

- Guideline: Define ; Build Intuition ; characterize

2. Metrics Define

1. Overview: Before we define a single metric

① Some data are not tractable or take

② Need to define them accurately

Active user → How to define 'active' → check activity should take into account

③ Need to think all measurements and then summarize them into a single metric

④ One metric v.s multiple metrics

• One

(1) Every team towards same goal ; (2) Good for PR and external report;

(3) Can create a single composite metric (use objective function)

• Multiple

(1) Can know how other things move and find the reason for the movement of

important metrics, or composite metrics

(2) Won't over optimize one thing when don't look at how other things move

• Consider how generally applicable the metric is

If run a whole suite of AB tests, ideally you'd have one or more metrics that can use across the entire suite. It's much better to use a metric that's less optimal for your test if you can use it across the suites that you can do comparison.

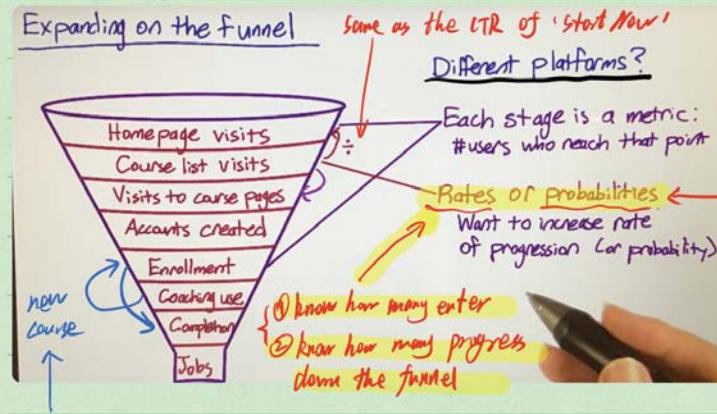
Custom metrics will increase the risk.

2. High level concept for metrics

Business Objective: { Helping students to find jobs
Financial sustainability

3. Expanding on the funnel

- Homepage visits
- Exploring the site
 - # users who view course list
 - # users who view course details
- Create an account
 - # users who enroll in a course
 - # users who finish Lesson 1, Lesson 2, etc
 - # users who sign up for coaching at various levels
- Completing a course
 - # users who enroll in a second class
 - # users who get jobs



- (Count)
• only use number for a few
key points
- prob. that unique users
progress down the funnel
(Binary case)

Time line for a particular user → more like a swirl (might jump forward and back)

For viewing a list of courses { rate: how easy it is to do it from the Course Overview Page
prob: whether a student reach the course list at all
(Binary probability)

4. Choosing metrics

Choosing Metrics



1. Click-through-rate on "Start Now" button
2. Click-through-probability on "Start Now" button
3. Probability of progressing from course list to course page
4. Probability of progressing from course page to enrolling
5. Probability that enrolled student pays for coaching

Choosing Metrics

- Update a description on the course list page

3 Continued progression down funnel

- Increase size of "Start Now" button

1 Rates better for usability test

- Explain benefits of paid service

5 User retention or usage

1. Click-through-rate on "Start Now" button
2. Click-through-probability on "Start Now" button
3. Probability of progressing from course list to course page
4. Probability of progressing from course page to enrolling
5. Probability that enrolled student pays for coaching

Size related to how easy to be found

for whom used this coaching, do they use more often now that it's more clear that coaches can do

∴ it's of usability \Rightarrow rates is better for it.

5. Difficult Metrics

- ① { Don't have access to data
Take too long

② Hard metrics example

- Rate of returning for 2nd course : Has data, but too long
- Average happiness of shoppers : Doesn't have data
- Probability of finding information via search : Doesn't have data

③ Techniques for approxi of hard metrics

- Survey, retrospective analyses, focus groups :
- Brainstorm new metrics and validating possible metrics

(a) external data

Eg: company collect market data such as market share,
company run survey on users
Academic research

(b) internal data get baseline/theory based on past metrics and related changes

existing data: retrospective analysis, running experiments
gather new data: user experiment research, survey, focus group problem:

