

Applied Statistics:

Statistical Testing & Experiments

Welcome & Thank you

Four Sessions

1) 9:00 AM - 11:00 AM

2) 11:30 AM - 1:00 PM

3) 2:30 PM - 4:30 PM

4) 5:00 PM - 6:30 PM

(Q) will these videos be available later?

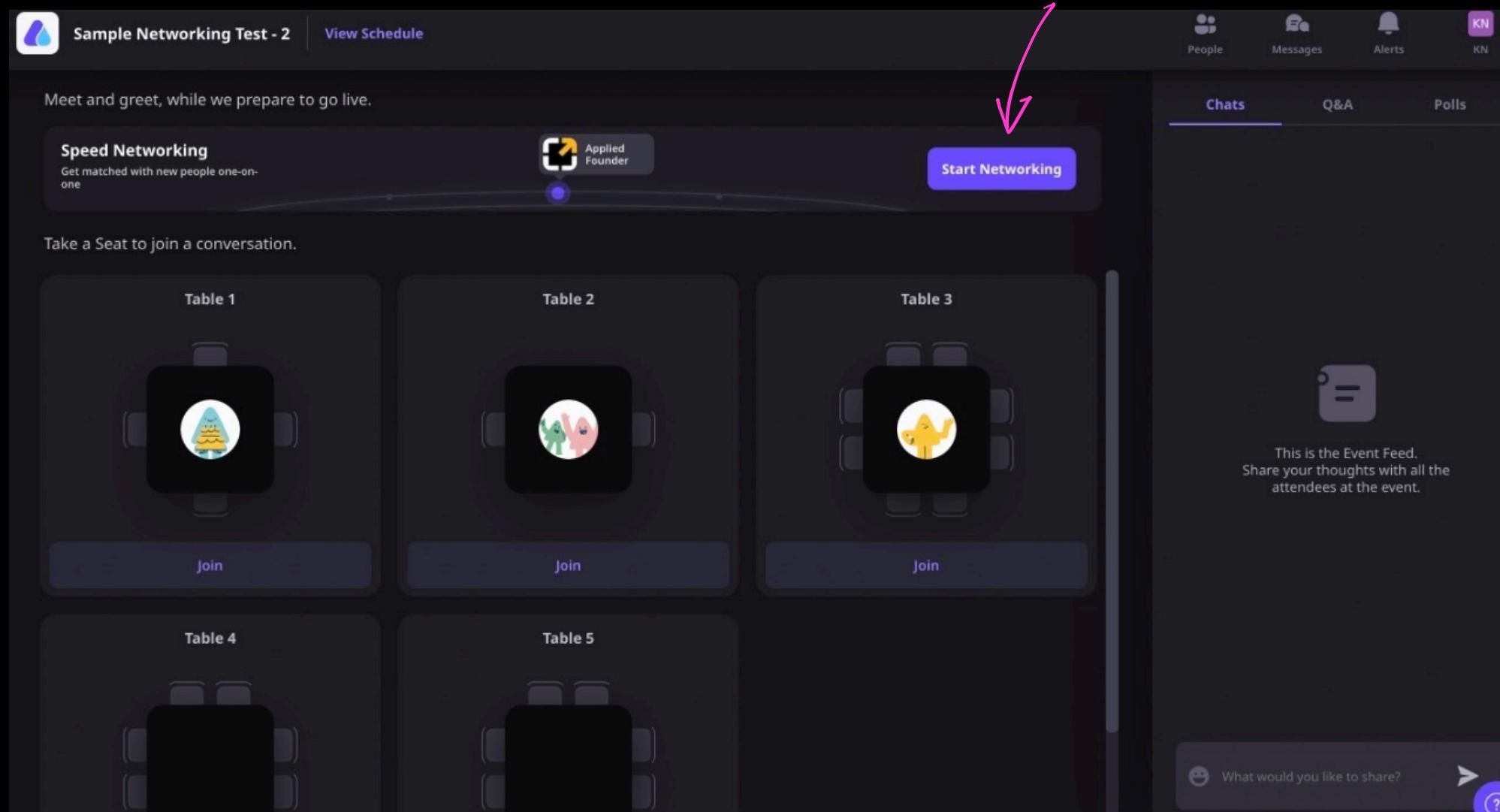
- Yes, as shorter videos over next few days
- YouTube; AppliedAI; Diploma

Pre-reqs:

- ① Probability (11 & 12th class)
- ② Python programming basics

Interactive - Learning :

- Quiz Question (timed: 2min / 5min)
↳ scenario-based ; Math ; code ; MSQ
- Think - Pair - Share
↳ (Speed Networking)



Sample Networking Test - 2 | View Schedule

Meet and greet, while we prepare to go live.

Speed Networking
Get matched with new people one-on-one

Applied Founder

Start Networking

Chats Q&A Polls

Take a Seat to join a conversation.

Table 1 Table 2 Table 3

Join Join Join

Table 4 Table 5

This is the Event Feed.
Share your thoughts with all the attendees at the event.

What would you like to share? 

Instructor - Participant interaction:

- "Q&A"; voting (end of each session)
- "Raise hand" + Mic

Any Technical glitches?

↳ Chat window

→ We will reconnect shortly

Participants from various backgrounds

↳ We will try & accommodate

Questions in Q&A

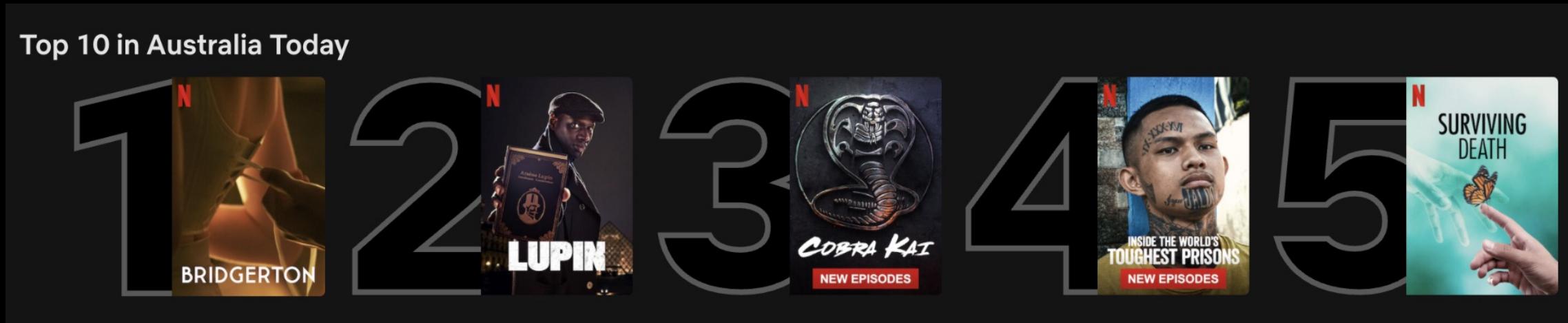
↳ organization

Statistical Testing

8

Experiments

Real-world Problem @ Netflix



<https://netflixtechblog.com/what-is-an-a-b-test-b08cc1b57962>

Does showing Top-10 (in a country)
impact time spent on Netflix

Netflix

↳ more time spent \Rightarrow more renewals

↓

more revenue

A/B tests

↳ Randomized Controlled Trials (RCT)

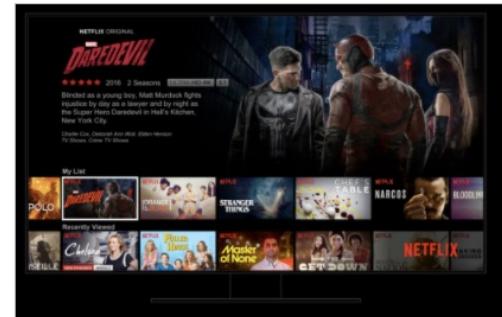
→ bucket testing

→ Split-run test

Netflix Members

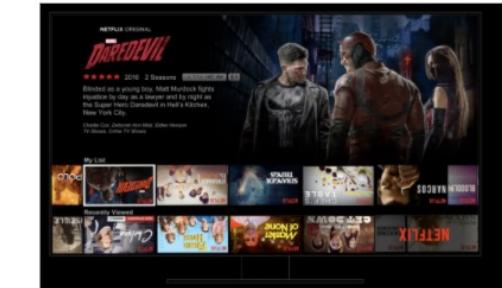


Version 'A' (Control)



