CSI 60: User Interfa	ce Design
Quantitative Evaluation	03/10/10
	Berkelev
	UNIVERSITY OF CALIFORNIA

Last Chances...

- 1. Check out a camera and tripod after class
- 2. Bring official DSP letter for special midterm accommodations to us

Charles Thacker wins Turing Award!

"Thacker created and collaborated on what would become the fundamental building blocks of the PC business. The **Alto** computer, developed in 1974, incorporated bitmap (TV-like) displays which enable modern **graphical user interfaces** (GUIs), including **What You See Is What You Get** (WYSIWYG) **editors**. These components have dominated computing during the last two decades. "



//en.wikipedia.org/wiki/File:Xerox_Alto.j

Topics

- Managing study participants (qualitative and quantitative studies)
- 2. Why do we conduct quantitative studies?
- 3. Designing controlled experiments



<section-header>

The Three Belmont Principles

Respect for Persons

Have a meaningful consent process: give information, and let prospective subjects freely chose to participate

Beneficience

Minimize the risk of harm to subjects, maximize potential benefits Justice

Use fair procedures to select subjects (balance burdens & benefits)

To ensure adherence to principles, most schools require Institutional Review Board approval of research involving human subjects.

Ethics: Stanford Prison Experiment

1971 Experiment by Phil Zimbardo at Stanford 24 Participants – half prisoners, half guards (\$15 a day) Basement of Stanford Psychology bldg turned into mock prison Guards given batons, military style uniform, mirror glasses,... Prisoners wore smocks (no underwear), thong sandals, pantyhose caps

Experiment quickly got out of hand

Prisoners suffered and accepted sadistic treatment Prison became unsanitary/inhospitable Prisoner riot put down with use of fire extinguishers Guards volunteered to work extra hours

Zimbardo terminated experiment early

Grad student Christina Maslach objected to experiment Important to check protocol with ethics review boards [from Wikipedia]



Ethics Vas it useful? ...,that's the most valuable kind of information that you can have - and that certainly a society needs it" (Zimbardo): Vas it ethical? Could we have gathered this knowledge by other means?

http://www.prisonexp.org/slide-42.htm

Ethics (more recently)

"In 2001, a faculty member from the business school of a major university designed a study to see how restaurants would respond to complaints from putative customers. As part of the project, the researcher sent letters to restaurants falsely claiming that he and/or his wife had suffered food poisoning that ruined their anniversary celebration. The letters disclaimed any intention of contacting regulatory agencies and stated that the only intent was to convey to the owner what had occurred "in anticipation that you will respond accordingly." Restaurant owners were understandably upset and some employees lost their jobs before it was revealed that the letter was a hoax."

CITI Human Subject Training Material

Beneficience: Example

MERL DiamondTouch: User capacitively coupled to table through seating pad. No danger for normal users, but possibly increased risk for participants with pacemakers.

Inform subjects in consent!



//www.merl.com/projects/images/DiamondTouch.jpg

Privacy and Confidentiality

Privacy: having control over the extent, timing, and circumstances of sharing oneself with others.

Confidentiality: the treatment of information that an individual has disclosed with the expectation that it will not be divulged

Examples where privacy could be violated or confidentiality may be breached in HCI studies?

Treating Subjects With Respect

Follow human subject protocols

Individual test results will be kept confidential Users can stop the test at any time Users are aware (and understand) the monitoring technique(s) Their performance will not have implications on their life Records will be made anonymous

Use standard informed consent form

Especially for quantitative tests Be aware of legal requirements

Conducting the Experiment

Before the experiment

Have them read and sign the consent form Explain the goal of the experiment in a way accessible to users Be careful about the demand characteristic (Participants biased towards experimenter's hypothesis) Answer questions

During the experiment Stay neutral Never indicate displeasure with users performance

After the experiment

Debrief users (Inform users about the goal of the experiment) Answer any questions they have

Managing Subjects

Don't waste users time

Use pilot tests to debug experiments, questionnaires, etc... Have everything ready before users show up

Make users comfortable

Keep a relaxed atmosphere Allow for breaks Pace tasks correctly Stop the test if it becomes too unpleasant

If you want to learn more...

Online human subjects certification courses: E.g., http://phrp.nihtraining.com/users/login.php

The Belmont Report: Ethical Principles and Guidelines for the protection of human subjects of research

1979 Government report that describes the basic ethical principles that should underly the conduct of research involving human subjects

http://ohsr.od.nih.gov/guidelines/belmont.html

Why Quantitative Studies?

Qualitative Studies

Qualitative: What we've been doing so far

Contextual Inquiry: try to understand user's tasks and conceptual model Usability Studies: look for critical incidents in interface

Qualitative methods help us:

Understand what is going on Look for problems Roughly evaluate usability of interface

Quantitative Studies

Ouantitative Use to reliably measure some aspect of interface Compare two or more designs on a measurable aspect

Approaches Collect and analyze user events that occur in natural use mouse clicks, key presses Controlled experiments

Examples of measures Time to complete a task Average number of errors on a task Users' ratings of an interface * Ease of use, elegance, performance, robustness, speed,...

 \ast You could argue that users' perception of speed, error rates etc is more important than their actual values

Comparison

Qualitative studies

Faster, less expensive \rightarrow esp. useful in early stages of design cycle In real-world design, quantitative study not always necessary

Quantitative studies Reliable, repeatable result \rightarrow scientific method

Best studies produce generalizable results

Pilot User Study Assignment (after midterm)

You will conduct a **qualitative** study We don't have enough time or subjects for quantitative studies But you should do a little quantitative analysis What are your measures? Compute summary statistics (mean, stdev) Do you have independent, dependent, and control variables?