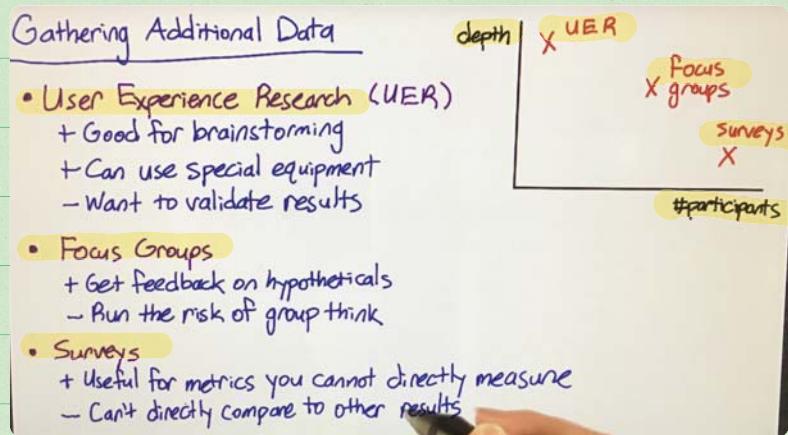
Show correlation but not Causation (actual controlled experiment)

The only way to validate is run experiment

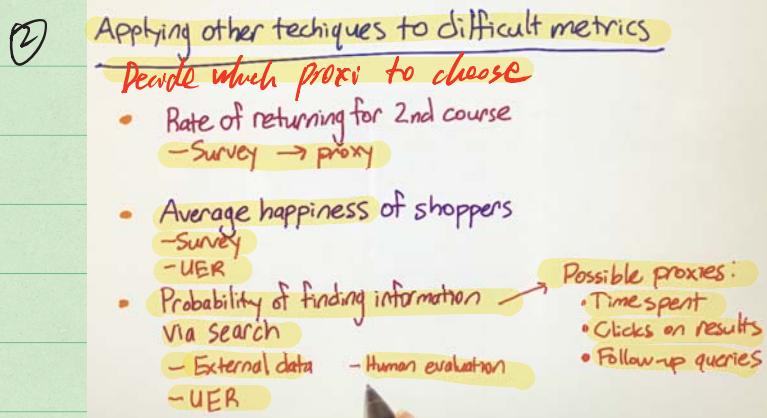
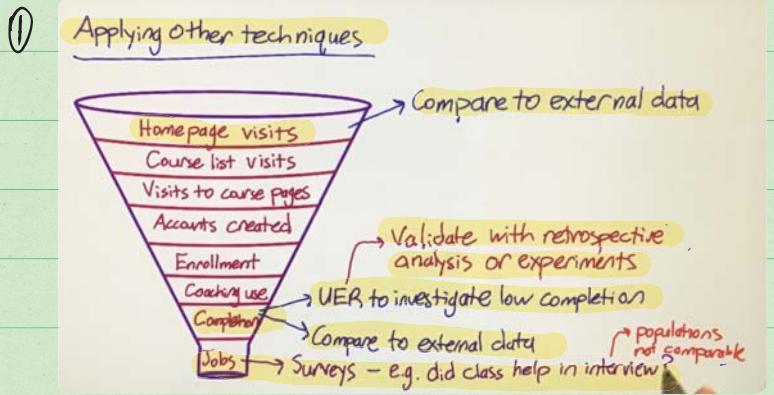
- Run an experiment to determine whether or not we can use a metric for evaluating experiments

(Want to make sure the metrics will move as we make changes)

6. Gathering Additional data



7. Use other techniques



③

Which techniques would you use?

- Measure user engagement
Course completion too long-term

4 2 5

- Decide whether to extend inventory

3 1

- Which ads get most views

1 2 5

1. External data

2. UER

3. Focus group

4. Survey

5. Retrospective analysis

6. Experiments

Use historical data to predict or set baseline

II. Data capture and build intuition

1. Data capture

Note: CTR may different based on different browsers.

