

A/B Testing 3

Mini 4 / Spring 2024

Carnegie Mellon University
Tepper School of Business

THE INTELLIGENT FUTURE



GUEST SPEAKERS



- Mon, April 8
 - Guest speaker: **Sachal Lakhavani (AI Product Manager @ Meta; Tepper MBA alum)**
 - 30 min Q&A about A/B testing & Product management (Section A: 5-5:30pm, Section E: 6:30-7pm)
 - **Everyone is expected to prepare for at least one question (I might cold call!)**
 - 15 min 1:1 Q&A (Section A: 5:30-5:45pm, Section E: 7-7:15pm)
- Mon, April 15
 - Guest speaker: **Michael Degnan (Senior VP, Head of Enterprise Innovation @ PNC)**
 - 1 hour lecture on New Product Development & Innovation
 - 30 min Q&A about Innovation & Product management

FINAL PROJECT PRESENTATIONS



- Schedule
 - Section A: Wed, April 17 (Teams 1-3) & Mon, April 22 (Teams 4, 5, 7)
 - Section E: Mon, April 15 (all teams)
- Deliverables
 - 20 min presentation & slides (no min/max requirement)
 - Submit the slides to Canvas by the beginning of Wed, April 17 (Section A) or Mon, April 15 (Section E)
- Logistics
 - Each team should ask at least one question for other teams' presentations
 - **Attendance is expected for all teams' presentations**

FINAL PROJECT PRESENTATIONS



- Product selection
 - Approved on Canvas; Let me know if you are not happy with the current selection
- Required
 - What is the Problem a product manager is facing (WFH or GenAI)
 - Environment analysis (e.g., 3C)
 - Customer needs discovery (e.g., interview, UGC, survey; please define your target segment)
 - **(At least one) Feature proposal** (show us what features you considered; connection with customer needs is essential; please consider the AARRR framework)
 - Design of A/B testing to assess the proposed feature
- Optional
 - Frameworks from the pre-class videos (e.g., RICE)
 - Qualitative/quantitative insights (e.g., from market reports, data, your own experiences)

FINAL PROJECT PRESENTATIONS

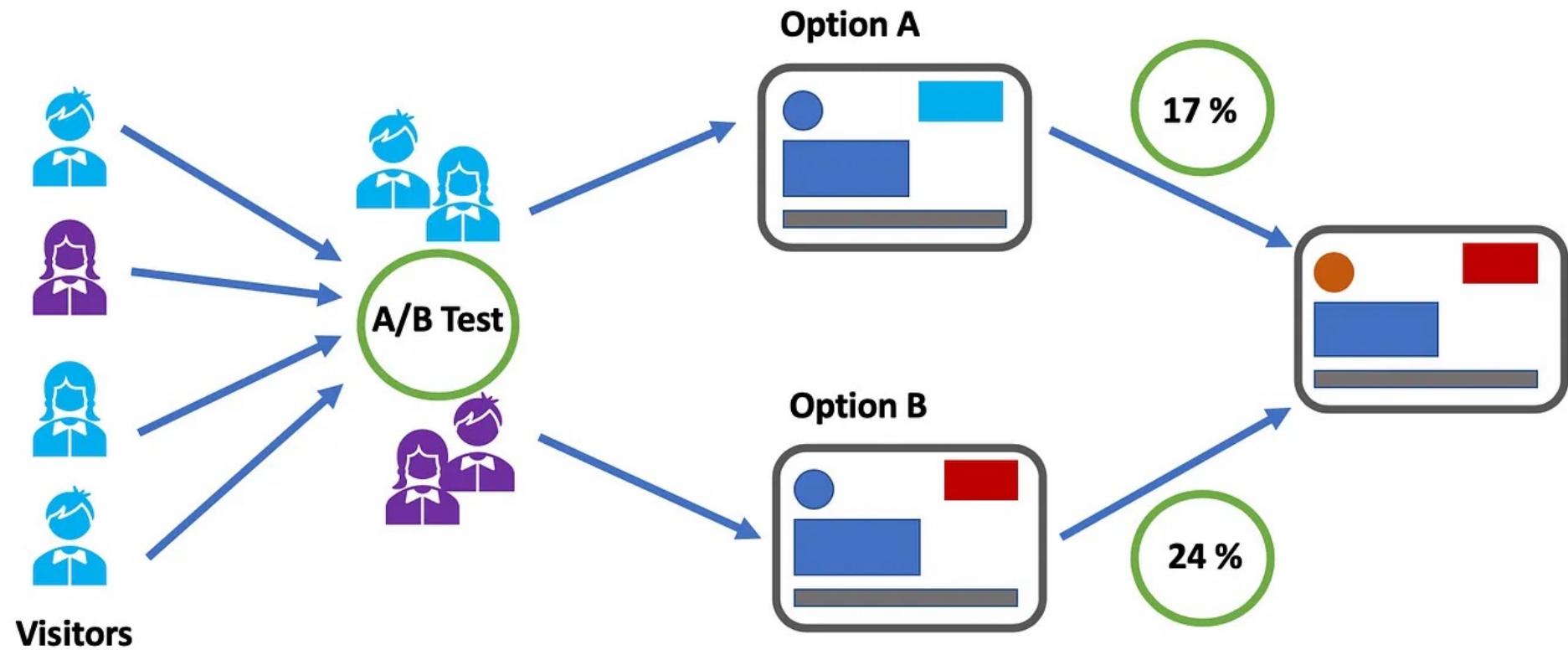


- Grading
 - Contents:
 - 1) Customer needs discovery
 - 2) Product feature
 - 3) A/B testing design (hypothesis, randomization, etc.)
 - Presentation: Clarity, Storytelling, Time limit (20 min)
 - Discussion: Q&A with the class
 - Teamwork: Peer evaluation of team members after the last day of class
 - All of you will give feedback to other teams and pick the best team(s) for bonus points (logistics to be elaborated on presentation days)

RECALL: LAST WEEK



- A/B Testing
 - Examples in tech firms: Amazon, Ebay, Netflix, etc.
 - Compare.com case: customer churn
 - How to define hypotheses for A/B testing
 - Uber pool case: how to randomize users into control and treatment conditions
 - Booking.com case: the role of culture and management, challenges in experimentation at scale

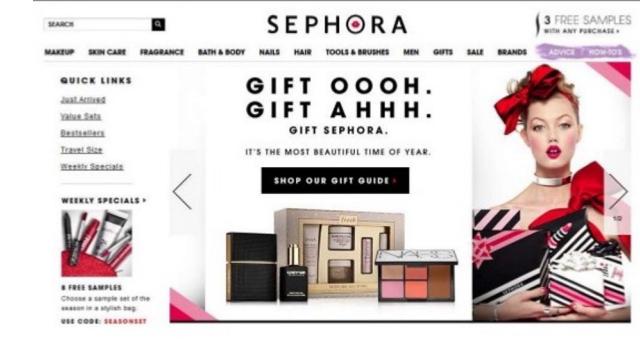




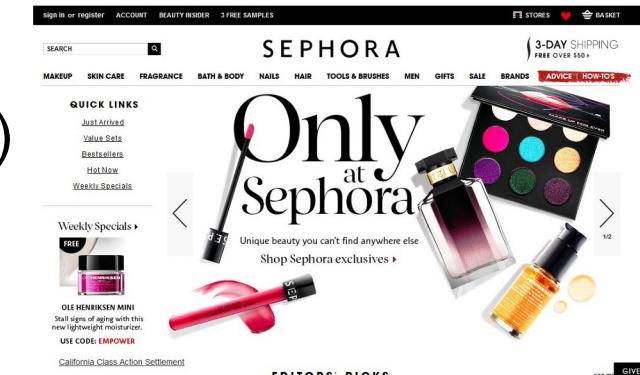
INTUITION: WHY A/B TEST?

EXAMPLE: IMPROVING SEPHORA'S WEBPAGE (NO A/B TEST)

- Context: Sephora.com launches new landing site by replacing the old landing site for all users and scrapping the old one.
- Metric: “Click-through-rate” of web page.
- Results
 - Old Webpage 10% click-through-rate (before new site)
 - New Webpage 12% click-through-rate (after new site)
- Question: Should we stay with the new webpage? Why or why not?



Old Webpage



New Webpage



LET'S SAY WE TRIED THE NEW WEBPAGE, THEN BACK TO OLD WEBPAGE, THEN THIS HAPPENED...



