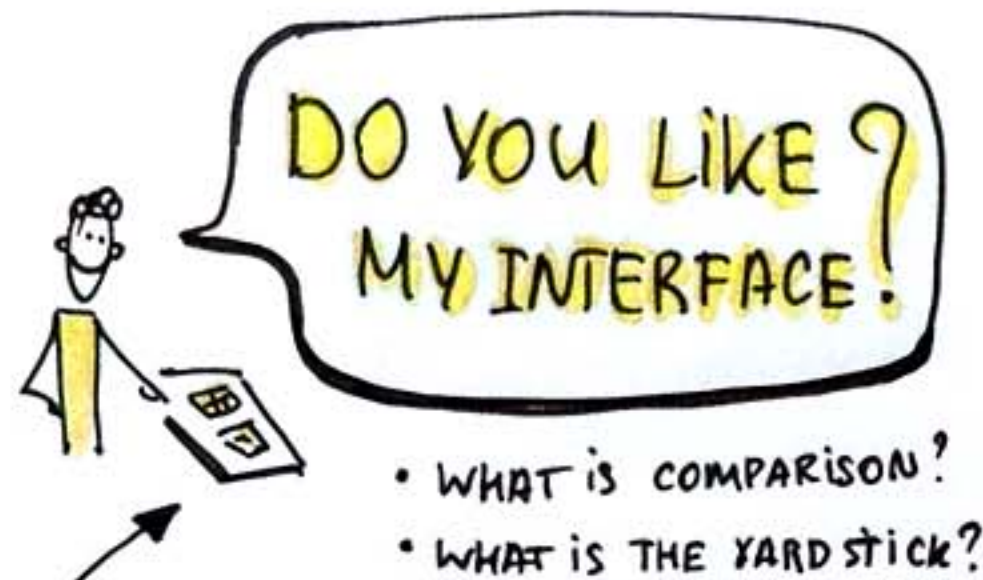


7.1 DESIGNING STUDIES



HOW TO IMPROVE

- ▣ **BASERATES**
how often does Y occur?
- ▣ **CORRELATIONS**
do X and Y co-vary?
- ▣ **CAUSES**
does X cause Y?

SKETCHNOTESPACE.COM
ANNA IURCHENKO



MANIPULATION

- diff conditions
- independent variables



MEASURES

- dependent variables
- accuracy
- recall
- emotional response



PRECISION

- internal validity
- if you ran this again will you see the same results
- # of people



GENERALIZABILITY

- external validity
- does this apply to this particular users?

“CONTROLLED COMPARISON ENABLES CASUAL INFERENCE”