It's not because the actual CTR different, but because the technology you use in order to gather the clicks (different failure rate for JAVA commands)

2. Define metrics : from high level to detailed

① Need to decide details like time interval

Defining a metric

High-level metric: Click-through probability = $\frac{\# \text{ users who click}}{\# \text{ users who visit}}$

Def #1: For each <time interval>, $\frac{\# \text{ cookies that click}}{\# \text{ cookies}}$

Def #2: $\frac{\# \text{ pageviews with click within } <\text{time interval}>}{\# \text{ pageviews}}$

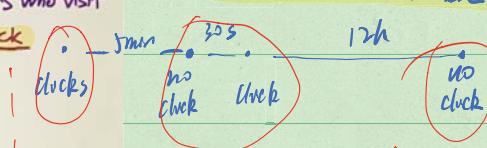
Def #3: $\frac{\# \text{ clicks}}{\# \text{ pageviews}}$ (Click-through rate)

pageview pageview click

30s 1s click

Def 1: per minute = 1 Def 2: per minute = $\frac{1}{2}$

• A user with same cookie



$$\text{per minute} = \frac{1}{2}, \text{ per hour} = \frac{1}{2}, \text{ per day} = 1$$

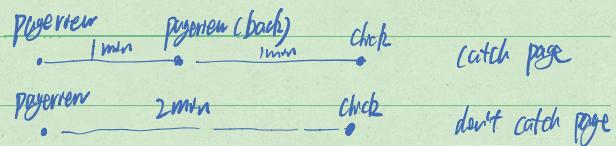
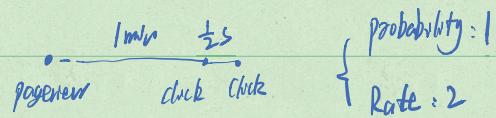
Def 2 is easier to compute

② Problem for metrics

Which metrics have which problems?

	1: Cookie Prob	2: PageView Prob	3: Rate
Double click	✓	✓	□
Back button catches page	✓	□	□
Click-tracking bug	□	□	□

Just record pageview but don't record click



$$\text{cookie: } 1 \quad \text{pageview: } \frac{1}{2}$$

3. Filtering and segmenting

Eg: Just want to make change for English users

Filter makes sense because you don't want to dilute results.

Filter to affected traffic, can increase power and sensitivity of experiment

① Goal: de-bias data

when filter out spam and fraud, don't introduce bias into your data

Eg: (a) metrics that can only be measured on logged-in user:

don't have data for non-logged users → biased

(b) Filtering out long or weird session of user behavior

Before do that, need to check and make sure it's not actually your website,

metrics, or even your logging that's causing these sessions to come up.

② How to check: compute baseline value for metric

- slice data: compute metric on a bunch of disjoint sets. eg: country, language

- compute a metric on all the different slices:

Show if we move traffic disproportionately: indication of bias

- look at week over week or day over day traffic pattern changes: useful to spot something unusual.

• Key for filtering: building intuition

Have to know what changes are going to be expected v.s. unexpected

When see data for real experiment: can know if have a problem

Do I believe this, what's really going on.

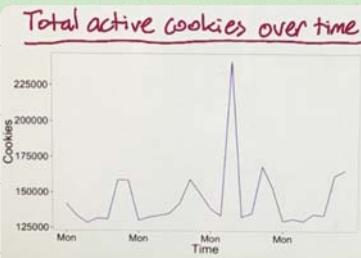
③ Check if it's seasonality effect

Segmenting and filtering data

Good for evaluating definitions and building intuition



④ Dig into what cause of spike: look at different segments



⑤ Detect suspect click tracking issue: clicks counted twice on mobile

which graphs would confirm this problem.

Which graphs would confirm this problem?

