

Topics: → simple, visual

- { ① QQ-plot → Is the data Gaussian?
  - ↳ Why, How, code, Variations
- { ② Transformations → convert data to Gaussian dist.
  - ↳ Box-Cox : code ✓
- { ③ Mode; Bi-modal distributions

~~OPS:~~

Ques → tab

↳ "END! ... "

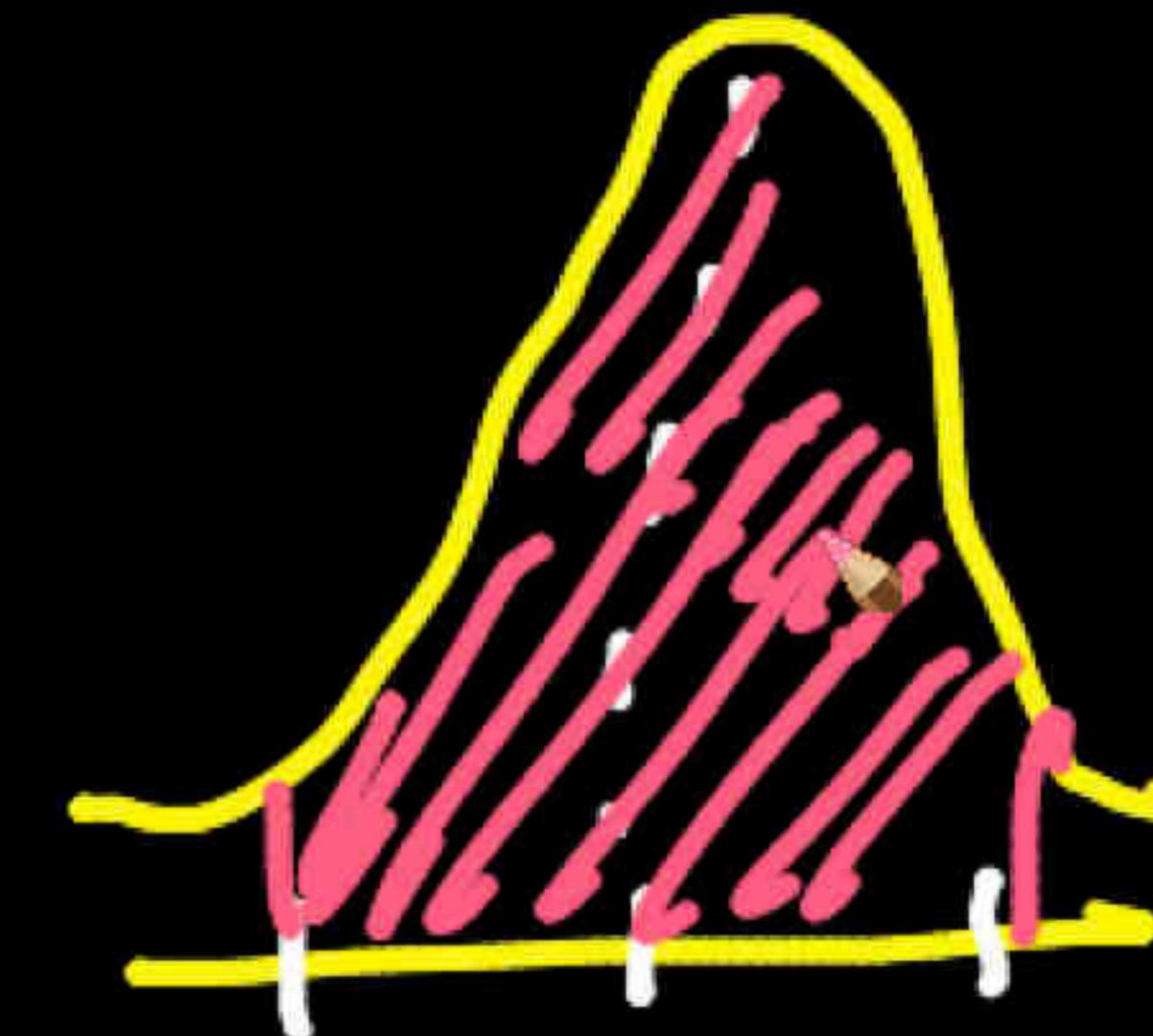
chat → interactivity, Y/N



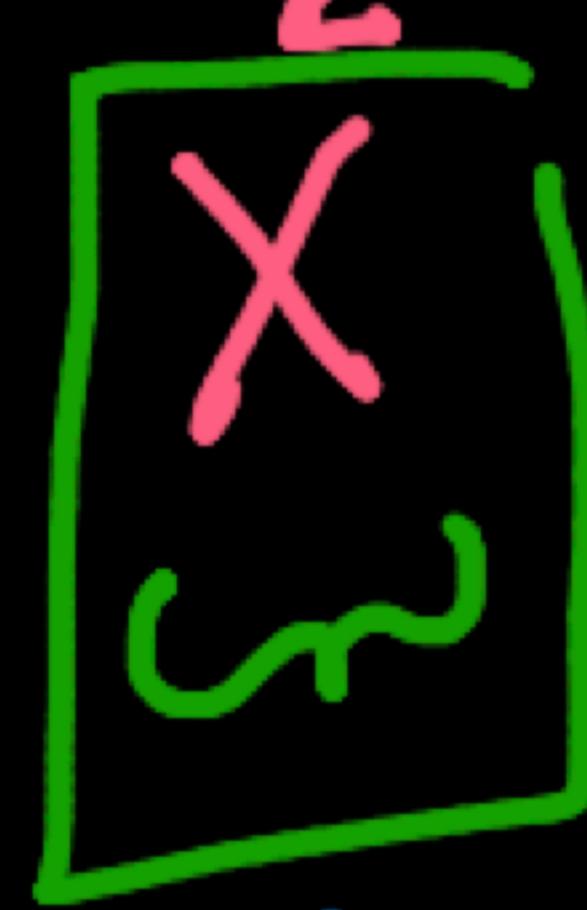
C.I using CLT



$$\left[ \mu - 1.96\sigma \leq \bar{x} \leq \mu + 1.96\sigma \right]$$


$$\bar{X}_n \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right) \text{ as } n \rightarrow \infty$$

height



height

$x_1, x_2, \dots, x_n \rightarrow 200$  (let)

obs

m,s

QQ-plot



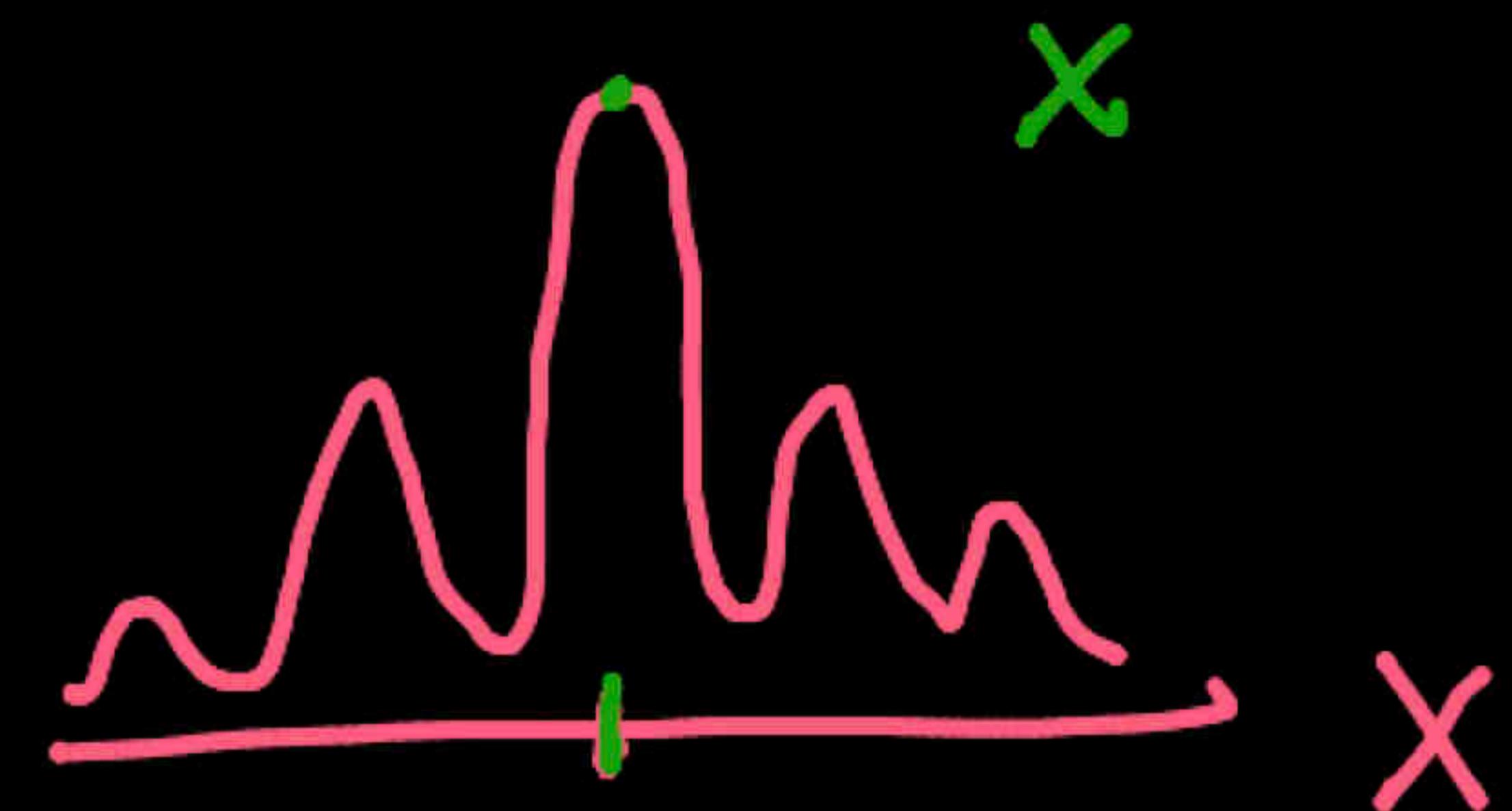
Is  $X$  Gaussian dist?

$\mu, \sigma$

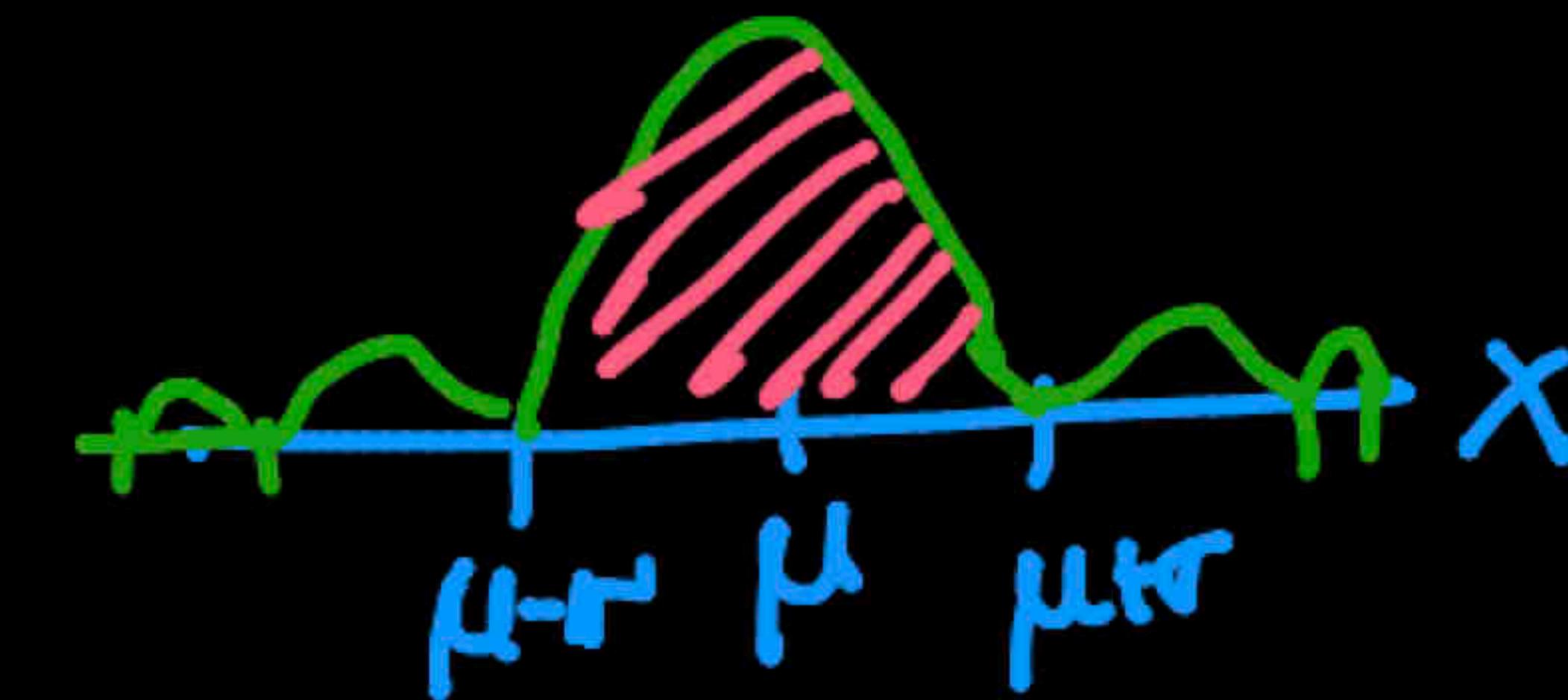
→ PDF

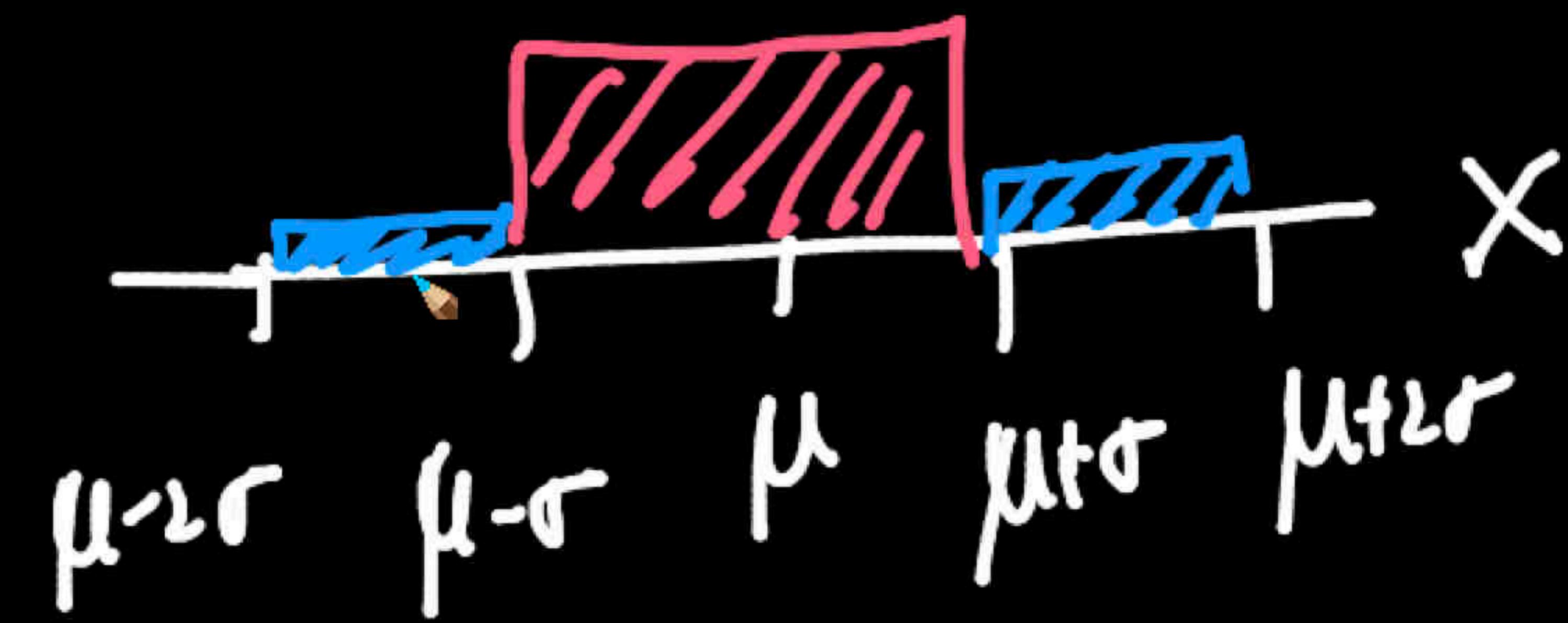


→ mean = mode = median



→ 68-95-99 rule





A logo consisting of a white outline of a house on a black background. The roof is orange, and the door is white with the text "KS test" above it. Two white diagonal lines extend from the bottom right corner of the house outline.

A large, hand-drawn style pink arrow pointing upwards and to the right, indicating the direction of the next section.

$N(\mu, \sigma)$

# ↳ equations

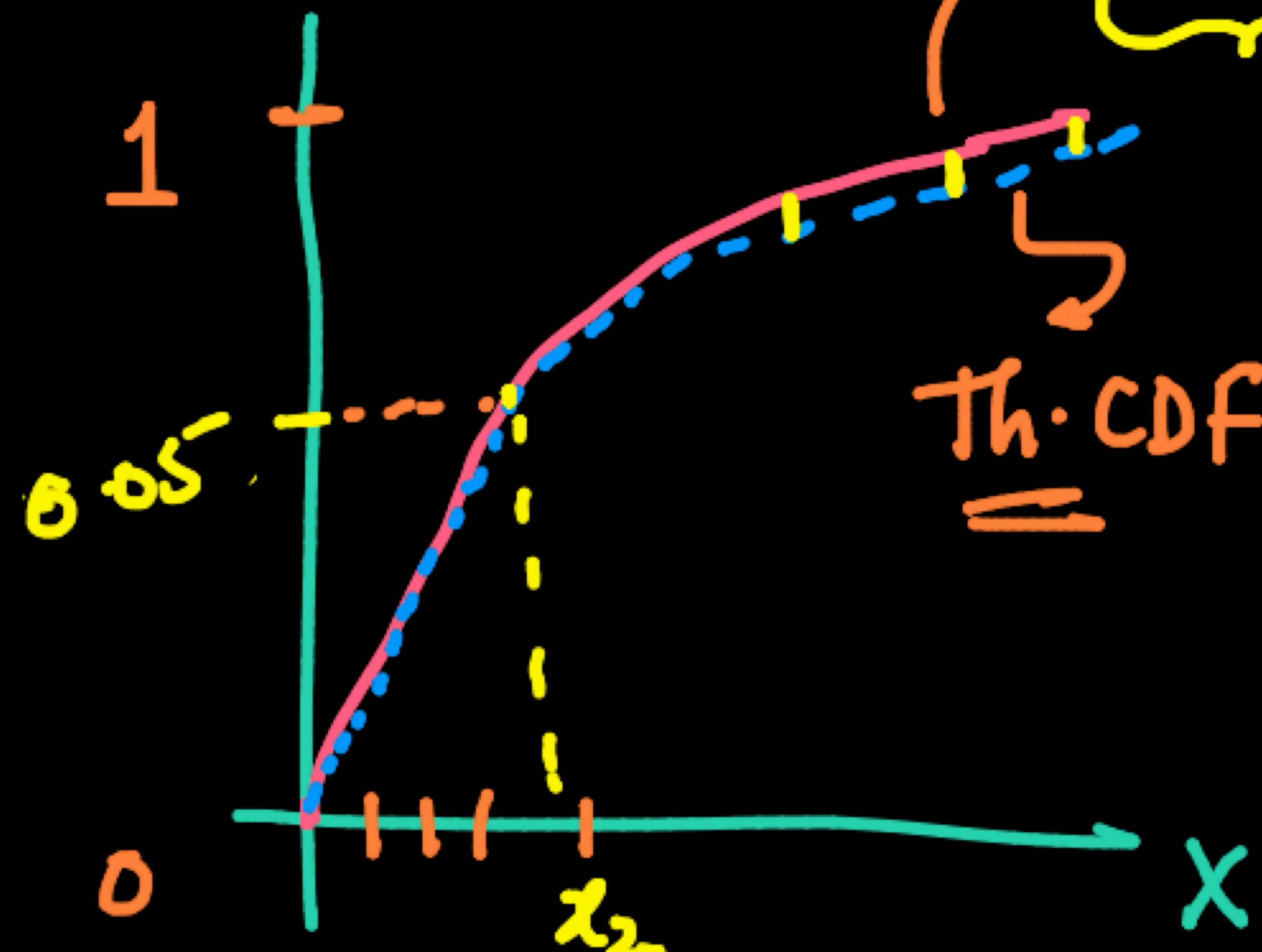
$x_1 x_2 \dots$

xn✓

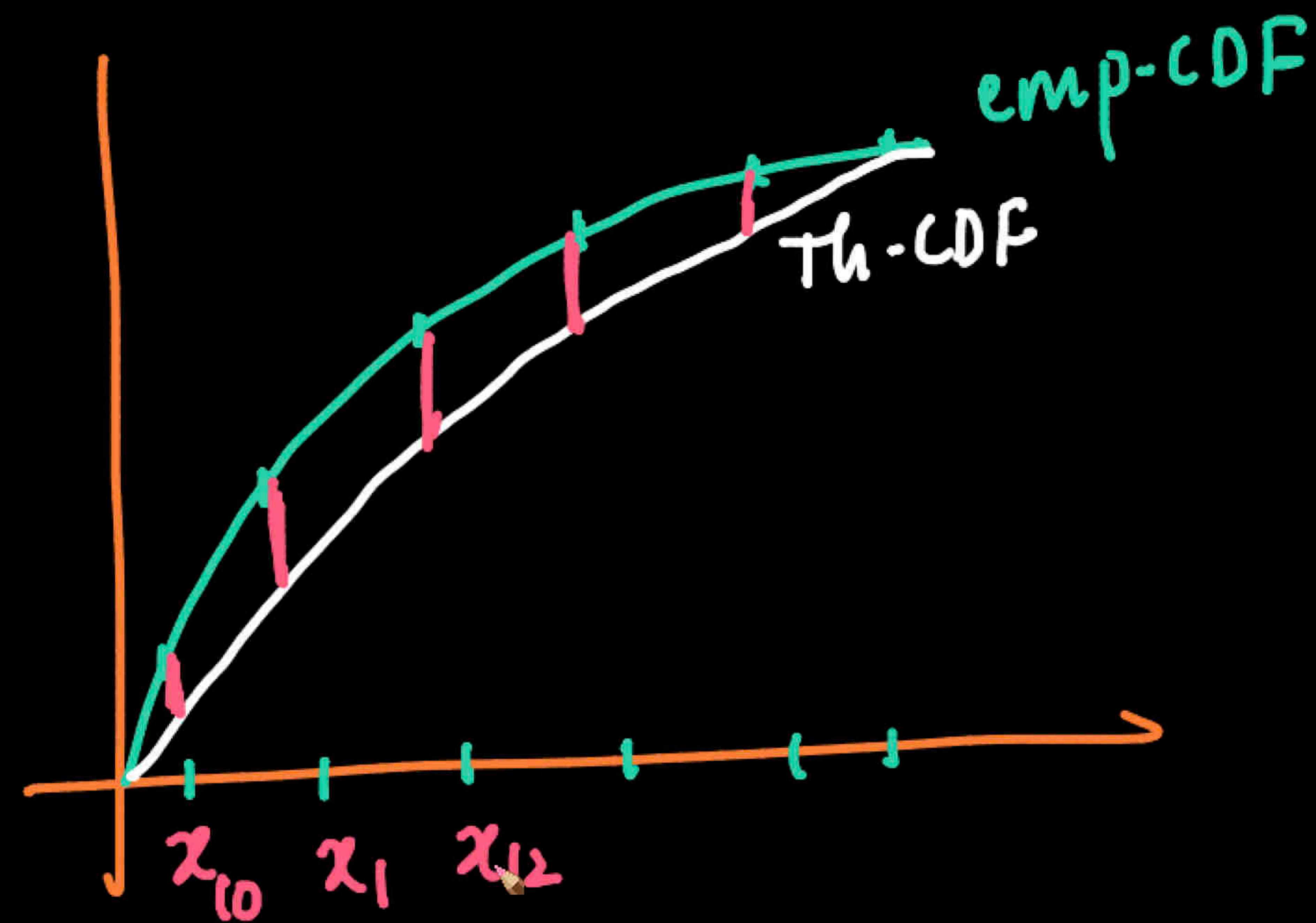
CDF

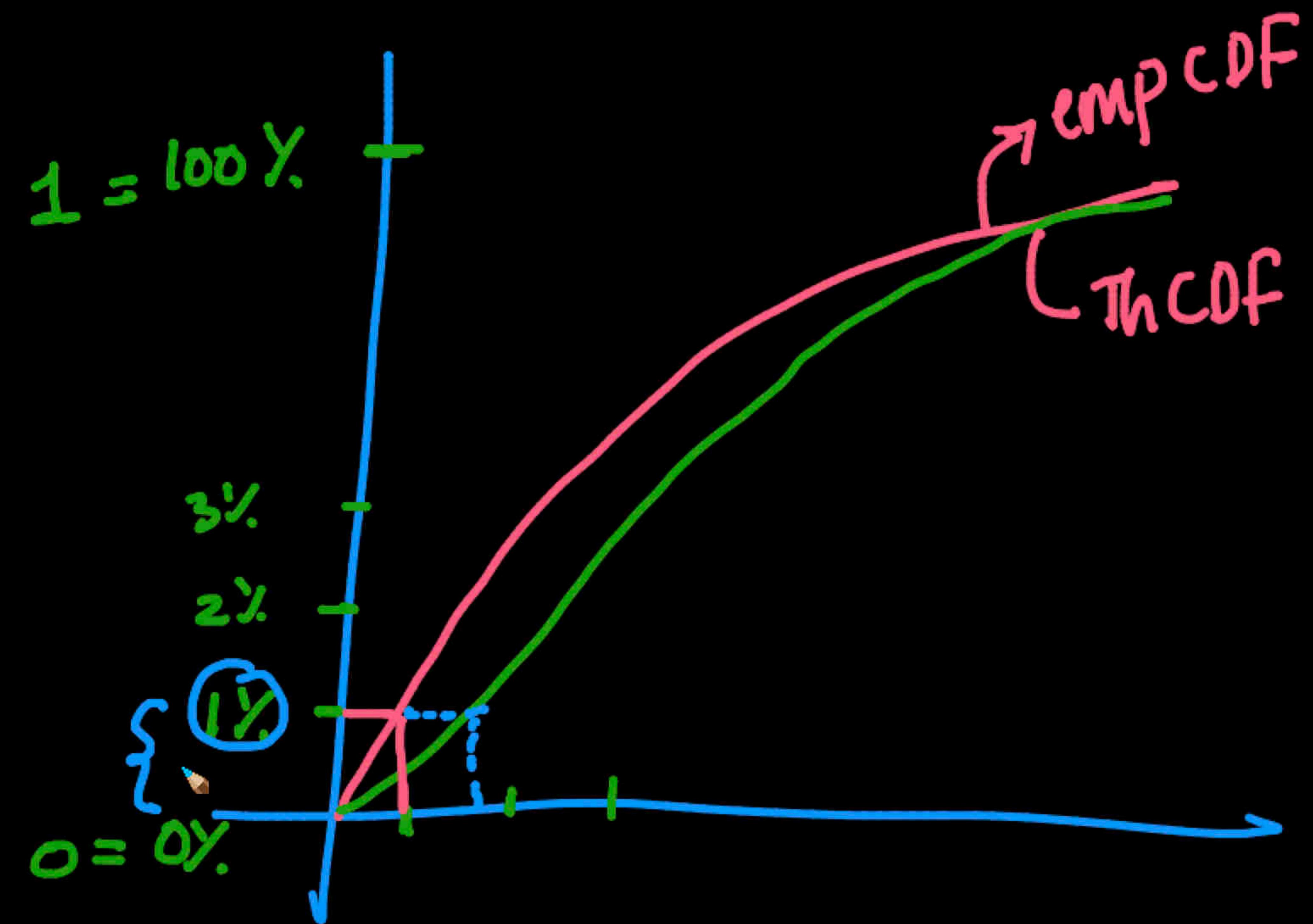
MIS

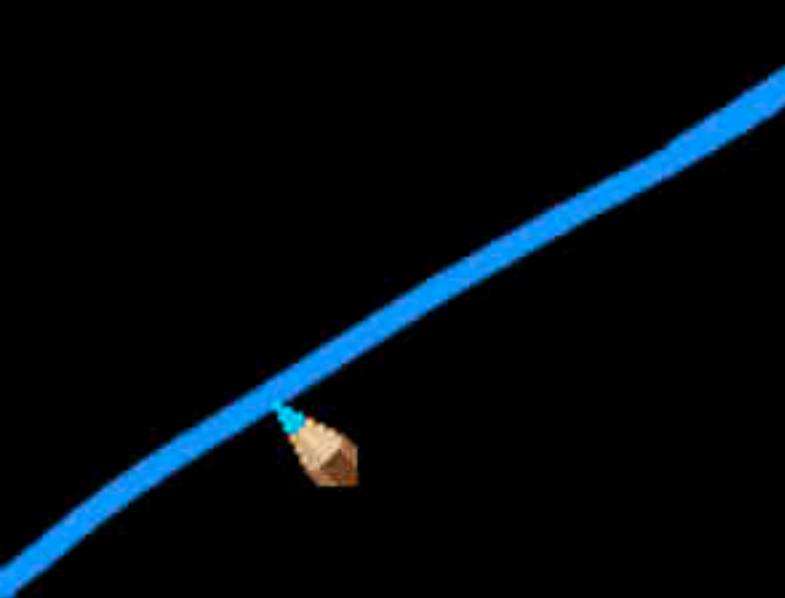
$$\mu \approx m \quad \tau \approx s$$



$$P(X \leq x_i)$$









$X \Rightarrow x_1, x_2, \dots, x_n \rightarrow M, S$

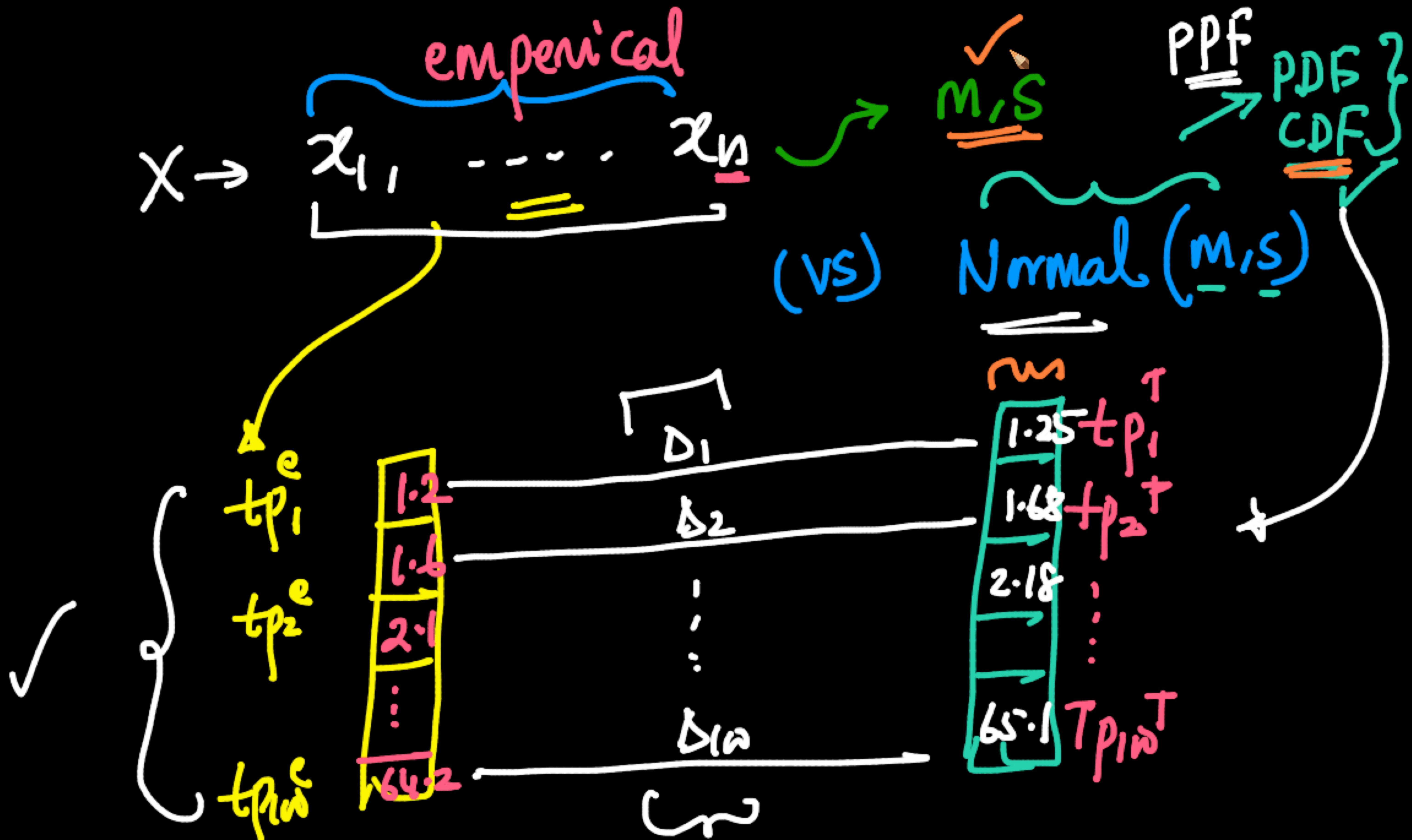


$$\left( \frac{x_i - M}{S} \right)$$

Standardize (not mandatory)