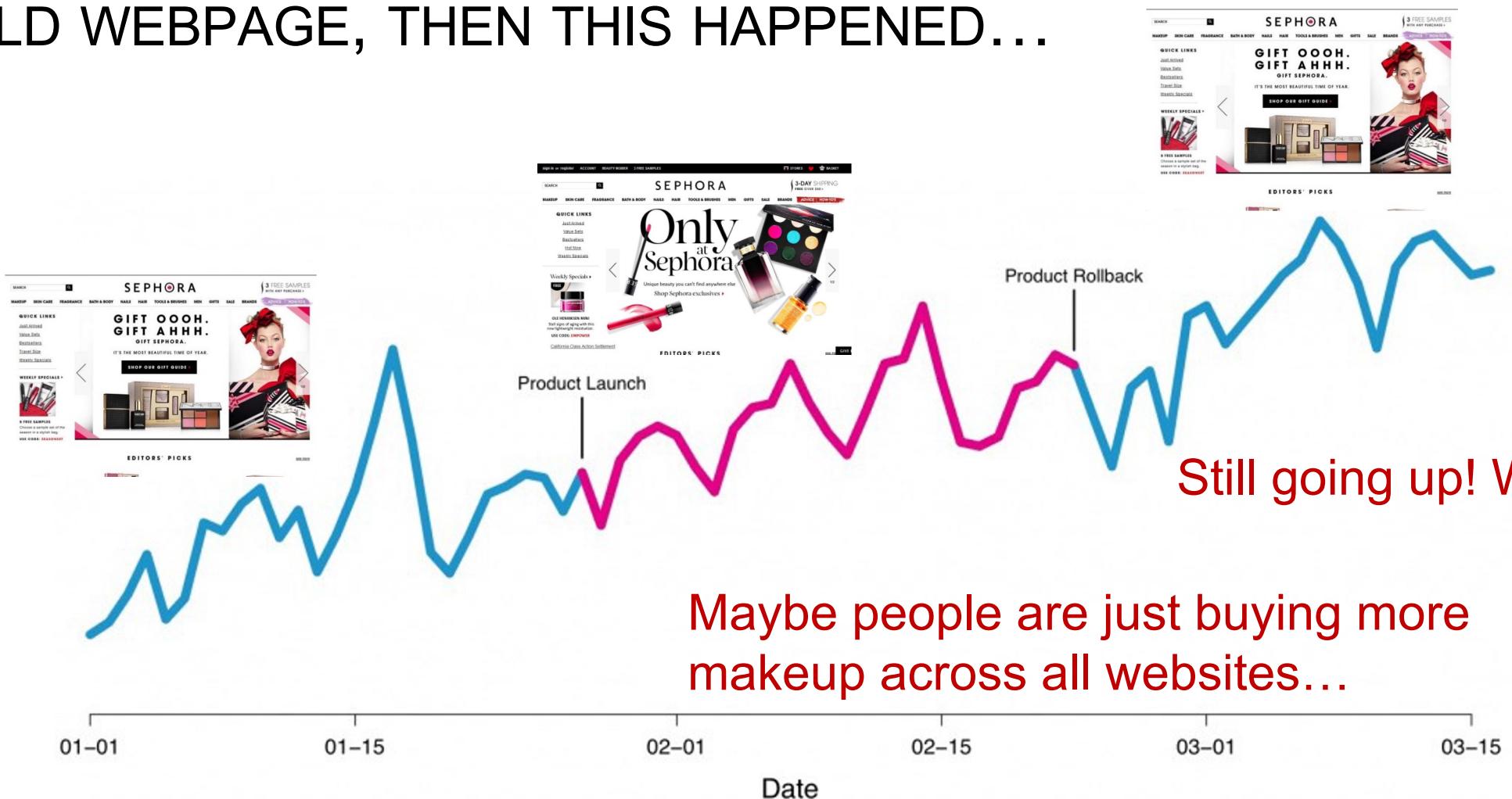
Carnegie Mellon University

Tepper School of Business

Source: <https://medium.com/airbnb-engineering/experiments-at-airbnb-e2db3abf39e7>

JOIN THE INTELLIGENT FUTURE

LET'S SAY WE TRIED THE NEW WEBPAGE, THEN BACK TO OLD WEBPAGE, THEN THIS HAPPENED...



Carnegie Mellon University

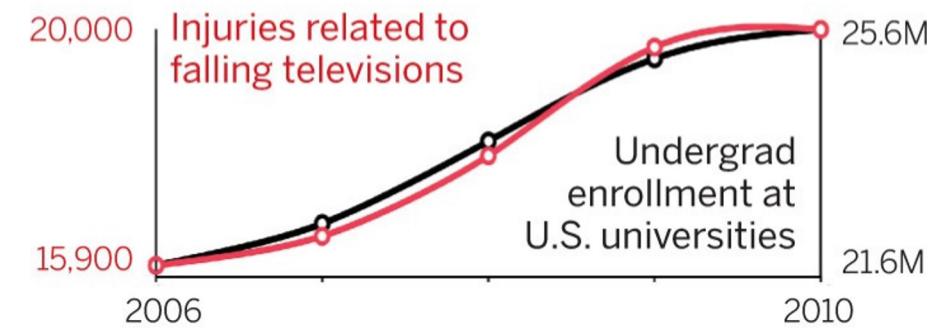
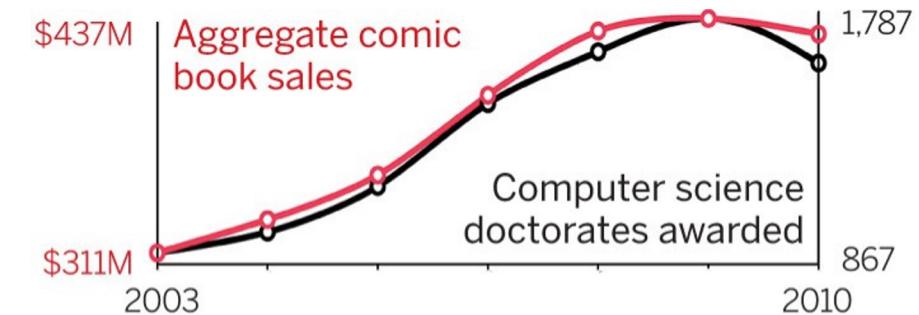
Tepper School of Business

JOIN THE INTELLIGENT FUTURE

CORRELATION IS NOT CAUSATION



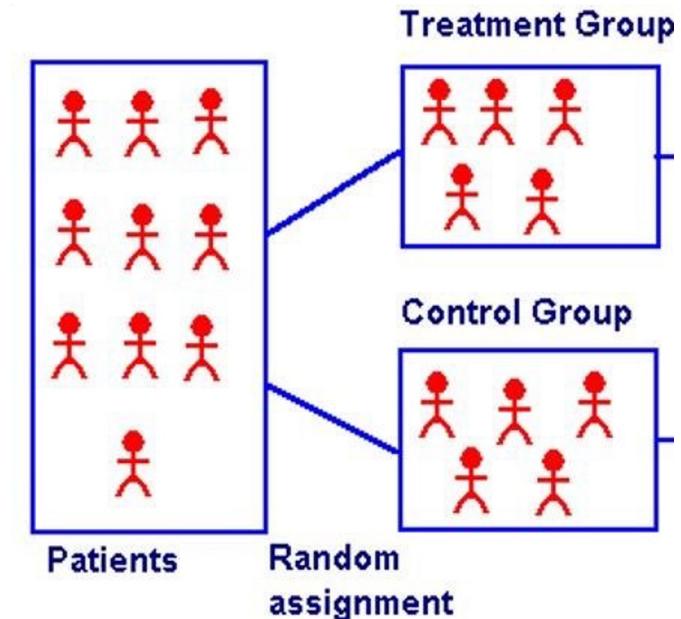
- We can not say the new webpage is better than the old webpage
- The click through rate could simply be higher for any webpage during the time we tested
- We need to control for “unrelated correlation”
- How? Using a “randomized controlled trial,” also known as an “A/B Test”



KEY IDEA: TREATMENT AND CONTROL GROUPS



- A/B Test is a “**Randomized controlled trial**”
 - Control Group (Group A) - This group sees no change from the current setup.
 - Treatment Group (Group B) - This group is exposed to the new web page
 - Goal of A/B Test
 - Compare the click-through-rates of the two groups using statistical inference.
 - Test our Product Feature “Hypothesis”
 - “[Specific repeatable action] will create [expected, measurable result]”



KEY IDEA: WE “RANDOMIZE” USERS INTO CONTROL (A) AND TREATMENT (B) TO MEASURE CAUSATION (OF OUR PRODUCT FEATURE)

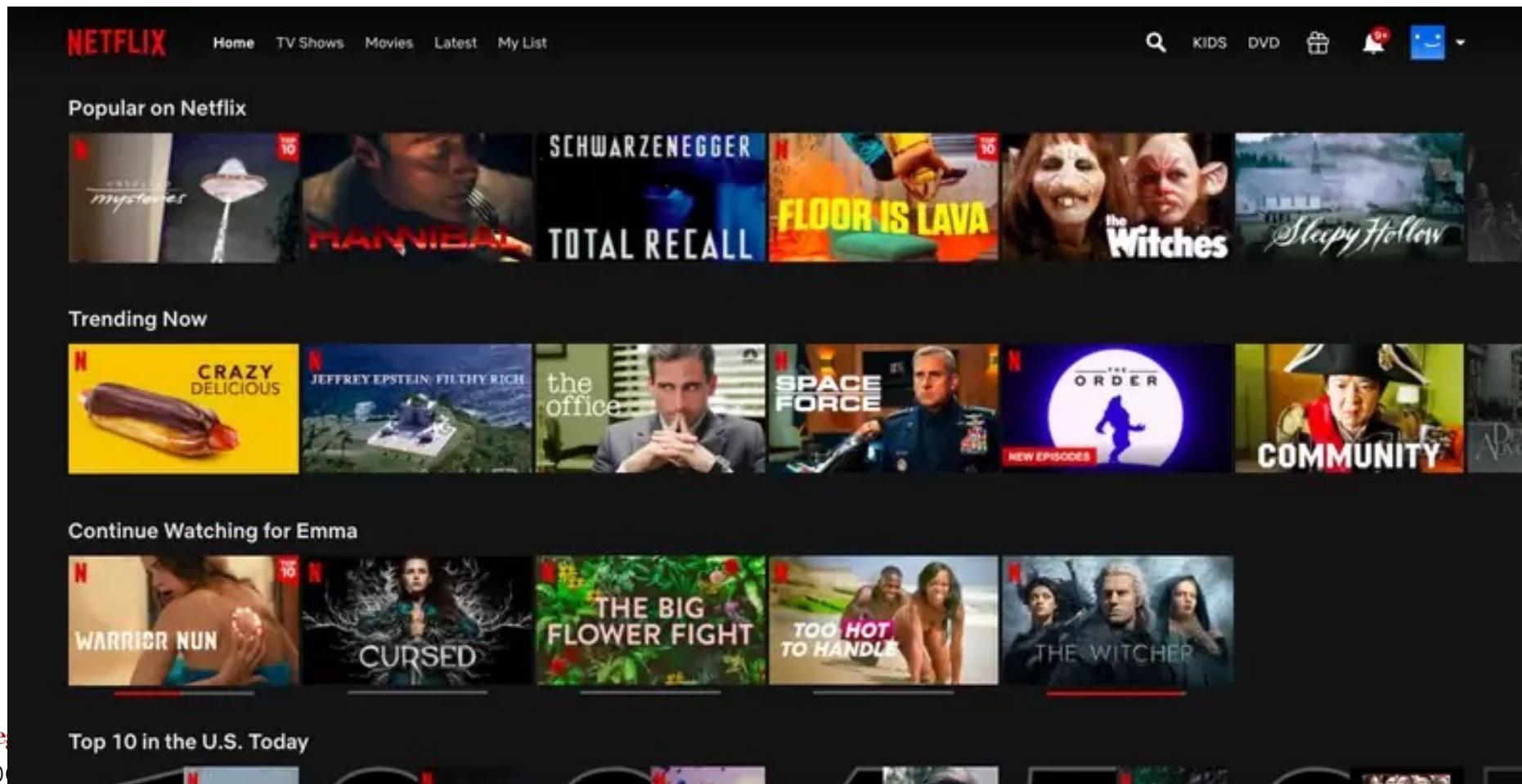


- Why? The world is not a vacuum...
 - More is happening than just the experiment (treatment vs control group) and effect.
- Many “random” situations lead to “biases”
 - Users overall are changing (e.g., buying more makeup).
 - Users characteristics are changing: Different demographics, new vs returning
 - Users have changing goals and intentions: Browsing around, buying immediately, bored and scrolling.
 - Users are finding our webpage differently than before: email, newsletters, web searches, social media
- Key Point: Randomization of users into treatment (A) and control (B) helps balance these out.



