

# Topics

90 min

- recap of CLT ✓
- Code for CLT & bootstrapping
- Samplings ✓

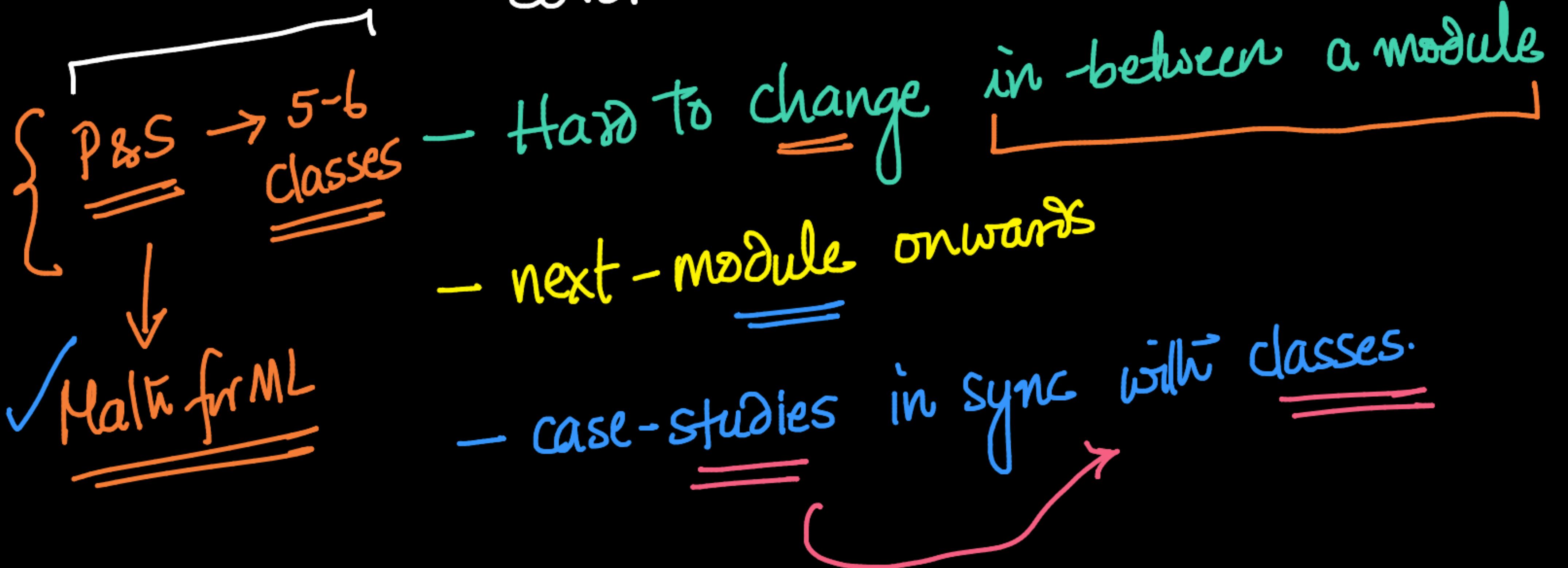
50-60 min

- Hypothesis Testing
- intuition
- terminology
- Math & code

## Assessments

→ ... list all problems :- .

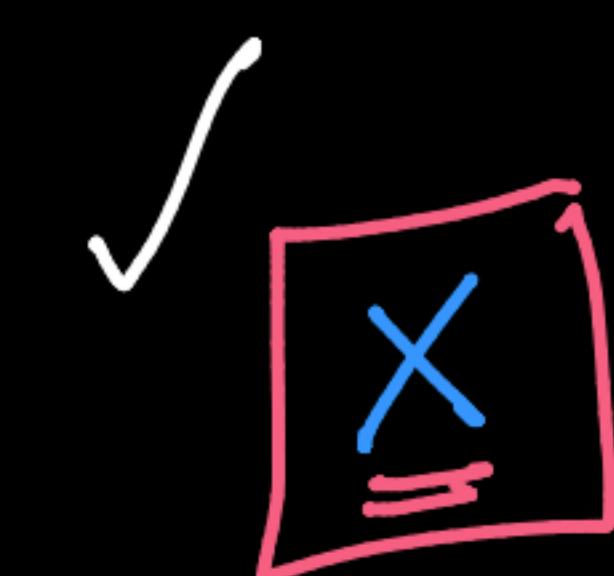
- Instructor assigns based on concepts
- Covered



END. ....  
{  
—  
—  
—  
—>

= Questions =

→



$[x_1, x_2, \dots, x_{75}]$

pop mean ( $\mu$ ) & std-dev ( $\sigma$ ) are finite

CLT

sample size  $\hat{n}$

$\{x_1^1, x_2^1, \dots, x_n^1\}$  — sample 1 →  $M_1$

$\{x_1^2, x_2^2, \dots, x_n^2\}$  —  $S_2$  →  $M_2$

⋮  
⋮

$\{x_1^K, \dots, x_n^K\}$  —  $S_K$  →  $M_K$

sample means:  $m_1, \dots, m_K$



Normal (mean =  $\mu$ ; std-dev =  $\frac{\sigma}{\sqrt{n}}$ )

=

as  $n \rightarrow \infty$

~~Why CLT~~

$X = \text{rec. time of patients medicine 1}$

 $\{x_1, \dots, x_{10}\}$ 

{95% C.I.  
on the  
mean-rec-line}

$s_1 - \dots$

$s_2 -$

$\vdots$

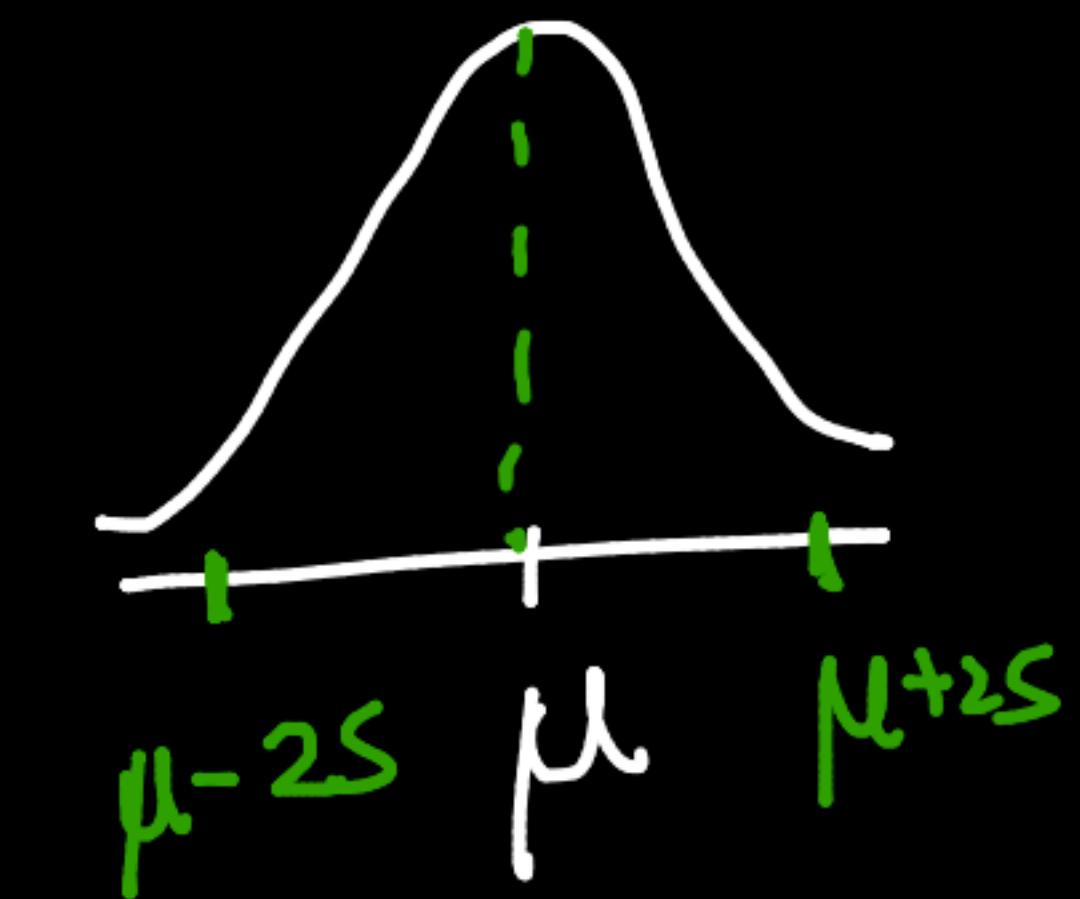
$s_k -$

$\rightarrow M_1$

$\rightarrow M_2$

$\vdots$

$\rightarrow M_K$



$$s = \text{SD} = \sqrt{\frac{\sum}{n}}$$

95 C.I. → bootstrapping: means: medians; P99

CLT: means ✓



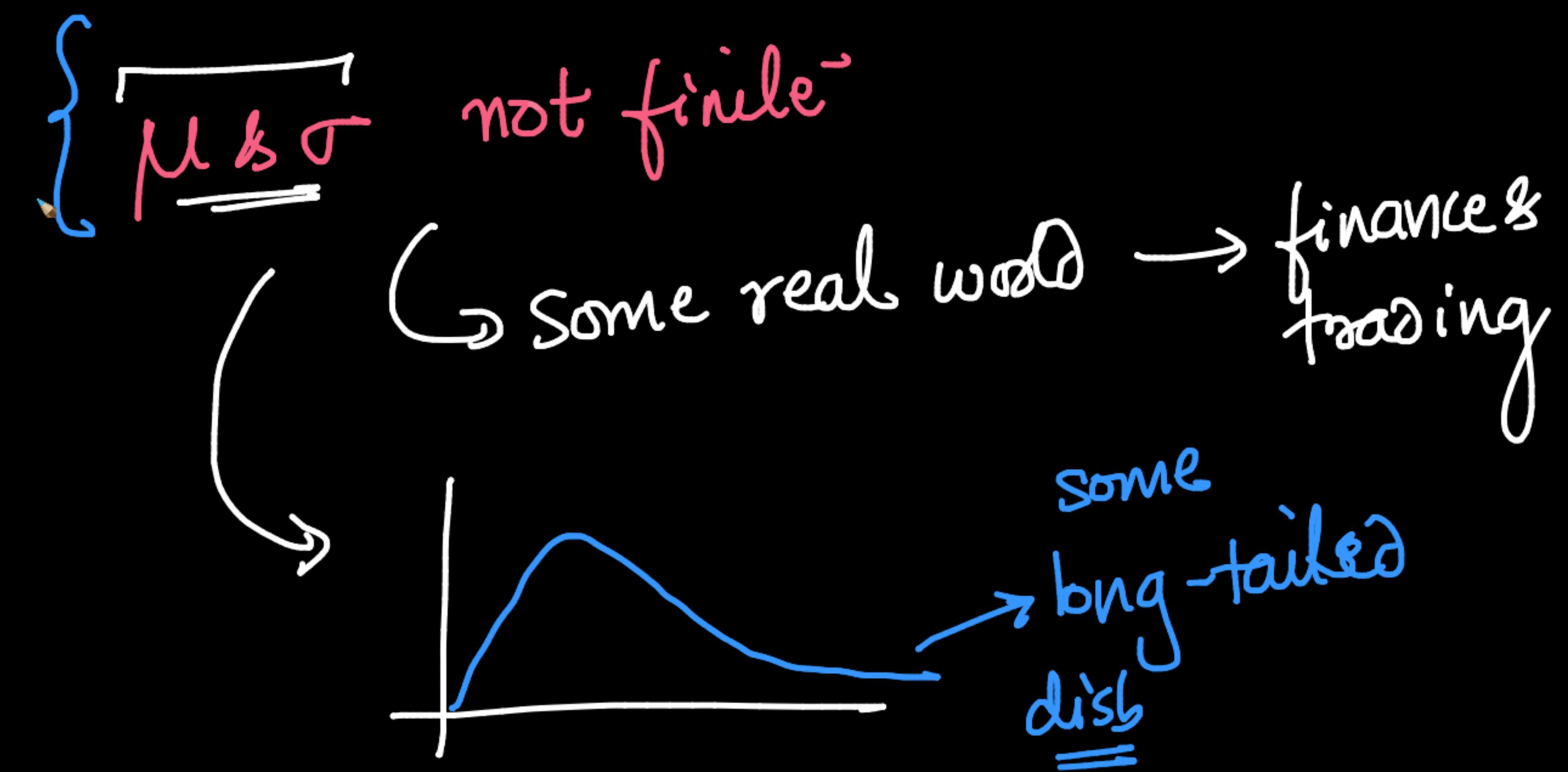
95% C.I. of mean

[ $\mu, \sigma$  are finite]

pop-

CLT:





Main page  
Contents  
Current events  
Random article  
About Wikipedia  
Contact us  
Donate  
  
Contribute  
Help  
Learn to edit  
Community portal  
Recent changes  
Upload file

Tools  
What links here  
Related changes  
Special pages  
Permanent link  
Page information  
Cite this page  
Wikidata item

Print/export  
Download as PDF  
Printable version

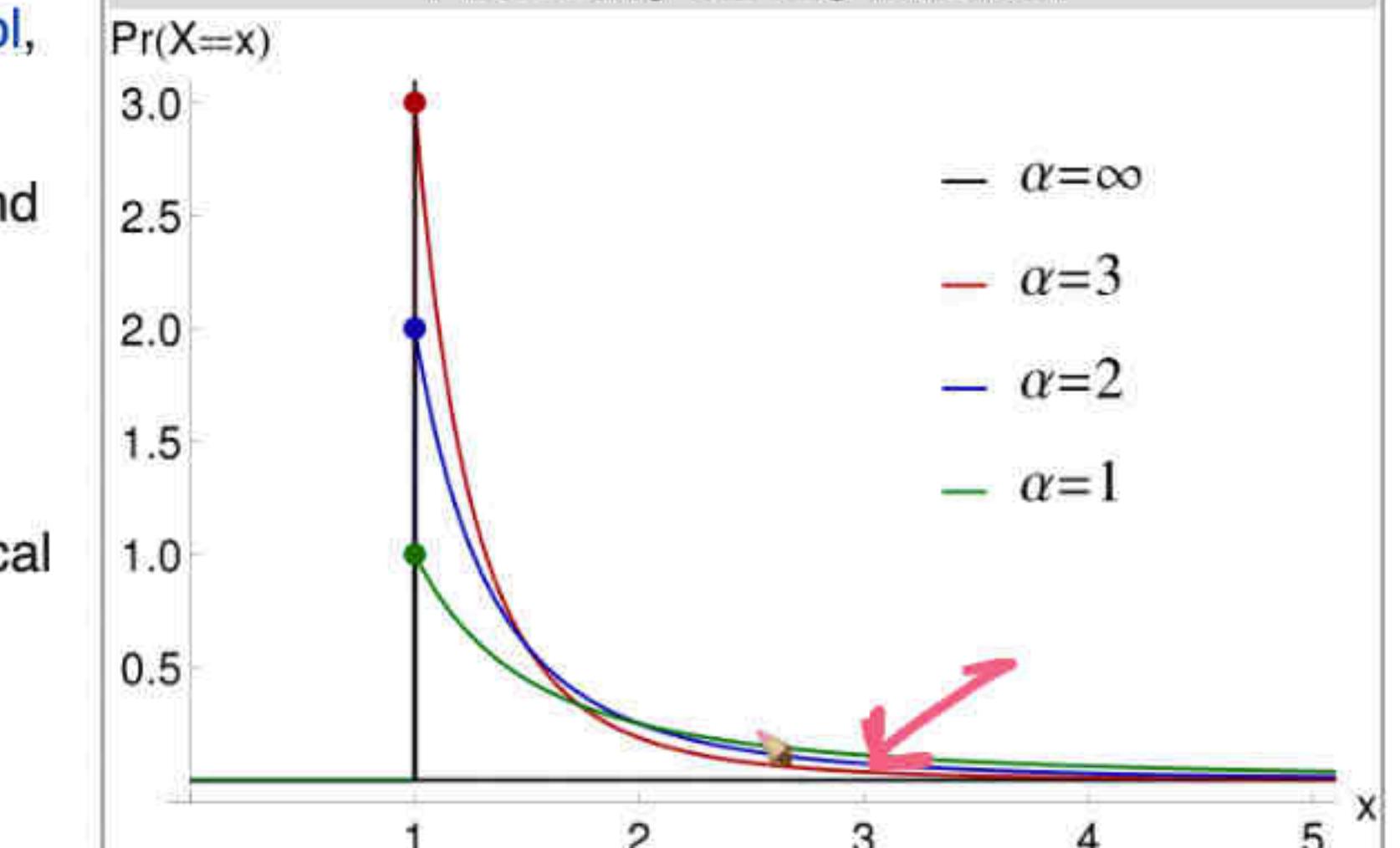
The **Pareto distribution**, named after the Italian civil engineer, economist, and sociologist Vilfredo Pareto,<sup>[1]</sup> (Italian: [pa're:to] US: /pə'retəʊ/ pə-RAY-toh),<sup>[2]</sup> is a power-law probability distribution that is used in description of social, quality control, scientific, geophysical, actuarial, and many other types of observable phenomena. Originally applied to describing the **distribution of wealth** in a society, fitting the trend that a large portion of wealth is held by a small fraction of the population.<sup>[3][4]</sup> The **Pareto principle** or "80-20 rule" stating that 80% of outcomes are due to 20% of causes was named in honour of Pareto, but the concepts are distinct, and only Pareto distributions with shape value ( $\alpha$ ) of  $\log_4 5 \approx 1.16$  precisely reflect it. Empirical observation has shown that this 80-20 distribution fits a wide range of cases, including natural phenomena<sup>[5]</sup> and human activities.<sup>[6][7]</sup>

## Contents [hide]

- 1 Definitions
  - 1.1 Cumulative distribution function
  - 1.2 Probability density function
- 2 Properties
  - 2.1 Moments and characteristic function
  - 2.2 Conditional distributions
  - 2.3 A characterization theorem
  - 2.4 Geometric mean
  - 2.5 Harmonic mean
  - 2.6 Graphical representation
- 3 Related distributions

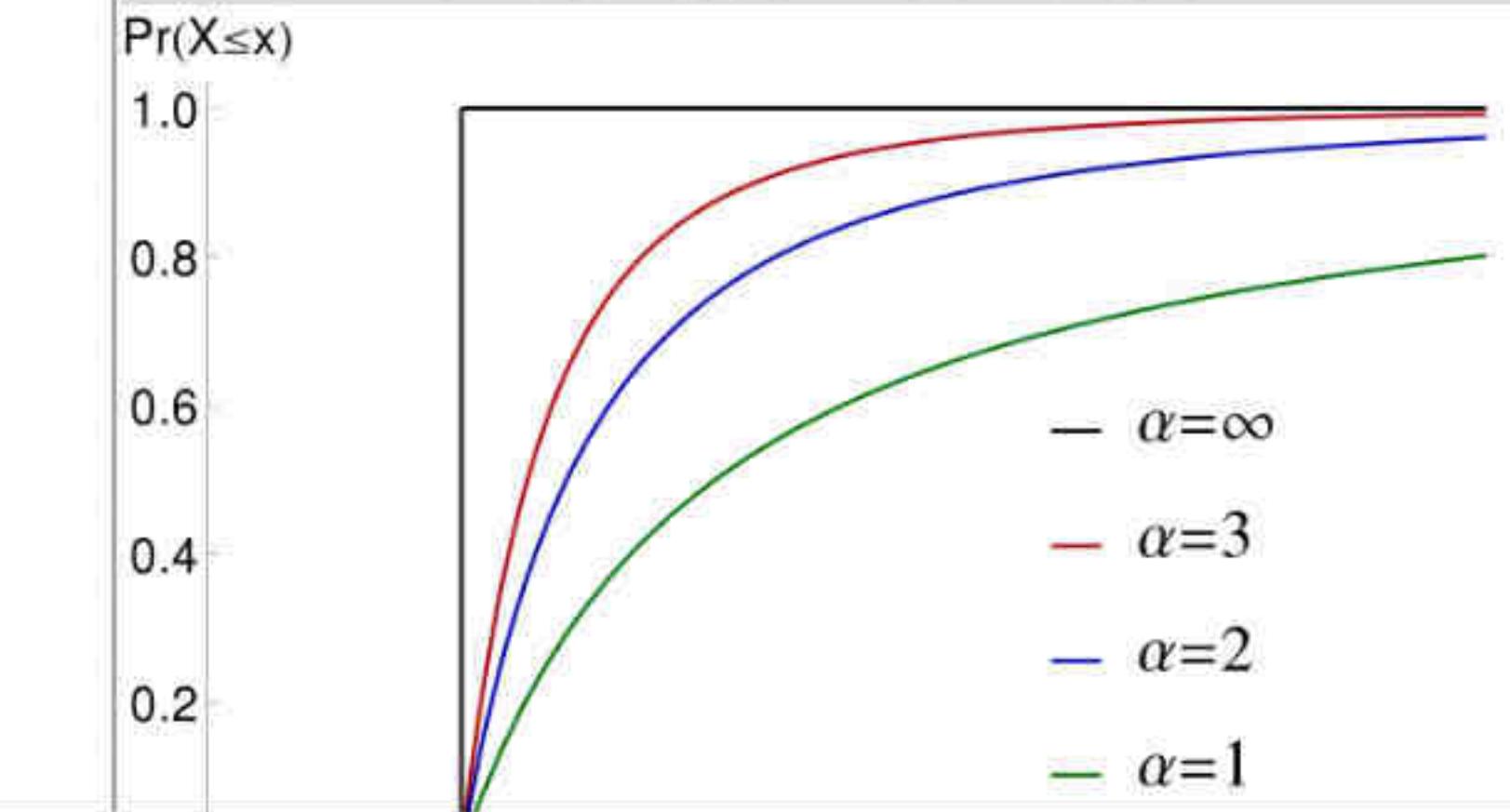
## Pareto Type I

### Probability density function



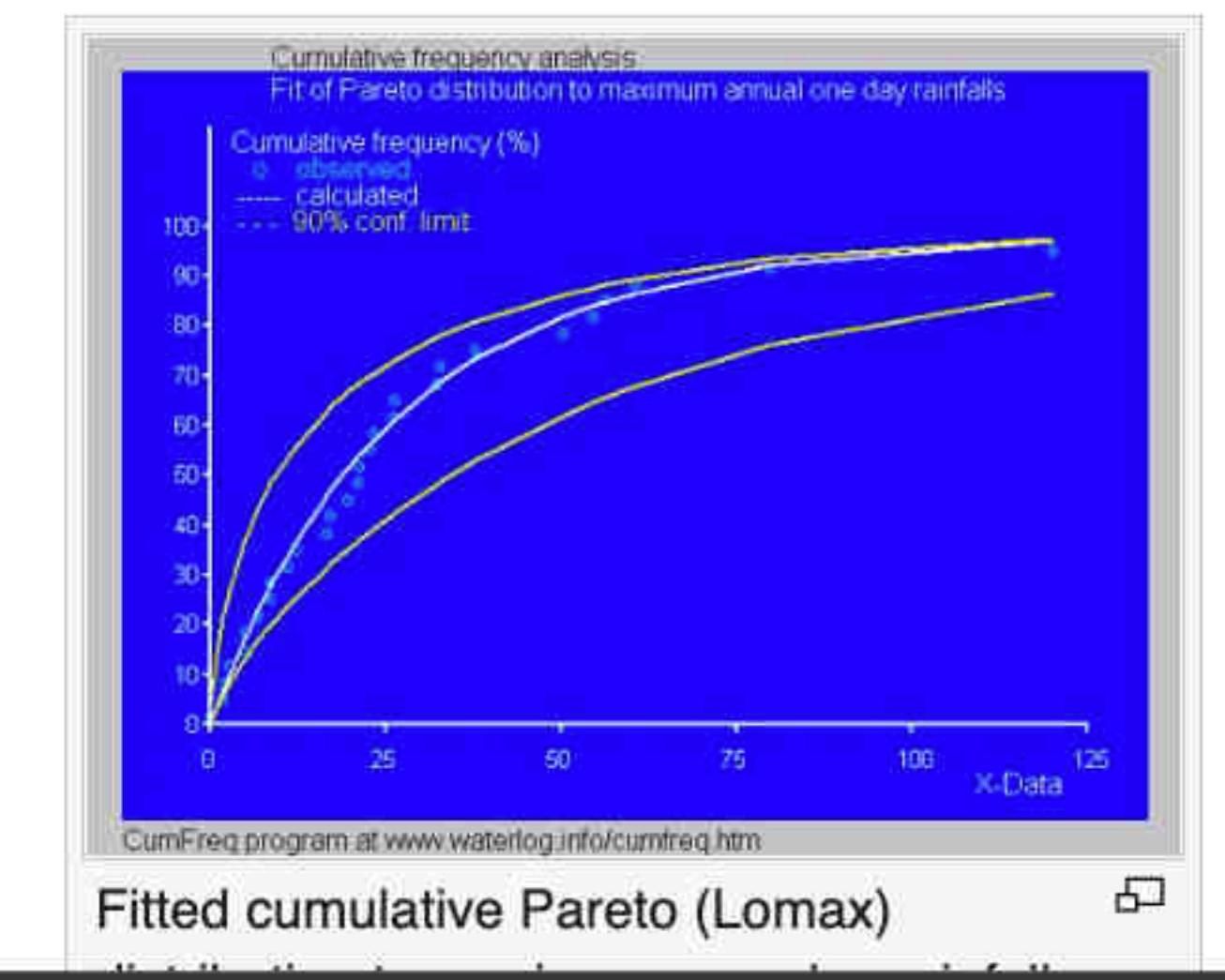
Pareto Type I probability density functions for various  $\alpha$  with  $x_m = 1$ . As  $\alpha \rightarrow \infty$ , the distribution approaches  $\delta(x - x_m)$ , where  $\delta$  is the Dirac delta function.

