# Online Experimentation and A/B Testing

Data Science Dojo

datasciencedojo

data science for everyone

# Agenda

- **Introduction**
  - What is A/B testing?
  - Some interesting A/B tests
- **Fundamentals**
  - Terminology.
  - Hypothesis testing
  - Metrics for A/B testing
  - Steps in Experimentation

datasciencedojo
data science for everyone

# INTRODUCTION

In God we trust. All others bring data.

W. E. Deming

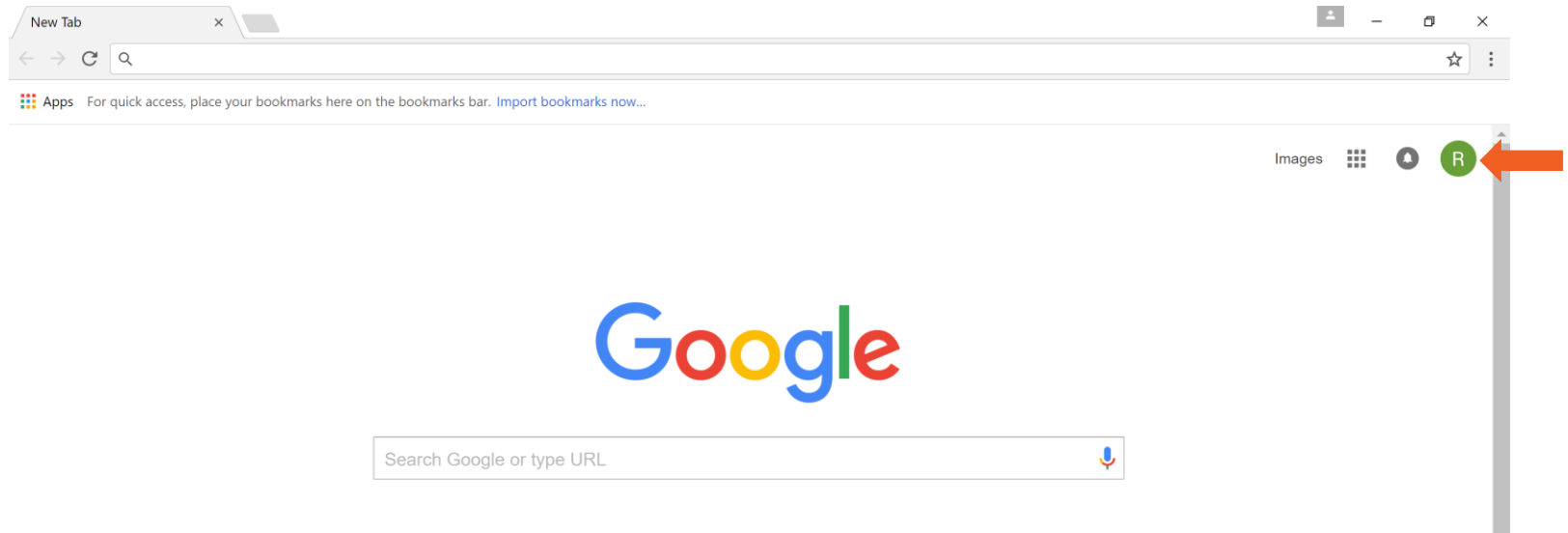# Background in Clinical Trials

## Case

- Try out a new product, treatment, or intervention
