

Clustering standard errors at the session level*

Duk Gyoo Kim[†]

May 11, 2025

Abstract

Session-specific features of a laboratory experiment, if those exist, do not disappear by clustering standard errors at the session level. Randomly ordering or counterbalancing sessions to deal with sampling issues, cannot justify clustering the standard errors at the session level. Unlike empirical studies, for laboratory experimental studies, the experimental design reflected on the researchers' intention should primarily determine the clustering level. In a typical controlled laboratory experiment where subjects make choices in the same environment repeatedly, clustering at a participant level is intended by the experimental design, and standard errors could be larger (that is, a statistical inference could be more conservative) when clustered at the individual or decision-group level than the session level. It implies that clustering standard errors at the session level can lead to false-positive treatment effects if it is mistakenly chosen. Having a small per-session sample to increase the number of sessions could yield undesirable heterogeneities that are hard for the experimenter to control or observe.

Keywords: Lab experiment, Cluster-robust standard errors, Statistical inference

JEL codes: C18, C90

1 Introduction

This paper aims to convince the experimental economists and the readers interested in lab-experimental studies that the session-level clusters should be used only in particular situations with proper justification. A session is typically defined as a group of individuals who participate in the same laboratory experiment simultaneously. For an

*I thank Guillaume Fréchette, Franziska Heinicke, Sang-Hyun Kim, Wooyoung Lim, Yoshiyasu Rai, Euncheol Shin, Donggyu Sul, Wladislaw Mill, and the participants at Mannheim/ZEW Experimental seminar, the Korean Economic Review International Conference, and 2020 ESA Global Online Around-the-Clock Conference for their helpful comments, and Elisa Casarin for her research assistance.

[†]School of Economics, Yonsei University, kim.dukgyoo@yonsei.ac.kr

experiment adopting a between-subject design, a subject participated in a session of one treatment¹ without knowing² the treatment condition of each session. A set of observations from an individual is a proper subset of the entire sample from a session, which is a proper subset of the entire sample from the same treatment. Thus, adding individual- or session-fixed effects on the regression does not help us examine a treatment effect due to perfect multicollinearity. Discussions on multi-way clustering with the non-nested clustering units (e.g., Petersen, 2008; Bertrand et al., 2004; Cameron et al., 2011) are not practically helpful because session-level clusters nest individual-level clusters.

Obtaining accurate standard errors of the treatment effect is fundamental for proper statistical inference. Although many studies discuss the proper use of cluster-robust standard errors (e.g., Cameron et al., 2008; Abadie et al., 2017; de Chaisemartin and Ramirez-Cuellar, 2024), to the best of my knowledge, only a few studies, including Moffatt (2016) explicitly discuss it within the context of laboratory experiments.³ Perhaps it is why we see some researchers report standard errors clustered at the session level and some at the individual level. Among all 322 published papers using lab-experimental data at the *Experimental Economics* from March 2010 to March 2020, 124 papers mentioned cluster-robust standard errors. Standard errors of 40 papers are clustered at the participant level, and those of 34 papers are at the session level.⁴

My research question is when we should cluster standard errors at the session level for analyzing experimental data from a controlled laboratory. With some caveats, the preview of the answer is that we should avoid considering the session as a clustering unit. It is often argued that standard errors should be clustered at the session level concerning the session-specific effects. For example, Keith Marzilli Ericson, a co-editor of the *Journal of Public Economics*, points out that many lab-experimental papers fail to randomly assign participants to treatment, with claiming that once researchers "[d]o session-level randomization,"⁵ then the statistical "[i]nference should cluster standard

¹On the contrary, a within-subject design assigns a participant to two or more treatments. In this case, considering session-level clusters is even less persuasive as the design's primary purpose is to examine individual changes.

²It is the main feature of a laboratory experiment in economics. Unlike medical experiments where knowing the treatment condition (for example, taking a new drug for controlling a high blood pressure) cannot affect the observational outcomes (unless subjects can choose their blood pressure level at will), knowing treatment conditions may arise issues in sample selection (i.e., choosing a session that a participant believes to maximize payoffs) and in experimenter demand effects (i.e., choosing accordingly decisions after grasping the purpose of the study). Since the treatment condition is known only after the subject participated in an experiment, the treatment is randomly assigned from the subject's perspective.

³Moffatt (2016) explains that researchers can consider different (subject-level as the lowest and session-level as the highest) clustering. When analyzing example data, he uses subject-level clustering only.

⁴Some papers use exogenously given clusters, such as classes and cohorts. Other papers used cluster-adjusted standard errors when analyzing empirical data, not experimental data. A few papers consider a fixed independent group as a clustering unit, which I will discuss in Section 3.

⁵In a typical setting, one session is conducted at one time, so session-level randomization practically

errors at the session level."⁶ Also, it is not uncommon that reports from referees point out that the standard errors should be clustered at the session level. Most of the time, their reasoning, including ones that Ericson made on his blog post, is that there might be some "static" session effects (Fr  chette, 2012).⁷ This reasoning—using session-level cluster adjustment for session effects—is not on solid ground. Concerns for static session effects are a reason for randomizing or counterbalancing the sessions so that the session-specific idiosyncratic features can be integrated out; such concerns are not a reason for clustering standard errors at the session level. My claim is not new: In the context of randomized field experiments, de Chaisemartin and Ramirez-Cuellar (2024) similarly claim that clustering standard errors at the unit-of-randomization level may lead to a severe downward bias of the variance estimator of the treatment effect. I am worried that many researchers seem to use session-level clustered standard errors to remedy session effects, without further justifying why a session should be the cluster level.

In line with Abadie et al. (2020), who claim to consider design-based uncertainty instead of sampling-based one for statistical inference, I claim that the experimental design reflecting the researchers' intention should determine the clustering level, and only when the design intends the strongly positive observational relationship within a session, standard errors should be clustered at the session level. Figure 1 summarizes my arguments.

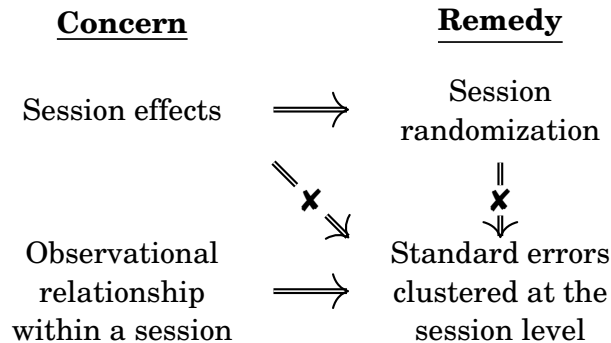


Figure 1: Clustering at the session level is not a remedy for session effects.

It is worth noting that the critical difference between typical empirical data and controlled-laboratory data is on whether the researchers' intentions involve the data-generating process. Namely, in experimental economics, the researchers who analyze the experimental data are *data generators* as well. Thus, it is indispensable to consider the

implies randomly ordering treatment conditions over the sessions.