WHAT TO CONSIDER IN A/B TESTING?

Let's suppose you are running a A/B test on Netflix main page layouts.



Carne
Teppu

Let's suppose you are running a A/B test on Netflix main page layouts.



WHY IS IT BETTER THAN HISTORICAL DATA ANALYSIS?

> WE JUST DISCUSSED IT.

QUESTION 1: WHAT IS THE CHALLENGE IN TESTING AT SCALE (TESTING MANY DIMENSIONS AT THE SAME TIME)?

QUESTION 2: WHAT ARE THE PROS/CONS OF RUNNING TESTS FOR ALL CUSTOMERS VS. A SUBSET OF CUSTOMERS (WHO WOULD YOU TEST)?

QUESTION 3: WHAT METRICS YOU WOULD MEASURE AS AN OUTCOME?

CAITLIN SMALLWOOD (VP OF SCIENCE @ NETFLIX) IN 2021

Carnegie Mellon University

Tepper School of Business

JOIN THE INTELLIGENT FUTURE



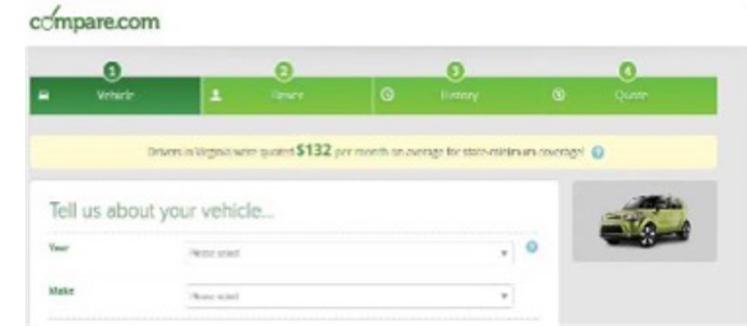
Carnegie Mellon University
Tepper School of Business

From 10:46

JOIN THE INTELLIGENT FUTURE

3 STEPS FOR A/B TESTING

- Step 1) Hypothesis Definition
 - How do we define success or failure (of our product feature)?
 - What metrics do we need to track and measure to test our hypothesis?
- **Step 2) Experimental Design and Launch A/B Test**
 - Who do we ask / who should we get data from?
 - How do we ensure we are not getting “biased” results?
 - How many samples do I need? How long to run A/B test?
- Step 3) Hypothesis Testing and Interpretation
 - Note: A/B Test = Hypothesis Test
 - What is statistically significant for our test of our hypothesis?
 - Do we go with product feature A or B?

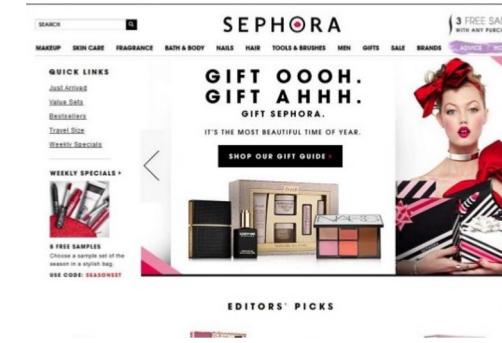


“Drivers in Connecticut
were quoted \$_____”

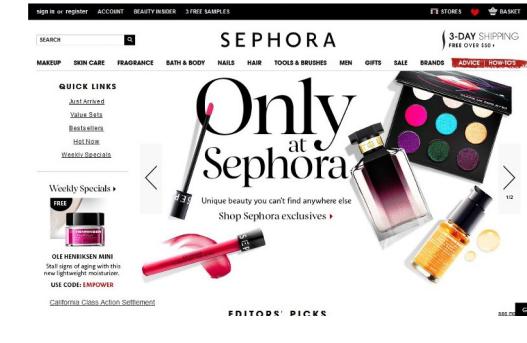


BACK TO SEPHORA: WEBPAGE A VS WEBPAGE B (NOW WITH A/B TEST)

- Now Sephora ran Webpages A and B
 - Users randomly assigned A or B
- Outcome
 - Webpage A: 10% click-through-rate
 - Webpage B: 12% click-through-rate



Old Webpage A



New Webpage B

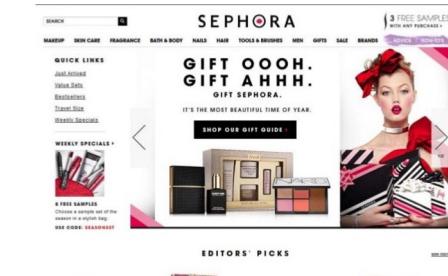
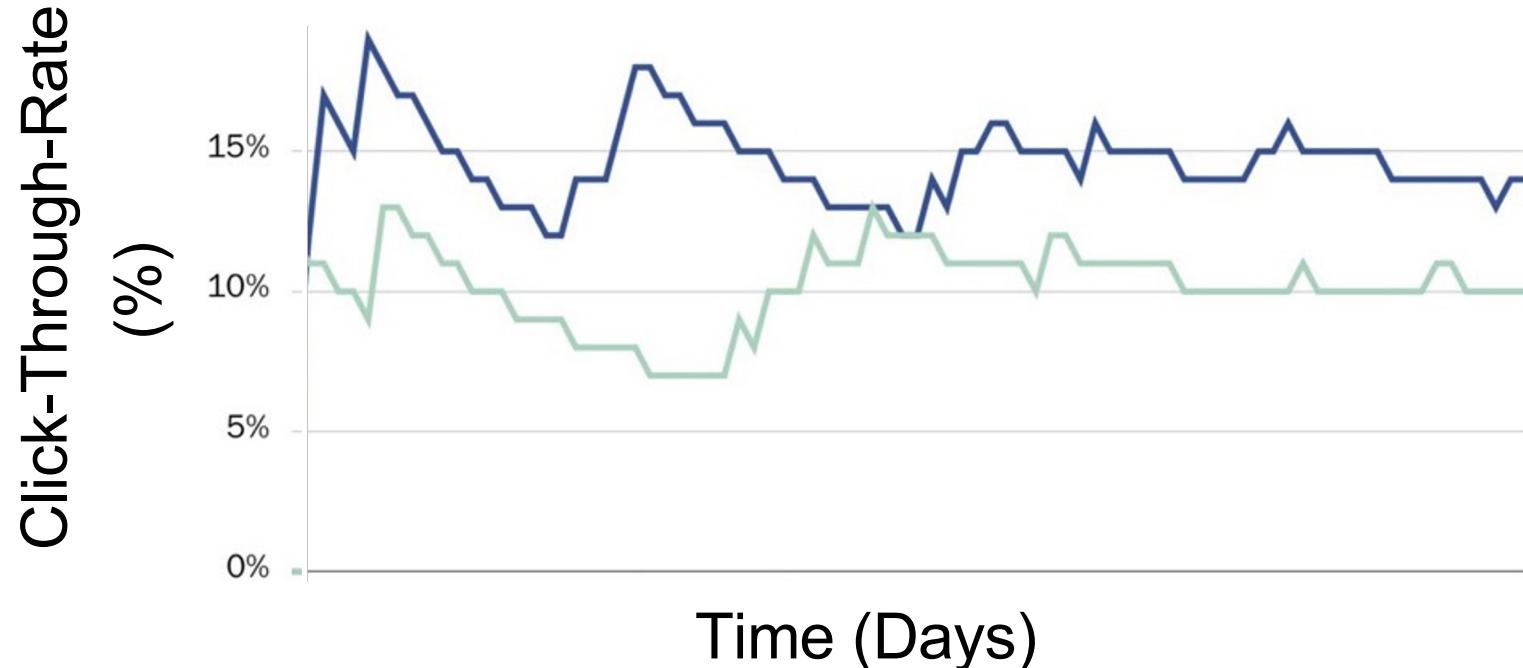
- Question: Do we implement Webpage B?
Why or why not?
- Answer: Not yet! We don't know if this is statistically significant.

Webpage	Number of Samples	Click-Through-Rate
A	1000	10%
B	1000	12%

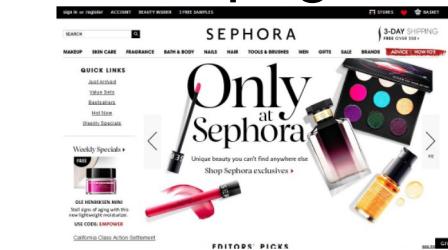
WHY HYPOTHESIS TEST?

ANSWER: BECAUSE A AND B HAVE RANDOMNESS.

WE NEED TO DETECT SIGNAL OVER NOISE



Webpage A



Webpage B

Notice: a lot of “bouncing around” of metric for A and B.

KEY POINT: METRICS NEARLY ALWAYS HAVE UNEXPLAINABLE RANDOMNESS.

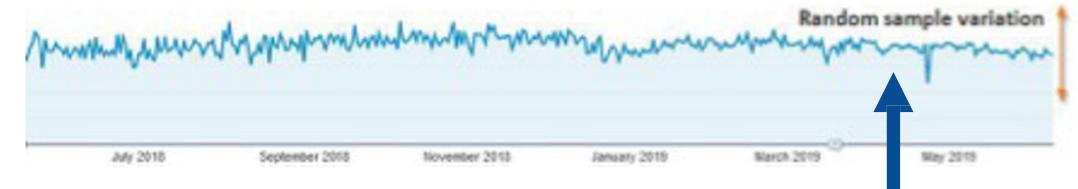


WE NEED (STATISTICAL) HYPOTHESIS TESTING TO CHECK SIGNIFICANCE.

- Why randomness?
 - We do not observe all reasons why people are bouncing at our webpage
 - Example: Economy bouncing back, Valentines Day sales for make up
- Common mistake
 - This is not “measurement error,” our measurements are fine
- Statistical Randomness of Sample vs Population
 - The data we collected are only samples of overall population.
 - We don’t know if we got samples that randomly had more click-through-rate in general.
 - We need statistics to help account for this.

Webpage	Number of Samples	Click-Through-Rate
A	1000	10%
B	1000	12%

Click-Through-Rate (CTR) over Time



12% +/- 3% Click Through Rate

HYPOTHESIS TESTING CONCEPTS



■ Null hypothesis

- *The hypothesis, often referred to as H_0 , that A and B are not different and observed differences during experiment are due to random fluctuations.*

■ P-value

■ Confidence level.

