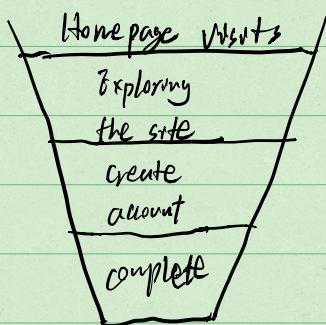


Case : Audacity that create online finance course

I. Case overview

Customer funnel



Experiment:

Change "start Now" from blue to pink

- Hypothesis: This change will increase how many students explore Audacity's course

First step to second step.

II . Refining the hypothesis (choose a metric)

1. Which metric to use ?

• λ : Total number of courses completed (Take too long)

• λ : Number of clicks (need to use percentage)

$$\checkmark: \frac{\text{Number of Clicks}}{\text{Number of page views}} = \text{Click through rate (CTR)}$$

$$\checkmark: \frac{\text{Unique visitors who click}}{\text{Unique visitors to page}} = \text{Click through prob. (CTP)} \quad (\text{Better than CTR})$$

2. Updated hypothesis:

Changing the "start Now" button from orange to pink will increase CTP of the button

3. Choose between CTR and CTP

{ **Rate:** Measure usability of a site. User can click many buttons on a page. It shows how often do they actually find that button

Prob: Measure the total impact. If just want to know how often user go to the second page, use prob. because you don't want to count if user double-click, reload, etc

4. Repeated measures for CTR

e.g. visitors = 1000, unique clicks = 100 $\Rightarrow CTR \approx 10\%$

III. Review of Statistics

1. Binomial distribution

① Two types of outcomes

② Independent events

③ Identical distribution (P same for all)

• Example of not identical:

x: drawing 20 cards from shuffled deck : red and black

{ For number : mean: np std: $\sqrt{np \cdot (1-p)}$

{ For prob : mean: p , std: $\sqrt{\frac{p(1-p)}{n}}$

• Check if binomial

V: roll a die 50 times: 6 or other

x: click on a search result page : click or not click

the next click will be influenced by previous click (find necessary information or not)

V: students complete of course after 2 months: complete or not

x: purchase of item with a week : purchased or not

related to number of items in previous purchase.

2. Confidence interval

① Benefit for Binomial distribution: have formula for sample SE and can calculate

② Meaning of 95% CI :

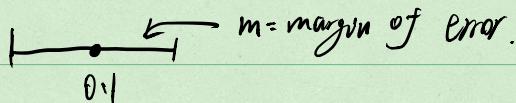
If we repeated sampling over and over again, we would expect the interval we construct around our sample mean to cover the true value in population 95% of the time.

③ calculating CI:

Center of interval: $\hat{P} = \frac{x}{N}$ # users who click
users

$$\hat{P} = \frac{100}{1000} = 0.1$$

To use normal: check $N \cdot \hat{P} > 5$ and $N \cdot (1-\hat{P}) > 5$



(\pm distribution)

Z : Z -score for $N(0,1)$

Normal approximation: $Z \approx \pm 1.96$

$Z_{0.05}$, two side: -1.96, 1.96

Binomial distribution: $Z \approx \sqrt{\frac{P \cdot (1-P)}{n}}$

$$\Rightarrow m = 0.019 \Rightarrow \begin{array}{c} | \\ 0.081 \quad 0.1 \quad 0.119 \end{array}$$

3. Hypothesis Testing

Need to calculate: PC (results due to chance)

$$H_0: P_{\text{cont}} = P_{\text{exp}} \text{ or } P_{\text{exp}} - P_{\text{cont}} = 0$$

$$H_A: P_{\text{exp}} - P_{\text{cont}} \neq 0$$

Measure P_{cont} and P_{exp} . calculate $P(C(P_{\text{exp}} - \hat{P}_{\text{cont}}) | H_0)$

compute the probability that this difference would have arisen by chance, if H_0 were true

Want to Reject Null if the prob. is small enough (< 0.05)

