

Design and Implementation of Online Behavioral Experiments

nodeGame.org **Stefano Balietti** MZES and Heidelberg Introduction: From Lab to Online

@balietti @nodegameorg stefanobalietti.com@gmail.com

Your Instructor: Stefano Balietti

http://stefanobalietti.com

Currently

- Fellow in Sociology Mannheim Center for European Social Research (MZES)
- Postdoc at the Alfred Weber Institute of Economics at Heidelberg University

Previously

- Microsoft Research Computational Social Science New York City
- Postdoc Network Science Institute, Northeastern University
- Fellow IQSS, Harvard University
- PhD, Postdoc, Computational Social Science, ETH Zurich

Your Instructor: Stefano Balietti

http://stefanobalietti.com

Main Research Interests:

- Consensus formation, social influence, and polarization
- Inequality, redistribution, and preferences thereof
- Incentives schemes for peer review systems
- Optimal experimental design
- Philosophy of science, in particular Paul Feyerabend

- 1. Review *fundamental* concepts of experimental design
- 2. Understand *challenges* of *opportunities* for online experiments

- 1. Review fundamental concepts of experimental design
- 2. Understand *challenges* of *opportunities* for online experiments

- 1. Review *fundamental* concepts of experimental design
- 2. Understand *challenges* of *opportunities* for online experiments
- 3. Acquire *solid* knowledge of JavaScript and Node.JS programming

- 1. Review *fundamental* concepts of experimental design
- 2. Understand *challenges* of *opportunities* for online experiments
- Acquire *solid* knowledge of JavaScript and Node.JS programming
- 4. Implement your own experiment with the nodeGame framework

- 1. Review *fundamental* concepts of experimental design
- 2. Understand *challenges* of *opportunities* for online experiments
- 3. Acquire *solid* knowledge of JavaScript and Node.JS programming
- 4. Implement your own experiment in the nodeGame framework
- 5. Learn how to run your experiment with digital clouds and online labor markets (Amazon Mechanical Turk)

Experimentation: from Lab to Online

What Is an Experiment?





 An experiment is a *methodological procedure* carried out with the goal of verifying, falsifying, or establishing the validity of a hypothesis.

 A *test* under *controlled conditions* that is made to demonstrate a known truth, examine the validity of a hypothesis, or determine the efficacy of something previously untried.

 An experiment is an empirical method that *arbitrates between* competing hypotheses.

What Is a Synchronous Experiment?

- It is an experiment where participants
 - . (i) interact in a common environment,
 - . (ii) at the same time,
 - (iii) and their actions have an *immediate* effect on the decision process, and experimental outcome (e.g., monetary payoff)

· Synchronous experiments can be **turn-based** or **real-time**

Internal Validity

The extent to which a causal inference based on a study is warranted by an experiment.



Control



Deliberate choice of the variables that will be kept constants and those that will be varied: **treatment variables**

Independence



Treatment variables must *not* be correlated

•

Randomization



- · Participants are **randomly** assigned to treatments
 - Treatments are executed in controlled order
 - Stratification might be required

"There is a property common to almost all the moral sciences, and by which they are distinguished from many of the physical; that is, that *it is seldom in our power to make experiments with them*"

John Stuart Mill, 1836





Friedman, Cassar (2004) Economists go to the laboratory. Who, what, when and why.

"Economics unfortunately cannot perform the controlled experiments of chemists or biologists because [it] cannot easily control other important factors. Like astronomers or meteorologists, [it] generally must be content largely to observe."

Samuelson and Nordhaus 1985, Economics Textbook (italics added)





"There is no laboratory in which economists can test their hypotheses" 1993

Nobel prize for Vernon Smith

"for having established laboratory experiments as a tool in empirical economic analysis, especially in the study of alternative market mechanisms"



2002

"The key idea of induced value theory is that proper use of a reward medium allows an experimenter to induce pre-specified characteristics, and the subjects innate characteristics become largely irrelevant." (Vernon Smith, 1976)



"The key idea of induced value theory is that proper us medium allows an experimenter to induce pre-specified characteristics become largely irrelevant." (Vernon Smith, 1976)



"The key idea of induced value theory is that proper use medium allows an experimenter to induce pre-specified characteristics become largely irrelevant." (Vernon Smith, 1976)

Monotonicity: Participants always prefer more of the reward medium (*m*)

Salience: Changes in *m* are directly related to a participant's actions

Dominance: changes in a subject's utility from participating in the experiment are primarily due to changes in m.