### Cumulative distribution function



way that a larger portion of the wealth of any society is owned by a smaller percentage of the people describe distribution of income.<sup>[4]</sup> This idea is sometimes expressed more simply as the [Pareto principle](#) or the "80-20 rule" which says that 20% of the population controls 80% of the wealth.<sup>[23]</sup> However, the 80-20 rule corresponds to a particular value of  $a$ , and in fact, Pareto's data on British income taxes in his *Cours d'économie politique* indicates that about 30% of the population had about 70% of the income.  
[\[citation needed\]](#) The [probability density function](#) (PDF) graph at the beginning of this article shows that the "probability" or fraction of the population that owns a small amount of wealth per person is rather high, and then decreases steadily as wealth increases. (The Pareto distribution is not realistic for wealth for the lower end, however. In fact, [net worth](#) may even be negative.) This distribution is not limited to describing wealth or income, but to many situations in which an equilibrium is found in the distribution of the "small" to the "large". The following [examples](#) are sometimes seen as approximately Pareto-distributed:

- The sizes of human settlements (few cities, many hamlets/villages)<sup>[24][25]</sup>
- File size distribution of Internet traffic which uses the TCP protocol (many smaller files, few larger ones)<sup>[24]</sup>
- Hard disk drive error rates<sup>[26]</sup>
- Clusters of Bose–Einstein condensate near absolute zero<sup>[27]</sup>
- The values of oil reserves in oil fields (a few large fields, many small fields)<sup>[24]</sup>
- The length distribution in jobs assigned to supercomputers (a few large ones, many small ones)<sup>[28]</sup>
- The standardized price returns on individual stocks<sup>[24]</sup>
- Sizes of sand particles<sup>[24]</sup>
- The size of meteorites
- Severity of large casualty losses for certain lines of business such as general liability, commercial auto, and workers compensation.<sup>[29][30]</sup>
- Amount of time a user on Steam will spend playing different games. (Some games get played a lot, but most get played almost never.)<sup>[21]</sup> [original research?]

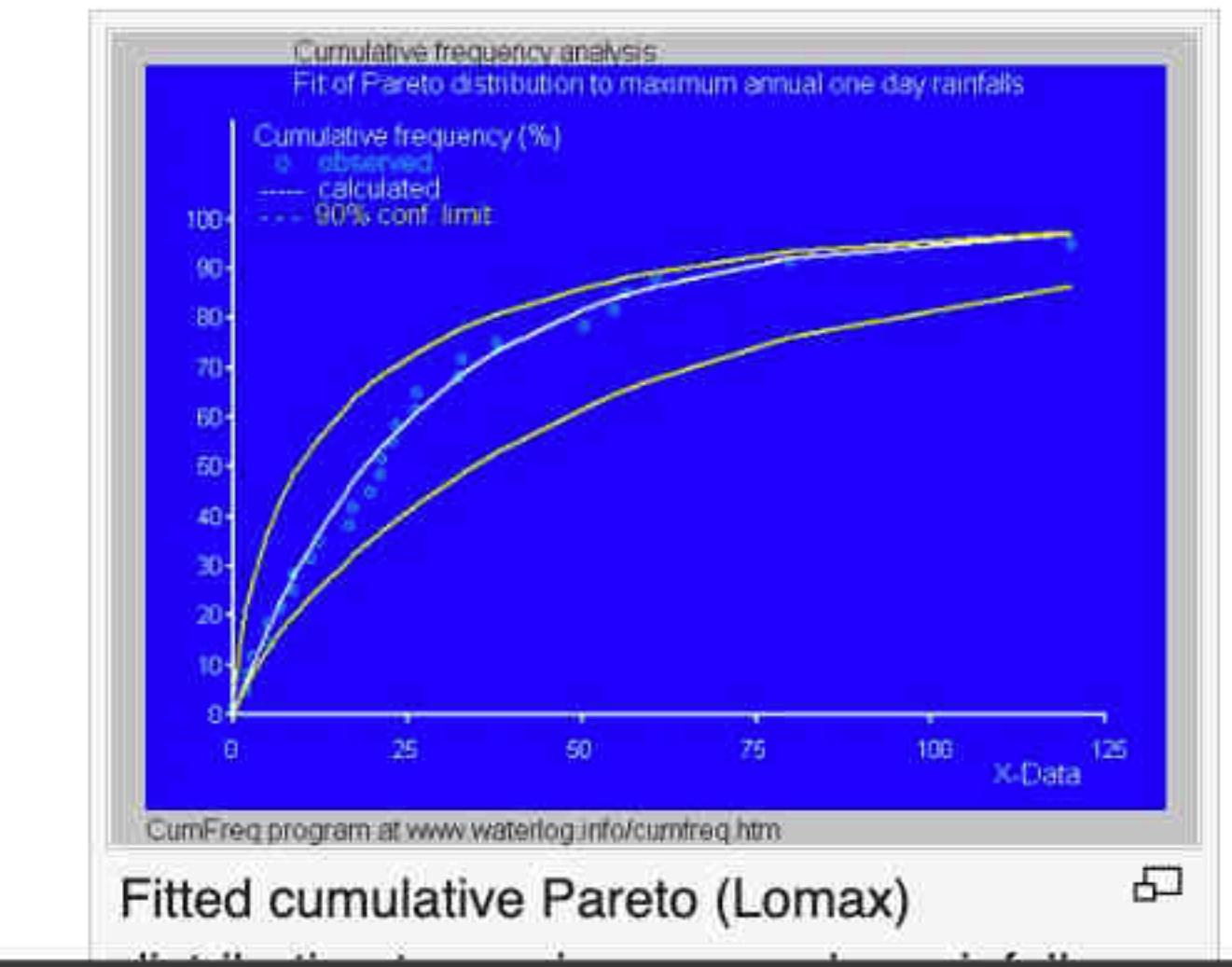


$\mu, \sigma$

{ Pareto dist  
 $d \leq 2$   
 $r = \infty$

way that a larger portion of the wealth of any society is owned by a smaller percentage of the people describe distribution of income.<sup>[4]</sup> This idea is sometimes expressed more simply as the [Pareto principle](#) or the "80-20 rule" which says that 20% of the population controls 80% of the wealth.<sup>[23]</sup> However, the 80-20 rule corresponds to a particular value of  $a$ , and in fact, Pareto's data on British income taxes in his *Cours d'économie politique* indicates that about 30% of the population had about 70% of the income.  
[\[citation needed\]](#) The [probability density function](#) (PDF) graph at the beginning of this article shows that the "probability" or fraction of the population that owns a small amount of wealth per person is rather high, and then decreases steadily as wealth increases. (The Pareto distribution is not realistic for wealth for the lower end, however. In fact, [net worth](#) may even be negative.) This distribution is not limited to describing wealth or income, but to many situations in which an equilibrium is found in the distribution of the "small" to the "large". The following [examples](#) are sometimes seen as approximately Pareto-distributed:

- The sizes of human settlements (few cities, many hamlets/villages)<sup>[24][25]</sup>
- File size distribution of Internet traffic which uses the TCP protocol (many smaller files, few larger ones)<sup>[24]</sup>
- [Hard disk drive error rates](#)<sup>[26]</sup>
- Clusters of Bose–Einstein condensate near absolute zero<sup>[27]</sup>
- The values of [oil reserves](#) in oil fields (a few [large fields](#), many [small fields](#))<sup>[24]</sup>
- The length distribution in jobs assigned to supercomputers (a few large ones, many small ones)<sup>[28]</sup>
- The standardized price returns on individual stocks<sup>[24]</sup>
- Sizes of sand particles<sup>[24]</sup>
- The size of meteorites
- Severity of large [casualty](#) losses for certain lines of business such as general liability, commercial auto, and workers compensation.<sup>[29][30]</sup>
- Amount of time a user on [Steam](#) will spend playing different games. (Some games get played a lot, but most get played almost never.)<sup>[21]</sup> [\[original research?\]](#)

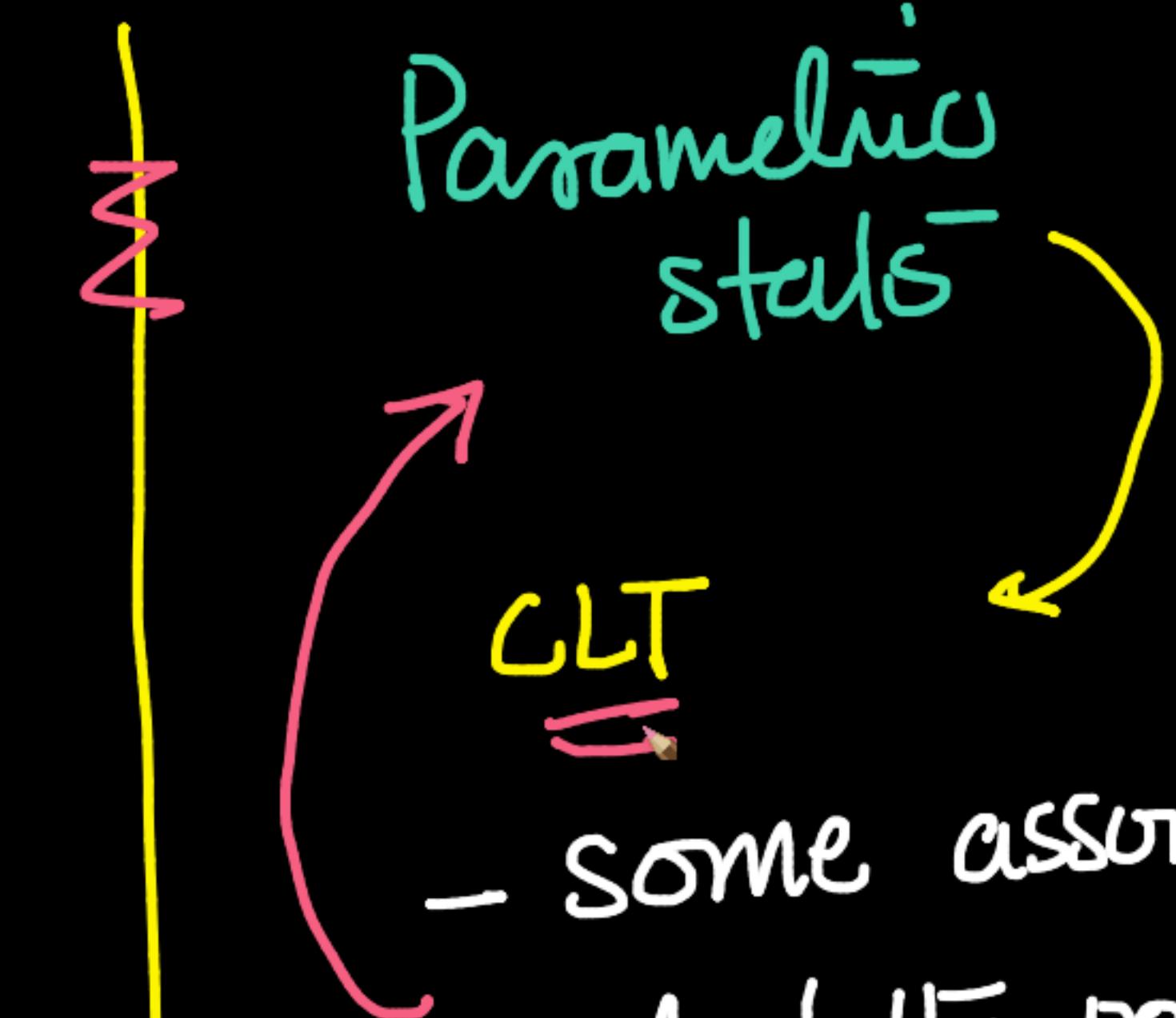




Non-parametric stats

Bootstrapping

{ we made no assumptions  
about underlying population



Parametric stats

CLT

- some assumptions  
about the population  
CLT: ( $\mu$  &  $\sigma$  are finite)

8 colab.research.google.com/drive/1rn1Wz-7i9WknfREQY\_ncl\_EschmEAAJBiW#scrollTo=ZmAP

Pareto distribution – Wikipedia

1

+ Code + Text

RAM Disk

1

```
-rw-r--r-- 1 root root 1825101 May 12 16:05 uber_dataset.z  
drwxr-xr-x 2 root root 4096 May 12 16:04 MACOSX
```

```
import pandas as pd

df = pd.read_csv("./uber_travel_data.csv")
df.sample(100).head()
```

| sourceid | source | dstid | destination   | travel_time |
|----------|--------|-------|---|-------------|
| 2936366  | 187    | 77    | 0 Mehrauli - Badarpur Road, Tughlakabad Instit...   | 4987        |
| 982706   | 63     | 172   | Jawaharlal Nehru Stadium Marg, CGO Complex, Pr... nullShiva Road, Pocket 14, Sector 8D, Rohini, ... | 3093        |
| 3711028  | 234    | 259   | 113, Press Colony, Press Colony, Mayapuri, New...   | Nan         |
| 114690   | 7      | 131   | P 13, Baird Place, Delhi Cantonment, New Delhi 4400 Gali Bahuji, Pahari Dhiraj, Sadar Bazaar,...    | 1840        |
| 2886810  | 183    | NaN   | 243   | Nan         |



```
[ ] df.shape
```

QQ-plot and CLT and Bootstrap.ipynb - Colab Research

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=7mABOWWF1cNT

+ Code + Text

-rw-r--r-- 1 root root 182510 May 12 16:03 uber\_dataset.zip  
drwxr-xr-x 2 root root 4096 May 12 16:04 \_\_MACOSX

RAM Disk

import pandas as pd

df = pd.read\_csv("./uber\_travel\_data.csv")  
df.sample(100).head()

|         | sourceid | source  | dstid | destination                                       | travel_time |
|---------|----------|---|-------|---|-------------|
| 2936366 | 187      | 60, Vijay Vihar Phase II, Pocket C, Sector 1, ... | 77    | 0 Mehrauli - Badarpur Road, Tughlakabad Instit... | 4987        |
| 982706  | 63       | Jawaharlal Nehru Stadium Marg, CGO Complex, Pr... | 172   | nullShiva Road, Pocket 14, Sector 8D, Rohini, ... | 3093        |
| 3711028 | 234      | 113, Press Colony, Press Colony, Mayapuri, New... | 259   |   | NaN         |
| 114690  | 7        | P 13, Baird Place, Delhi Cantonment, New Delhi    | 131   | 4400 Gali Bahuji, Pahari Dhiraj, Sadar Bazaar,... | 1840        |
| 2886810 | 183      |   | NaN   | 243   | NaN         |

[ ] df.shape

14 / 14

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=M8Qmeb2C8IcH

+ Code + Text RAM Disk

[ 6 ] 114690 7 P 13, Baird Place, Delhi Cantonment, New Delhi 131 4400 Gali Bahuji, Pahari Dhiraj, Sadar Bazaar, ... 1840

2886810 183 NaN 243 NaN 2715

{x}

df.shape

(4542026, 5)

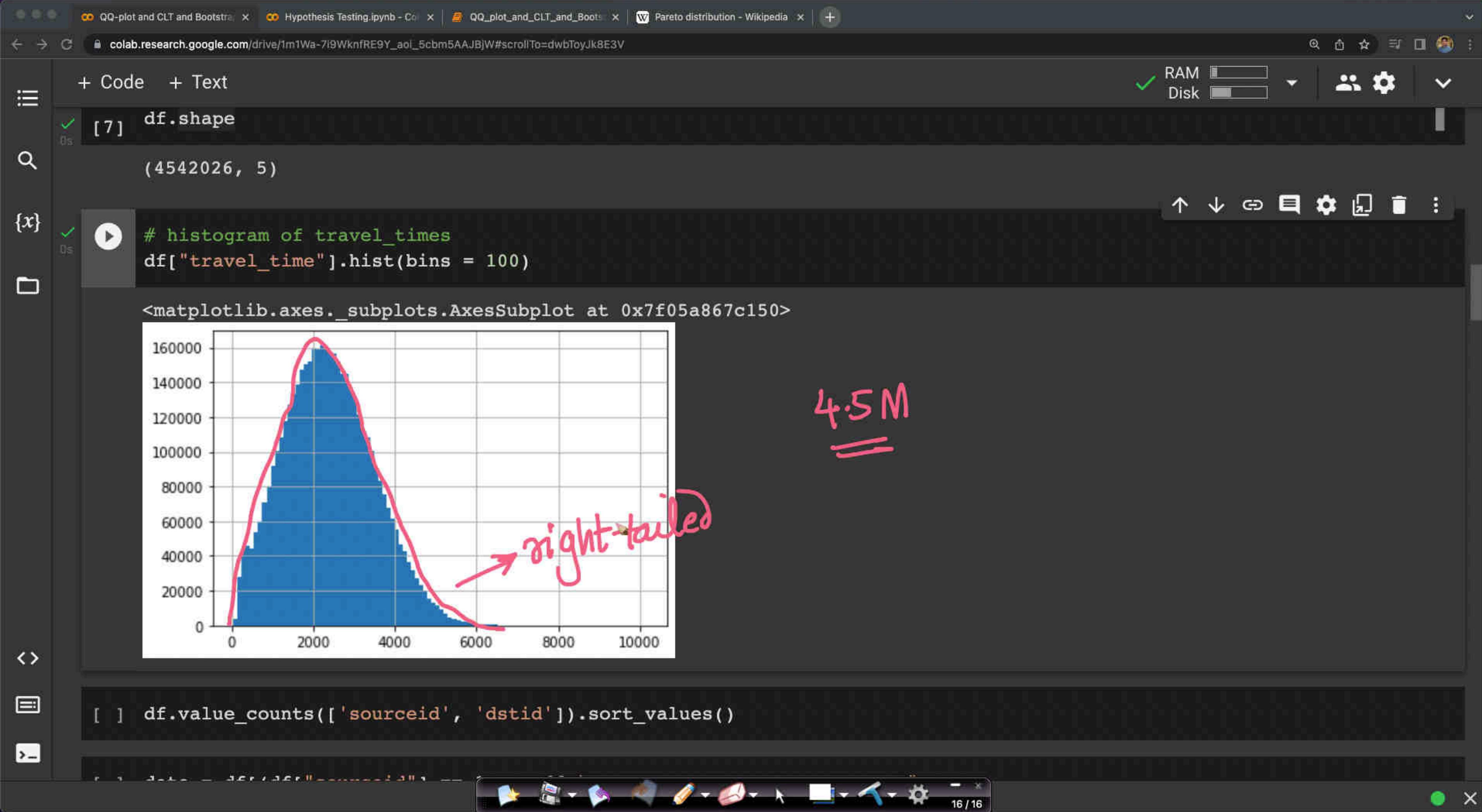
[ ] # histogram of travel\_times  
df["travel\_time"].hist(bins = 100)

[ ] df.value\_counts(['sourceid', 'dstid']).sort\_values()

[ ] data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel\_time"]  
data.shape

[ ] data.hist(bins = 100)

CLT for C.I on mean of travel time



QQ-plot and CLT and Bootstrap x Hypothesis Testing.ipynb - Colab Notebook x QQ\_plot\_and\_CLT\_and\_Bootstrap x Pareto distribution - Wikipedia x +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=dwbToyJk8E3V

+ Code + Text

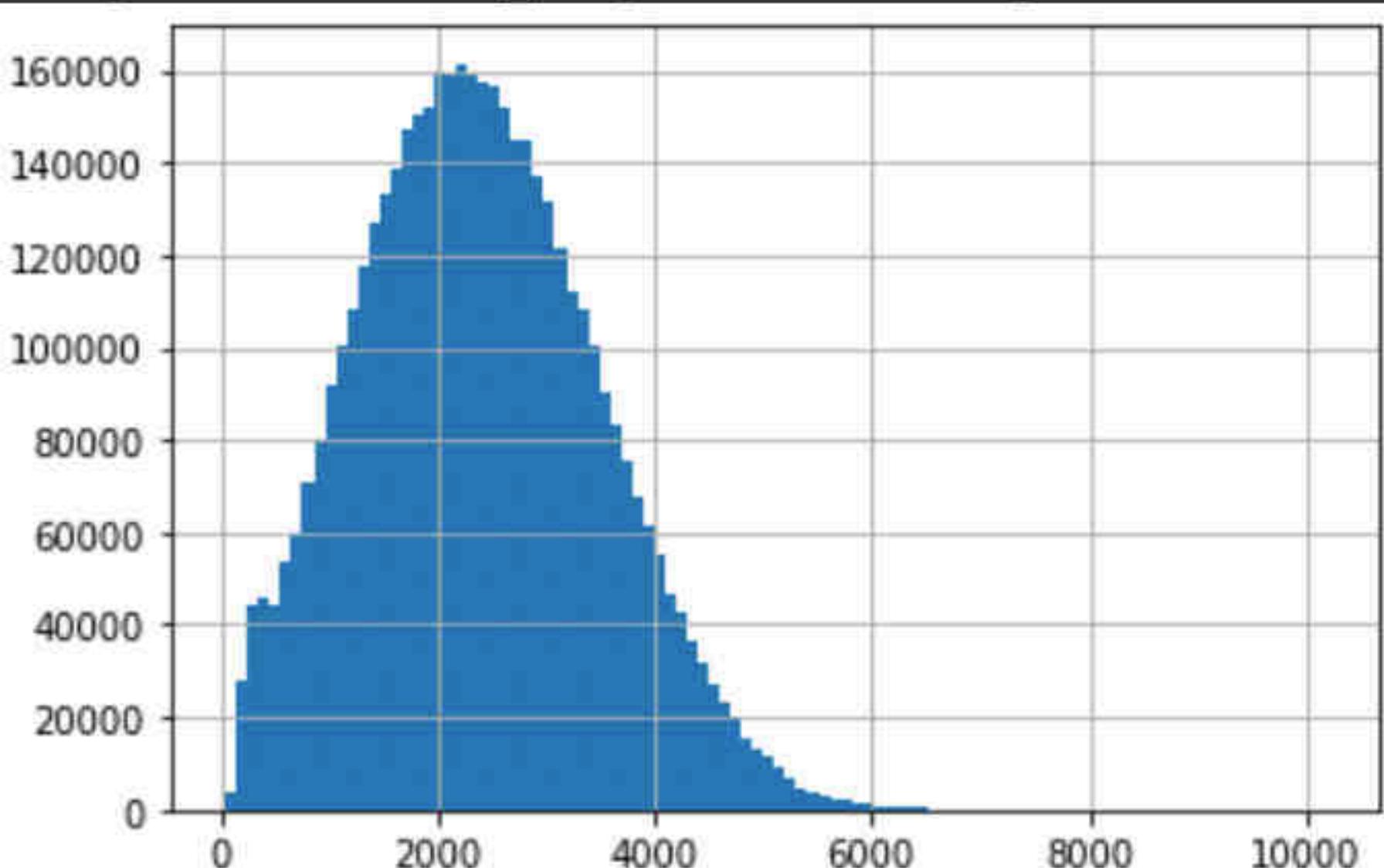
[ 7 ] df.shape

(4542026, 5)