4. How to compare two samples: pooled standard error

① Reason for using it:

" Have two samples with different number of observations "

" Need to choose a standard error that gives us a good comparison "

of both samples \Rightarrow calculate pooled standard error

② Calculation:

$$X_{\text{cont}}, X_{\text{exp}}; N_{\text{cont}}, N_{\text{exp}}$$

$$\hat{P}_{\text{pool}} = \frac{X_{\text{cont}} + X_{\text{exp}}}{N_{\text{cont}} + N_{\text{exp}}}$$

$$SE_{\text{pool}} = \sqrt{\hat{P}_{\text{pool}} \cdot (1 - \hat{P}_{\text{pool}}) \cdot \left(\frac{1}{N_{\text{cont}}} + \frac{1}{N_{\text{exp}}}\right)}$$

$$\hat{d} = \hat{P}_{\text{exp}} - \hat{P}_{\text{cont}}, \quad H_0: d=0, \quad d \sim N(0, SE_{\text{pool}})$$

• If $\hat{d} > 1.96 \cdot SE_{\text{pool}}$ or $\hat{d} < -1.96 \cdot SE_{\text{pool}}$, we reject null

IV. Design the experiment

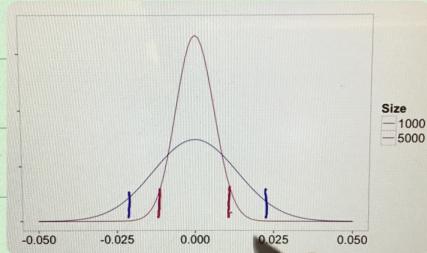
1. Size v.s. power trade-off

① How page width affect sensitivity

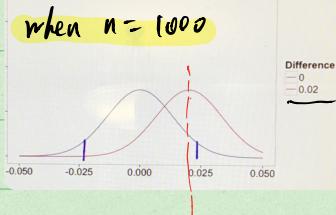
$\alpha: P(\text{reject null} \mid \text{null is true})$

$\beta: P(\text{fail to reject null} \mid \text{null false}) \quad \alpha=0.05 \Rightarrow \text{area outside boundary: } 0.05$

↑ that is, if new mean in the range, we fail to reject null



when $n=1000$



when size is small, it's already practically significant

but not statistically significant.

⇒ β is pretty high

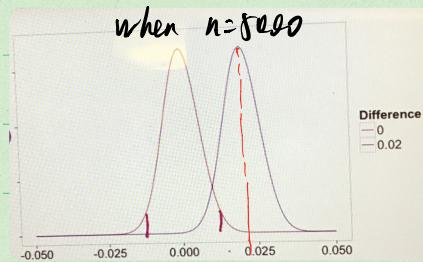
$\alpha: \text{prob to launch a bad experiment}$

$\beta: \text{prob to fail to launch a good experiment that actually did have a difference you care about.}$

• Small sample: α low, β high

(80%)

• FB: sensitivity what a high level of sensitivity at the practical significant boundary



More data:

When there is a true difference, much less likely to fall within the range of

\Rightarrow More likely to reject null when there is a difference

- Large sample: given δ same, β lower \Rightarrow power higher

2. practically / substantive significant

① Def: From business perspective, what size of change matters to us

② Reason: investment, want for next change that do better

In previous case, practically significant is 2%

③ Purpose: repeatability

- Statistical significant is about repeatability: When set up experiment, want to get the guarantee that these results are repeatable, so it's statistically significant
- Also want to see a change in experiment that you are interested in from a business standpoint so it's practically significant, it's also statistically significant.
- Need to size experiment appropriately, such that the statistically significant bar is actually lower than the practically significant bar.