- Randomly selected study group
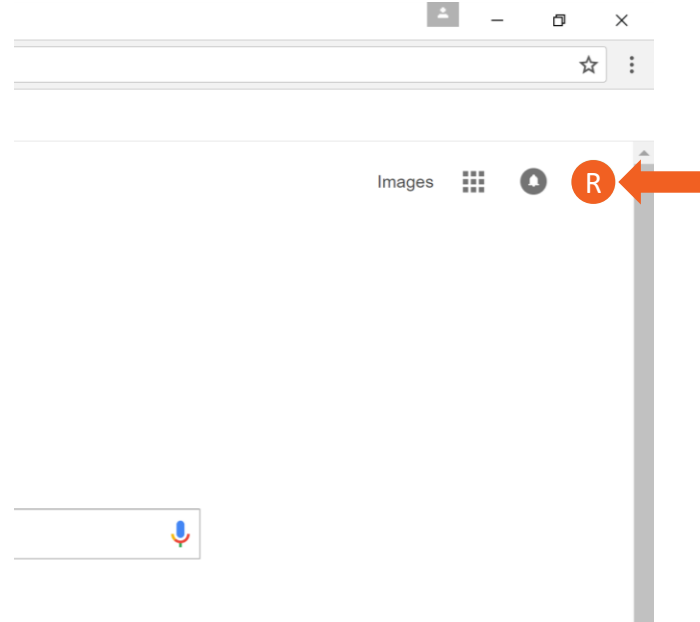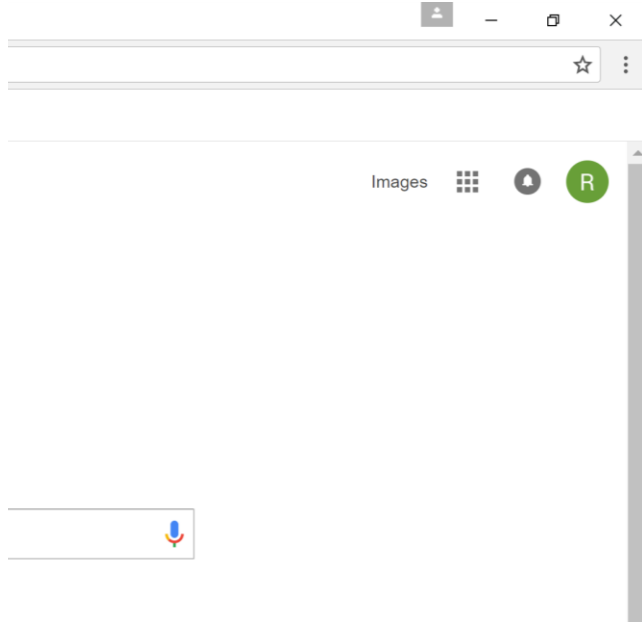- Keep all other factors constant

## Control

- The old product or the typical treatment
- Randomly selected control group
- Similarly, keep all other factors constant

datasciencedojo
data science for everyone

# What is A/B Testing?

# What is A/B Testing?

# What is A/B Testing?

- Would changing the color of the icon give me:
  - Higher CTR
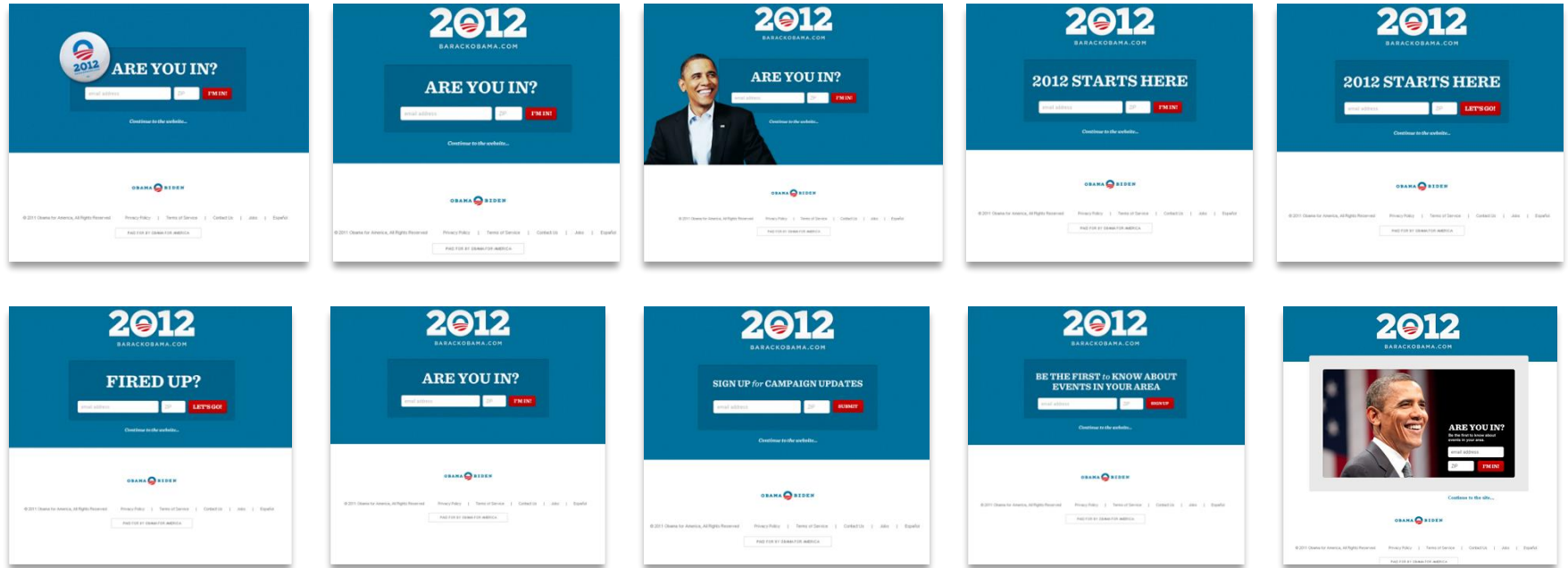  - Lower bounce rate
  - Higher revenue per user