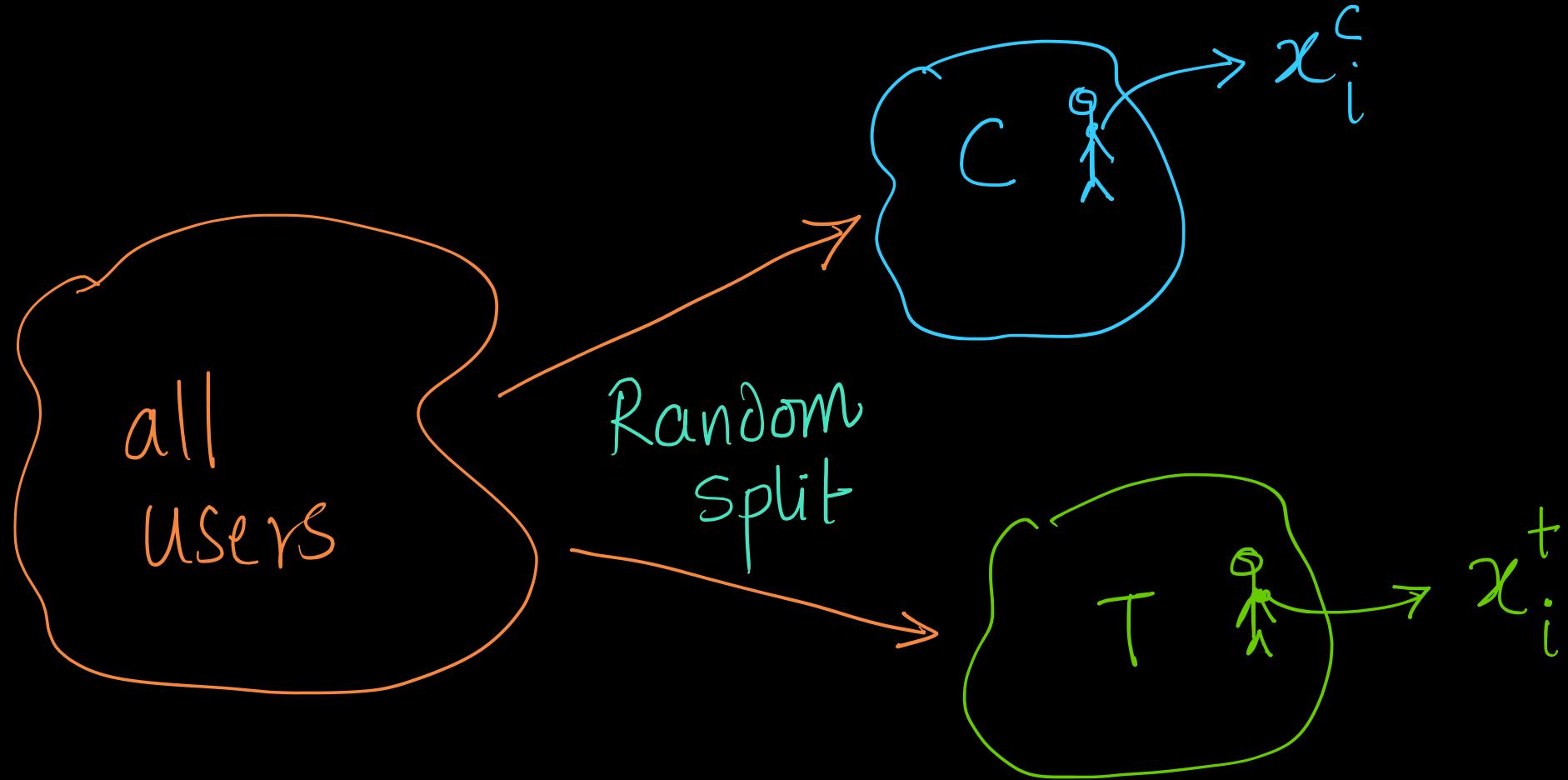
No Top-10

Version 'B' (Test)



Compare member behavior

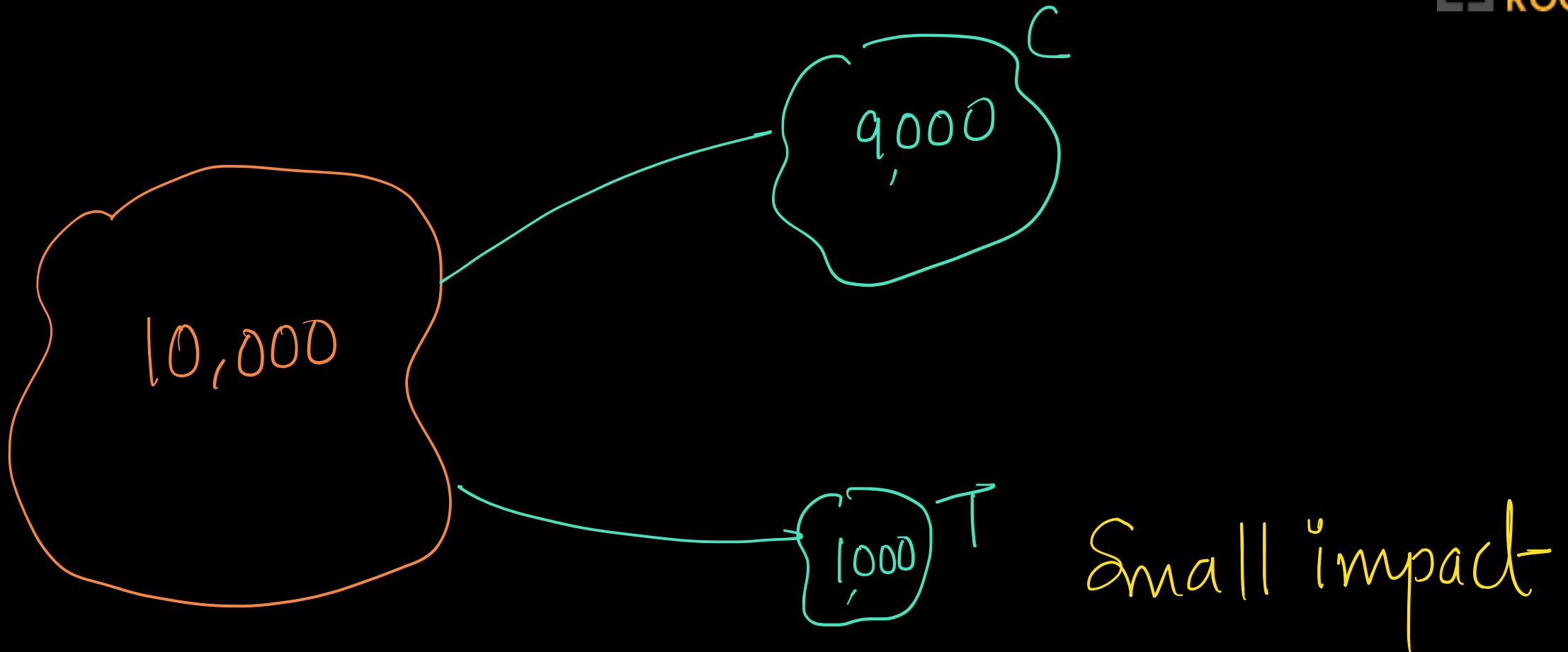
Top-10

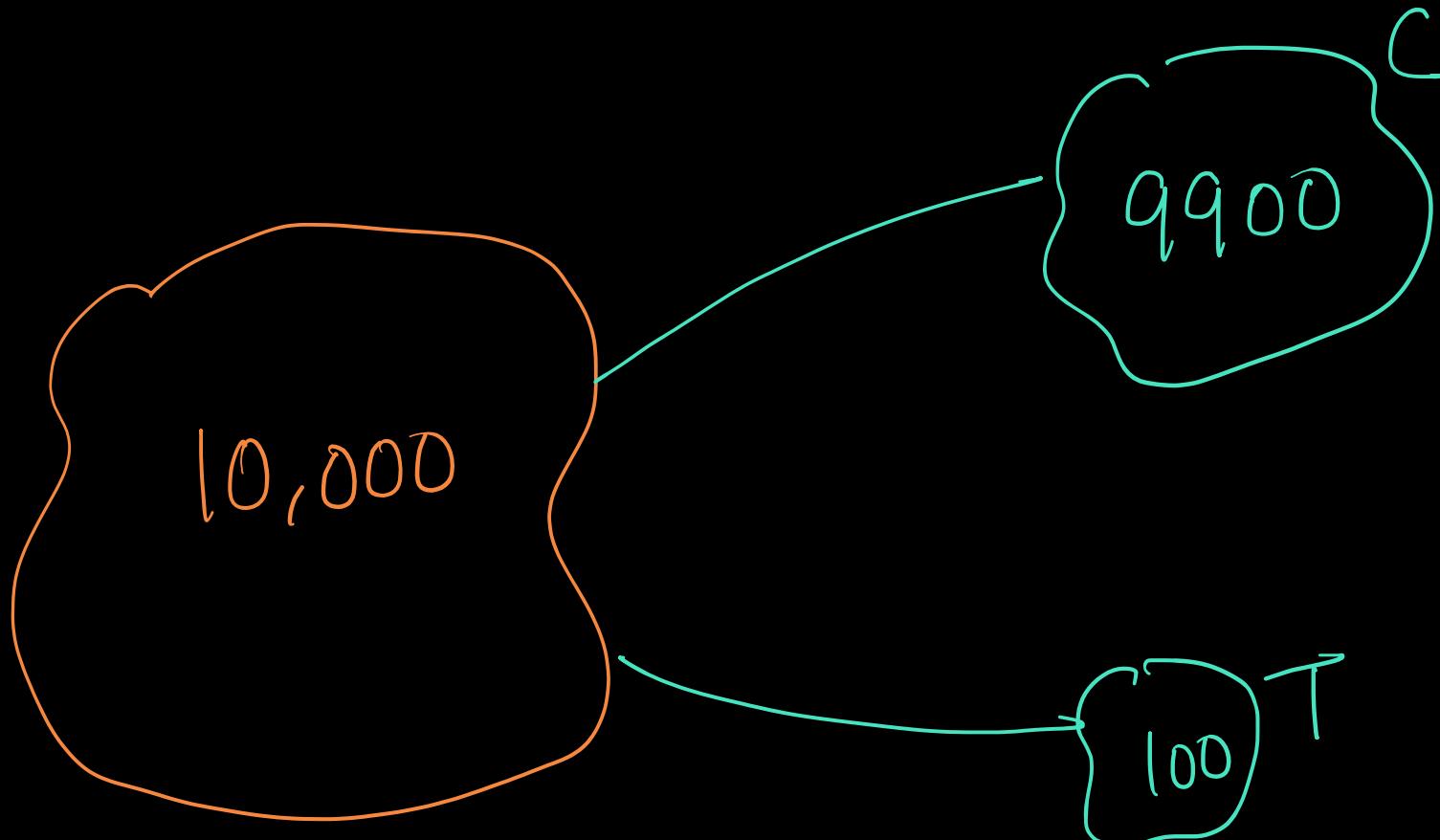


x_i^c : time spent on Netflix

Cost of experimentation

Treatment (B) could be terrible





minimal
impact

(Q)

How do we randomly split users into 90:10 ratio using code

→ Write the pseudo-code & paste it in chat [3 min]

(sol)

For each user u :

$\gamma = \text{random}()$ # uniform $\gamma \cdot v$

if $\gamma \leq 0.9$:

assign u to control

else:

assign u to treatment

collect data:

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

$X^t: x_1^t, x_2^t, \dots, x_m^t$

n & m can be

different

Is x_t "greater" than x_c ?