⁶More details are in his blog post (Ericson, 2018, [Design Issues in Economics Lab Experiments: Randomization](#)).

⁷The static session effects can be understood as a realization from a noise distribution that affects the observational outcomes in level. The dynamic session effects can be understood as the observational relationships across subjects within a session due to the feedback from interactions with other participants.

researchers' intention of the experimental design. When the experiment adopts a random rematch protocol, the researchers intend to disconnect (or minimize) the dynamic relationships between decision rounds. If the session size is sufficiently large, the random rematch will approximate perfect stranger matching. Even when the session size is small, the random rematch prevents subjects from considering dynamic strategies because they do not know who the previously matched players in what decision round were. When the researchers allow the subjects to play the same game repeatedly, unless the researchers expect the subjects to play in a completely random manner, they intend the positive observational dependence within a subject. However, the researchers do not *design* the direction of the possible dynamic session effects. Of course, the researchers may expect some interactions among subjects will affect their decisions, but the directions of such interactions are not designed. Take the public goods game with a perfect stranger match, for example. Suppose one subject observes more contributions from others in one decision round. Would he respond to his observation by increasing his contribution later because he wants to be a conditional cooperator, or by decreasing it because he observes that the public goods are well provided without his contributions and the free-riding incentives become salient? When the experimental design allows subjects to play the public goods game repeatedly, is the direction of the interactions designed as well? If answers to both questions are negative, then the observational relationship within a session is not intended by the experimental design. Moreover, as I will elaborate later in Section 5, the negative relationship between a subject's decision and the decisions of the previous group members would substantially *exacerbate* the type-1 error.

To minimize destructive discussions, I emphasize two things that I am *not* claiming. First, I am not claiming that we should not worry about static session effects. The experimenter's crucial responsibility is to maintain every session's environment as homogeneous as possible, except for the treatment conditions being examined. Since it is challenging, if not impossible, to make every session environment identical, the experimenter must make sure both the control-group participants and the treatment-group participants are from the same population by randomizing or counterbalancing the session orders. In this regard, I entirely agree with what Ericson wrote in his blog: "Your subject population could be changing over time (perhaps early subjects are more eager, or have lower value of time). Or news events could change beliefs and preferences. The list of potential stories can be long; some can be ruled out, others cannot." Indeed, the list of potential stories is long: Perhaps one experimenter manages sessions better than another experimenter. Subjects participating in an early morning session may have distinctive characteristics than other subjects. An exogenous aggregate shock (e.g., COVID-19 pandemic) may arise between sessions. Some sessions may be conducted in more dis-

turbing situations due to unexpected constructions, delays caused by technical glitches, or unexpectedly high/low temperatures, to name a few. Thus, it is legitimate for readers, editors, and referees to demand more sessions if they are concerned about potential static session effects. For similar reasons, a sequential modification of the experimental design—earlier sessions conducting X and Y and (perhaps several months) later sessions conducting X' and Z—may significantly undermine the internal validity of the research. Although I am wholly sympathetic to the concerns about the static session effects, it is a reason for being careful about sampling subjects from the same population pool by randomizing the sessions, a reason for making session environments as homogeneous as possible, and a reason for checking and controlling for session-particular features, but not the reason for clustering standard errors at the session level.

Second, I am not claiming that clustering standard errors at the session level is futile, especially when the experiment exploits strongly *positive* interactions among subjects in a session. A session-level cluster can undoubtedly address the "dynamic" session effects or the observational dependence within the session. It is sometimes tightly aligned with the experimental design, especially when the subjects made decisions only once or the session-(or "market")-level interactions are of the main interest.⁸ Although Fréchette (2012) argues for using standard errors clustered at the session level when there is "only one observation per subject so that we do not need to keep track of the periods" (p. 488),⁹ it should not be merely extrapolated to a case where subjects make several decisions under the same environment. Thus, this paper can be understood as an extension of his paper. Again, the current paper focuses less on the studies where the experimental design intends strong and positive dynamic interactions within a session, where I believe the session-level clustering is appropriate, but it focuses on the discussions about the proper cluster level when individuals in the lab make repeated decisions.

The rest of this paper is organized as follows. Sections 2 and 3, without formal expositions, illustrate why standard errors need to be clustered and why clustering at the session level should be considered carefully. I target the potential readers interested in laboratory experiments but not equipped with solid econometrics background. Those who do not need justification for the use of standard error clustering may skip these

⁸For example, Engelmann and Hollard (2010) have participants who made only a small number of decisions and focus more on the interaction within a session. Cipriani et al. (2017)'s interest is on the session-level information contagion, so the interactions within a session are inherited from the design. Corgnet et al. (2018) similarly justify their use of session-level clustering because each experimental market features a zero-sum game where an increase in one trader's earnings mechanically reduces other traders' possible gains within a session. Bracha et al. (2015) and Carpenter (2016) experimentally examine the attributes of labor supply, which is the accumulation of an individual's decisions, so it is pertinent to regard the labor supply as one observation per subject.

⁹This restriction is judicious because Fréchette (2012) focuses on the discussions about the session effects, not the relative importance of subject-specific effects and the session effects.

sections. Section 4 presents a simple econometric model to pinpoint the issues in choosing proper cluster levels. Section 5 shows some Monte-Carlo simulation results. Section 6 discusses practical issues regarding cluster-robust standard errors for the laboratory data, and Section 7 concludes.

2 Why do we cluster standard errors?

Clustered standard errors should be considered when observations within a cluster are related to each other. In other words, if the observations within a cluster are similar, then the errors within a cluster will be more correlated than those of the entire sample. Thus, without "penalizing" the observational similarity, we will have downward-biased standard errors, leading to false-positive treatment effects more often. Throughout the entire paper, I consider situations with no true treatment effects. Thus, by a false positive, I mean a Type-I error to mistakenly reject the true null hypothesis (no treatment effect). In more practical terms, using standard errors clustered at the session level may yield some "stars" in the regression table when they are not supposed to appear.

To elaborate on what I mean by "penalizing" similarity, consider the following. There are ten observations: five from a control group experiment, and the other five from the treatment group experiment. Assume that except for the treatment condition, everything is homogeneous and appropriately controlled. A researcher tests if the mean control-group observation is different from the mean treatment-group observation.

	Control session					Treatment session				
ID	1	2	3	4	5	6	7	8	9	10
Obs.	2	1.2	1.5	1.7	2.5	2.6	1.8	1.5	1.3	2.1

A standard t-test does not reject the null hypothesis that the two means are the same (p-value=0.8085). The standard error of the mean difference is 0.3194.¹⁰

Now, suppose that the researcher's half-sleeping RA mistakenly duplicated the observations several times.

¹⁰For this and following results, I used Stata and MATLAB. Refer to the README file in the Open Science Framework repository (<https://osf.io/sp3kt>) to learn how to replicate these results and accordingly, how to run the similar analysis.