datasciencedojo
data science for everyone

# Obama 2012 Campaign
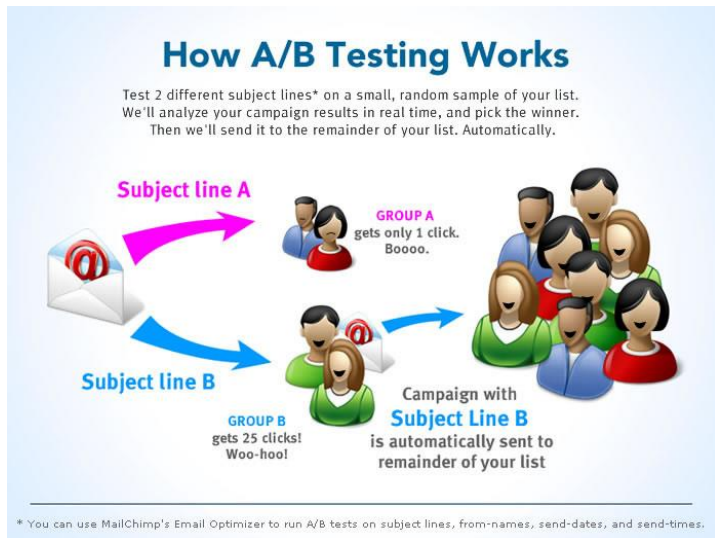
datasciencedojo
data science for everyone

# Obama 2012 Campaign

## Maximize Sign-Ups And Donations



**Source:** http://www.nathanielward.net/2011/06/see-ab-testing-in-action-on-barack-obamas-reelection-website/
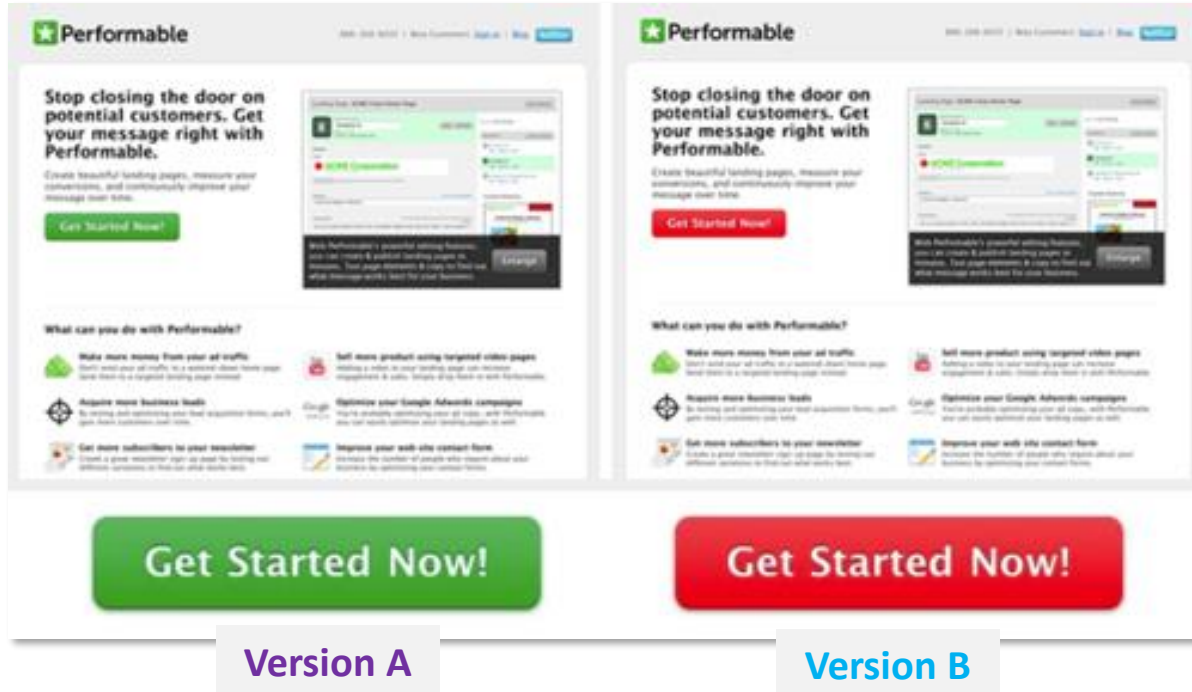
11