"The key idea of induced value theory is that proper use medium allows an experimenter to induce pre-specified characteristics become largely irrelevant." (Vernon Smith, 1976)

Monotonicity: Participants always prefer more of the reward medium (*m*)

Salience: Changes in *m* are directly related to a participant's actions

Dominance: changes in a subject's utility from participating in the experiment are primarily due to changes in m.





Can you think of examples where failing to induce values can affect experimental results?

The extent to which the results of a study can be generalized or extended to others.





Galileo's hand-made drawing of the moon surface. Sidereus Nuncius (1610)



Galileo's hand-made drawing of the moon surface. Sidereus Nuncius (1610)



Galileo's hand-made drawing of the moon surface. Sidereus Nuncius (1610)

Some Limitations of Traditional Experimental Research





- Small groups of individuals, generally 4-12 (bounded by size uni. Lab)
- Behavior of individuals in social networks under-researched
- Limited in the number of hypotheses that can be tested
- Observations are limited in time
- Generally high costs for an academic budget
- Limited *external validity:*
 - Sample WEIRD student
 - Artificial conditions

Goals for Present and Future Behavioral Research



Goals for Present and Future Behavioral Research

İİİİİ

ŤŤŤŤ

İİİ

İİ



Size

- The size of a group influence the behavior of group members.
- Can we detect non-linear scaling laws?
- Do we observe interactions with the other dimesions?
- Can we perform online experiments with hundreds or thousands of participants?

Time

- Do real-time games qualitatively differ from turn-based games?
- Can we observe people playing games for long periods of time?
 Can we realize highly controlled, yet highly realistic experiments?
- Explore the full space of model parameter of a theory
 Perform large number of repetitions to produce more accurate effect-size estimates.
 Can we detect non-linear effects of treatments?

Granularity





express

Examples of Online Experiments

Survey Experiments

· Game-Based Asynchronous Experiments

Game-Based Synchronous Experiments

Income Inequality has increased dramatically in the United States since 1980 Incomes of poorer and middle-income families have grown very little while top incomes have grown a lot.

How would YOU be doing if inequality had not increased?

The slider below shows how much each group would make if incomes had grown by the same percentage since 1980 for all groups: the poor, the middle class, and the rich. Use the slider to answer the questions below.

A household making \$25,800 today would instead be making \$35,200 if inequality had not changed since 1980. In other words, if growth had been evenly shared, this household would have earned 37% more.




Survey Experiments

Income Inequality has increased dramatically in the United States since 1980. Incomes of poorer and middle-income families have grown very little while top incomes have grown a lot.

Insert

a new block of questions, or provide additional info

Measure preferences for outcome variable at the *individual* level

How would YOU be doing if inequality had not increased?

The slider below shows how much each group would make if incomes had grown by the same percentage since 1980 for all groups: the poor, the middle class, and the rich. Use the slider to answer the questions below.

A household making \$25,800 today would instead be making \$35,200 if inequality had not changed since 1980. In other words, if growth had been evenly shared, this household would have earned 37% more.

Relatively simple to implement

Kuziemko et al. (2015) How Elastic Are Preferences for Redistribution

Game-Based Asynchronous Experiments

Create a common environment where a single variable is changed (e.g network structure)

Observe variation on outcome variable at the *group* level

Generally complex design Can extend for long periods of time



Centola (2010) The Spread of Behavior in Online Social Networks

Game-Based Synchronous Experiments

Create a common environment where a single variable is changed (e.g., group size or payoff)

Observe variation on outcome variable at the *group* level

Can be really complex or stylized



Mao et al. (2016) An Experimental Study of Team Size and Performance on a Complex Task



Mao et al. (2016) An Experimental Study of Team Size and Performance on a Complex Task

Game-Based Synchronous Experiments

Create a common environment where a single variable is changed (e.g., group size or payoff)

Observe variation on outcome variable at the *group* level

Can be really complex or stylized

Time left 00:01

Round: 1 of 10 Points: 0



Balietti, Goldstone and Helbing (2016) "Peer Review and Competition in the Art Exhibition Game" PNAS 113(30) pp. 8414–8419

Art Exhibition Game

E SLARIDI ×	anonisum									Fri Oct 07
C localhost:8080										Q,
AME 📋 BurdenShare_TE	🗋 BS_T2 📋 Ult_Req_Test 🐰	Google Maps Lite 📋 SPA 📋 Pen	sion 💻 Bad Conditions 🗋 NEWPR 🗋 NE	WU 🗋 R Data Analysis 📋	PR_AUTH CLAST_PR CLOCAL	8080 1	2 3	4 5	6	
Time left 00:50	Round: 1 of 12	Points: 0								Continue
cale head horizontally				1		No past e	exhibitions	vet.		
Scale head vertically										
ve and Eyebrow height	1	1								
Eye spacing										
cale eyes horizontally	0									
Scale eyes vertically										
Eyebrow length										
Eyebrow angle										
Eyebrow from eye										
Eyebrow spacing										
Upper lip										
Lower lip										
Zoom in										
20011111										