	Control session					Treatment session				
ID	1	2	3	4	5	6	7	8	9	10
Obs.	2	1.2	1.5	1.7	2.5	2.6	1.8	1.5	1.3	2.1
Obs.2	2	1.2	1.5	1.7	2.5	2.6	1.8	1.5	1.3	2.1
Obs.3	2	1.2	1.5	1.7	2.5	2.6	1.8	1.5	1.3	2.1
Obs.50	2	1.2	1.5	1.7	2.5	2.6	1.8	1.5	1.3	2.1

The (un-clustered) standard error of the mean difference is 0.040, and the null hypothesis is rejected ($p=0.0487$). This inference is obviously wrong because it ignores the perfect correlation between observations at the participant level. The standard error clustered at the participant level is 0.3014, and the treatment effect becomes insignificant again. Table 1 summarizes the three regression results on the treatment dummy and a constant. The estimated coefficient of the treatment dummy captures the treatment effect, the mean difference between the control and the treatment. The first column shows that the treatment effect is not statistically significant when the original observations are only considered. The second and the third columns show the regression results using the half-asleep RA's duplicated dataset. The second column shows that the un-clustered standard errors incorrectly lead to a statistically significant treatment effect, but in the third column, the treatment effect is insignificant with the clustered standard errors.

Obs	(1)	(2)	(3)
Treatment	0.0800 (0.25)	0.0800** (1.98)	0.0800 (0.27)
_cons	1.780*** (7.88)	1.780*** (62.19)	1.780*** (8.49)
Cluster SE	–	–	ID
N	10	500	500

t statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table 1: A false-positive effect when SEs are unclustered.

Although the example above is too unrealistic because of the perfect correlation between observations within a cluster, we can draw one clear takeaway message. A researcher must consider clustering standard errors when observations within a cluster are expected to be positively related. That is the way of providing more robust statistical results.

A naturally following question is what the proper cluster level would be. Unlike other empirical studies where the clustering units can be non-nested, potential clusters in a between-subject experiment—individual or session—are nested: The set of individual-level observations is a proper subset of the set of session-level observations. In the following sections, I claim that if the lab experiment asks the participants to make decisions in a similar environment repeatedly, clustering at the participant level is intended by the experimental design, so it is unnatural to cluster standard errors at the session level.

3 Illustration: Is the session-level clustering robust?

If standard errors clustered at the session level are larger than those at the individual level, it means that the session-level observations are more positively correlated than the individual's repeated choices. This may not be the case when the subjects are asked to make decisions in the same environment repeatedly.

To illustrate my claims, I use hypothetical data. This choice allows me to capture the features of actual data that are most relevant to my argument, yet it does not rely on previous studies whose design and approach to data analysis might have been driven by other considerations.

Imagine a particular type of controlled lab experiment on a group decision making,¹¹ adopting a between-subject design, random rematch, anonymity, and no opportunity for communication. To be more illustrative, suppose that six subjects per session have ten repeated decision rounds choosing an integer between 1 and 50, and the payoff of each round is determined by the subject's decision, a randomly-matched pair's decision, and some luck. At the beginning of a new round, the subjects are randomly rematched with another subject in the session. Their decisions are made anonymously, and they are not allowed to communicate with each other. Each subject participates in only one session. Suppose a researcher collected data from four (two control and two treatment) sessions,¹² as shown in Figure 2.

Each vertical line of Figure 2 shows a vector of an individual's decisions over ten rounds. A researcher wants to examine the mean treatment effect. If we do not cluster standard errors, the mean control-group observation is significantly different from the

¹¹For an experiment where a single player makes a streak of decisions under some uncertainties, it is straightforward to cluster standard errors at the individual level. Here I focus on experiments involving strategic decisions.

¹²Admittedly, the example here contains too few samples (six subjects in each session of four). The mere purpose is to display the entire observations in a simple scatterplot, and I do not intend to use the small sample properties. Increasing the exemplary data tenfold, that is, six subjects in each session of 40, does not change the results. The dataset used for Figure 2 is available at the [Open Science Framework repository](https://osf.io/nr3j6) (osf.io/nr3j6).

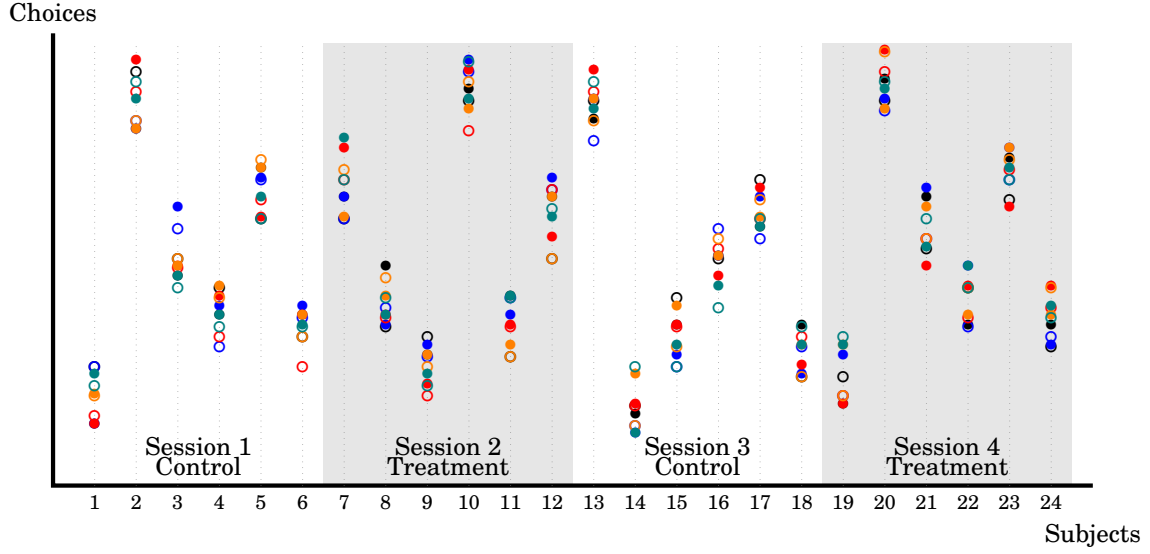


Figure 2: Data from Four Sessions

mean treatment-group observation ($\bar{y}_C=21.55$, $\bar{y}_T=24.15$, $t=1.9808$, $p\text{-value}=0.0488$). The standard error of the difference is 1.313. The corresponding regression results are shown in model (1) of Table 2.

