

Experimental Design & Recommendations

Concepts in Experimental Design

- ↳ correlation ≠ causation
- ↳ Other factors may influence variables

How do we deduce causality?

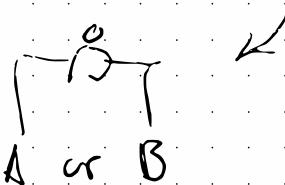
- ↳ Experiments (control confounding vars)

Experiments

- ↳ comparing results of 2+ groups
- ↳ control by random assignment of to groups
- What if this isn't possible?
 - ↳ observational studies
 - ↳ smoking, studies not possible due to ethics
 - ↳ establishing causality becomes difficult

Types of experiments

A/B tests : compare treatments vs controls (Between subj. design)



- ↳ sometimes it's possible to test of A & B on the same person

Sampling

- ↳ collecting population data is difficult
 - ↳ sample is used to est. pop. sized parameters
(select at random)

Doesn't work well if people aren't distributed uniformly

Measuring Outcomes

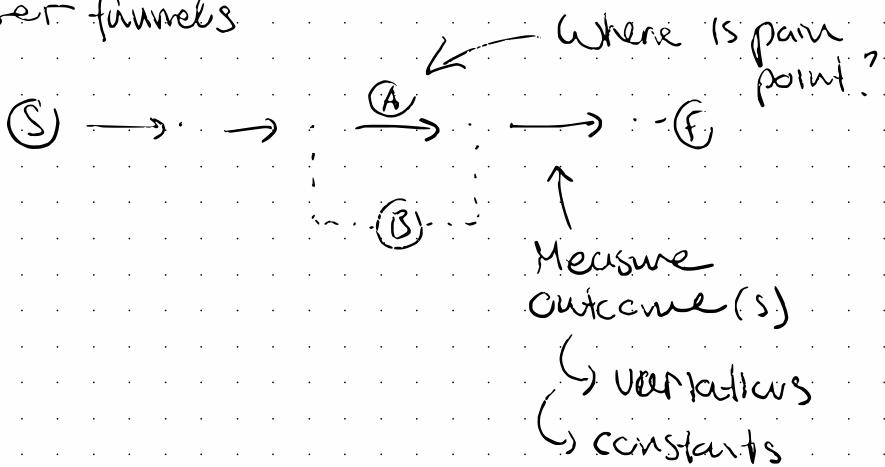
- ↳ What metrics are best to choose for experiments

depends on goals, keep in mind that your metrics can get hacked

Creating Metrics

- ↳ Think about user journey (have empathy)

- ↳ User funnels



Controlling variables

- ↳ vars that create (in)dependent

Should remain constant

- ↳ removes confounding variable

Checking validity (useful for handling data in general.)

construct

- ↳ Metrics are aligned w/ goals

internal validity

- ↳ "Degree to which an exp.'s claims of causality can be supported"

external validity

- ↳ Degree to which study can be generalised

Checking bias

Systematic errors that effect interpretation of results

↳ Sampling

↳ Methodological

↳ Novelty

↳ Order (primary, secondary)

Ethics

- ↳ Participant risk ↓
- ↳ Clear benefit(s) to experiment
- ↳ Informed consent (Difficult)
- ↳ Have strong privacy/security

S - specific

M - measurable

A - Achievable

R - relevant

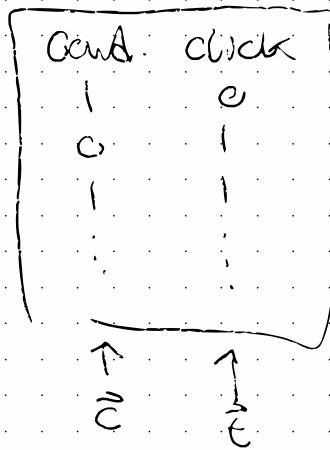
T - timely

Statistical Considerations in testing

- ↳ How much data?
- ↳ & other decisions

Start. sig.

↳ tests have tables b.ve so.



Steps,

- ① validate cond is even split
- ② Check if circles for sig diff

$$\text{obs} = \text{len}(\vec{c})$$

$$n_{\text{control}} = \text{obs} - \sum \vec{c}_i = n(0 \in \vec{c})$$

→ check control

$$\text{H}_0: p = 0.5 \quad \text{sd} | p = 0.5 = \sqrt{p(1-p) * n}$$

$$\text{H}_1: p \neq 0.5$$

$$z = \frac{n_{\text{control}} + 5 - (p_{\text{H0}} * n)}{\text{sd}}$$

observed 1
correction norm. const.
 expect val

then... compare z w/ $1.96 * \text{norm.cdf}(z)$

p-value

p-value is prob. of data given parameter
ie prob of seeing data if null is true

↳ If evenly dist, p val should not be
significantly different.