Designing Controlled Experiments

Steps in Designing an Experiment

- I. State a lucid, testable hypothesis
- 2. Identify variables (independent, dependent, control, random)
- 3. Design the experimental protocol
- 4. Choose user population
- 5. Apply for human subjects protocol review
- 6. Run pilot studies
- 7. Run the experiment
- 8. Perform statistical analysis
- 9. Draw conclusions







Experiment Design

- Testable hypothesis Precise statement of expected outcome
- Independent variables (factors) Attributes we manipulate/vary in each condition Levels – values for independent variables

Dependent variables (response variables) Outcome of experiment (measurements)

Usually measure user performance

Experiment Design

Control variables

Attributes that will be fixed throughout experiment Confound – attribute that varied and was not accounted for Problem: Confound rather than IV could have caused change in DVs Confounds make it difficult/impossible to draw conclusions

Random variables

Attributes that are randomly sampled Increases generalizability

Variable Types

Nominal: categories with labels, no order

Ordinal: categories with rank order

Continuous: interval (w/o zero point), ratio (w/ zero point)

Common Metrics in HCI

Performance metrics:

- Task success (binary or multi-level)
- Task completion time
- Errors (slips, mistakes) per task
- Efficiency (cognitive & physical effort)
- Learnability

Satisfaction metrics:

• Self-report on ease of use, frustration, etc.





Satisfaction Metric: Likert Scales

Respondents rate their level of agreement to a statement

I: Strongly Disagree 2: Disagree 3: Neither agree nor disagree 4:Agree 5: Strongly agree

"Overall, I am satisfied with the ease of completing the tasks in this scenario"

Likert data is ordinal, not continuous (matters for analysis)!

Variables	
Independent variables	
Dependent variables	
Control variables	Grouped Bar
Random variables	Divided Bar

Variables	
Independent variables Chart type Leaf Node vs Non-Leaf Node Comparison Data Density (# of Leaf Nodes) Dependent variables Response Time Estimation Error Leas Set for the set	TreeMap
Control variables Color scheme, rendering style	
Random variables Location, environment, Attributes of subjects Age, sex,	Grouped Bar



Experimental Protocol

What is the task? (must reflect hypothesis!) What are all the combinations of conditions? How often to repeat each combination of conditions? Between subjects or within subjects Avoid bias (instructions, ordering, ...)

Number of Conditions

Consider all combinations to isolate effects of each IV (factorial design) (3 chart types) * (3 leaf/non-leaf combinations) * (3 densities) = 27 combinations) * (3 densities) = 27

Adding levels or factors can yield lots of combinations!

Reducing Num. of Conditions

Vary only one independent variable leaving others fixed

Problem: ?

Reducing Num. of Conditions

Vary only one independent variable leaving others fixed

Problem: Will miss effects of interactions

Other Reduction Strategies

Run a few independent variables at a time If strong effect, include variable in future studies Otherwise pick fixed control value for it

Fractional factorial design

Procedures for choosing subset of independent variables to vary in each experiment

Choosing Subjects

Pick balanced sample reflecting intended user population Novices, experts Age group Sex

....

Example 12 non-colorblind right-handed adults (male & female)

Population group can also be an IV or a controlled variable What is the disadvantage of making population a controlled var?

Between Subjects Design

Wilma and Betty use one interface

Dino and Fred use the other





Between vs. Within Subjects

Between subjects

- Each participant uses one condition +/- Participants cannot compare conditions
- + Can collect more data for a given condition
- Need more participants

Within subjects

All participants try all conditions

- + Compare one person across conditions to isolate effects of individual diffs + Requires fewer participants
- Fatigue effects
- Bias due to ordering/learning effects

Within Subjects: Ordering Effects

In within-subjects designs ordering of conditions is a variable that can confound results Why?

Turn it into a random variable

...

Randomize order of conditions across subjects Counterbalancing (ensure all orderings are covered) Latin square (partial counterbalancing)

Run the Experiment

Always pilot it first! Reveals unexpected problems Can't change experiment design after starting it

Always follow same steps - use a checklist

Get consent from subjects

Debrief subjects afterwards

Results: Statistical Analysis

Descriptive Statistics

Continuous data: Central tendency (mean, median, mode), Dispersion, Shape of distribution Categorical data: Frequency distributions









Are the Results Meaningful?

Hypothesis testing Hypothesis: Manipulation of IV effects DV in some way Null hypothesis: Manipulation of IV has no effect on DV Null hypothesis assumed true unless statistics allow us to reject it

Statistical significance (p value)

Statistical tests

T-test (I factor, 2 levels) Correlation ANOVA (1 factor, > 2 levels, multiple factors) MANOVA (> 1 dependent variable)





T-test

Compare means of 2 groups Null hypothesis: No difference between means

Assumptions

Samples are normally distributed Very robust in practice Population variances are equal (between subjects tests) Reasonably robust for differing variances Individual observations in samples are independent Extremely important!







15



Are we done with our analysis? No! Multiple IVs effect DV non-additively! We had 3 IVs (chart type, density, node/leaf combo), so we should investigate interaction effects!







Draw Conclusions

What is the scope of the finding? Does the experiment reflect real use? External validity Are there other parameters at play? Internal validity

Summary

Quantitative evaluations Repeatable, reliable evaluation of interface elements To control properly, usually limited to low-level issues Menu selection method A faster than method B

Pros/Cons

Objective measurements Good internal validity → repeatability But, real-world implications may be difficult to foresee Significant results doesn't imply real-world importance 3.05s versus 3.00s for menu selection

NextTime

Midterm review! No new readings – revisit old readings. Heuristic Evaluation & Low-Fi Prototypes due!