	(1)	(2)	(3)
Treatment	2.600** (1.98)	2.600 (0.63)	2.600*** (4.43)
_cons	21.55*** (23.22)	21.55*** (7.26)	21.55*** (36.75)
Cluster SE	—	ID	Session
N	240	240	240

t statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table 2: False-positive effects when SEs are clustered at a session level

In this hypothetical data, standard deviations of the individual-level observations are small, which implies that they made similar choices over the rounds. The session-level standard deviations are as large as the standard deviation of the entire sample. If we cluster the standard error of the mean difference at the subject level, the difference is no longer statistically significant (model (2) in Table 2). However, clustering standard errors at the session level does not handle the false-positive treatment effect (model (3) in Table 2.) In other words, a researcher might mistakenly claim a significant treatment

effect when clustering standard errors at the session level.

Unless the experiment encourages every subject to make completely arbitrary decisions, the observational similarity at the participant level is intended by the experimental design when the experiment asks a participant to make repeated decisions. Two of the primary reasons for the repetitions are to increase the number of observations and to allow subjects to learn the equilibrium of the game in the course of getting feedback. Thus, when the learning effects are not of their primary interest, researchers often focus on the observations from the latter half decision rounds. Those observations are likely "less noisy," meaning that the individual's decisions are similar over rounds. Roughly put, the observations become similar to the half-sleeping RA's duplicated data.

Instead of a random rematch, if the experiment involves fixed independent groups of the participants over the repeated decision rounds, then clustering standard errors at the group level could also be considered. If the experiment features repeated games (e.g., [Duffy and Fehr, 2018](#)) or asks each group to achieve a collective goal (e.g., [Hortala-Vallve et al., 2013](#)), it is appropriate to have a fixed group to interact over time. In this case, both individual-level clusters and group-level clusters can be intended by the experimental design. If there are two or more ways of defining a cluster, and those ways are equally justifiable by the experimental design, then a researcher, given that he/she wants to report more robust statistical results, must choose a cluster within which observations are more related. One rule of thumb is to check the standard deviation of the observations within a potential cluster. For illustration, consider a public goods experiment with a fixed group of three subjects. Suppose that a researcher has collected data shown in Table 3.

Group 1							Group 2							Group 3							Group 4						
ID	1	2	3	4	5	6	ID	7	8	9	10	11	12	ID	7	8	9	10	11	12	ID	7	8	9	10	11	12
r01	0	10	5	2	5	3	r01	7	6	5	1	1	2	r01	7	6	5	1	1	2	r01	7	6	5	1	1	2
r02	0	10	4	2	5	3	r02	3	4	4	5	3	3	r02	3	4	4	5	3	3	r02	3	4	4	5	3	3
r03	0	10	4	3	5	1	r03	2	2	3	6	5	3	r03	2	2	3	6	5	3	r03	2	2	3	6	5	3
r04	0	10	3	3	5	0	r04	0	1	4	3	5	4	r04	0	1	4	3	5	4	r04	0	1	4	3	5	4
r05	0	10	3	3	6	0	r05	1	0	1	5	6	3	r05	1	0	1	5	6	3	r05	1	0	1	5	6	3
r06	0	10	3	3	5	0	r06	0	1	0	7	5	6	r06	0	1	0	7	5	6	r06	0	1	0	7	5	6
r07	1	10	3	3	5	0	r07	1	0	0	9	7	8	r07	1	0	0	9	7	8	r07	1	0	0	9	7	8
r08	0	10	3	2	4	0	r08	0	0	0	10	10	9	r08	0	0	0	10	10	9	r08	0	0	0	10	10	9
r09	0	10	3	1	4	0	r09	0	0	0	10	10	10	r09	0	0	0	10	10	10	r09	0	0	0	10	10	10
r10	0	10	1	2	5	0	r10	0	0	0	10	10	10	r10	0	0	0	10	10	10	r10	0	0	0	10	10	10
Std.	0.32	0.00	1.03	0.70	0.57	1.25	Std.	2.22	2.07	2.06	3.17	3.08	3.19	Std.	2.22	2.07	2.06	3.17	3.08	3.19	Std.	2.22	2.07	2.06	3.17	3.08	3.19
Std.(Group)=4.25							Std.(Group)=1.95							Std.(Group)=2.05							Std.(Group)=3.06						

Table 3: Strong dependence at the participant level (L) or the decision-group level (R)

If the individual choices vary little, as illustrated on the data from Groups 1 and 2 in Table 3, standard deviations of the participant-level observations (varying from 0.00 to 1.25) are smaller than those of the group-level observations (1.95 to 4.25). It implies that individual observations are more related to each other than group observations, so in this case, the standard error clustered at the individual level should be used. Meanwhile, if a group's choices vary less than individual choices, as illustrated on the data from Groups 3 and 4, the researchers may consider standard errors clustered at the independent-group level. I imagined situations where a group collectively reaches to complete free-riding or complete cooperation. Such a case may happen when group members' previous actions influence a subject's action more than the subject's own previous actions.¹³

The discussion above may be extrapolated to justify session-level clusters. If the session-level observations are more positively correlated than the individual's or the decision group's repeated choices, it could mean that the session-level clustered standard errors yield more robust statistical results. I am skeptical about this data-driven approach,¹⁴ and I will discuss it after introducing cluster-robust inference in the following section.

4 Cluster-Robust Inference

In this section, I present a prototype parametric¹⁵ model for cluster-robust inference of the mean treatment effect. I assume only one treatment (and one control) and that the experimenter controls session effects appropriately, so the model does not include them. An econometrician has $N = (S + S) \times I \times R$ observations in total, where S is the number of controlled and treated sessions, I is the number of per-session subjects, and R is the number of repetitions of the same game.¹⁶

For simplicity, set the dependent variable as the deviation from the mean of control-

¹³Some papers, e.g., [Robbett \(2014\)](#) and [Gallo and Yan \(2015\)](#), used the term "session" as a fixed independent group. In this case, it would be appropriate to cluster standard errors at the session (or independent-group) level.

¹⁴I should clarify that I do not mean to avoid any data-driven approach. Instead, I claim that the experimental design, or the intention of the researchers who design the experiment, should be prioritized over the purely statistical data features. If the experimental design well justifies two different clustering levels, then researchers could use a level that renders more (statistically) conservative reports. Thus I suggest observing some statistical features *within* the design-driven approach.

¹⁵Some researchers prefer non-parametric tests that take the session-level aggregate data as one independent data point. This approach may be free from the concern about the clustering issues as well as parametric assumptions, but the current paper does not address the comparative advantages of non-parametric methods.

