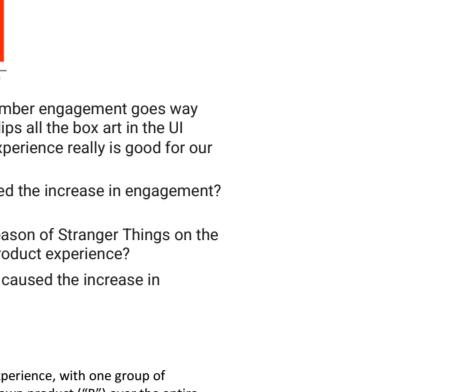
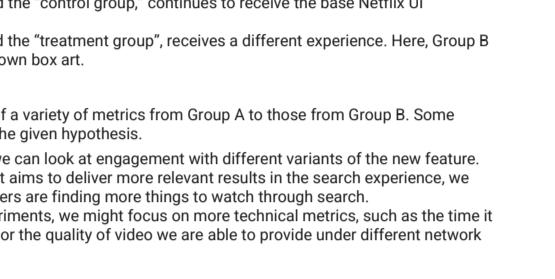


Lecture Objective:

- Business sense towards A/B testing.
- How to launch A/B testing to get "valid" and "reliable" results.
 - Developing hypothesis
 - Designing A/B tests
 - Evaluating test results
 - Making decisions
 - Launch recommendations and pitfalls with A/B testing set ups.

Example Cases: A/B testing Use Cases→ Certain Group

Population ↓ after ↓

Population with filter ↓

sample → samples sample

A B

Control Test

[normal] [upside down]

Group A B

40 to 60 80 to 100

4000 4000

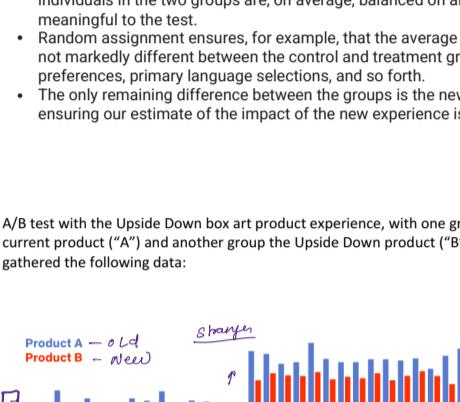
30 days

What is an A/B Test?

An A/B test is a simple controlled experiment. Let's say - (Not ever close to reality) Let say we want to learn if a new product experience that flips all the box art in the UI upside down, how confident would you be that the new product experience really is good for Netflix members.

Without A/B Test

Let's say Experiment was implemented, and flipped the switch to the Upside Down experience on the 16th day of a month. How would you act if we gathered the following data



The data look good: Release of a new product experience and member engagement goes way up! But if you had these data, plus the knowledge that Product B flips all the box art in the UI upside down, how confident would you be that the new product experience really is good for our members?

Can we really know that the new product experience is what caused the increase in engagement? What other explanations are possible?

What if you also knew that Netflix released a hit title, like a new season of Stranger Things on the same day as the (hypothetical) roll out of the new Upside Down product experience?

The key point is that we don't know if the new product experience caused the increase in engagement