Normal

$$N(0, 1) = Z(0, 1)$$



PPF

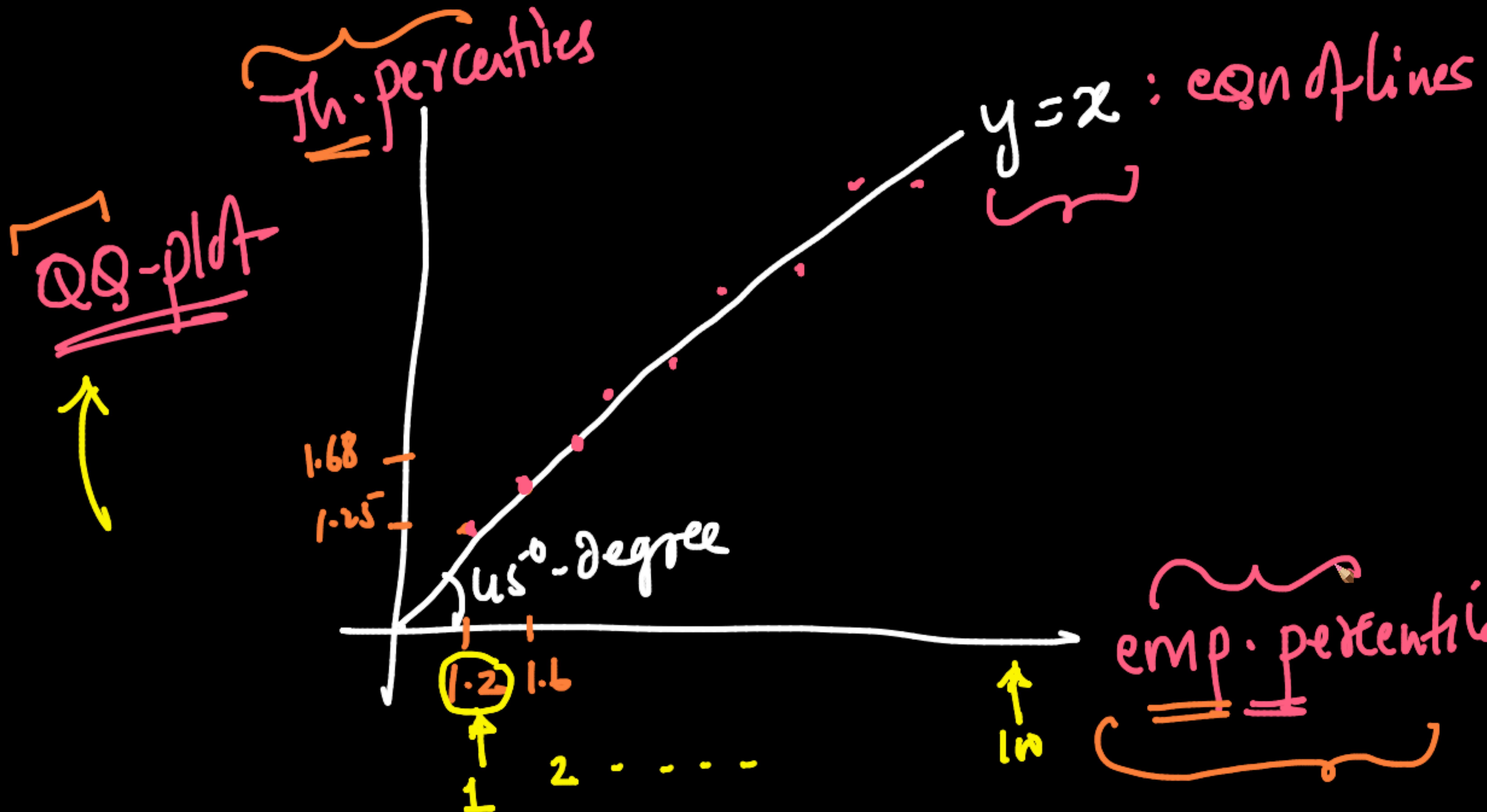
Normal( $m, s$ )

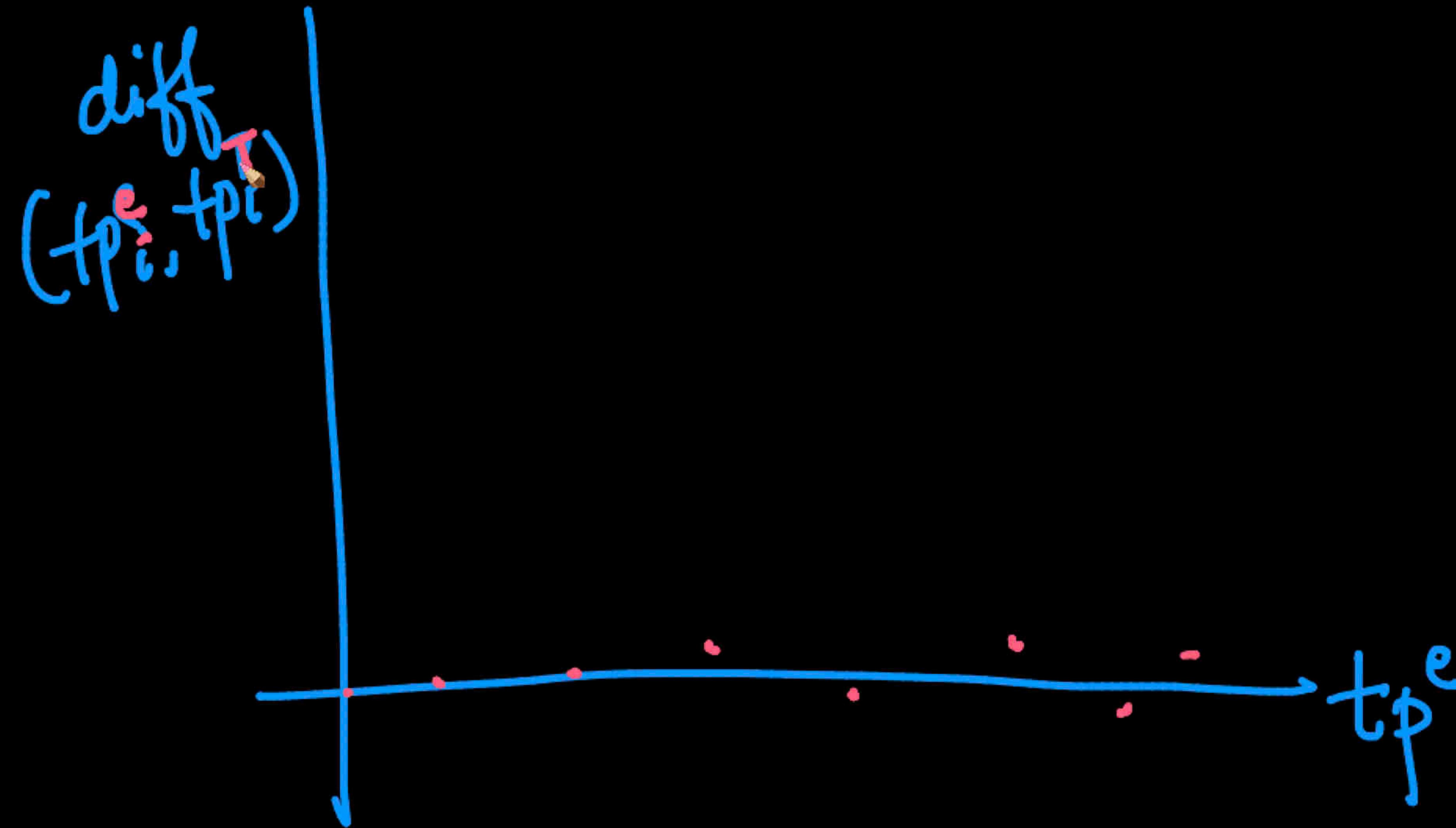
$$P(x \leq \underline{\text{val}}) = \underline{\text{Y.}}$$

$\lceil$        $\underline{\text{P}}$

PDF(val)

$(tp_i^e, tp_i^t)$  for all  
 $i : 1 \text{ to } n$





Dis-adv:

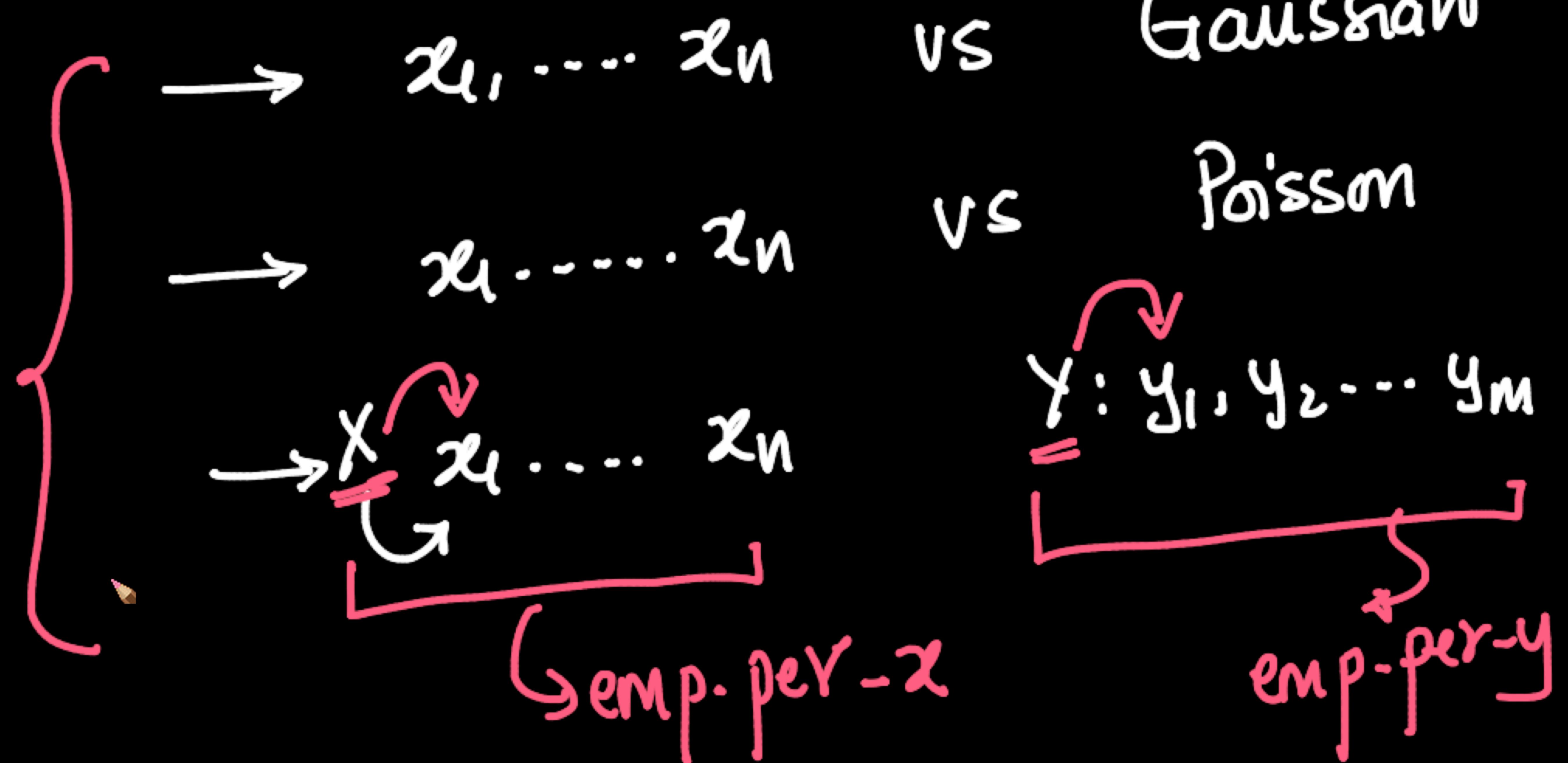
↳ [ visual  
only ] ✓

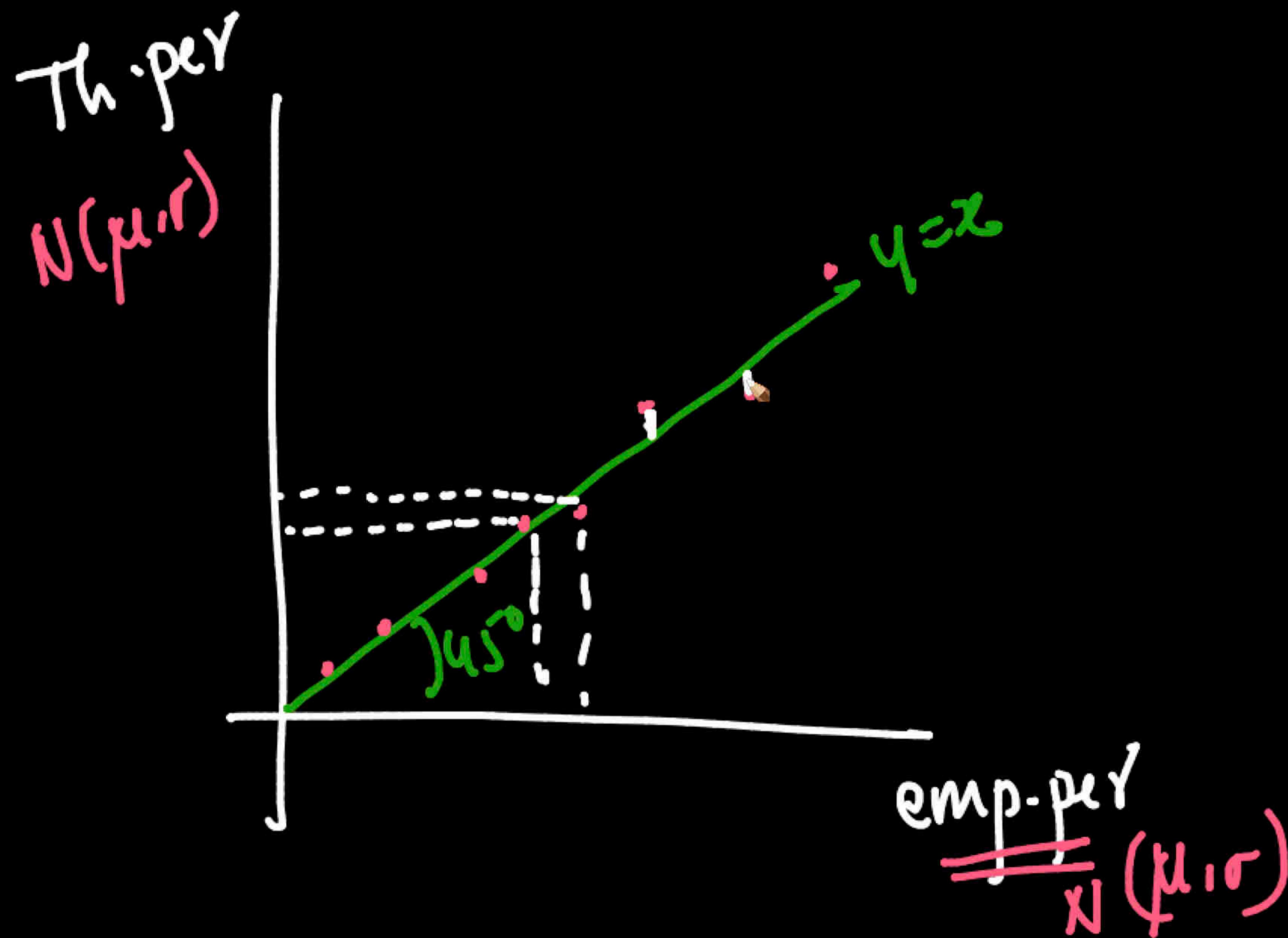
↳ how confident

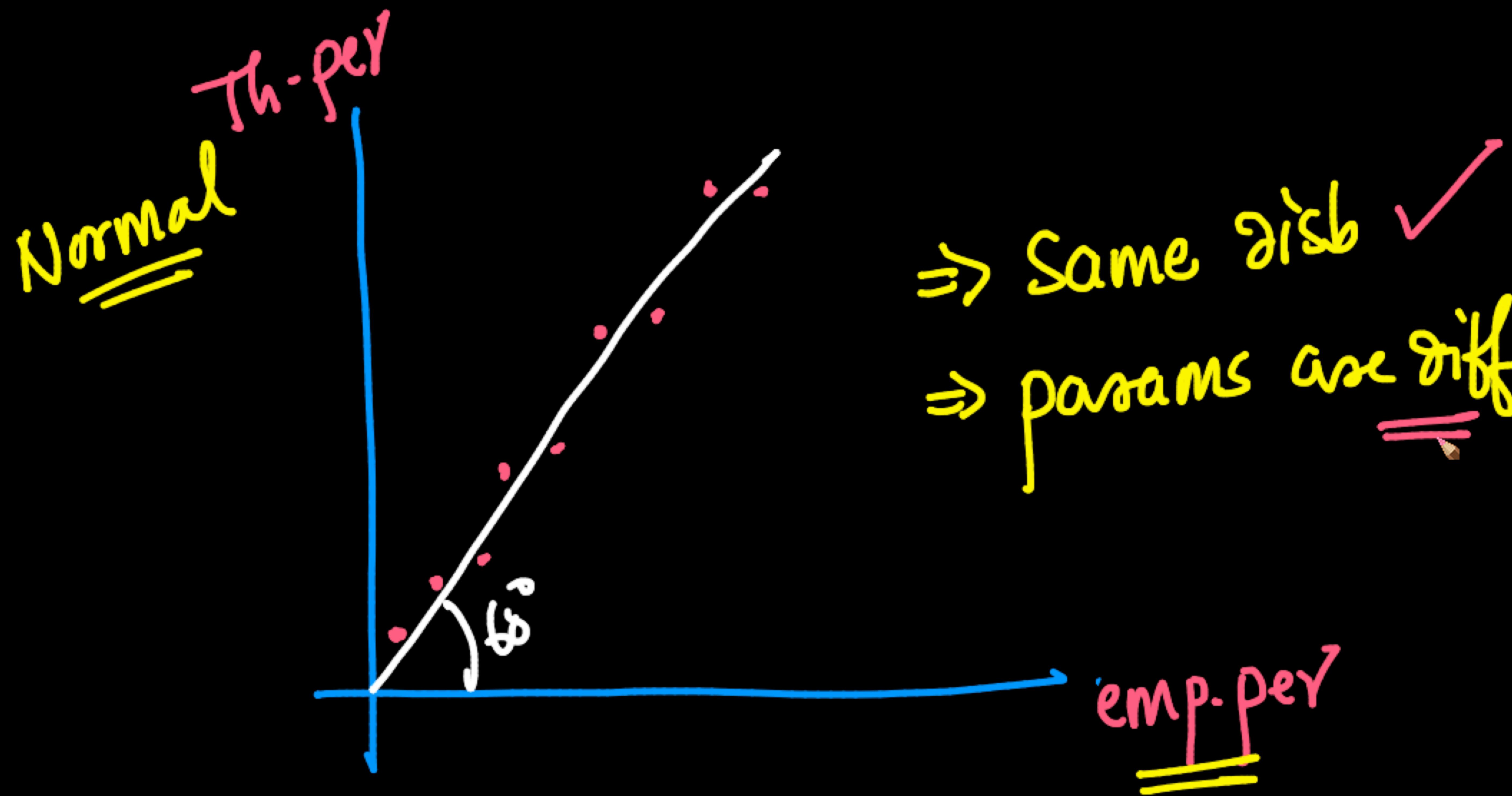
↳ Hyp-testing

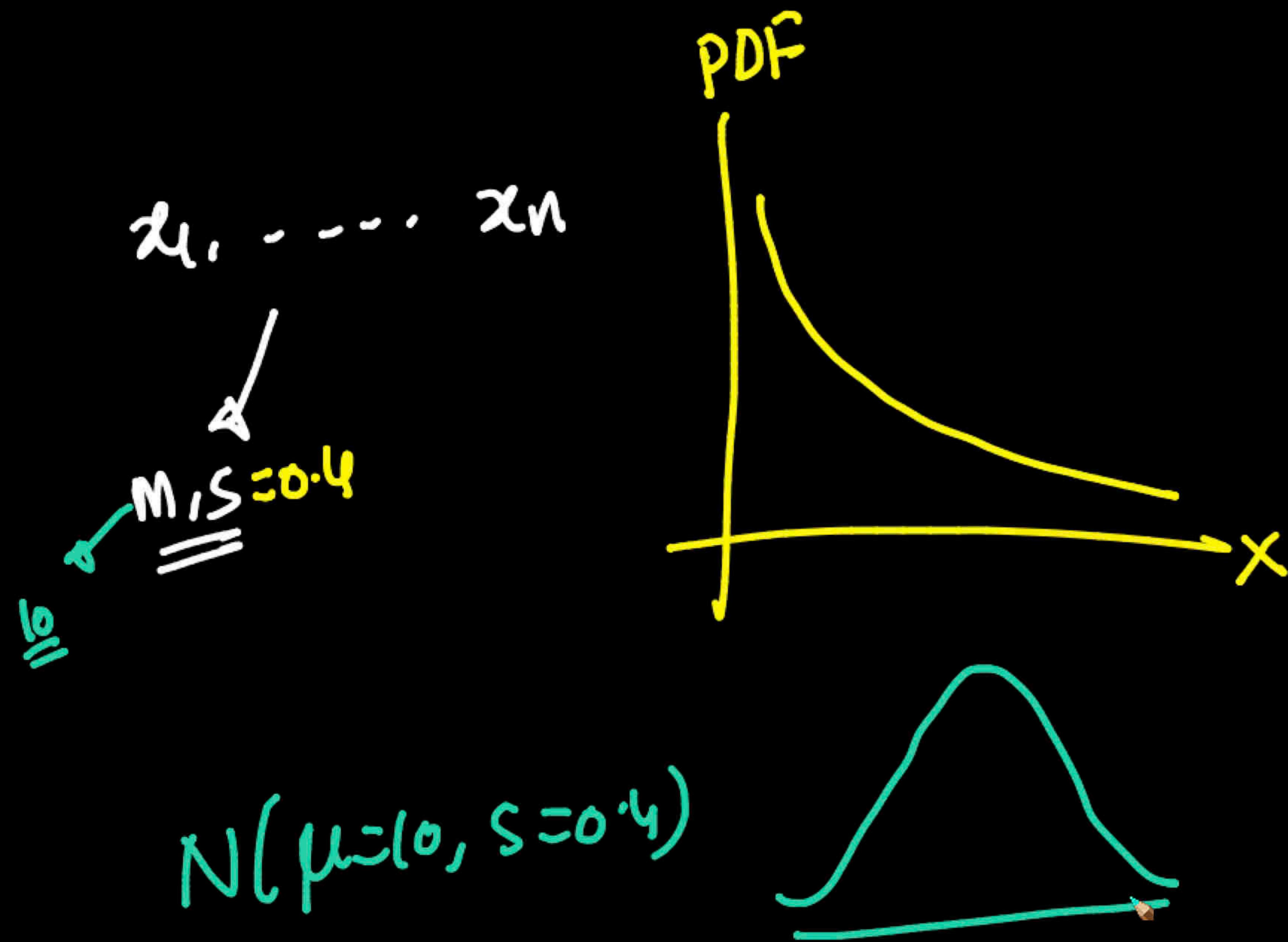
$$\tilde{x}_1 \dots \tilde{x}_n \sim N(\mu, \sigma)$$

# QQ-plot









+ Code + Text

Reconnect



```
[ ] [-0.01101/05, 0.04/20004, 0./1049555, 0.5/500/01, 0.52028192,  
0.46012344, 0.92880202, 0.67717555, 0.04629136, 0.47546512,  
1.4513086 , 1.45343272, 1.54991229, 0.62037232, 0.01407354,  
0.46979478, 0.05595689, -1.73249288, 0.23003225, 1.29352827,  
0.80189453, 1.61966331, 0.69681106, 1.03215339, -1.21549361,  
0.93475221, 1.30143537, 0.7254352 , 0.22841529, 1.50249735,  
-0.02415314, -0.18205881, 0.95388083, 0.66182587, 0.08282857,  
1.53432986, 1.07818559, 1.04804152, 0.62920033, 0.2221568 ,  
1.11689153, 0.70328342, 1.48907562, -0.85967934, 0.37330663,  
0.10042743, 0.43601618, -0.84872277, -0.18902961, 1.16872747,  
0.49445364, 0.97912906, 0.16970087, 1.43121388, 0.67825154,  
0.8233865 , 1.20263091, 0.49206124, 0.34548617, 1.58287164]
```

[ ] x = np.array(data)



x.shape

(500, )



&lt;&gt; sns.distplot(x)



```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecate  
warnings.warn(msg, FutureWarning)  
<matplotlib.axes._subplots.AxesSubplot at 0x7f146d788c50>
```



+ Code + Text

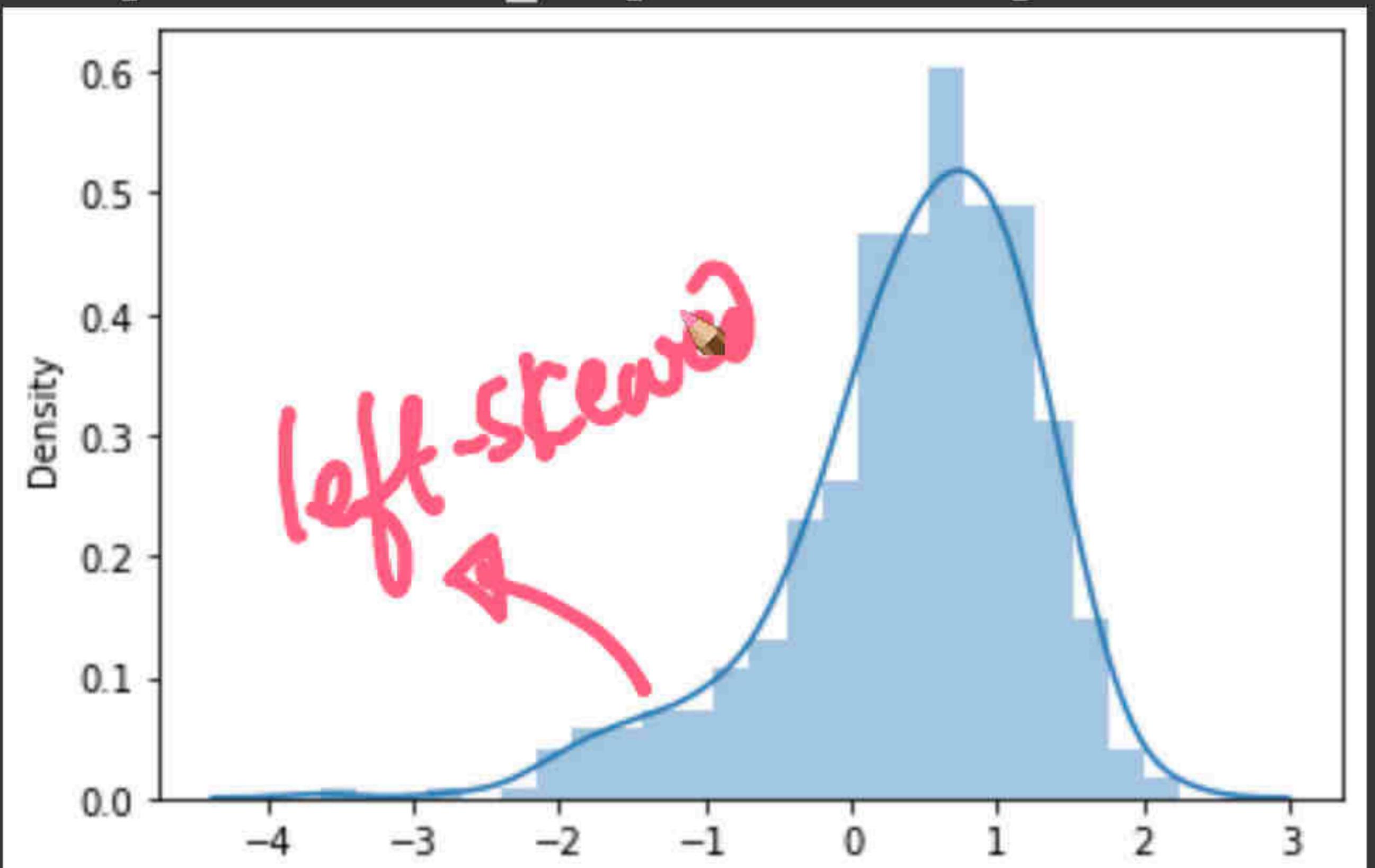
Reconnect



sns.distplot(x)



```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecation warning.warn(msg, FutureWarning)  
<matplotlib.axes._subplots.AxesSubplot at 0x7f146d788c50>
```



```
import statsmodels.api as sm  
fig = sm.qqplot(x, line='45', fit=True)  
plt.grid()
```