¹⁶For expositional simplicity, I assume that the number of the subjects and the repetitions are the same for each session and that the number of controlled sessions is equal to the number of treated sessions, but these assumptions do not affect the main messages.

group observations. Then, the treatment effect is captured by β in

$$y_i = \beta T_i + \varepsilon_i, \quad (1)$$

where $i = 1, \dots, N$ is an index for observations, and $E[\varepsilon_i] = 0$. T_i has a value 1 if the observation is from the treated session and 0 otherwise. $\beta = 0$ implies that the treatment-group and the control-group means are the same.¹⁷ With a slight abuse of notation, T is a set of treated observations such that for $i \in T$, $T_i = 1$. The OLS estimator is

$$\hat{\beta} = \frac{\sum_i T_i y_i}{\sum_i T_i^2} = \frac{\sum_{i \in T} y_i}{SIR}, \quad (2)$$

and the variance of the estimator is

$$V[\hat{\beta}] = E[(\hat{\beta} - \beta)^2] = \frac{V[\sum_{i \in T} \varepsilon_i]}{S^2 I^2 R^2} \quad (3)$$

$V[\sum_{i \in T} \varepsilon_i] = \sum_{i \in T} \sum_{j \in T} \text{Cov}[\varepsilon_i, \varepsilon_j] = \sum_{i \in T} \sum_{j \in T} E[\varepsilon_i \varepsilon_j]$ is of our interest. If errors are uncorrelated, that is, $E[\varepsilon_i \varepsilon_j] = 0$ for $i \neq j$, it becomes $\sum_{i \in T} E[\varepsilon_i^2]$, and its sample analog, $\sum_{i \in T} (y_i - \hat{\beta} T_i)^2 = \sum_{i \in T} u_i^2$, yields the heteroskedasticity-robust standard error. We are concerned that this is not the case, at least within a cluster. Let C_i denote the cluster that i belongs to. If $E[\varepsilon_i \varepsilon_j] \neq 0$ for i and $j \in C_i$,

$$V_{clu}[\hat{\beta}] = \frac{\sum_{i \in T} \sum_{j \in T} E[\varepsilon_i \varepsilon_j] \mathbf{1}_{j \in C_i}}{S^2 I^2 R^2}, \quad (4)$$

where $\mathbf{1}_A$ is an indicator whose value is 1 when condition A holds and 0 otherwise. Given that the number of clusters is sufficiently large,¹⁸ we can use the variance estimate

$$\hat{V}_{clu}[\hat{\beta}] = \frac{\sum_{i \in T} \sum_{j \in T} u_i u_j \mathbf{1}_{j \in C_i}}{S^2 I^2 R^2} \quad (5)$$

Two remarks are (1) if the cluster is the entire set, $\hat{V}_{clu}[\hat{\beta}]$ becomes zero because $\sum_{i \in T} u_i = 0$, and (2) if clusters are defined in a far-fetched manner so that $u_i u_j$ is negative

¹⁷Provided that the experiment is appropriately conducted so that the treatment condition is orthogonal to other control variables such as subjects' characteristics, and treatment- and control-group participants are drawn from the same population distribution, the treatment effect should not be affected by other control variables. In other words, adding other control variables does not alter main claims of this paper.

¹⁸Another concern would regard the asymptotic refinement of the clustered standard errors when the number of clusters is small (Cameron et al., 2008). Kézdi (2004) shows simulation results that 50 clusters are often large enough for accurate inference. A typical laboratory experiment has fewer sessions than 50, while it has more subjects than 50. A common practice of using the standard errors at the session level seems to ignore this concern. Bootstrap-based tests (e.g., Roodman et al., 2019) instead of t tests should be considered when considering the session-level clusters, but this point is beyond what the current paper concerns. Nonetheless, the observations from the simulation results in Section 5 remain valid even when I increase the number of sessions. See Appendix A.

for many pairs of i and j , the cluster-robust variance estimate could even be smaller than the heteroskedasticity-robust one.

Without loss of generality, lexicographically order the observations such that $i = s \times n \times r$, $s = 1, \dots, 2S$, $n = 1, \dots, I$, and $r = 1, \dots, R$. Then $\sum_{i \in T} \sum_{j \in T} u_i u_j \mathbf{1}_{j \in C_i}$ is the summation of entities on the block diagonal matrices of $u_i u_j$, $i, j \in T$. Figure 3 illustrates the difference between standard errors clustered at the individual level and the session level. The main difference is that there are more off-diagonal (but still within a larger block diagonal) entities when clustering standard errors at the session level (see hatched areas in Figure 3). If standard errors clustered at the session level are larger than those at the individual level, it implies that the summation of those off-diagonal entities is positive. It happens when the signs of u_i and u_j are, in general, the same for $j \in C_i$. Since the residual is the deviation from the conditional mean, the same signs imply error correlations.

If the experimental design intends the strong and positive correlation between, for example, the first choice of individual i and the last choice of individual j in the same session, then the session-level cluster might be used. Perhaps someone's initial choice profoundly affects other's later choices so that those observations are related. Many questions then follow. Is that relationship stronger than the relationship between a subject's own choices? Is that relationship stronger than the relationship between the last observations in one session and those in another session with the same treatment condition? It is undoubtedly possible that errors are weakly but positively correlated within a session, but considering a larger-size cluster comes at a price. Given the same number of observations, larger-size clusters have a stronger downward bias due to fewer clusters. Although statistical analysis software uses finite-cluster corrections,¹⁹ it is unclear whether the standard error's downward bias will be appropriately corrected when a session is used as a clustering unit. While the experimenters may be concerned about the observational relationship within a session for any laboratory experiments, they should want to double-check whether the experimental design is well associated with such a relationship from the beginning.