FIDELITY OF IMPL. APPROACH

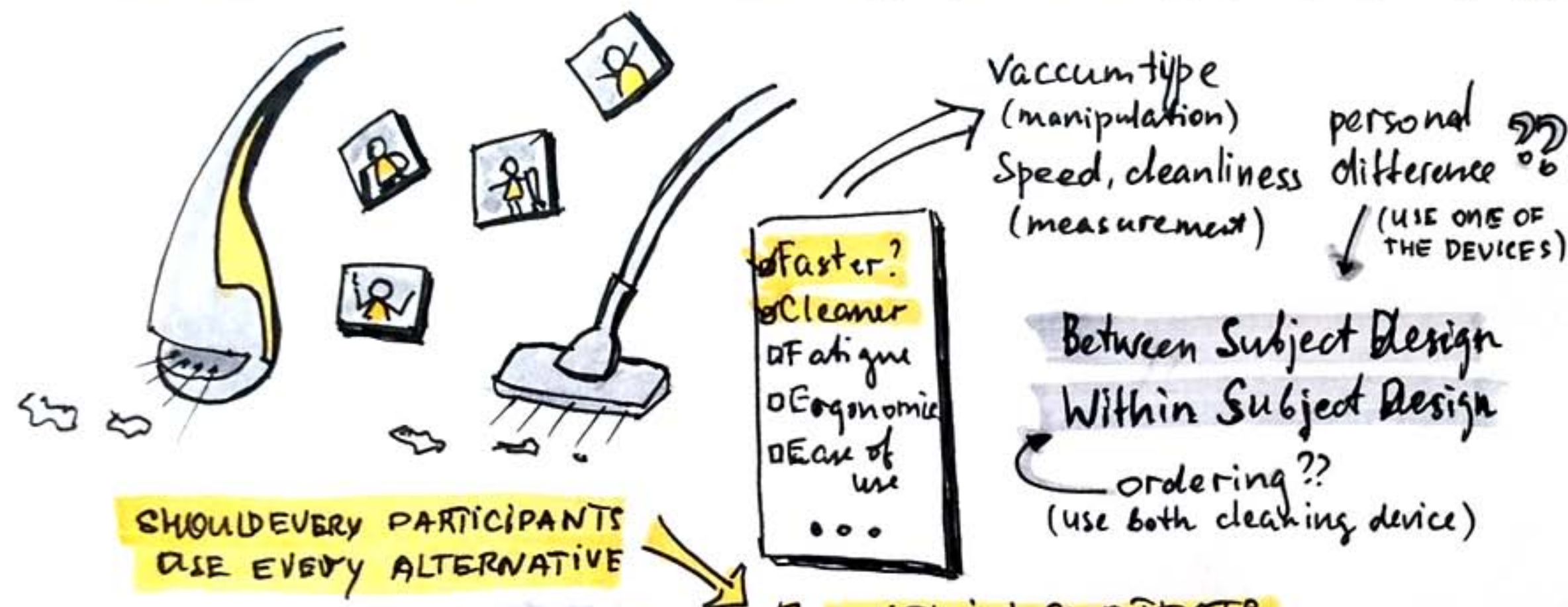


STRATEGIES FOR FAIRER COMPARISONS

- insert your new approach into production setting
 - scale things down so you're looking at a piece of a larger system
 - train people up
 - make a version of a production thing in the same style as new approach
- ← you can learn stuff 😊

7.2 ASSIGNING CONDITIONS

SKETCHNOTESPACE.COM
ANNA IURCHENKO



THREE or MORE ALTERNATIVES

1	2	3
2	3	1
3	1	2

LATIN SQUARE

- each person will use all three conditions, order is changed

ASSIGNMENT SHOULD BE RANDOM (HAWTHORNE EFFECT)



1. WITHIN SUBJECTS

- (everyone tries all the options)
- not worry about learning

2. BETWEEN-SUBJECTS

- (each person tries one)
- more people, attention to fair assign.

3. COUNTERBALANCING

- minimize variation in a between subject design



MAKE CLEAR GOALS

- narrow & scope to the purpose of your study
- WHO, WHERE, WHEN
- SCENARIO (REALISTIC)

Question

Data to collect

Setup

7.3 IN-PERSON STUDIES

- planning
- EXECUTING
- ANALYSING

CONCRETE TASKS

- WRITE THEM DOWN
- THINK OF ORDER
- WHAT TO DO IF USER CAN'T ACCOMPLISH?
- TRAINING?

REMINDE USERS
You are testing the site, not the users

PILOT



COLLEAGUES

REAL USERS (ONE)

CAPTURING

THINKING ALOUD

GREETINGP.

DEBRIEF

• WHAT YOUR CONS

COUNTERBALANCE ASSIGNMENT

- USE PRE-TEST
- EACH PARTICIPANT HAS AN EQUAL CHANCE OF LANDING IN EITHER CONDITIONS

DAVID MARTIN
DOING PSYCHOLOGY EXPERIMENTS

A 35 59
B 32 57
①
(TYPERS EXAMPLE)

REGRESSION IS A DANGER

7.4 RUNNING WEB EXPERIMENTS

ROLL OUT DIFFERENT VERSIONS OF A USER INTERFACE, GET FEEDBACK, ITERATE QUICKLY

DUSTIN CURTIS
[dcurt.is]
BLOG

TYPOGRAPHY EXPERIMENT

Ron Kohavi

color change on MSN search

↑ 3% (AD)
CONVERSION RATE?



SMALL DISTRACTIONS (extra fields) CAN YIELD BIG CHANGES

COMMITMENT ESCALATION

★★★★★ → [] WHY?

BY COMBINING
ITERATIVE DESIGN
AND CONTROLLED
EXPERIM.

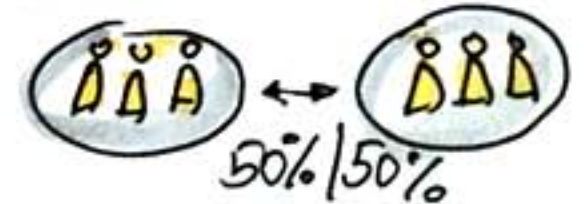
YOU CAN
DIAL IN THESE
PHENOMENA
TO ↑ EFFECT

PRINCIPLES

0.1%

RAMP-UP & AUTO ABORT
(simple analyses to find egregious problem)

EQUAL AMOUNT
OF PEOPLE



RUN IT LONG
ENOUGH



1. CONSISTENT
2. DURABLE
3. INDEPEND.

rules of random assignments

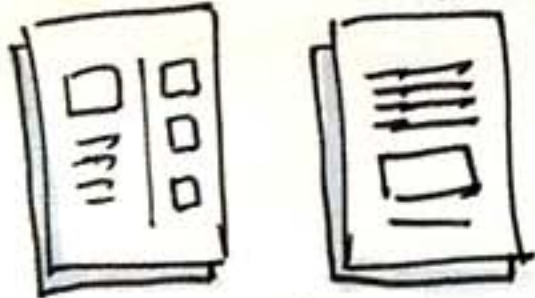
DESIGNER
ROLE SHIFTS TO
BEING ABOUT CREATING
MULTIPLE ALTERNATIVES

MEASURE
THINGS THAT
MATTER

WAYS DESIGN MAKES A DIFFEREN.