Click-through-rate over time



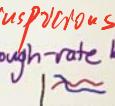
Click-through-probability over time



Both rate and probability on same graph



Click-through-rate by platform



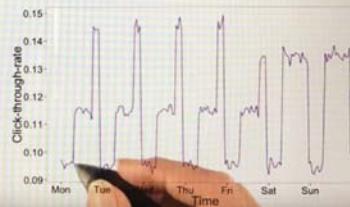
Click-through-probability by platform



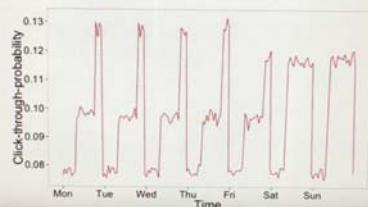
Both rate and probability by platform



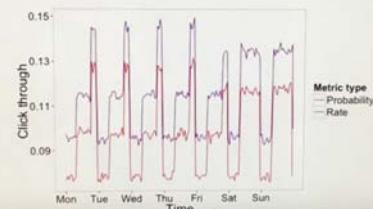
Click-through-rate over time - No



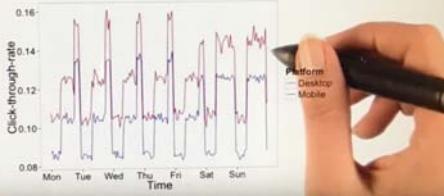
Click-through-probability over time - No



Both rate and probability on same graph - No



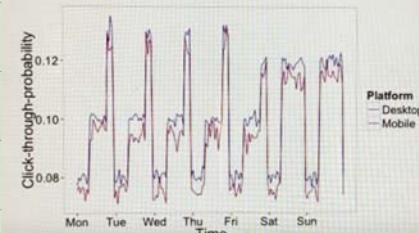
Click-through-rate by platform - Suspicious!



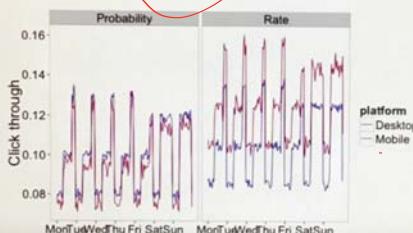
Don't know how much higher it's weird, and
the difference is consistent.

May have different behavior

Click-through-probability by platform - No



Both rate and probability by platform - Yes



prob. eliminates effect of count twice

P: same but R: much higher

4. Summary Metrics

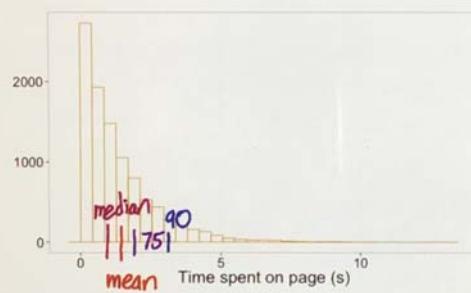
① Use metrics to summarize individual events

Categories of summary metrics

- Sums and counts
e.g. # users who visited page
- Means, medians, and Percentiles
e.g. mean age of users who completed a course or median latency of page load
- Probabilities and rates
- Probability has 0 or 1 outcome in each case
- Rate has 0 or more
- Ratios
e.g. $\frac{P(\text{revenue-generating click})}{P(\text{any click})}$

② Choose metric based on distribution

Means, medians, and percentiles



Possible metrics

- median
- mean
- 75th percentile
- 90th percentile
- Percent of users who spend at least 3 seconds

5. Sensitivity and robustness of metric (sensitivity)

- ① Idea: Choose a metric that picks up changes you care about but doesn't move a lot when nothing interesting happens
eg. median is more robust than mean (robustness)

② Measure sensitivity and robustness

(a) run experiments / using experiments you already have

Data should move in a way that intuitively make sense

- eg. increase quality will leads to increase in load time:
we could see if the metrics actually response to that

(b) A vs. A experiment:

to determine if they are too sensitive

Don't change anything, just compare people who saw the same thing to each other. See if metrics pick up any spurious difference between the two.

(Make sure not to call things significant that may not really meaningful)

(c) Look back to experiments run earlier.