- *The probability of failing to reject (i.e., retaining) the null hypothesis when it is true.*
- *Confidence level.* Commonly set to 95%, this implies that 5% of the time we will incorrectly conclude that there is a difference when there is none (Type I error). All else being equal, increasing level reduces our statistical power.

■ Alpha or Significance Level

- $1 - \text{Confidence Level}$
- There is a 95% chance of the new feature beating the original feature

■ Statistical Power

- *The probability of correctly rejecting the null hypothesis, H_0 , when it is false. Power measures our ability to detect a difference when it indeed exists.*

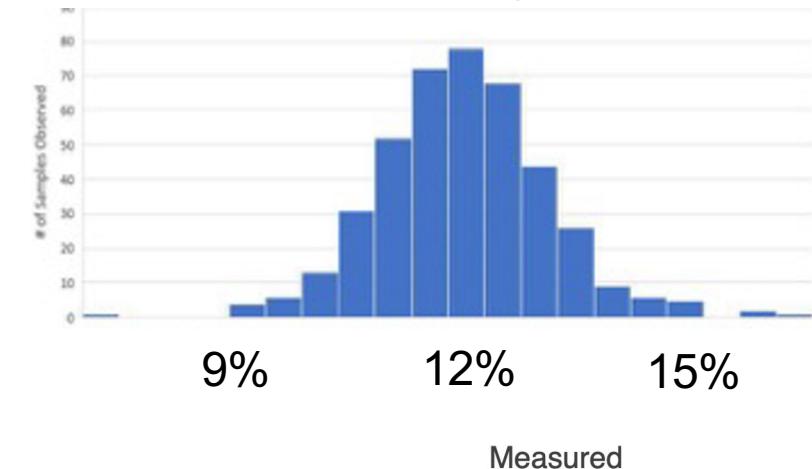
■ Standard error

- SD of the sampling distribution; The smaller the Std-Err, the more powerful the test.

Carnegie Mellon University

Tepper School of Business

Webpage A
Click-Through-Rate



		Reality	
		H_0 is true	H_0 is false
Measured	Do not reject H_0	Correct decision 😊	Type I False Positive (α)
	Reject H_0	Type II False Negative (β)	Correct decision 😊

JOIN THE INTELLIGENT FUTURE

NULL HYPOTHESIS: EXAMPLE OF RESTAURANT TECH PRODUCT FEATURE



Hypothesis definition for product feature

“[Specific repeatable action] will create [expected, measurable result]”

Example: “Restaurants that add the food pickup feature will increase their overall number of orders per day”

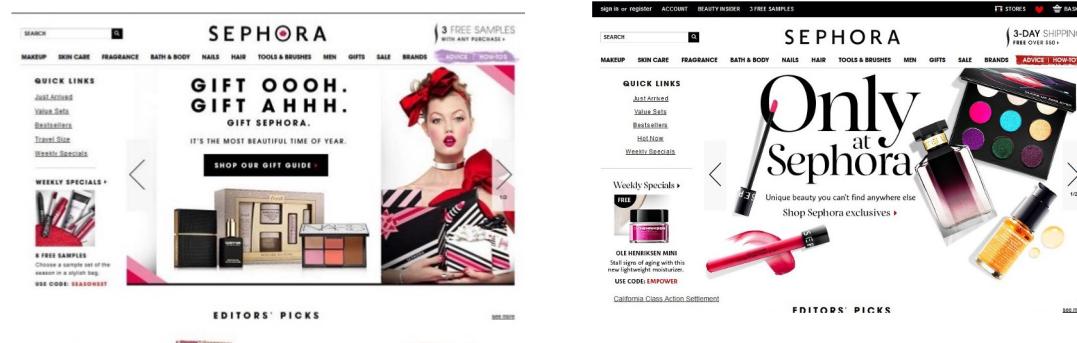
- H_0 (null): Restaurants that add the food pickup feature will neither increase or decrease total number of orders per day
- H_1 (ours): Restaurants that add the food pickup feature will increase their overall number of orders per day



QUESTION: WHAT IS NULL HYPOTHESIS AND HYPOTHESIS TO TEST?

Context: Testing two version of website – A (old) and B (new)

Metric: Click-through-rate (CT) for header banner on website

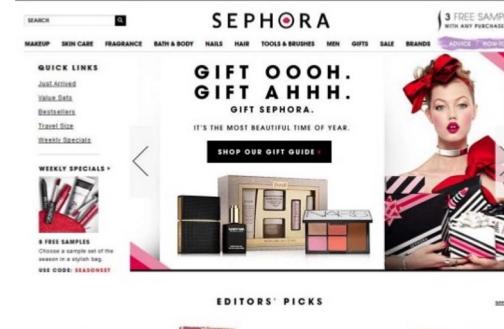


A

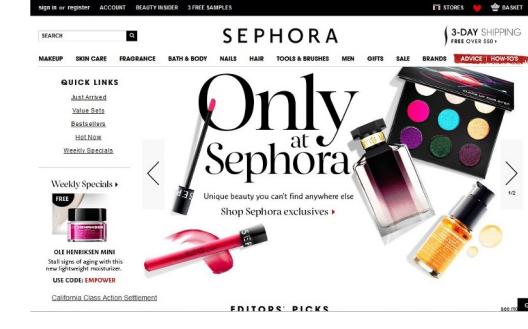
B

- H_0 : <null hypothesis>
- H_1 : <hypothesis to test>

ANSWER: WHAT IS NULL HYPOTHESIS AND OUR HYPOTHESIS?



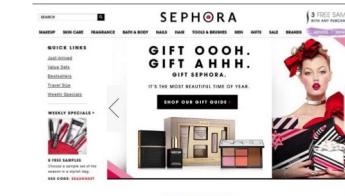
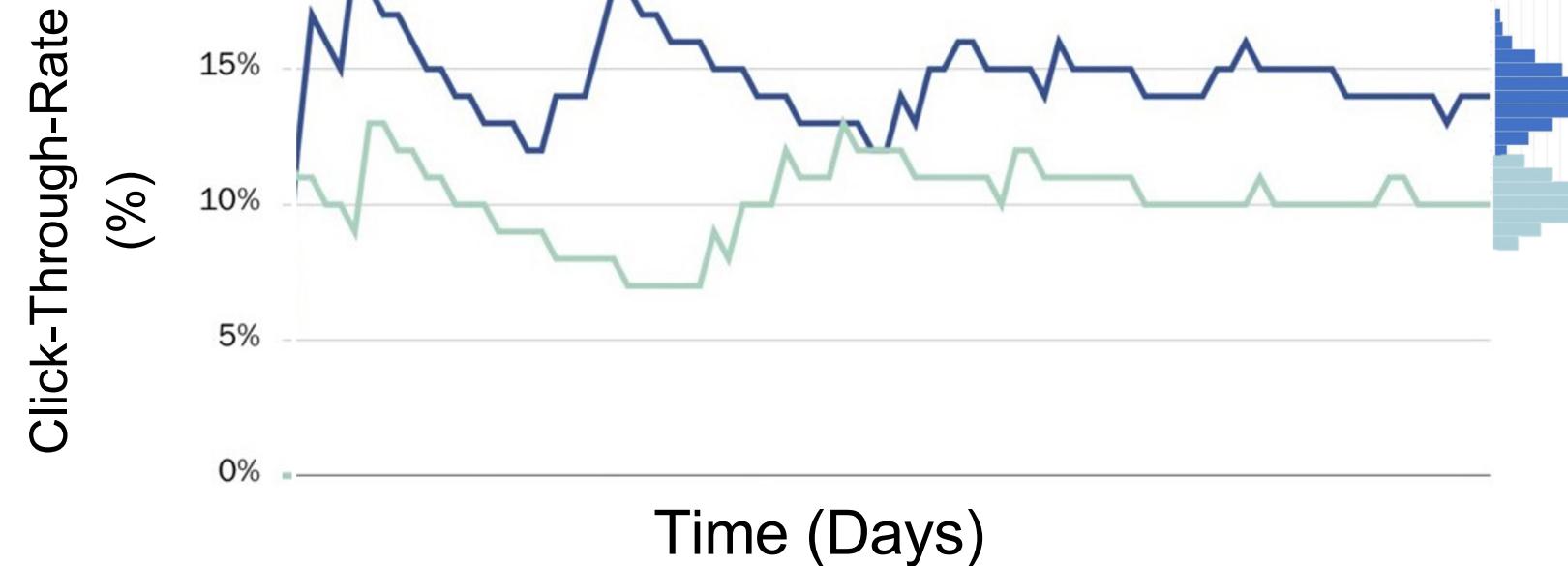
A



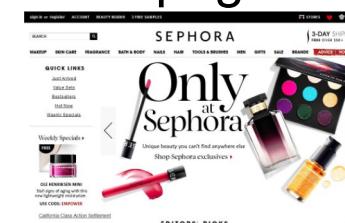
B

- H_0 : the click through rate is the same for A and B.
- H_1 : the click-through-rate is higher (or lower) for webpage A than webpage B. There is an “effect.”

INTUITION: “IS B BETTER THAN A?” AND THEIR DISTRIBUTIONS



Webpage A



Webpage B

Notice: There is a lot of “bouncing around” of metric for A and B.

Key Point: This randomness leads to a distribution of our tracked metric for A and B (blue and teal distributions).