QQ plot and Transfo

statsmodels.graphics

scipy.stats.probplot

Power transform -

scipy.stats.boxcox -

Normal distribution

+

▼

colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=sIplfvVWgUa9

🔍 🔍 ⭐



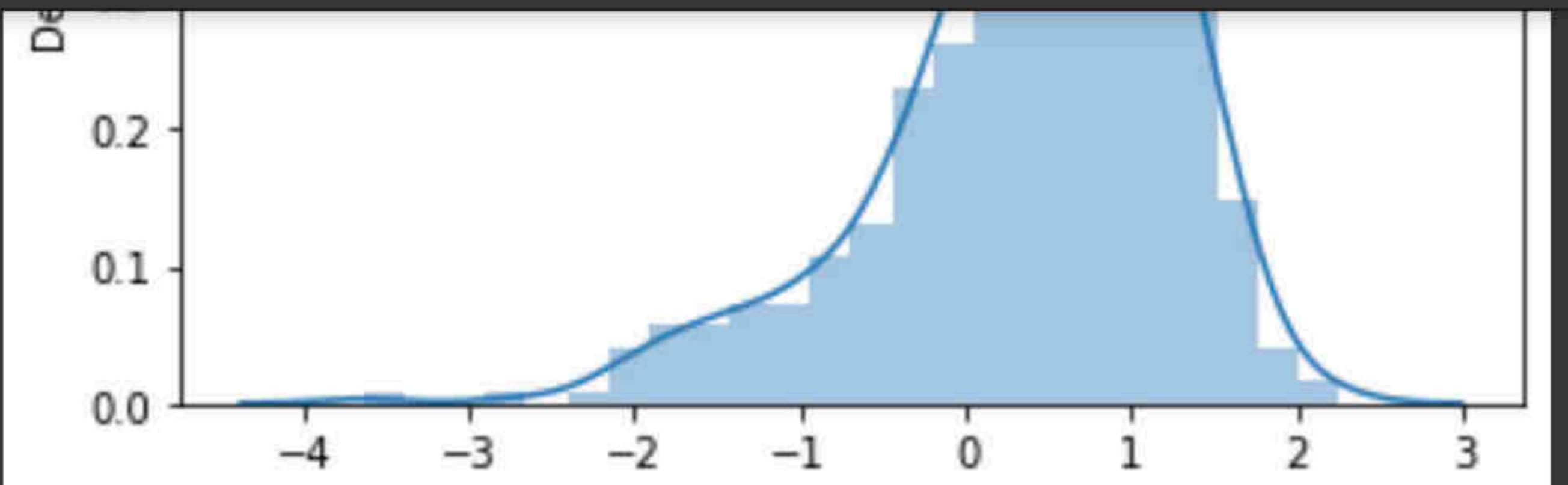
⋮

+ Code + Text

Reconnect



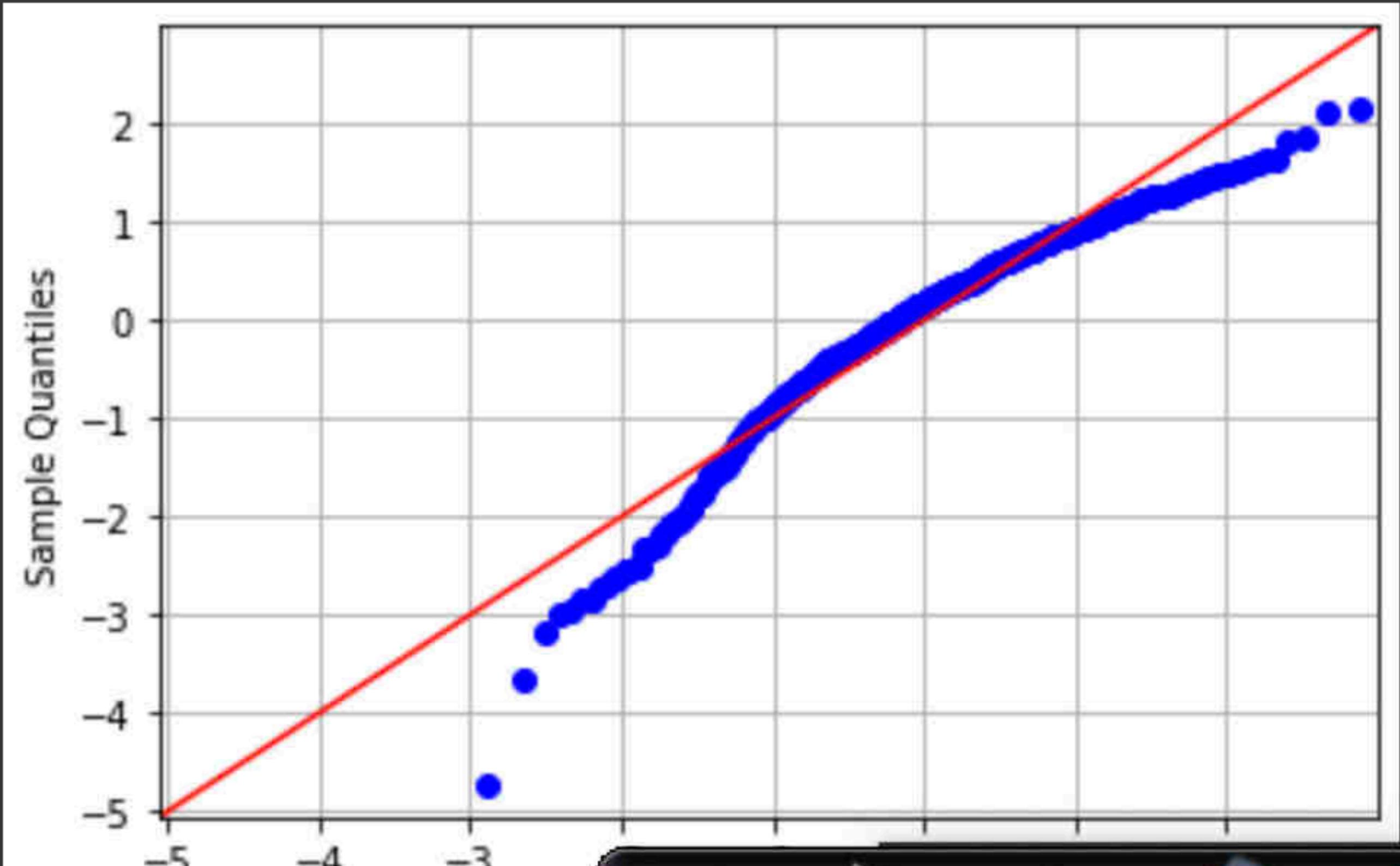
▼

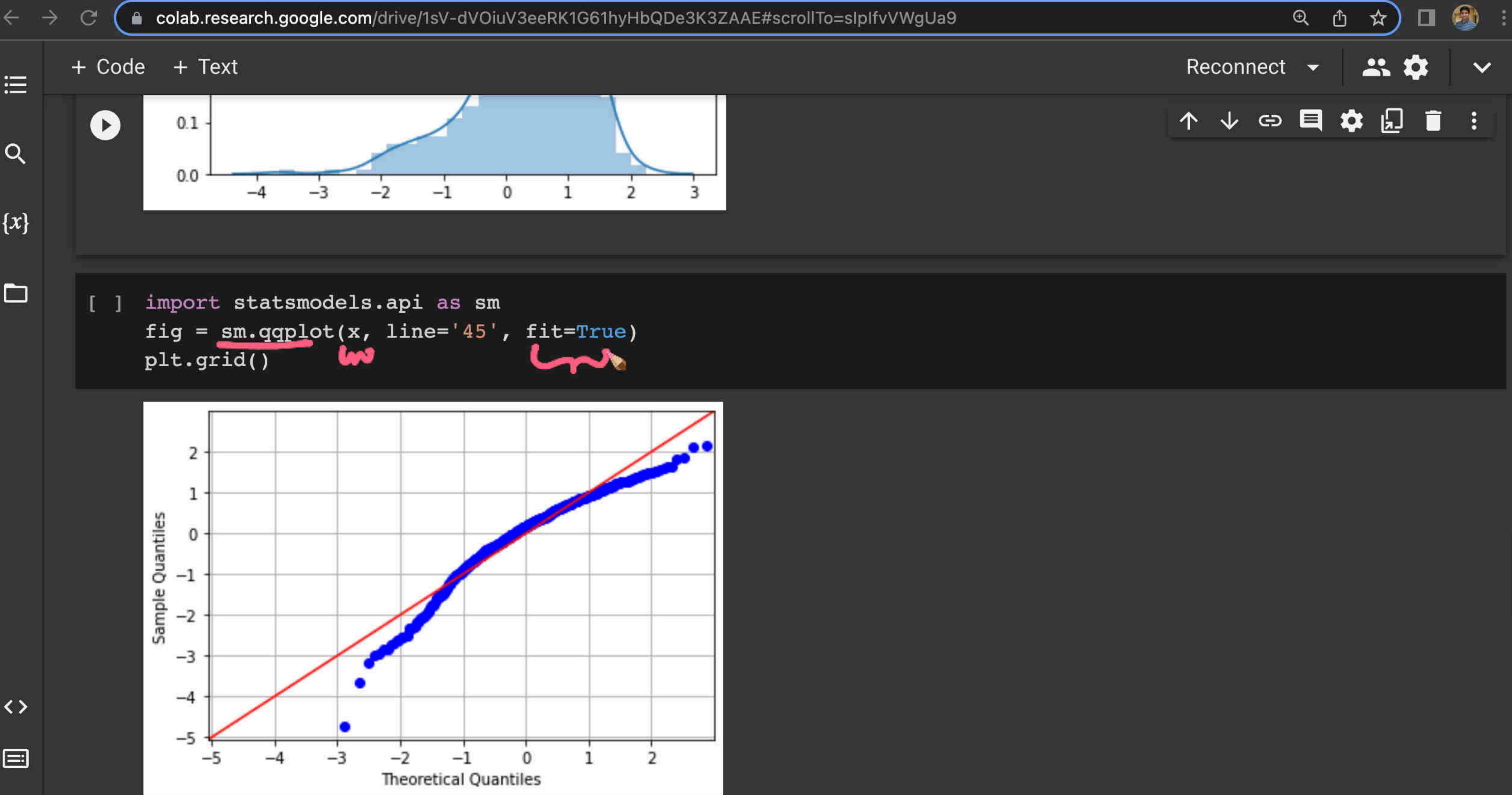


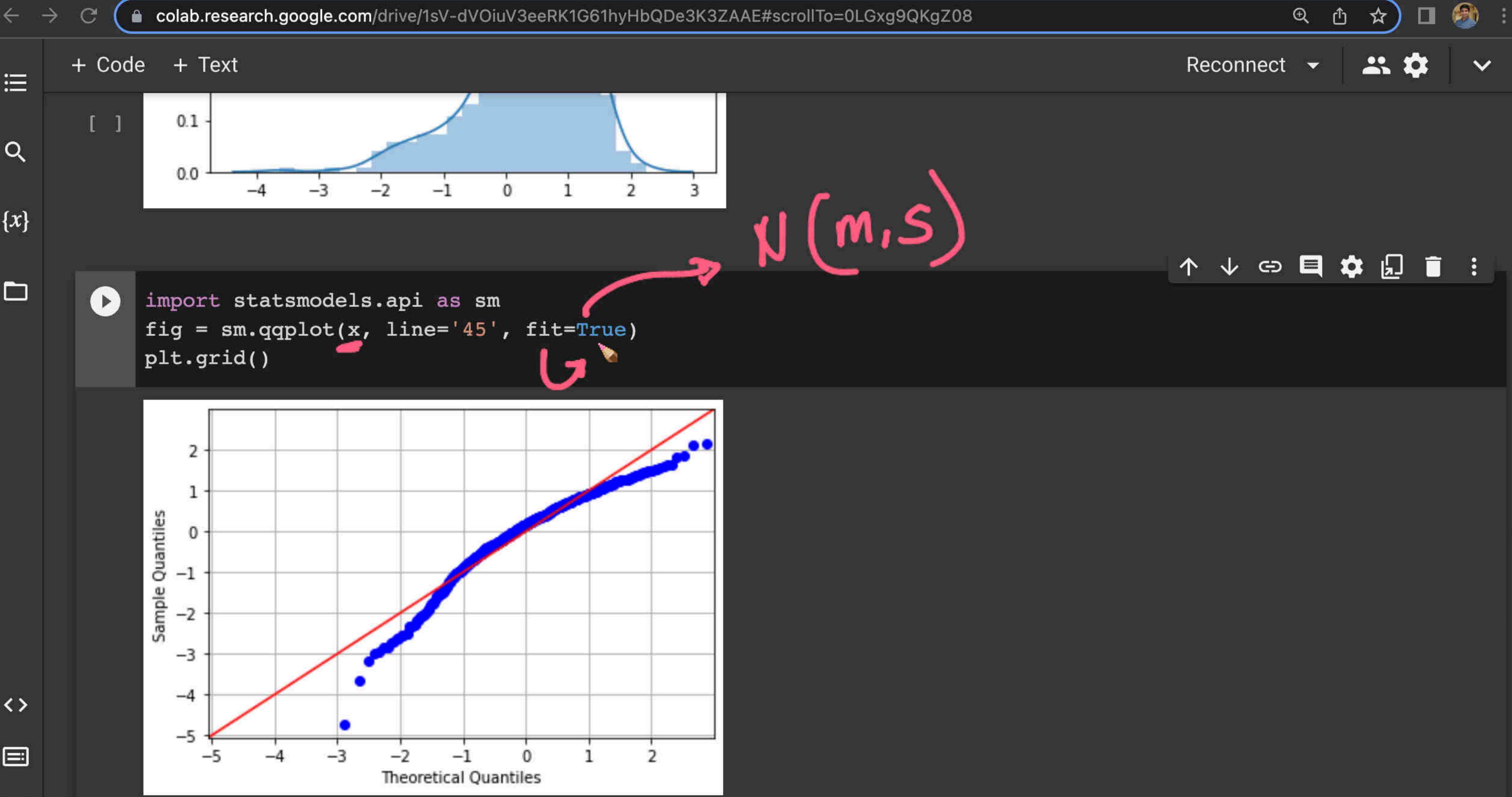
scipy.stats



```
import statsmodels.api as sm
fig = sm.qqplot(x, line='45', fit=True)
plt.grid()
```







QQ plot and Transf

statsmodels.graphics

scipy.stats.probplot

Power transform - V

scipy.stats.boxcox - X

Normal distribution X

colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=OLGxg9QKgZ08

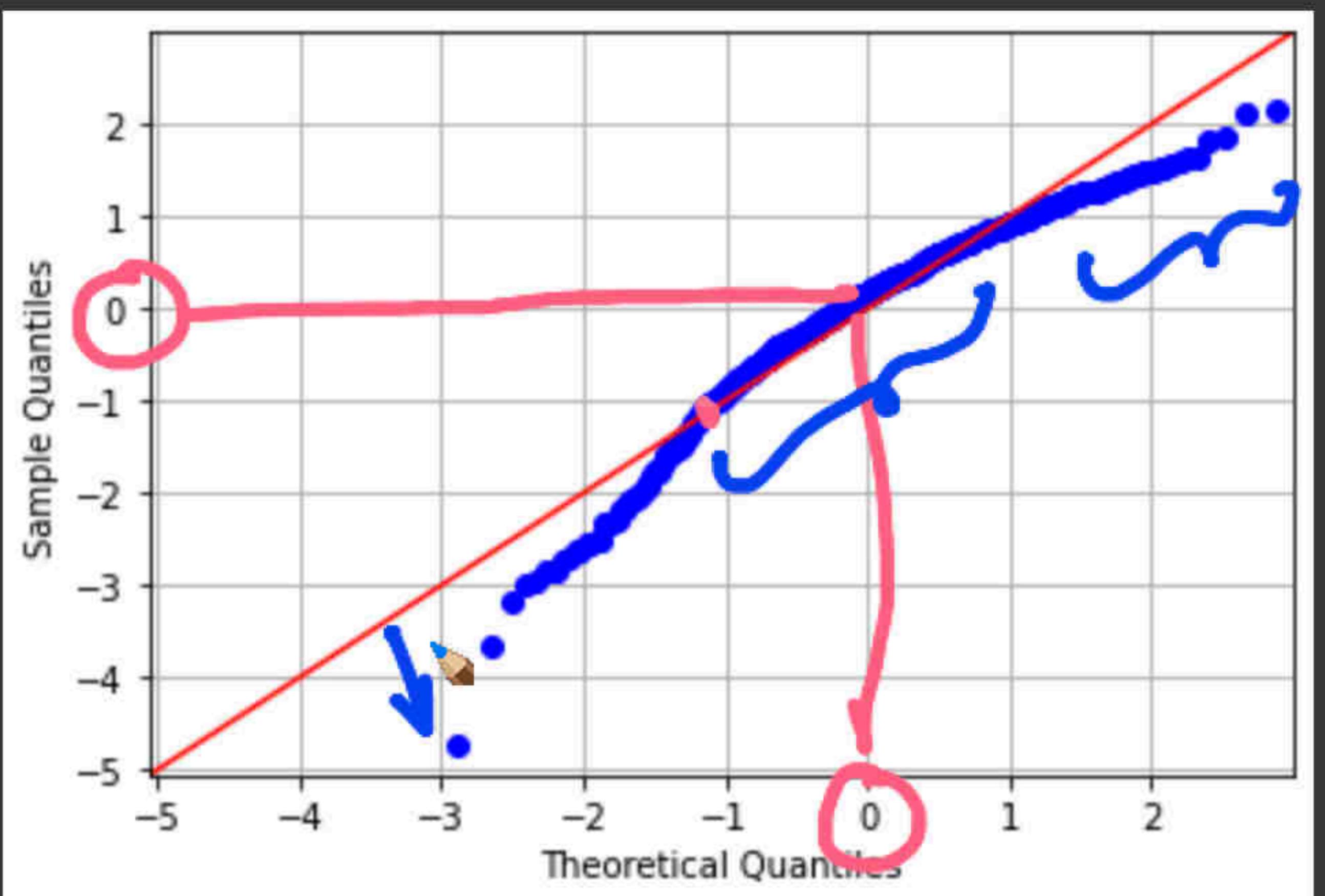


+ Code + Text

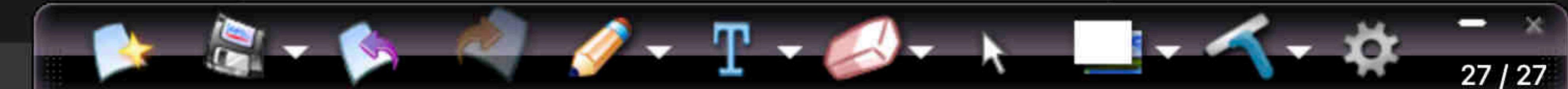
Reconnect



```
import statsmodels.api as sm
fig = sm.qqplot(x, line='45', fit=True)
plt.grid()
```

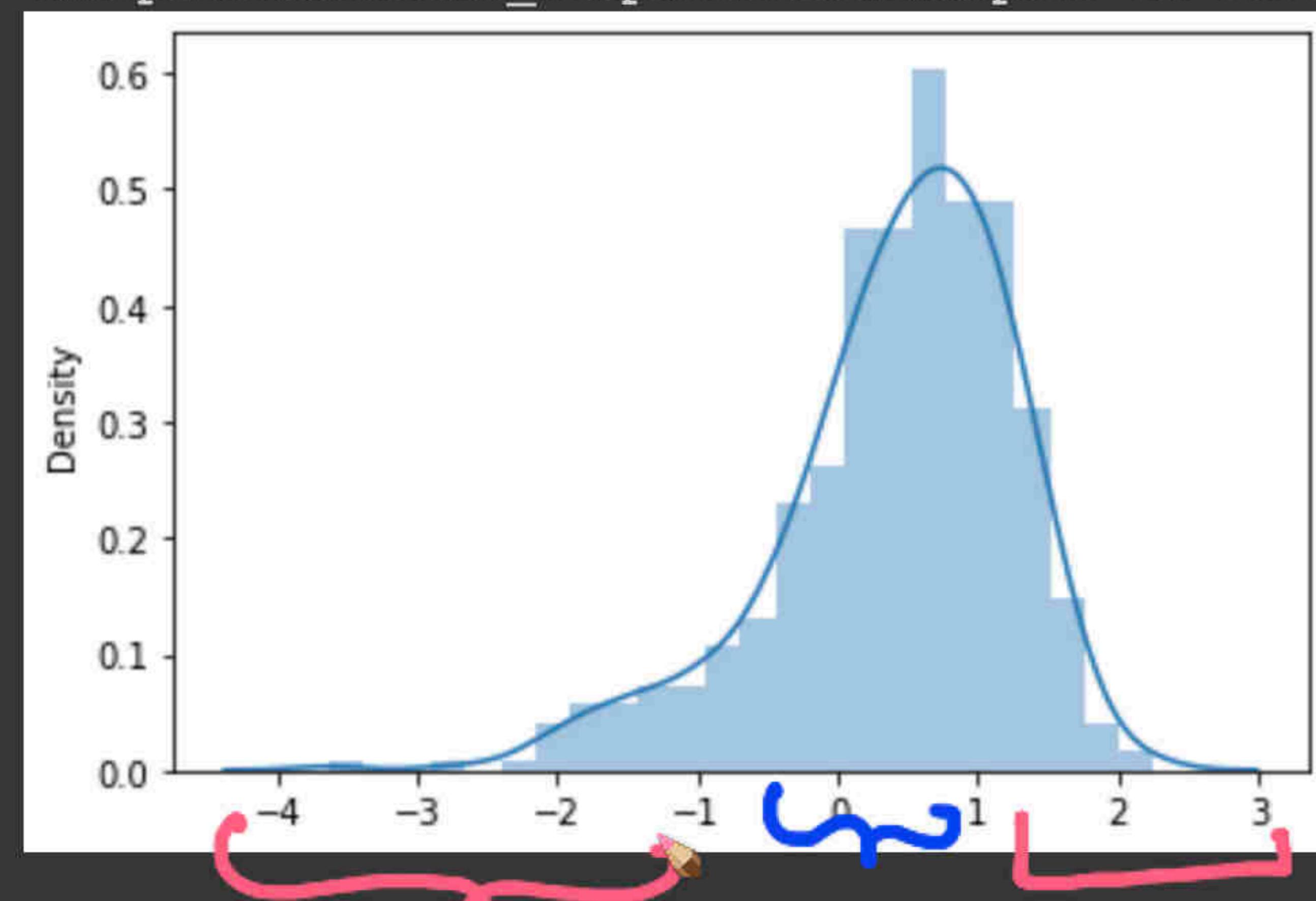


```
[ ] # Let us compare CDFs also
mu = np.mean(x)
s = np.std(x)
print(mu,s)
```



+ Code + Tex

## Reconnect



QQ plot and Transf...

statsmodels.graphics

scipy.stats.probplot

Power transform -

scipy.stats.boxcox -

Normal distribution

+

▼

colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=OLGxg9QKgZ08

🔍 📁 ⚡

⋮

+ Code + Text

Reconnect

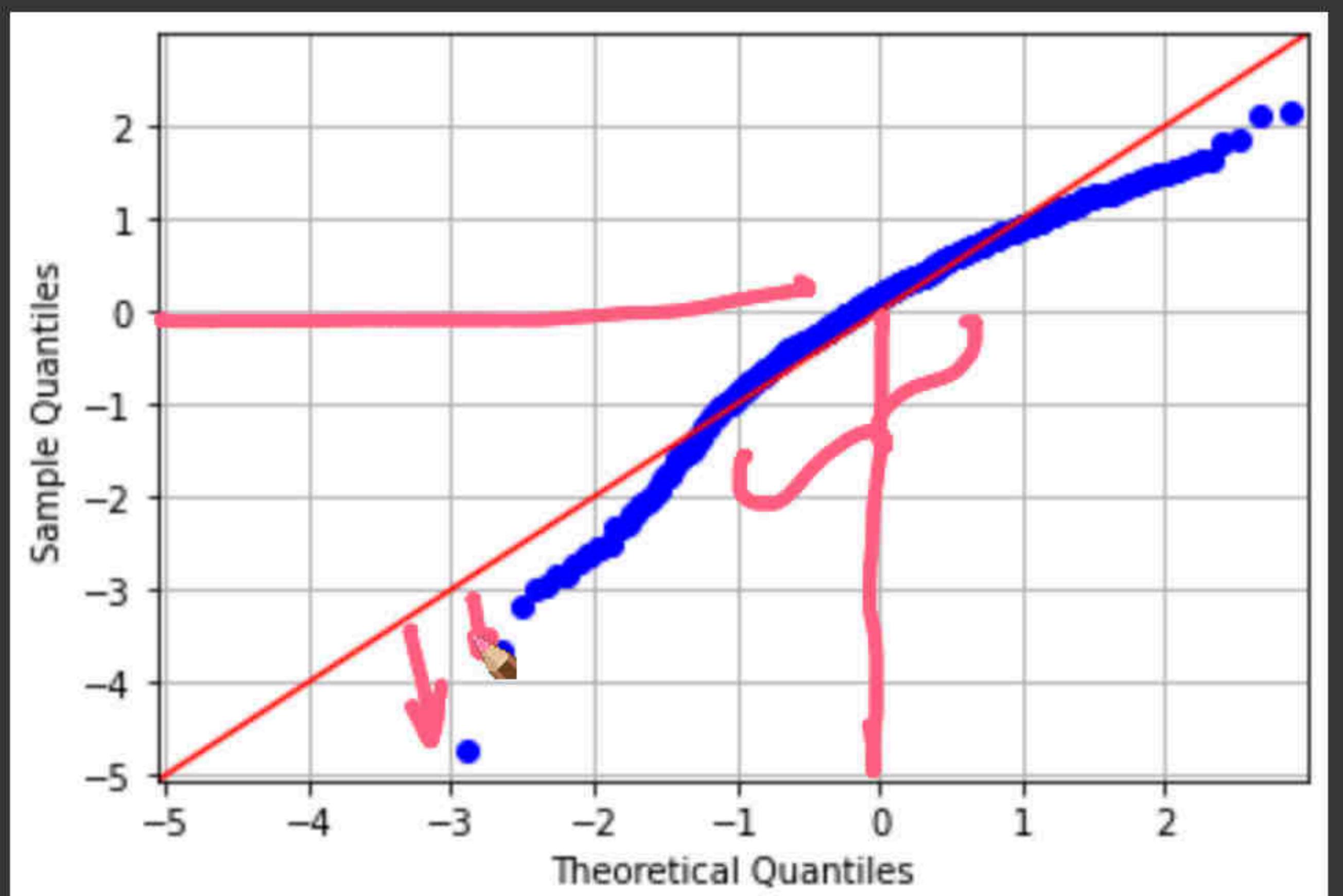
👤⚙️

▼

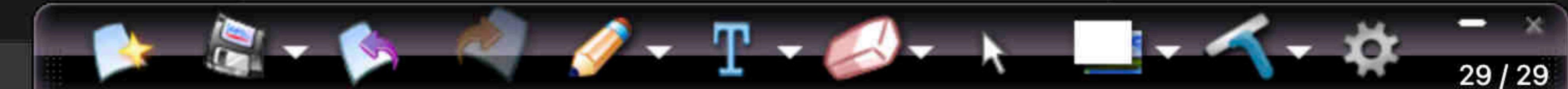
↑ ↓ ↻ ⏷ ⚙️ 🗑️ ⏷



```
import statsmodels.api as sm
fig = sm.qqplot(x, line='45', fit=True)
plt.grid()
```



```
[ ] # Let us compare CDFs also
mu = np.mean(x)
s = np.std(x)
print(mu,s)
```



29 / 29

×

optional  
Task:

{ Code  
scipy...  
qqplot

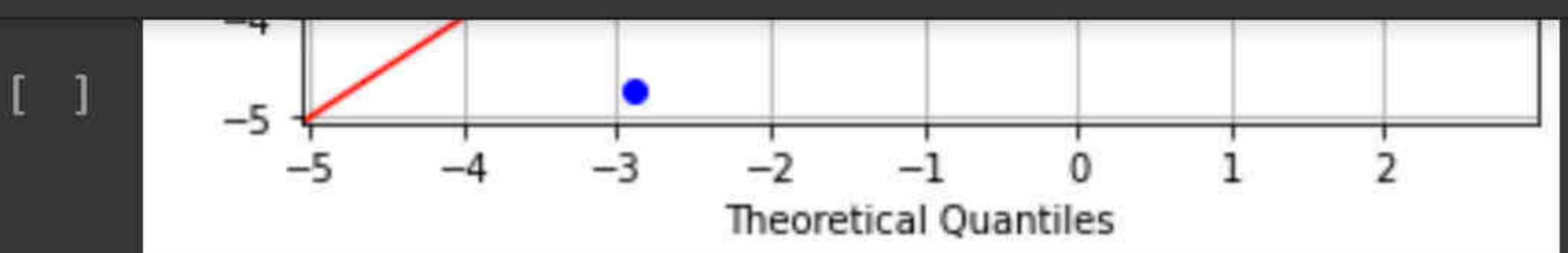
$x_1 \dots x_n$

↓  
WAF to plot QQ-Plot

scipy.stats  
lineplot

+ Code + Text

Reconnect



{x}



```
# Let us compare CDFs also
mu = np.mean(x)
s = np.std(x)
print(mu, s)
```

3 → M, S

0.40812679178 0.8571215209422999

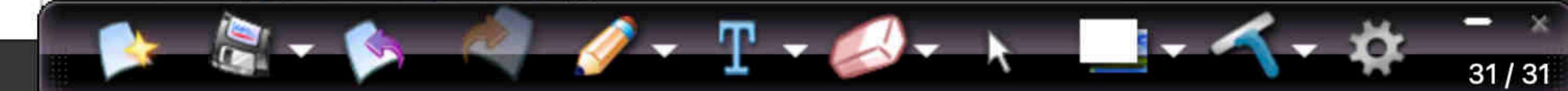
```
[ ] # normally distributed data with mean=mu, std-dev=s
y = stats.norm.rvs(loc=mu, scale=s, size=500)
```

↔ N

$y_1, y_2, \dots, y_{500}$

$\overbrace{N(M, S)}$

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



+ Code + Text

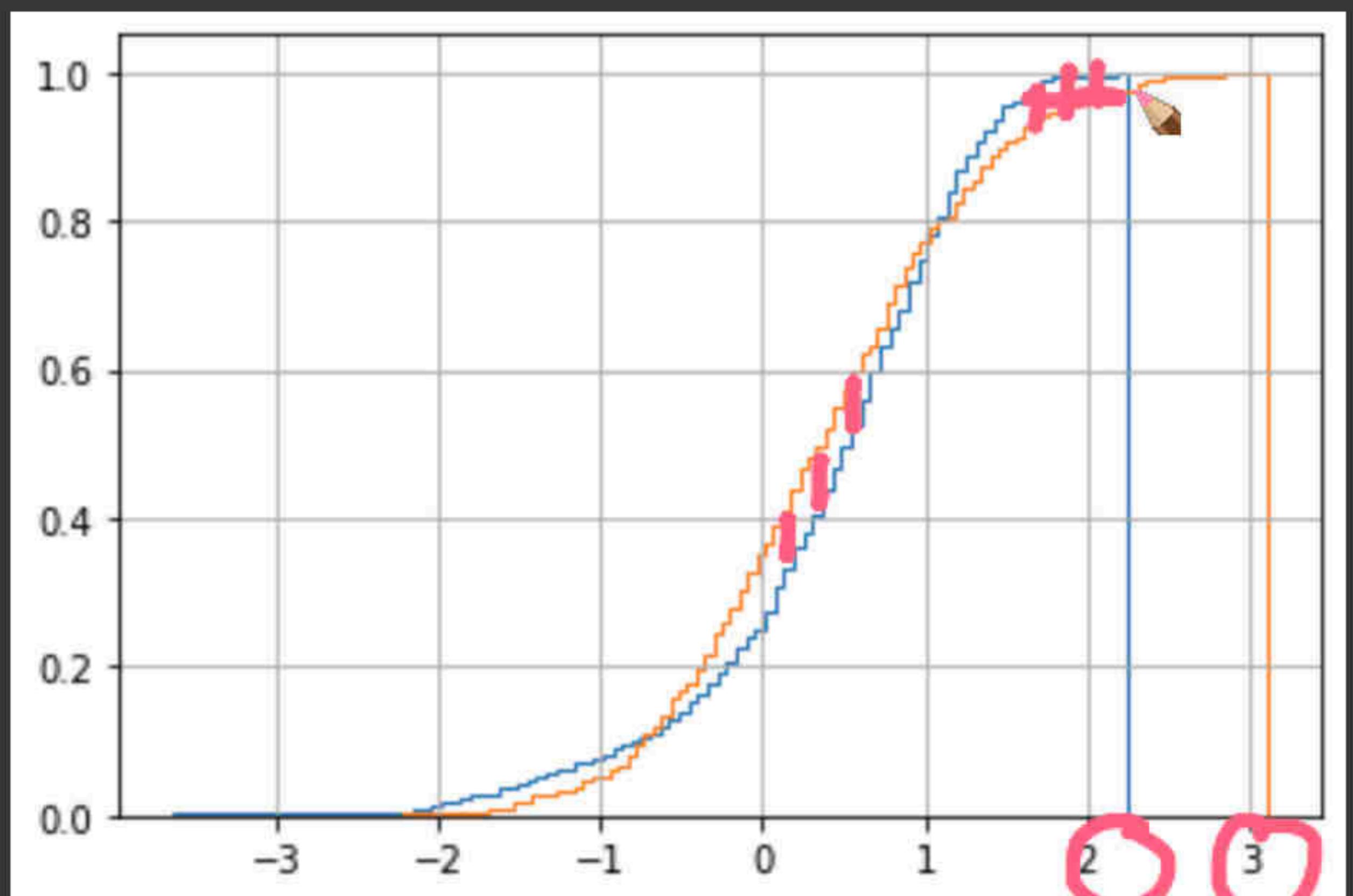
Reconnect



```
# normally distributed data with mean=mu, std-dev=s
[ ] y = stats.norm.rvs(loc=mu, scale=s, size=500)
```