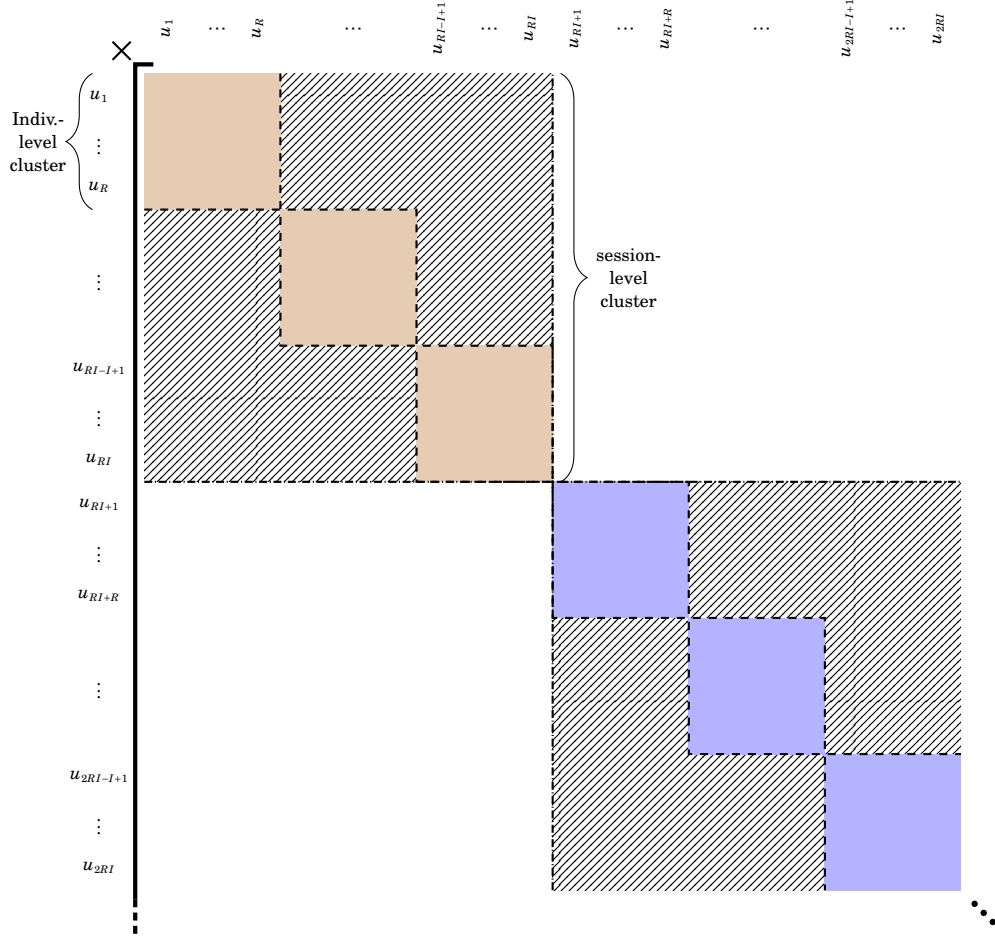
5 Simulations

For backing up the illustrations in Section 3, this section presents some Monte-Carlo simulation results.²⁰ For all simulation results, I consider $S = 4$ (four sessions per control

¹⁹For example, Stata uses $\frac{G}{G-1} \frac{N-1}{N-k} u_i$ instead of u_i , where G is the number of clusters, N is the number of observations, and k is the number of regressors.

²⁰The code and the instructions are available at the [Open Science Framework repository](https://osf.io/nr3j6) (osf.io/nr3j6).

Figure 3: Individual-level vs. session-level clusters



This figure illustrates a part of N -by- N matrix where entity at (i, j) is $u_i u_j$. The cluster-robust standard error of the treatment effect is the sum of the entities on block-diagonal sub-matrices. Clustering standard errors at the session level, compared to the individual level, involves more off-diagonal entities.

and treatment each), $I = 18$ (18 subjects per session), and $R = 5$ (the last five repetitions of the game). Those numbers are in the range of typical laboratory experiments. Further, I assume that the group size is three (or six groups per session) and a random rematch (six groups are randomly shuffled every round). For this simulation, I have in mind a standard public goods game where a subject can choose a contribution level between 0 and 50 or a Tullock contest where a subject can invest up to 50 tokens to win the prize. Simulations are conducted in the following way.

1. Generate the treatment indicator, the session number, subject id, and the group number.²¹
2. For each iteration, $2 * S * I * R$ observations are generated in the following way.
 - (a) In the first round of the experiment, each subject draws a choice from a discrete uniform distribution between 0 and 50.
 - (b) From the second round and beyond, subjects tend to (i) stick to their previous choice and (ii) consistently respond to their previous group choices. Thus, the observation in the next round is a linear combination of three numbers: the number chosen by the subject in the previous round (with linear coefficient ρ_{ind}), the average number chosen by the group members in the previous round (with ρ_{ss}), and the randomly generated number from the same discrete uniform distribution (with $1 - \rho_{ind} - \rho_{ss}$).
3. Regress observations on the treatment dummy and a constant.
4. Calculate the heteroskedasticity-robust standard error²², the standard error clustered at the session level, and the standard error clustered at the individual level. Count if the two-tailed p-value of the t-statistic ($\frac{\hat{\beta}}{SE}$) is less than 0.05.
5. Repeat Steps 2–4 for 10,000 times.

At least four points are worth mentioning. First, for this simulation, the population mean of the control is the same as the population mean of the treatment. Since the primary purpose of this exercise is to check the claim that the standard errors clustered at the session level may lead to a false-positive result (that is, reporting a statistically significant treatment effect when there is supposed to be no treatment effect), it is important to set no fundamental differences between the control and the treatment. Second,

²¹A fixed match is not considered in this simulation, but the simulation code can also serve the purpose. Check the instructions for the simulation.

²²Since the simulation data is generated in a homogeneous manner, the heteroskedasticity-robust standard error is practically identical to the OLS standard error. I omit the simulation results for the OLS standard errors.

the sign of ρ_{ss} captures the direction of the responses to the previous observations from the matching group.²³ For example, a positive ρ_{ss} can be interpreted that the subject tries to imitate the previous average observation (such as conditional cooperation in the public goods game and learning the optimal investment level by observing other's investments), and a negative ρ_{ss} implies that the subjects deviate what the average players do (such as more free-riding after observing sufficient contributions from others and more investment than the average level to win the contest). Third, the magnitude of ρ_{ss} is naturally limited as the number of subjects increases. If the experiment adopts a perfect stranger match with sufficiently many subjects, whatever the subjects had learned from the previous game has nothing to do with the new game. When a session consists of 18 subjects and the size of a group is three, the probability of meeting at least one member of the previous group again is $2/17 \approx 0.1176$. With having this probability in mind, I vary ρ_{ss} from -0.20 to 0.20 . Whichever the sign of ρ_{ss} , the larger value implies the larger observational dependence across subjects within a session. Fourth, since the first-round data is generated from the discrete uniform distribution, and the second round and beyond depend on the initial realizations, a learning effect toward a particular decision point (for example, a Nash equilibrium) is not considered.

ρ_{ind}	ρ_{ss}	Mean			St.Dev.			Pr(p-value<0.05)		
		SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}
0.8	-0.2	0.8660	1.2595	1.6713	0.0304	0.3659	0.0808	0.1929	0.0972	0.0150
0.8	-0.1	0.8282	1.4068	1.6743	0.0300	0.4092	0.0745	0.2653	0.0972	0.0283
0.8	0	0.8106	1.6224	1.7067	0.0302	0.4731	0.0701	0.3413	0.0987	0.0507
0.8	0.1	0.8201	1.9268	1.7795	0.0314	0.5628	0.0701	0.4194	0.0984	0.0853
0.8	0.2	0.8597	2.3330	1.8932	0.0344	0.6836	0.0768	0.4805	0.0999	0.1298
0.5	-0.2	0.9215	1.0738	1.3314	0.0215	0.3121	0.0719	0.1054	0.0969	0.0206
0.5	-0.1	0.8479	1.1073	1.2725	0.0208	0.3216	0.0671	0.1500	0.0972	0.0324
0.5	0	0.7821	1.1649	1.2249	0.0206	0.3383	0.0623	0.2081	0.0975	0.0495
0.5	0.1	0.7262	1.2596	1.1933	0.0211	0.3664	0.0583	0.2778	0.0973	0.0777
0.5	0.2	0.6836	1.4062	1.1828	0.0223	0.4089	0.0556	0.3579	0.0972	0.1155
0.2	-0.2	1.1077	1.0444	1.2374	0.0195	0.3038	0.0682	0.0472	0.0982	0.0277
0.2	-0.1	1.0214	1.0468	1.1725	0.0182	0.3042	0.0643	0.0685	0.0986	0.0365
0.2	0	0.9398	1.0555	1.1101	0.0172	0.3066	0.0605	0.0956	0.0983	0.0506
0.2	0.1	0.8631	1.0735	1.0514	0.0166	0.3119	0.0569	0.1291	0.0969	0.0692
0.2	0.2	0.7921	1.1072	0.9980	0.0166	0.3217	0.0538	0.1802	0.0963	0.0917
0	0	1.0956	1.0442	1.0983	0.0183	0.3036	0.0605	0.0501	0.0994	0.0516

