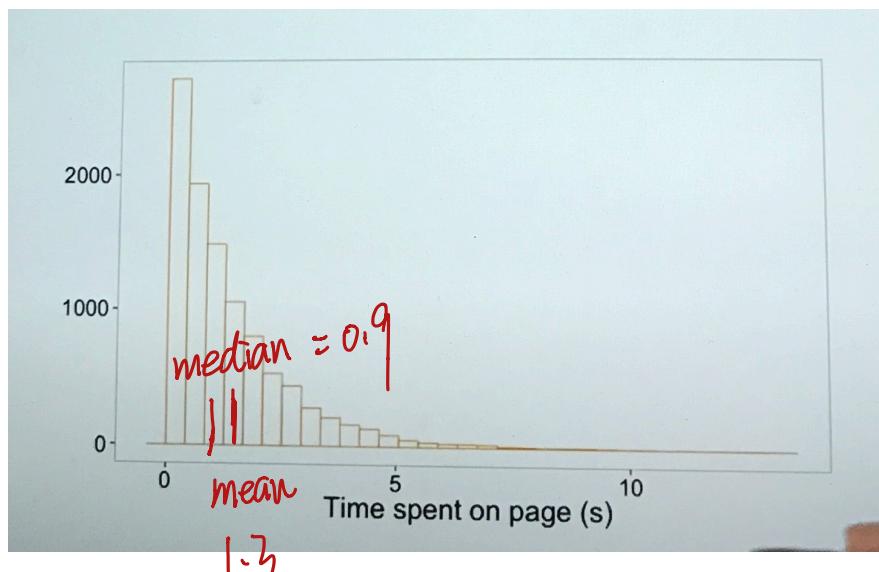


$$\text{e.g. } \frac{P(\text{revenue-generating click})}{P(\text{any click})}$$

①为什么是指示器.

Means, median, and percentage (指看分布) ②什么是频率.

⇒ 用户在网页的某个页面花多长时间.



How to measure sensitivity and robustness?

1. running experiments or using experiment you already have
- + a versus a experiments to determine if they're too sensitive.
↳ It's an experiment where you don't change anything. You just compare people who saw the same thing to each other.

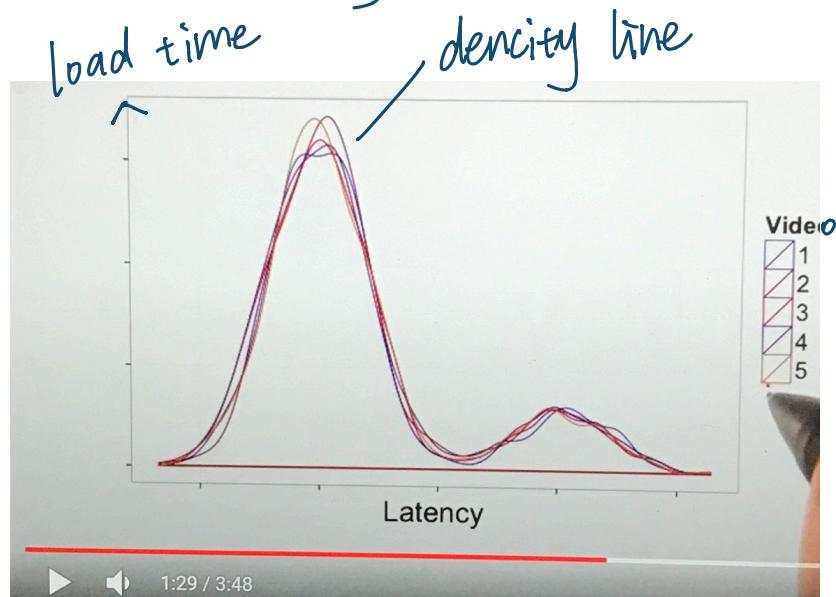
2. Sort of a retrospective analysis of your logs: 对你记录的回放性能分析

If you don't have experiment data or you can't run new experiments, you can look back at the changes you know you made to your site, and see if the metrics you're interested in actually moved in conjunction with those changes.