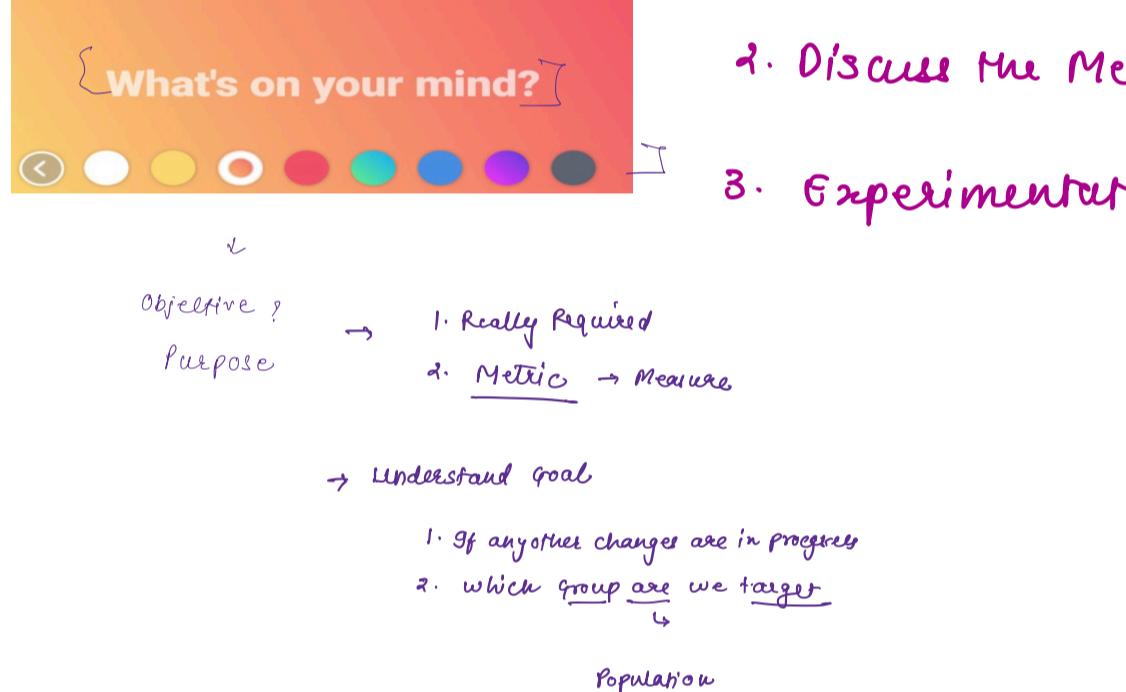
What if instead we'd run an A/B test with the Upside Down box art product experience, with one group of members receiving the current product ("A") and another group the Upside Down product ("B") over the entire month?

To run the experiment, Lets take a subset of members, usually a simple random sample, and then use random assignment to evenly split that sample into two groups.

- Group "A," often called the "control group," continues to receive the base Netflix UI experience,
- Group "B," often called the "treatment group," receives a different experience. Here, Group B receives the Upside Down box art.

Then compare the values of a variety of metrics from Group A to those from Group B. Some metrics will be specific to the given hypothesis.

- For a UI experiment, we look at engagement with different variants of the new feature.
- For an experiment that aims to deliver more relevant results in the search experience, we can look at if members are finding more things to watch through search.
- In other types of experiments, we might focus on more technical metrics, such as the time it takes the app to load, or the quality of video we are able to provide under different network conditions.



Holding everything else constant

- Because there are control ("A") and treatment ("B") groups using random assignment, individuals in the two groups are, on average, balanced on all dimensions that may be meaningful to the test.
- Random assignment ensures, for example, that the average length of Netflix membership is not markedly different between the control and treatment groups, nor are content preferences, primary language selections, and so forth.
- The only remaining difference between the groups is the new experience we are testing, ensuring our estimate of the impact of the new experience is not biased in any way.

A/B test with the Upside Down box art product experience, with one group of members receiving the current product ("A") and another group the Upside Down product ("B") over the entire month, and gathered the following data:



In this case, we're led to a different conclusion: the Upside Down product results in generally lower engagement (not surprisingly!), and both groups saw an increase in engagement concurrent with the launch of the big title

A/B tests let us make causal statements. We've introduced the Upside Down product experience to Group B only, and because we've randomly assigned members to groups A and B, everything else is held constant between the two groups. We can therefore conclude with high probability that the Upside Down product caused the reduction in engagement

Using a framework based on hypothesis generation, A/B testing, and statistical analysis allows us to carefully quantify uncertainties, and understand the probabilities of making different types of mistakes.

There are two types of mistakes we can make in acting on test results.