{x}

# histogram of travel\_times  
df["travel\_time"].hist(bins = 100)

`<matplotlib.axes._subplots.AxesSubplot at 0x7f05a867c150>`



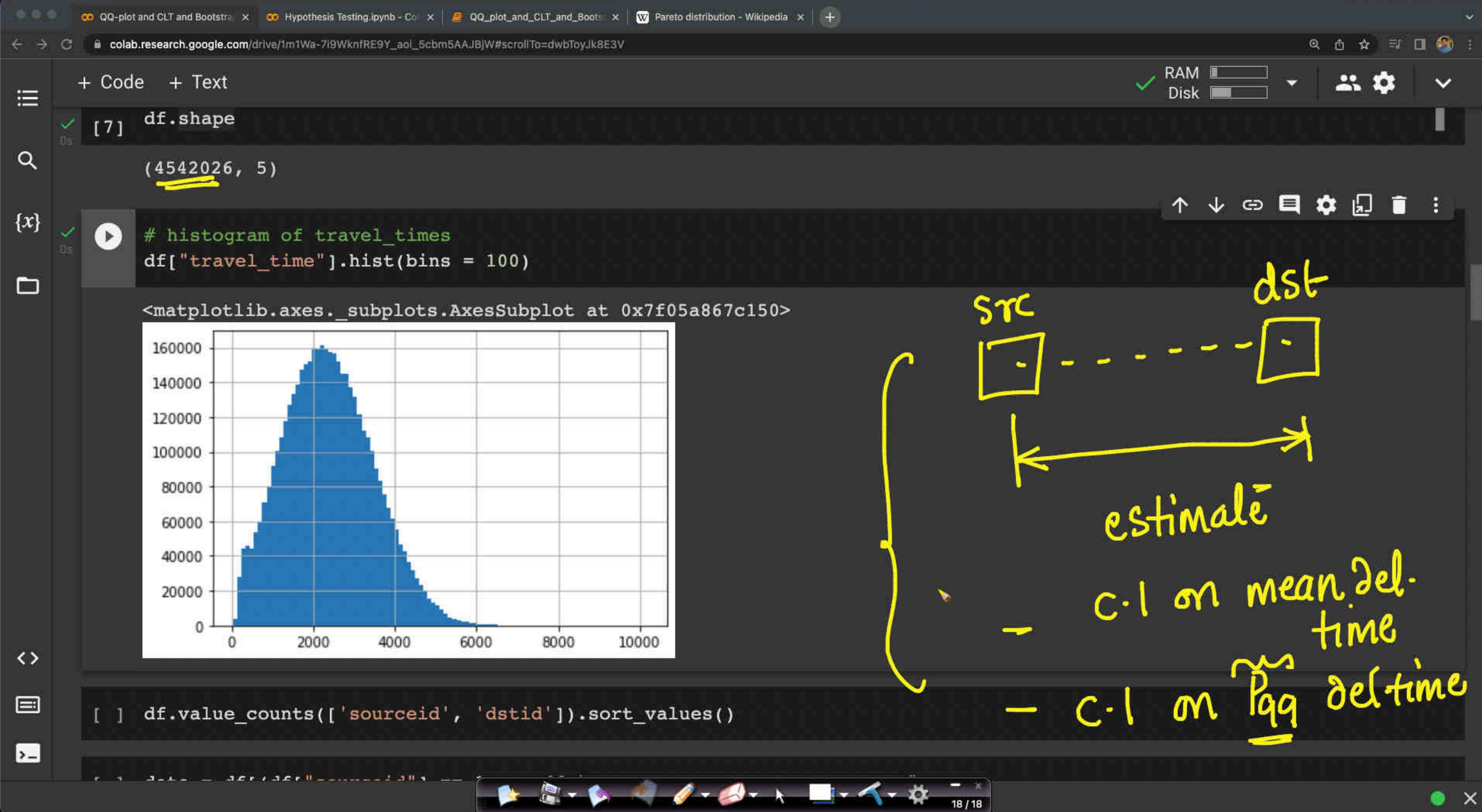
A histogram of travel times with 100 bins. The x-axis ranges from 0 to 10,000, and the y-axis ranges from 0 to 160,000. The distribution is highly right-skewed, with the highest frequency bin around 2,000 travel time units.

finite  $\mu$  &  $\sigma$ : population

[ ] df.value\_counts(['sourceid', 'dstid']).sort\_values()

RAM Disk

17 / 17



QQ-plot and CLT and Bootstrap x Hypothesis Testing.ipynb - Colab Notebook x QQ\_plot\_and\_CLT\_and\_Bootstrap x Pareto distribution - Wikipedia x +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=dwbToyJk8E3V

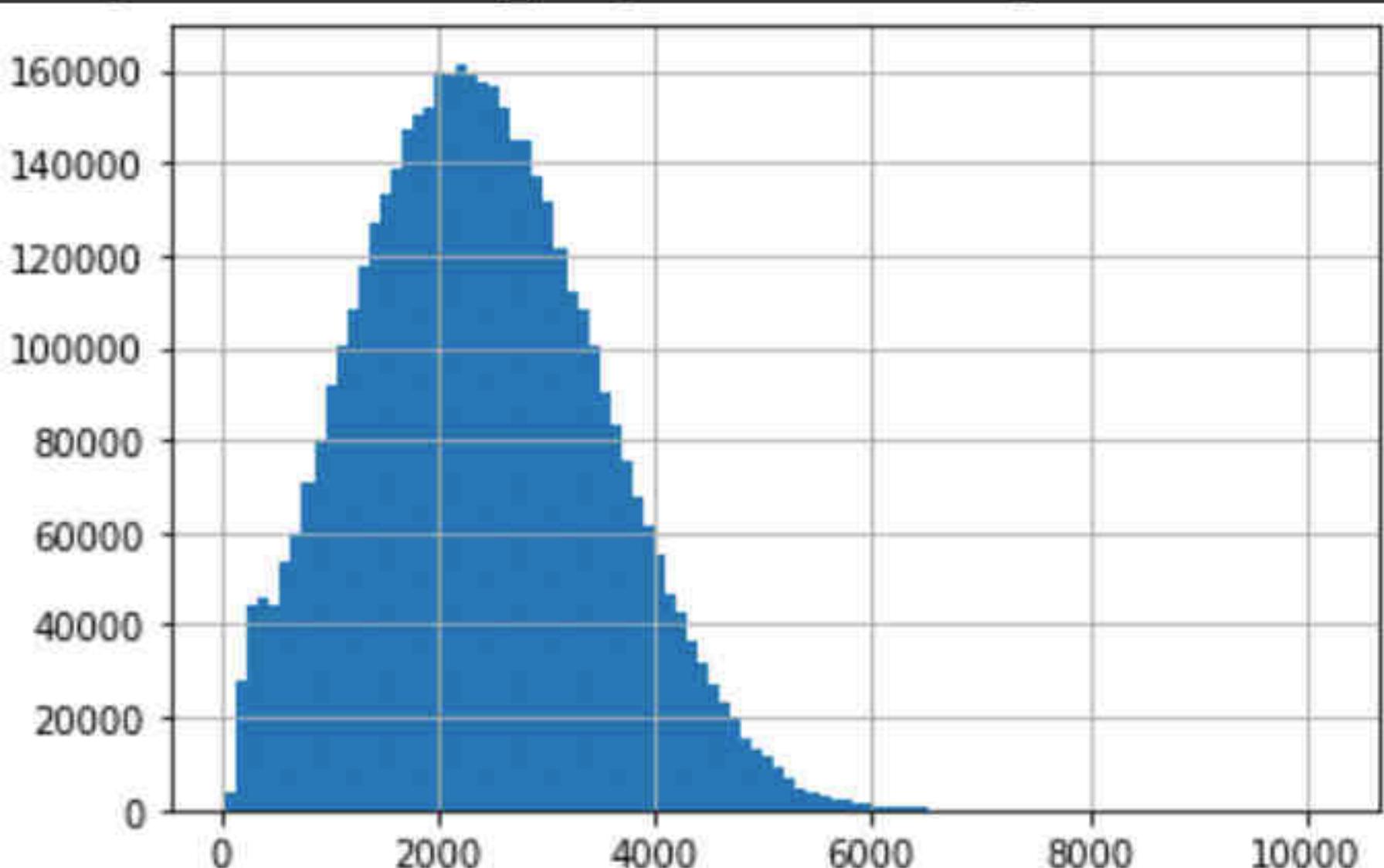
+ Code + Text

[ 7 ] df.shape

(4542026, 5)

{x} # histogram of travel\_times  
df["travel\_time"].hist(bins = 100)

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f05a867c150>



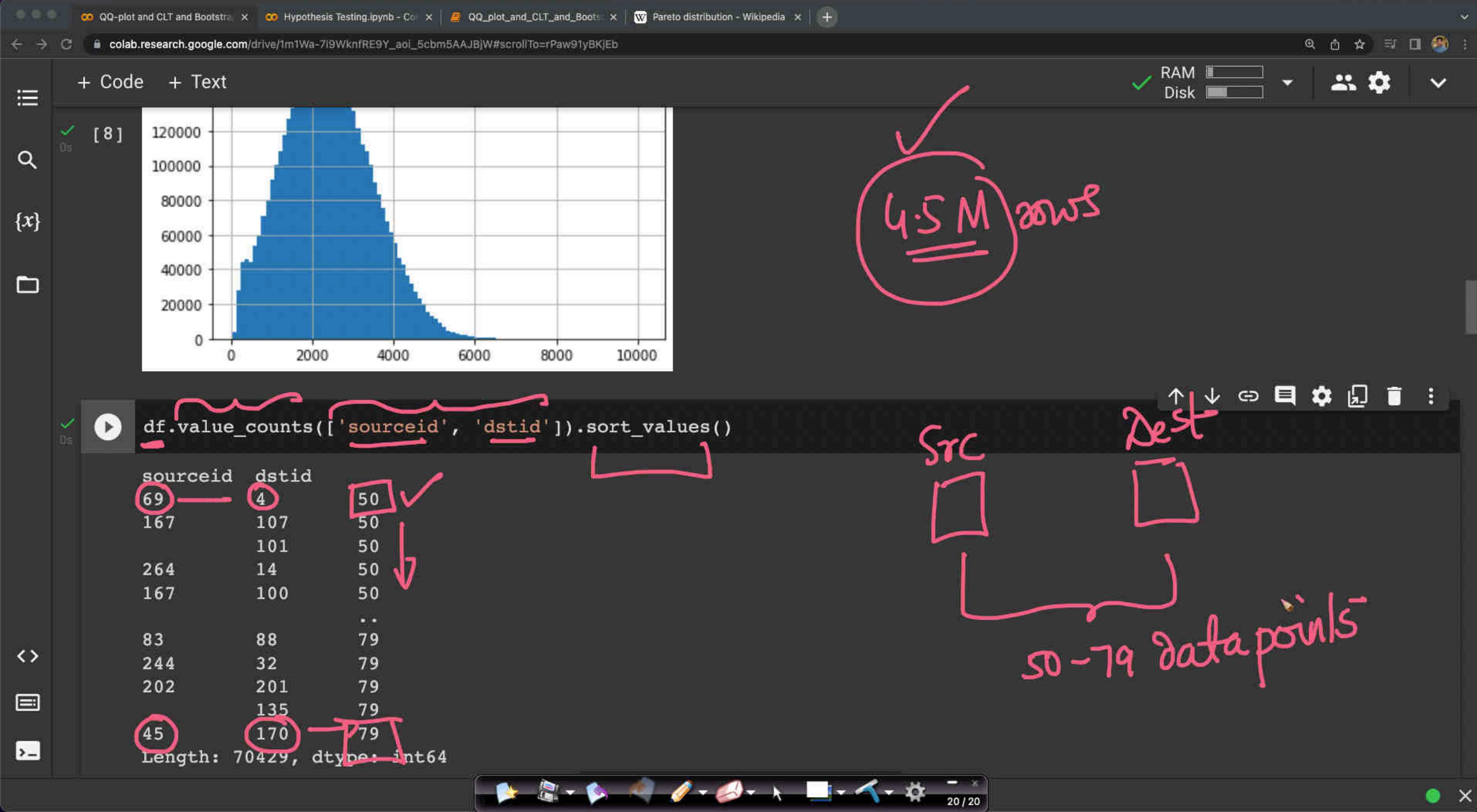
4.5 Million travel-times

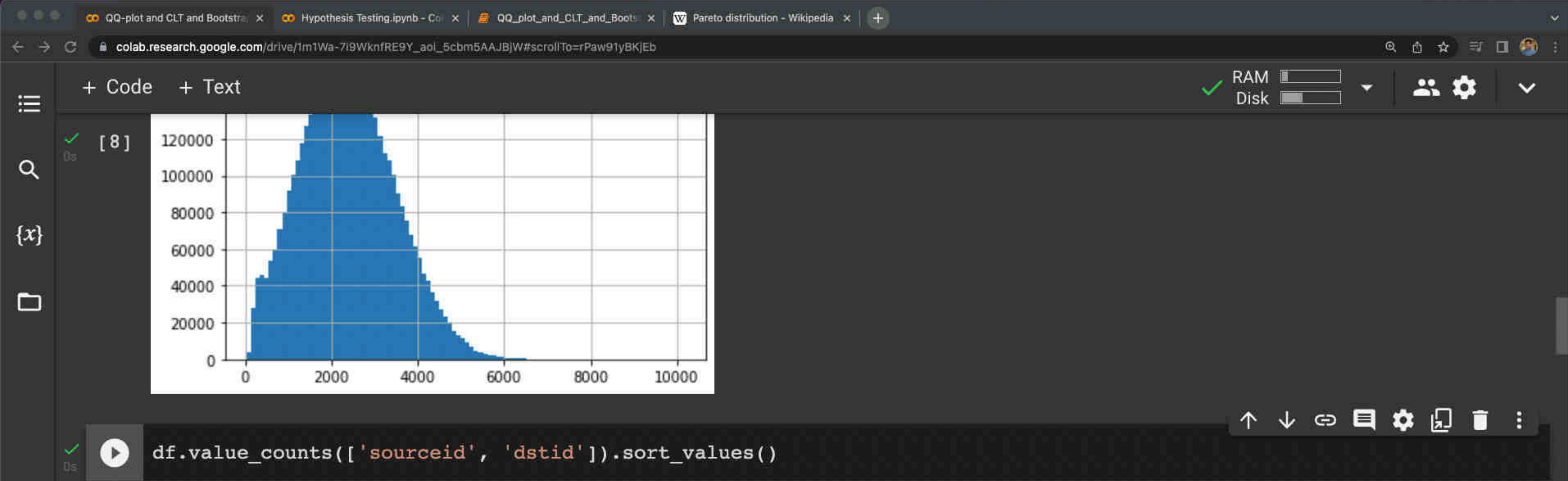
approx it as pop-disb

[ ] df.value\_counts(['sourceid', 'dstid']).sort\_values()

RAM Disk

19 / 19





```
sourceid  dstid
69         4      50
167        107     50
              101     50
264        14      50
167        100     50
              ..
83         88      79
244        32      79
202        201     79
              135     79
45         170     79
Length: 70429, dtype: int64
```

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=OgEYaa0eKUMI

+ Code + Text

RAM Disk

45 170 79  
Length: 70429, dtype: int64

{x} `[ ] data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)]["travel_time"]`

src=1 dst=5  
75 obs

[ ] data.shape  
(75, )

[ ] data.hist(bins = 100)

▶ CLT for C.I on mean of travel\_time

[ ] ↳ 6 cells hidden

▶ 95% C.I on 99th percentile value for travel\_time via bootstrapping

[ ] ↳ 7 cells hidden

▶ box-cox transform

22 / 22

QQ-plot and CLT and Bootstrap x Hypothesis Testing.ipynb - Colab Notebook x QQ\_plot\_and\_CLT\_and\_Bootstrap x Pareto distribution - Wikipedia x +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=rhun6Z8uLhrm

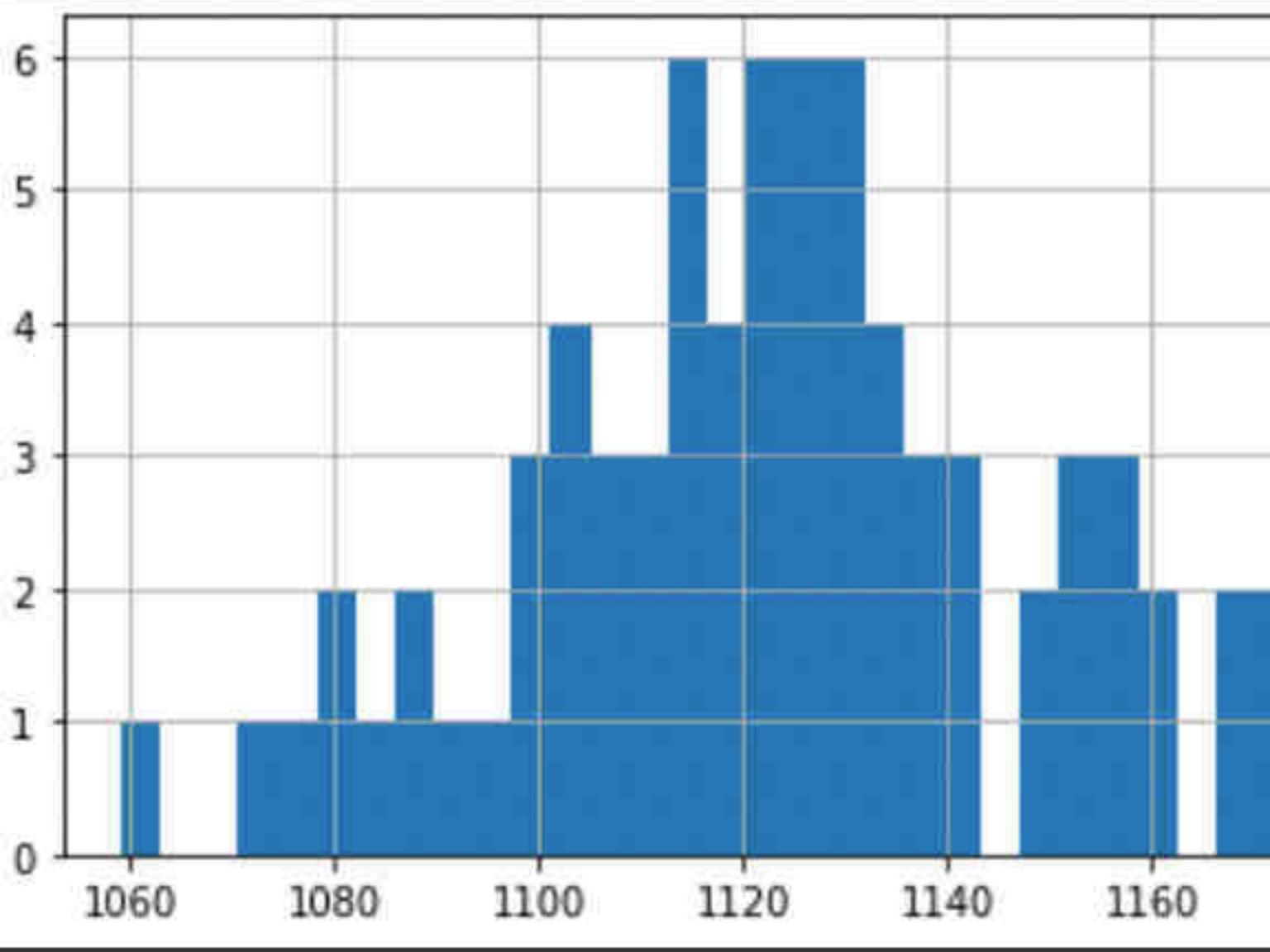
+ Code + Text RAM Disk

[10] data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel\_time"]  
data.shape

{x} (75, ) ✓

data.hist(bins=30)

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f05a5fd56d0>



Handwritten annotations:

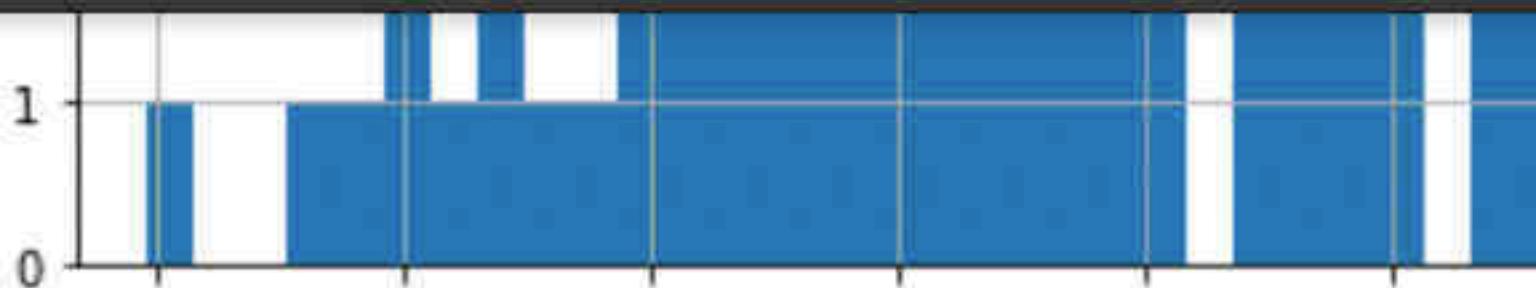
- A yellow circle highlights the output '(75, )' from the code.
- A green box encloses the number '1'.
- A green box encloses the number '5'.
- A green circle encloses the number '75' with the word 'Sample' written above it.
- A green arrow points from the '75' circle to the word 'random'.