Carnegie Mellon University

Tepper School of Business

JOIN THE INTELLIGENT FUTURE

INTUITION: THE “NORMAL” DISTRIBUTION OF A AND B



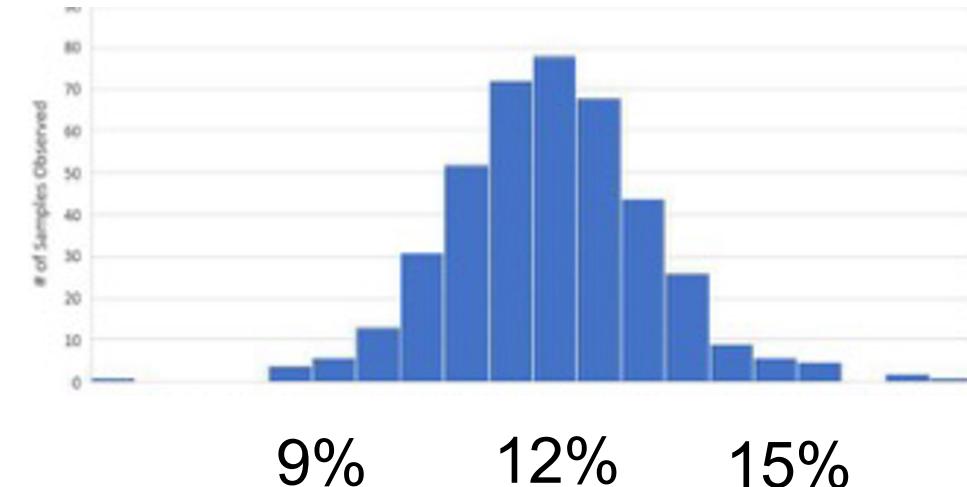
- “Normal” is the “statistical model” we will be using for click-through-rates (and A/B testing)

- A.k.a., the “Gaussian” Distribution

- Technical Note:

- We actually need “Binomial Distribution” - the number of Heads in a sequence of Bernoulli Trials with replacement
 - But this is approximated with a Gaussian

Webpage A
Click-Through-Rate



INTUITION FOR HYPOTHESIS TEST: OVERLAPPING DISTRIBUTIONS A AND B



Overlap of Two Distributions

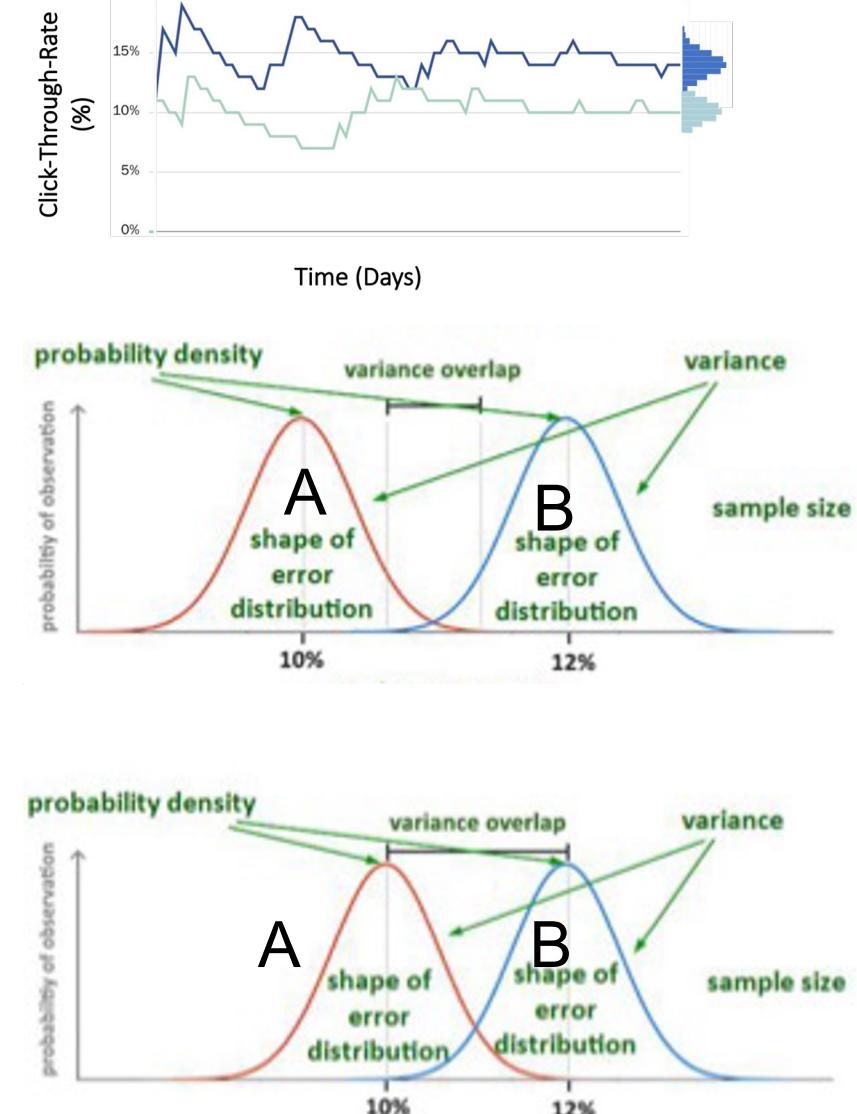
- Webpage A and Webpage B

Intuition: Less overlap of A and B

- More likely an “effect” of Webpage B versus A
- “Webpage B increases click-through-rate (CTR)”

Intuition: More overlap of A and B

- Less likely there is an “effect” of B vs A
- “Webpage B doesn’t change CTR”





HOW TO HYPOTHESIS TEST: CONVERTING DATA FROM A AND B

■ Data on A and B

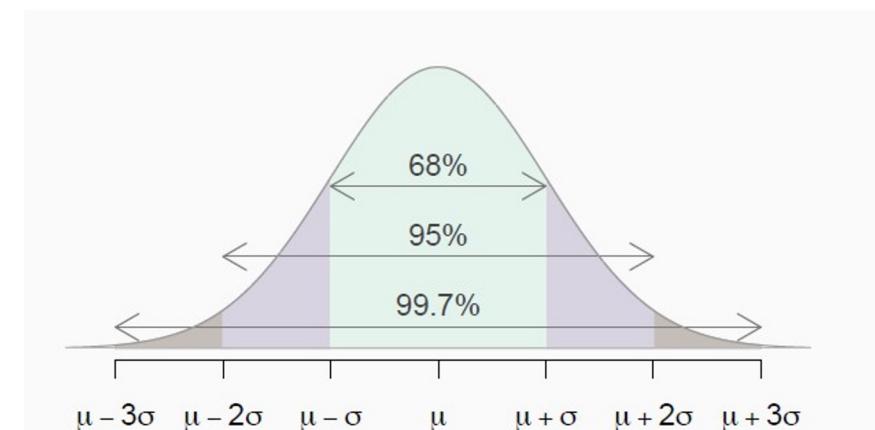
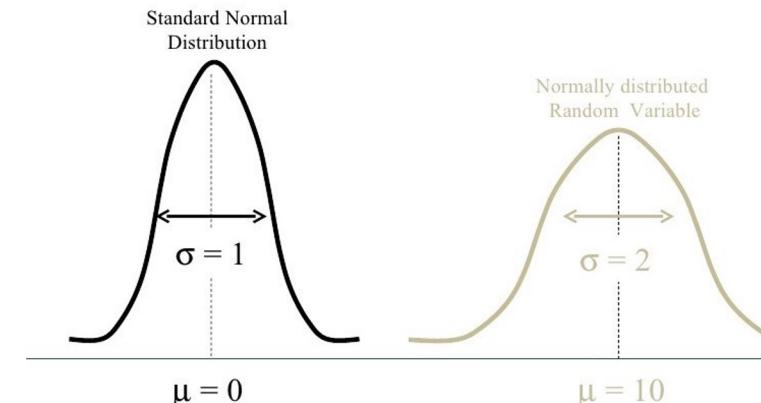
- Users that saw either Webpage A or Webpage B
- “Distribution” of user data for each webpage

■ Take the difference between webpage A and B distributions.

- This creates a new distribution which we can convert to a “standard” Normal distribution

■ Statistical Test

- This allows us to use a “Z-Test” to test our hypotheses

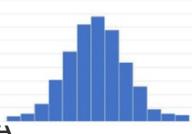


HOW TO HYPOTHESES TEST: MEAN AND STANDARD DEVIATION OF COLLECTED DATA



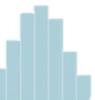
■ Webpage A:

- Mean μ_A
- Standard Deviation: σ_A
- Data sample i of CTR: a_i
- Sample Size: n_A



■ Webpage B:

- Mean μ_B
- Standard Deviation: σ_B
- Data sample i of CTR: b_i
- Sample Size: n_B



■ Z-Statistic

- Difference of Means

Carnegie Mellon University
Tepper School of Business

Mean and Standard Deviation of
Webpage A and Webpage B

$$\hat{\mu}_A = \frac{1}{n_A} \sum_i a_i \text{ and } \hat{\sigma}_A = \sqrt{\frac{1}{n_A} \sum_i (\hat{\mu}_A - a_i)^2}$$

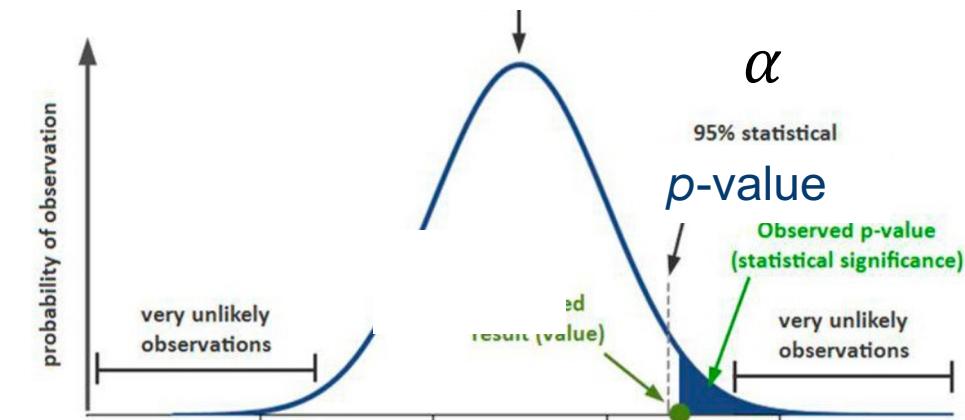
$$\hat{\mu}_B = \frac{1}{n_B} \sum_i b_i \text{ and } \hat{\sigma}_B = \sqrt{\frac{1}{n_B} \sum_i (\hat{\mu}_B - b_i)^2}$$

“Z-Transformation of our Data”
(Difference between Webpage A and B)

$$Z_{AB} = \frac{\hat{\mu}_B - \hat{\mu}_A}{\sqrt{\frac{\hat{\sigma}_B^2}{n_B} + \frac{\hat{\sigma}_A^2}{n_A}}}$$

STATISTICAL SIGNIFICANCE

- How?
 - Choose significance level α for rejecting the null hypothesis H_0 that webpage A and B have same CTR
- Reject H_0 if our p -value is $\leq \alpha$
 - This means we find H_1 true
 - Or more accurately, H_1 is more consistent with the data
- H_0 : “the click through rate is the same for the two webpages”
- H_1 : “the click-through-rate is higher (or lower) for webpage B than A”