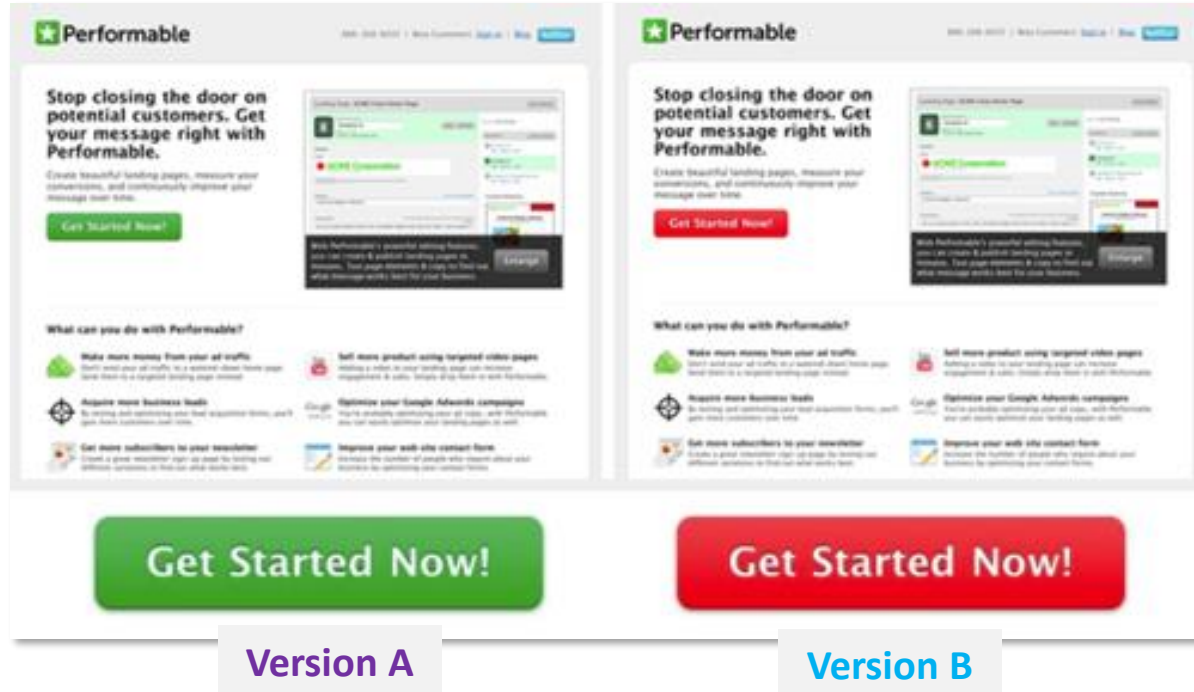# A/B Testing On Newsletters And Email



## Run tests on many things

➢ Subject lines
➢ **From** names
➢ **Send** dates
➢ **Send** time

# Testing Call-to-Action Button



Version A          Version B

Which button increased **number of clicks?**

13

# Testing Call-to-Action Button
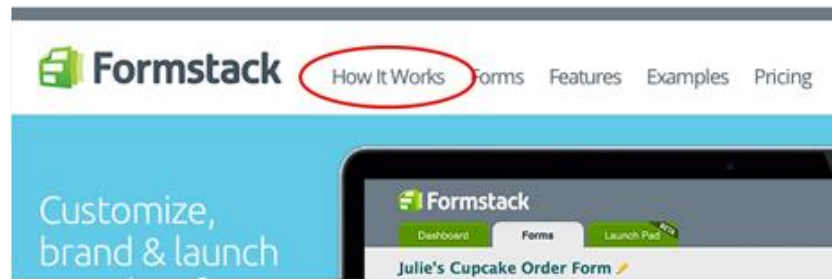


Version A                    Version B

**Red button** increased clicks by **21%**

14

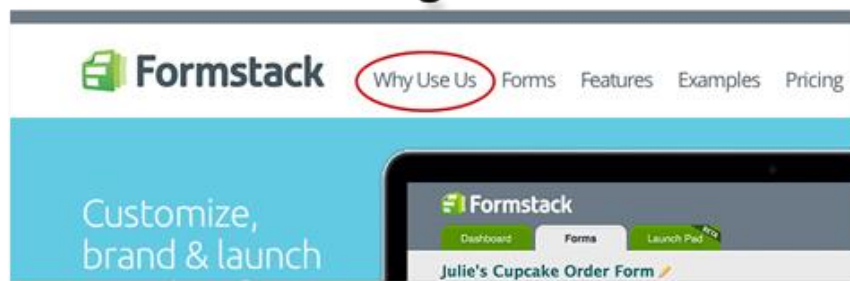# Testing Navigation Bar



Version A

Version B

'Which experience increased **clicks** by 47.7%

# Testing Navigation Bar



Version A

Version B

**'How It Works'** increased clicks by 47.7%

datasciencedojo

data science for everyone

# Jocelyn or Michael?



**Conversion Rate:** Who gives better conversion rate?

# Jocelyn or Michael?



**Michael** increased conversions by 21%

# AwayFind - Mobile notifications for priority messages



**Version A**



**Version B**

Which version increased **sign-ups** by 38%?

# AwayFind - Mobile notifications for priority messages



Version A

Version B

**Version B!**
A longer yet clearer message is more effective.

20

# Online Form



**Version A**

**Version B**

Which radically redesigned form increased **B2B leads** by **368.5%?**

21

# Online Form



Version A

Version B

Version A!

# CALIFORNIA CLOSET



**Version A**

**Version B**

Which version **increased leads** by 115%

25

# CALIFORNIA CLOSET



**Version A**

**Version B**

Version A increased leads by 115%.
**This is why you should test...!**

# FUNDAMENTALS

# Why We Use A/B Testing

## Problem

- Users are complex and our intuition is often wrong

- Rolling out a feature to all the users at the same time is risky

## A/B testing purpose

- Know what the users want subconsciously or otherwise.

- Helps to fail fast and move on

Impact is always expected to be positive, but outcome is often humbling

datascïencedojo
data science for everyone