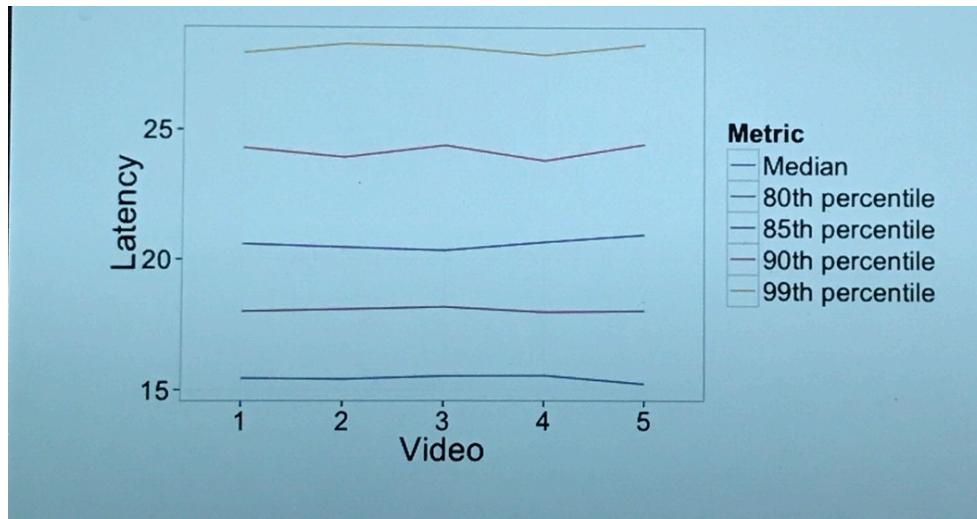
or look at the history of the metric and see if you can find a cause for any major changes that you see.

Measuring sensitivity and robustness

① - choose summary metric for latency of a video



有2个高峰. 1个加载时间相对长



90th 与 99th 是 Z shape

↳ not robust enough

操作直线

Mid and 80th is more

sensitive enough.

② look at the actual experiment.

↳ distribution for experimental video

Absolute & relative difference

1. 首先如何计算这个 comparison:

现在获得的是实验组的数据 + control data

方法: 1. difference

又称 percent change

2. relative change 而不是 absolute change

然后 choose one practical significance boundary to get

stability over time.

applicable situation: seasonality, system is changing

缺点: Variability.

Absolute vs. relative difference

Suppose you run an experiment where you measure the number of visits to your homepage, and you measure 5000 visits in the control and 7000 in the experiment. Then the absolute difference is the result of subtracting one from the other, that is, 2000. The relative difference is the absolute difference divided by the control metric, that is, 40%.

Relative differences in probabilities

For probability metrics, people often use percentage points to refer to absolute differences and percentages to refer to relative differences. For example, if your control click-through-probability were 5%, and your experiment click-through-probability were 7%, the absolute difference would be 2 percentage points, and the relative difference would be 40 percent. However, sometimes people will refer to the absolute difference as a 2 percent change, so if someone gives you a percentage, it's important to clarify whether they mean a relative or absolute difference!

Variability: {
 确定实验规模
 分析置信区间 和得出结论

使用 count 或 probability, only dealing with the variability of a single measurement, or a constrained one, in case of probability.

使用 ratios or percentiles, compute the variability empirically. 根据经验

Calculating variability.

To calculate a confidence interval, you need:

- Variance (or standard deviation)

- Distribution

Binomial distribution: 在哪方面利用了这是一个二项分布的事实:

$$SE = \sqrt{\frac{\hat{P}(1-\hat{P})}{N}} \leftarrow \text{利用是二项分布的事实得出标准误差的公式}$$

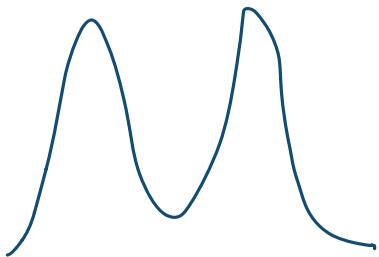
$$m = Z^* \cdot SE \leftarrow \text{我们假设这是一个正态分布.}$$

△ N 越大, 二项分布越趋近于正态分布.

type of metric	distribution	estimated variance
probability	binomial (normal)	$\frac{\hat{P}(1-\hat{P})}{N}$
	根据中心极限定理	
mean	normal (if sample size is larger)	$\frac{\hat{\sigma}^2}{N}$
△ 注意样本方差与实际指标方差的区别:		
样本方差: 所有数据点的方差. take each data point and collect the variance of them.		
实际指标方差: (Variance of the actual metric): if you were to collect a new sample, how would you expect this metric to vary		
median/percentile	depends on the distribution of the underlying data	depends on the distribution of the underlying data