Table 4: Monte-Carlo simulation results

Table 4 shows the simulation results with different ρ_{ind} and ρ_{ss} . Three columns un-

²³I assume that ρ_{ind} is always non-negative because it captures the subject's decision consistency. $\rho_{ind} \approx 0$ means that the subject merely ignores what he/she previously chose, and $\rho_{ind} < 0$ implies that the subject intentionally oscillates the decisions.

der "Mean" show the average value of 10,000 simulated heteroskedasticity-robust standard errors (SE_{het}), standard errors clustered at the session level (SE_{clu}^{ss}), and the standard errors clustered at the individual level (SE_{clu}^{ind}), respectively. The following three columns under "St.Dev." show the standard deviations of those standard errors, and the last three columns show the actual test sizes for a nominal test size of 0.05. Since there are no population differences between the control and the treatment, such size is supposed to converge to the significance level (0.05) as the number of repetitions increases, when $\rho_{ind} = \rho_{ss} = 0$. Note that in the last row of Table 4, the size is close to 0.05 when using heteroskedasticity-robust standard errors and standard errors clustered at the individual level, but not standard errors clustered at the session level.²⁴

One can observe that the standard deviation of SE_{clu}^{ss} is distinctively larger than those of SE_{het} and SE_{clu}^{ind} . Figure 4 shows a histogram of one of the results ($\rho_{ind} = 0.8, \rho_{ss} = 0.2$) summarized in Table 4, which clearly illustrates that the standard errors at the session level (Clu-Ss) vary substantially more than heteroskedasticity-robust standard errors (Het) and the standard errors clustered at the individual level (Clu-Ind). It implies that the test statistics would substantially vary when the standard errors are clustered at the session level, although the data are obtained through an identical process. As a result, the statistical test finds a significant treatment effect more often when clustering the standard errors at the session level than at the individual level, except for some cases with $\rho_{ss} = 0.20$.

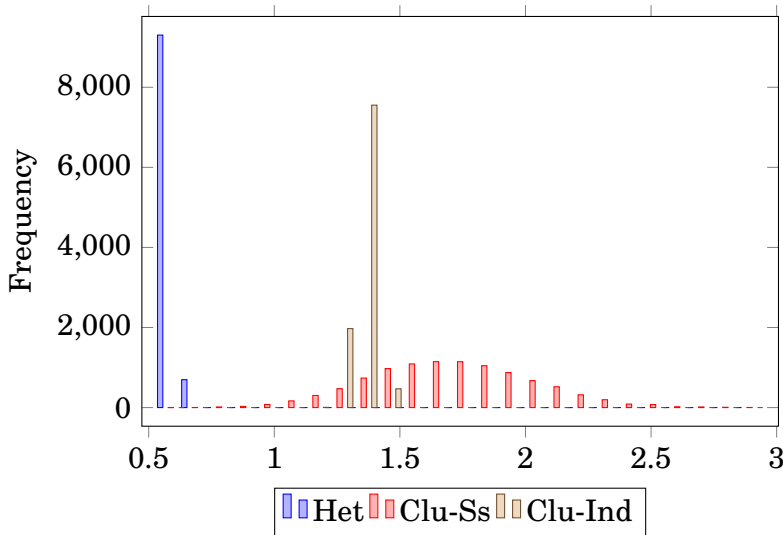


Figure 4: Simulated (iter=10,000) standard errors clustered at the session level vary.

²⁴The oversized test with the standard errors clustered at the session level could be due to the insufficient number of clusters. The small number of sessions by itself could be the reason to avoid clustering standard errors at the session level, but decreasing the session size to increase the number of sessions is not desirable. See Section 6 for relevant discussions about the session size and the number of sessions.

Another noticeable result from this simulation is that the test sizes using standard errors clustered at the session level are sometimes *even larger* than the test sizes using the heteroskedasticity-robust standard errors, especially with weak or negative dependence across observations such as in the case of $\rho_{ind} = 0.2$ and $\rho_{ss} = -0.2$. This never happens with standard errors clustered at the individual level. It means that the attempt to find more robust statistical results could undesirably lead to the opposite outcomes when clustering standard errors at the session level.

6 Discussions

6.1 Standard deviation as a rule of thumb

Suppose clustering observations can be done in two or more ways, equally justifiable by the experimental design. In that case, a researcher ought to choose a cluster within which the observations are more related to each other. I propose to check the within-cluster standard deviations of the observations. Recall that the residuals of the simple regression are the deviations from a conditional mean. A sufficiently smaller within-cluster standard deviation than the standard deviation of the entire sample may imply that the residuals flock together within sessions, and hence the errors are correlated within the cluster. Thus, when both session-level and individual-level clusters are equally justifiable by the experimental design, my rule of thumb is to compare within-cluster standard deviations. Consider I individual-level clusters, and S session-level clusters, where an individual-level cluster is a proper subset of a session-level cluster. Let std_I and std_S respectively denote the standard deviation of the individual-cluster observations and that of the session-cluster observations. If $std_I < std_S$ in general, then consider clustering the standard errors at the individual level.²⁵

If std_S is distinctively smaller than the standard deviation of the entire sample of the same treatment, then the session-level clustering might lead to larger standard errors. If this is the case, especially when the experimental design does not intend the observational relationship within a session, a researcher may want to check whether the sessions are sufficiently randomized. A relevant situation is illustrated in Figure 5, which displays a scatterplot of observations from eight (four control and four treatment) sessions. Almost all residuals from sessions 1, 4, and 8 are positive, and almost all residuals from sessions 2, 3, and 6 are negative. Thus, the products of pairs of those residuals within a session have positive values, and the standard error clustered at the

²⁵Consider, as an extreme case, std_I to be zero such that $std_I = 0 < std_S$. The zero standard deviation means that an individual's decisions are identical (or perfectly correlated) across decision rounds. That is, the observations are the same as the half-sleeping RA's duplicated data illustrated in Section 2.

session level will be larger than the heteroskedasticity-robust one. However, if thinking conversely, one may wonder whether the samples are balanced because otherwise, it is hard to explain the stark differences between identically-treated sessions. This distinctive variation across sessions may be due to the failure of session randomization or the session size being too small.

