

Autumn 2019

Prof. James A. Landay

Stanford University

dt+UX DESIGN THINKING FOR USER EXPERIENCE DESIGN + PROTOTYPING + EVALUATION

Usability Testing

刘哲明
 Prof. James A. Landay
 Computer Science Department
 Stanford University

Autumn 2019
 November 13, 2019

1

Hall of Fame or Shame?

- Kitchen Stories

2

Hall of Fame!

- Kitchen Stories
- Like
 - Large pictures of recipes
 - Photos & videos
 - Shopping list that marks off as you purchase
- Wish
 - ?

3

dt+UX DESIGN THINKING FOR USER EXPERIENCE DESIGN + PROTOTYPING + EVALUATION

Usability Testing

刘哲明
 Prof. James A. Landay
 Computer Science Department
 Stanford University

Autumn 2019
 November 13, 2019

4

Outline

- Why do usability testing?
- Choosing participants
- Ethical considerations
- Designing & conducting the test
- Using the results
- Experimental options & details

5


Why do Usability Testing?

- Can't tell how good UI is until?
 - people use it!
- Expert review methods are based on evaluators who may?
 - know too much
 - not know enough (about tasks, etc.)
- Hard to predict what real users will do

6

Choosing Participants



- Representative of target users. How so?
 - job-specific vocab / knowledge
 - tasks
- Approximate if needed
 - system intended for doctors?
 - get **medical students** or nurses
 - system intended for engineers?
 - get engineering students
- Use incentives to get participants
 - t-shirt, mug, free coffee/pizza



7

Ethical Considerations

- Usability tests can be distressing
 - users have left in tears
- Testing/fieldwork can be coercive if there is a power imbalance (e.g., in under resourced communities)





People may feel no option but to speak to you or give you their time even though they may not get anything of value in return.

8

Ethical Considerations


- You have a responsibility to alleviate these issues
 - make voluntary with informed consent (form)
 - avoid pressure to participate
 - let them know they can stop at any time
 - stress that you are testing the system, not them
 - make collected data as anonymous as possible
- Often must get human subjects approval (IRB)



9

Usability Test Proposal


- A report that contains
 - objective
 - description of system being testing
 - task environment & materials
 - participants
 - methodology
 - tasks
 - test measures
- Get approved & then reuse for final report
- Seems tedious, but writing this will help “debug” your test



10

Selecting Tasks


- Tasks from low-fi design can be used
 - may need to shorten if
 - they take too long
 - require background that test user won't have
- Don't train unless that will occur in real deployment
- Avoid bending tasks in direction of what your design best supports
- Don't choose tasks that are too fragmented?
 - fragmented = do not represent a complete goal someone would try to accomplish with your application
 - e.g., phone-in bank test



11

Two Types of Data to Collect

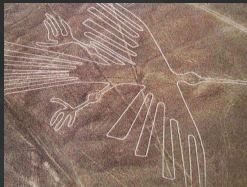
- Process data
 - observations of what users are doing & thinking
 - **qualitative**
- Bottom-line data
 - summary of what happened
 - time, errors, success
 - i.e., the dependent variables
 - **quantitative**




12

Which Type of Data to Collect?

- Focus on process data first
 - gives good overview of where problems are




http://www.redicecreations.com/id_img24592nazca_bird.jpg

13

Which Type of Data to Collect?


- Focus on process data first
 - gives good overview of where problems are
- Bottom-line doesn't tell you ?
 - where to fix
 - just says: "too slow", "too many errors", etc.
- Hard to get reliable bottom-line results
 - need many users for statistical significance



14

The "Thinking Aloud" Method


- Need to know what users are thinking, not just what they are doing
- Ask users to talk while performing tasks
 - tell us *what they are thinking*
 - tell us *what they are trying to do*
 - tell us *questions that arise as they work*
 - tell us *things they read*



15

Thinking Aloud (cont.)

- Prompt the user to keep talking
 - "tell me what you are thinking"
- Only help on things you have pre-decided
 - keep track of anything you do give help on
- Make a *recording* & take good notes
 - make sure you can tell what they were doing
 - use a digital watch/clock
 - record audio & video
 - or even event logs

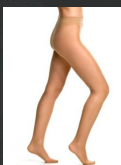


16

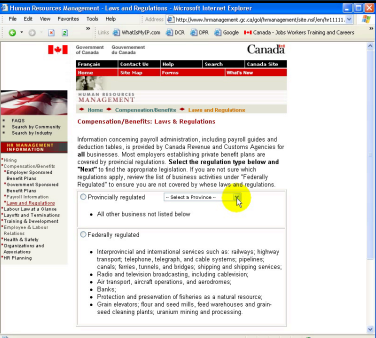
Will thinking out loud give the right answers?

- Not always
- If you ask, people will always give an answer, even if it has nothing to do with facts
 - panty hose example

→ Try to avoid specific questions (especially that have binary answers)



17



18

<http://bit.ly/cs147-quiz3-20-13>

Closed notes & no web lookup


5 minutes

Do not communicate about this quiz with anyone inside or outside this room

19

Using the Test Results


- Summarize the data
 - make a list of all critical incidents (CI)
 - positive & negative
 - include references back to original data
 - try to judge why each difficulty occurred
- What does data tell you?
 - UI work the way you thought it would?
 - users take approaches you expected?
 - something missing?



20

Using the Results (cont.)