What does this mean?

①

$$\text{mean}(x^t) > \text{mean}(x^c)$$

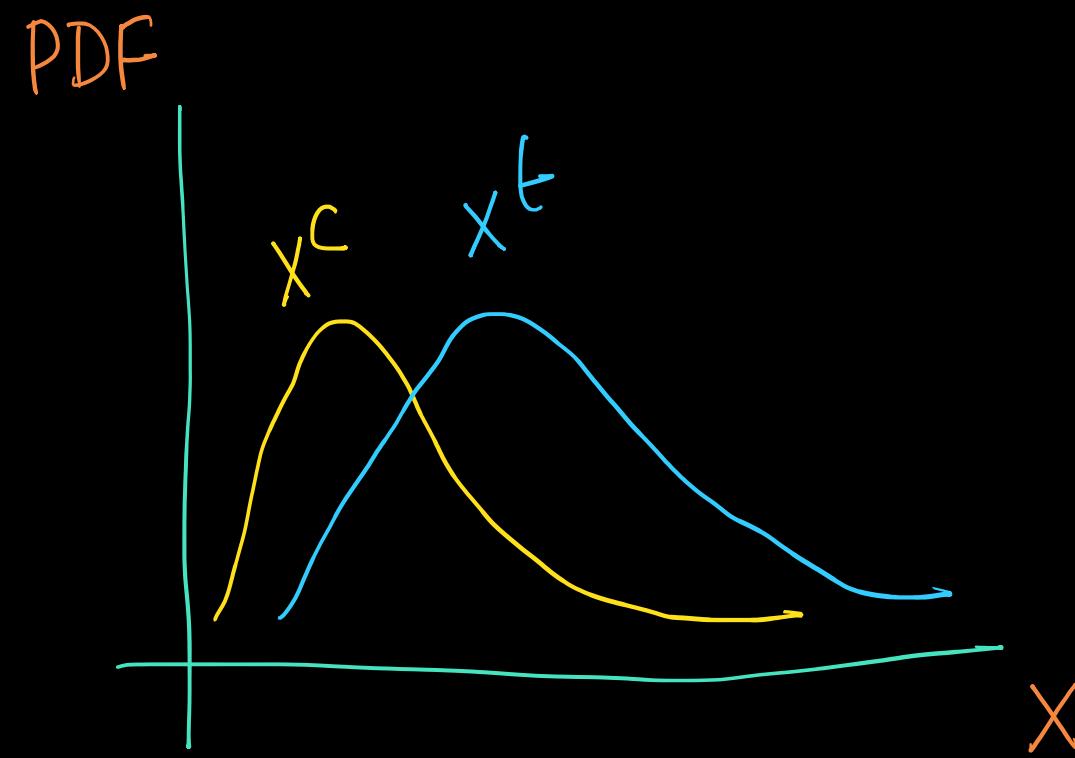
②

$$\text{median}(x^t) > \text{median}(x^c)$$

③

$$\text{tp99}(x^t) > \text{tp99}(x^c)$$

4



many possible & valid methods

Simplest:

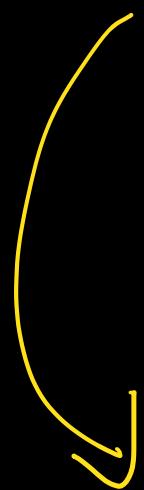
$$\text{median}(x^t) \text{ vs } \text{median}(x^c)$$

④ robust to outliers

empirical observation:

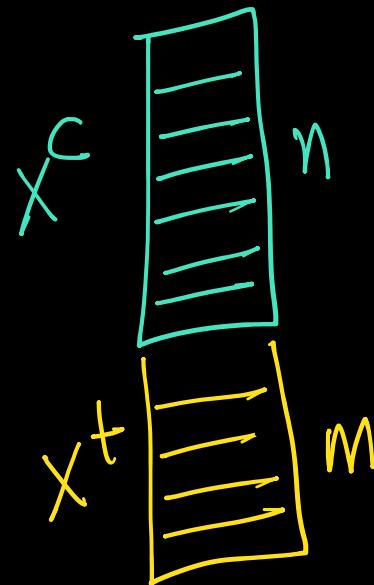
$$\text{median}(x_t) - \text{median}(x_c)$$

$$= 20 \text{ minutes (let)}$$

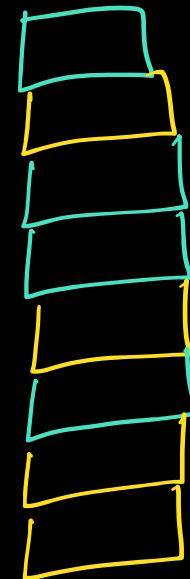


median-diff-observed

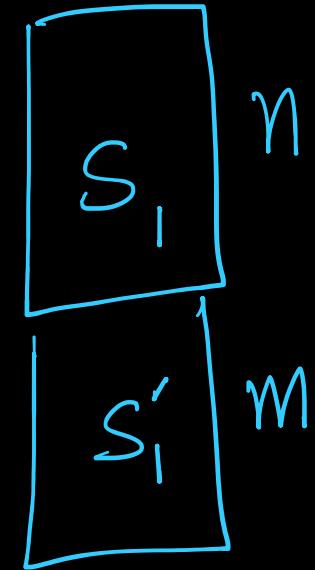
Permutation testing:



Combine



shuffle
/Permute



Split 1

$$\text{median}(S'_i) - \text{median}(S_i)$$

$$\text{median} - \text{diff-1}$$

why combine, shuffle & split

→ simulating if Control & Treatment
(*)
data come from the same
distribution

median-diff - 1

One possible median-diff if
C & T are the same with
no difference

Repeat k -times

median-diff-1 (M_1)

median-diff-2 (M_2)

.

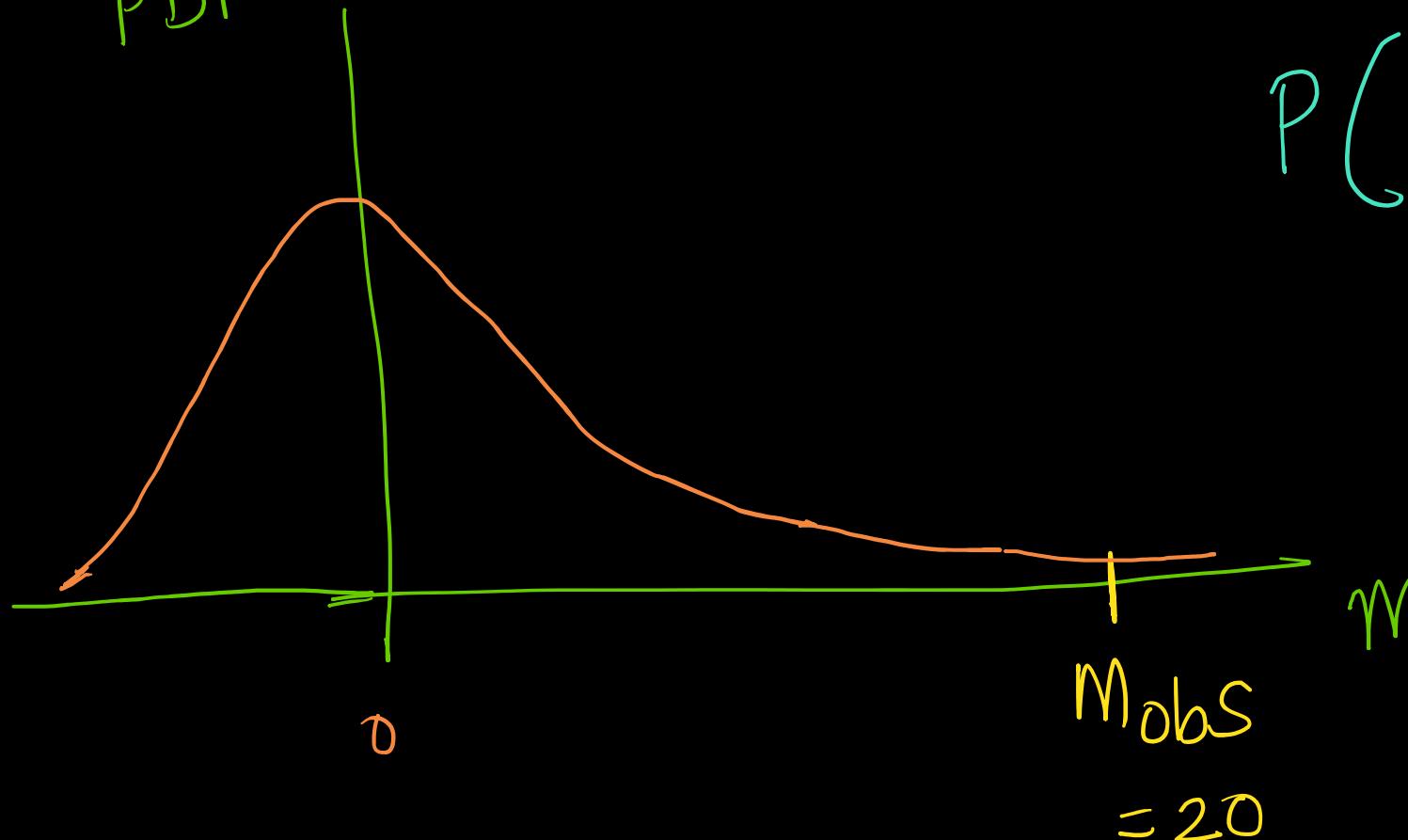
:

