substantive problem in many applications.^{[6][7]}

Box–Cox transformation [edit]

The one-parameter Box–Cox transformations are defined as

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

and the two-parameter Box–Cox transformations as

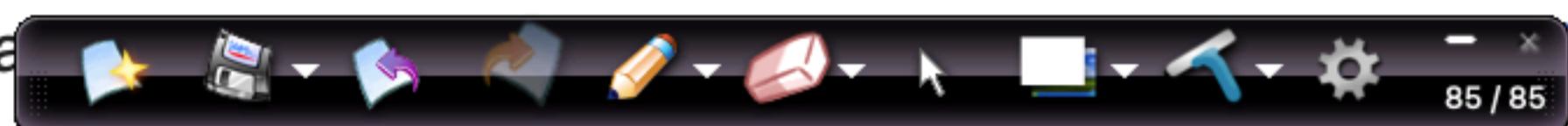
$$y_i^{(\lambda)} = \begin{cases} \frac{(y_i + \lambda_2)^{\lambda_1} - 1}{\lambda_1} & \text{if } \lambda_1 \neq 0, \\ \ln(y_i + \lambda_2) & \text{if } \lambda_1 = 0, \end{cases}$$

as described in the original article.^{[8][9]} Moreover, the first transformations hold for $y_i > 0$, and the second for $y_i > -\lambda_2$
^[8]

The parameter λ is estimated using the [profile likelihood](#) function and using goodness-of-fit tests.^[10]

Confidence interval [edit]

Confidence interval



85 / 85

...y constructed using Wilks's theorem on the profile likelihood function to find all the possible values of λ that fulfill the following restriction:^[11]

$$\ln(L(\lambda)) > \ln(L(\hat{\lambda})) - \frac{1}{2}\chi^2_{1,1-\alpha}.$$