P-VALUE



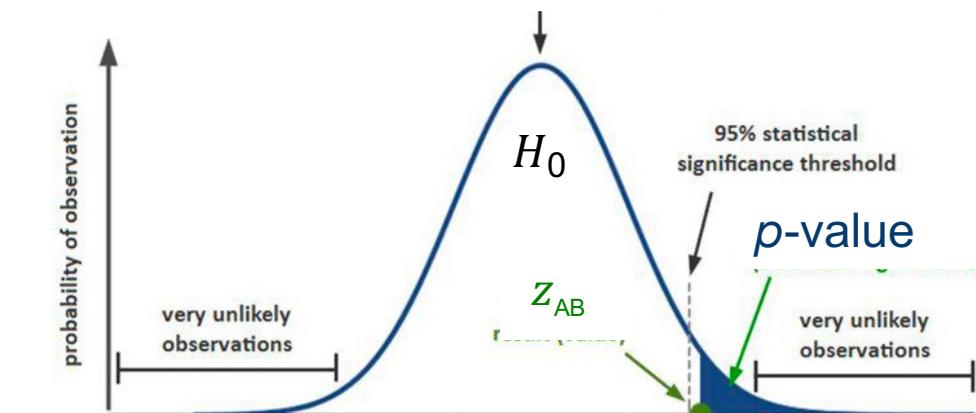
- p-value – Measure of the strength of evidence the sample data provides against the null hypothesis H_0

- We get p-value from z-statistic z_{AB}

$$p\text{-value} = P(Z \geq z_{AB}; H_0)$$

- p -value is the probability of the null hypothesis H_0 giving a more extreme difference between webpage A and B than the observed difference z_{AB}

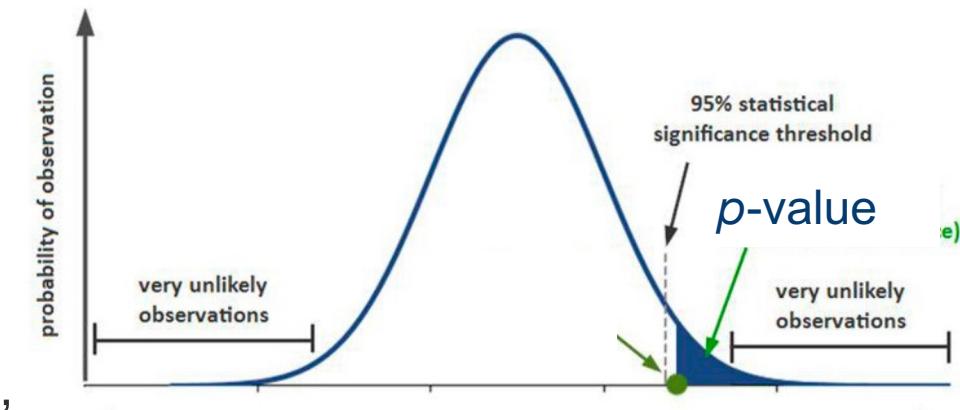
- It's the amount of "tail" in the difference of means distribution, the standard Normal specified for H_0



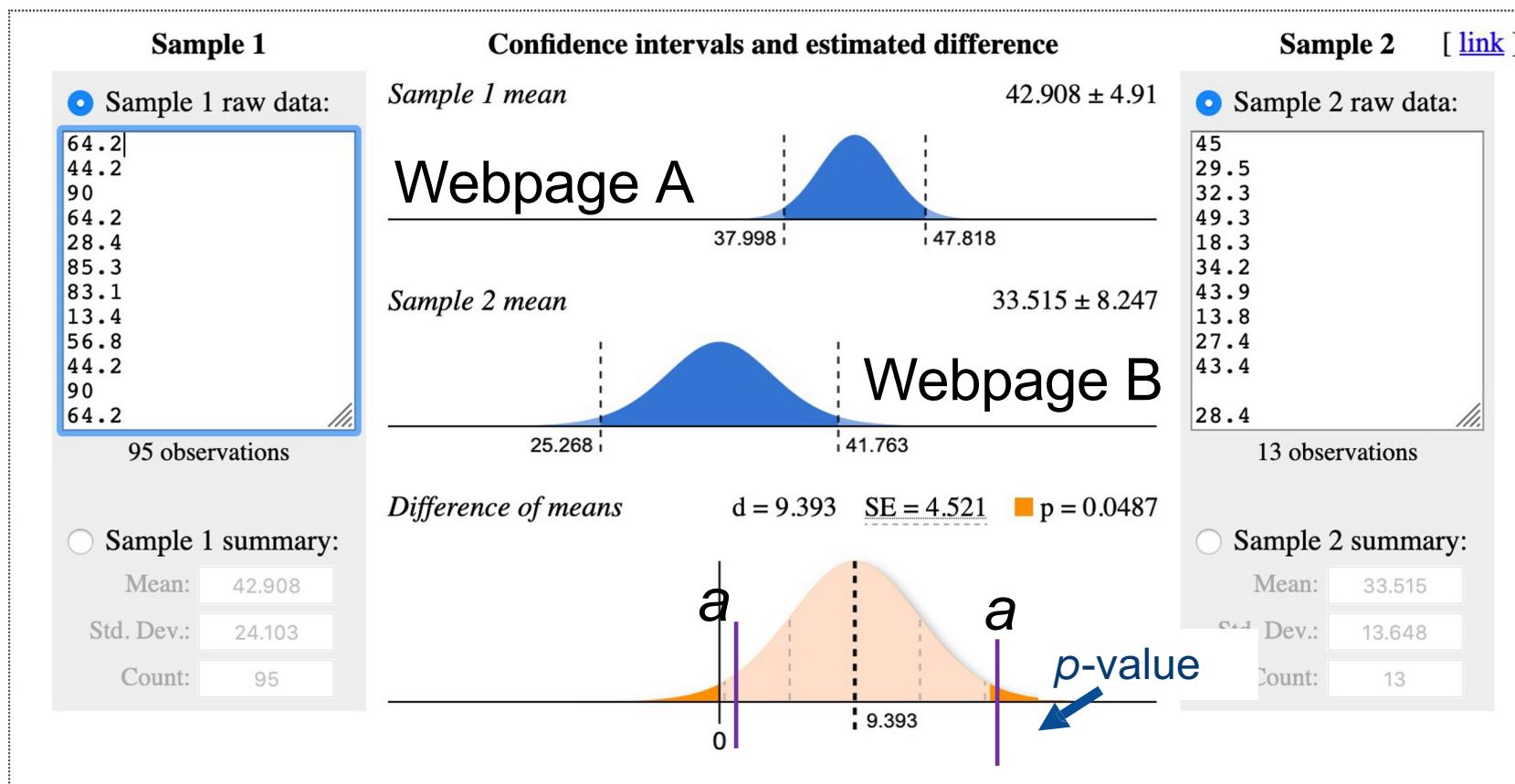
STATISTICAL SIGNIFICANCE: COMPARING P -VALUE WITH ALPHA

- We saw comparing p-value to alpha performs hypothesis test. But how do we interpret?

- p -value: We do not choose this. This is calculated from data.
- Alpha or “Significance Level”
 - We choose this. Default α is often 5%
 - You may have heard, “ $p < 0.05$ ”
 - There is a 5% or less chance of “false positive”
 - There is a 95% chance of the new webpage A having higher CTR than the old webpage B. This interpretation is called the “confidence level”



STATISTICAL SIGNIFICANCE: OVERLAPPING DISTRIBUTIONS



Verdict: Sample 1 mean is greater

Webpage A is statistically significant in its CTR from Webpage B

QUESTION: DOES P-VALUE < 0.05 ALWAYS MEAN STATISTICALLY SIGNIFICANT?

■ We choose the alpha. It just happens that many academic fields have adopted alpha of 0.05

- This was introduced in 1925
- It is an arbitrary choice.
- Physicists for example used a p-value of 0.0000003 for the Higgs Boson

- p < 0.05 is very abused
 - With “internet scale” data, we can often “find” statistical significance.

K. Pearson, "On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling," *Philosophical Magazine*, vol. 50, no. 5, p. 157–175, 1900.

R. Fischer, *Statistical Methods for Research Workers*, Edinburgh: Oliver & Boyd, 1925.

Carnegie Mellon University

Tepper School of Business

JOIN THE INTELLIGENT FUTURE

Title: Redefine Statistical Significance

Authors: Daniel J. Benjamin^{1*}, James O. Berger^{2†}, Magnus Johannesson^{3*}, Brian A. Nosek^{4,5}, E.-J. Wagenmakers⁶, Richard Berk^{7, 10}, Kenneth A. Bollen⁸, Björn Brembs⁹, Lawrence Carin¹¹, Colleen M. Cawthon¹², David Chalmers¹³, Christopher J. Chaitin¹⁴, Mette Christensen¹⁵, Thomas D. Cook^{16, 18}, Kristof De Boeck¹⁷, Zdenek Drizal¹⁸, Anna Dreber¹⁹, Kenny Easwaran¹⁹, Charles Efferson²⁰, Ernst Fehr²¹, Andy P. Field¹⁹, Malcolm Foster²², Edward I. George¹⁹, Richard Gonzalez²⁴, Steven Goodman²³, Edwin Green²⁵, Donald P. Green²⁶, Andrew Gromov²⁷, James D. Hadfield²⁸, Larry J. Hedges²⁹, Leanne Hedges²⁹, Teck-Hui Ho¹⁹, Harald Hox³⁰, John Hsu³¹, John A. Ioannidis³², James Joseph³³, Daniel J. Hruschka³⁴, Gonda Imamura³⁵, John P. A. Ioannidis³⁷, Minjieng Joos³⁶, Michael Kirchner³¹, David Labus³², John List³⁴, Roderick Little³⁴, Arthur Lupia⁴³, Edward Machery⁴⁶, Scott E. Maxwell³⁷, Michael McCarthy³⁴, Don Moore³⁹, Stephen L. Morgan³⁰, Marco Munafò^{31, 32}, Shinichi Nakagawa³³, Brendan Nyhan⁴⁰, Daniel Oberauer⁴¹, Michael O’Boyle⁴², Michael O’Reilly⁴³, Daniel Oberauer⁴⁴, Judith Rousseau³⁹, Victoria Savalei⁴⁰, Felicia D. Schiehler⁴⁵, Thomas Sellke³², Béatrice Sinclair⁴⁶, Dustin Tingley⁴⁸, Trisha Van Zandt⁴⁹, Simine Vazire⁴⁶, Duncan J. Watts⁵⁰, Christopher Winslade⁴⁸, Robert T. Wolpert¹, Yu Xie⁴⁹, Cristobal Young⁵⁰, Jonathan Zimmerman⁵¹, Valen E. Johnson^{7, 24}

Affiliations:

¹Center for Economic and Social Research and Department of Economics, University of Southern California, Los Angeles, CA 90089-3332, USA.