:

median-diff- K (M_K)

median-diff-
observed
(M_{obs})

PDF



$$P(M \geq M_{obs})$$

$$= 0.02$$

$$M_{obs}$$

$$\approx 20$$

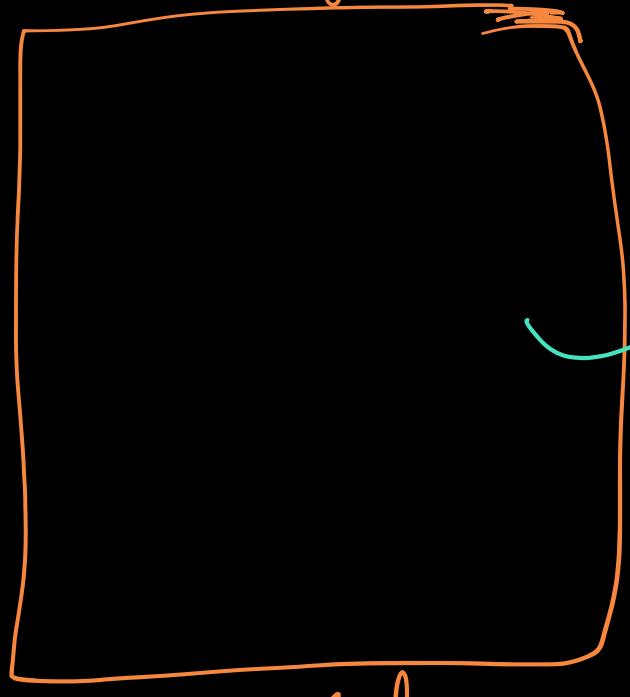
$$M$$

Conditional-prob:

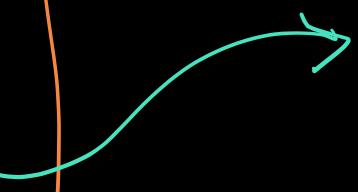
$P(\text{observing a median-diff} \geq m_{\text{obs}} \mid$
no difference in C8T)

$$= 0.02 = 2\%$$

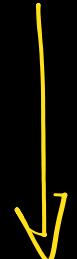
No diff b/w C & T



w odd

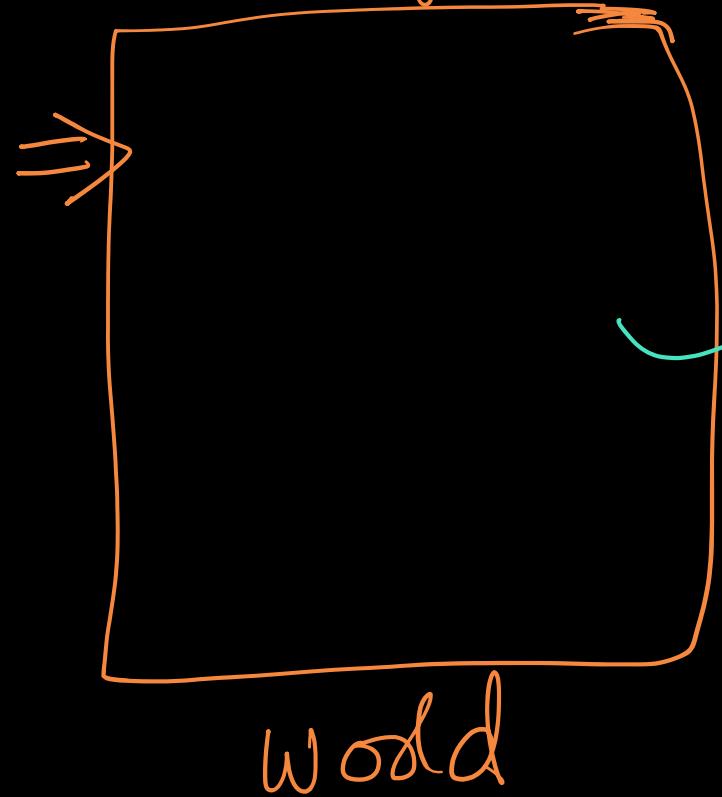


$P(m \geq 20)$ is very
small (2%).



But we observed
this in our
A/B test

No diff b/w C & T



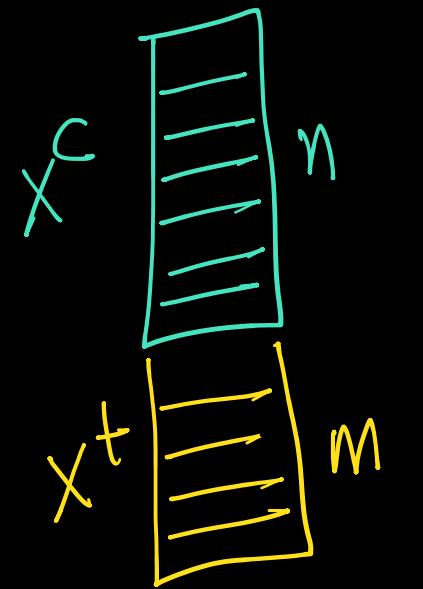
this word is not
likely

\Rightarrow CGT's medians differ

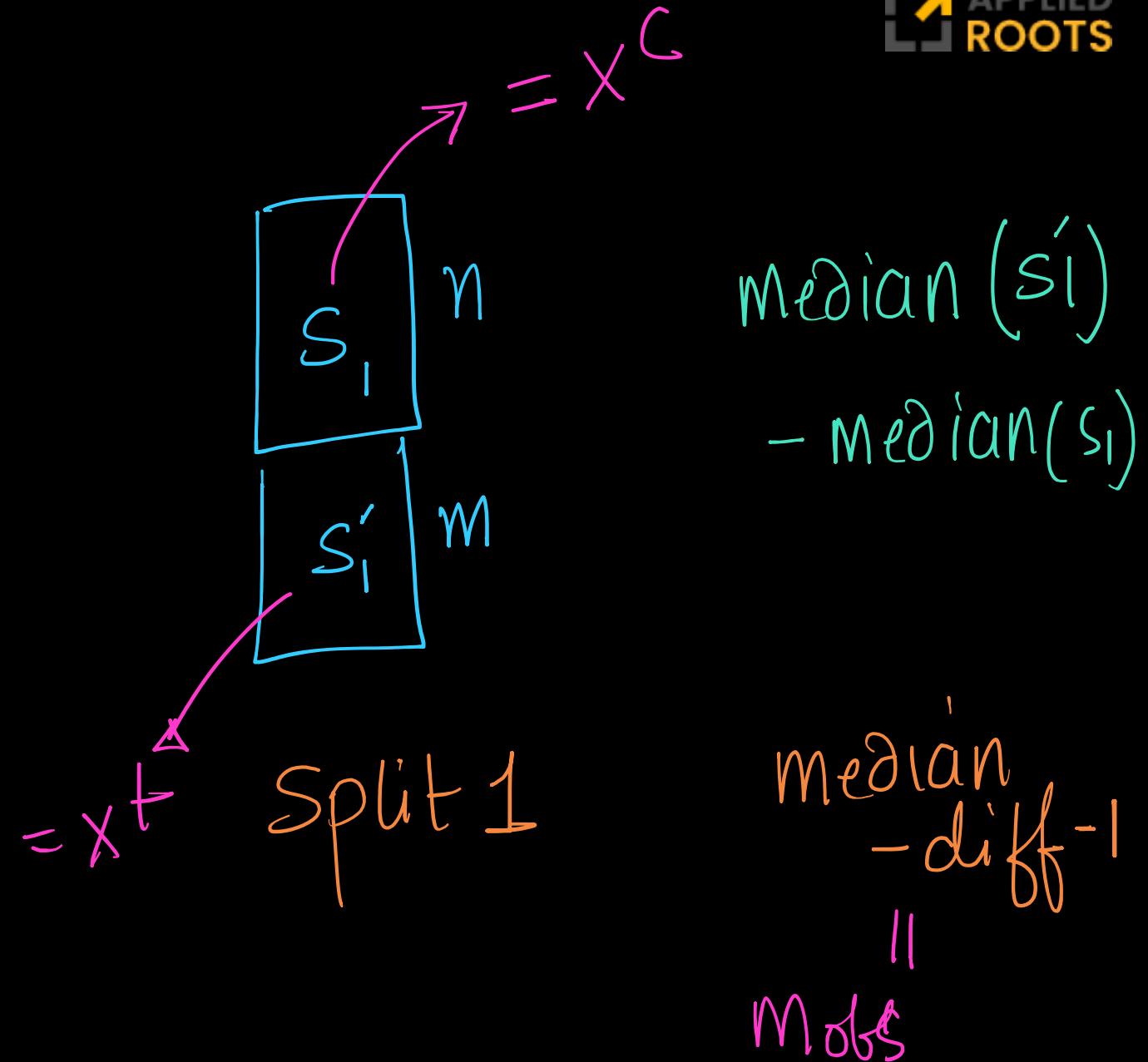
(Q) what would happen if we
don't shuffle / permute

→ provide an answer with justification

No-shuffle



Combine



(Q) How would you shuffle in code?