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



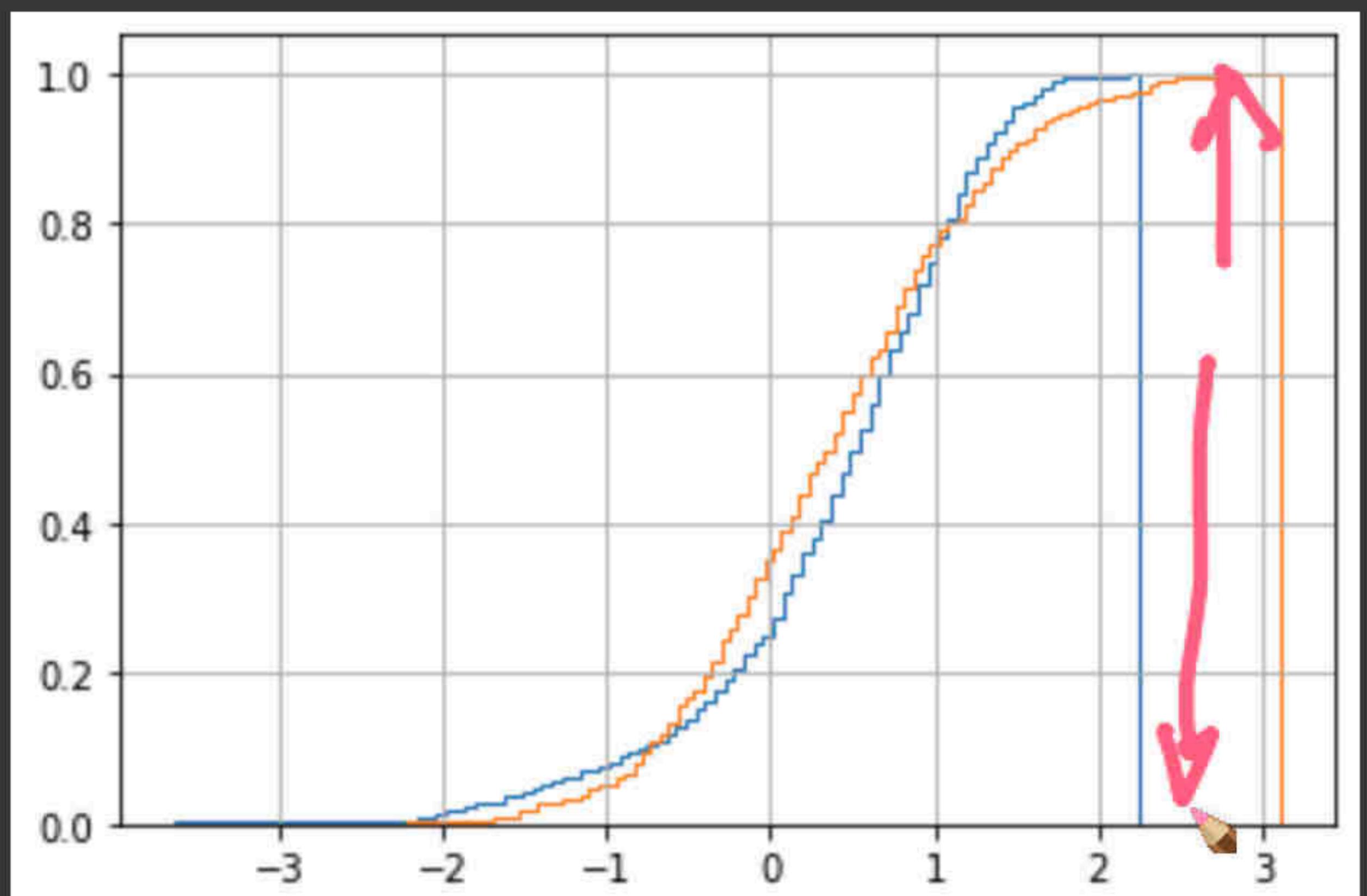
## ► QQ Plot: Gaussian vs Gaussian

+ Code + Text

```
# normally distributed data with mean=mu, std-dev=s
[ ] y = stats.norm.rvs(loc=mu, scale=s, size=500)
```



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

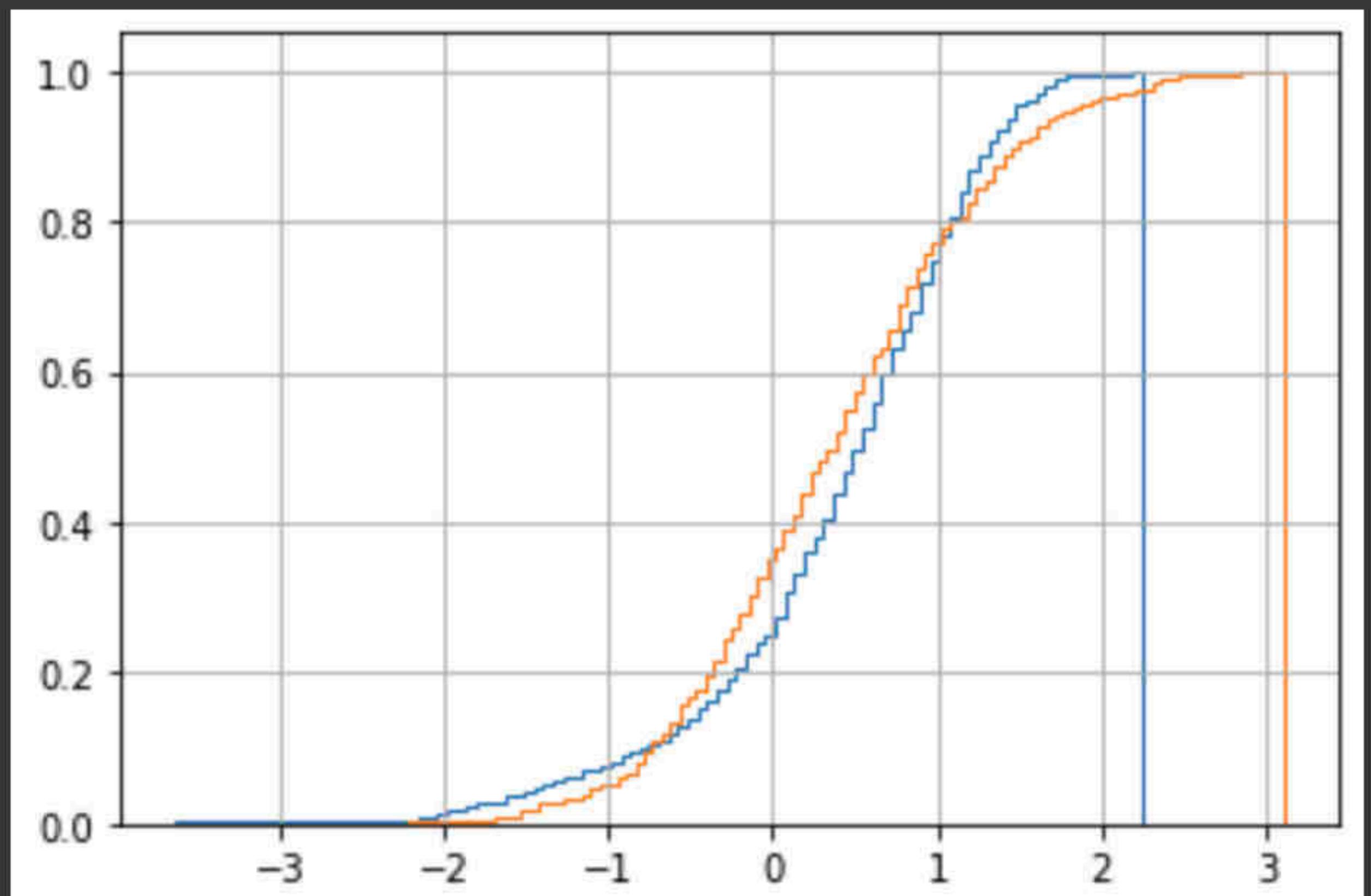


## ► QQ Plot: Gaussian vs Gaussian

```
# normally distributed data with mean=mu, std-dev=s
[ ] y = stats.norm.rvs(loc=mu, scale=s, size=500)
```



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



## ► QQ Plot: Gaussian vs Gaussian

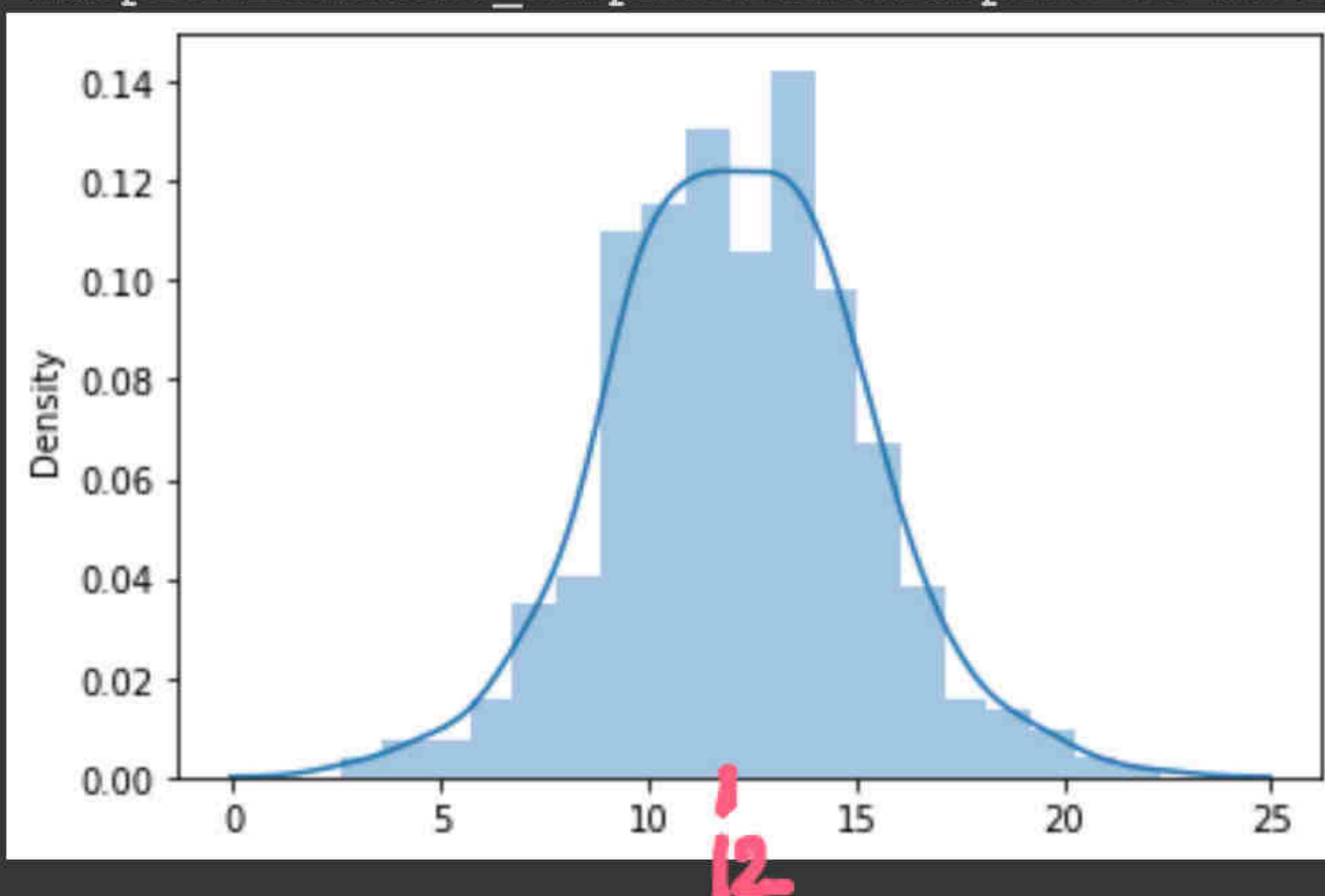
+ Code + Text

Reconnect

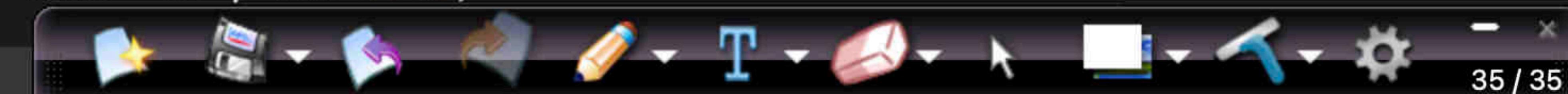
[ ] 7.13412165, 10.26586752, 14.90528691, 13.65151123, 11.46129682,  
14.28189362, 9.25854953, 10.40331241, 14.82192117, 9.85303781])

{x} [ ] sns.distplot(x)

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecate
  warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f146d5e9690>
```



sized  
500 sample

[ ] fig = sm.qqplot(x, line='45', fit=True)  
plt.grid()

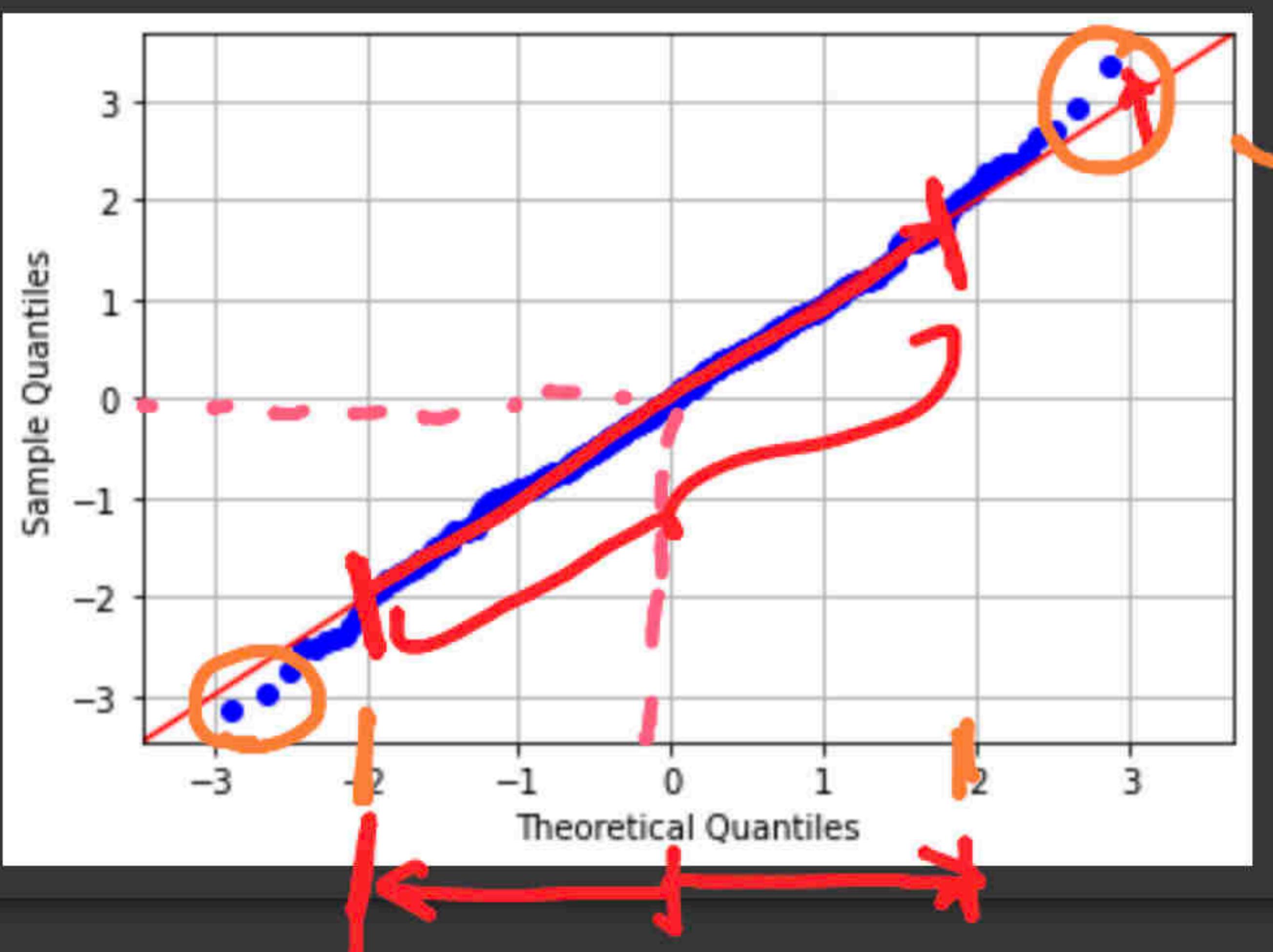
+ Code + Text

Reconnect

User Settings

Up Down Left Right Chat Gear Copy Delete More

```
fig = sm.qqplot(x, line='45', fit=True)  
plt.grid()
```



500 obs } → Scarcity  
of data

5M obs

QQ Plot between Non Gaussian and Non-gaussian

[ ] ↳ 4 cells hidden

+ Code + Text

Reconnect

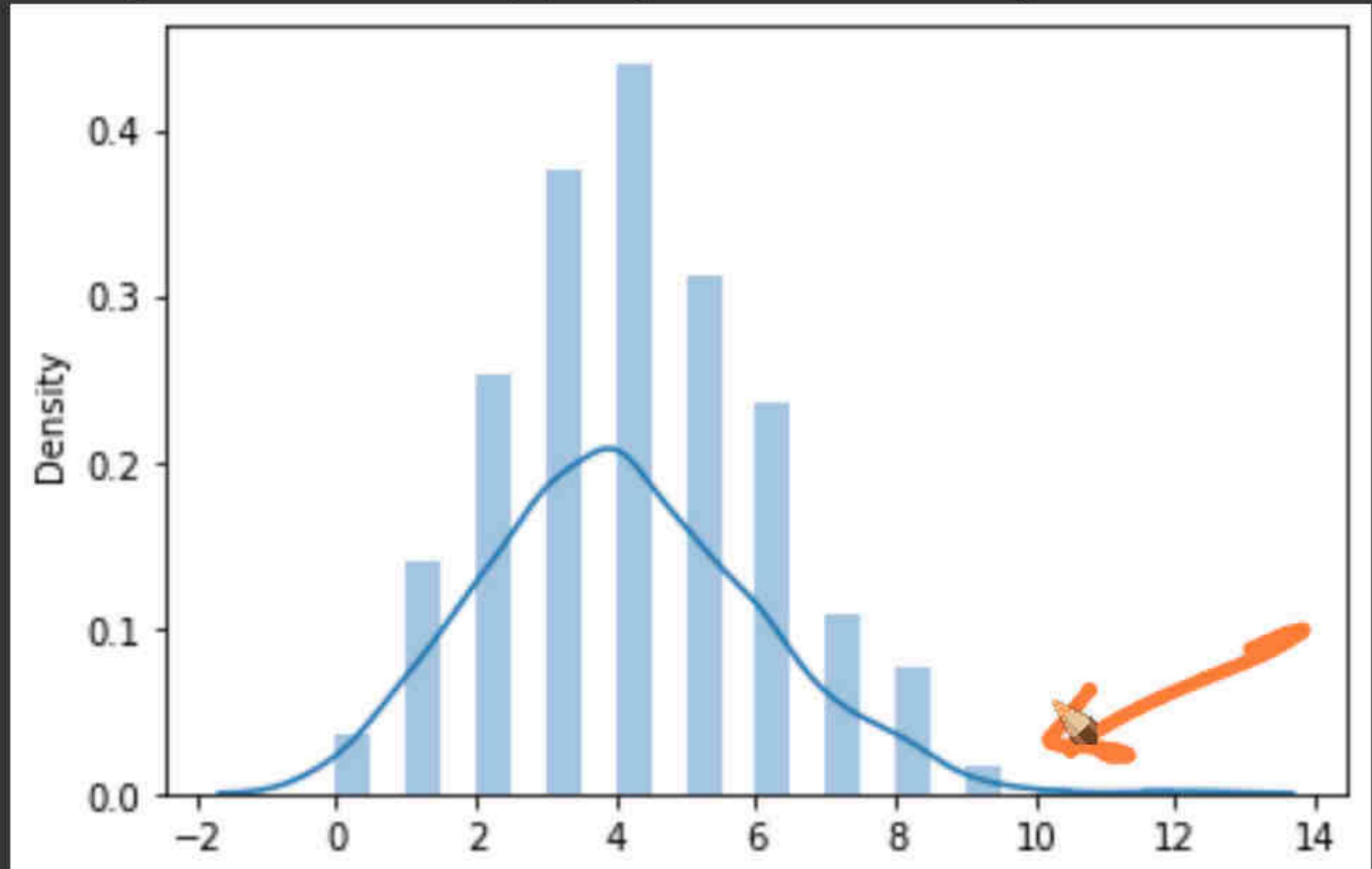
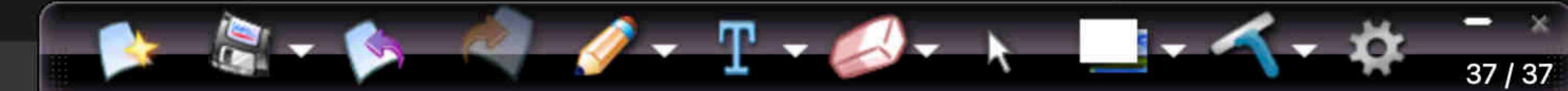


4, 6, 8, 3, 3, 2, 2])



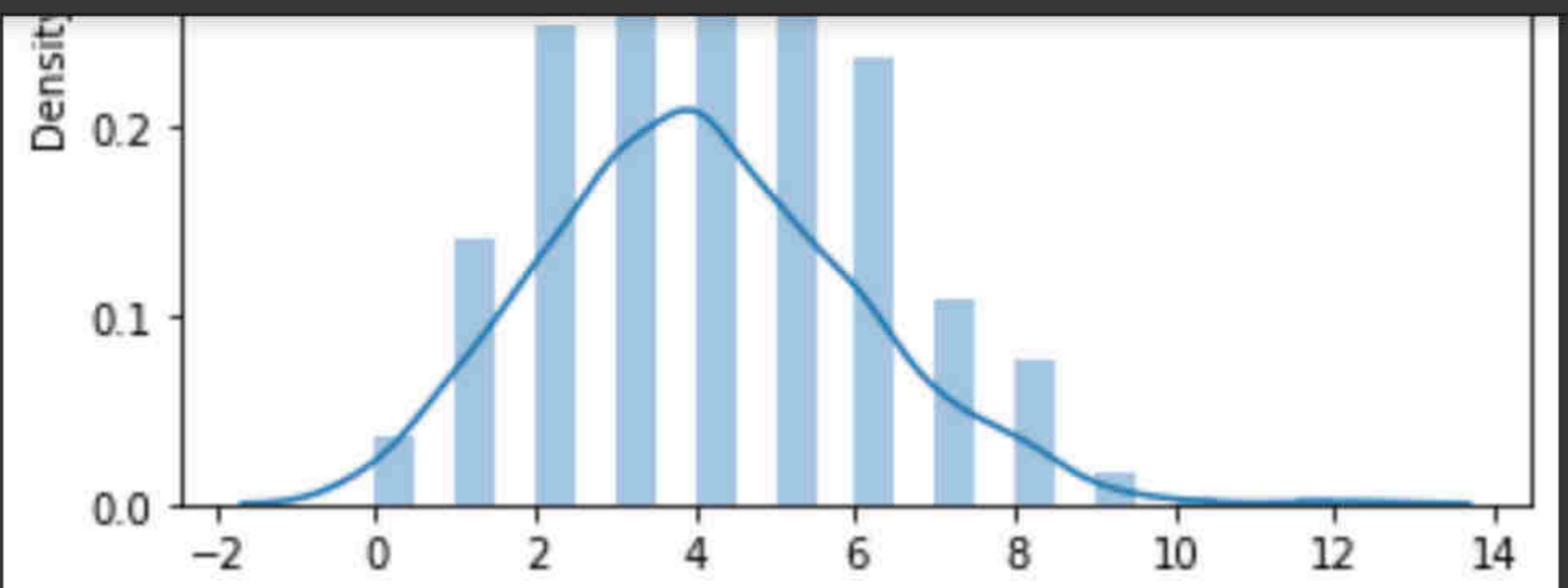
[ ] sns.distplot(x)

{x}

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecate  
warnings.warn(msg, FutureWarning)  
<matplotlib.axes.\_subplots.AxesSubplot at 0x7f146d913f90>[ ] import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
plt.grid()

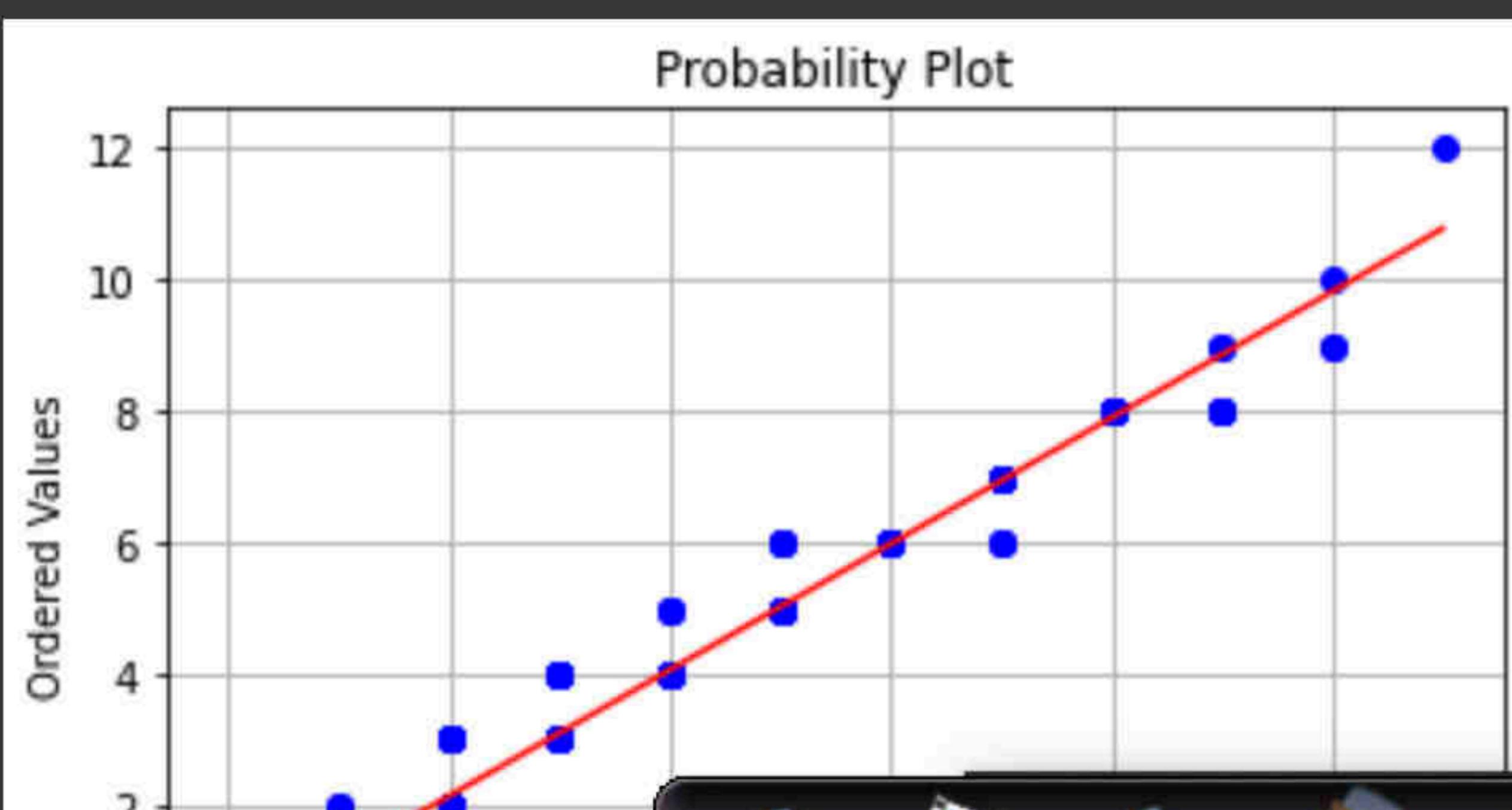
+ Code + Text

Reconnect



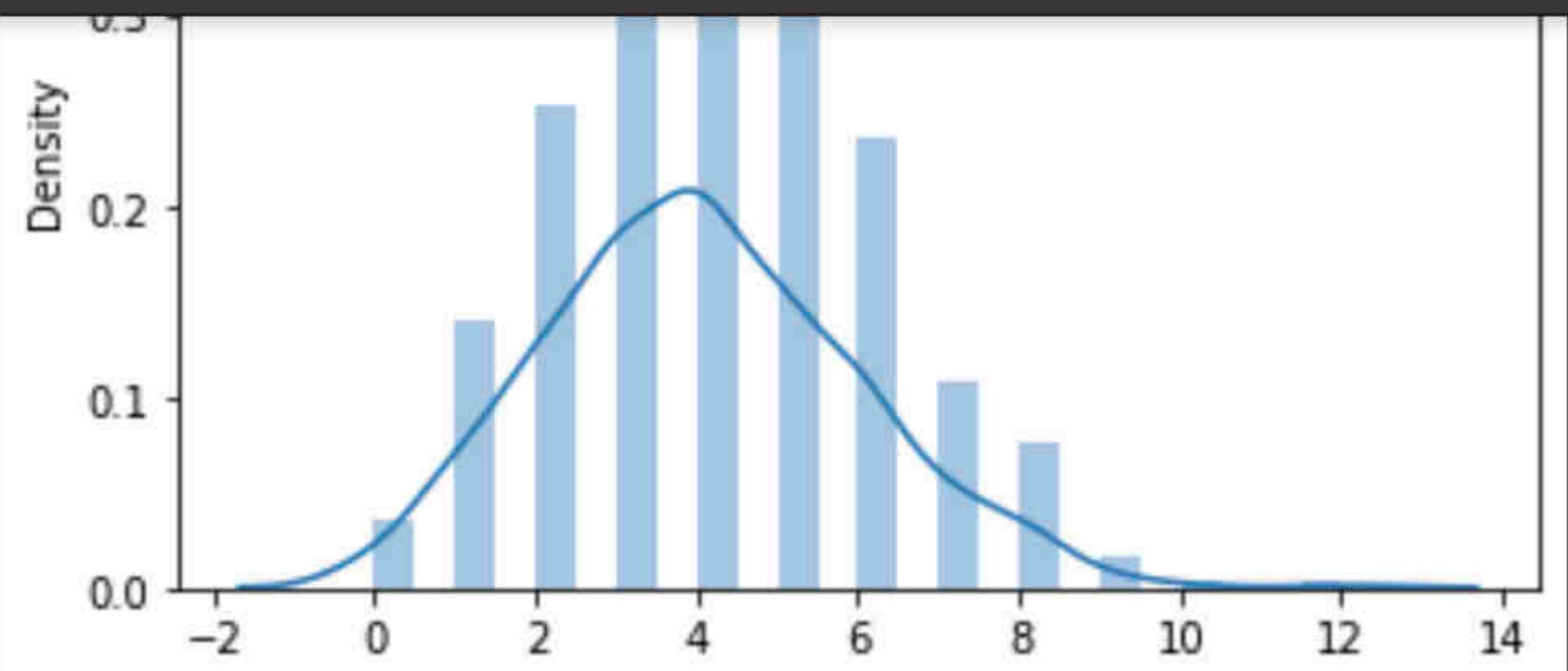
# Poisson ( $\lambda$ )