- position & color of call to action
- whitespace?
- position of testimonials
- position of heading
- number of columns
- # of visual element competing for attract.
- the age, sex, appearance of people on the photo

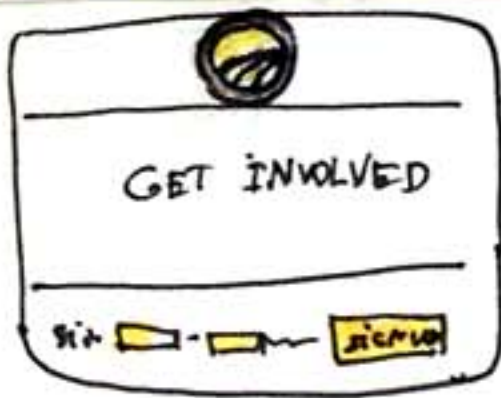
AB



METRICS? CTR?

MULTIVARIATE
TEST (variations
of two diff. parts
of the page)

CALL TO ACTION EXPERIM



baseline 8.2%

- Variations
- 7.5% button (sign up, learn more...)
 - 8.91% MEDIA (family img, change img, obama img...)



Small changes, HUGE IMPACT
(text on the button)

SKETCHNOTESPACE.COM
ANNA IURCHENKO

7.5 ANALYZING EXPERIMENTS



COMPARING RATES

ANALYZE DATA IN 3 STEPS

- 1 HOW IT LOOKS LIKE? (plot your data)
- 2 OVERALL NUMBERS (deviation, average)
- 3 IS THE DIFF. 'REAL'? (compute significance)

PEARSON'S CHI-SQUARED TEST

compare the rates of an expected value to an observed value

$$\chi^2 = \frac{(\text{OBSERVED} - \text{EXPECTED})^2}{\text{EXPECTED}}$$

Sign \uparrow p

IMPROVED CLICK-THROUGHS

[example]

- 10% CTR
- CHANGE TEXT
- 1 week, 119 clicked "Learn More" out of 1000

Learn More

YES

CAN WE SAY WITH CONFIDENCE THAT "Learn More" have \uparrow CTR

$$\frac{(119 - 100)^2}{100} + \frac{(881 - 900)^2}{900} = 4.01$$

df = 1 (degree of freedom) $p \leq 0.05$



STATISTICAL TESTING

formalise "we're pretty sure"

helps generalize from small samples



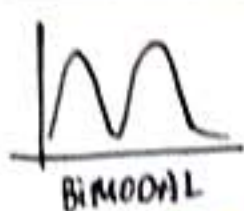
1908, WILLIAM GOSSET, 'Guinness'

T-test (small sample) 2 conditions

ANOVA

(compare > 2 conditions)

Data often ain't 'Normal'



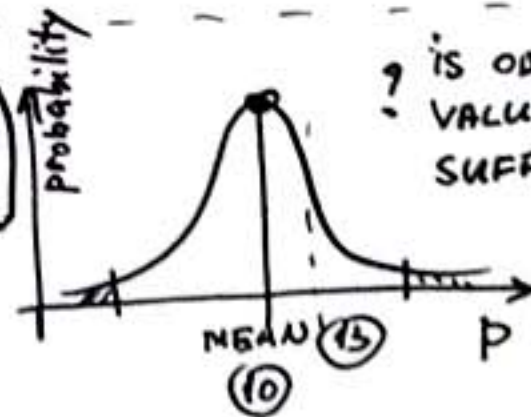
BIMODAL



- A/A tests
- randomised testing



SKETCHNOTESPACE.COM
ANNA TURCHENKO



? IS OBSERVED VALUE IS WEIRD SUFFICIENTLY

NORMAL VARIANCE

THE NULL HYPOTHESIS

• our opening bid in any stat. test is that there is no relation between measured phenomena



- Practical statistic for HCI, Jacob Wobbrock, dept.s.washington.edu
- Doing Psychology Experiments, David Martin
- Statistics as Principled Argument, Robert P. Abelson
- Learning to use stat. test, Judith Green