# Multi-Variate Testing

# A/B Testing vs Multivariate Testing

|  | **A/B Testing** | **Multivariate Testing** |
|---|---|---|
| Common use | Compare two very different designs with each other | Several minor variations are up for debate:<br>➢ Two colors of button with three different headlines<br><br>ⓘ Also called full factorial testing |
| Advantages | ➢ Simple in design<br>➢ Small sample size may be ok | A lot of different combinations tried at once. |
| Limitations | Trying only one alternative | ➢ Bigger sample size<br>➢ Complex<br>➢ Need better understanding of interactions |

datasciencedojo
data science for everyone

# TERMINOLOGY

# Control and Treatment

## Control

Default experience, the way things are now.

**Example:** Current look and feel of your 'Buy Now' button

**Buy Now**

## Treatment

The change we want to make.

**Example:** Change the button from green to blue

**Buy Now**

## Illustration



**Clinical Research Participants (100 people)**

Random placement into each group (like a coin toss)

**Experimental Group (50 people)**
Receive new medication or therapy

**Control Group (50 people)**
Receive a standard treatment with known effects

datasciencedojo
data science for everyone

# METRICS FOR A/B TESTING

# Metrics Used For A/B Testing

➢ **Search engines**

   Queries/UU, Session length, Sessions/UU, Page views, Bounce rate

➢ **Online Retailers**

   Conversion rate, revenue/UU, Avg Cart Value and so on

➢ **Other websites**

   CTR, signup for newsletter

## Each business is different

datasciencedojo

data science for everyone

# Null vs Alternate Hypothesis

- Null Hypothesis ($H_o$)
  - Control and treatment are similar (in terms of the parameter we are estimating)
- Alternate Hypothesis ($H_a$)
  - Treatment is different from control

datasciencedojo
data science for everyone

# Null vs Alternate Hypothesis



Control               Treatment

- **Null Hypothesis ($H_o$)**
  - Green and blue buttons have the same CTR
- **Alternate Hypothesis ($H_a$)**
  - Each button has a different CTR