- Update tasks & rethink design
 - rate severity & ease of fixing CIs
 - fix both severe problems & make the easy fixes



21

Measuring Bottom-Line Usability


- Situations in which numbers are useful
 - time requirements for task completion
 - successful task completion %
 - compare two designs on speed or # of errors
- Ease of measurement
 - time is easy to record
 - error or successful completion is harder
 - define in advance what these mean
- Do not combine with thinking-aloud. Why?
 - talking can affect speed & accuracy



22

Analyzing the Numbers

- Example: trying to get task time ≤ 30 min.
 - test gives: 40, 5, 20, 90, 10, 15
 - mean (average) = 30
 - median (middle) = 17.5
 - looks good!
- Did we achieve our goal?
- Wrong answer, not certain of anything!
- Factors contributing to our uncertainty?
 - small number of test users ($n = 6$)
 - results are very variable (standard deviation = 32)
 - std. dev. measures dispersal from the mean



23

Analyzing the Numbers (cont.)

- This is what basic statistics can be used for
- Crank through the procedures and you find
 - 95% certain that typical value is between 5 & 55

24

Analyzing the Numbers (cont.)

| Web Usability Test Results | | |
|---|----------------------------|--------------------------------------|
| Participant # | Time (minutes) | |
| 1 | 20 | |
| 2 | 15 | |
| 3 | 40 | |
| 4 | 90 | |
| 5 | 10 | |
| 6 | 5 | |
| | | |
| number of participants | 6 | |
| mean | 30.0 | |
| median | 17.5 | |
| std dev | 31.8 | |
| | | |
| standard error of the mean | = stddev / sqrt (#samples) | 13.0 |
| | | |
| typical values will be mean +/- 2*standard error | | --> 4 to 56! |
| | | |
| what is plausible? = confidence (alpha=5%, stddev, sample size) | 25.4 | --> 95% confident between 4.6 & 55.4 |

25

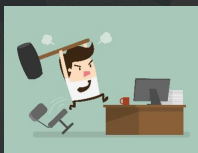
Analyzing the Numbers (cont.)

- This is what basic statistics can be used for
- Crank through the procedures and you find
 - 95% certain that typical value is between 5 & 55
- Usability test data is *highly variable*
 - need lots to get good estimates of typical values
 - 4x as many tests will only narrow range by 2x
 - breadth of range depends on sqrt of # of test users
 - this is when online methods become useful
 - easy to test w/ large numbers of users

26

Measuring User Preference

- How much users like or dislike the system
 - can ask them to rate on a scale of 1 to 10
 - or have them choose among statements
 - "best UI I've ever...", "better than average"...
 - hard to be sure what data will mean
 - novelty of UI, unrealistic setting ...
- If many give you low ratings → trouble
- Can get some useful data by asking
 - what they liked, disliked, where they had trouble, best part, worst part, etc.
 - redundant questions are OK



27

Comparing Two Alternatives

- *Between groups* experiment
 - two groups of test users
 - each group uses only 1 of the systems
- *Within groups* experiment
 - one group of test users
 - each person uses both systems (cheaper)
 - can't use the same tasks or order (learning)
 - best for low-level interaction techniques
 - e.g., new mouse, new swipe interaction, ...



28

Comparing Two Alternatives

- Between groups requires many more participants than within groups
- See if differences are statistically significant
 - assumes normal distribution & same std. dev.
- Online companies can do large AB tests
 - look at resulting behavior (e.g., buy?)

29

Instructions to Participants

- Describe the purpose of the evaluation
 - "I'm testing the product; I'm not testing you"
- Tell them they can quit at any time
- Demonstrate the equipment
- Explain how to think aloud
- Explain that you will not provide help
- Describe the task
 - give written instructions
 - one task at a time



31

Autumn 2019

Prof. James A. Landay

Stanford University

Reporting the Results

- Report what you did & what happened
- Images & graphs help people get it!
- Video clips can be quite convincing



33

Heuristic Evaluation vs. User Testing

- HE is much faster
 - 1-2 hours each evaluator vs. days-weeks
- HE doesn't require interpreting user's actions
- User testing is far more accurate (by def.)
 - takes into account actual users and tasks
 - HE may miss problems & find "false positives"
- Good to alternate between HE & user testing
 - find different problems
 - don't waste participants

34

Summary

- User testing is important, but takes time/effort
- Use ????? tasks & ????? participants
 - *real tasks* & *representative* participants
- Be ethical & treat your participants well
- Want to know what people are doing & why? collect
 - process data
- Bottom line data requires ????? to get statistically reliable results
 - *more participants*
- Difference between between & within groups?
 - between groups: everyone participates in one condition
 - within groups: everyone participates in multiple conditions

35

Further Reading on Ethical Issues With Community-based Research

- Children and Families "At Promise, Beth B. Swadener, Sally Lubeck, editors, SUNY Press, 1995, <http://www.sunypress.edu/tc-2023-children-and-families-at-promise.aspx>
- "Yours is better!" Participant Response Bias in HCI, Proceedings of CHI 2012, by Nicola Dell, et al., <http://research.microsoft.com/pubs/163718/CHI2012-Dell-ResponseBias-proc.pdf>
- "Strangers at the Gate: Gaining Access, Building Rapport, and Co-Constructing Community-Based Research", Proceedings of CSCW 2015, by Christopher A. Le Dantec & Srah Fox, <http://dl.acm.org/citation.cfm?id=2675133.2675147&coll=DL&dl=ACM>
- "Imperialist Tendencies" blog post by Jan Chipchase, <http://janchipchase.com/content/essays/imperialist-tendencies/>
- "To Hell with Good Intentions" by Ivan Illich, speech to the Conference on InterAmerican Student Projects (CIASP), April 20, 1968, <http://www.swara.org/illich-hell.htm>

36

Next Time

- Lecture
 - Midterm ("closed-book")
- Studio
 - Hi-fi prototype planning session

37