CLT for C.I on mean of travel\_time

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=lcDaEZiY\_Nj0

+ Code + Text

[13] 

{x}

CLT for C.I on mean of travel\_time

finished μ & σ

```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)

for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

```
[ ] import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```

QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

+ Code + Text

RAM Disk



## CLT for C.I on mean of travel\_time

{x}

```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
✓ r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)

for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

▶ import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_means, bins=100)  
plt.grid()  
plt.show()



QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

RAM Disk

+ Code + Text

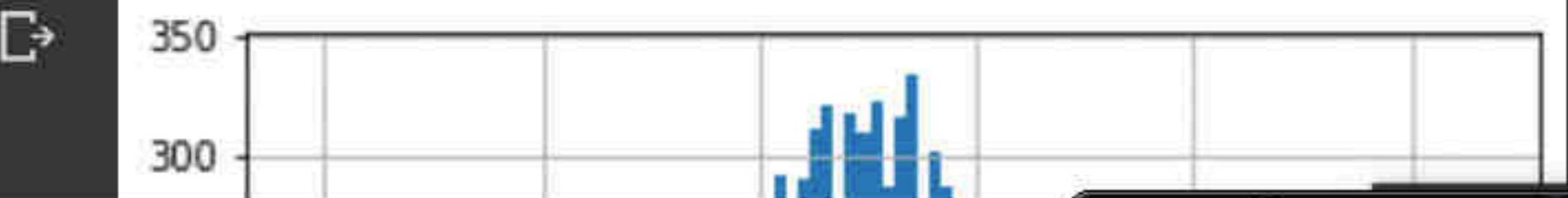
## CLT for C.I on mean of travel\_time

{x}

```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)

for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

▶ import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_means, bins=100)  
plt.grid()  
plt.show()



QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

+ Code + Text

RAM Disk

## CLT for C.I on mean of travel\_time

75 obs

```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
✓ r = 10000
✓ data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50 ✓
bs_means = np.empty(r)

for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

 $\underline{\underline{50}} = n = \text{sample-size}$ 

```
▶ import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```



QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

+ Code + Text

RAM  
Disk

## CLT for C.I on mean of travel\_time

{x}

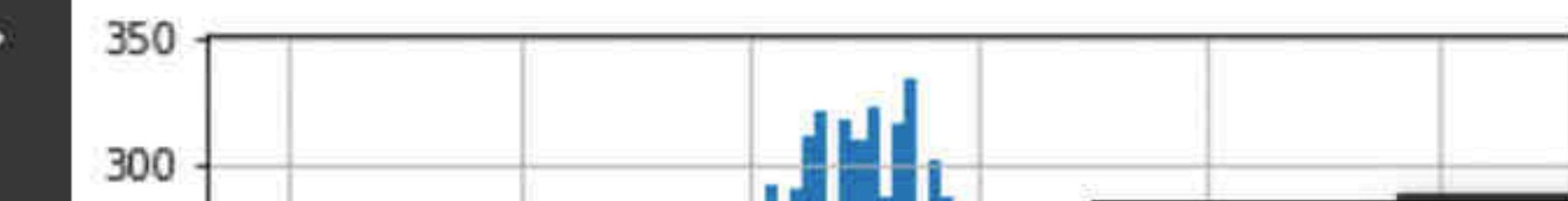
```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)
for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

M<sub>1</sub>M<sub>2</sub>

⋮

M<sub>10,000</sub>

```
▶ import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```



QQ-plot and CLT and Bootstrap.ipynb - Colab | Hypothesis Testing.ipynb - Colab

Hypothesis Testing.ipynb - Colab

QQ\_plot\_and\_GLT\_and\_Bootstrap

- Wikipedia x | +

## - Code + Text

✓ RAM Disk

1

↑ ↓ ⌂ ☰ ✎ ↻ 🗑 ⋮

### CLT for C.I on mean of travel\_time

100

```
[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)

for i in range(r):
    bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

```
▶ import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs_means, bins=100)  
plt.grid()  
plt.show()
```



QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

RAM Disk

+ Code + Text



{x} CLT for C.I on mean of travel\_time

[ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50  
# bs\_means is a list of 'r' bootstrap sample means  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel\_time"]  
size = 50  
bs\_means = np.empty(r)  
  
for i in range(r):  
 bs\_sample = np.random.choice(data, size=size)  
 bs\_means[i] = np.mean(bs\_sample)

<> [ ] import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_means, bins=100)  
plt.grid()  
plt.show()

Sample-means ~ Gaussian

QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

+ Code + Text

✓ RAM  
Disk

User Settings

Up Down Left Right Home End

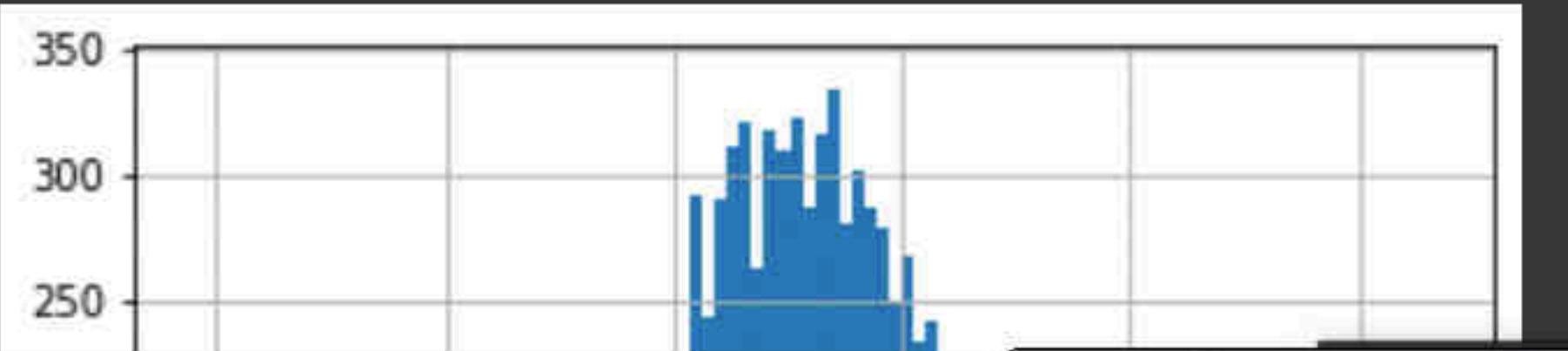
## CLT for C.I on mean of travel\_time

```
{x} [ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50
# bs_means is a list of 'r' bootstrap sample means
r = 10000
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]
size = 50
bs_means = np.empty(r)

✓ for i in range(r):
{   bs_sample = np.random.choice(data, size=size)
    bs_means[i] = np.mean(bs_sample)
```

 $S_1$  $S_2$ 

```
▶ import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```

I  
.  
. $S_Y$ 

QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

+ Code + Text

✓ RAM Disk

User Settings

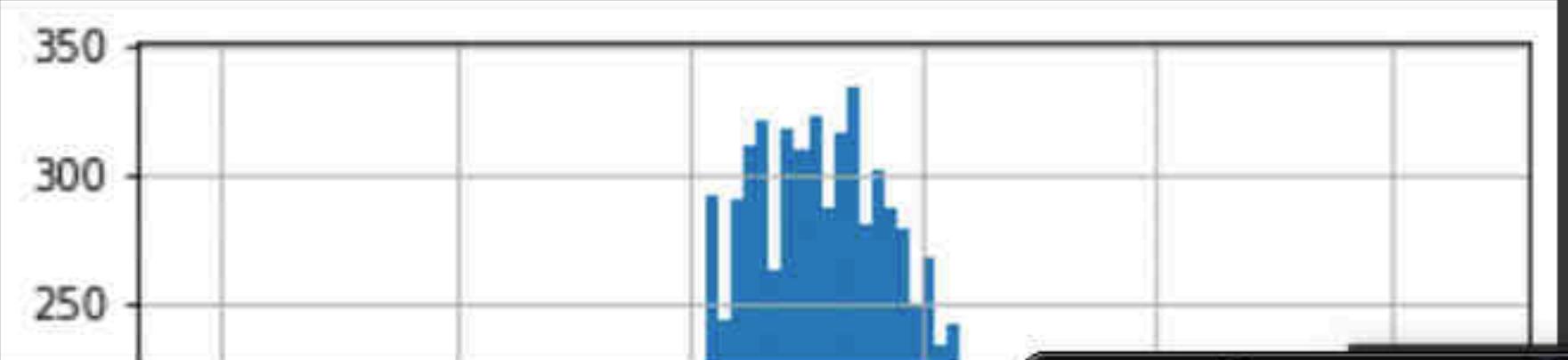
## CLT for C.I on mean of travel\_time

```
{x} [ ] # Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50  
# bs_means is a list of 'r' bootstrap sample means  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 50  
bs_means = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size) →  $S_1, S_2, \dots, S_r = 10,000$   
    bs_means[i] = np.mean(bs_sample)
```

```
▶ import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs_means, bins=100)  
plt.grid()  
plt.show()
```

$S_1, S_2, \dots, S_r = 10,000$

$M_1, M_2, \dots, M_r$



QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

+ Code + Text

RAM Disk

bs\_means = np.empty(r)

for i in range(r):

    bs\_sample = np.random.choice(data, size=size)

    bs\_means[i] = np.mean(bs\_sample)

import matplotlib.pyplot as plt

plt.figure()

plt.hist(bs\_means, bins=100)

plt.grid()

plt.show()

M<sub>1</sub>, M<sub>2</sub> - - - M<sub>k</sub>

350  
300  
250  
200  
150  
100  
50  
0

1110 1115 1120 1125 1130 1135

QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

+ Code

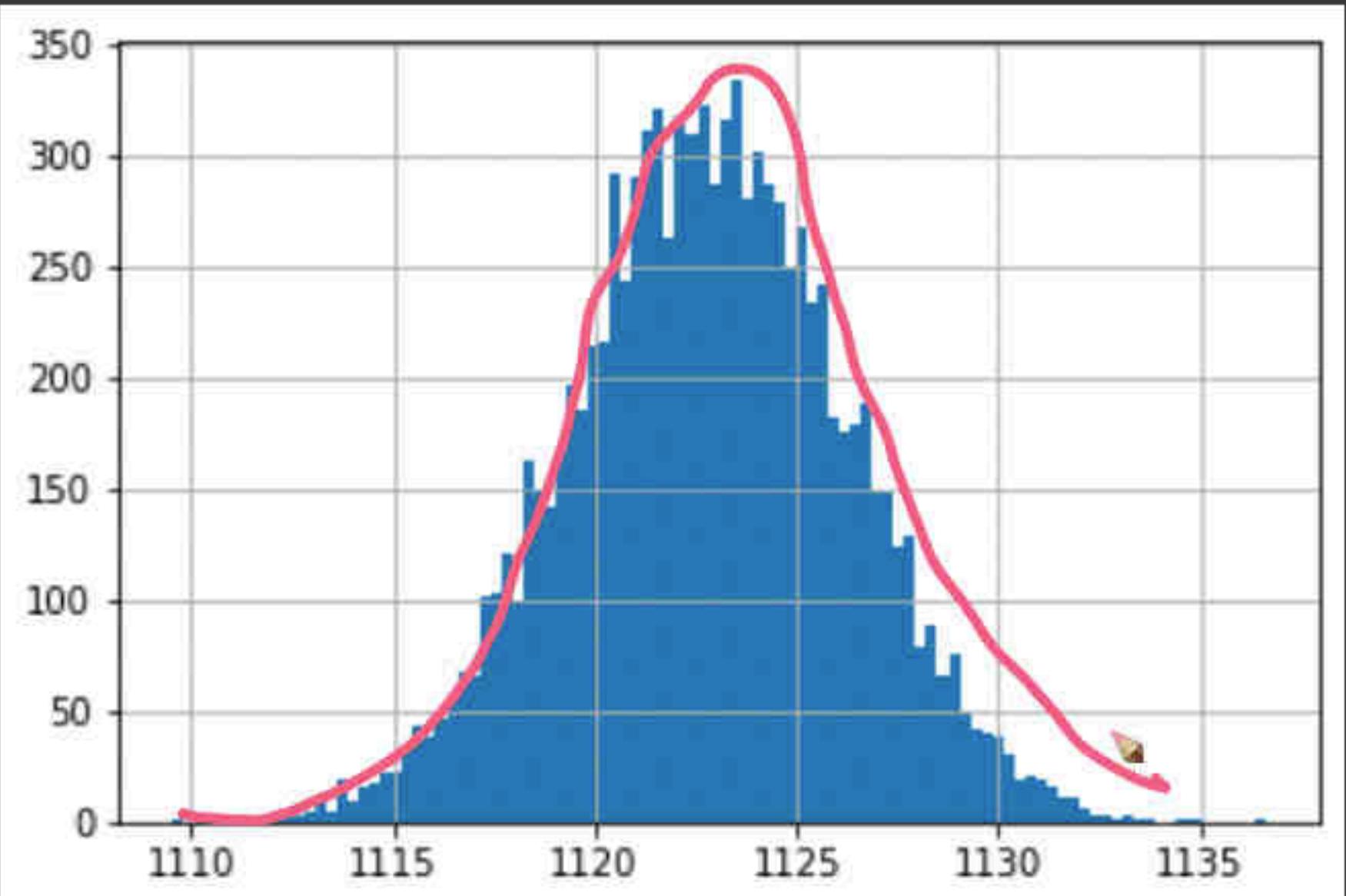


```
bs_sample = np.random.choice(data, size=size)
bs_means[i] = np.mean(bs_sample)
```

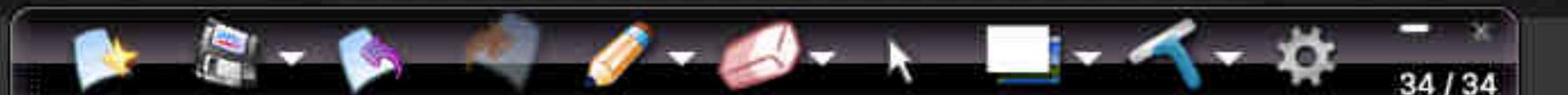
RAM  
Disk

{x}

```
[ ] import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```



```
[ ] # QQ-plot with normal distribution
```



∞ QQ-plot and CLT and Bootstrap × ∞ Hypothesis Testing.ipynb - Colab Notebooks × ∞ QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb × W Pareto distribution - Wikipedia × +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

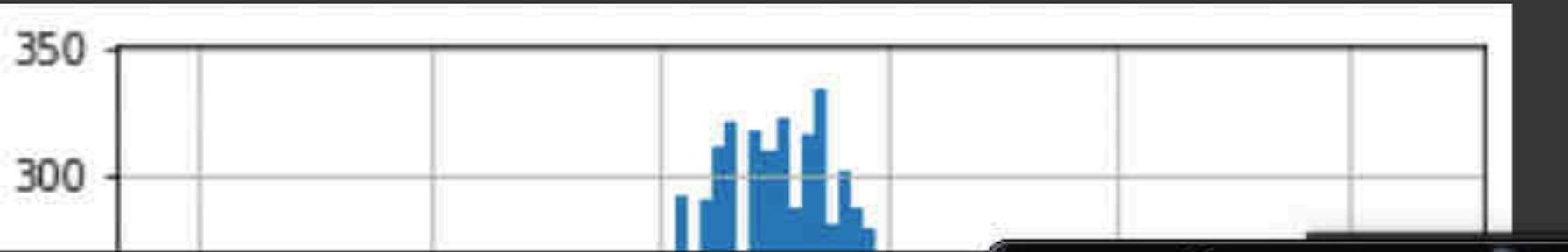
+ Code + Text RAM Disk

CLT for C.I on mean of travel\_time

{x}

# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50  
# bs\_means is a list of 'r' bootstrap sample means  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel\_time"]  
size = 50  
bs\_means = np.empty(r)  
  
for i in range(r):  
 bs\_sample = np.random.choice(data, size=size)  
 bs\_means[i] = np.mean(bs\_sample)

import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_means, bins=100)  
plt.grid()  
plt.show()



350  
300

35 / 35

∞ QQ-plot and CLT and Bootstrap × ∞ Hypothesis Testing.ipynb - Colab Notebook × ∞ QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb × W Pareto distribution - Wikipedia × +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

+ Code + Text RAM Disk

CLT for C.I on mean of travel\_time

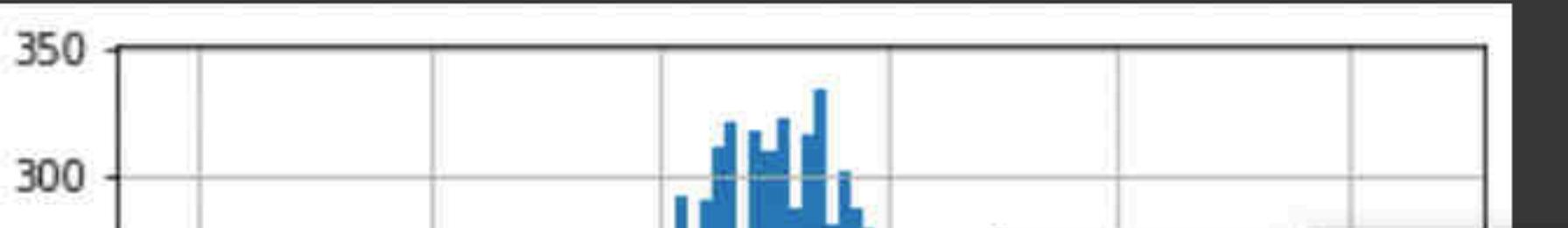
{x}

# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=50  
# bs\_means is a list of 'r' bootstrap sample means  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel\_time"]  
size = 50 ←  
bs\_means = np.empty(r)  
  
for i in range(r):  
 bs\_sample = np.random.choice(data, size=size)  
 bs\_means[i] = np.mean(bs\_sample)

import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_means, bins=100)  
plt.grid()  
plt.show()

CLT:  
Sample-Means  
 $\sim N(\mu; \frac{\sigma}{\sqrt{n}})$

as  $n \rightarrow \infty$



350  
300

36 / 36

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

+ Code + Text

```
[ ] plt.hist(bs_means, bins=100)
plt.grid()
plt.show()
```

{x}

□

350  
300  
250  
200  
150  
100  
50  
0

1110 1115 1120 1125 1130 1135

→ Std-dev =  $\frac{\sigma}{\sqrt{n}}$

$n = 50$

μ

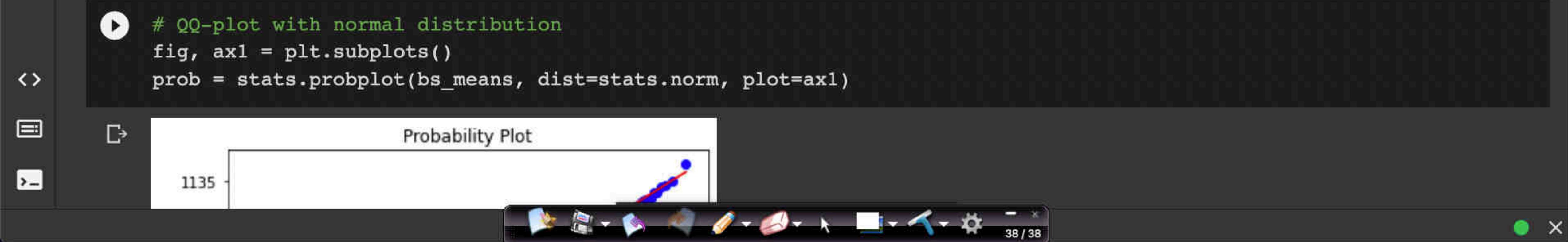
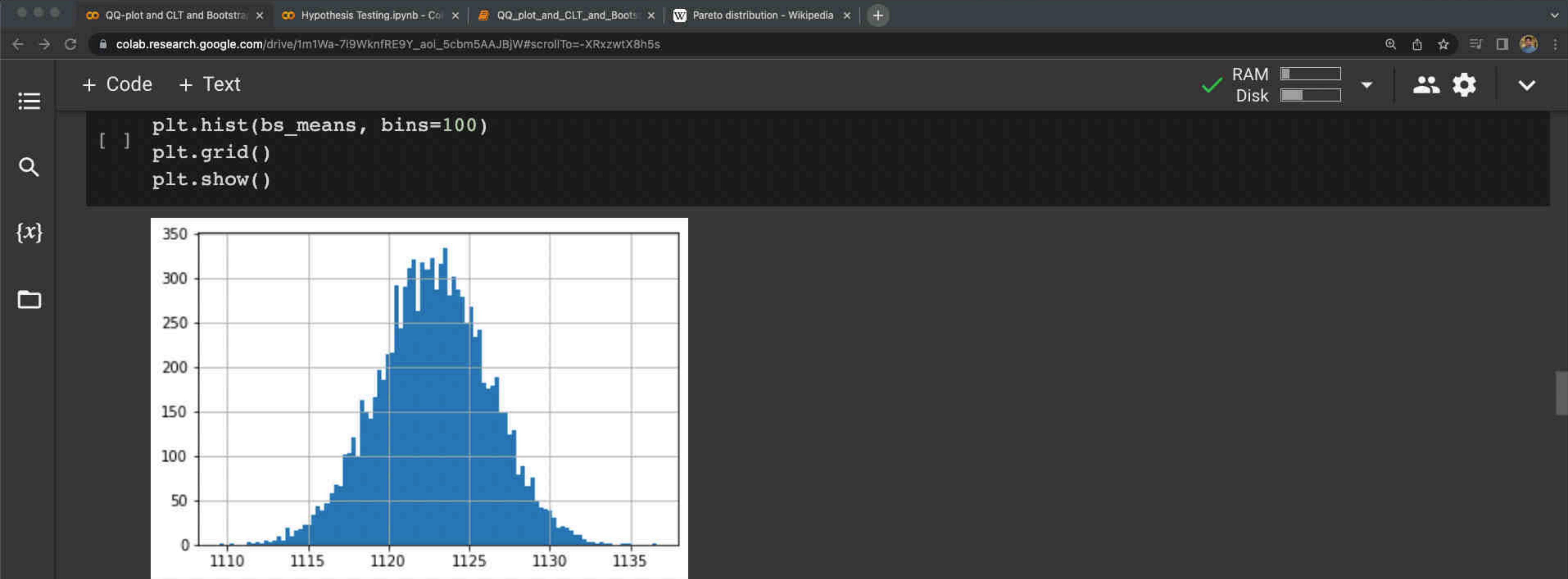
```
▶ # QQ-plot with normal distribution
fig, ax1 = plt.subplots()
prob = stats.probplot(bs_means, dist=stats.norm, plot=ax1)
```

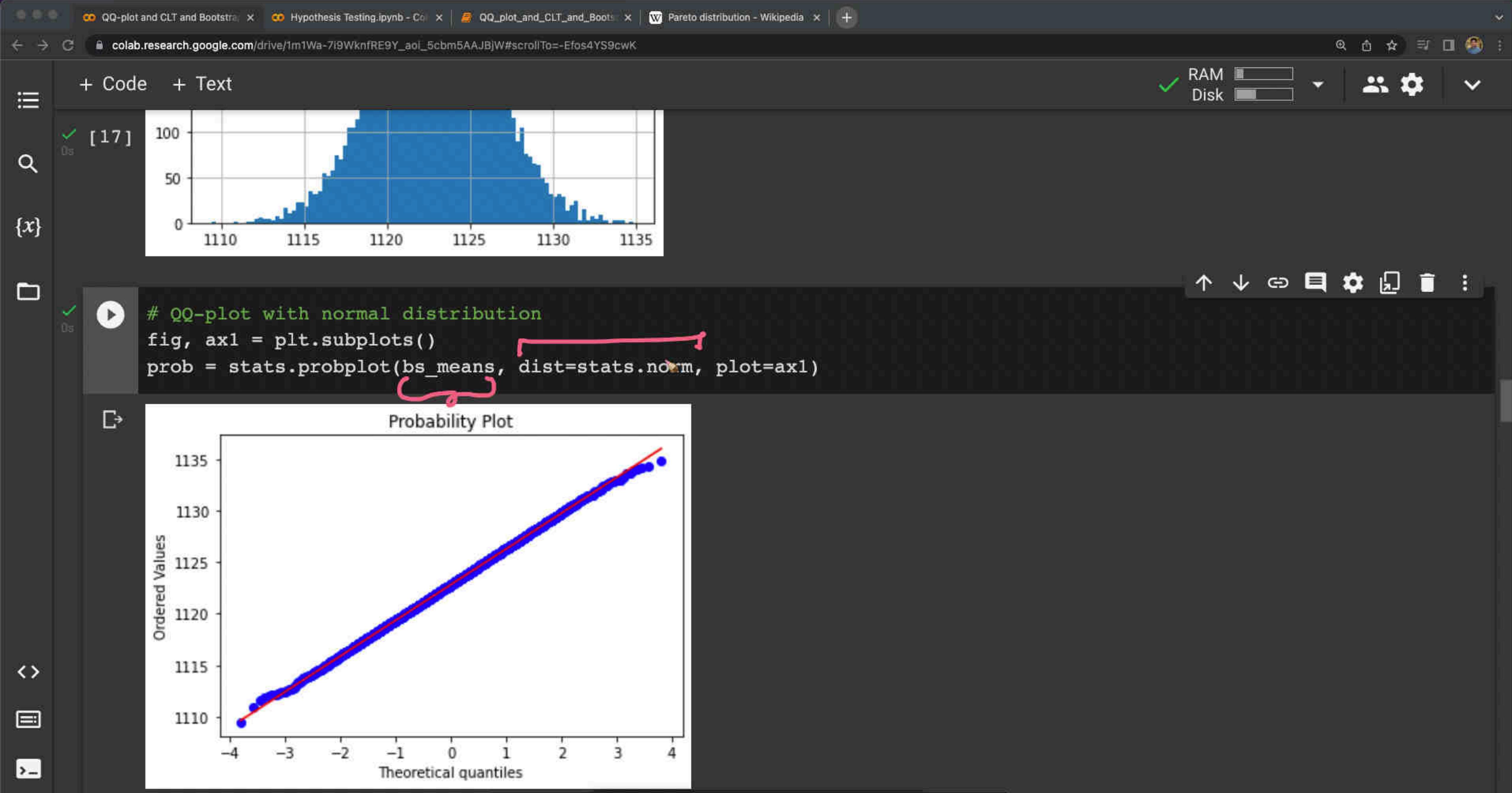
Probability Plot

1135

RAM Disk

37 / 37

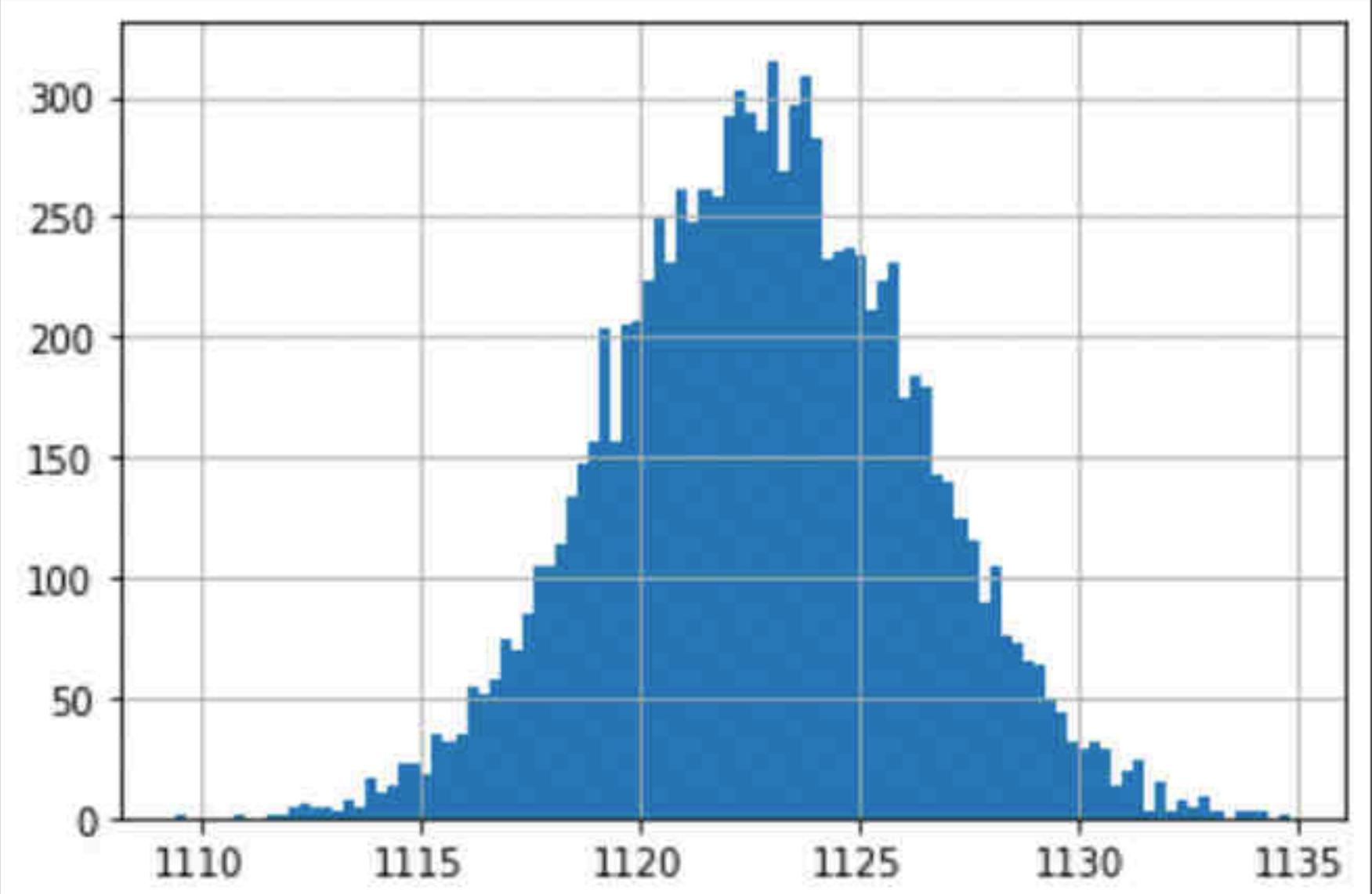




+ Code + Text

RAM Disk