```
import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
plt.grid()  
plt.show()
```

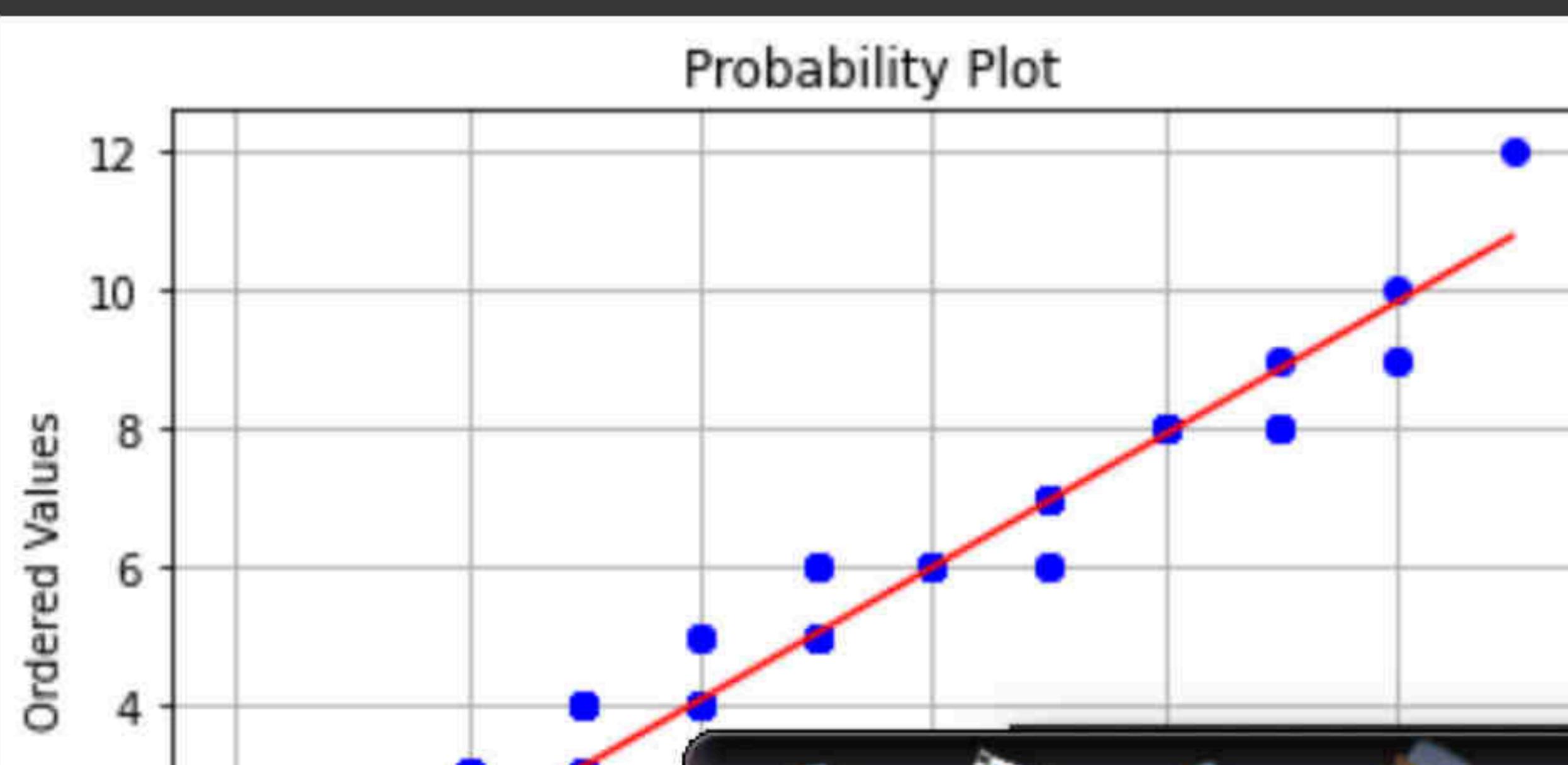


+ Code + Text

Reconnect



```
import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
{  
    plt.grid()  
    plt.show()
```



QQ plot and Transf...

statsmodels.graphics

scipy.stats.probplot

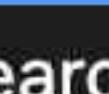
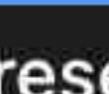
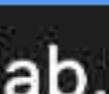
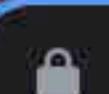
Power transform -

scipy.stats.boxcox -

Normal distribution

+

▼



colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=vS6zL7-ei4yz

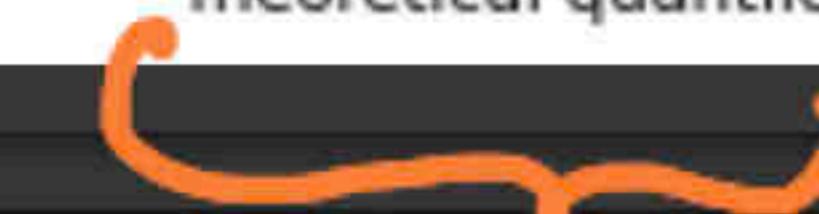
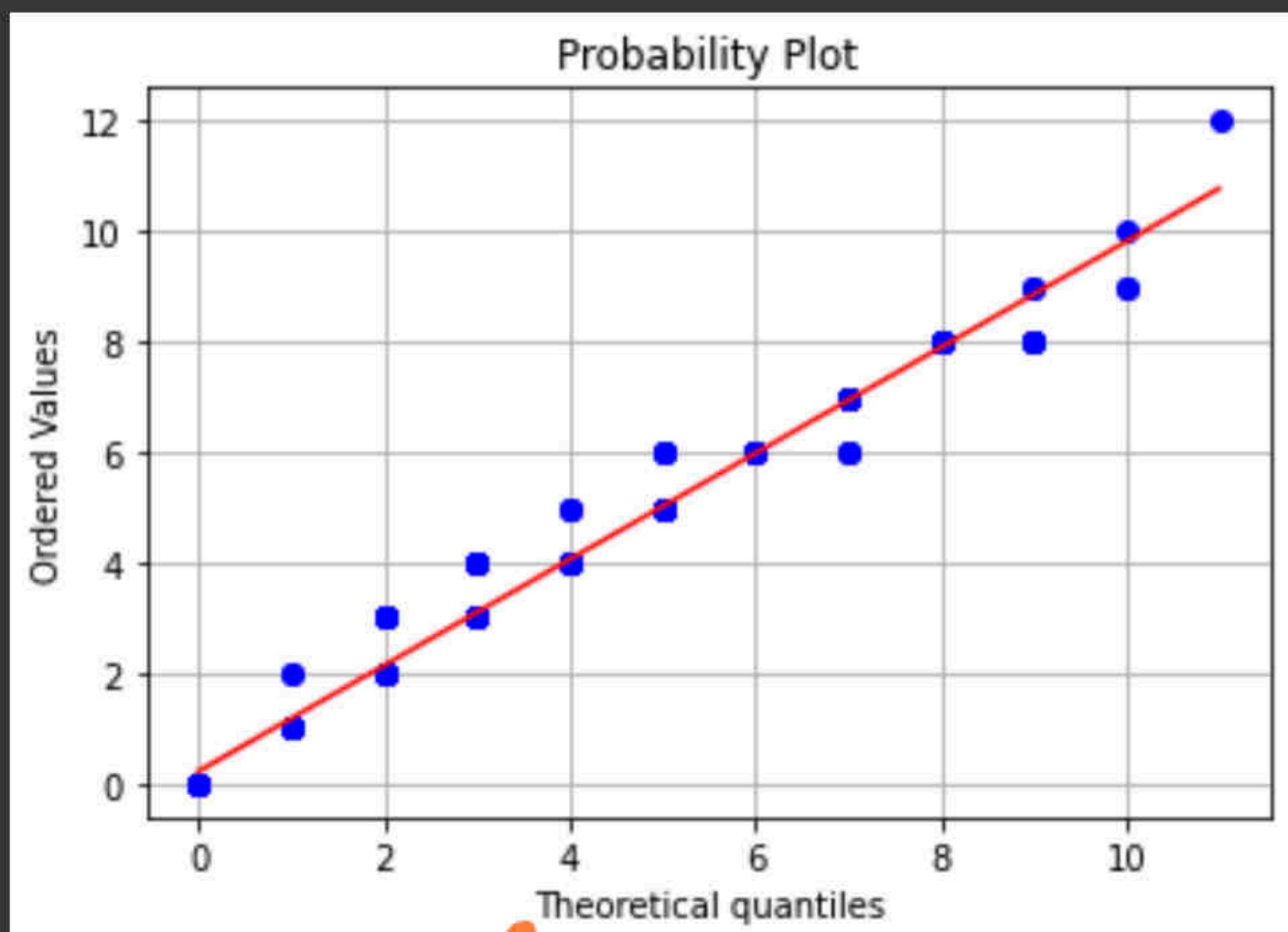


+ Code + Text

Reconnect



```
import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
plt.grid()  
plt.show()
```



QQ plot and Transfor

statsmodels.graphics

scipy.stats.probplot

Power transform -

scipy.stats.boxcox -

Normal distribution

+

▼



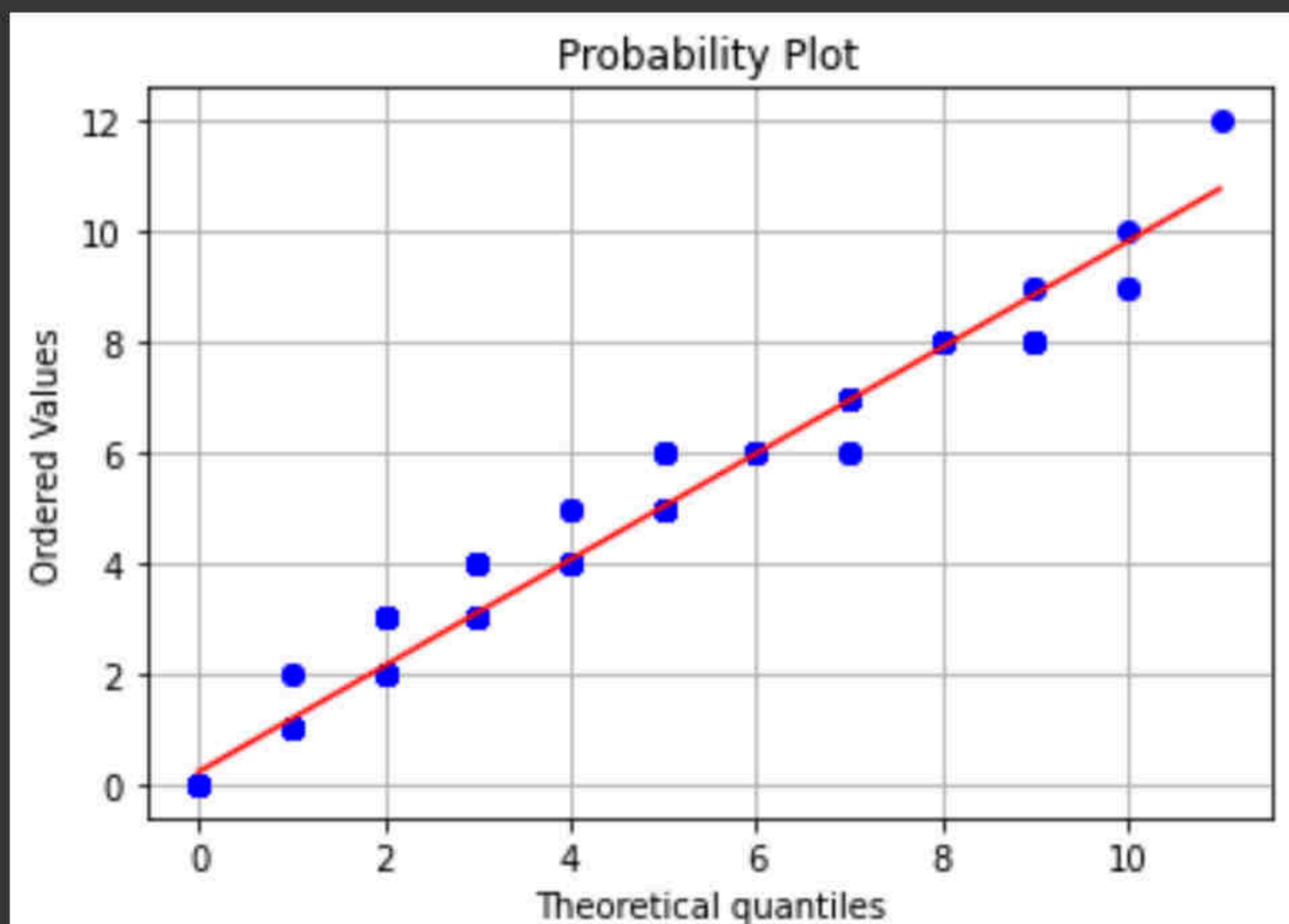
+ Code



+ Text



```
import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
plt.grid()  
plt.show()
```



[ ]

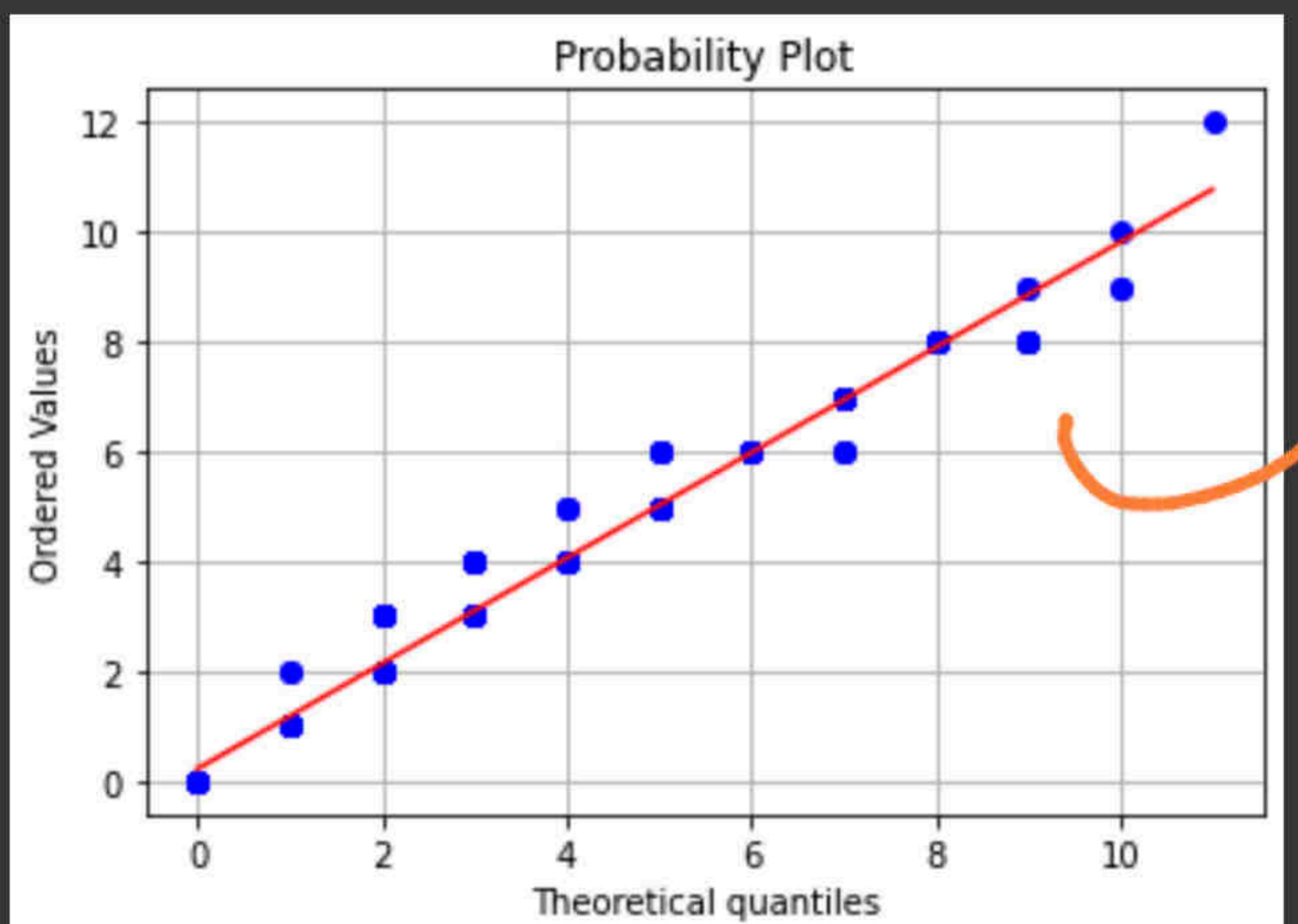


+ Code + Text

Reconnect