Balietti, Goldstone and Helbing (2016) "Peer Review and Competition in the Art Exhibition Game" PNAS 113(30) pp. 8414–8419



COMPLEXITY OF EXECUTION

- Recruitment
- Payment

Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

• <u>Recruitment</u>

Payment

• Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

Recruiting Platforms

- Amazon Mechanical Turk
 <u>https://www.mturk.com</u>
- TurkPrime
 <u>https://www.turkprime.com</u>
- Figure Eight (ex Crowd Flower, ML) https://www.figure-eight.com
- Prolific.ac <u>https://www.prolific.ac</u>

 Reddit (for bursty access) <u>https://www.reddit.com/r/WebGames</u>

- YouGov (country-specific) <u>https://yougov.co.uk</u>
- PollFish (also built-in surveys)
 <u>https://www.pollfish.com</u>

- Psychological Research on the Net <u>http://psych.hanover.edu/research/exponnet.html</u>
- The Web Experiment List
 <u>http://www.wexlist.net</u>
- Online Social Psychology Studies
 <u>http://www.socialpsychology.org/expts.htm</u>

Citizen Science Initiatives

- Volunteer Science
 <u>https://volunteerscience.com</u>
- Science@Home <u>https://www.scienceathome.org</u>
- Zoouniverse

https://www.zooniverse.org

- Lab in the Wild <u>https://labinthewild.org</u>
- CitizenLab

https://www.citizenlab.co

- Recruitment
- Payment

• Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)



- Recruitment
- Payment

• Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)



Figure 1: 1a: The communication network among Amazon Mechanical Turk workers. 1b-1f: Subnetworks for Reddit HWTF (magenta; 660 workers, 1837 edges), MTurkGrind (red; 392 workers, 1331 edges), TurkerNation (green; 200 workers, 740 edges), Facebook (blue; 133 workers, 357 edges), and MTurkForum (black; 312 workers, 244 edges).

Ming et al. (2016) "The Communication Network Within the Crowd"

- Recruitment
- Payment

Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

- Recruitment
- Payment

• Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)



- Anonymity might increase deviant behavior
- (but get access to specialized populations)

- The same person uses the same computer to participate repeatedly.
- The same person uses different computers to participate repeatedly.
- Different persons use the same computer to participate.

- Recruitment
- Payment

Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

- Recruitment
- Payment

Anonymity

Dropouts

- Interference
 - <u>Data Quality</u>
 - Farming
 - <u>Cheating</u>

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

Improve Data Quality and Attention

Content responsive faking:

(a) responses influenced by the item content but (b) are not completely accurate.

purposeful faking (e.g., malingering in medical surveys) socially desirable response of both intentional and nonintentional varieties

Content nonresponsivity:

responding without regard to item content.

Large variance estimation (2.5 – 50%; 10.6 for students) Inattentive only on some questions

Meade, A. W. and S. B. Craig (2012). Identifying Careless Responses in Survey Data. Psychological methods 17(3), 437.

Improve Data Quality and Attention

• Engagement is key

ui design, interactive elements, colors, length, feedback...

- Nonsensical or "bogus" items "30th February"
- Instructed response items

"To monitor quality, please respond with a two for this item"

- Unlikely sets "Select all the good qualities you possess"
- Consistency indexes Evaluating differences among responses to items highly similar in content
- Too quick response times and other obvious patterns in responses

Meade, A. W. and S. B. Craig (2012). Identifying Careless Responses in Survey Data. Psychological methods 17(3), 437.

Commitment Device

- Ask participants if they paid attention to the *previous* set of questions
- It prompts participants to pay attention to the subsequent set of questions
- Its purpose is fulfilled regardless of whether the respondents answer honestly

Commitment Device Example

Alesina, Miano, and Stancheva (2018) "Immigration and Redistribution" NBER 24733

Before proceeding to the next set of questions, we want to ask for your feedback about the responses you provided so far. It is vital to our study that we only include responses from people who devoted their full attention to this study. This will not affect in any way the payment you will receive for taking this survey. **In your honest opinion, should we use your responses, or should we discard your responses since you did not devote your full attention to the questions so far?**

- Yes, I have devoted full attention to the questions so far and I think you should use my responses for your study

- No, I have not devoted full attention to the questions so far and I think you should not use my responses for your study

AMT Bot Scare

Dennis et al. (2018)