Other summary metrics may be harder to analyze:

e.g. median - could be non-normal if data is non-normal



e.g. latency is a bimodal distribution.



mean 是正态分布，而 mid 不是
为了估计单位数的方差，你需要
要假设基础数据的分布情况
根据你的假设，你通常轻松

估计其分析方差。

→ 对于 demographic data

各个变量的方差之和

count / difference

normal (maybe)

$\text{Var}(X) + \text{Var}(Y)$

rates → 实验组

poisson

\bar{x} 平均率

e.g. $\frac{\hat{P}_{\text{exp}}}{\hat{P}_{\text{cont}}}$ instead of $\hat{P}_{\text{exp}} - \hat{P}_{\text{cont}}$
 $\hat{P}_{\text{exp}} \rightarrow$ 对照组

ratios

depends on 分子和分母
的分布

depends on 分子和
分母的分布。

Confidence interval for a mean

Measure: Mean number of homepage visits per week

$$N_1 = 87,029$$

$$N_2 = 113,407$$

$$N_3 = 84,843$$

$$N_4 = 104,994$$

$$N_5 = 99,327$$

$$N_6 = 92,052$$

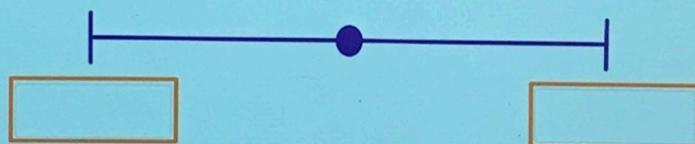
$$N_7 = 60,684$$

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

$$\sigma = \text{SD}(N_1, \dots, N_7)$$

$$SE = \frac{\sigma}{\sqrt{7}}$$

use spreadsheet
or programming language
95% confidence interval



$$\bar{N} = 91762.29 \quad \sqrt{7} = 2.646$$

$$G = 17014.80$$

$$SE = 6470.99$$

$$(\bar{X} - Z \cdot \frac{G}{\sqrt{n}}, \bar{X} + Z \cdot \frac{G}{\sqrt{n}})$$

Nonparametric answer

Nonparametric method: 不需对分布情况作假设，也能分析数据。

e.g. sign test (符号检验).

持续20天的AB test, 有15天实验组测量数据高于对照组。

使用二项分布的方法计算结果的实际发生概率，并检查实验组和对照组之间是否有差异。

缺点：不能估计效应的大小

优点：easy to do.

当得出经验方差后：

choose 1: summary statistic distribution is nice and normal

→ use a normal confidence interval with the variants you estimated empirically.

choose 2: data is funnier / want to be really robust

→ compute a non-parametric confidence interval.

实证.

Empirical variability

针对复杂的指标，使用实证预测方差，而不是通过分析计算。

计算指标方差时，需要对数据潜在的分布有一个假设。

对于复杂指标，分布会很复杂，所以使用实证。

e.g. google 使用 A/B 测试从而通过实证预估所有指标的差异性。

A/B : {
A → 对照
B → 实验

A/A : {
A → 对照
A → 第一个对照

A/A → 你测量的任何区别是因为潜在的差异性。e.g. system of user population, what users are doing.

边际效益递减

there's a diminishing return as you run more A/A test. The key rule is that standard deviation is going to be proportional to the square root of the number of sample. 标准偏差与样本量的平方根

端正书

bootstrap = 自助法 — you take that big sample , and randomly
divide it up into a bunch of small samples and do the comparison
within those random subsets

A/A . v.s. bootstrap

使用 A/A is because if your experiment system is itself
complicated , it's actually a very good test of your system.

使用 bootstrap : push towards more complicated metrics or
running more and more features

Calculating variability empirically

Look at A/A tests on click-through-probability

Uses of A/A tests:

1. compare results to what you expect (sanity check)
2. estimate variance and calculate confidence

3. directly estimate confidence interval

Calculating variability empirically

Compare results to what you expect:

20 experiments, each on 0.5% of traffic 50 users in each group

20 more, each on 1% 100 users per group

10 more, each on 5% 500 users per group

How many experiments will show a statistically significant difference at the 95% level?