→ provide the idea/pseudo-code
DO NOT give the python function.

1, 2, 3, ----- n+m



Permutation (using recursion)



3, 2, 6, 4, -----

Hypothesis Testing (Terminology)

H_0 : Null Hypothesis

$\hookrightarrow X^C$ & X^T are the same

H_a : Alternative hypothesis

$$x^T > x^C$$

Test statistic : $\text{median}(x^t) - \text{median}(x^c)$

p-value:

$P \left(\text{observing a value as extreme} \mid \text{as the empirically observed value} \right) \mid H_0$

$$\hookrightarrow P(M \geq M_{\text{obs}} \mid H_0) = 2\%$$

Significance-level: (α)

α : 5% (typically)

if p-value < α :

reject H_0

else Accept H_0

(e.g.)

$$P(m > M_{obs}) = 0.02 < \alpha = 5\%.$$



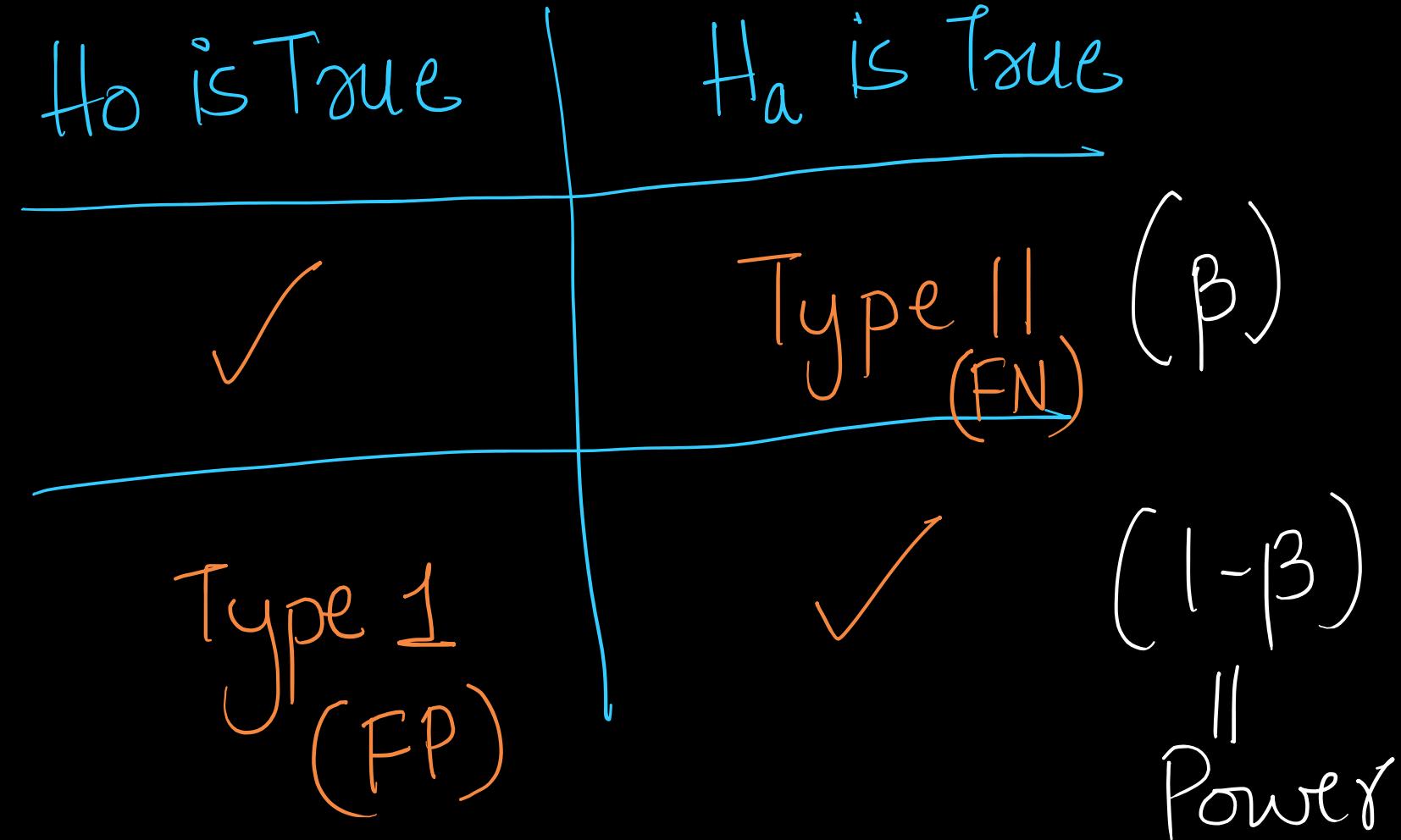
reject H_0 @ 5% significance level

* Accept H_0 @ 1% significance level

Errors:

accept H_0

accept H_a



Power

↳ how often do we accept H_a
when H_a is true

(Q) Why did we choose the
Ho as $X^c \approx X^t$?

→ Provide short explanation in the chat -

if we choose

$$H_0: X^c \neq X^t$$

→ how do we simulate this H_0 ?

$$H_0 : X^c \approx X^t$$



Can be simulated using Permutation
+ Resamplings

Code for Permutation Testing

<https://colab.research.google.com/drive/1Lp5LW2TiJtSFBZd37uqMmAghpMLvQRp?usp=sharing>

(Q) what happens if $M = N = \text{Small}$
 (100)

& everything else stays the same

$d = 5\%$

$K = 100$

(Q) What happens when k-increases
to 100?

$$M = N = 100$$

$$\alpha = 5 \%$$

(Q) What if

$$M = 10000$$

$$N = 1000$$

$$\alpha = 5\%$$

$$K = 100$$

(Q) Medical domain

$$m \approx n \approx 30$$

(Q) if $n = 30; m = 20$

what is the maximum value of K ?

$$n+m = 50$$

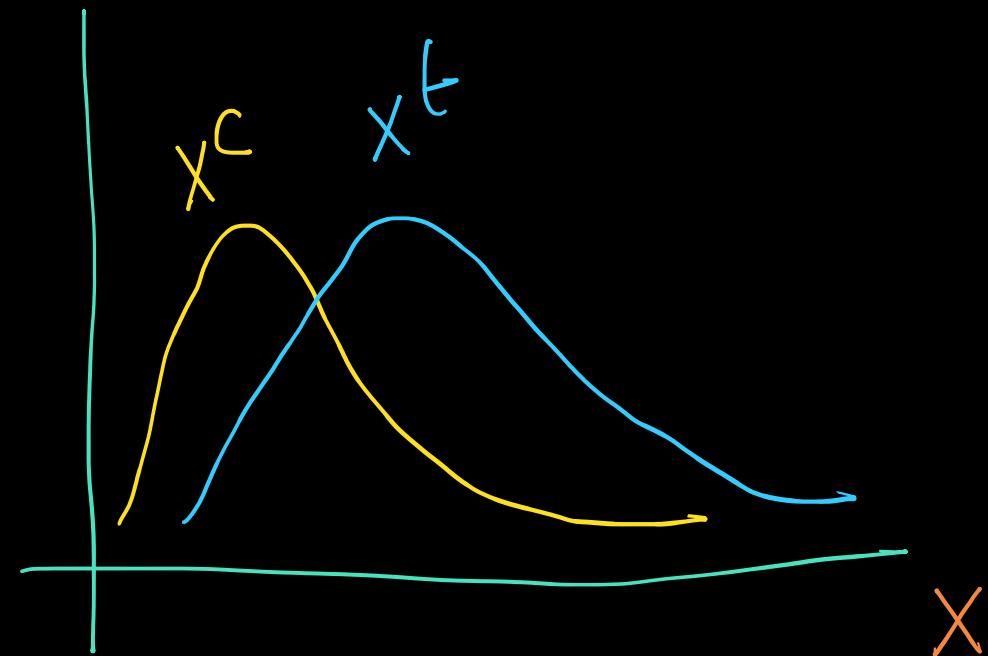
— — — — — →
1 2 3 4 - - - - 50

50!

No - test is perfect

Permutation test cannot compare

PDF



Assumption :

Observations are exchangeable

True for Randomized trials

- Super easy to experiment with
- No distributions or complex assumptions
- Test - statistic flexibility

Variations:



default



new

which art-work generates
more clicks?

