

ACS P230 / Part II unit / CDH - Alan Blackwell & Advait Sarkar

Overview Practical experimental course lectures provide overview and sample of current research This introduction general principles, research approaches, current trends Specialist lectures: six specialist topics Design and run your own study discussion and feedback each week Final presentation of your results





	None	Casual	Student	Professional
HCI	9	9	8	2
ML	6	9	12	1



Lecture topics

- Week 2 Mixed initiative interaction (AB)
 information gain, cognitive ergonomics, agency & control
- Week 3 Labelling (AS)
 attribution, subjectivity, reliability, consistency
- Week 4 Program synthesis (AB)
 end-user programming, attention investment
- Week 5 Visual analytics (AS)
 visualisation, tool chains, design case studies
- Week 6 Addressing data bias (Dr Daniela Massiceti, Microsoft Research)
- Week 7 Interpretability (Prof Neil Lawrence)
- Week 8 Your research presentations

Practical work plan – for ACS and Part II

- Week I select research question (part II select replication study)
- Week 2 discuss potential study approaches (part II discuss preparations)
- Week 3 review and feedback on study proposals (part II may start work)
- Week 4 & 5 review logistical issues / practical progress
- Week 6 discuss preliminary findings
- Week 7 discuss research implications
- Week 8 final presentation

7

Assessment

- Everyone: final research report (80%)
 - Based on your practical work
 - Presented as a research paper
- Part II: phased submissions of work-in-progress drafts (20%)
 - Flat mark for each submission. Text can be revised/reused in the final report!
- ACS: optional (but recommended) work-in-progress drafts (0%)
 Advisory grades will be provided as feedback, for revision in final report
- ACS: reflective diary (20%)
 - Summarise lectures
 - Document discussions
 - Record development of your own thinking
 - Make 8 weekly entries ...
 - ... plus a final summative review



Continuous feedback

- Week 2 Research question (200 words) + a sample diary entry for ACS
- Week 3 Study design (400 words)
- Week 4 Another sample diary entry for ACS
- Week 5 Draft literature review for final report (400 words)
- Week 6 Draft introduction to report (200 words)
- Week 7 Draft results section for report (400 words)
- Week 8 Draft discussion section for report (200 words)

Week	Part II	ACS
2	Selection of study for replication, with summary of the key research question	Research question (200 words) + a sample diary entry
3	Detailed work plan for data collection and analysis	Study design (400 words)
4	Literature review (summary of original publication, and other work that cites it)	Another sample diary entry
5	Introduction to the replicated study	Draft literature review for final report (400 words)
6	Results and data analysis	Draft introduction to report (200 words)
7	Discussion of results and conclusion	Draft results section for report (400 words)
8		Draft discussion section for report (200 words)









First wave: HCI as engineering "human factors"

- The "user interface" (or MMI "man-machine interface") is a separate module, designed independently of the main system.
- Design goal is efficiency (speed and accuracy) for a human operator to achieve well-defined functions.
- Use methods from cognitive science to model users' perception, decision and action processes and predict usability
 - > At this point, relatively closely aligned with AI





