```
[17] import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs_means, bins=100)  
plt.grid()  
plt.show()
```



# disb of Sample-Means

```
# QQ-plot with normal distribution
fig, ax1 = plt.subplots()
prob = stats.probplot(bs_means, dist=stats.norm, plot=ax1)
```

QQ-plot and CLT and Bootstrap.ipynb - Colab

Theoretical quantiles

```
# compute C.I on the mean given that bs_means follows Gaussian distribution: CLT
print(np.mean(bs_means))
print(np.std(bs_means))

1122.823726
3.460407709638272
```

95% C.I on mean

std-dev=3

asym.

```
[ ] print(np.mean(bs_means)-2*np.std(bs_means))
print(np.mean(bs_means)+2*np.std(bs_means))

1115.9224421865092
1129.6414778134913
```

1 5

1122.82 1122.82  
-2x3.41 +2x3.41

```
[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs_means,2.5))
print(np.percentile(bs_means,97.5))

# what if r is say 100 and not 10,000?

1115.98
1129.4205000000002
```

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

Theoretical quantiles

```
# compute C.I on the mean given that bs_means follows Gaussian distribution: CLT
print(np.mean(bs_means))
print(np.std(bs_means))

1122.823726
3.460407709638272
```

```
[ ] print(np.mean(bs_means)-2*np.std(bs_means))
      print(np.mean(bs_means)+2*np.std(bs_means))
```

1115.9224421865092  
1129.6414778134913

qq's C.I ✓

68-95-99 rule ✓

```
[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs_means,2.5))
print(np.percentile(bs_means,97.5))

# what if r is say 100 and not 10,000?
```

```
1115.98
1129.4205000000002
```

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

Theoretical quantiles

```
# compute C.I on the mean given that bs_means follows Gaussian distribution: CLT
print(np.mean(bs_means))
print(np.std(bs_means))

[ ] print(np.mean(bs_means)-2*np.std(bs_means))
print(np.mean(bs_means)+2*np.std(bs_means))

[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs_means,2.5))
print(np.percentile(bs_means,97.5))

# what if r is say 100 and not 10,000?

```

RAM Disk

43 / 44

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

Theoretical quantiles

```
# compute C.I on the mean given that bs_means follows Gaussian distribution: CLT
print(np.mean(bs_means))
print(np.std(bs_means))

1122.823726
3.460407709638272
```

[ ] print(np.mean(bs\_means)-2\*np.std(bs\_means))
print(np.mean(bs\_means)+2\*np.std(bs\_means))

1115.9224421865092
1129.6414778134913

```
[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs_means,2.5))
print(np.percentile(bs_means,97.5))

# what if r is say 100 and not 10,000?
```

1115.98
1129.4205000000002

RAM Disk

95% C.I on the mean (CLT)

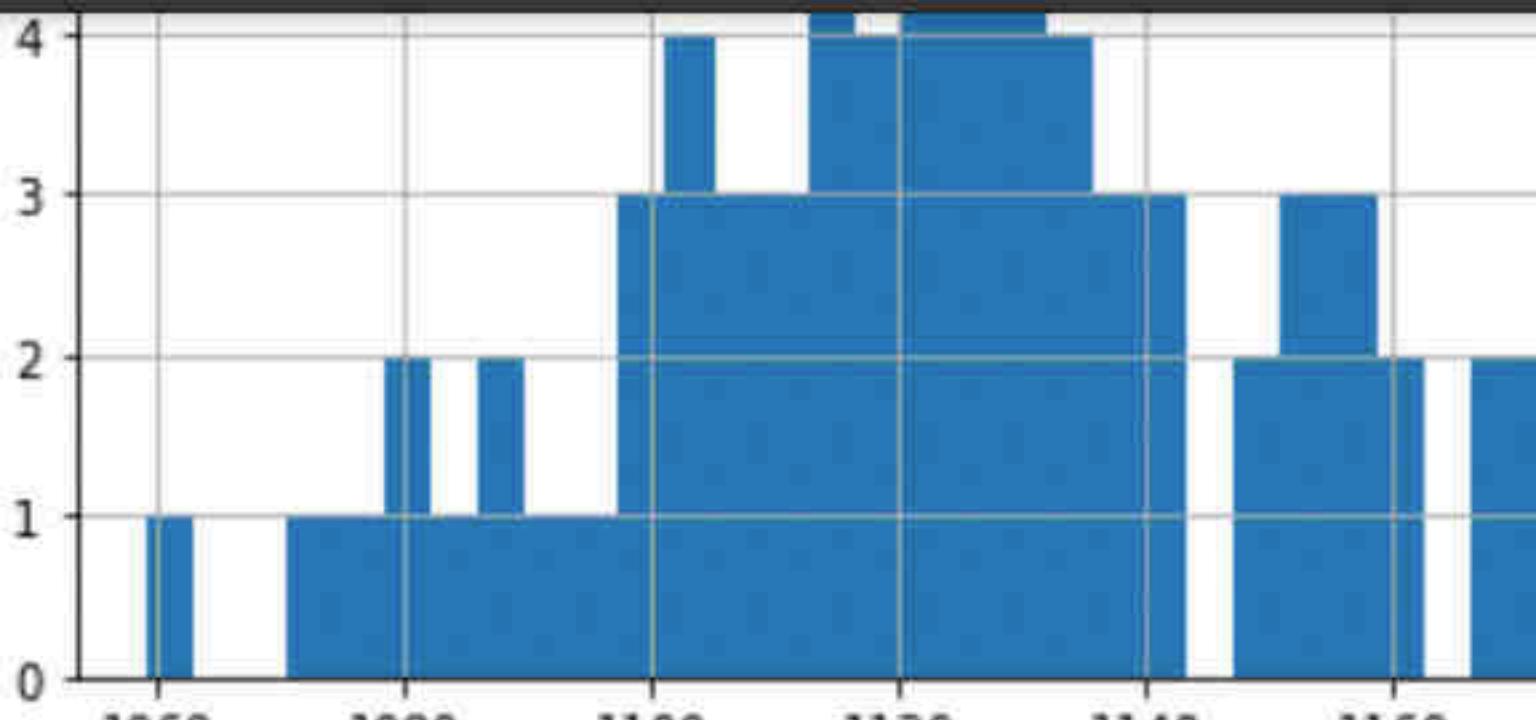
44 / 45

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=sYnwKETg\_VhM

+ Code + Text

RAM Disk

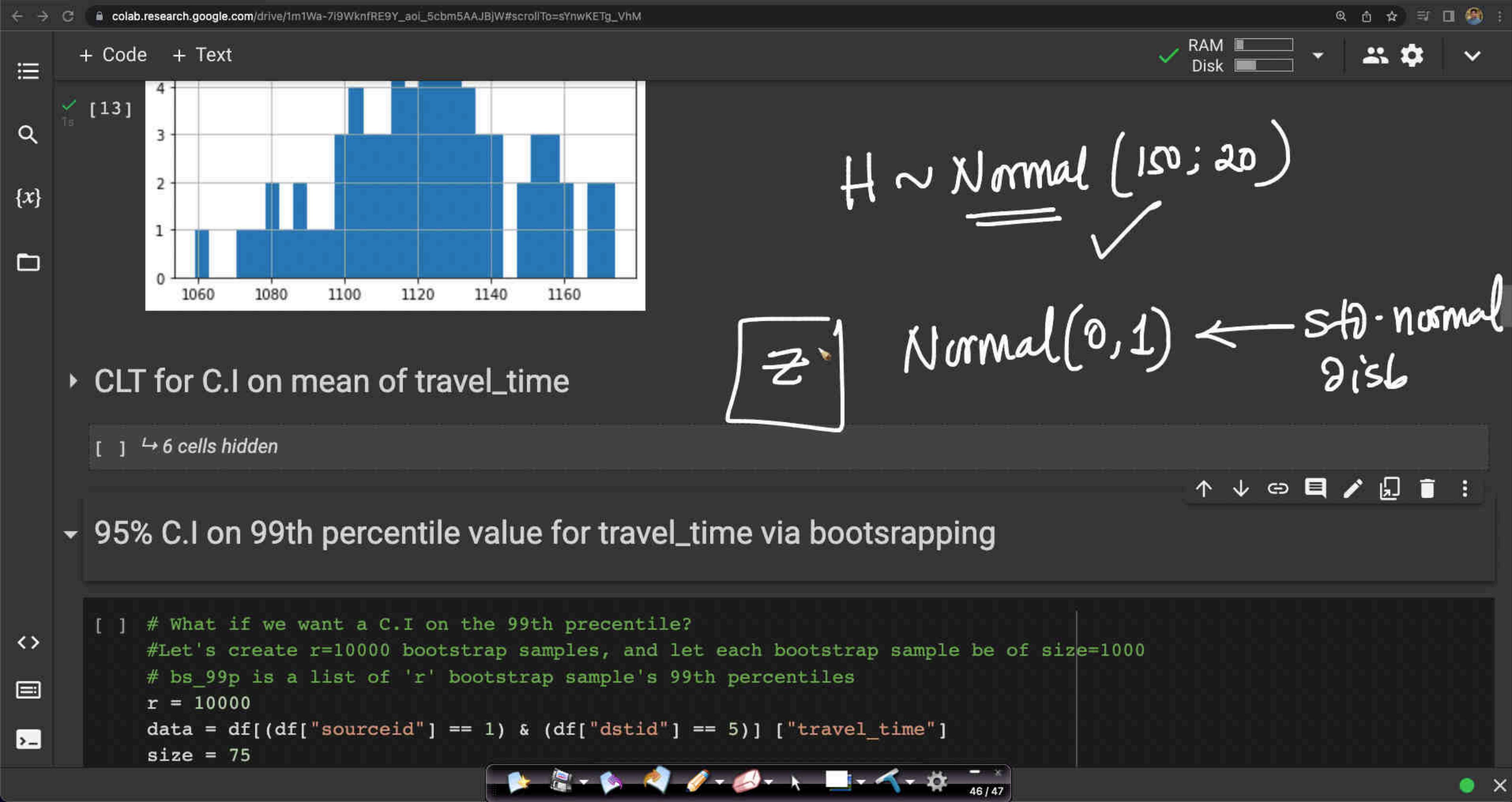
[13] 

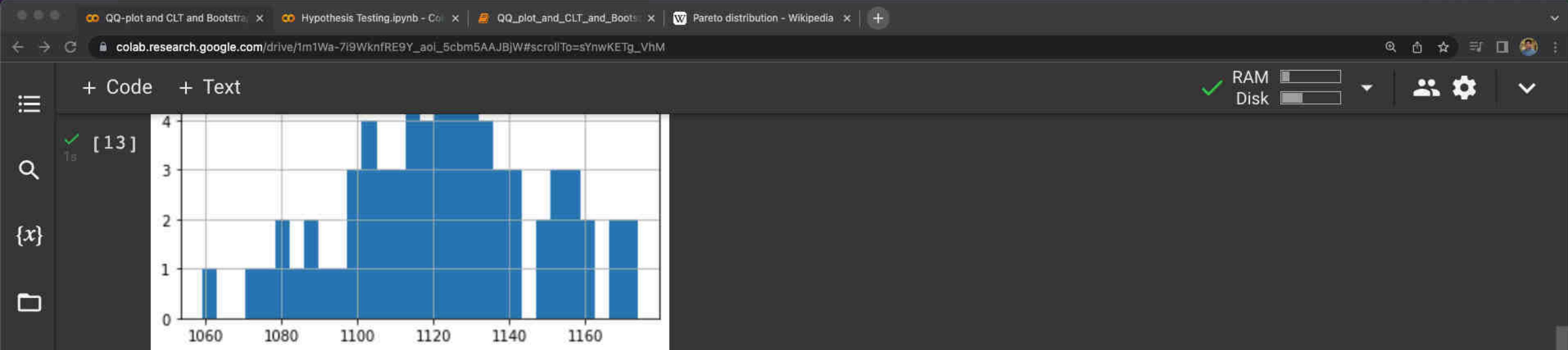
X Normal mean = 0  
std = 1

CLT for C.I on mean of travel\_time

95% C.I on 99th percentile value for travel\_time via bootstrapping

```
[ ] # What if we want a C.I on the 99th percentile?  
# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75
```





## ▶ CLT for C.I on mean of travel\_time

[ ] ↳ 6 cells hidden



## ▼ 95% C.I on 99th percentile value for travel\_time via bootstrapping

```
[ ] # What if we want a C.I on the 99th percentile?  
# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75
```

QQ-plot and CLT and Bootstrap.ipynb - Colab

```
[19] # compute C.I on the mean given that bs_means follows Gaussian distribution: CLT
print(np.mean(bs_means))
print(np.std(bs_means))

M = 1122.823726
S = 3.460407709638272
```

[ ] print(np.mean(bs\_means)-2\*np.std(bs\_means))
print(np.mean(bs\_means)+2\*np.std(bs\_means))

```
1115.9224421865092
1129.6414778134913
```

[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs\_means,2.5))
print(np.percentile(bs\_means,97.5))

# what if r is say 100 and not 10,000?
1115.98
1129.4205000000002

CLT

$M$

$S$

$\mu \sim N$

$std = \frac{\sigma}{\sqrt{n}}$

$S = 3.46 = \frac{\sigma}{\sqrt{n}}$

$S \times \sqrt{n} \approx \sigma$

POPSOP

$\mu$

$\sigma$

$\sqrt{n}$

$\approx$

$\sigma$

$\approx$

std · eraser

05% CI on 99th percentile value for t-distribution bootstrapping

∞ QQ-plot and CLT and Bootstrap x ∞ Hypothesis Testing.ipynb - Colab x ∞ QQ\_plot\_and\_CLT\_and\_Bootstrap x W Pareto distribution - Wikipedia x +

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

+ Code + Text RAM Disk

[19] # compute C.I on the mean given that bs\_means follows Gaussian distribution: CLT

```
print(np.mean(bs_means))
print(np.std(bs_means))
```

1122.823726  
3.460407709638272

n is large

[ ] print(np.mean(bs\_means)-2\*np.std(bs\_means))
print(np.mean(bs\_means)+2\*np.std(bs\_means))

1115.9224421865092  
1129.6414778134913

[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs\_means,2.5))
print(np.percentile(bs\_means,97.5))

# what if r is say 100 and not 10,000?

1115.98  
1129.4205000000002

95% C.I on 99th percentile value for travel time via bootstrapping

$\bar{M} \approx M$

$S = \text{std} * \sqrt{n}$



+ Code + Text

```
[19] # compute C.I on the mean given that bs_means follows Gaussian distribution: CLT  
print(np.mean(bs_means))  
print(np.std(bs_means))
```

```
1122.823726  
3.460407709638272
```

bias estimates

```
[ ] print(np.mean(bs_means)-2*np.std(bs_means))  
print(np.mean(bs_means)+2*np.std(bs_means))
```

```
1115.9224421865092  
1129.6414778134913
```

```
[ ] # could we just use the 2.5th percentile and 97.5th percentile value  
print(np.percentile(bs_means,2.5))  
print(np.percentile(bs_means,97.5))  
  
# what if r is say 100 and not 10,000?
```

```
1115.98  
1129.4205000000002
```

▶ 95% CI on 99th percentile value for travel time via bootstrapping

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=-XRxzwtX8h5s

+ Code + Text RAM Disk

[19] # compute C.I on the mean given that bs\_means follows Gaussian distribution: CLT

```
print(np.mean(bs_means))
print(np.std(bs_means))
```

1122.823726  
3.460407709638272

[ ] print(np.mean(bs\_means)-2\*np.std(bs\_means))
print(np.mean(bs\_means)+2\*np.std(bs\_means))

1115.9224421865092  
1129.6414778134913

[ ] # could we just use the 2.5th percentile and 97.5th percentile value
print(np.percentile(bs\_means, 2.5))
print(np.percentile(bs\_means, 97.5))

# what if r is say 100 and not 10,000?

```
1115.98
1129.4205000000002
```

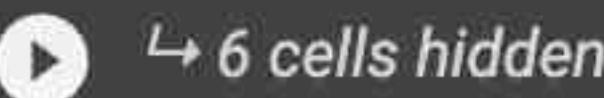
95% C.I on 99th percentile value for travel time via bootstrapping

2.5% ← → 2.5%  
M-2S M M+2S  
2.5% percentile  
97.5% percentile

[https://colab.research.google.com/drive/1m1Wa-7i9WknfBEGY\\_apI-5cbm5AAIBiW#scrollTo=IcDaEZIV](https://colab.research.google.com/drive/1m1Wa-7i9WknfBEGY_apI-5cbm5AAIBiW#scrollTo=IcDaEZIV)

+ Code + Text

✓ RAM Disk





$\{x\} \rightarrow$  95% C.I on 99th percentile value for travel\_time via bootstrapping

```
[ ] # What if we want a C.I on the 99th percentile?  
#Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

of size=1000

question Pgq [ l, u ]

```
[ ] len(bs_99p)
```

[ ] bs\_99p

QQ-plot and CLT and Bootstrap.ipynb

Hypothesis Testing.ipynb

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb

Pareto distribution - Wikipedia

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=IcDaEZiY\_Nj0

+ Code

+ Text

RAM  
Disk

User Settings

Up Down Left Right Home End

{x} ▾ 95% C.I on 99th percentile value for travel\_time via bootstrapping

P<sub>qq</sub> → CLT X

```
[ ] # What if we want a C.I on the 99th percentile?  
# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

```
[ ] len(bs_99p)
```

```
10000
```

```
[ ] bs_99p
```

QQ-plot and CLT and Bootstrap.ipynb - Colab | Hypothesis Testing.ipynb - Colab

QQ-plot and CLT and Bootstrap.ipynb - Colab | Hypothesis Testing.ipynb - Colab

QQ-plot and CLT and Bootstrap.ipynb - Colab | QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb - Colab

QQ-plot and CLT and Bootstrap - Colab | Hypothesis Testing.ipynb - Colab | QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb - Colab | Pareto distribution - Wikipedia

[colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\\_aol-5cbm5AAJBjW#scrollTo=tmxw002-9](https://colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y_aol-5cbm5AAJBjW#scrollTo=tmxw002-9)

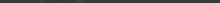
+ Code + Text

[ 1 ] ↳ 6 cells hidden

✓ RAM Disk

1

- ▼ 95% C.I on 99th percentile value for travel\_time via bootstrapping



```
# What if we want a C.I on the 99th percentile?  
#Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

[ ] len(bs 99p)

10000

[ ] bs 99p

QQ-plot and CLT and Bootstrap.ipynb - Colab | Hypothesis Testing.ipynb - Colab

Hypothesis Testing.ipynb - Co

QQ\_plot\_and\_CLT\_and\_Boots x | W Pareto distribution

Wikipedia X | +

Code + Text

✓ RAM Disk

▼

[ ] ↳ 6 cells hidden

$\{x\} \rightarrow$  95% C.I on 99th percentile value for travel\_time via bootstrapping

↑ ↓ ⌂ 目 齿 阅 ⌂ ⋮

```
# What if we want a C.I on the 99th percentile?  
#Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

[ ] len(bs\_99p)

10000

[ ] bs\_99p

QQ-plot and CLT and Bootstrap | Hypothesis Testing.ipynb - Colab | QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb | Pareto distribution - Wikipedia | +

[colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\\_aol\\_5cbm5AAJBjW#scrollTo=tmxw002-9n-A](https://colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y_aol_5cbm5AAJBjW#scrollTo=tmxw002-9n-A)

+ Code + Text

✓ RAM Disk

[ ] ↳ 6 cells hidden

{x} ▾ 95% C.I on 99th percentile value for travel\_time via bootstrapping

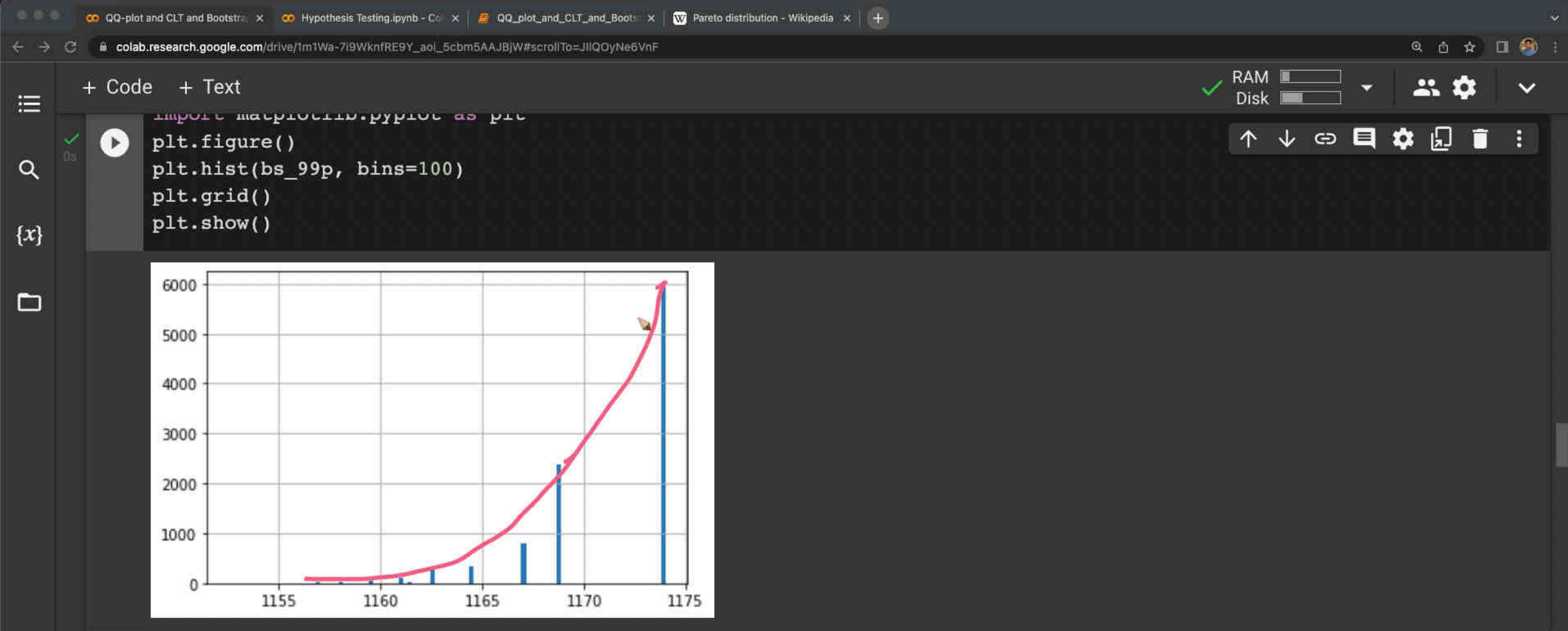
A set of small, light-gray navigation icons located at the bottom of the page. From left to right, they include: a double arrow pointing up and down, a circular arrow, a speech bubble, a gear, a square with a diagonal line, a trash can, and three vertical dots.