→ Check click for sig diff.

for each c in \vec{C} , there is p of click

$$\hookrightarrow P(c_1) - P(c_0) = P_{\text{diff}}$$

null hypothesis is overall avg

$$H_0: P(c_1) = P(c_0)$$

$$H_1: P(c_1) \neq P(c_0)$$

$$\text{std err} = \sqrt{P_{\text{null}}(1-P_{\text{null}}) * \left(\frac{1}{n_{\text{cont}}} + \frac{1}{n_{\text{exp}}}\right)}$$

$$z = \frac{P_{\text{diff}}}{\text{std err}} \Rightarrow 1 - \text{norm.cdf}(z)$$

p value

Note, to get p value you need to establish H_0 & H_1

- ↳ establish in terms of a diff (Δ)
- ↳ normalise by variance or std err
- ↳ find where H_1 lies on a distribution depending on tests.

simulation

$\theta \in \{\text{exp, cont}\}$

- ↳ draw from Binomial(n_θ , p_{null} , many times)

$$\text{drifts} = \frac{\text{exp circles}}{n_{\text{exp}}} - \frac{\text{ctrl circles}}{n_{\text{ctrl}}}$$

↳ list of normalised drifts

if $[P_{\text{diff}} < \text{diff}; \text{then false else true}]$.mean()

↳ gives p-value

Practical Significance

↳ even if difference is significant,
it doesn't mean that it should be
implemented

Experiment sizes

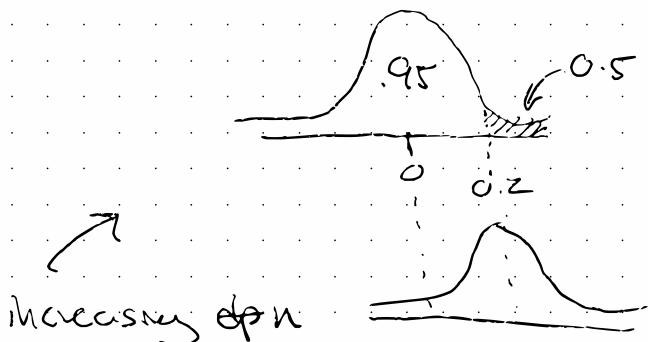
Statistical Power: given true mean,
the prob. of rej H_0

i.e.

$$\alpha = 10\% \text{ curr}$$

$$\beta = 12\% \text{ curr (desired)}$$

$$1 - \beta = 28$$



will narrow curves & increase power of test

find n people we will need for experiment

Note, power = $1 - \beta$ rate of type 2 error
 $P(\text{Rej } H_0 | H_0 \text{ is true})$

$$\delta e_{\text{null}} = \Phi_{\text{null}}(1 - \text{Power})$$

Dummy tests

↳ comparison between same exp. group

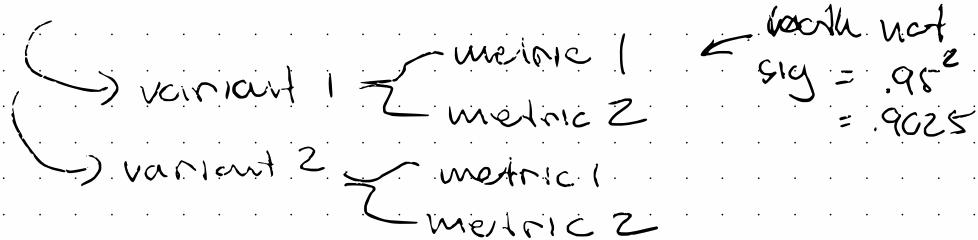
A/A test

↳ Make sure tests is working

↳ get more info about metrics

Analyzing Multiple Metrics

↳ multiple metrics \rightarrow increases chance
for false pos



\Rightarrow Apply bonferroni correction to account for
this

Early Stopping

- ↳ Stopping & viewing test results can lead to more false tues
- ↳ need to google how to do early stopping if necessary

Introduction to recommendation engines

↳ used in many ways on the web these days.

knowledge based recommendations

↳ good recommendations?

↳ how about top song/movie/etc.