(how the metric change in past experiments)

(d) Retrospective analysis

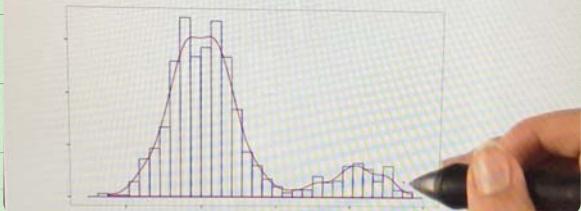
- look back at changes and see if metrics you're interested in actually moved in conjunction with those changes
- look at history of metrics to see if you can find a cause for any major changes that you see.

③ Measure sensitivity and robustness

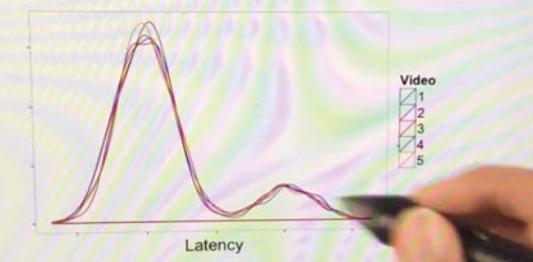
Measuring sensitivity and robustness

Choose summary metric for latency of a video

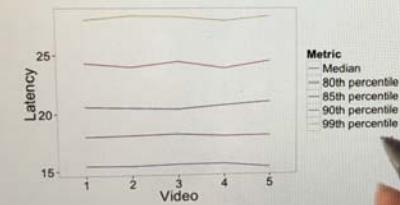
Distribution for a single video:



Distribution for similar Videos



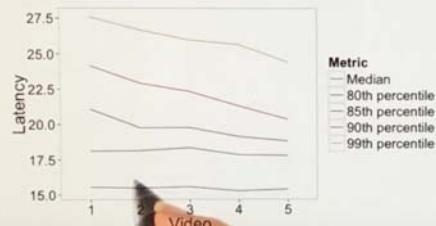
Distribution for similar videos



90th and 99th percentile: not robust enough

∴ 2nd-2nd for different videos

Distribution for experimental videos



90th and 99th percentile: not robust enough

Median and 80th percentile: not sensitive enough

∴ don't change when resolution changes

→ with lower resolution

6. Absolute / relative difference

start with absolute and then use relative (better for comparison)

III. Characterize

1. variability { use distribution
compute empirically

2. Calculating variability

① parametric way : get a nice distribution of data

To calculate CZ, you need { variance
distribution

$$\text{eg: For binomial: } SE = \sqrt{\frac{P \cdot (1-P)}{N}}, \quad m = 2^k \cdot SE$$

Other summary metrics may be harder to analyze
e.g. median — could be non-normal if data is non-normal



type of metric	distribution	estimated variance
probability	binomial (normal)	$\frac{P(1-P)}{N}$
mean	normal	$\frac{\sigma^2}{N}$
median/percentile	depends	depends
count/difference	normal (maybe)	$\text{Var}(X) + \text{Var}(Y)$
rates	poisson	\bar{X}
ratios	e.g. $\frac{\text{Pop}}{\text{Point}}$ instead of $\text{Pop} - \text{Point}$	

Confidence interval for a mean

Measure: Mean number of homepage visits per week

$$N_1 = 87,029$$

$$N_2 = 113,407$$

$$N_3 = 84,843$$

$$N_4 = 104,994$$

$$N_5 = 99,327$$

$$N_6 = 92,052$$

$$N_7 = 60,684$$

$$m = z^* \cdot SE$$

$$= 1.96 \cdot SE$$

$$= 12,605$$

$$\bar{N} = \frac{N_1 + \dots + N_7}{7} = 91,762$$

$$\sigma = SDC(N_1, \dots, N_7) = 17,015$$

$$SE = \frac{\sigma}{\sqrt{7}} = 6430 \quad 95\% \text{ confidence interval}$$



② Non-parametric way : use empirical variability

: for complicated metrics, the distribution may be very weird

Even for simple metrics, ~~analytical~~ estimate of the variance ended up ~~underestimate~~

∴ Use A/A experiments across the board to estimate empirical variability

- difference in measures are due to underlying variability:

System, population, user behavior

• When don't have enough sample for A/A test: bootstrap

- If bootstrap estimate is similar to analytical estimates, move on
- - - - - not -- -, do A/A test to see what's going on.