```
# What if we want a C.I on the 99th percentile?  
#Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=100  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)  
  
len(bs_99p)
```

[ ] len(bs\_99p)

10000

[ ] bs 99p



Code:

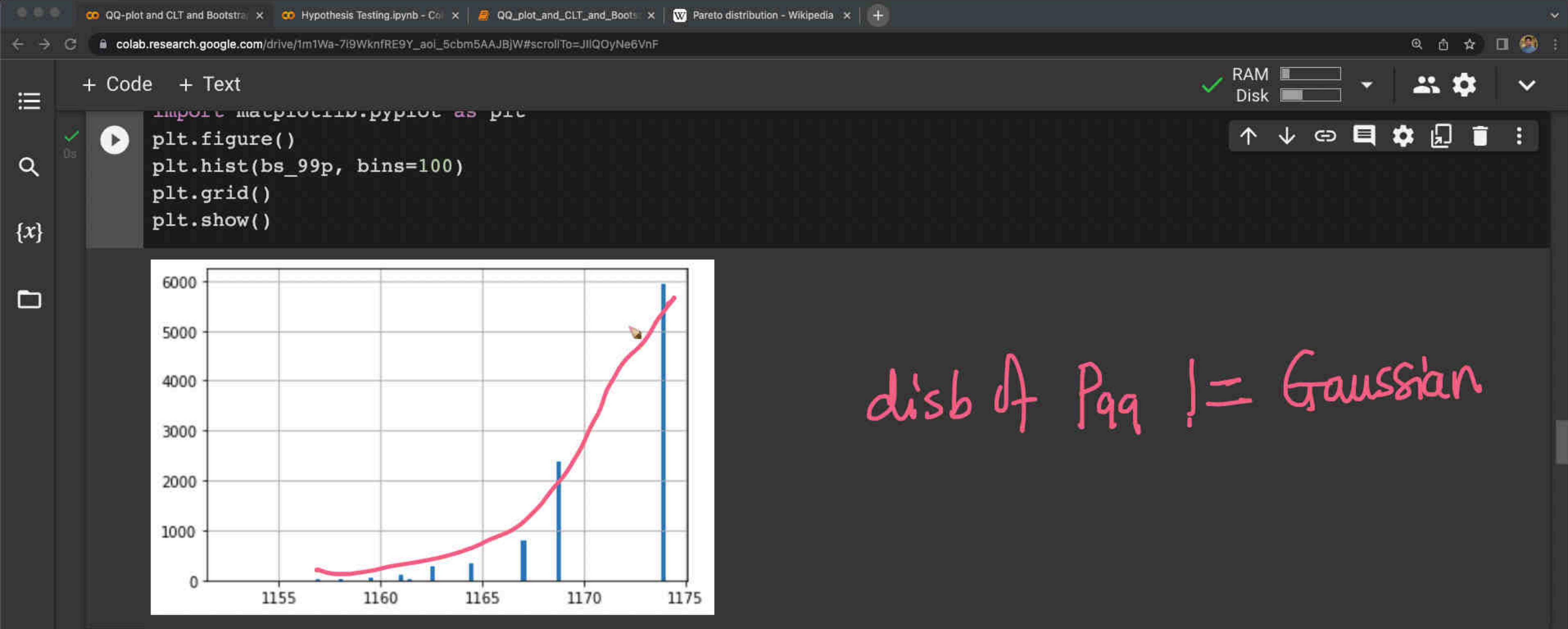
```
# QQ-plot with normal distribution
fig, ax1 = plt.subplots()
prob = stats.probplot(bs_99p, dist=stats.norm, plot=ax1)
```

Output:

Probability Plot

1185

57 / 57



# QQ-plot with normal distribution  
fig, ax1 = plt.subplots()  
prob = stats.probplot(bs\_99p, dist=stats.norm, plot=ax1)

Probability Plot

1185

58 / 58

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=zyEVchf-C6q

+ Code + Text

RAM Disk

[20] `for i in range(r):  
 bs_sample = np.random.choice(data, size=size)  
 bs_99p[i] = np.percentile(bs_sample, 99)`

[21] `len(bs_99p)`

10000

[23] `bs_99p`

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

#bs\_99p may or maynot be normally distributed.  
`print(np.percentile(bs_99p,2.5))  
print(np.percentile(bs_99p,97.5))`

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
`print(np.percentile(data,99))`

$x_1, x_2, \dots, x_{75}$

↓ SI with replacement

$x_1, \dots$

59 / 59

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebooks

Hypothesis Testing.ipynb - Colab Notebooks

QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb - Colab Notebooks

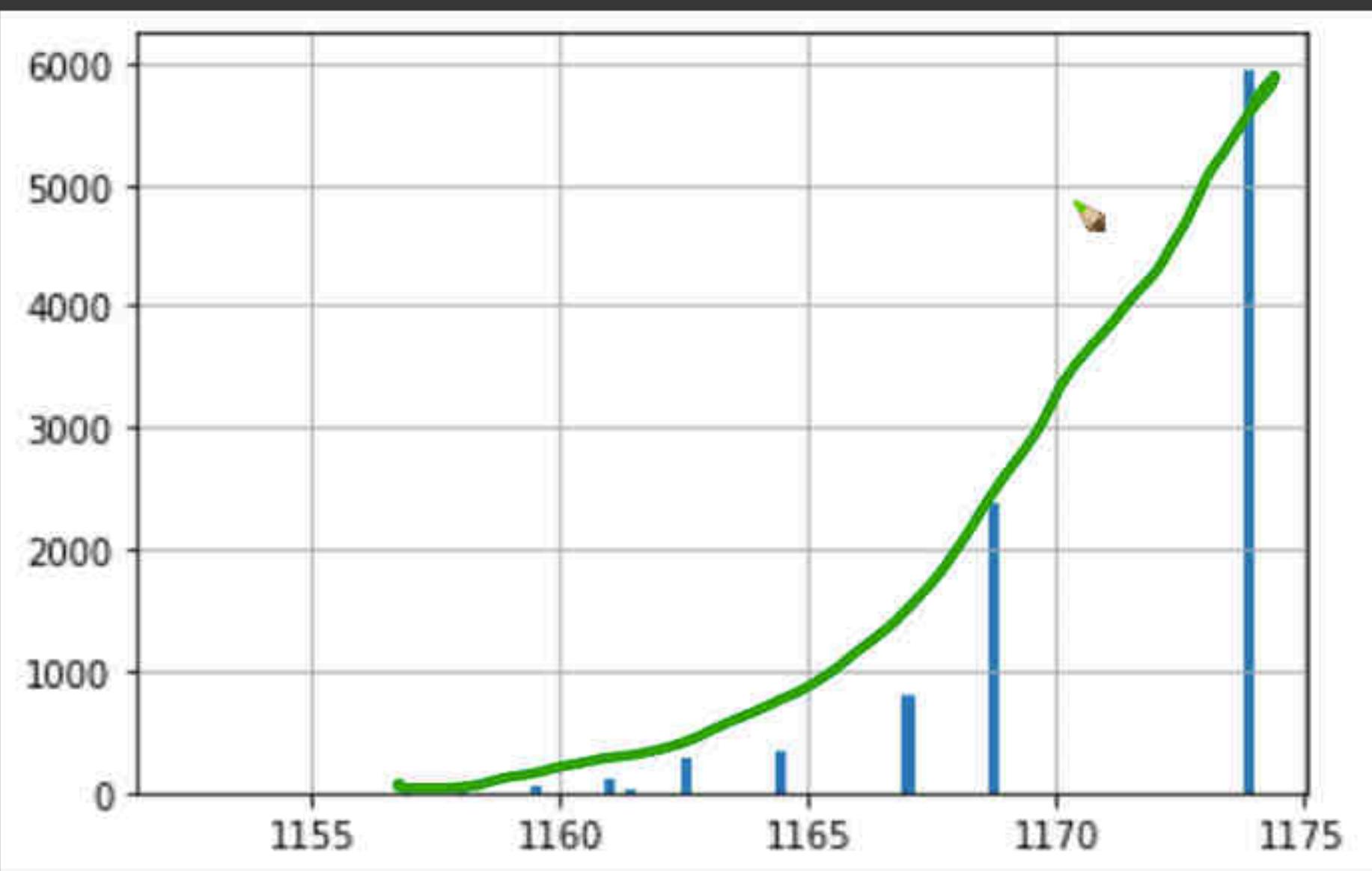
Pareto distribution - Wikipedia

+ Code + Text

✓ RAM  
Disk

👤⚙️

```
[22] # plot the pdf of bs_99p
[22] import matplotlib.pyplot as plt
[22] plt.figure()
[22] plt.hist(bs_99p, bins=100)
[22] plt.grid()
[22] plt.show()
```



```
▶ # QQ-plot with normal distribution
[23] fig, ax1 = plt.subplots()
[23] prob = stats.probplot(bs_99p, dist=stats.norm, plot=ax1)
```

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=zyEVchf-C6q

+ Code + Text RAM Disk

10000

{x} [23] bs\_99p

```
array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])
```

[ ] #bs\_99p may or maynot be normally distributed.

```
print(np.percentile(bs_99p,2.5))
print(np.percentile(bs_99p,97.5))
```

{1162.56  
1174.0}

[ ] # Point estimate of the 99th percentile of the 75 observed samples

```
print(np.percentile(data,99))
```

1174.0

[22] # plot the pdf of bs\_99p

```
import matplotlib.pyplot as plt
plt.figure()
plt.hist(bs_99p, bins=100)
```

RAM Disk

Play button

2.5%

97.5%

P<sub>99</sub>

P<sub>2.5</sub>:

1162.56 1174.0

61 / 61

60

 [QQ-plot and CLT and Bootstrap.ipynb](#) ×  [Hypothesis Testing.ipynb - Co](#)

Hypothesis Testing.ipynb - Colab

Pareto distribution - Wikipedia

17

[https://colab.research.google.com/drive/1m1Wa-7j8WkfREFQY\\_aj-5cbm5AAIBIW#scrollTo=zvYEyhf-C6g](https://colab.research.google.com/drive/1m1Wa-7j8WkfREFQY_aj-5cbm5AAIBIW#scrollTo=zvYEyhf-C6g)

◎ 佛學研究

## # Code # Text

✓ RAM Disk

V

10000

✓ [23] bs 99p

```
array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])
```

A set of small, light-gray navigation icons located at the bottom of the page. From left to right, they include: a double arrow pointing up and down, a circular arrow, a magnifying glass, a gear, a square with rounded corners, a document icon, and three vertical dots.

```
#bs_99p may or maynot be normally distributed.  
print(np.percentile(bs_99p,2.5))  
print(np.percentile(bs_99p,97.5))
```

1162.56  
1174.0

```
[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data, 99))
```

1174.0

75

[ 22 ]

```
[22] # plot the pdf of bs_99p
    import matplotlib.pyplot as plt
    plt.figure()
    plt.hist(bs_99p, bins=100)
```

[QQ-plot and CLT and Bootstrap](#) | [Hypothesis Testing.ipynb - Colab](#) | [QQ plot and CLT and Bootstrap](#) | [Pareto distribution - Wikipedia](#)

- Code + Text

```
# What if we want a C.I on the 99th percentile?  
# Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>r</sub>

[21] len(bs\_99p)

10000

✓ [ 23 ] bs\_99p

```
array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])
```

```
[ ] #bs_99p may or maynot be normally distributed.  
print(np.percentile(bs_99p,2.5))  
print(np.percentile(bs_99p,97.5))
```

QQ-plot and CLT and Bootstrap.ipynb - Colab

```
# What if we want a C.I on the 99th percentile?  
#Let's create r=10000 bootstrap samples, and let each bootstrap sample be of size=1000  
# bs_99p is a list of 'r' bootstrap sample's 99th percentiles  
r = 10000  
data = df[(df["sourceid"] == 1) & (df["dstid"] == 5)] ["travel_time"]  
size = 75  
bs_99p = np.empty(r)  
  
for i in range(r):  
    bs_sample = np.random.choice(data, size=size)  
    bs_99p[i] = np.percentile(bs_sample, 99)
```

$S_i$ : - - - - -

[21] len(bs\_99p)

10000

[23] bs\_99p

```
array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])
```

```
[ ] #bs_99p may or maynot be normally distributed.  
print(np.percentile(bs_99p,2.5))  
print(np.percentile(bs_99p,97.5))
```

Chrome File Edit View Bookmarks Profiles Tab Window Help

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=tmxw002-9n-A

+ Code + Text

[23]  $\text{bs\_99p}$

```
array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])
```

{ } [ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

[22] # plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

RAM Disk

95% C.I. m Pqq

bs-99p

95%

q-5 97-5

0-5 99% C.I. 99-5

6000

65 / 65

QQ-plot and CLT and Bootstrap.ipynb - Colab

QQ-plot\_and\_CLT\_and\_Bootstrap.ipynb - Colab

Pareto distribution - Wikipedia

RAM Disk

+ Code + Text

[23] ~~array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])~~

{x} [ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

[ ] # plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

6000

66 / 66

QQ-plot and CLT and Bootstrap.ipynb - Colab

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

[ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

[ ] # plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

6000

RAM Disk

q5% C.I [ l, u ]

on P99

75

QQ-plot and CLT and Bootstrap.ipynb - Colab

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

[ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

# plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

point estimate : Mean  $\bar{x}$

as 1. c.l: [98, 102.1]

6000

RAM Disk

68 / 68

QQ-plot and CLT and Bootstrap.ipynb - Colab

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

[ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

[ ] # plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

6000

RAM Disk

travel-times

4.5 Million

sample-size  
S1 → 100,000 → .  
S2  
:  
S10,000

69 / 69

QQ-plot and CLT and Bootstrap.ipynb - Colab

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

[ ] #bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

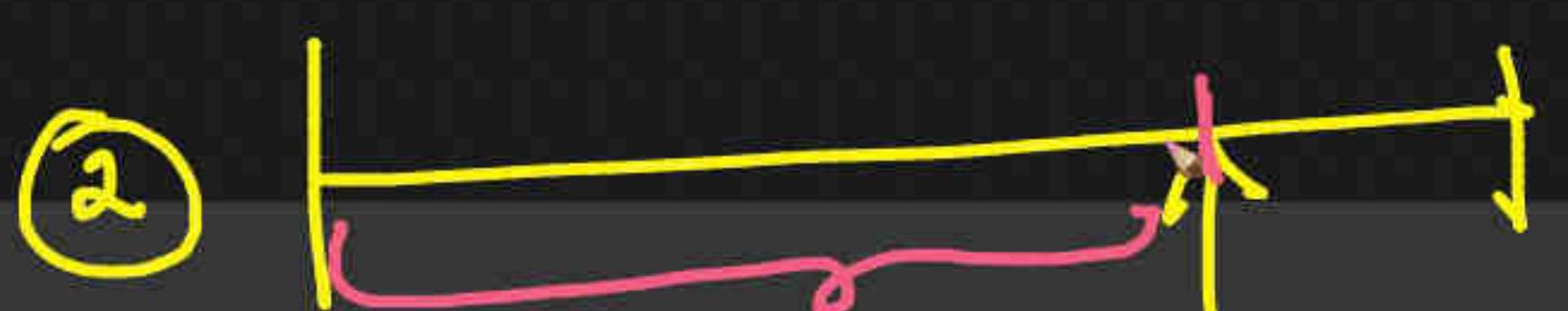
[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

# plot the pdf of bs\_99p  
import matplotlib.pyplot as plt  
plt.figure()  
plt.hist(bs\_99p, bins=100)  
plt.grid()  
plt.show()

6000

① Set  $x_1, \dots, x_{100}$

② 

99th position

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=zyEVchf-C6q

+ Code + Text RAM Disk

[21] len(bs\_99p)

10000

{x}

[23] bs\_99p

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

#bs\_99p may or maynot be normally distributed.  
print(np.percentile(bs\_99p,2.5))  
print(np.percentile(bs\_99p,97.5))

1162.56  
1174.0

95% C.I of P99

[ ] # Point estimate of the 99th percentile of the 75 observed samples  
print(np.percentile(data,99))

1174.0

point estimation of P99

[22] # plot the pdf of bs\_99p  
import matplotlib.pyplot as plt

qq-plot and CLT and Bootstrap.ipynb - Colab

prices of items sold

99 100 Max

1168.82

1162.56

1174.0

#bs\_99p may or maynot be normally distributed.

print(np.percentile(bs\_99p,2.5))

print(np.percentile(bs\_99p,97.5))

[ ] # Point estimate of the 99th percentile of the 75 observed samples

print(np.percentile(data,99))

1174.0

[22] # plot the pdf of bs\_99p

import matplotlib.pyplot as plt

QQ-plot and CLT and Bootstrap.ipynb - Colab

colab.research.google.com/drive/1m1Wa-7i9WknfRE9Y\_aol\_5cbm5AAJBjW#scrollTo=zyEVchf-C6q

+ Code + Text

RAM Disk

[21] len(bs\_99p)

10000

[23] bs\_99p

array([1174. , 1174. , 1174. , ..., 1167. , 1174. , 1168.82])

#bs\_99p may or maynot be normally distributed.

print(np.percentile(bs\_99p,2.5))

print(np.percentile(bs\_99p,97.5))

1162.56

1174.0

# Point estimate of the 99th percentile of the 75 observed samples

print(np.percentile(data,99))