 Useful overview papers: Dudley, J. J., & Kristensson, P. O. (2018). A Review of User Interface Design for Interactive Machine Learning. ACM Transactions on Interactive Intelligent Systems, 8(2), 1–37. https://doi.org/10.1145/3185517 Abdul, A., Vermeulen, J., Wang, D., Lim, B.Y., & Kankanhalli, M. (2018). Trends and Trajectories for Explainable, Accountable and Intelligible Systems. Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18, 1–18. https://doi.org/10.1145/3173574.3174156
Ashktorab, Z., Jain, M., Liao, Q.V., & Weisz, J. D. (2019). Resilient Chatbots: Repair Strategy Preferences for Conversational Breakdowns. Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, 1–12. <u>https://doi.org/10.1145/3290605.3300484</u>
Eiband, M., Völkel, S.T., Buschek, D., Cook, S., & Hussmann, H. (2019). When people and algorithms meet: User-reported Problems in Intelligent
Everyday Applications. Proceedings of the 24th International Conference on Intelligent User Interfaces, Part F1476, 96–106. https://doi.org/10.1145/3301275.3302262
Everyday Applications. Proceedings of the 24th International Conference on Intelligent User Interfaces, Part F1476, 96–106. https://doi.org/10.1145/3301275.3302262 Buolamwini, J., & Gebru, T. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. Proceedings of Machine Learning Research, 81, 77–91. http://proceedings.mlr.press/v81/buolamwini18a.html
Everyday Applications. Proceedings of the 24th International Conference on Intelligent User Interfaces, Part F1476, 96–106. https://doi.org/10.1145/3301275.3302262 Buolamwini, J., & Gebru, T. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. Proceedings of Machine Learning Research, 81, 77–91. http://proceedings.mlr.press/v81/buolamwini18a.html Alvarado, O., & Waern, A. (2018). Towards Algorithmic Experience: Initial Efforts for Social Media Contexts. Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18, 1–12. https://doi.org/10.1145/3173574.3173860
 Everyday Applications. Proceedings of the 24th International Conference on Intelligent User Interfaces, Part F1476, 96–106. https://doi.org/10.1145/3301275.3302262 Buolamwini, J., & Gebru, T. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. Proceedings of Machine Learning Research, 81, 77–91. http://proceedings.mlr.press/v81/buolamwini18a.html Alvarado, O., & Waern, A. (2018). Towards Algorithmic Experience: Initial Efforts for Social Media Contexts. Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18, 1–12. https://doi.org/10.1145/3173574.3173860 Chen, NC., Suh, J., Verwey, J., Ramos, G., Drucker, S., & Simard, P. (2018). AnchorViz: Facilitating Classifier Error Discovery through Interactive Semantic Data Exploration. Proceedings of the 2018 Conference on Human Information Interaction&Retrieval - IUI '18, 269–280. https://doi.org/10.1145/3172944.3172950







Controlled Experimental Methods

- Participants (subjects), potentially in groups
- Experimental task
- Performance measures (speed & accuracy)
- Trials
- Conditions / Treatments / Manipulations
 - modify the system
 - use alternative systems
 - Use different features of the system
- **Effect** of treatments on sample means
 - Within-subjects (each participant uses all versions)
 - Between-subjects (different groups use different versions)









Experimental Manipulations Compare productivity gains (effect size) of version with new feature to one without? Will system work without the new feature? Will the experimental task be meaningful if the feature is disabled? Must new feature be presented second in a within-subjects comparison (order effect) Is your system sufficiently well-designed for external validity of productivity measure?

- Is full implementation necessary?
 - > Can you simulate features with Wizard of Oz technique?

39

Measurement

- Speed (classically 'reaction time')
 - Time to complete task
- Accuracy (number of (non)errors).
 - Is outcome as expected
- Trade-off between speed and accuracy?
 - Or poor performance on both?
 - Check correlation between them
- Task completion:
 - Stop after a fixed amount of time (ideally < I hour)</p>
 - > Measure proportion of the overall task completed

Self-Report

- Did you find this easy to use? (Likert scale)
 - applied value: appeal to customers
 - theoretical value: estimate 'cognitive load'

Danger of bias

- Subjective impressions of performance inaccurate
- Suffer from experimental demand
 - > Participants want to be nice to the experimenter
 - Should disguise which manipulation is the novel one

May be necessary to capture affect measures:

- Did you enjoy it, feel creative/ enthusiastic?
- Alternative is to collect 'richer' data ...