↳ most popular
in genre?

knowledge based

↳ recommends top canid for
certain sentence

How is this done?

↳ create df sorted by "top" criteria
(highest avg, most rated, etc.)

↳ filter as necessary

Advantage of knowledge based?

- ↳ don't know a lot about users
- (X) not personalized to the user

Collaborative filtering

↳ "people who $x \times$ also $y \cdot z$ "

→ two types

model based

Neighborhood based

N based collaborative filtering

Similarity vs distance

→ Correlation
coeff.

→ distance
measurements

Movie, Person matrix

	8	9	9	9	People
8	□	8	9		
1		1	5	7	
6	□	6			
3	□	3	5		

for recommendation
to create person
remove:
- seen movies
- movies w/ low rat.

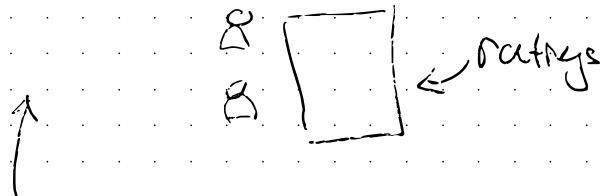
↳ creates sparse
matrix

- ↳ for each user, calc list of movies seen
- ↳ recommends differently depending on number of movies seen.

↖ ↘
 $\leftarrow 2 \quad \rightarrow 7=2$

calc corr. between
users

similar movies



You can even use
different similarity
metrics

content based recommendations

↳ "cold start problem"

What if the number of users you have is small?

→ find similar items by name,
appearance,
item type
etc.

↳ think of people w/ less than 9 ratings for a movie

get users & their top rated movies

↳ if movies x attribute matrix is A
then take $A \cdot A^T$ to get a dot prod matrix for movies

↳ $A \cdot A^T$ for which attributes are similar

→ for a top rated movie, find the most similar. (which haven't been seen already)

Types of ratings

- likes vs dislikes
- higher
- 3, 5, 10 star ratings
- granularity

→ include neutral? depends

Goals of recommendations

- ↳ relevance
 - ↳ novelty
 - ↳ serendipity
- exposing users to something new

Matrix Factorisation for Recommenders

How do we know if recommendations are good?

How to use ML for recommendations

How to recommend to new users

Quality of recommendations

↳ watch metrics over time

↳ AB test

↳ train/test cross validation

→ you can also do older vs newer data

SVD

$$A = U \Sigma V^T$$

singular
values

↑ orthonormal matrices
with the basis of the row/col spaces.

Why use SVD? (Latent Factors)

↳ Factors that aren't in the data
but are hidden underneath.

- think how PCA also reveals stuff
like this.

Item

$$A = \text{users} \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = U \Sigma V^T$$

$U \leftarrow$ orthonormal basis for
 $N \times k$ user directions of most
variance

$V^T \leftarrow$ On basis for
 $k \times M$ item dir of most var.

$\Sigma \leftarrow$ connects U w/ V^T w/ strengths
 $k \times k$ (latent factors) (Diagonal Matrix)

you can take a subset of PCs in order to recreate a full picture
^
estimate

$\Sigma \leftarrow$ will contain the variance found in each PC

$$\hookrightarrow \text{amount of var explained} = \frac{\sum \text{exp } P_{ci}^2}{\sum \text{all } P_{ci}}$$

Problem: SVD break down w/
missing values!

Funle SVD

\hookrightarrow possible to do SVD w/ missing vals

- ① fill V^T & U w/ random numbers
- ② choose non missing var ($\Sigma_{i,j}$)
- ③ get pred b w/ $U_{i,*} \cdot V_{*,j}^T$
- ④ get sq. error

\hookrightarrow minimize w/ AD

$$\frac{\delta}{\delta u_i} (y - uv)^2 = -2(y - uv)v_i$$

$$\frac{\delta}{\delta v_i} (y - uv)^2 = -2(y - uv)u_i$$

example

$$y_{ij} = 9$$

$$u_i = [0.8, 1.2, -2]$$

$$v_i = [-1.2, 1.1, -0.2]$$

$$u_i \cdot v_i^T = 4$$

$$\therefore \text{error} = (9 - 4)^2$$

$$= \dots$$

update u_i, v_i

$$u_{i\text{new}} = u_i + \alpha \text{error} v_i$$

$$v_{i\text{new}} = v_i + \alpha \text{error} u_i$$

↳ rinse & repeat till convergence

Some quick notes on AV w/ dd & new vals

1. take dataset
2. order by date
3. split data