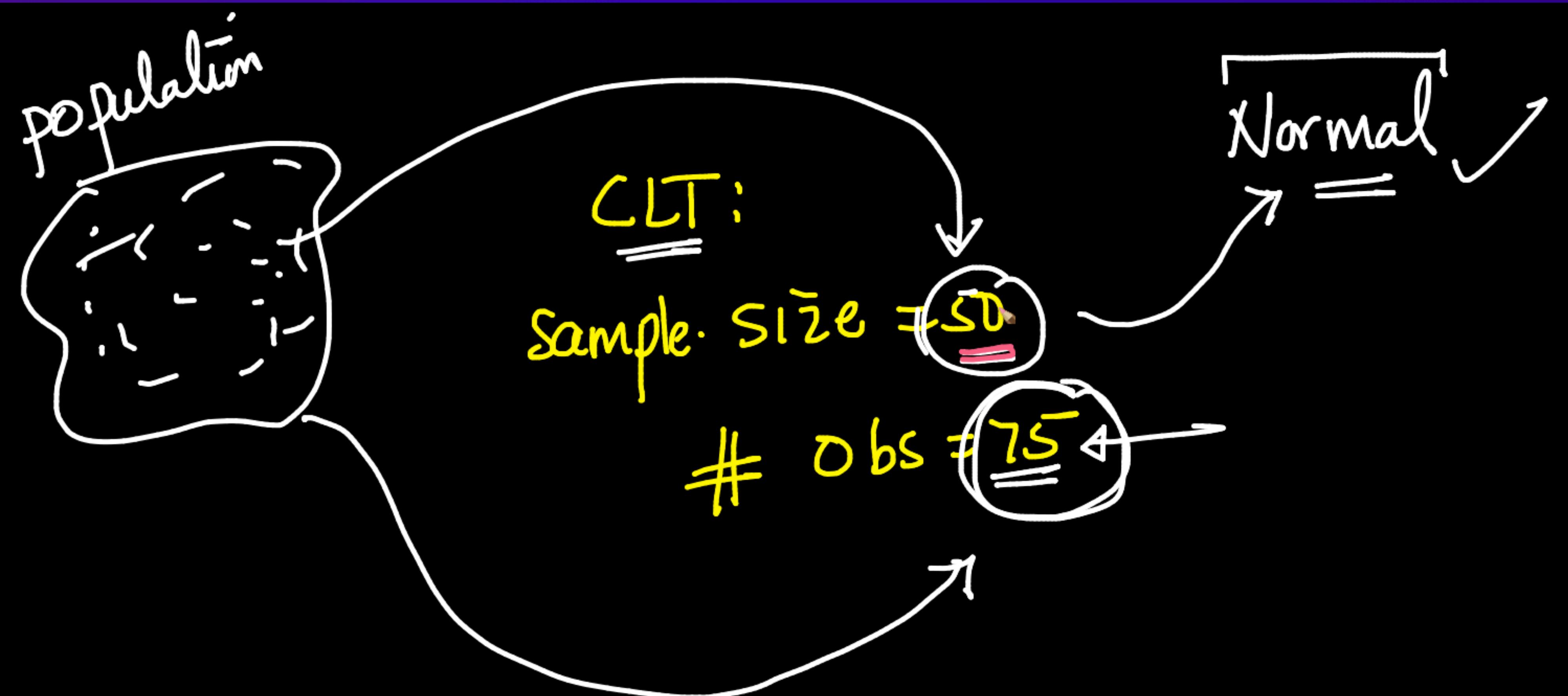
1174.0

[22] # plot the pdf of bs\_99p

import matplotlib.pyplot as plt

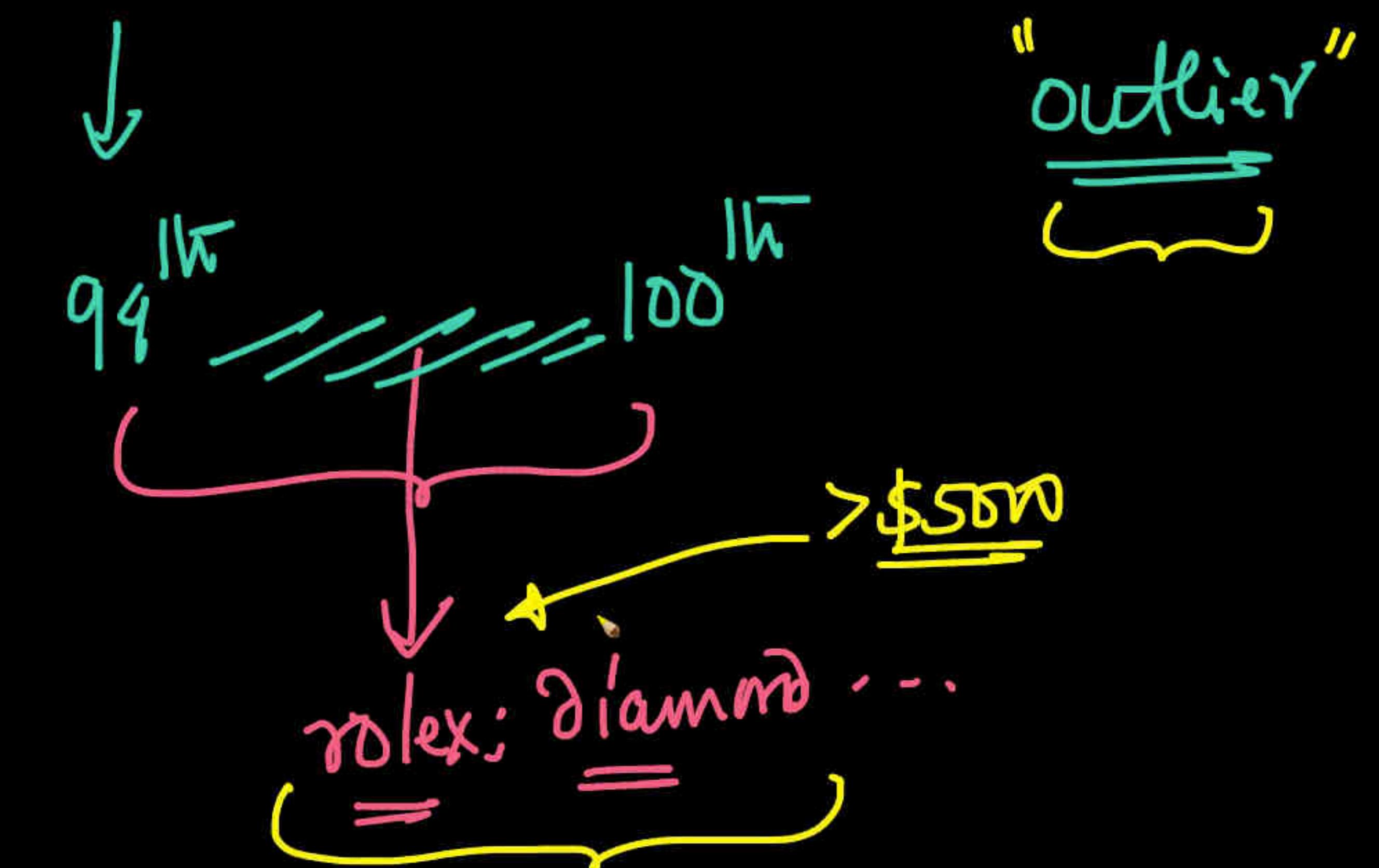
qq-q [W]

0.1 Y,



$\gamma = 10,000 \leftarrow$  computational reasonable

$\underline{n} = \text{sample size} \leq \underline{\# \text{obs}}$



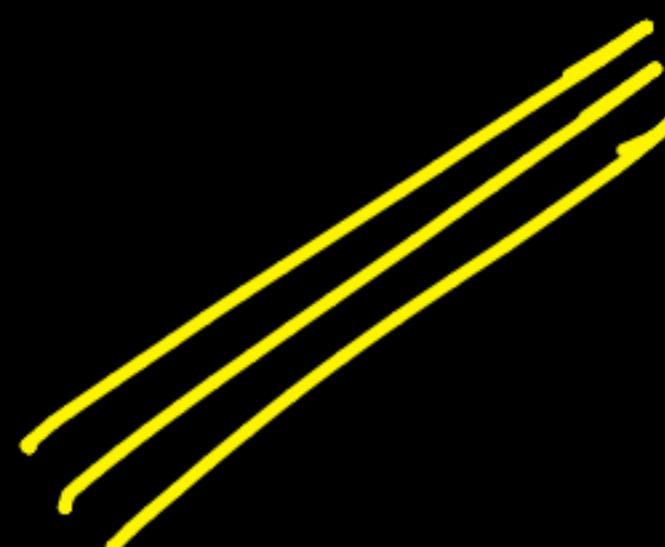
can be used ..

- 98.5 vs 99.5

vs 99  $\bar{h}$



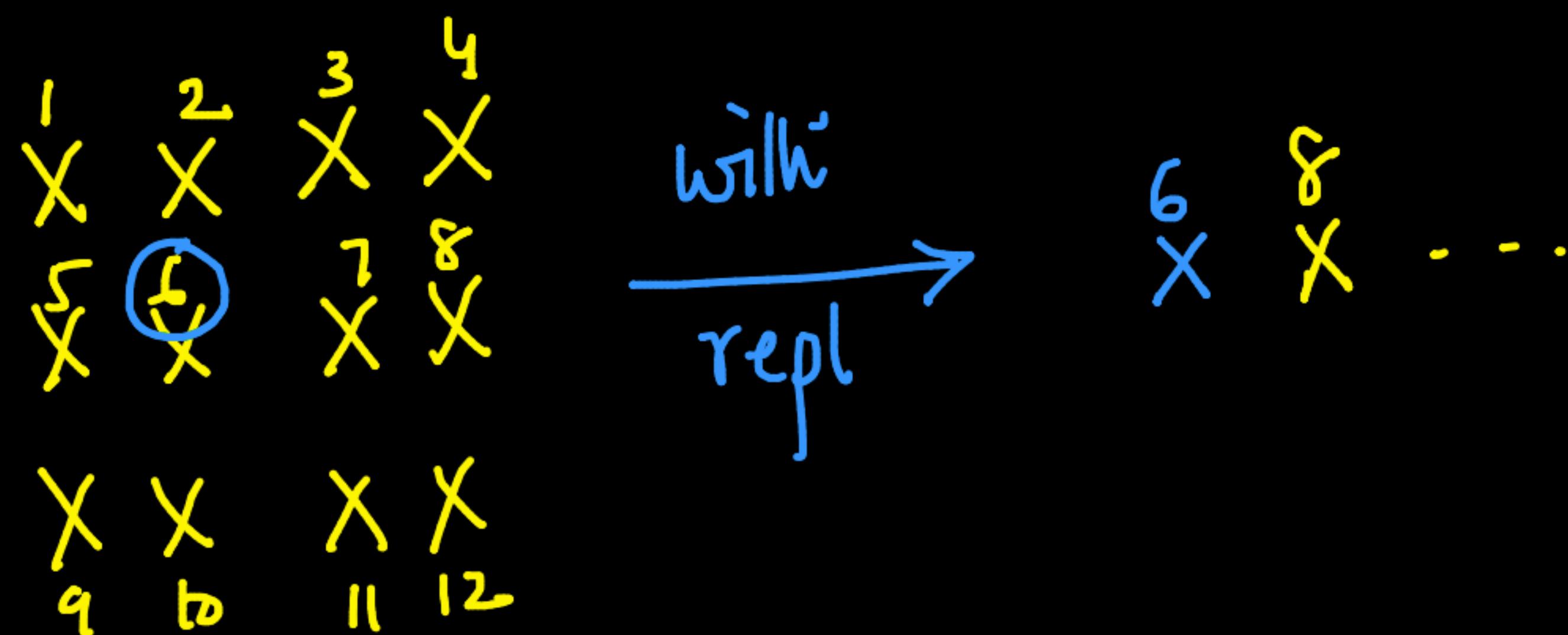
typically  
used values



Samplings:

Random Sampling ✓

Uniform dist  
[1,12]



[1 ... 12]

1 2 3 4  
~~2~~ ~~6~~ ~~5~~ ~~4~~  
5 6 ~~8~~ ~~7~~ ~~12~~

~~9~~ ~~2~~ ~~6~~ ~~7~~  
9 10 11

8-pls →

without  
repl

$x^7, x^8, x^{12}$

✓ { family incomes in India  
    across states ✓

Stratified Sampling

✓ random - sample

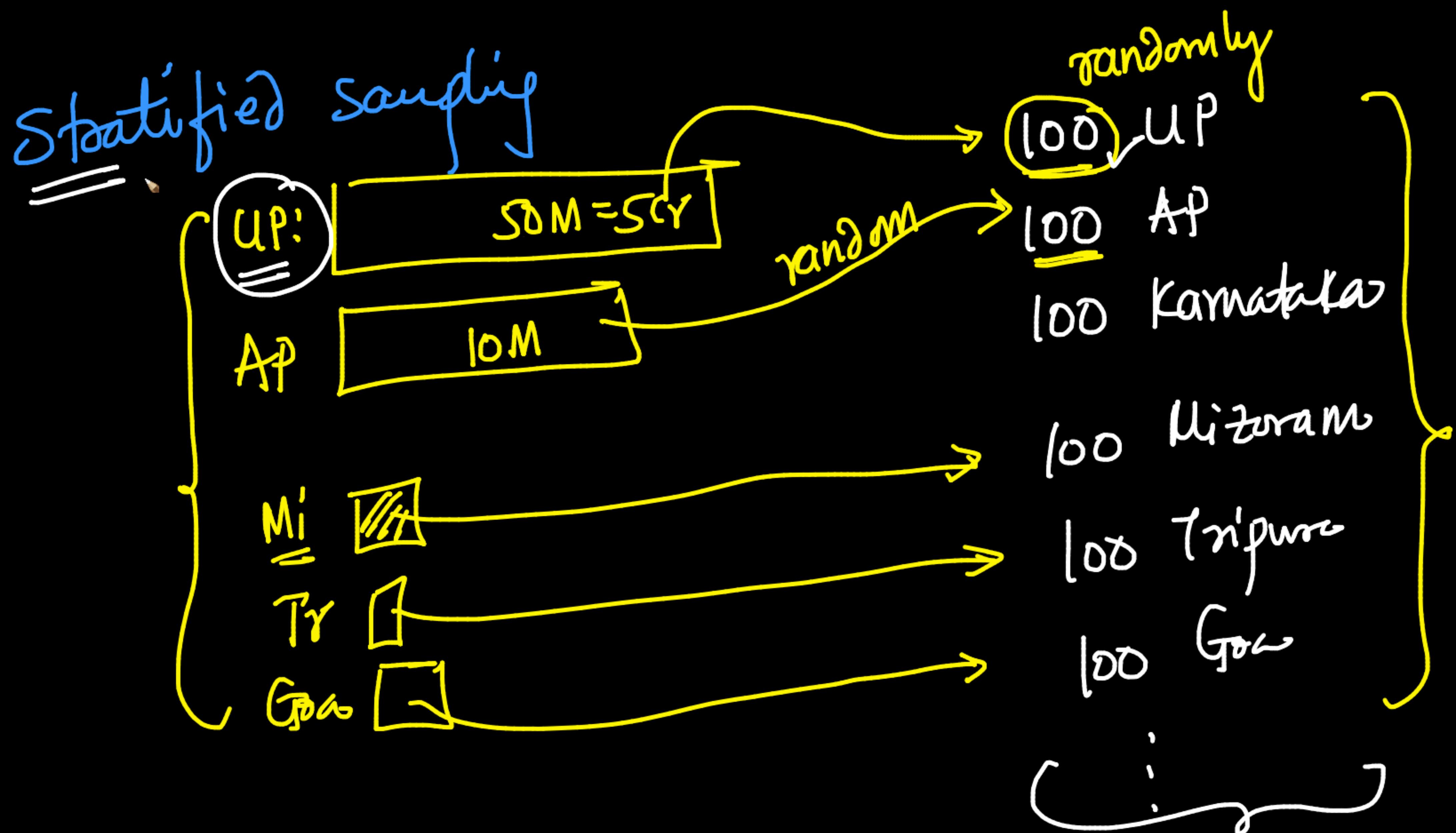
→ More obs from more  
    populous states

[UP] ✓

330M

DO NOT WANT

Mizoram



~~Systematic Sampling:~~

4 → 330M

Salaries of each family in Income

①

Sort all the families by salary

330M → inc

$F_1, F_2, \dots, F_{330}, \dots, F_{10}$

1M families

2

picking @ regular intervals

# Hypothesis testing

- Most useful & fundamental
- Tricky to grasp
  - ↳ counter intuition
- takes time to grasp
- Terminology ...

↳ C.I → ✓  
↳ p-value  
→ dish ✓

revise after  
class

# Toy - problem:

→ toss  $n$  times

coin is fair or biased towards heads

$$P(H) = 0.5$$

1

$$PCH) \geq 0.5$$

probabilistic

6/10

heads → biased or fair → NO  
towards heads

INTUITION

65/100

heads → biased towards heads

52/100

→ fair (very confident)

651/1000

→ biased (v.v.confident)

✓ Task: coin is fair or unbiased towards heads  
Experiment: Toss 100 times

Test-statistic:  $\sqrt{T}$  (T)  
Sensible way to measure the outcome of the expt which can help us answer & solve my task  
#heads in 100 tosses (expt)

let  $T_{obs}$  ← observed value of T after conducting the expt  
 $T_{obs} = 65$

default statement / assumption

✓ Null-hypothesis ( $H_0$ ): coin is fair ✓

✓ Alt. hyp. ( $H_a$ ): coin is biased towards heads =  
=  $\rightarrow (H_1)$

$P$  (observing as extreme or more than  $T_{\text{obs}} = 65$  |  $H_0$ ) = ?  $\rightarrow$  coin is fair

$$P(T > T_{\text{obs}} = 65 \mid H_0) = ? \quad \dots$$

$\hookrightarrow$  how to compute it



$H_0$  is true

$\checkmark [T \sim \text{Binomial}(n=100; p=0.5)]$

CDF

$$P(H) = 0.5$$

$T = \# \text{heads}$   
in 100 tosses

$$P\left(\frac{\text{65 or More heads in 100 tosses}}{||} H_0\right) = P(T \geq 65) = 1 - \text{CDF}(T=64)$$

$$\underline{\underline{P(T \leq 64)}}$$

QQ-plot and CLT and Bootstrap.ipynb · Hypothesis Testing.ipynb · QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb · Pareto distribution - Wikipedia

colab.research.google.com/drive/1uNuXZ3pBTxJxZTdJEx-AHzq3lx8teh5n#scrollTo=MRYFtPjZ9iCz

+ Code + Text

✓ RAM Disk



```
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
import seaborn as sns
```

```
[ ] # Expt: Toss the coin 100 times
# Test statistic: Count the number of heads
# H0: Coin is fair
# Ha: Coin is biased towards heads
# T ~ Binomial(n=100, p=0.5) under Null Hypothesis(=H0)

prob = stats.binom.cdf(k=64, n=100, p=0.5)
print(1-prob) # P(T >= 65 | H0 )
```

0.0017588208614850442

```
[ ] # Plot PDF of Binomial(100, 0.5)
X = stats.binom.rvs(n=100, p=0.5, size=100000)
plt.grid()
sns.distplot(X)
```

QQ-plot and CLT and Bootstrap.ipynb · Hypothesis Testing.ipynb · QQ\_plot\_and\_CLT\_and\_Bootstrap.ipynb · Pareto distribution - Wikipedia

colab.research.google.com/drive/1uNuXZ3pBTxJxZTdJEx-AHzq3lx8teh5n#scrollTo=MRYFtPjZ9iCz

+ Code + Text RAM Disk

import numpy as np  
import matplotlib.pyplot as plt  
from scipy import stats  
import seaborn as sns

[ ] # Expt: Toss the coin 100 times  
# Test statistic: Count the number of heads  
# H0: Coin is fair  
# Ha: Coin is biased towards heads  
#  $T \sim \text{Binomial}(n=100, p=0.5)$  under Null Hypothesis( $=H_0$ )

prob = stats.binom.cdf(k=64, n=100, p=0.5)  
print(1-prob) #  $P(T \geq 65 | H_0)$

0.0017588208614850442

0.175 ✓

$\rightarrow \text{CDF}(64)$

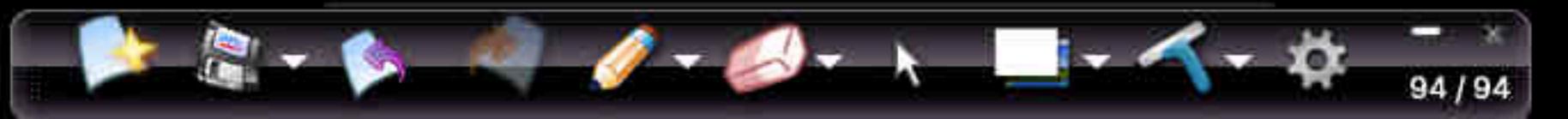
0.00175 - - -

# Plot PDF of Binomial(100, 0.5)  
X = stats.binom.rvs(n=100, p=0.5, size=100000)  
plt.grid()  
sns.distplot(X)

/usr/local/lib/python3.7/dist-packages/ 2019-05-29 11:20:54.545996: FutureWarning: 'distplot' is a deprecated function

91 / 91

$$\begin{aligned} P(T \geq 65) &= 1 - P(T < 65) \\ &= 1 - P(T \leq \underline{\underline{64}}) \\ &\quad \text{CDF}(T=64) \end{aligned}$$



$$P(T \geq T_{\text{obs}} = 65 \mid H_0) = 0.17\%$$

$\checkmark$  }  $\times_{H_0}$ : coin is fair  
 $\equiv$   
 $0.17\% : \text{v-low}$   
 $\equiv$   
reject  $H_0$

$\text{Ha!}$  } coin is biased towards heads  
 $\equiv$   
accept  $\text{Ha!}$

rule of thumb

default threshold



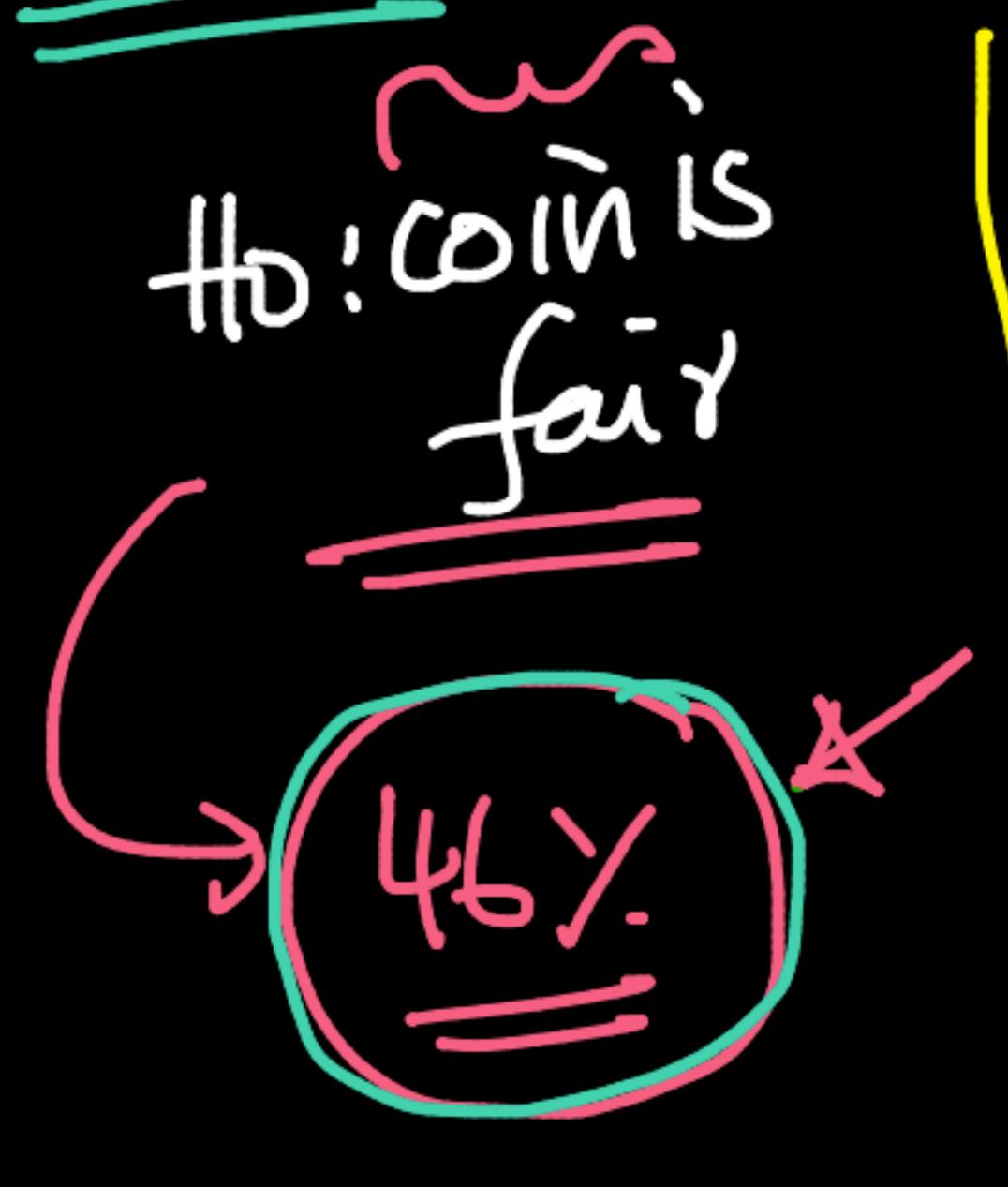
can't reject  $H_0$   
accept  $H_0$

binomial obs

$T_{obs} = 51$  heads

$$P(T \geq T_{obs} = 51 \mid H_0)$$

$H_0$ : coin is fair



read your eqn in English

$$= 0.46 = 46\%$$

$H_a$ : coin is biased toward heads

γwatch → did not click

46% ≠ p(heads)

↳  $P(T > 51 \text{ (Tobs)} | H_0) = 46\%$ .