3. Calculating number of page views

Need to know: use online calculator

- {
 - Base line conversion rate (c before experiment)
 - Minimum detectable effect (c practical significant level)
 - absolute or relative change (will introduce later)
- (statistical power (sensitivity)) : $1 - \beta = P(\text{reject null} \mid \text{null is false})$
- significance level : α

4. How number of page views varies changes:

① Higher CTP in control group (but less than 0.5) : Increase page views

$$SE = \sqrt{\frac{P \cdot (1-P)}{N}} \quad \sqrt{0.1 \times 0.9} = 0.3 \Rightarrow \sqrt{0.5 \times 0.5} = 0.5$$

∴ need to keep SE same, $P \uparrow$, $N \uparrow$

② Increase practical significant level (α_{min}) : Decrease page views

larger change is more easy to detect

③ Increase confidence level ($1 - \alpha$) ($\alpha \downarrow$) : Increase page views

need to be more certain that a change has occurred

before you reject the null \Rightarrow more conservative

α : $P(\text{reject null} \mid \text{null is true})$, β : $P(\text{fail to reject null} \mid \text{null is False})$

∴ when $\alpha \downarrow$, $\beta \uparrow$, then sensitivity ($1 - \beta$) \downarrow

④ $n \uparrow$, sensitivity \uparrow (\because more page view \Rightarrow narrow distribution)

⑤ Higher sensitivity : Increase page views

V : Analyze the results

1. Is the difference due to random variation:

need to calculate confidence interval

2. Calculate CZ:

$$N_{\text{cont}} = 10,072$$

$$N_{\text{exp}} = 9886$$

$$\alpha_{\min} = 0.02$$

$$X_{\text{cont}} = 974$$

$$X_{\text{exp}} = 1242$$

$$\text{Confidence level} = 95\%$$

$$\Rightarrow \hat{P}_{\text{pool}} = \frac{974 + 1242}{10,072 + 9886} = 0.111$$

$$SE_{\text{pool}} = \sqrt{0.111 \cdot (1 - 0.111) \cdot \left(\frac{1}{10,072} + \frac{1}{9886} \right)} = 0.00445$$

$$\hat{\delta} = \frac{X_{\text{exp}}}{N_{\text{exp}}} - \frac{X_{\text{cont}}}{N_{\text{cont}}} = 0.0289 \quad m = SE_{\text{pool}} \cdot 1.96 = 0.0087$$

difference in two groups

margin of error

$$\begin{array}{c} \hat{\delta} \\ \hline \end{array}$$
$$\hat{\delta} - m = 0.0202 \quad \hat{\delta} + m = 0.0376$$

$\because \hat{\delta} - m > \alpha_{\min} \therefore$ we should launch the change

3. Whether to launch an experiment

{ CZ all higher than α_{\min} : launch

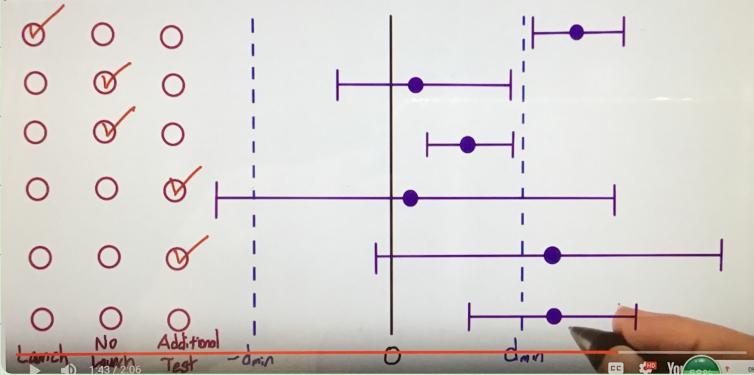
{ CZ all lower than α_{\min} : no launch

{ CZ very large, part of it is larger than α_{\min} : need additional test

(with more power)

(power = sensitivity = $1 - \beta$)

Confidence Interval Cases



4. What to do when last three cases happened

- Sometimes no time for a new A/B testing
- Need to communicate to decision-makers when they are going to have to make a judgement, and take a risk, because the data is uncertain. They're going to use other factors, like strategic business issues, or other factors besides the data.