²Department of Statistical Science, Duke University, Durham, NC 27708-0251, USA.

³Department of Economics, Stockholm School of Economics, SE-113 83 Stockholm, Sweden.

⁴University of Virginia, Charlottesville, VA 22908, USA.

⁵Center for Open Science, Charlottesville, VA 22903, USA.

⁶University of Amsterdam, Department of Psychology, 1018 VZ Amsterdam, The Netherlands.

⁷University of Pennsylvania, School of Arts and Sciences and Department of Criminology, Philadelphia, PA 19104-6286, USA.

⁸University of North Carolina Chapel Hill, Department of Psychology and Neuroscience, Department of Sociology, Chapel Hill, NC 27599-3270, USA.

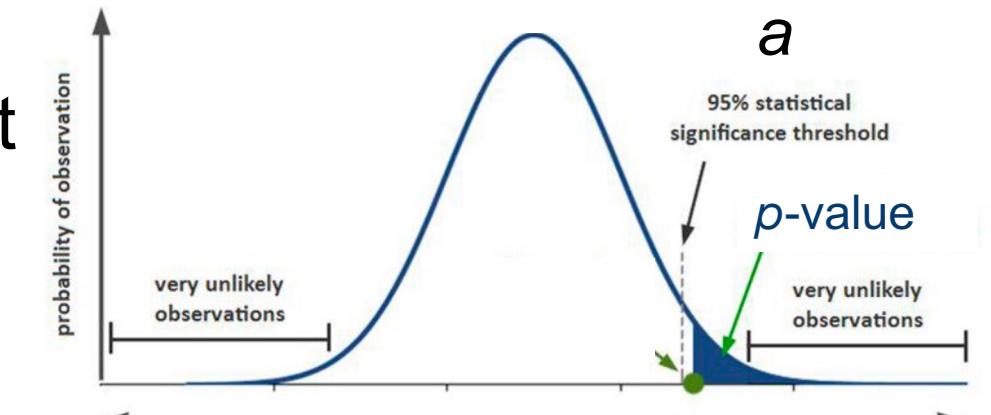
⁹Institute of Zoology - Neurogenetics, Universität Regensburg, Universitätsstrasse 31 93040 Regensburg, Germany.

QUESTION: WHAT CAN LEAD TO A LOW P-VALUE?



Observing a low p -value means either:

1. The null hypothesis is not true.
2. The null hypothesis is true, but we have observed a very rare outcome
3. The statistical model is inadequate so the calculated p -value is not an actual p -value.



HOW TO GET A LOW P-VALUE

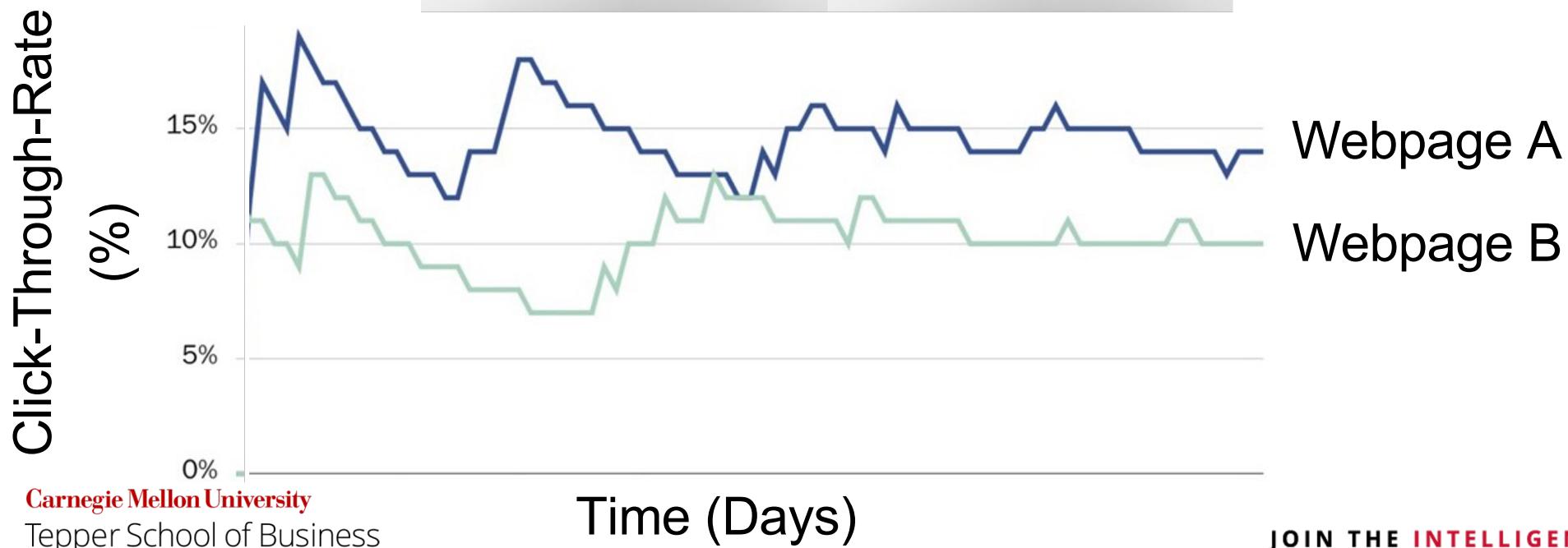


- The larger the variance of the data, the larger the p-value, all else being equal.
- The larger the sample size, the smaller p-value, all else being equal.
- The larger the observed discrepancy, the smaller the p-value, all else being equal.

Recommended interactive visualization of how sample size and variance affect statistical significance:

<https://www.evanmiller.org/ab-testing/t-test.html>

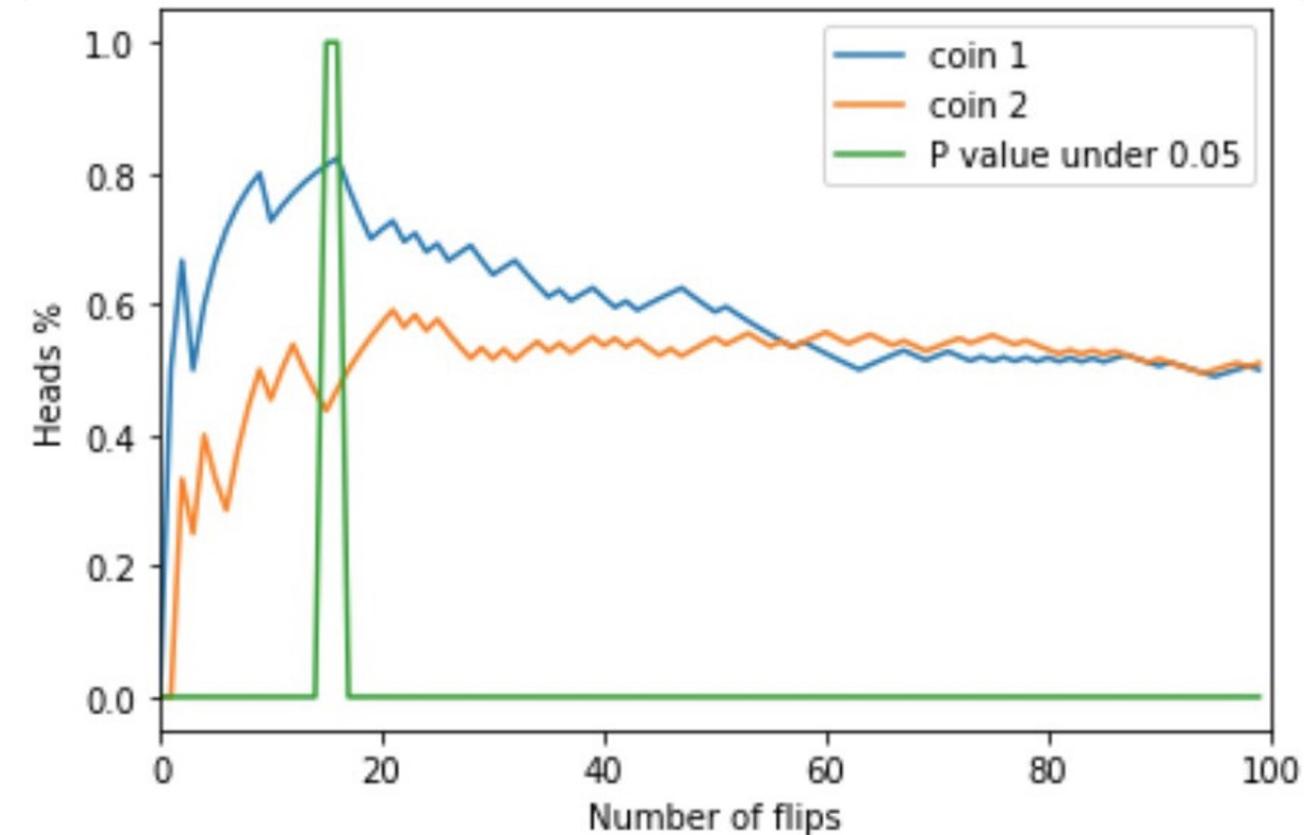
QUESTION: HOW LONG SHOULD WE RUN A/B TEST?