# Type I and Type II Error

Type I Error:
The probability of **falsely rejecting** null hypothesis

Type II Error:
The probability of **falsely accepting** null hypothesis

**Ground Truth**

|  | Ho is true | Ho is false |
|---|---|---|
| **Reject Ho** | Type I error | Correct decision |
| **Do not reject Ho** | Correct decision | Type II error |

**Experiment Outcome**

datasciencedojo
data science for everyone

# CAN YOU TELL ME IN SIMPLE WORDS

# The Cook and Smoke Detector

- Null Hypothesis ($H_o$): There is no fire
- Alternate Hypothesis ($H_a$): There is fire

# The Cook and Smoke Detector

- **Type I Error:** There is no fire but smoke detector goes off.
- The cook removes the alarm to prevent type I errors.
- This increases the chance of **Type II Error** (i.e., a fire without an alarm)

datasciencedojo
data science for everyone

# The Boy Who Cried Wolf

- Null Hypothesis ($H_o$): There is no wolf
- Alternate Hypothesis ($H_a$): There is a wolf

# The Boy Who Cried Wolf

- **Type I Error:** Villagers believe the boy when there is no wolf

- **Type II Error:** Villagers do not believe the boy when the wolf is really there

# Confidence Interval

**Problem:** On a 5-point scale, a product has an average rating of 4.32 and a standard deviation of 0.845 based on 62 ratings. What is the 95% confidence interval?

Standard deviation

Mean

$$95\% \text{ Confidence Interval} = \bar{X} \pm 2.0 \frac{\sigma}{\sqrt{n}}$$

Critical value

Sample size

# Confidence Interval

Mean $\bar{X} = 4.32$

Standard deviation $\sigma = 0.845$

Standard error (SE) $= \dfrac{\sigma}{\sqrt{n}} = \dfrac{0.845}{\sqrt{n}} = \dfrac{0.845}{\sqrt{62}} = 0.11$

Margin of error is 2.0 x 0.11 = 0.22

The confidence interval is:

4.32 − 0.22 = 4.10

4.32 + 0.22 = 4.54

Intuition – We are 95% confident that the average review for all customers is between 4.10 and 4.54.

datasciencedojo

data science for everyone

# Confidence Interval Intuition

Let's say we know for sure that the average product review for all customers is 4.4.

This is known as the *population mean* and is denoted by $\mu$.

Additionally, we decide to take 20 random samples of customer reviews from the population (i.e., all reviews).



$\mu = 4.4$

If we calculate the 95% confidence interval for all 20 samples we know that typically 19 of them will contain $\mu$.

datasciencedojo
data science for everyone

# Calculating Confidence Interval

Normal Distribution

| Confidence level | Z score |
|---|---|
| 90% | 1.645 |
| 95% | 1.960 |
| 98% | 2.326 |
| 99% | 2.576 |

Student's t Distribution

### Critical Values (t*)

| n − 1 | Confidence Level | | |
|---|---|---|---|
| | 0.900 | 0.950 | 0.990 |
| 10 | 1.812 | 2.228 | 3.169 |
| 20 | 1.725 | 2.086 | 2.845 |
| 30 | 1.697 | 2.042 | 2.750 |
| 40 | 1.684 | 2.021 | 2.704 |
| 50 | 1.676 | 2.009 | 2.678 |
| 60 | 1.671 | 2.000 | 2.660 |
| 70 | 1.667 | 1.994 | 2.648 |
| 80 | 1.664 | 1.990 | 2.639 |
| 90 | 1.662 | 1.987 | 2.632 |
| 100 | 1.660 | 1.984 | 2.626 |

Our example!