QQ-plot and CLT and Bootstrap.ipynb - Colab Notebook

colab.research.google.com/drive/1uNuXZ3pBTxJxZTdJEx-AHzq3lx8teh5n#scrollTo=wBXEHYCIKb9u

+ Code + Text

RAM Disk

[1] import numpy as np  
import matplotlib.pyplot as plt  
from scipy import stats  
import seaborn as sns

# Expt: Toss the coin 100 times  
# Test statistic: Count the number of heads  
# H0: Coin is fair  
# Ha: Coin is biased towards heads  
# T ~ Binomial(n=100, p=0.5) under Null Hypothesis(=H0)

prob = stats.binom.cdf(k=50, n=100, p=0.5)  
print(1-prob) # P(T >= 65 | H0 )

0.4602053813064103

# Plot PDF of Binomial(100, 0.5)  
X = stats.binom.rvs(n=100, p=0.5, size=100000)  
plt.grid()  
sns.distplot(X)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:261: FutureWarning: `distplot` is a deprecated function

RAM Disk

P(T >= 65 | H0 )

5-6 classes

assessment → off sync  
→ back to sync → next module

newatch → hyp.-testing

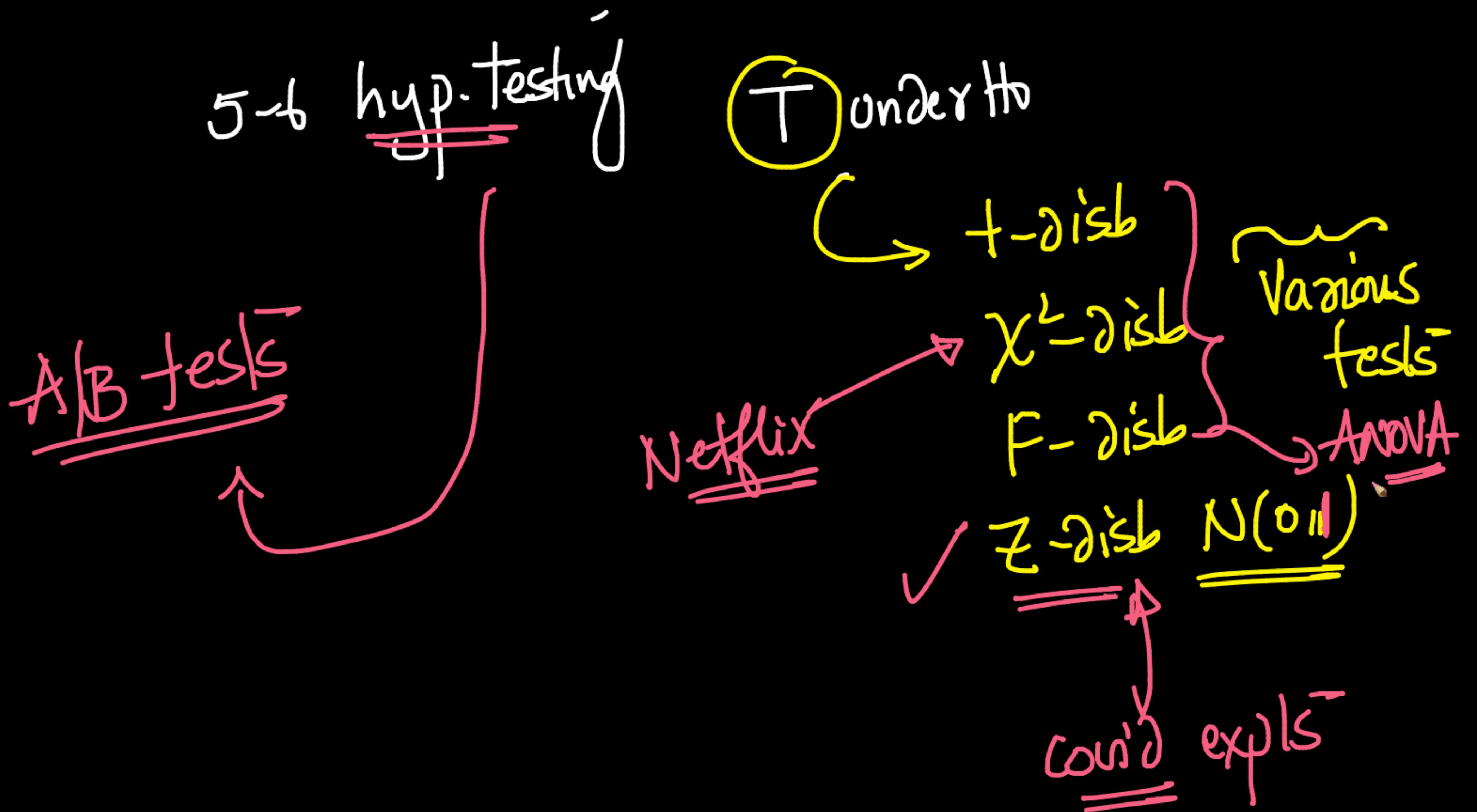
small <5%  convention

has some dish  
 $T \sim \text{Bin}(n, p) \rightarrow P(T, T_{\text{obs}} | H)$

Under  $H_0$

future classes - - -  
(permutation testing)

↓  
videos as post  
real



rec-times

$M_1: x_1^1, x_2^1 \dots x_{50}^1$  ✓

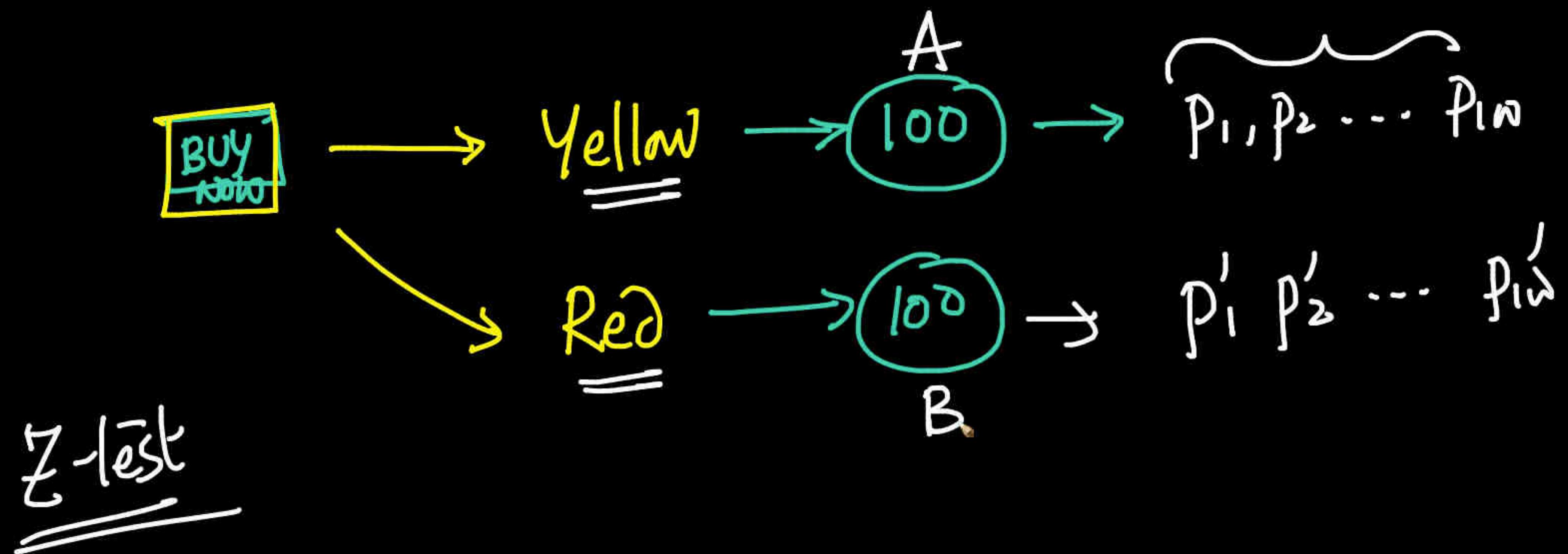
$M_2: x_1^2, x_2^2 \dots x_{60}^2$



$H_0:$  mean rec-time of  $\tilde{M}_1$  = mean rec-time  
of  $\tilde{M}_2$

$H_a:$  "

≠      1



conversations about popular titles. Second, we can help members choose some great content to watch by fulfilling the intrinsic human desire to be part of a shared conversation.

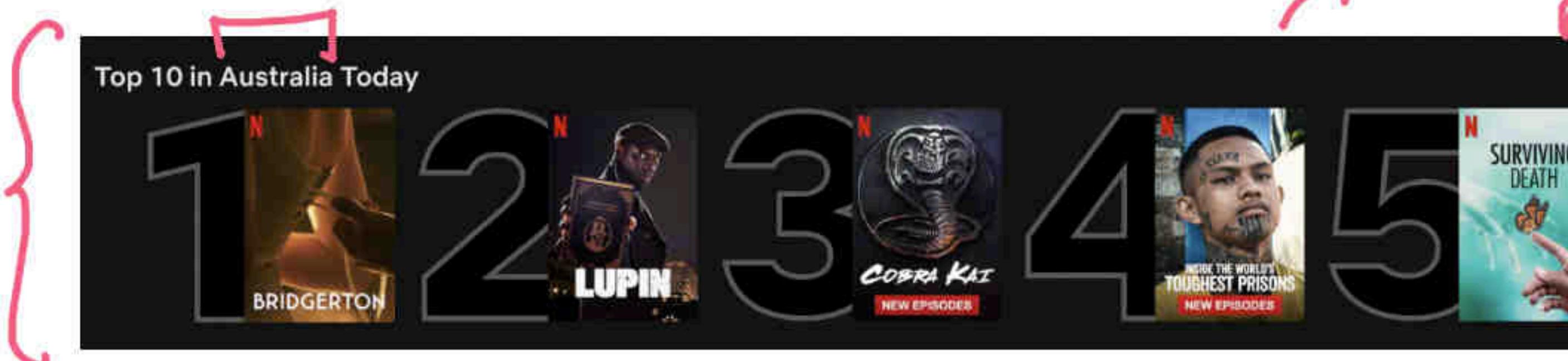


Figure 5: An example of the Top 10 experience on the Web UI.

We next turn this idea into a testable hypothesis, a statement of the form “If we make change X, it will improve the member experience in a way that makes metric Y improve.” With the Top 10 example, the hypothesis read: “*Showing members the Top 10 experience will help them find something to watch, increasing member joy and satisfaction.*” The primary decision metric for this test (and many others) is a measure of member engagement with Netflix: are the ideas we are testing helping our members to choose Netflix as their entertainment destination on any given night? Our research shows that this metric (details omitted) is correlated, in the long term, with how many members will retain their subscriptions. Other areas of the

## Netflix Technology Blog

265K Followers

Learn more about how Netflix designs, builds, and operates our systems and engineering organizations

Follow



## More from Medium

 Netflix Technology ... in Netflix Tech...

**Experimentation is a major focus of Data Science across Netflix**



 Diana Jerman in disney-streaming

**Disney Streaming Embraces 3 Key Tenets of Experimentation**



 Clare La... in Expedia Group Technol...

**2021's Top 10 Posts on Expedia Group Technology**



 Zhiheng Xu in The Airbnb Tech Blog

**Intelligent Automation Platform: Empowering Conversational AI and Beyond at Airbnb**



Task:

A:

100  $\rightarrow$  no top 10

$t_1, t_2, \dots, t_{10}$   
how much time  
Spent

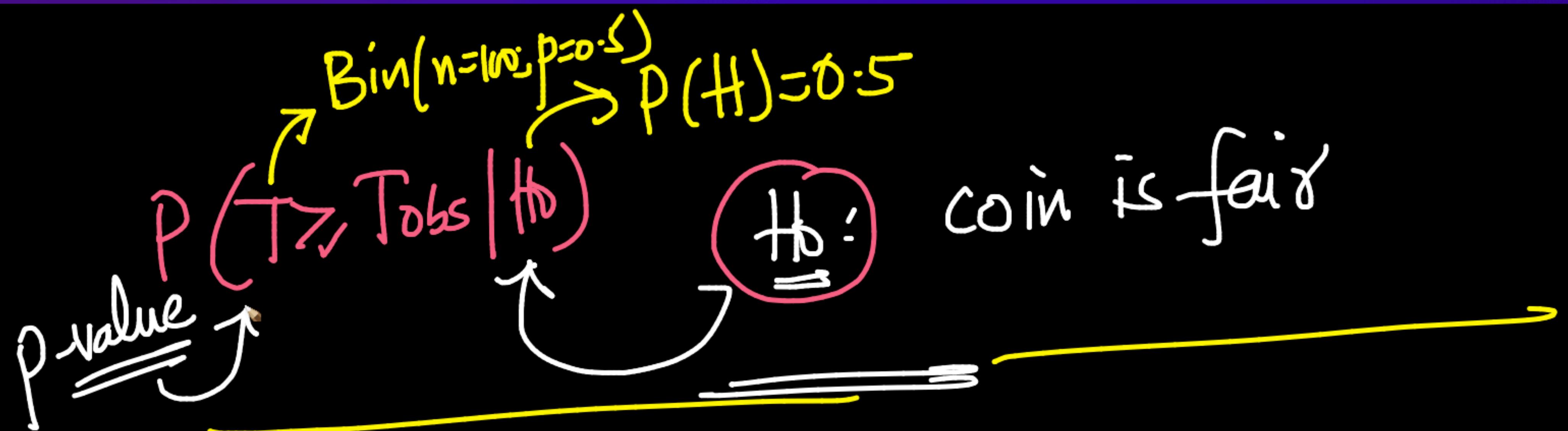
B:

100  $\rightarrow$  top 10

$t'_1, t'_2, \dots, t'_{10}$

{ H: Mean watchtime A = mean  $WT_B$   
i Ha! " "





$$P(T > T_{obs} | H_0)$$

$H_0:$  coin is biased towards heads

Bin(n=100; p=?)  
 $\rightarrow$  is impossible

$$\left\{ \begin{array}{l} P(H) = 0.5 \\ 0.52 \\ 0.53 \\ \dots \end{array} \right.$$

∞-many

$$P(T > \tilde{5} | H_0) = 0.46 - 46\% > 5\%$$

$H_0$

46%

Can't reject

reject  $H_0 \rightarrow$  accept  $H_a$

can't reject  $\underline{H_0} \Rightarrow$  accept  $H_0 \&$   
reject  $H_a$

$$\Pr(T \geq T_{\text{obs}} \mid H_0) = p\text{-val}$$

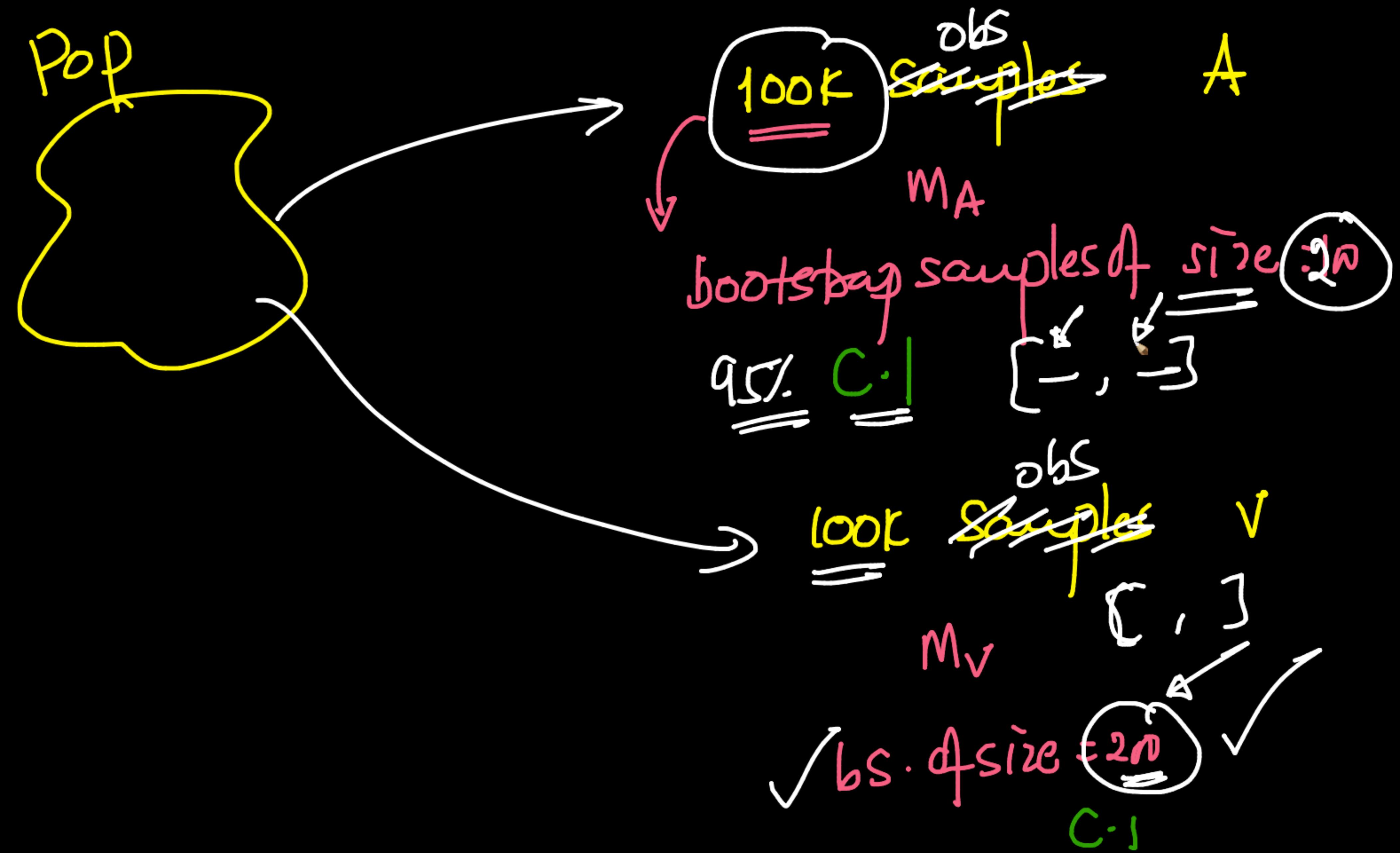
↓

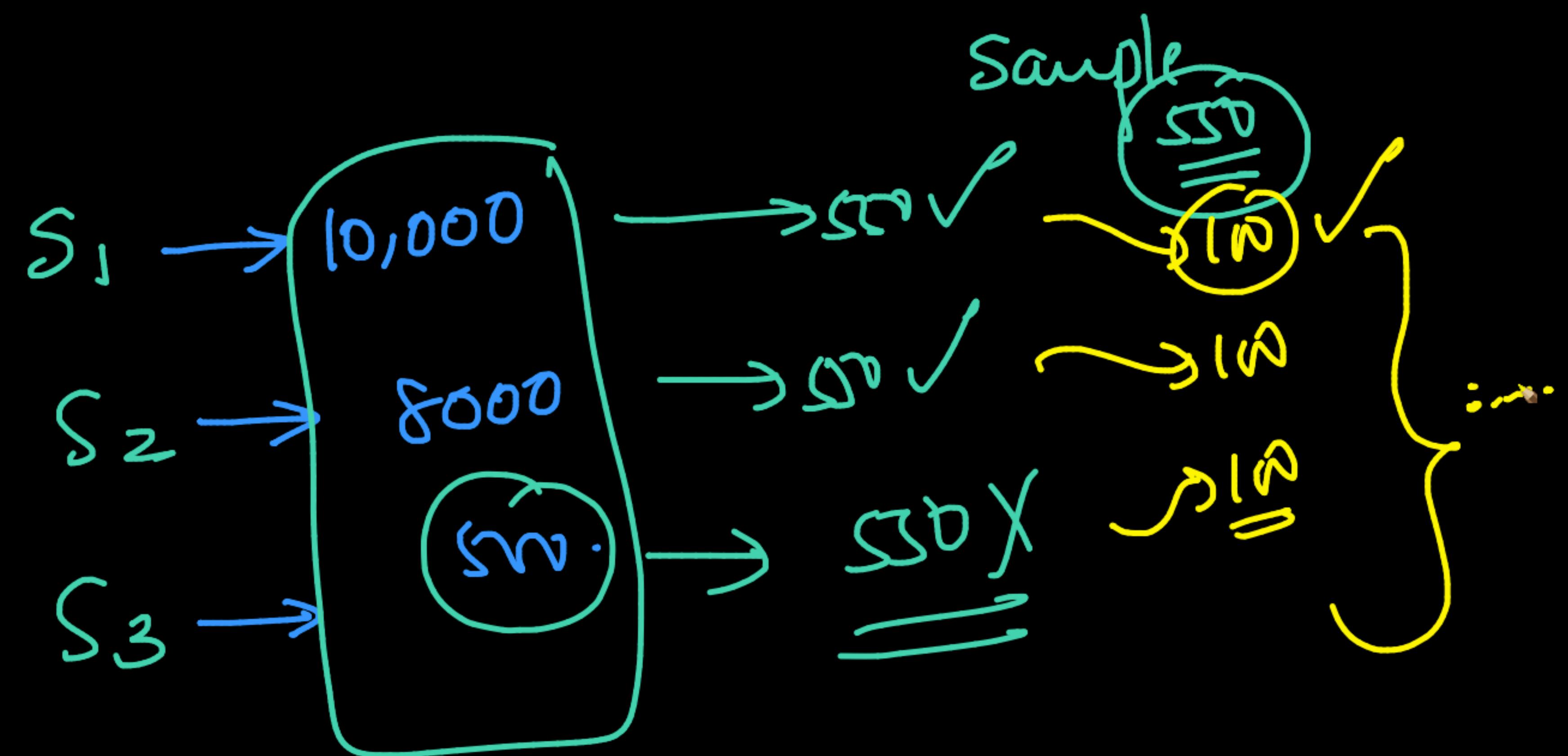


Significance level

reject  $H_0 \Rightarrow$  accept  $H_a$

o/w:- accept  $\underline{\underline{H_0}}$





conversations about popular titles. Second, we can help members choose some great content to watch by fulfilling the intrinsic human desire to be part of a shared conversation.

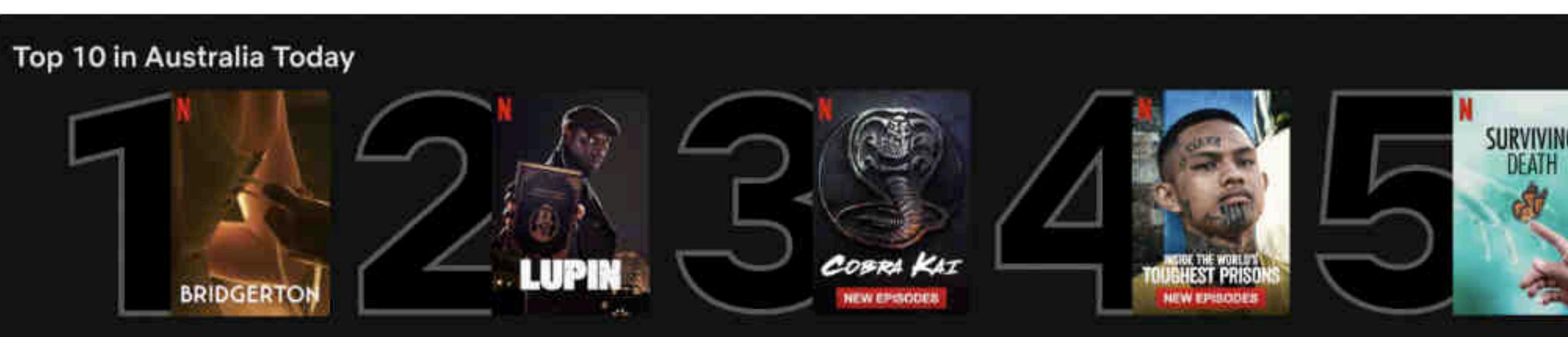


Figure 5: An example of the Top 10 experience on the Web UI.

We next turn this idea into a testable hypothesis, a statement of the form “If we make change X, it will improve the member experience in a way that makes metric Y improve.” With the Top 10 example, the hypothesis read: “*Showing members the Top 10 experience will help them find something to watch, increasing member joy and satisfaction.*” The primary decision metric for this test (and many others) is a measure of member engagement with Netflix: are the ideas we are testing helping our members to choose Netflix as their entertainment destination on any given night? Our research shows that this metric (details omitted) is correlated, in the long term, with what members will retain their subscriptions. Other areas of the business in which we run tests, such as the

## Netflix Technology Blog

265K Followers

Learn more about how Netflix designs, builds, and operates our systems and engineering organizations

Follow



## More from Medium

Netflix Technology ... in Netflix Tech...

**Experimentation is a major focus of Data Science across Netflix**



Diana Jerman in disney-streaming

**Disney Streaming Embraces 3 Key Tenets of Experimentation**



Clare La... in Expedia Group Technol...

**2021's Top 10 Posts on Expedia Group Technology**



Zhiheng Xu in The Airbnb Tech Blog

**Intelligent Automation Platform: Empowering Conversational AI and Beyond at Airbnb**