<https://netflixtechblog.com/selecting-the-best-artwork-for-videos-through-a-b-testing-f6155c4595f6>

- Binary variable 1: click 0: no click

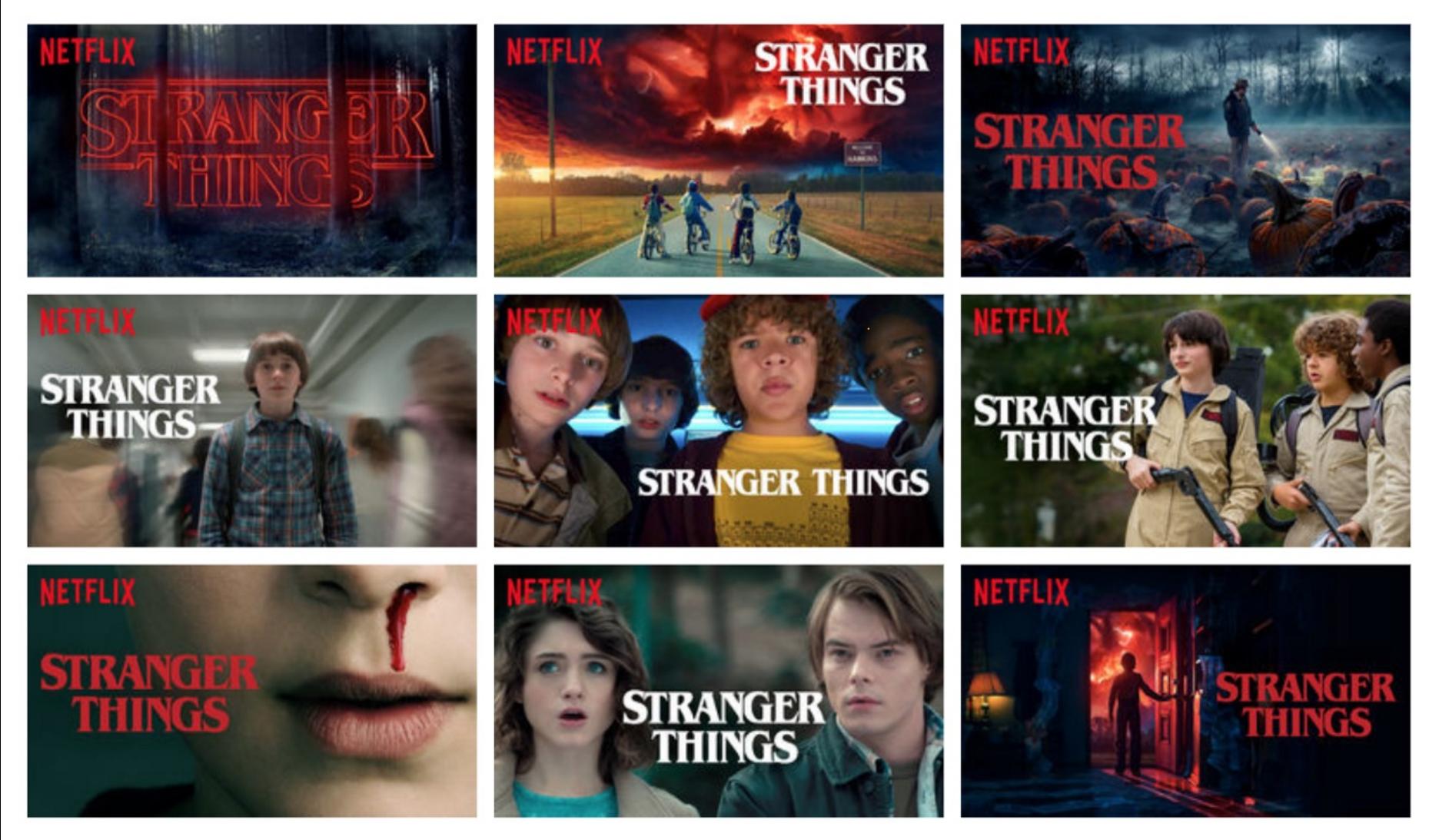
C	1	}	n
C	0		
:	:		
C	1		
T	1		
T	1	}	m
T	0		
T	0		
T	0		

Test - statistic

↳ click Through rate

Permute - resample - compare

What if we have 9 alternatives ?



- A/B Tests can be wasteful

$$9 \times 5\% \Rightarrow 45\%$$

[Multi-arm bandits (alex)]

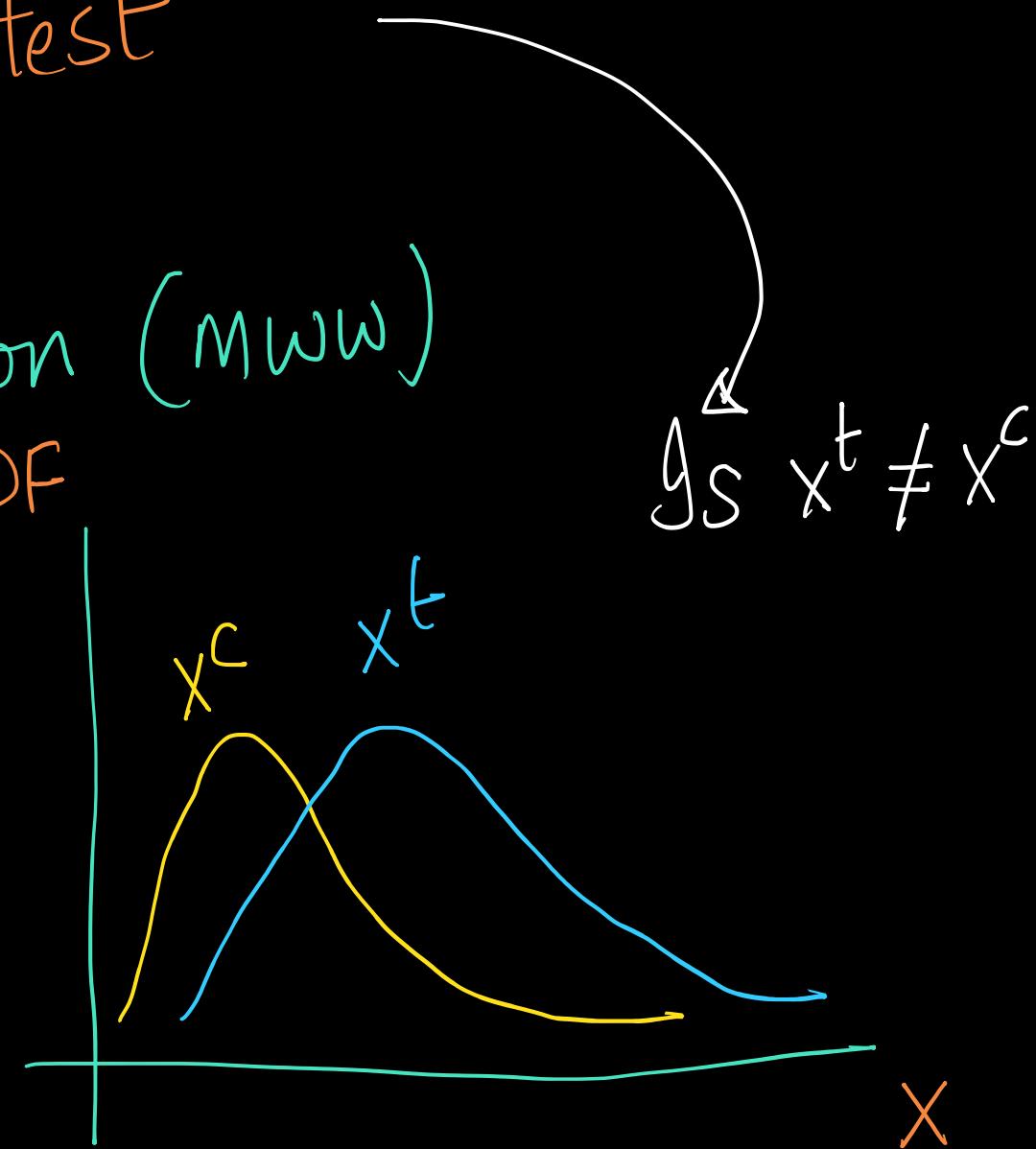
Mann - Whitney - U test

- Mann - whitney - wilcoxon (MWW)

PDF

- rank - sum test

- very - popular



Assumptions

- ① x^c & x^t 's observations are independent of each other
- ② x_i^c & x_i^t values are ordinal
 $<, =, >$

③ $H_0: X^c \approx X^t$

$$H_a: X^c \neq X^t$$

test - statisic (u)

$$u = \sum_{i=1}^n \sum_{j=1}^m s(x_i^c, x_j^t)$$

$$s(x_i^c, x_j^t) = \begin{cases} 1 & \text{if } x_i^c > x_j^t \\ 1/2 & \text{if } x_i^c = x_j^t \\ 0 & \text{if } x_i^c < x_j^t \end{cases}$$

$\hookrightarrow S_{ij}$

- compose every pair x_i^c & x_j^t

x^c	x^t
1	1
6	6
9	9

$$\begin{aligned}
 & \left(\frac{1}{2} + 0 + 0 \right) \\
 & + \left(1 + \frac{1}{2} + 0 \right) = 4.5 \\
 & + \left(1 + 1 + \frac{1}{2} \right) = \frac{m \times n}{2}
 \end{aligned}$$

compose every pair x_i^c & x_j^t .

x^c	x^t
1	2
6	12
9	18

$$\begin{aligned}
 & (0 + 0 + 0) \\
 & + (1 + 0 + 0) = 2 \\
 & + (1 + 0 + 0)
 \end{aligned}$$

compose every pair x_i^c & x_j^t .