COMMON PITFALL: EARLY STOPPING

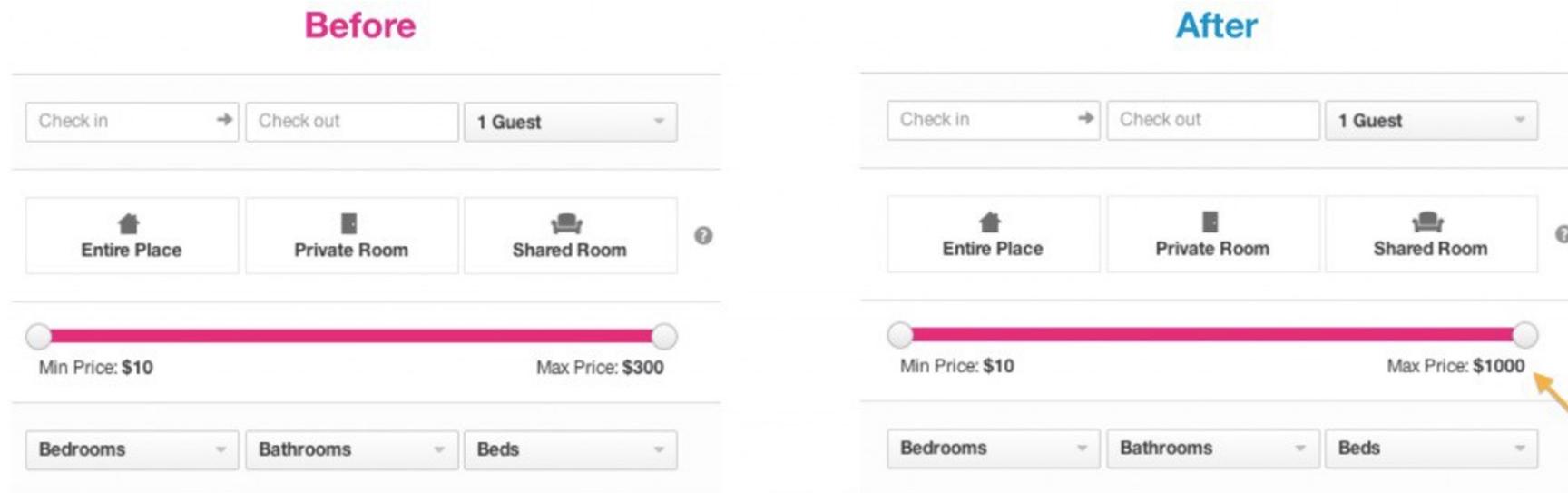


- Imagine two coins
 - Coin 1: 50% Heads
 - Coin 2: 50% Heads
- “If we stopped when we saw a “significant effect” we would have said the 2 coins were different





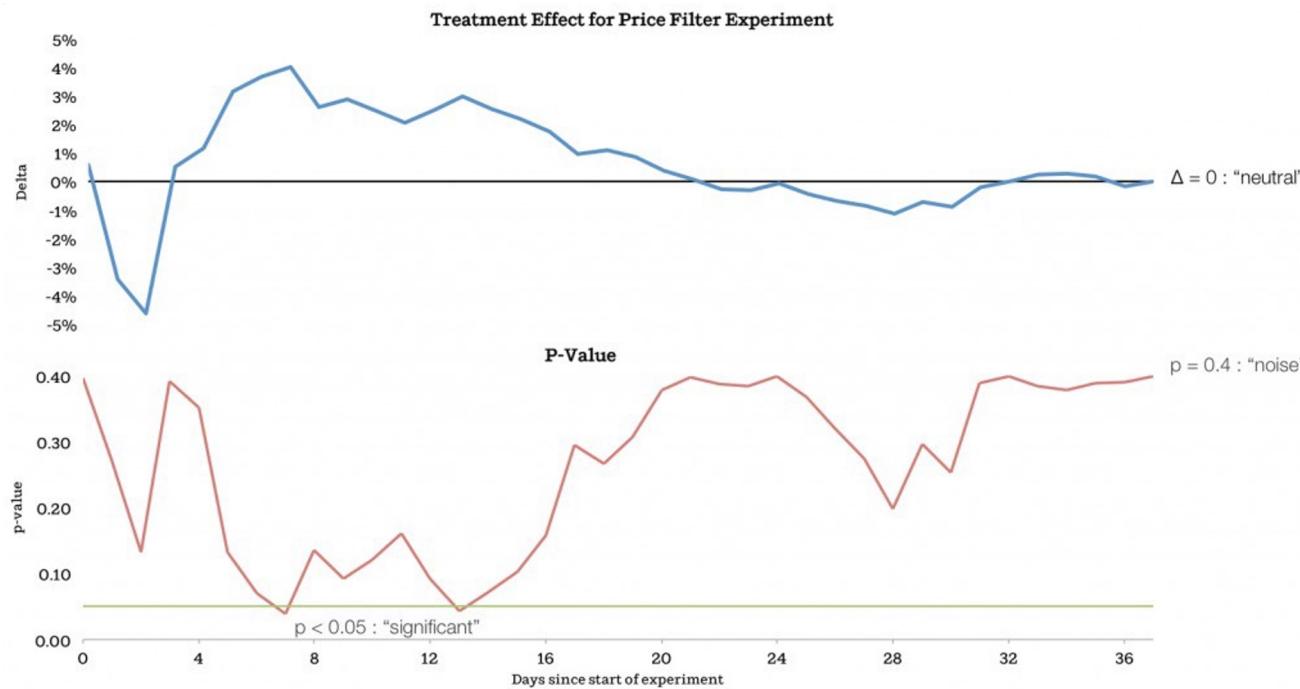
REAL EXAMPLE: AIRBNB



- Changing the Maximum Price Filter on Search Page

Source: <https://medium.com/airbnb-engineering/experiments-at-airbnb-e2db3abf39e7>

REAL EXAMPLE: AIRBNB



Takeaway: If Airbnb stopped when it was “significant” then this would have been (an incorrect) spurious correlation.

HOW MANY SAMPLES DO I NEED?



Equation are based on desired detection of effect size (from base rate) and standard deviations of A and B

Use sample size calculators:

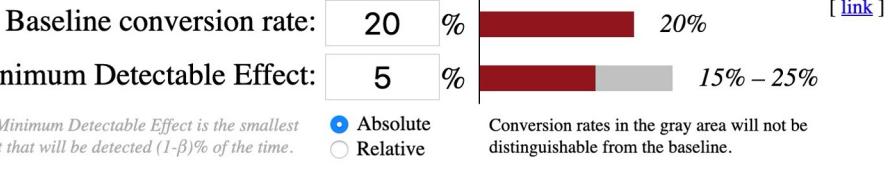
- Many online
- Link below is good one

<https://www.evanmiller.org/ab-testing/sample-size.html>

[Sample Size Calculator](#) | [Chi-Squared Test](#) | [Sequential Sampling](#) | [2 Sample T-Test](#) | [Survival Times](#) | [Count Data](#)

Need A/B sample sizes on your iPhone or iPad? Download [A/B Buddy](#) today.

Question: How many subjects are needed for an A/B test?



Sample size:

1,030

per variation

Statistical power $1-\beta$:  80% Percent of the time the minimum effect size will be detected, assuming it exists

Significance level α :  5% Percent of the time a difference will be detected, assuming one does NOT exist

See also: [How Not To Run an A/B Test](#)

IN-CLASS EXERCISE: HOW MANY SAMPLES WOULD YOU NEED AND HOW MANY DAYS NEEDED?



Go to sample size calculator

- <https://www.evanmiller.org/ab-testing/sample-size.html>

Question 1: How many samples needed for?

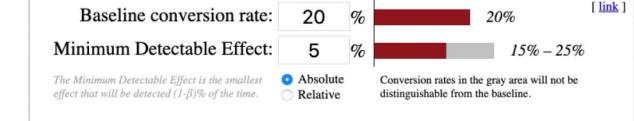
- Significance Level = 5%
- False Error Rate = 20% (Power = 80%)
- Baseline Metric = 40%
- Minimum Effect = 6% (Relative)
- Average Daily Volume (Unique Users) = 1,500
- % of unique users eligible for test = 10%

Question 2: How many days will you need for an A/B test with two variations?

[Sample Size Calculator](#) | [Chi-Squared Test](#) | [Sequential Sampling](#) | [2-Sample T-Test](#) | [Survival Times](#) | [Count Data](#)

Need A/B sample sizes on your iPhone or iPad? Download [A/B Buddy](#) today.

Question: How many subjects are needed for an A/B test?



Sample size:

1,030

per variation

Statistical power $1-\beta$: 80% Percent of the time the minimum effect size will be detected, assuming it exists

Significance level α : 5% Percent of the time a difference will be detected, assuming one does NOT exist

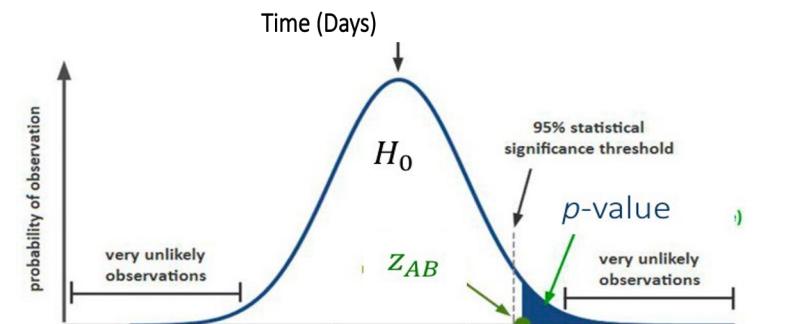
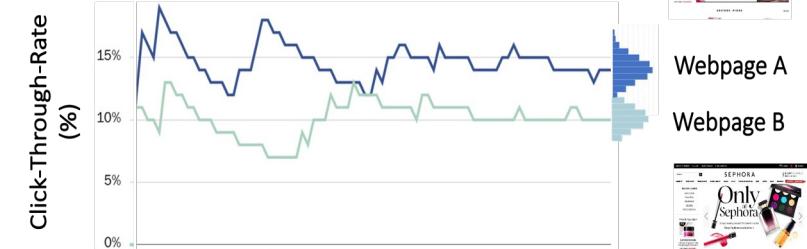
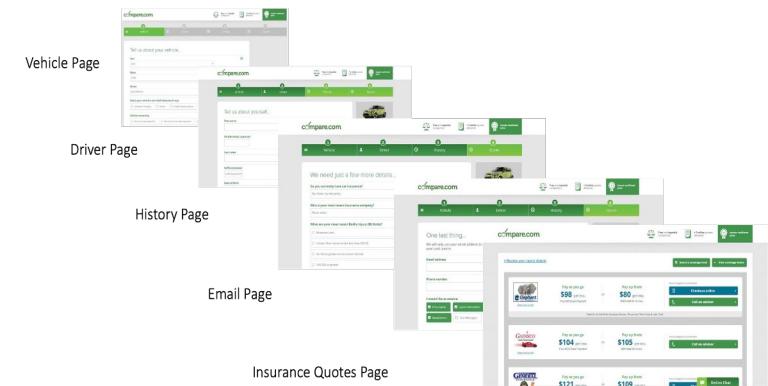
See also: [How Not To Run an A/B Test](#)

TODAY'S LEARNING OBJECTIVES



- What and Why A/B Test?
 - Focus on Intuition and Definitions
 - 3 Steps for A/B Testing

- How to A/B Test?
 - Statistical Significance
 - Concepts: p-values, confidence intervals, etc.



JOIN THE INTELLIGENT FUTURE