# Confidence Interval

- Range of plausible values of parameter being estimated (e.g., the mean), given the sample data

# A/A Test

- Comparing the identical experience on different random sets of users
- Used for validation of setup



**Control**



**Treatment**

# Steps in Experimentation

**Planning**
- Choose factors, levels, sample size(how long to run)
- What business question to answer
- Metrics and expected outcome

**Coding and Logging**
- Setup of test and instrumentation

**A/A Test**
- To make sure the setup is correct.

**Make a Decision**
- To ship or not to ship

**Analysis and interpretation**
- Some times this can be an art
- Newness effect
- Seasonality, segments etc.

**A/B and/or multivariate test**

datasciencedojo
data science for everyone

# Categories of Metrics

|  | Short-term | Medium-term | Long-term |
|---|---|---|---|
| **Examples** | ➢ CTR <br> ➢ PVs <br> ➢ Bounce Rate | ➢ PVs/user/day <br> ➢ CTR/user /day <br> ➢ Avg. session length | **Days with at least one visit:** <br> ➢ Total time on site <br> ➢ Repeat visits/user |
| **What is measured?** | Immediate or almost immediate impact | Engagement over hours up to a day | Loyalty |

datasciencedojo
data science for everyone

# A/B Testing Tools

# Questions?

# APPENDIX & COMMON PITFALLS

# Pitfalls in Online Experimentation

1. Picking an overall evaluation criteria (OEC) for which it is easy to beat the control

2. Incorrectly computing the confidence intervals

3. Using standard statistical formulas for computation of variance and power

4. Combining metrics over periods where proportions assigned to Control and Treatment vary or over subpopulations sampled at different rates

5. Neglecting to filter bots

6. Failing to validate each step of the analysis pipeline and the OEC components

7. Forgetting to control for all differences, and assuming that humans can keep the variants in sync

# Pitfall 1: Picking an Easy-to-Beat Overall Evaluation Criteria (OEC)

- Before running an experiment an OEC is selected
- OEC should be tied to a long term goals as opposed to short term goals:
  - For example, click-through rate (CTR) vs. long term revenue)
- Loyal/repeat users get more weight?
- Sometimes getting the true metric is hard:
  - For example, high CTR does not necessarily mean high conversion rate

datasciencedojo
data science for everyone

# Pitfall 1: Picking an Easy-to-Beat Overall Evaluation Criteria (OEC)

- Measuring click through on a small area of the page, ignoring the impact on other areas:
  - What if the small area on the page was bold/flashing/high contrast?
  - What happens to the whole page CTR?
- Is 'time on site' a good OEC?
  - What if the treatment has a reduced user's effectiveness?

datasciencedojo

data science for everyone

# Pitfall 2: Incorrect Computation of Confidence Intervals

- **Hypothesis Test:** Determines whether there is a statistically significant difference in the means of the control and the treatment

- **Confidence Interval:** Provides a plausible range of the size of the effect (i.e., difference between the means of the control and treatment means)

datasciencedojo

— data science for everyone —

# Pitfall 2: Incorrect Computation of Confidence Intervals

$$0.95 = 1 - \alpha = P(-z \le Z \le z) = P\left(-1.96 \le \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \le 1.96\right)$$

$$= P\left(\bar{X} - 1.96\frac{\sigma}{\sqrt{n}} \le \mu \le \bar{X} + 1.96\frac{\sigma}{\sqrt{n}}\right)$$

$$= P\left(\bar{X} - 1.96 \times 0.5 \le \mu \le \bar{X} + 1.96 \times 0.5\right)$$

$$(\bar{x} - 0.98; \bar{x} + 0.98) = (250.2 - 0.98; 250.2 + 0.98) = (249.22; 251.18).$$



**Confidence interval implies:** We are 95% confident that average for all cups of coffee from the machine is between 249.22 and 251.18 grams.

# Pitfall 2: Incorrect Computation of Confidence Intervals

- Confidence interval should be formed out of absolute difference

- Do not form a confidence interval around percent change. Percentage change involves dividing by a random variable.

- Some techniques to compute CI are mentioned when the OEC is a linear/non-linear combination of metrics that have the same/different basis/experimental unit.

datasciencedojo
data science for everyone

# Pitfall 3: Standard Statistical Formulas for Computation of Variance and Power

- Variance of the metric is needed to compute the statistical significance

- Variance estimates using standard statistical formula for some families of metrics are inaccurate

- This happens when the experimental unit used in random assignment is different from the experiment unit used in the calculation of the metric.