Out of 20 experiments, we expect to see 1 significant difference

Calculating variability empirically

Estimate variance and calculate confidence interval:

Since we expect a normal distribution:

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

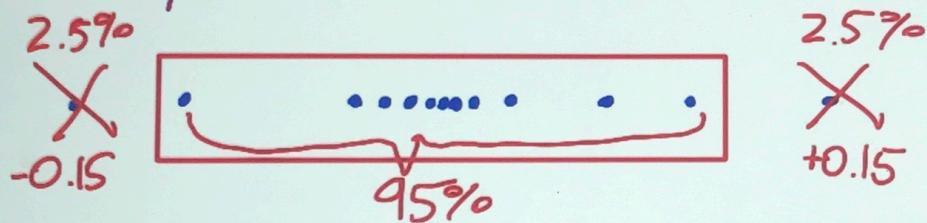
$$= 0.059 \cdot 1.96 = 0.116 \text{ empirically}$$

$$\text{Analytically: } SE = \sqrt{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}}) \left(\frac{1}{N_{\text{cont}}} + \frac{1}{N_{\text{exp}}} \right)}$$

Slightly different margin of error for each experiment

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 to 0.06

Empirical standard deviation: $0.059 * 1.65 = 0.097$

\uparrow
z-score for 90% confidence

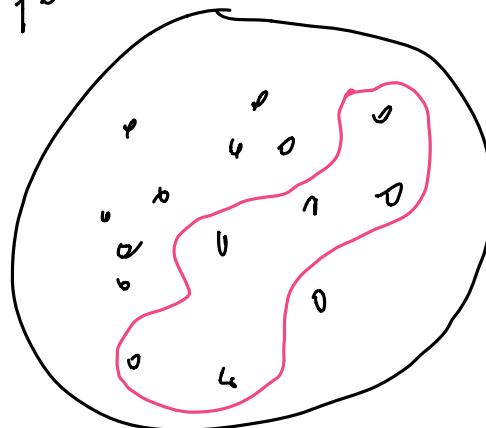
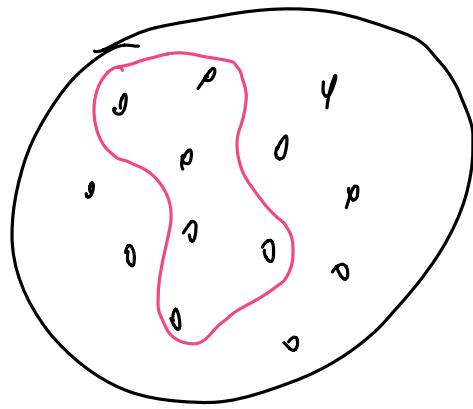
两个数据不太相似是因为我们只有 20 个 data point

Bootstrapping:

你运行一次实验，在 excel sheet 中，每个实验每一组我们只显示一个数字
一定的 CTR，实际上这些数字是由单根据很多 data point 计算得来
的。包括很多页面浏览量与点击数据。

→ take random sample from each side 计算 CTR based
on use as a "simulated experiment"

repeat to get many experiments
one experiment



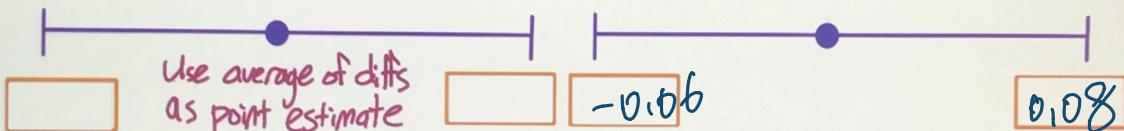
Calculating a confidence interval empirically:

- 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

Confidence level = 95%



[View Intro](#)

0.0043

[VIEW ANSWER](#)

[SUBMIT ANSWER](#)

A/A test data

$$m = 0.0364 \times 1.96 =$$

Sort the difference 小到大

第二個小值，第二個大值

17

$$95\% \times 40 = \underline{38} \text{ min 40 max } \overline{40}$$

Variability summary

1. different metric have different variability
if variability is so high for a metric, not practical to use
2. To compute the variability, we need to understand the distribution of the underlying data.
 |
 | analytical techniques
 | empirical techniques