x^c	x^t
1	0
6	1
9	2

$$\begin{aligned}
 & (1+1/2+0) \\
 & + (1+1+1) \\
 & + (1+1+1)
 \end{aligned}
 = 7.5$$

if $n & m \geq 20$; under H_0

$$U \sim \text{Normal}(\mu_U, \sigma_U)$$

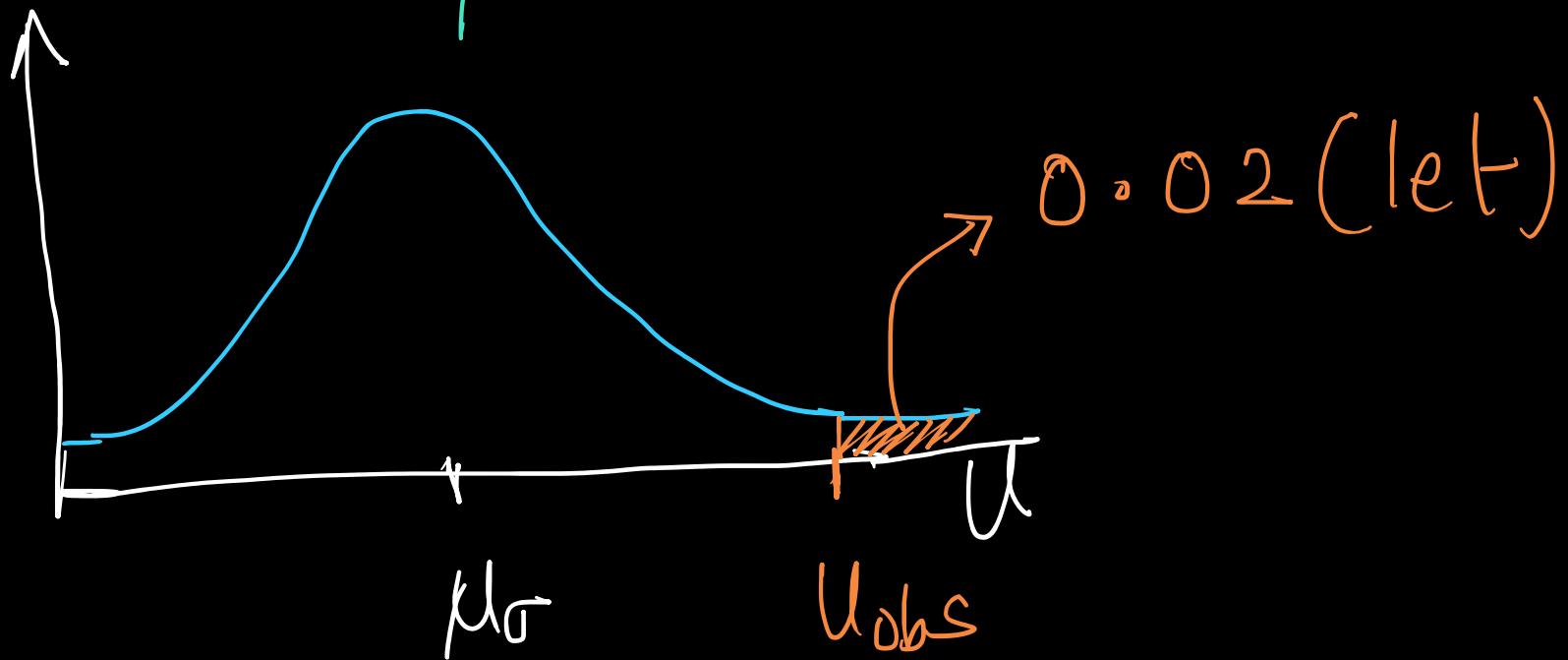
$$\frac{n+m}{2}$$

$$\sqrt{\frac{nm(n+m+1)}{12}}$$

Hypothesis testing:

- ① $H_0: X^c \approx X^t$ $H_a: X^c \neq X^t$
- ② Compute U_{obs}

③

 $U \sim \text{Norm}(\mu_U, \sigma_U)$ under H_0 

(4)

if $\alpha = 5\%$ reject H_0

Note :

→ We know the distribution of

test-statistic

$P(\text{test-statistic} > \text{obs} | H_0)$ is easy
to compute

Permutation testing

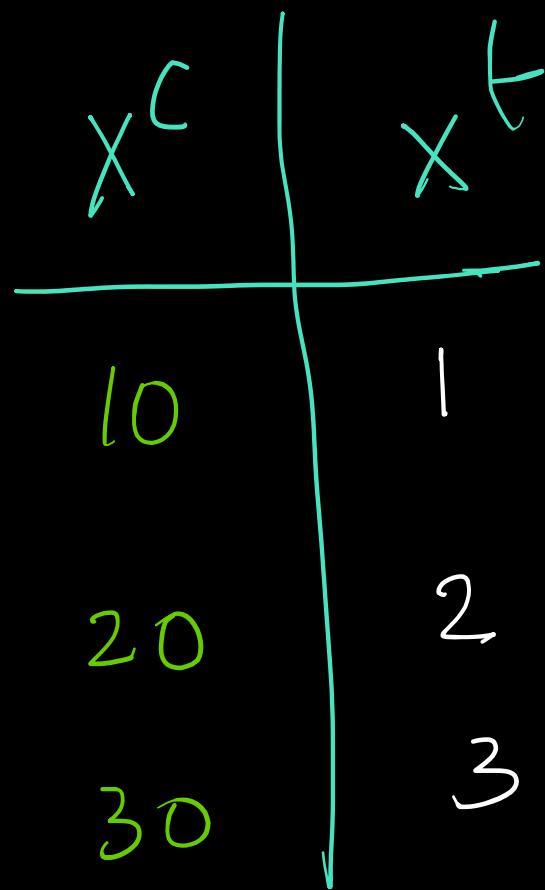
$$\hookrightarrow P(\text{test statistic} > \text{obs} | H_0)$$

Simulation
(resampling)

MWW

- widely used
- non-parametric
- every pair is compared
- cannot say $x^c > x^t$ or $x^t > x^c$

(Q) What is the max value of
test-statistic t & when does
that occur?



$$\begin{aligned} & (1+1+1) \\ & + (1+1+1) \\ & + (1+1+1) \\ & = 9 \end{aligned}$$

$m \times n$

Python - Code

(Q) Small sample sizes

Permutation vs MWU
(median)

Hypothesis Testing (recap)

- H_0 & H_a
- assumptions
- test-statistic & its distribution
- accept or reject

Multiple - Testing

$$X^1 \sim X^2$$

$\alpha = 5\%$

$$X^1 \sim X^3$$

n is large

$$X^1 \sim X^4$$

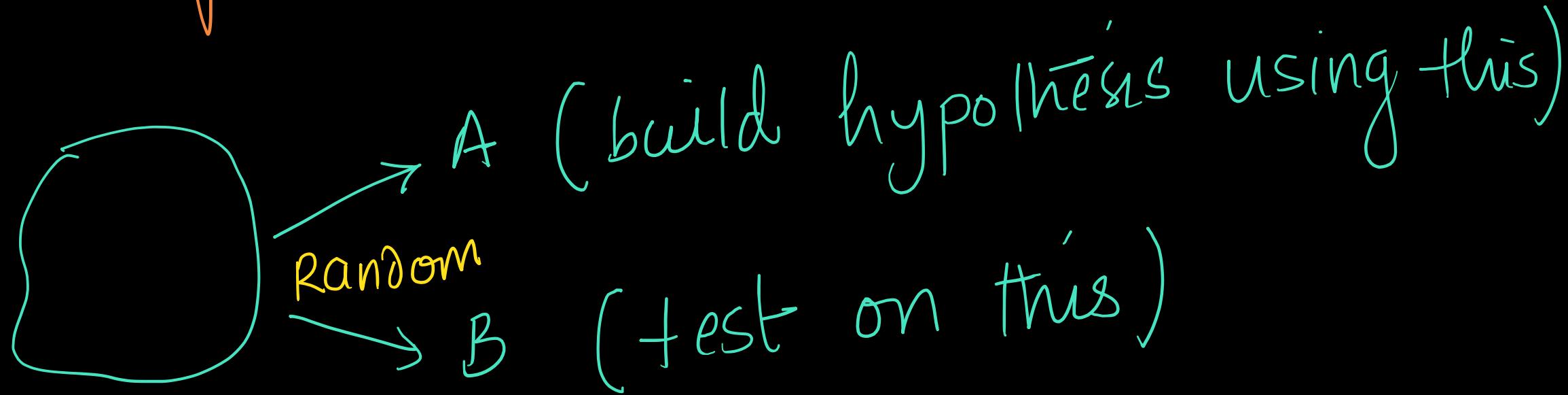
"ERRORS"