"MTurk Workers' Use of Low-Cost 'Virtual Private Servers' to Circumvent Screening Methods: A Research Note" https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3233954

Ahler et al. (2018) The Micro-Task Market for "Lemons": Collecting Data on Amazon's Mechanical Turk http://gsood.com/research/papers/turk.pdf

TurkPrime Blog

https://blog.turkprime.com/after-the-bot-scare-understanding-whats-been-happening-with-data-collection-on-mturk-and-how-to-stop-it https://blog.turkprime.com/moving-beyond-bots-mturk-as-a-source-of-high-quality-data

Tools to filter Farmers' IP addresses

https://github.com/MAHDLab/rIP

https://itaysisso.shinyapps.io/Bots

- Recruitment
- Payment

Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

Fairness for Crowd Workers

- IRB, IRB, IRB (Institutional Review Board)
- Crowd worker have flexibility, but almost no rights
- Low pay, menial tasks, no accrual of sick days, or contributions to pension, no unions, etc.
- In the US, 20 million people (Pew Research)
- Primary income for:
 - 1.25-2.5 million (Economic Policy Institute, 2018)
 - 2-10 million (JPMorgan Chase Institute, 2017)



Mary Gray & Sid Suri http://www.ghostwork.info http://www.inthecrowd.org

- Recruitment
- Payment

• Anonymity

Dropouts

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)

- Recruitment
- Payment

• Anonymity

<u>Dropouts</u>

- Interference
 - Data Quality
 - Farming
 - Cheating

- Technical Barriers for implementation
- Cross-Device Compatibility
- Accessibility (impaired users)



Dropouts



Dropouts (Attrition)

Arechar et al. (2018) "Conducting interactive experiments online"

"The most severe problem for online interactive studies, and the largest discrepancy with laboratory experiments, is attrition (participant dropout)."

Major challenge to internal validity, if dropout rates vary with treatment, selection bias may arise

Dropouts (Attrition)



Arechar et al. (2018) "Conducting interactive experiments online"

Just a little bit more than 50% of the groups that started finished without dropouts

What Keeps Participants In?

Intrinsic Motivation

Volunteer Science, Galaxy Zoo, Science@Home (Deci, 2005)

Monetary incentives

about 75% of US Turkers reported that MTurk was their primary or secondary source of income (Paolacci et al, 2009)

Curiosity

Solving and unrelated, but intriguing puzzle (Law et al., 2016)

Personalized Feedback

Lab-in-the-Wild (Ye et al., 2017)

Two Simple Techniques To Limit Dropouts

• Implement warm-up phase (nodeGame levels)



Implement seriousness checks

Reips (2002) Standards for Internet-Based Experimenting

What to Do When Dropouts Happen?



What to Do When Dropouts Happen?



- Terminate experiment immediately
- Continue with less participants (notify? Does behavior change?)
- Replace with bots (notify? Behavior is different)
- Replace with other humans (difficult to implement technically and humanly)

What to Do When Dropouts Happen?



- Terminate experiment immediately
- Continue with less participants (notify? Does behavior change?)
- Replace with bots (notify? Behavior is different)
- Replace with other humans (difficult to implement technically and humanly)

→ Try to design games that are **dropout-robust**, Make Leaving Harder, Inform about total time for experiment, Inform about waiting times during the experiment, Wait for reconnections, Notify
Challenges in Online Experiments

- Recruitment
- Payment

• Anonymity

Dropouts

• Fairness

- Interference
 - Data Quality
 - Farming
 - Cheating

- <u>Technical Barriers for</u> <u>implementation</u>
- <u>Cross-Device Compatibility</u>
- Accessibility (impaired users)

List of Platforms for Online experiments

Group-Behavior

- nodeGame ③
 <u>https://nodeGame.org</u>
- Wextor (pioneer) <u>https://www.wextor.eu</u>
- Otree (large base) https://www.otree.org
- Lioness (new) https://lioness-lab.org
- Breadboard (networks) http://breadboard.yale.edu

- Empirica (very new) <u>https://empirica.ly</u>
- TurkServer (discontinued) https://turkserver.readthedocs.io

Individual

- JSPsych (many plugins) https://www.jspsych.org
- PsiTurk (groups possible) <u>https://psiturk.org</u>

Other Resources

Browser Stack

https://www.browserstack.com

CSS Frameworks
 <u>https://getbootstrap.com</u>

nodeGame: the Good Parts



- Powerful API to customize experiments
- Access to low-level details
- Integrated JS database
- Fast and highly scalable
- Game Levels
- Modular design (games, widgets, window)
- Well-documented (and active Forum)
- Integration with Amazon Mechanical Turk