41

Think-aloud

- "Tell me everything you are thinking"
 - 'concurrent verbalisation'

Problems:

- Hard tasks become even harder while speaking aloud
- > During the most intense (interesting) periods, participants simply stop talking,
- Alternative:
 - > make video recording, or eye-tracking trace
 - playback for participant to narrate
 - 'retrospective verbal report'

Qualitative Data

- Protocol analysis methods, e.g.
 - verbal protocol transcript of recorded verbal data
 - video protocol recording of actions

Hypothesis-, or theory-driven

- > Create 'coding frame' for expected/hypothetical categories of behaviour
- > Segment the protocol into episodes, utterances, phrases etc
- Classify these into relevant categories (considering inter-rater reliability)
- Compare frequency or order statistically

Grounded theory

- > Open coding, looking for patterns in the data
- Stages of thematic grouping and generalization
- > Constant comparison of emerging framework to original data
- More interpretive, danger of subjective bias

Experiment Design

- Arrangement of participants, groups, tasks, trials, conditions, measures, and hypothesized effects of treatments
- Within-subjects designs are preferred
 because so much variation between individuals
- This leads to order effects:
 - > first condition may seem worse, because of learning effect
 - Iast condition may suffer from fatigue effect
 - task familiarity can't use the same task twice

Precautions:

- Prior training to reduce learning effects
- Minimise experimental session length to reduce fatigue effects
- > Use different tasks in each condition, but 'balance' with treatment and order
- These are typically combined in a 'latin square' where each participant gets a different combination

⁴³

Analysis

- For an easy life, plan your analysis before collecting data!
- Will quantitative data be normally distributed?
 - t-test to compare two groups
 - ANOVA to compare effect of multiple conditions (which include latin square of task and order)
 - > Pearson correlation to compare relationship between measures
- Distributions of task times are often skewed:
 - > a small number of individuals complete the task quite slowly
 - b don't exclude 'outliers' who have difficulty with your system
 - > log transform of time is usually found to be normally distributed
- Subjective ratings are seldom normally distributed
 - chi-square test of categories
 - non-parametric comparison of means



Evaluation

- Rather than testing hypothesis, or comparing treatments
 - ask 'is my system usable'?
- More typical of commercial practice, for short-term goals, rather than general understanding
 - Formative evaluation assesses options early in design process
 - > Summative evaluation identifies usability problems in a system you have built
 - Repeated for iterative refinement in user-centred design
- Weaker research, because no direct contribution to theory
 - > However applied research venues require evidence of any claims made for new tools

Field Study Methods

- Laboratory studies are not adequate for:
 - organizational context of system deployment
 - interaction within a user community

Typical methods:

- 'contextual inquiry' interviews
- focus group' discussions
- 'case studies' of projects or organisations
- 'ethnographic' field work as participant-observer
- All result in qualitative data, often transcribed, and analysed using grounded theory approaches







Representative tasks and measures

- Identify user activities you plan to observe
 - assigned tasks (controlled experiment)
 - or user's goal (observational study)
- Will these explore an interesting research question?
- What measures are relevant to that question?
- Will qualitative data analysis be necessary?
- Will there be a threat to external validity?
 - From task, measure or analysis



- Do you wish to carry out a comparison, an evaluation, or an open exploratory study?
- If you plan to conduct a controlled experiment, will it be possible to use a withinsubjects design?
- What data analysis method will you use?
- What would you need to do in order to complete a pilot study?
- > What ethical issues are raised by your planned research?
- A good starting point is to choose a published study that you would like to emulate / replicate

Theoretical goal

- > What do you expect to learn from conducting your study?
- What contribution will it make to the research literature in interaction with machine learning?
- Where would you publish the results?
- A good starting point is to review contributions that were made in published studies you would like to emulate
 - Warning be careful of studies done without prior training in HCI, and not published in peer-reviewed HCI venues.



Week	Part II	ACS
2	Selection of study for replication, with summary of the key research question	Research question (200 words) + a sample diary entry
3	Detailed work plan for data collection and analysis	Study design (400 words)
4	Literature review (summary of original publication, and other work that cites it)	Another sample diary entry
5	Introduction to the replicated study	Draft literature review for final report (400 words)
6	Results and data analysis	Draft introduction to report (200 words)
7	Discussion of results and conclusion	Draft results section for report (400 words)
8		Draft discussion section for report (200 words)
		+ don't forget diary entries every week