⋮

$$X^1 \sim X^n$$

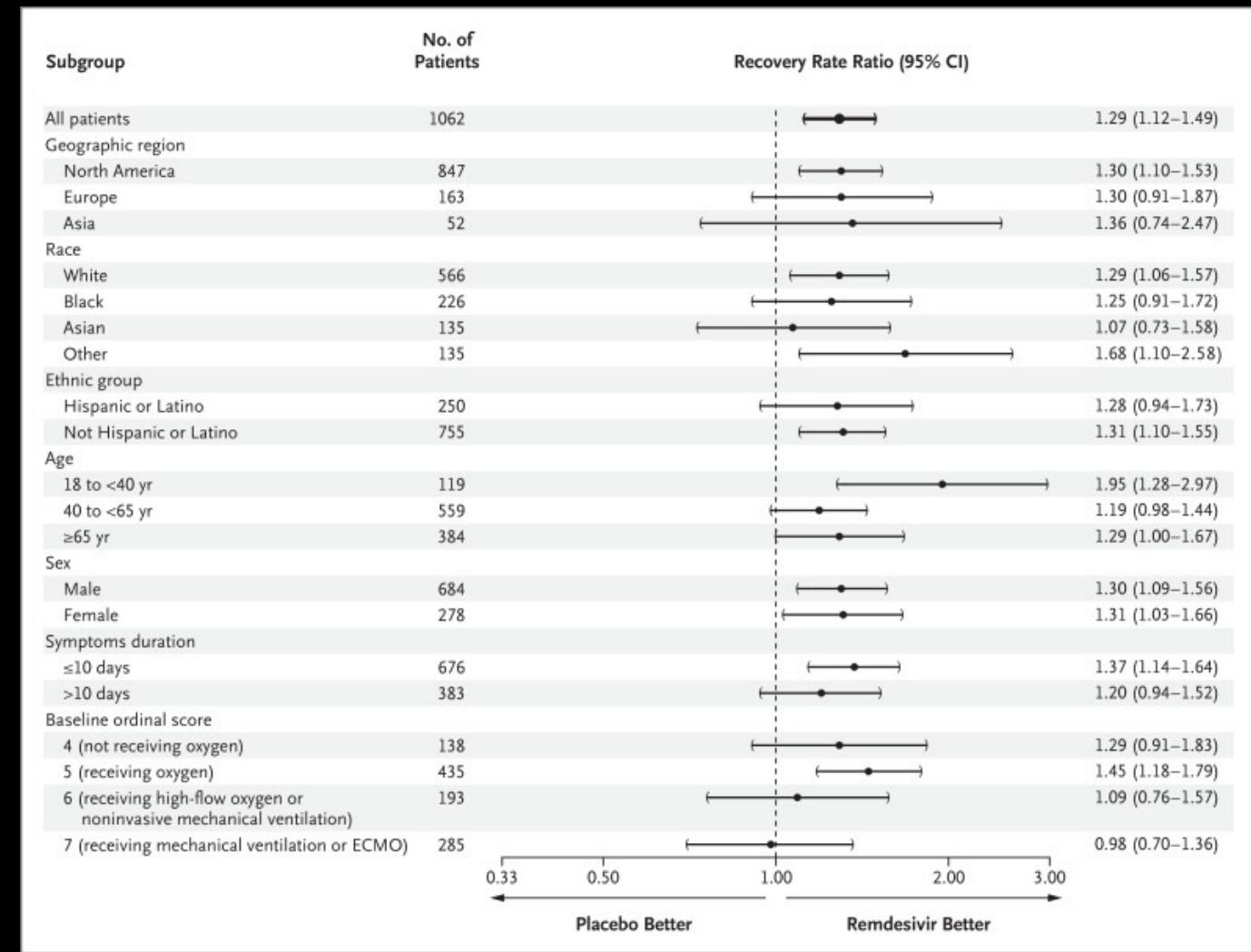
HOLM - BONFERRI correction

out-of-sample randomized tests



Bootstrapping & Confidence Intervals

<https://pubmed.ncbi.nlm.nih.gov/32445440/>



Why has health care moved to C.I.?

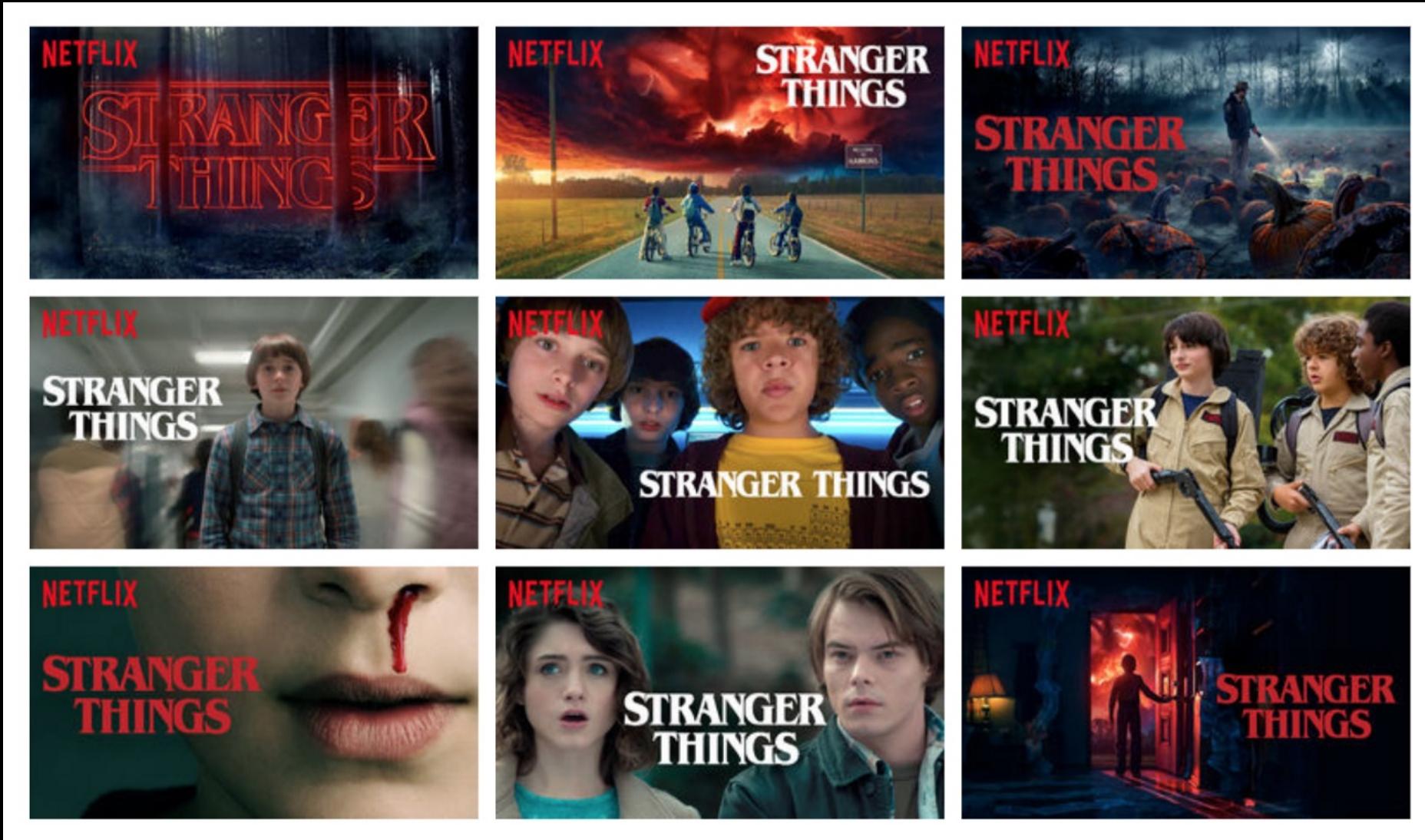
→ Intervals over point estimates



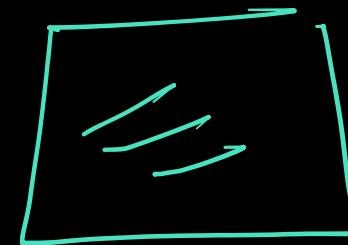
(distribution)

1000's of simultaneous experiments

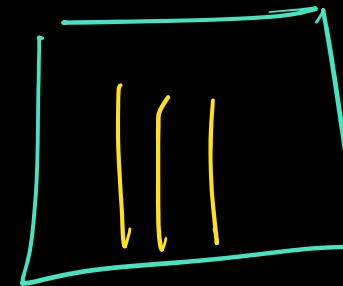
Multi-arm - Bandits



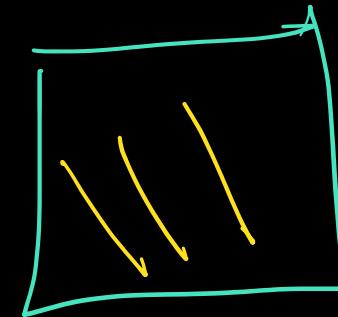
Google Ads



C1



C2



C3



<https://towardsdatascience.com/beyond-a-b-testing-multi-armed-bandit-experiments-1493f709f804>

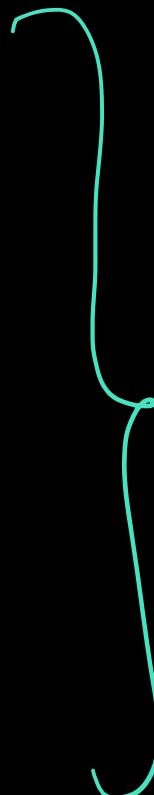
MAB

- ϵ -greedy
- Bayesian Thompson Sampling
- Contextual Bandits

Chi-Square Test

t-test

Dive-deep into
Multi-arm-bandits



future Sessions

Additional References for t-test & Chi-square test

<https://www.khanacademy.org/math/statistics-probability/significance-tests-one-sample#error-probabilities-and-power>

<https://www.khanacademy.org/math/statistics-probability/inference-categorical-data-chi-square-tests/chisquare-goodness-of-fit-tests/v/chi-square-distribution-introduction>

<https://youtu.be/WXPBoFDqNVk>

<https://machinelearningmastery.com/how-to-code-the-students-t-test-from-scratch-in-python/>