```
import scipy.stats as stats  
stats.probplot(x, dist='poisson', sparams=(4,), plot=plt)  
plt.grid()  
plt.show()
```



$X \sim \text{Poisson} (\lambda=4)$



QQ plot and Transfor

statsmodels.graphics

scipy.stats.probplot

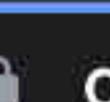
Power transform -

scipy.stats.boxcox -

Normal distribution

+

▼



colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=zL\_38krDiMDk

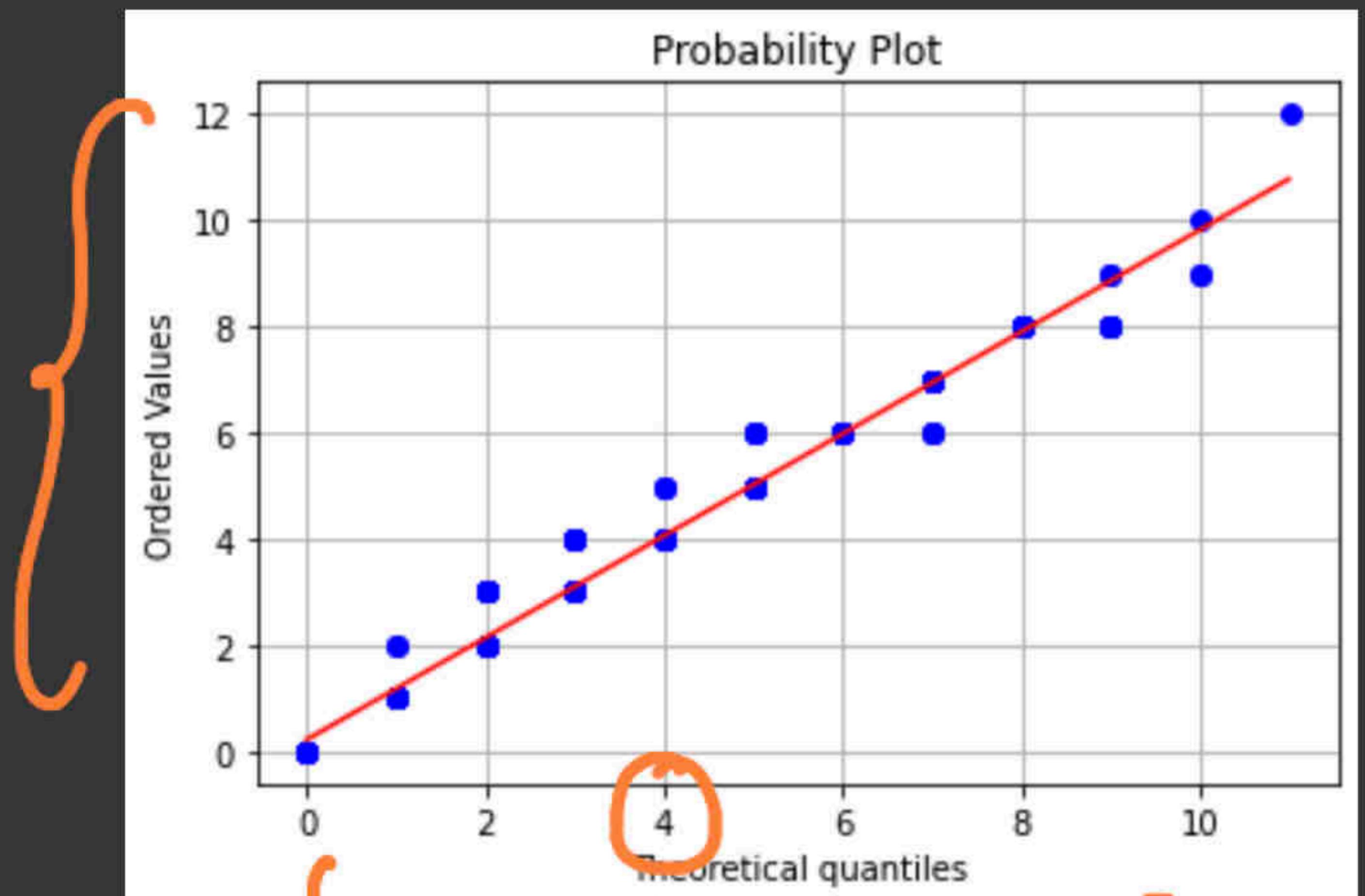


+ Code + Text

Reconnect



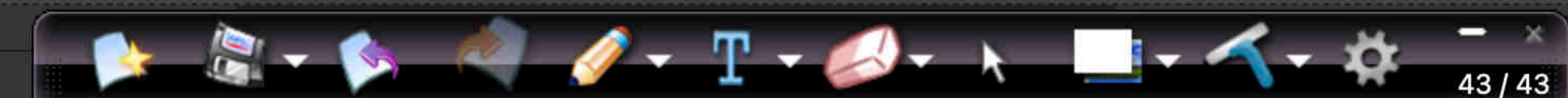
plt.show()

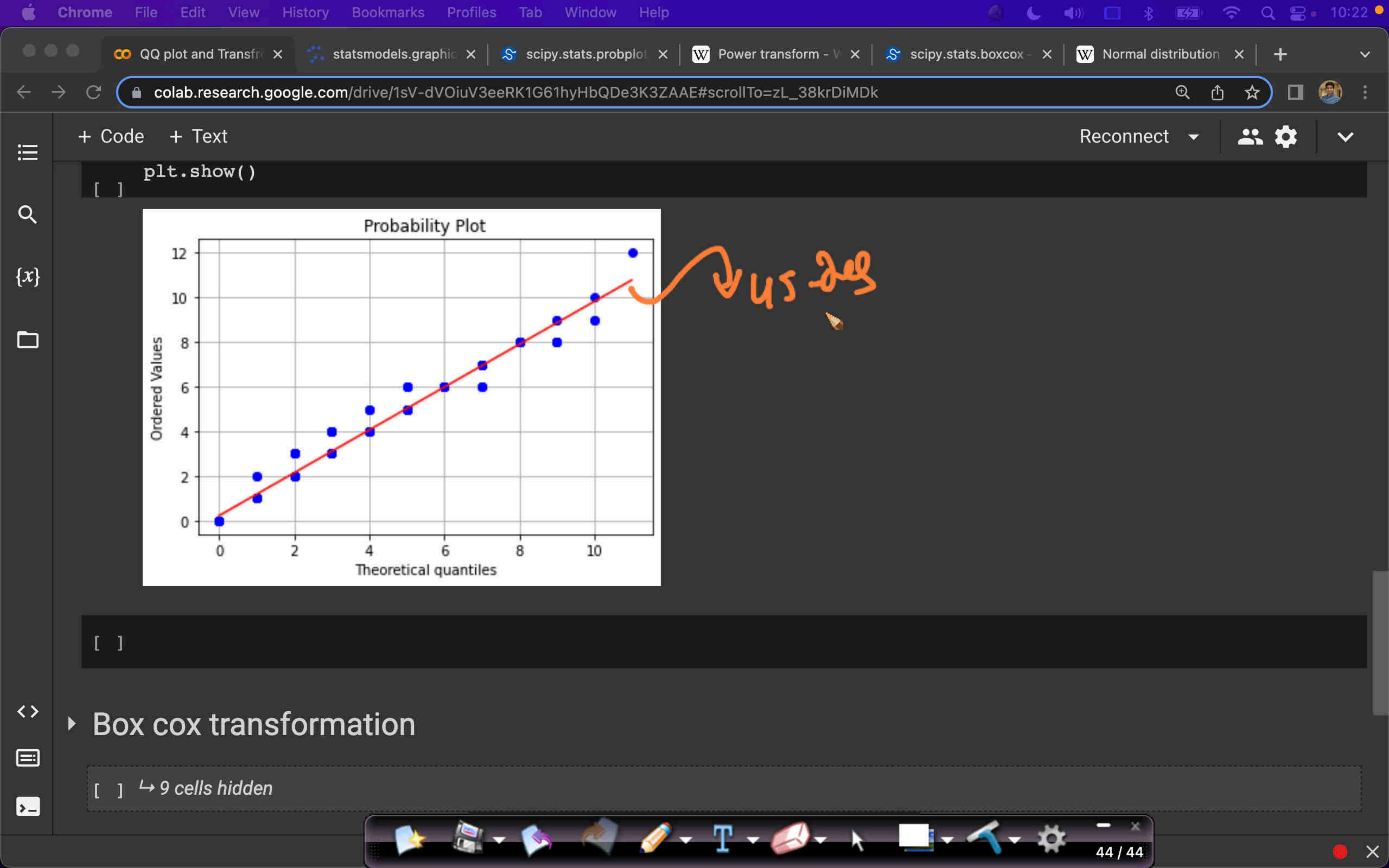


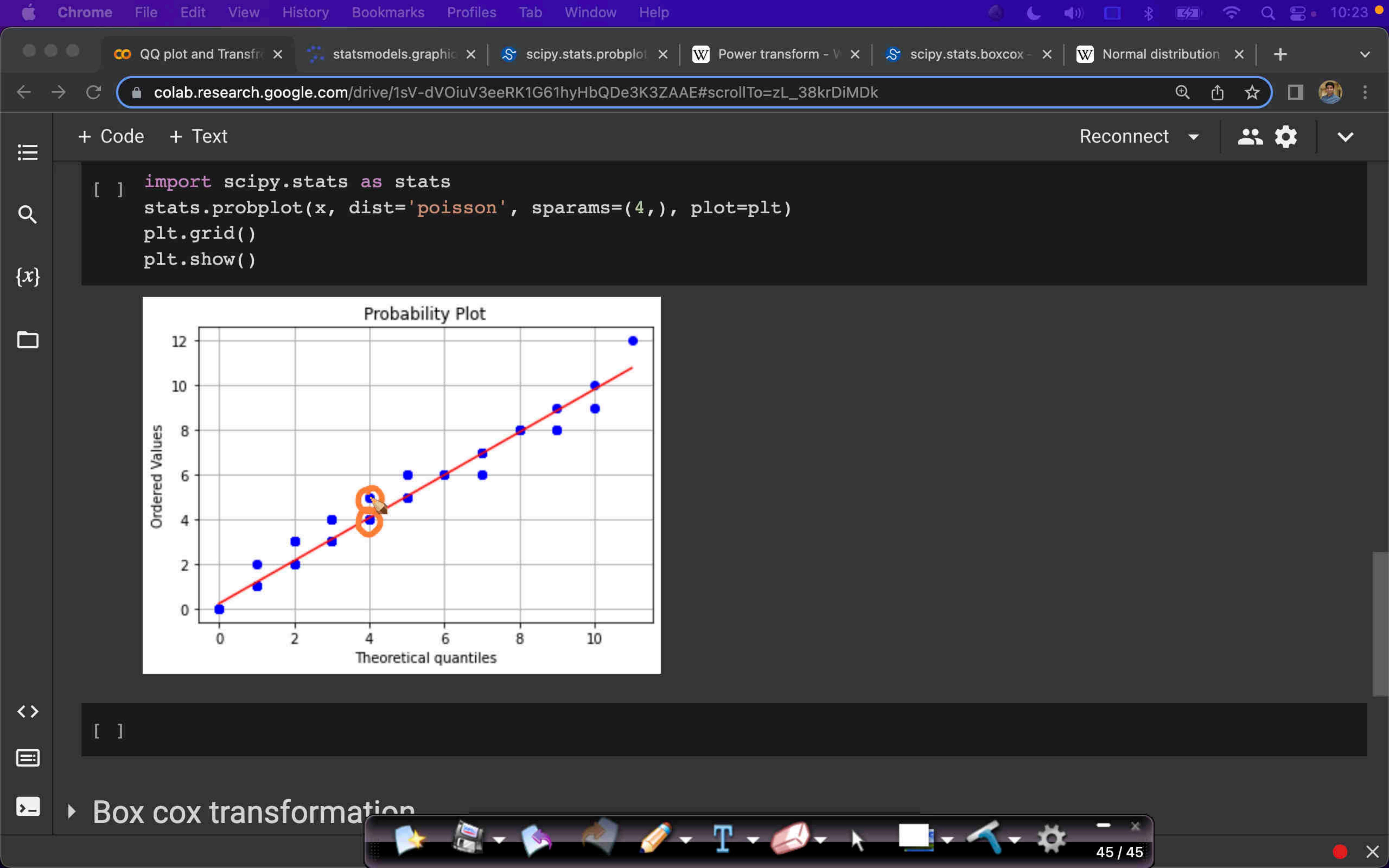
[ ]

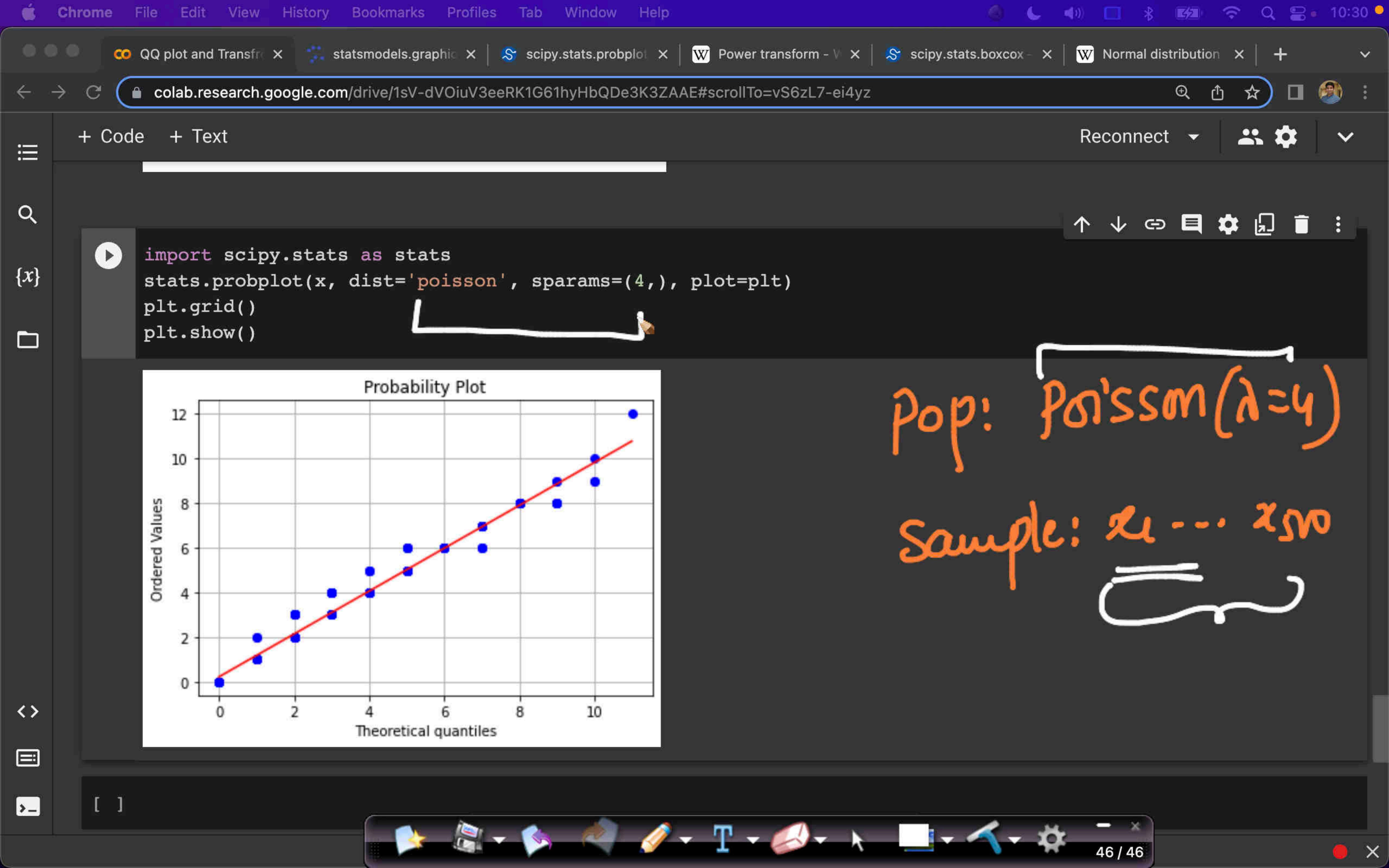
## Box cox transformation

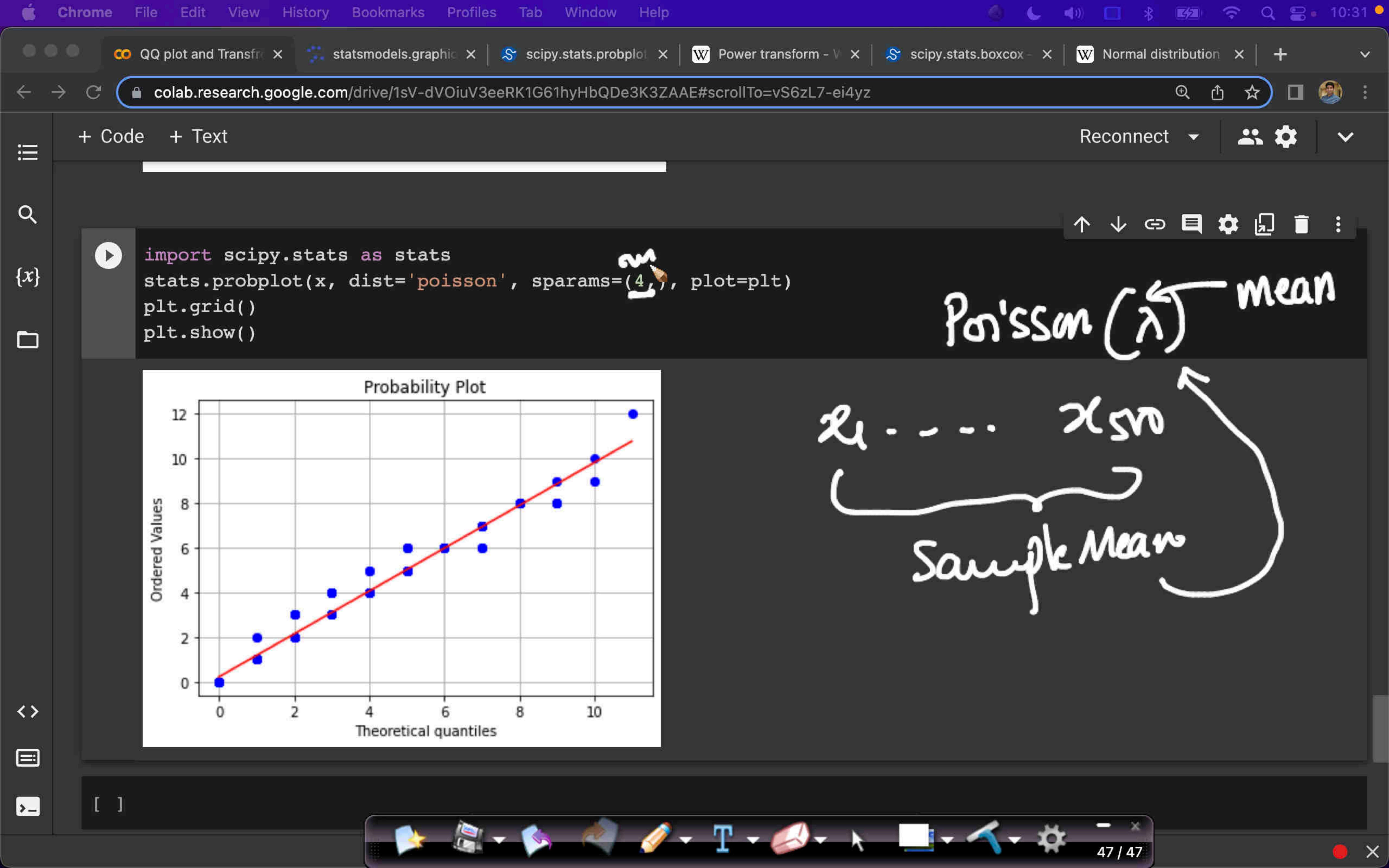
[ ] ↳ 9 cells hidden











Task:

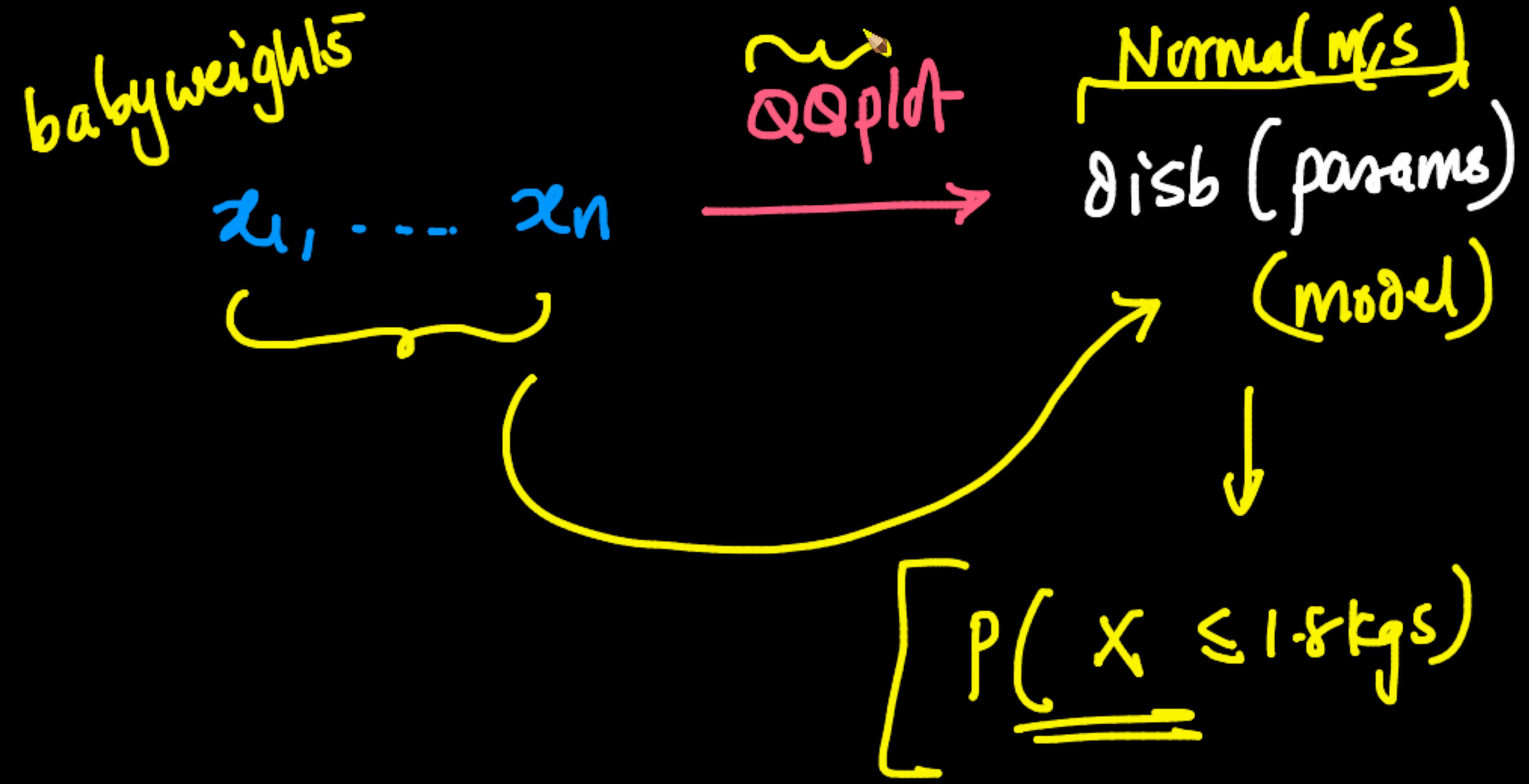
$x_1 \dots x_{500}$

`Normal()`

QQ-plot

scipy.stats  
np  
!

line plot



# Transformations

Task:

$$\begin{matrix} X \\ x_1 \ x_2 \dots x_i \ x_n \end{matrix} \xrightarrow{f()} \begin{matrix} X' \\ x'_1, x'_2, \dots x'_i, x'_n \end{matrix}$$

Not-Gaussian

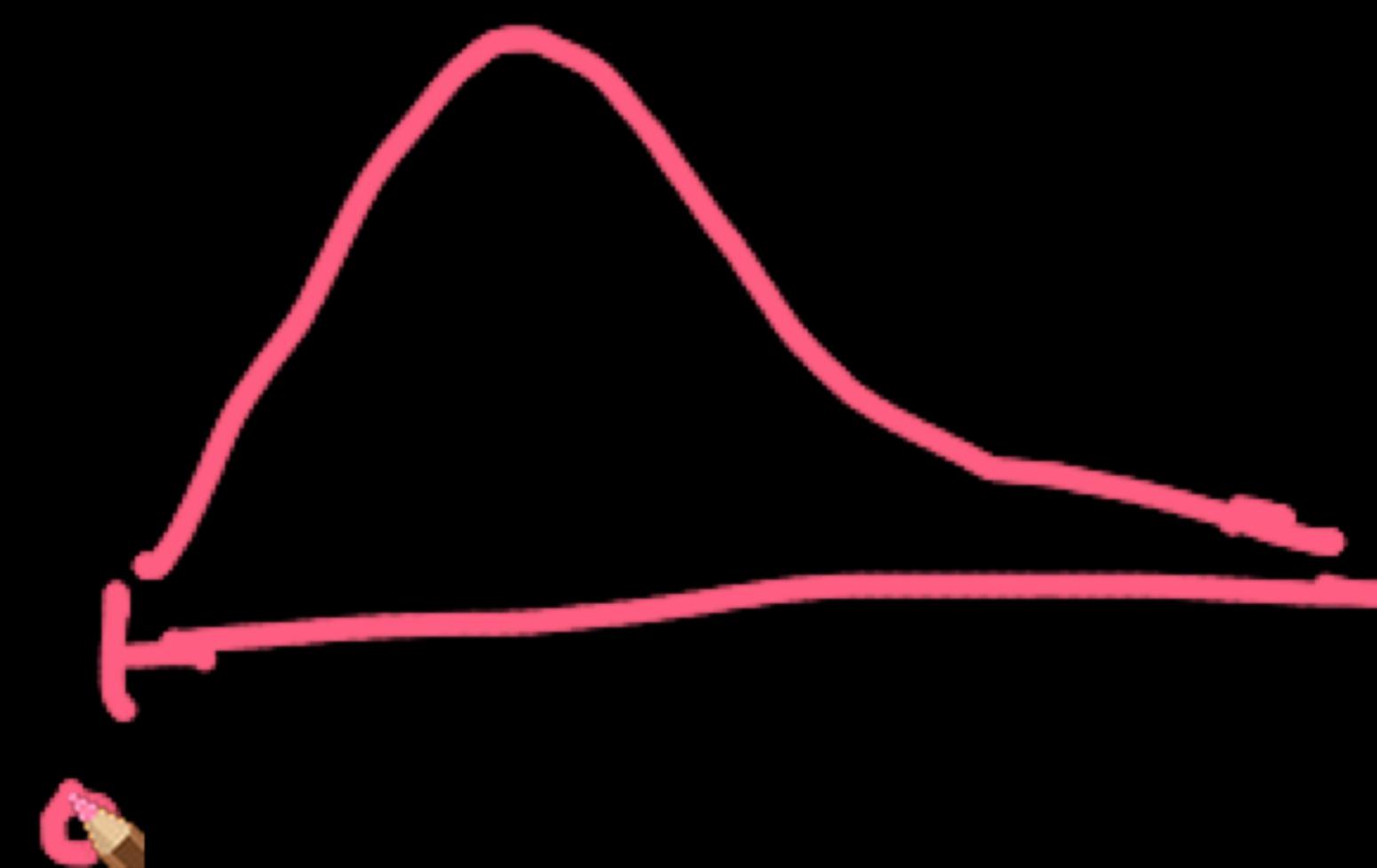
[ML: Th-proof]

→ tons of Math  
easy to explain

(Q) guess a transformation:

$$x_1, x_2, \dots, x_i, \dots, x_n \xrightarrow{\ln()} x'_1, x'_2, \dots, x'_i, \dots, x'_n$$

log-normal



$$\underline{\log_e} = \underline{\ln}$$

$f()$  →  $\log()$  ✓  
↳  $\text{sqrt}()$  → not popular

⋮

$f(x_i) = \underbrace{x_i}_{\text{for all } i : 1 \rightarrow n}$   
↳ QQ-plot

Power transform

Box-Cox

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

$\underline{x_1}, \underline{x_2}, \dots, \underline{x_n} \rightarrow \text{Boxcox}$

$$\underline{\lambda}$$

optim  $\underline{\lambda}$

Box-Cox:

$\lambda = 0 \Rightarrow x_1 \dots x_n \sim \text{log-normal}$

$\lambda \neq 0 \Rightarrow$

# 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

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

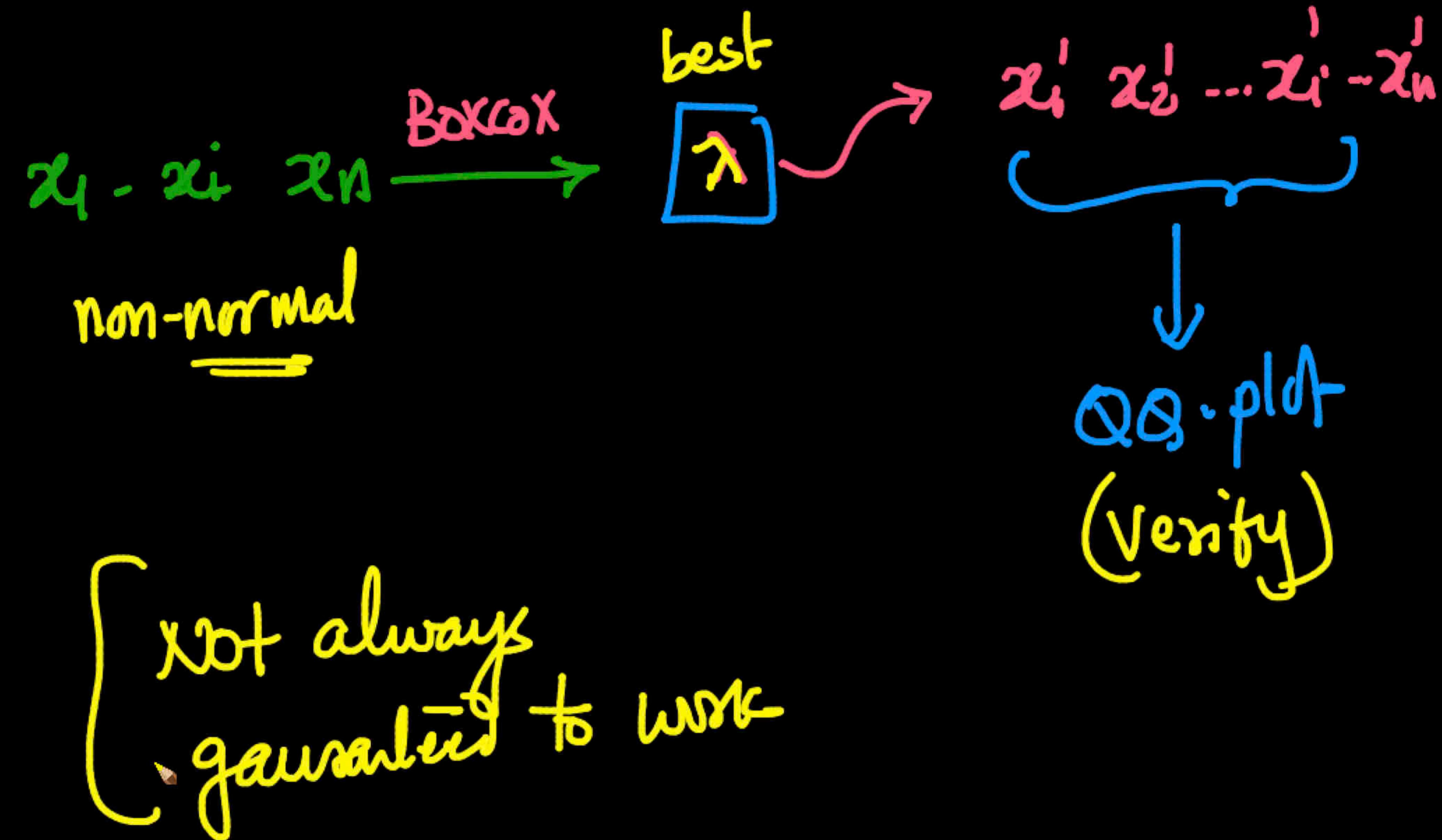
as described in the original article.<sup>[8][9]</sup> Moreover, the first transformations hold for  $y_i > 0$ , and the second for  $y_i > -\lambda_2$ .<sup>[8]</sup>

The parameter  $\lambda$  is estimated using the [profile likelihood](#) function and using goodness-of-fit tests.<sup>[10]</sup>

# Confidence interval [edit]

Confidence interval for the Box–Cox transformation can be asymptotically constructed using Wilks's theorem on the profile likelihood function to find all the possible values of  $\lambda$  that fulfill the following restriction:<sup>[11]</sup>





QQ plot and Transfo

statsmodels.graphics

scipy.stats.probplot

Power transform - V

scipy.stats.boxcox - X

Normal distribution X

+

▼

colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=dJYEyZ2gIIYb



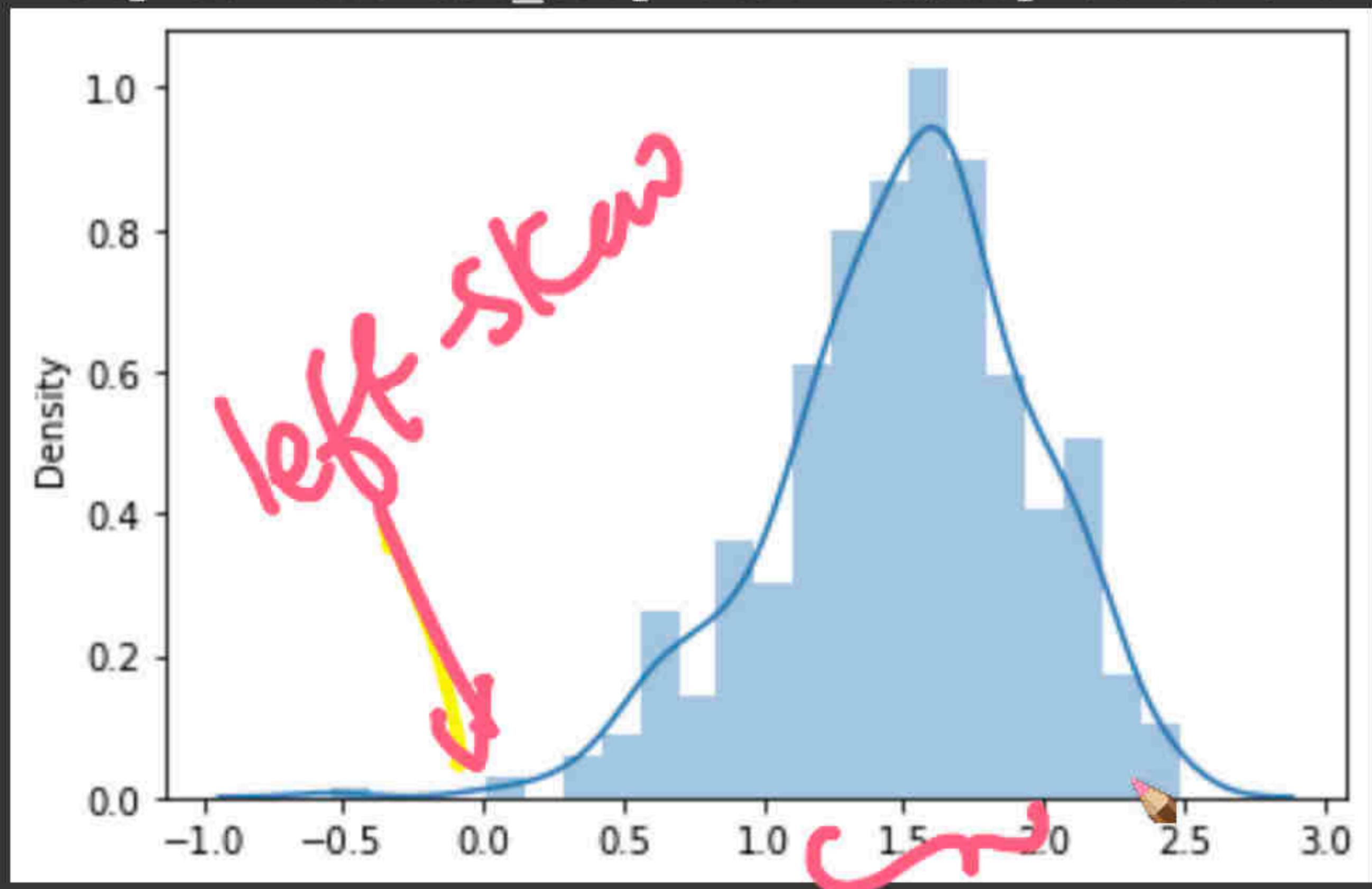
+ Code + Text

Reconnect ▾

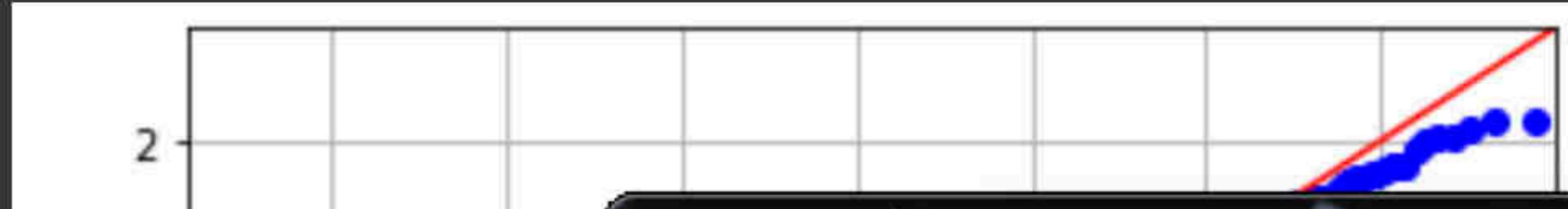


▼

```
[ ] /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecate
  warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f146d53c650>
```



```
▶ fig = sm.qqplot(x, line='45', fit=True)
  plt.grid()
```



QQ plot and Transfor

statsmodels.graphics

scipy.stats.probplot

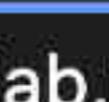
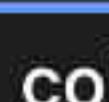
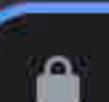
Power transform - V

scipy.stats.boxcox - X

Normal distribution X

+

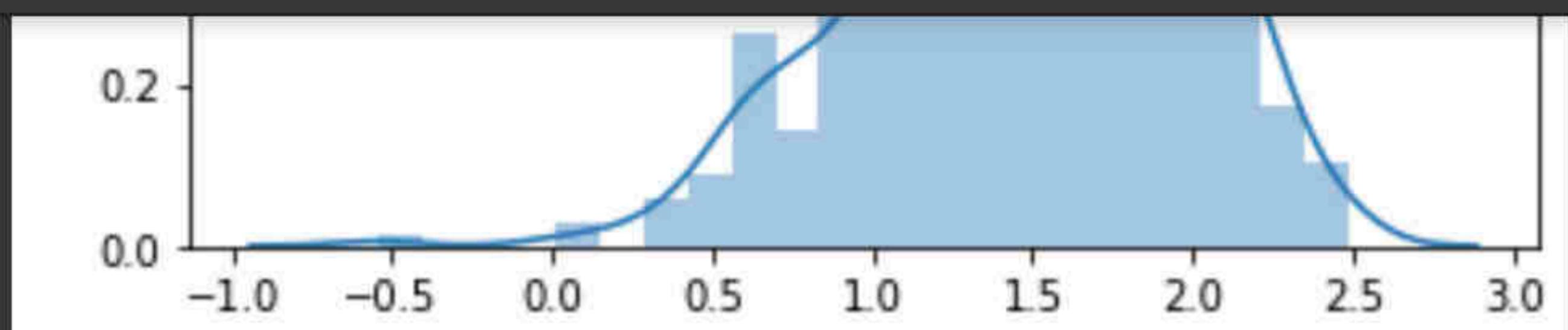
▼



colab.research.google.com/drive/1sV-dVOiuV3eeRK1G61hyHbQDe3K3ZAAE#scrollTo=dJYEyZ2gIIYb

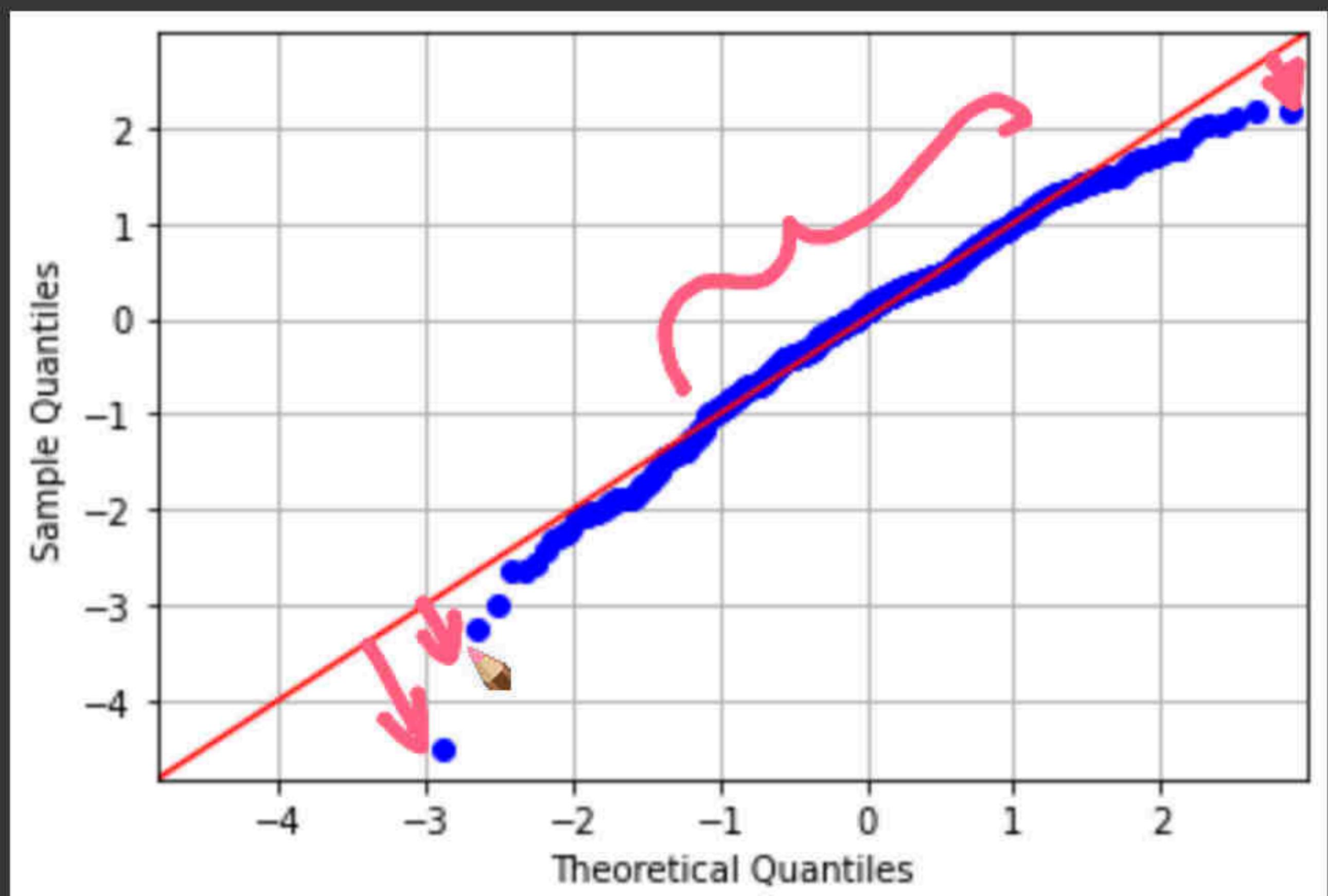
+ Code + Text

Reconnect

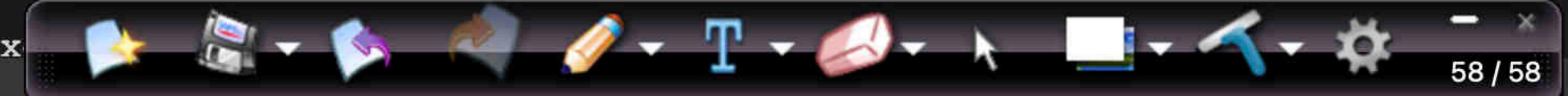


{x}

```
[ ] fig = sm.qqplot(x, line='45', fit=True)
plt.grid()
```



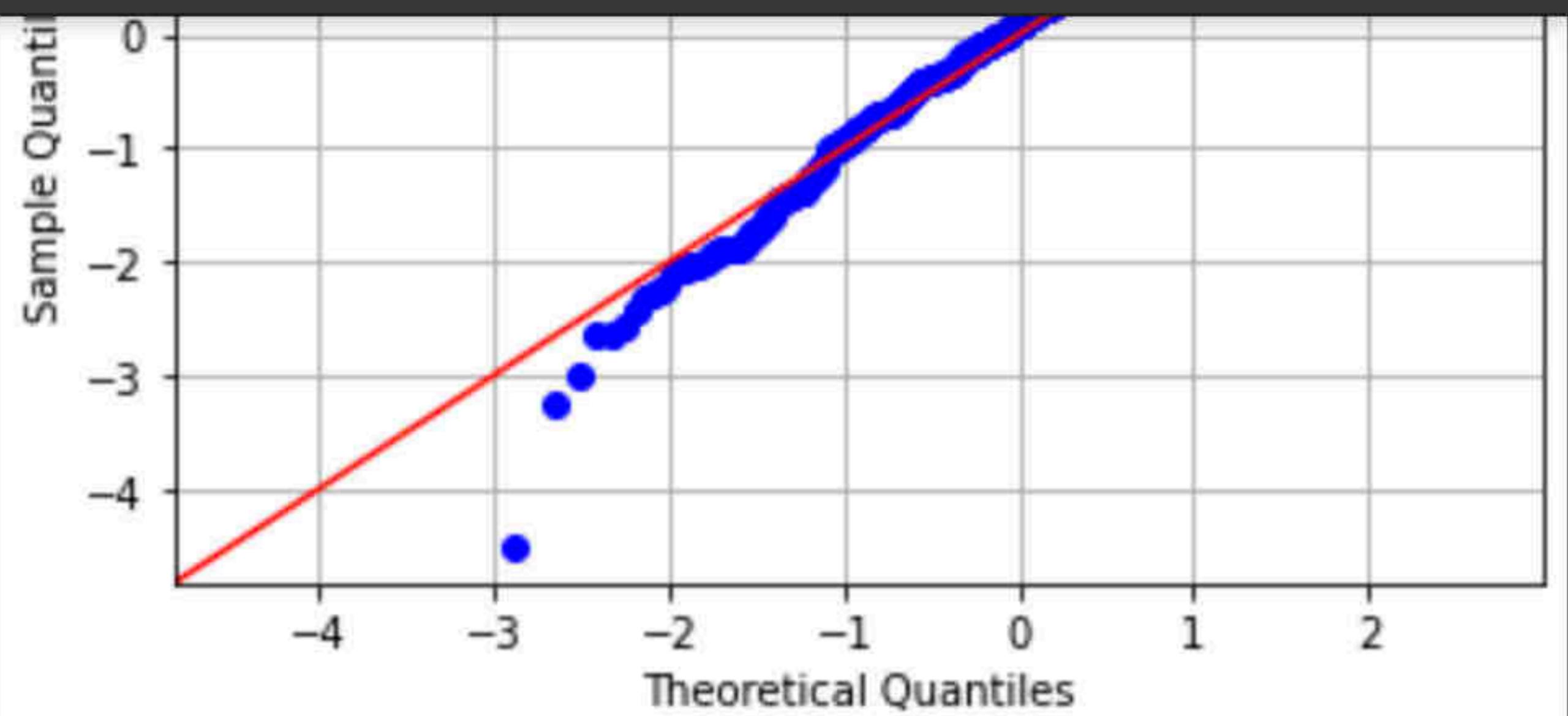
1 xt. 1 = stats.box



58 / 58

+ Code + Text

Reconnect



{x}