A false positive* (also called a Type I error) occurs when the data from the test indicates a meaningful difference between the control and treatment experiences, but in truth there is no difference. This scenario is like having a medical test come back as positive for a disease when you are healthy

A false negative (also called a Type II error), occurs when the data do not indicate a meaningful difference between treatment and control, but in truth there is a difference.

**Case Study:**

How would you test if Facebook incorporating colored backgrounds to statuses improves user engagement?

**1. Clarify the goal****2. Discuss the Metrics****3. Experimentation**

→ Objective ?
Purpose → 1. Really Required
a. Metric → Measure

→ Understand goal

- If any other changes are in progress
- which group are we targeting

Population ↓

Population - filters

→ Define the Metric

→ Engagement → Active users

+ Engagement %

↳ Likes
↳ shares etc

3. hypothesis

Null H: No significant change in engagement

Alternate → significant increase in engagement

4. Test: 2 groups (Ratio)

↓

Z - proportion test

5. find out Control group / Test group

p ↓

p-filterout C Target group

size ← Random sample + web app

R1 R2

size ← 100 → 200

2. 100% FN ↓

400 → place 1 → 200 ↓

400 → place 2 → 200 ↓

4000 → place 3 → 200 ↓

4000 → place 4 → 200 ↓

4000 → place 5 → 200 ↓

4000 → place 6 → 200 ↓

4000 → place 7 → 200 ↓

4000 → place 8 → 200 ↓

4000 → place 9 → 200 ↓

4000 → place 10 → 200 ↓

4000 → place 11 → 200 ↓

4000 → place 12 → 200 ↓

4000 → place 13 → 200 ↓

4000 → place 14 → 200 ↓

4000 → place 15 → 200 ↓

4000 → place 16 → 200 ↓

4000 → place 17 → 200 ↓

4000 → place 18 → 200 ↓

4000 → place 19 → 200 ↓

4000 → place 20 → 200 ↓

4000 → place 21 → 200 ↓

4000 → place 22 → 200 ↓

4000 → place 23 → 200 ↓

4000 → place 24 → 200 ↓

4000 → place 25 → 200 ↓

4000 → place 26 → 200 ↓

4000 → place 27 → 200 ↓

4000 → place 28 → 200 ↓

4000 → place 29 → 200 ↓

4000 → place 30 → 200 ↓

4000 → place 31 → 200 ↓

4000 → place 32 → 200 ↓

4000 → place 33 → 200 ↓

4000 → place 34 → 200 ↓

4000 → place 35 → 200 ↓

4000 → place 36 → 200 ↓

4000 → place 37 → 200 ↓

4000 → place 38 → 200 ↓

4000 → place 39 → 200 ↓

4000 → place 40 → 200 ↓

4000 → place 41 → 200 ↓

4000 → place 42 → 200 ↓

4000 → place 43 → 200 ↓

4000 → place 44 → 200 ↓

4000 → place 45 → 200 ↓

4000 → place 46 → 200 ↓

4000 → place 47 → 200 ↓

4000 → place 48 → 200 ↓

4000 → place 49 → 200 ↓

4000 → place 50 → 200 ↓

4000 → place 51 → 200 ↓

4000 → place 52 → 200 ↓

4000 → place 53 → 200 ↓

4000 → place 54 → 200 ↓

4000 → place 55 → 200 ↓

4000 → place 56 → 200 ↓

4000 → place 57 → 200 ↓

4000 → place 58 → 200 ↓

4000 → place 59 → 200 ↓

4000 → place 60 → 200 ↓

4000 → place 61 → 200 ↓

4000 → place 62 → 200 ↓

4000 → place 63 → 200 ↓

4000 → place 64 → 200 ↓

4000 → place 65 → 200 ↓

4000 → place 66 → 200 ↓

4000 → place 67 → 200 ↓

4000 → place 68 → 200 ↓

4000 → place 69 → 200 ↓

4000 → place 70 → 200 ↓

4000 → place 71 → 200 ↓

4000 → place 72 → 200 ↓

4000 → place 73 → 200 ↓

4000 → place 74 → 200 ↓

4000 → place 75 → 200 ↓