# Pitfall 3: Standard Statistical Formulas for Computation of Variance and Power

- Variance, Power and Sample size estimates may be wrong if care is not taken
- How to correct this?
  - <u>Bootstrap method</u>: Estimate variance using bootstrap samples and compare with the variance from standard formula
- This should be done for all metrics and especially for the one with different experiment and randomization units

# Pitfall 4: Simpson's Paradox

- Unintuitive but not uncommon
- Simpson's paradox: 'A correlation or trend present in different groups is reversed when the groups are combined'.

|  | Treatment A | Treatment B |
|---|---|---|
| Small Stones | Group 1 | Group 2 |
|  | 93% (81/87) | 87% (234/270) |
| Large Stones | Group 3 | Group 4 |
|  | 73% (192/263) | 69% (55/80) |
| Both | 78% (273/350) | 83% (289/350) |

# Pitfall 4: Simpson's Paradox

- 1 million visitors/day
- On Friday the treatment ran with 1% traffic
- On Saturday, the allocation was raised to 50%.
- If we consider Friday and Saturday separately T has a better CTR
- T's CTR is worse when aggregated over days

**Table 1: Conversion Rate for two days.**
**Each day has 1M customers, and the Treatment (T) is better than Control (C) on each day, yet worse overall**

|   | Friday C/T split: 99%/1% | Saturday C/T split: 50%/50% | Total |
|---|---|---|---|
| C | $\frac{20,000}{990,000} = 2.02\%$ | $\frac{5,000}{500,000} = 1.00\%$ | $\frac{25,000}{1,490,000} = 1.68\%$ |
| T | $\frac{230}{10,000} = 2.30\%$ | $\frac{6,000}{500,000} = 1.20\%$ | $\frac{6,230}{510,000} = 1.20\%$ |

It is possible to have $\frac{a}{b} < \frac{A}{B}$ and $\frac{c}{d} < \frac{C}{D}$ while $\frac{a+c}{b+d} > \frac{A+C}{B+D}$

# Pitfall 4: Simpson's Paradox – A Scenario in Controlled Experiments

➢ Sampling of users with non uniform sampling to make sure all browsers have a representative sample

➢ Overall results show treatment is better than control but when segmented by browser, control looks better than treatment for each browser

# Pitfall 5: Ignoring Bot Traffic

➢ For experimentation, we are interested in removing bots/fraud clicks that are not uniformly distributed across the control and treatment

➢ Uniformly distributed bots will only reduce the power of the experiment

# Pitfall 5: Ignoring Bot Traffic

Failing to exclude bot traffic and fraud clicks may invalidate the results of an experiment

datasciencedojo
— data science for everyone —

# Pitfall 6: Failing to Validate Each Step of Analysis

It is important to keep a check on the health of the pipeline

➢ Assignment of users to experiment variants

➢ Calculation of metrics

➢ Any abnormal shift in metrics

➢ Movement of metrics that are not expected to move

➢ Broken instrumentation

datasciencedojo
— data science for everyone —

# Pitfall 6: Failing to Validate Each Step of Analysis

**Logging Tests:**

- Compare with real historical data
- Compare with generated data
- Look for **unexpected patterns**
  - Volume of data over time
  - New and repeat users over time
  - Abnormal shift in any of the metrics
- A/A Tests
- Rich Instrumentation

# Pitfall 7: Failing to 'Control' the Control

- **Don't allow any difference** between the Control and the Treatment besides what is actually being tested

- If the **Treatment** has some **updates**, **Control** should have them too and vice versa

datasciencedojo
data science for everyone

# Pitfall 7: Failing to 'Control' the Control

- If the site is receiving **frequent updates**, these updates should be **applied equally** to the control and the treatment

- Forgetting to **control for all differences**, and assuming that humans can keep the variants in sync.

# HUMOR

# Have you heard the latest statistics joke?

Probably....

datasciencedojo
data science for everyone

How many statisticians does it take to change a light bulb?

**1 – 3.    α=0.05 (.95 confidence)**

datasciencedojo

data science for everyone

# Did you hear about the statistician who was thrown in jail?

He now has zero degrees of freedom.

datasciencedojo
data science for everyone

A statistician's wife has twins. He was delighted, and he called to tell his minister the good news.

"Excellent!", said the minister. "Bring them to church on Sunday and we'll baptize them."

"No," replied the statistician. "Let's just baptize one. We'll keep the other as control."

*Three statisticians go out hunting together. After a while they spot a solitary rabbit.*

*The first statistician takes aim and overshoots. The second aims and undershoots.*

*The third shouts out "We got him!"*