③ Calculating variability empirically

* Use of A/A test

(a) Compare results to what you expect using A-A test

- | | |
|---|--------------------------|
| { 20 experiments, each on 0.5% of traffic | : 50 users in each group |
| 20 more, each on 1% | : 100 users — |
| 10 more, each on 5% | : 500 users |

How many experiments will show a statistically significant difference at 95% level? 1/20

Calculating variability empirically



Range becomes tighter when size ↑

↑ We may still regard it as normal

(b) Estimate variance and σ^2 using A-A test

Calculating variability empirically

Estimate variance and calculate confidence interval:

Since we expect a normal distribution:

$$\begin{aligned} m &= SD \cdot z^* \quad \text{get SD empirically} \\ &= 0.059 \cdot 1.96 = 0.116 \quad \text{empirically} \end{aligned}$$

$$\text{Analytically: } SE = \sqrt{P_{\text{fail}}(1-P_{\text{fail}})(\frac{1}{N_{\text{start}}} + \frac{1}{N_{\text{end}}})}$$

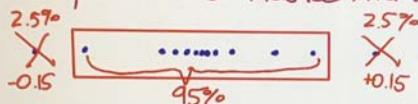
Slightly different margin of error for each experiment

每个都有点不一样

(c) Directly estimate confidence interval use A-A test

Calculating variability empirically

Directly estimate confidence interval:



Since we have 20 data points, dropping the highest and the lowest gives a 90% confidence level: $-0.1 \text{ to } 0.06$

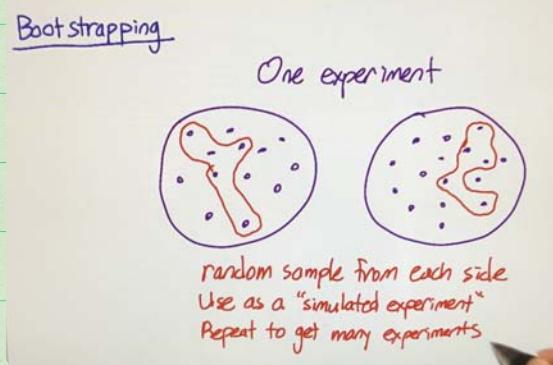
$$\text{Empirical standard deviation: } 0.059 \cdot 1.65 = 0.097$$

z-score for 90% confidence

Not accurate for this case because

number of A-A test too small

(d) bootstrap for A/A test (when no enough sample)

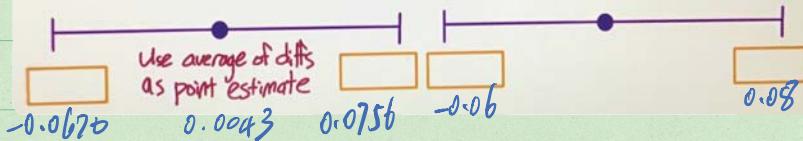


Calculating a confidence interval empirically:

- For each experiment, calculate the difference in click-through-probability between the two groups

Calculate the standard deviation of the differences, and assume metric is normally distributed.

Calculate an empirical confidence interval, making no assumptions about the distribution



$$m = 0.036 * 1.96 = 0.0713 \quad \text{Find first/last two after sort}$$

SD z-score

3. Variability Summary

- ① If too high, may not practical to use in evaluating experiments, even if the metric makes a lot of business or product sense
- ② To compute variability, need to understand distribution of underlying data
Can use both analytical and empirical techniques for compute variability

IV. Lessons learned

① There are many ways to define and aggregate metrics (related to sensitivity)

· Need to build intuition, understand data/system and work with engineer to understand how data is captured.

② Variability: start with analytical way

if distribution needed (revenue per query): easier to compute empirically

③ About metrics

Necessity of invariance / sanity checking

V. Conclusion:

① Establishing high level concept → ② Go to detail to compute the metrics

→ ③ How to build intuition about metrics, how appropriate they'll be for different experiments.

→ ④ Characterize variability of metrics