```
xt, l = stats.boxcox(x)
ValueError
<ipython-input-46-80d99d196c18> in <module>()
----> 1 xt, l = stats.boxcox(x)

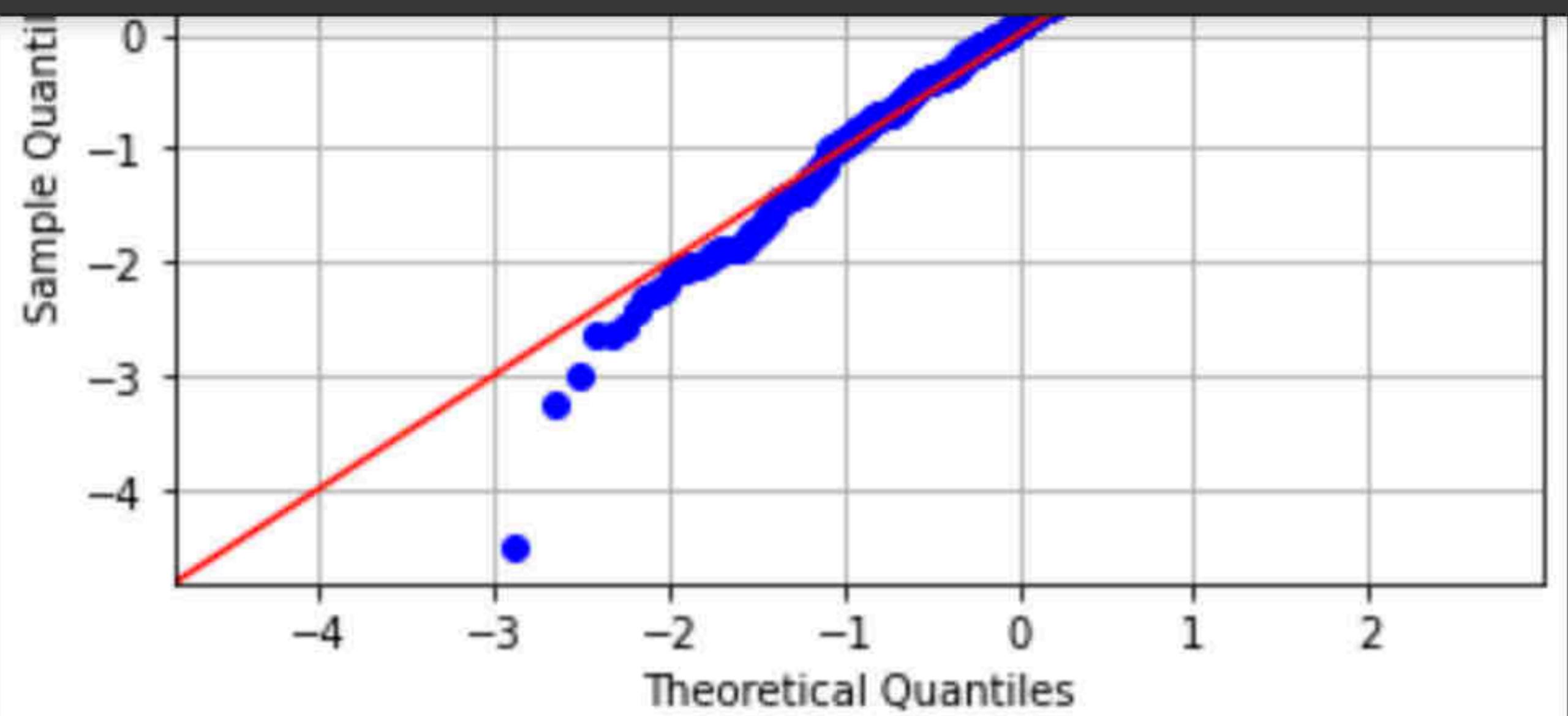
/usr/local/lib/python3.7/dist-packages/scipy/stats/morestats.py in boxcox(x, lmbda, alpha)
    1043
    1044      if any(x <= 0):
-> 1045          raise ValueError("Data must be positive.")
    1046
    1047      if lmbda is not None: # single transformation
```

ValueError: Data must be positive.



+ Code + Text

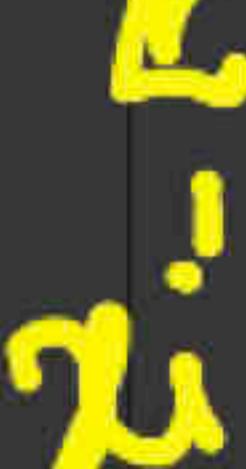
Reconnect



{x}



[ ] xt, l = stats.boxcox(x)



ValueError

Traceback (most recent call last)

&lt;ipython-input-46-80d99d196c18&gt; in &lt;module&gt;()

----&gt; 1 xt, l = stats.boxcox(x)

/usr/local/lib/python3.7/dist-packages/scipy/stats/morestats.py in boxcox(x, lmbda, alpha)

1043

1044 if any(x &lt;= 0):

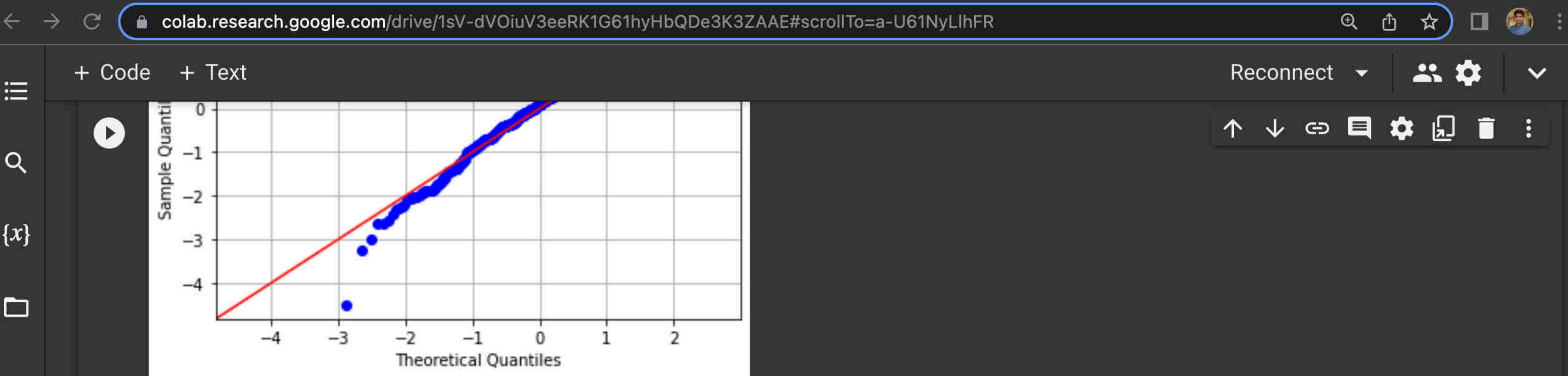
-&gt; 1045 raise ValueError("Data must be positive.")

1046

1047 if lmbda is not None: # single transformation

ValueError: Data must be positive.





```
[ ] xt, l = stats.boxcox(x)
```

```
-----  
ValueError                                Traceback (most recent call last)  
<ipython-input-46-80d99d196c18> in <module>()  
----> 1 xt, l = stats.boxcox(x)  
  
/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
```

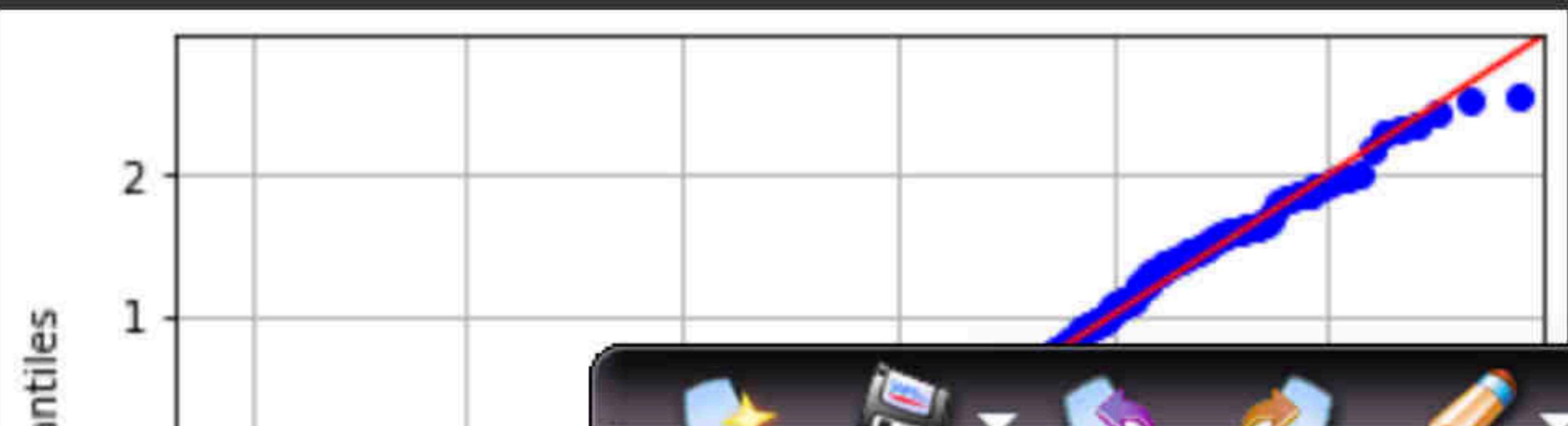
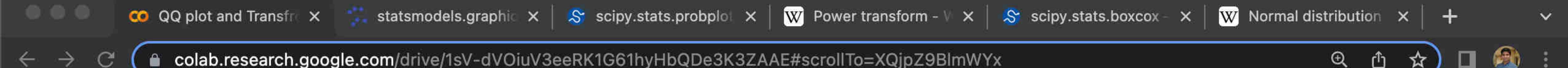
ValueError: Data must be positive.



$\tilde{x} \rightarrow x_1, x_2, \dots, x_n$

-0.54

$x+1 \rightarrow x_1+1, x_2+1, \dots, x_{n+1}$  > 0

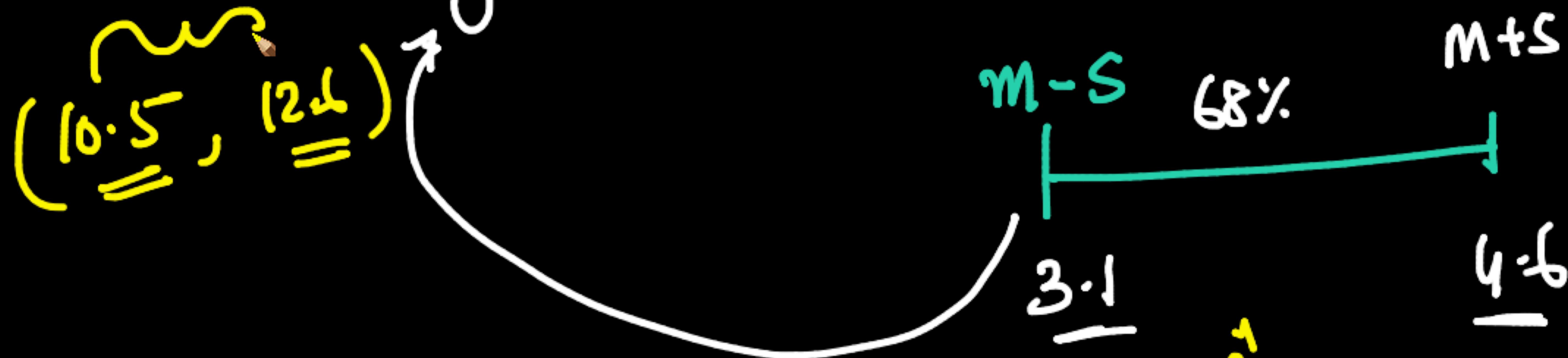


$$x \xrightarrow{+1} (x+1) \quad \text{boxcox}$$
$$x_i \xrightarrow{+1} (x_i+1)$$
$$\lambda = 1.86$$

$$x_i \xrightarrow{f(\cdot)} x'_i$$
$$\left\{ \frac{(x_i+1)^{\lambda} - 1}{\lambda} = x'_i \right.$$

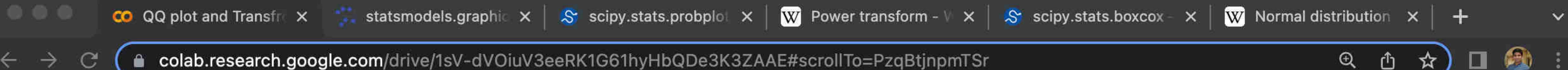
$\tilde{x}_i$  days →  $x_i'$  Gaussian

68-95-99



$$x_i \xrightarrow{f} x_i' \xleftarrow{f^{-1}}$$

$$x_i \xleftarrow{\log} e^x \xrightarrow{\log} x_i$$

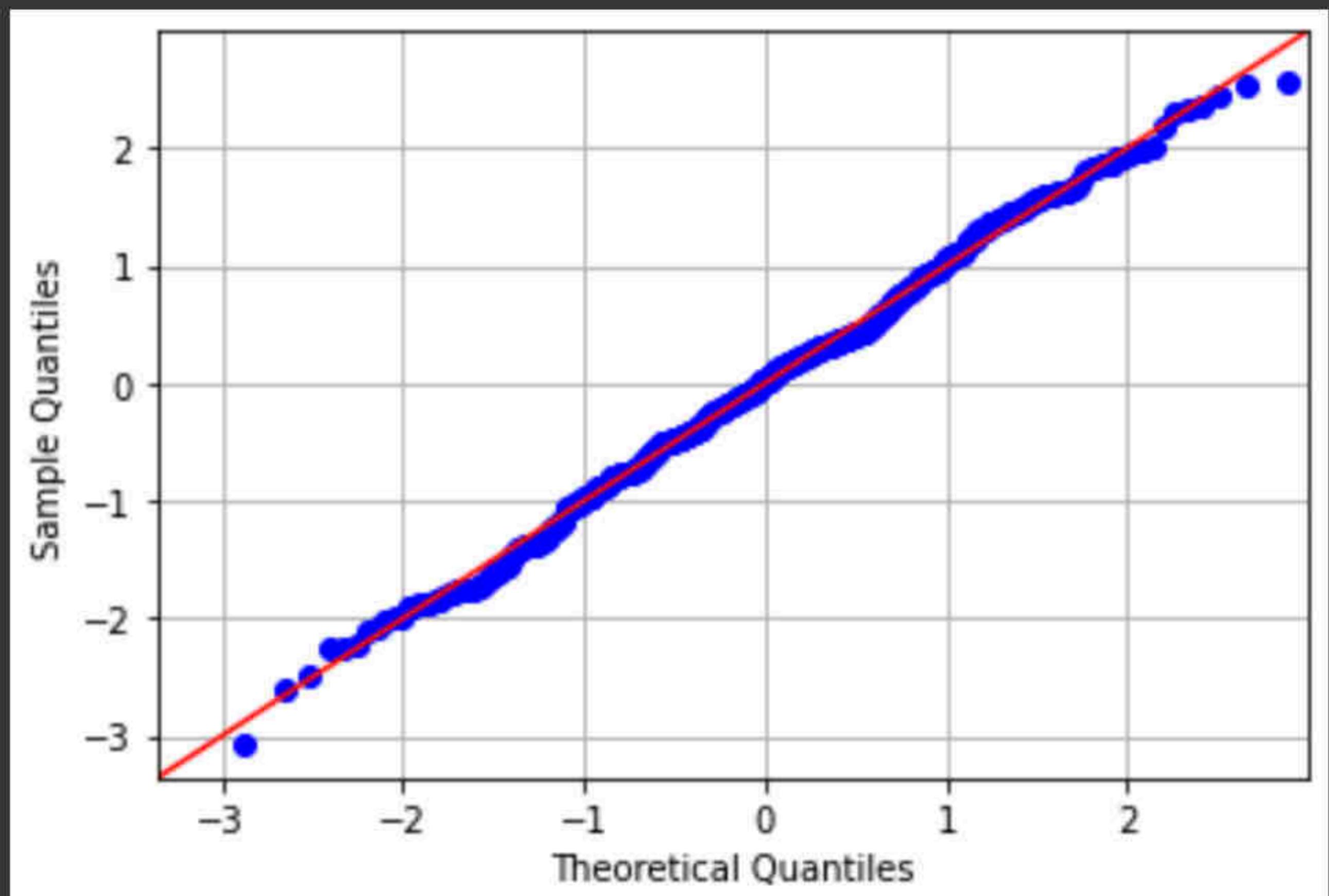


+ Code + Text

```
[ ] print(1)
```

1.8162672858204911

```
{x} fig = sm.qqplot(xt, line='45', fit=True)  
plt.grid()
```



## 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


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

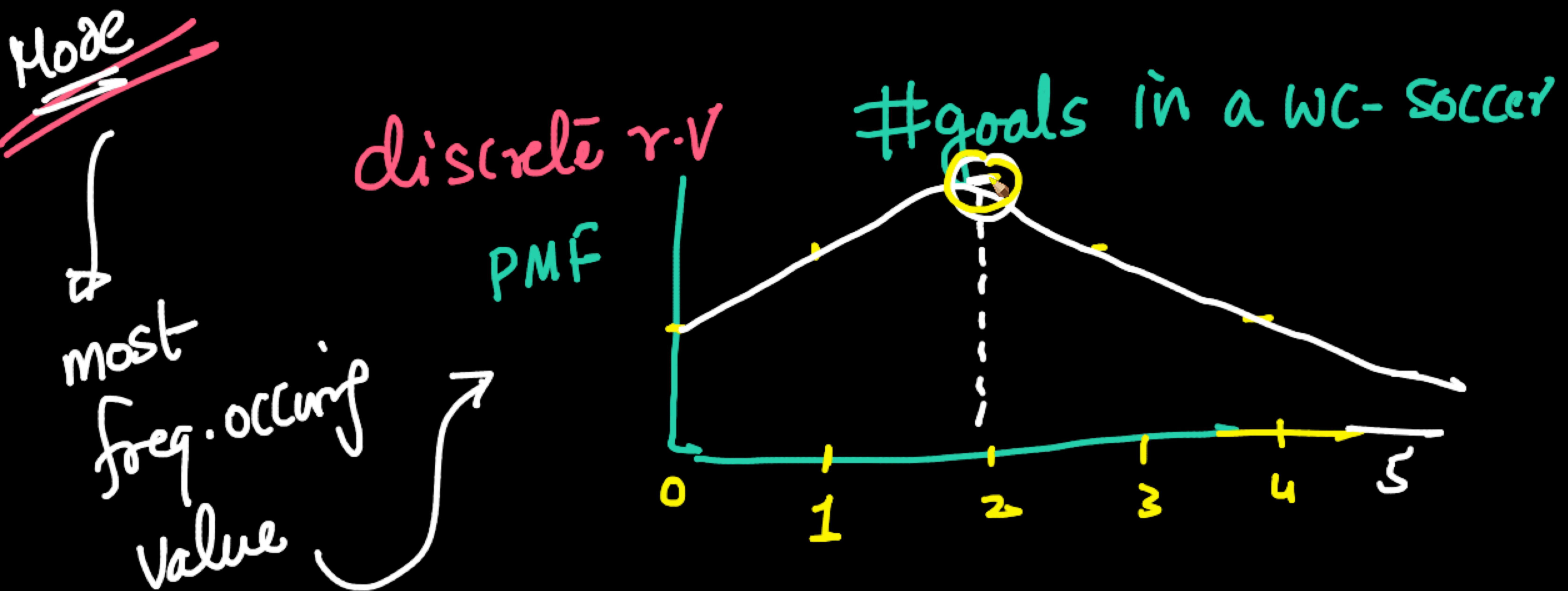
as described in the original article.<sup>[8][9]</sup> Moreover, the first transformations hold for  $y_i > 0$ , and the second for  $y_i > -\lambda_2$ .<sup>[8]</sup>

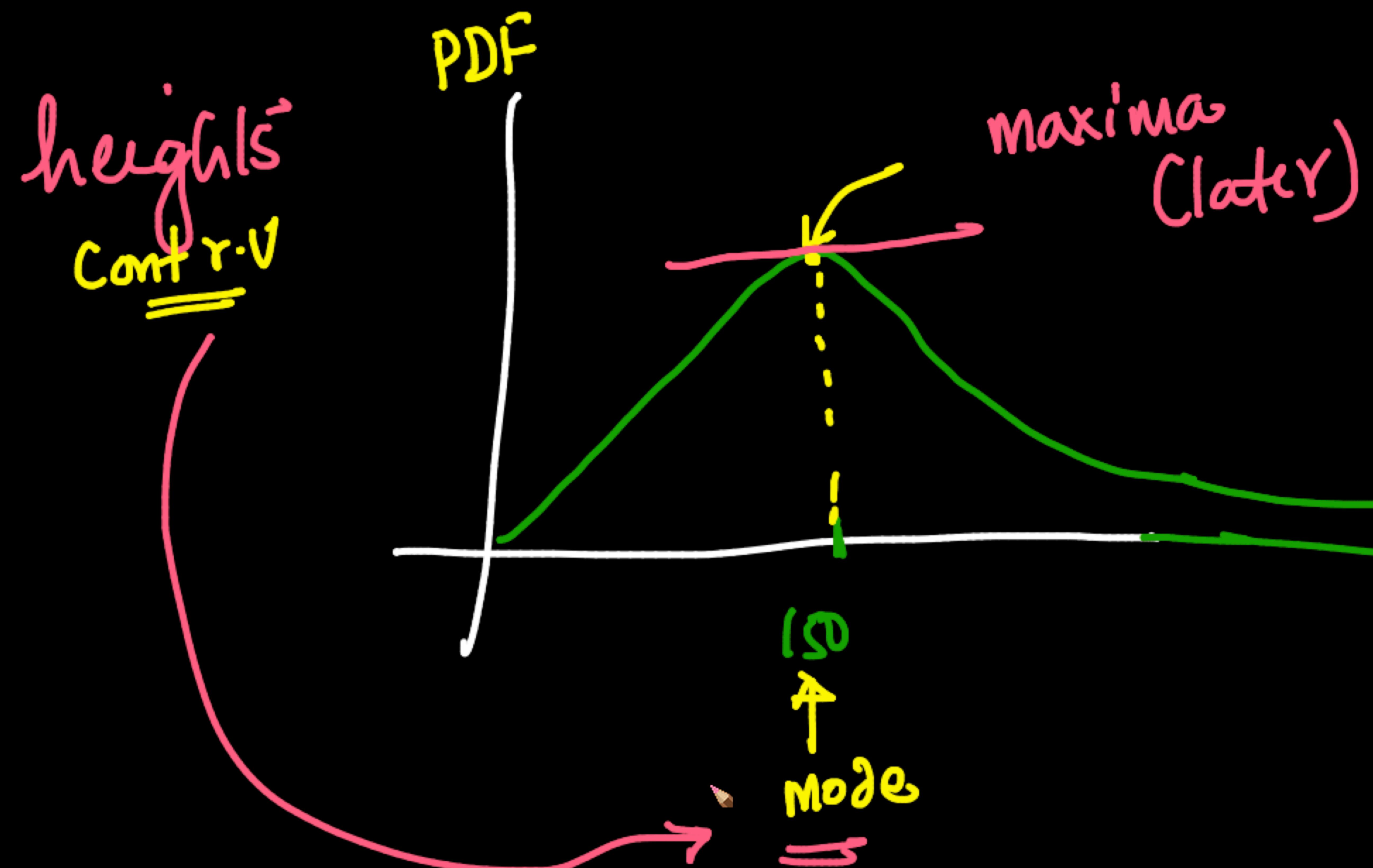
The parameter  $\lambda$  is estimated using the [profile likelihood](#) function and using goodness-of-fit tests.<sup>[10]</sup>

## Confidence interval [ edit ]

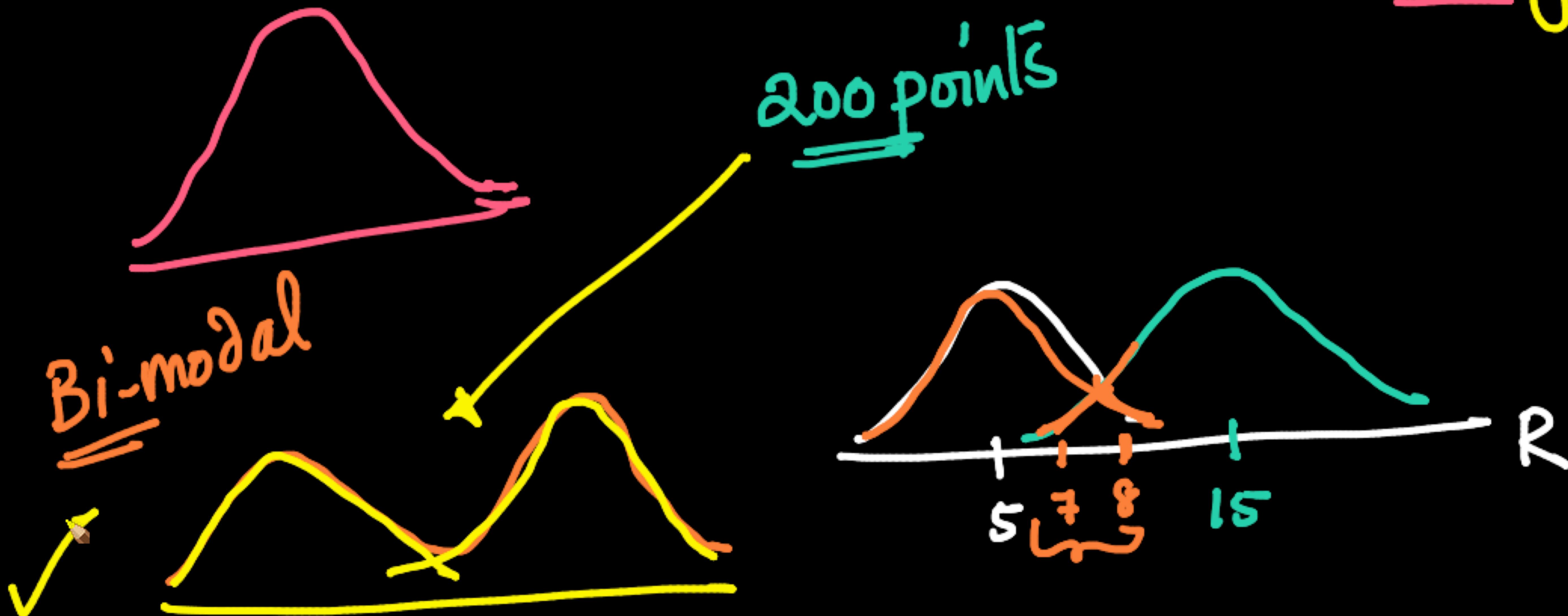
Confidence interval for the Box–Cox transformation can be asymptotically constructed using Wilks's theorem on the [profile likelihood](#) function to find all the possible values of  $\lambda$  that fulfill the following restriction.<sup>[11]</sup>

$$\chi_1 = -0.5 \xrightarrow{\text{square}} 0.25 \rightarrow \text{Loxcox}$$
$$\chi_{10} = 0.5 \xrightarrow{\quad} 0.25$$

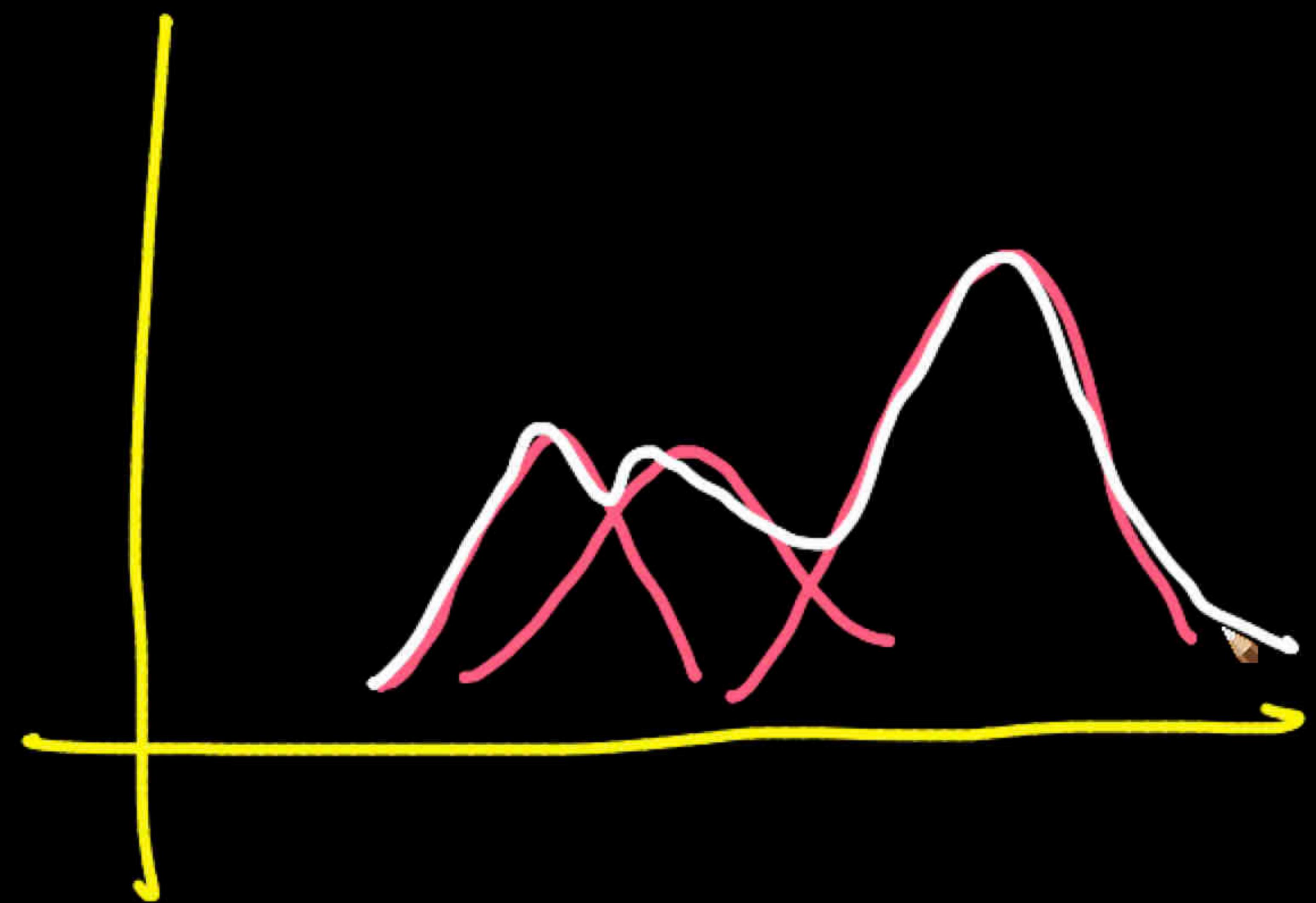


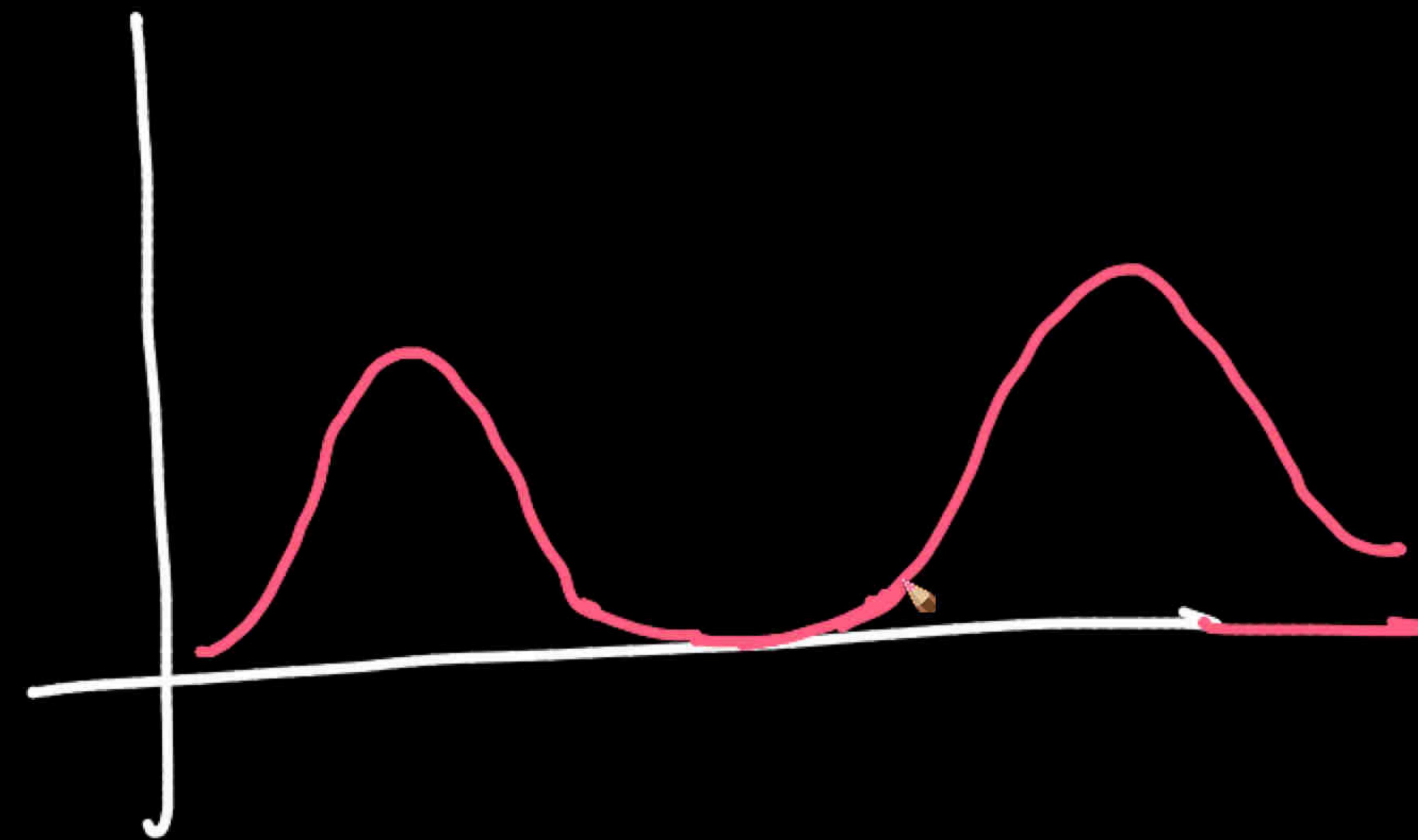


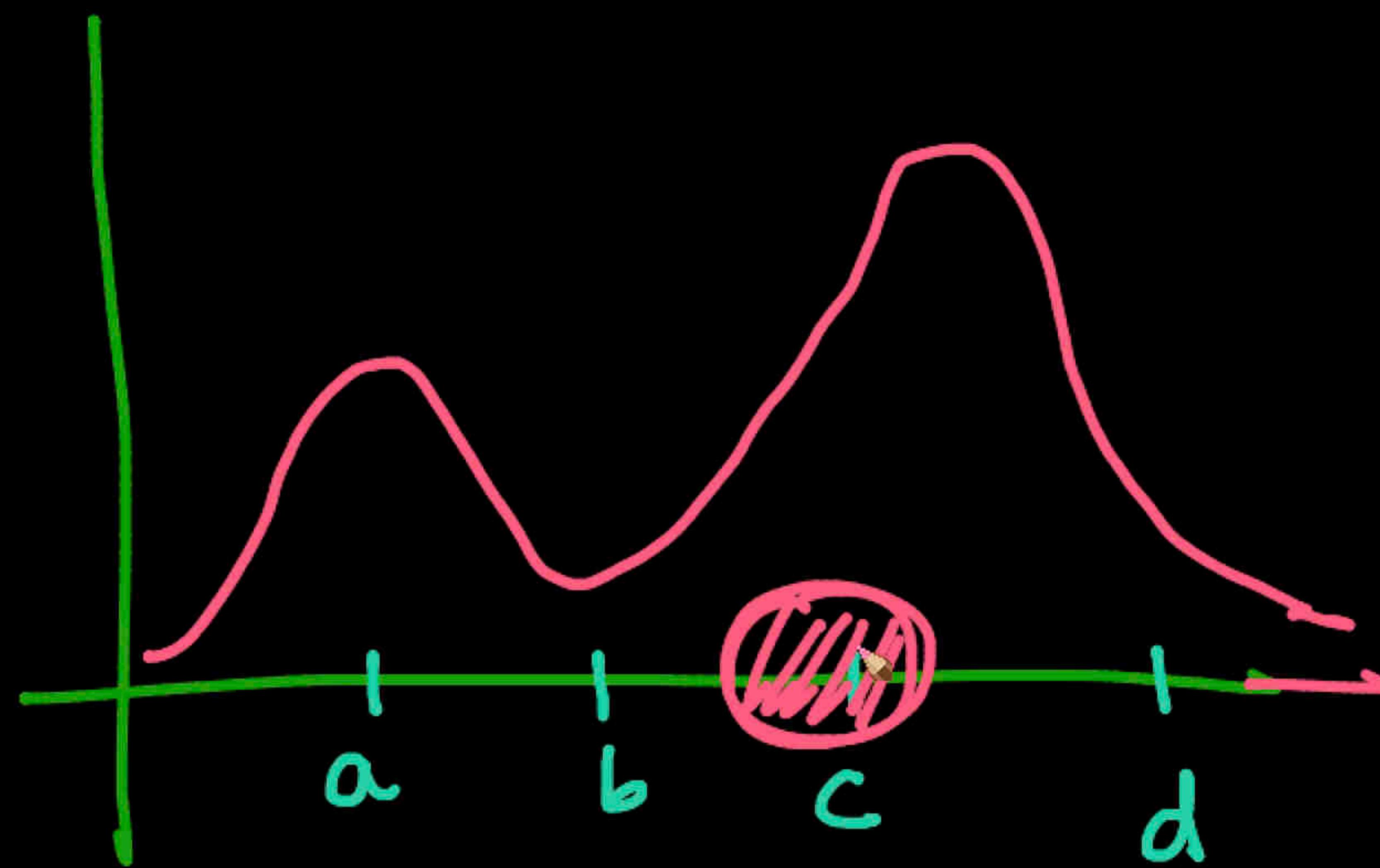
could recovery times  
kids (<5yrs)  
elderly 60+  
∞

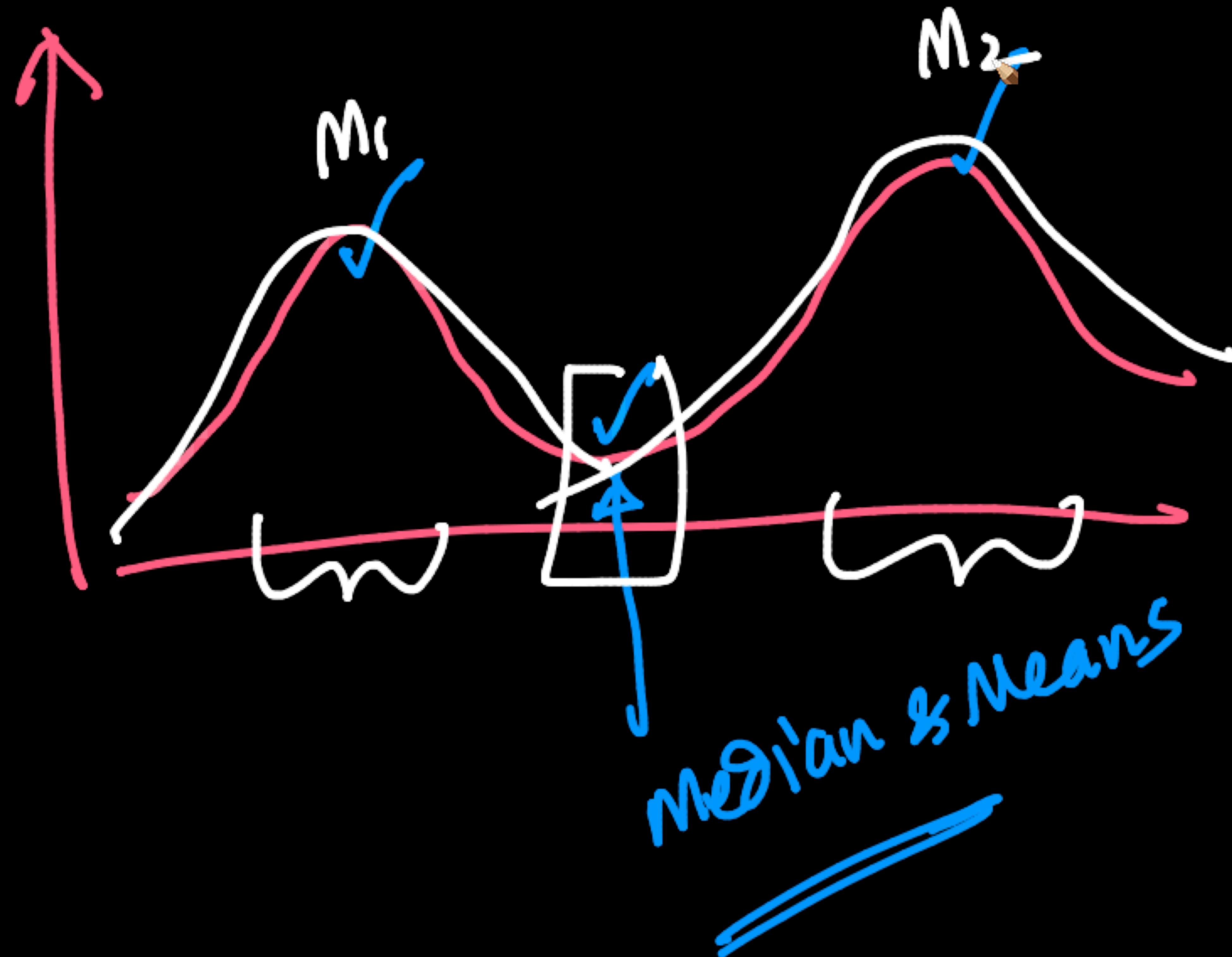


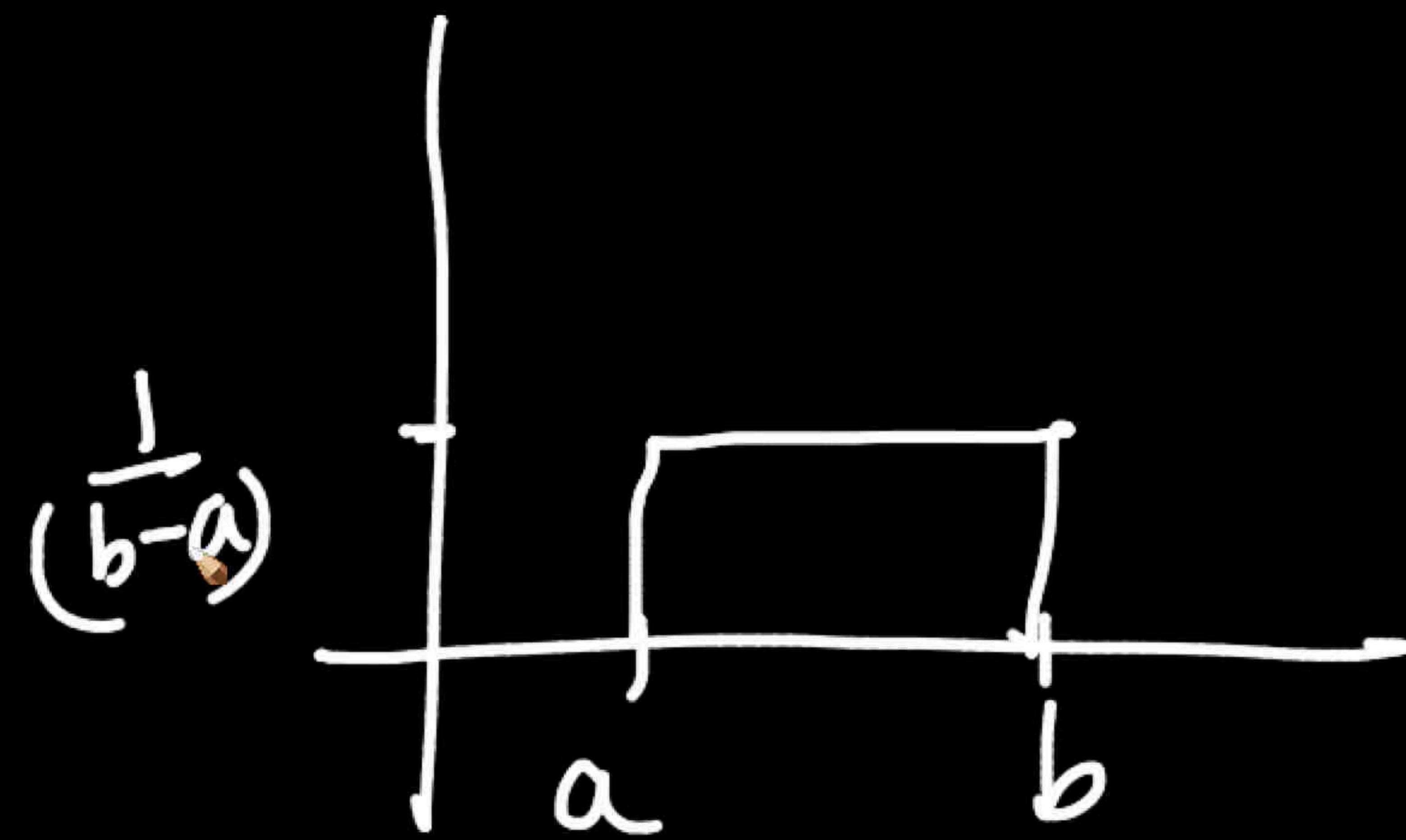
Indians  
Chinese  
Norwegians

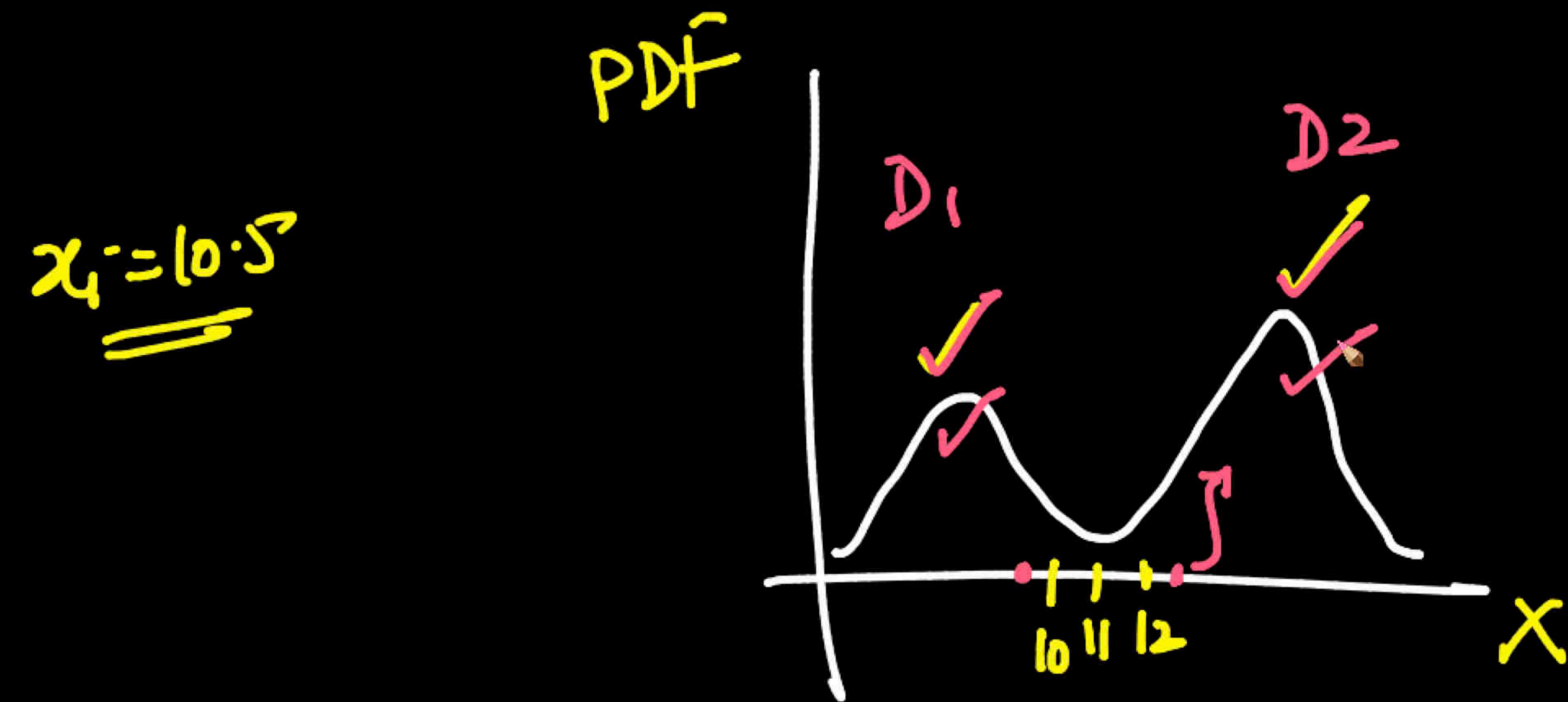






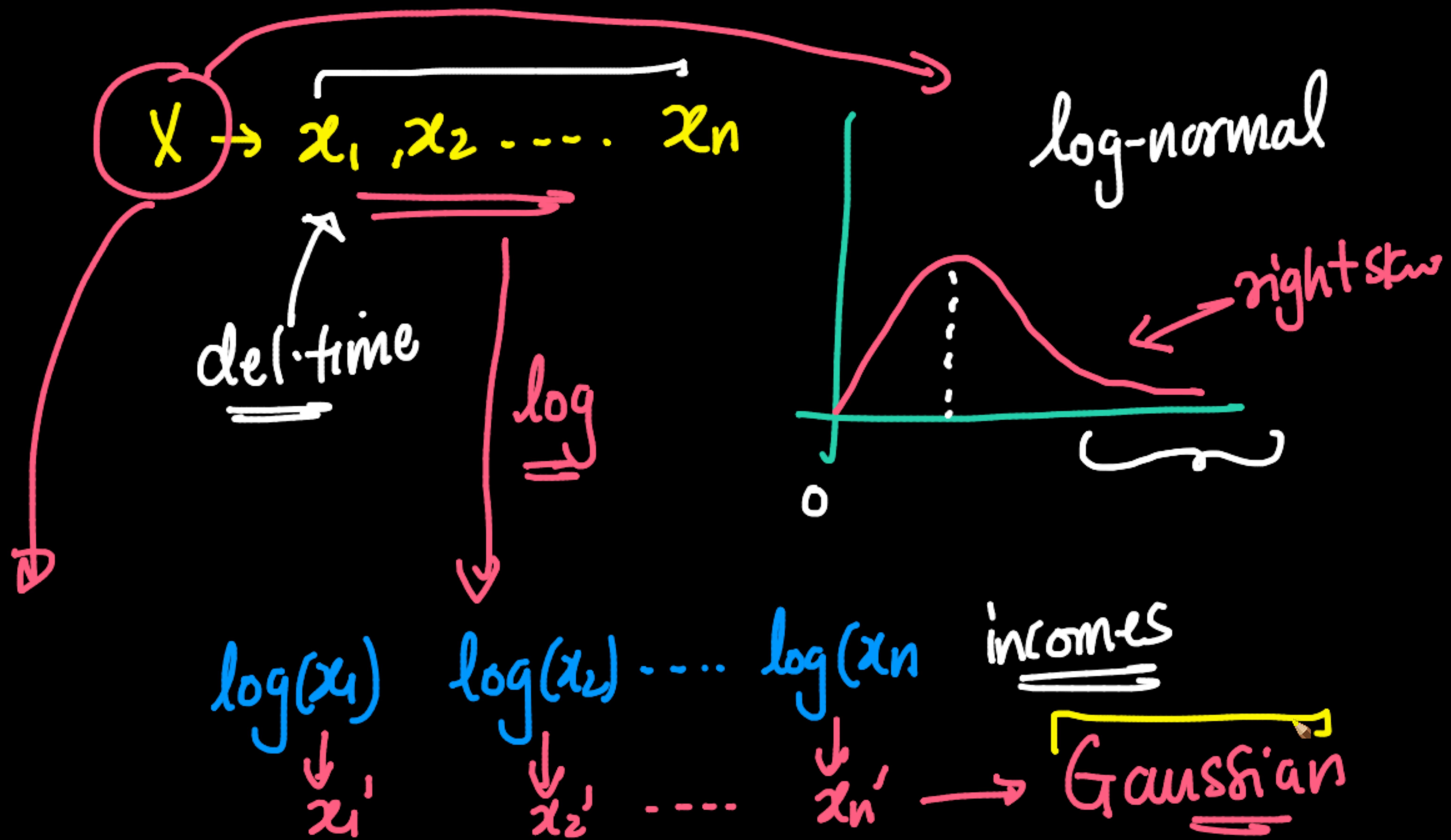






Mixture-models

(ML)





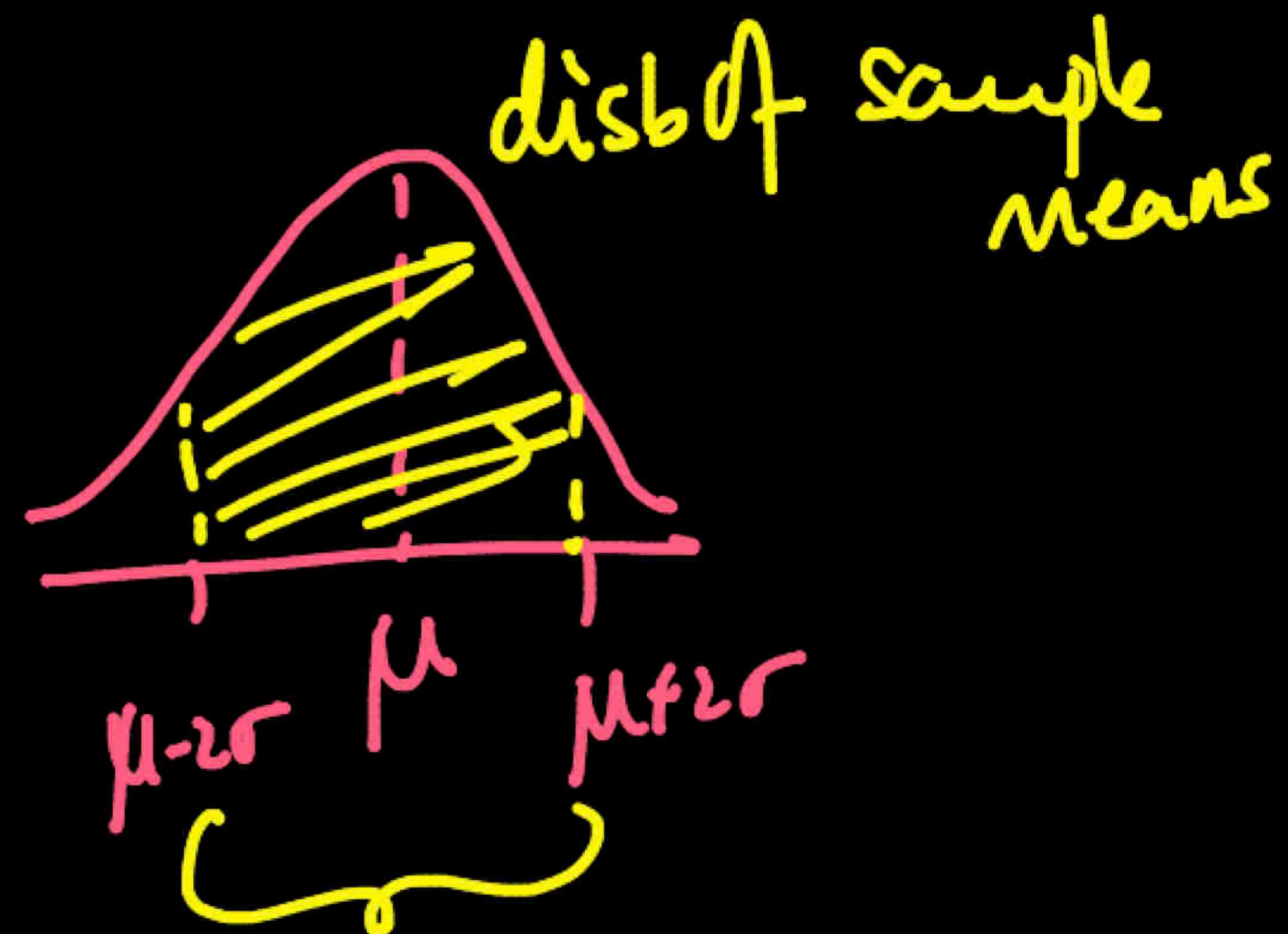
$x_1, x_2, \dots, x_m$

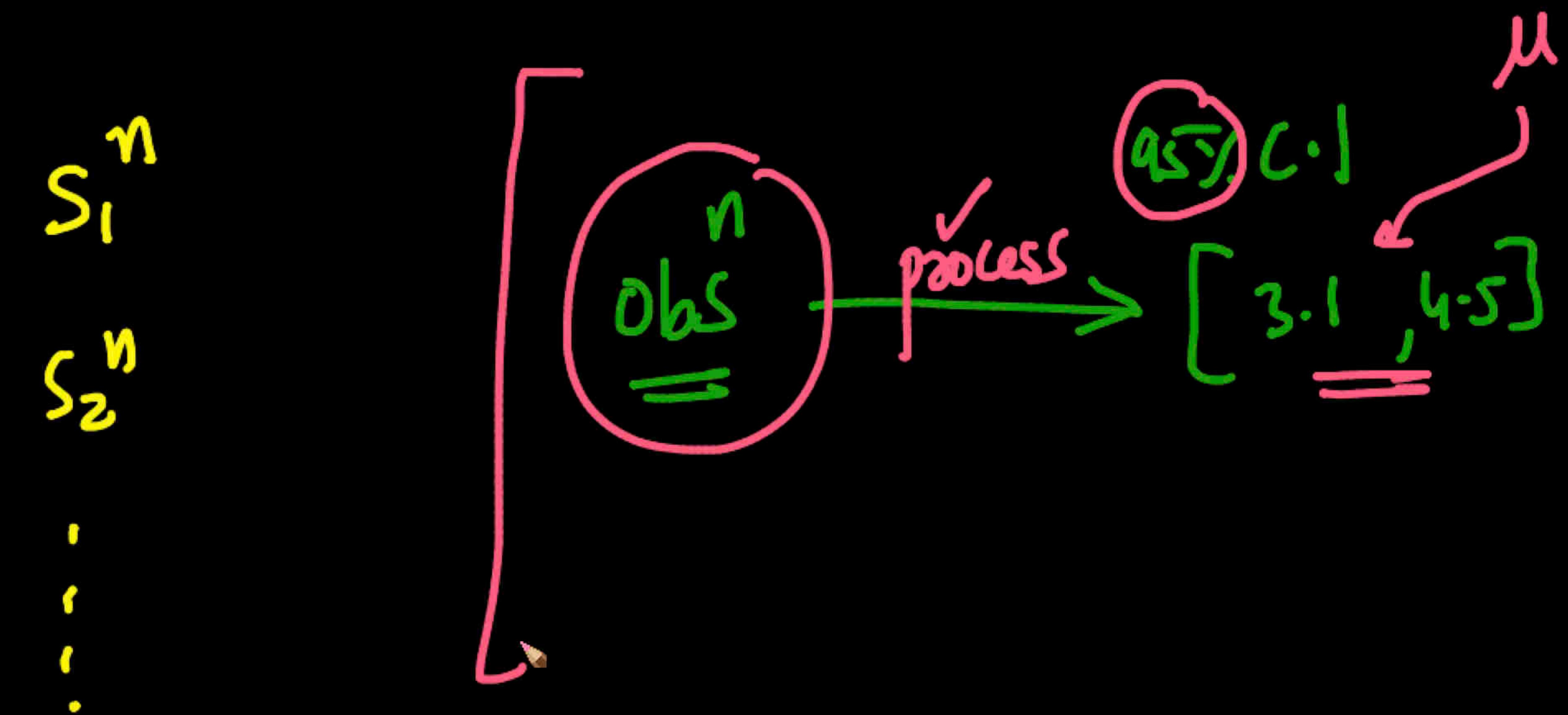
$S_1^n \rightarrow M_1$

$S_2^n \rightarrow M_2$

$\vdots$

$S_\gamma^n \rightarrow M_\gamma$





Chrome File Edit View History Bookmarks Profiles Tab Window Help

Rem QQ scip Pow scip Norm Log A Fir New Lect Deri New

colab.research.google.com/drive/1KdfJLe\_oy2i36WI4OuoJxSyUHhCawJJH

Comment Share Settings User

File Edit View Insert Runtime Tools Help Last edited on 8 July

+ Code + Text Connect Editing

Number of surgeries in a hospital per day.

{x}

[ ] import numpy as np  
from scipy.stats import poisson, expon  
import matplotlib.pyplot as plt  
import pandas as pd  
import seaborn as sns

# Number of surgeries in a hospital per day : Poisson  
# Sample of 200 observations

```
V = [ 2, 5, 6, 3, 2, 6, 2, 10, 4, 4, 4, 8, 3, 3, 3, 3, 8, 7,  
      6, 9, 4, 5, 5, 5, 5, 4, 10, 9, 4, 0, 8, 4, 6, 5, 5,  
      6, 1, 6, 7, 2, 3, 5, 7, 2, 3, 4, 5, 7, 4, 5, 4, 3,  
      3, 8, 1, 6, 8, 6, 6, 3, 9, 8, 2, 3, 6, 5, 8, 8, 5,  
      7, 7, 6, 7, 6, 6, 9, 3, 6, 12, 4, 4, 3, 2, 2, 2, 4,  
      8, 5, 3, 6, 2, 5, 3, 2, 6, 9, 2, 7, 4, 4, 3, 4, 5,  
      8, 9, 4, 6, 6, 2, 2, 7, 2, 10, 8, 2, 6, 3, 7, 2, 3,  
      8, 11, 4 ]
```

82 / 82