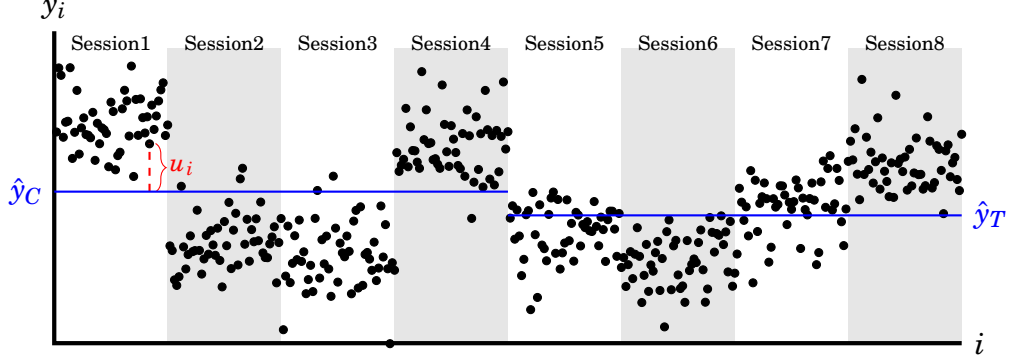


Figure 5: Small session-level standard deviations may question balanced sampling.

Each dot represents subject i 's decision. The first half and the last half sessions are with the same treatment condition, respectively. The figure illustrates a situation with a stronger dependency of the observations within a session than that within a treatment condition.

6.2 Further thoughts on the session-level clustering

I have claimed that the session-level cluster should be cautiously used either when the experimental design intends the observational relationship within a session or when the session-level residuals flock together. The latter reasoning is inconclusive as it relies on the mechanical aspects of the data, not on the experimental design. If researchers consider clustering standard errors at the session level because it generates large standard errors, why not consider clustering at the date-of-session or time-of-session level, why not at the experimenter level, and why not at the experiment level for meta-analysis²⁶ if those do the same or a better job? Furthermore, if we are willing to embrace let-the-data-tell-us approaches, why don't we consider direct tests on the level of clustering using bootstrapping or random forest algorithms (Ibragimov and Müller, 2016; MacKinnon et al., 2023)?

An ad-hoc definition of a session also obscures the session-level clustering. Suppose 24 subjects show up at the lab, and the experimenter decides to split the subjects into two subgroups without informing them, but across the subgroups, the experiment proceeds

²⁶Embrey et al. (2017) provide a meta-analysis of prior experimental research on the finitely repeated prisoner's dilemma and report the standard errors clustered at the study level. Detailed discussions and robustness checks on the clustering level for a meta-analysis are in the paper's Appendix A.4.

identically. In this case, would a session consist of 24 subjects, or would two sessions consist of 12 subjects each? This concern becomes more relevant to the fixed-group experiment. Suppose there are 24 subjects in one session, but only 12 subjects show up in another session due to severe no-shows. If a fixed group of six subjects repeatedly play a game, one session is a cluster of four decision groups, and the other one is a cluster of two decision groups. If the former session's observations are seemingly less related because of more (potentially heterogeneous) groups, the session with fewer participants affects standard errors clustered at the session level more. Is having a different weight on each session justified?²⁷

Another practical issue is the trade-off between the session size and the number of sessions. Given that the total number of participants is practically limited, considering session-level clusters pushes researchers toward having more sessions with fewer subjects per session. This approach is problematic in several aspects. First, many experiments adopt a random rematch design to minimize the strategic interactions between the games. If the number of subjects per session is small, then the indirect interactions are unavoidable. If a subject plays ten games with a randomly paired partner in a session of 40 subjects, the probability that a subject does not meet any match again is 28.34%, but with 16 subjects per session, that probability plummets to 1.89%. Such a low probability implies that, with fewer subjects per session, the fundamental reason for adopting a random rematch is compromised: Although the subjects do not know whether the current match is new, they know that it is highly likely to have met before or would meet again. Second, fewer subjects per session can prevent us from having a balanced sample: Given that the subjects are drawn from the same population distribution, small-sized sessions feature more (un)observable variations in sessions.²⁸ Suppose each session consists of only four subjects each, and the female proportion dramatically varies from 0% to 100%. How can a researcher be sure whether the session effects are controlled, and if not, how does she distinguish the gender-ratio effect from others unless having more sessions with sufficiently large variations of the gender ratio? What is even worse, if the substantial variations across sessions are due to unobservable characteristics, not like observable gender ratios? A vicious cycle of demanding more sessions to control issues with small-size sessions may be established.

My argument here is simple. Suppose a researcher considers either 12 sessions with 6 subjects per session or 4 sessions with 18 subjects per session. If a researcher adopts

²⁷Moreover, Müller (2020) points out that when clusters are of different sizes, the p-values from typical statistical packages, such as Stata, are not reliable.

²⁸Tversky and Kahneman (1974) point out that most people are unaware that smaller samples are subject to higher variance in characteristics. If both a large and a small hospital recorded the days when more than 60% of the newborns were boys, which hospital is more likely to record more such days? Only 22% of the subjects correctly answered the small hospital.

a random rematch to disconnect or minimize the dynamic relationships between decision rounds, then the random rematch can be understood as a practical proxy for a perfect-stranger match. Thus, having 4 sessions with 18 subjects per session is more aligned with the researcher’s intention. However, if the researcher intends to facilitate the interactions across subjects over decision rounds, it would be better to have many (small-sized) sessions, but it begs the question of using the random rematch from the beginning.

7 Conclusions

Session-specific idiosyncratic features can and should be integrated out when the researchers carefully randomize the sessions. If the purpose of clustering standard errors is to make more robust standard errors to minimize false-positive treatment effects, then one must consider clusters within which observations are more related, but across which observations vary. In a controlled laboratory experiment where participants repeatedly make choices in the same environment, individual-level clusters should be considered first, as the observational similarity within an individual is intended by the experimental design. Takeaway messages are summarized below:

1. The experimenter’s crucial responsibility is to ensure the participants in both the control and the treatment sessions are from the same population distribution and make each session environment as homogeneous as possible.
2. If the experiment asks participants to make repeated decisions in a similar environment, the experimental design intends clusters at the participant (or independent decision-group) level. Thus, it is natural to cluster standard errors at the participant (or independent decision-group) level.
3. The standard deviation of individual-level observations, when the individuals are asked to make decisions in the same environment repeatedly, tends to be smaller than that of session-level observations. Thus, clustering standard errors at the participant level may yield more conservative statistical results.
4. If the experimental design equally justifies two ways of clustering observations, a researcher would choose a cluster within which observations are more positively related.
5. Although not justifiable by the experimental design, clustering standard errors at the session level may be considered if the session-level observations are more

strongly and positively related than those of participant- or group-level observations. It begs further questions of why a session should be a level for clustering, among several other potential levels, and whether the sessions have balanced samples.

References

- Abadie, Alberto, Susan Athey, Guido W. Imbens, and Jeffrey Wooldridge**, “When Should You Adjust Standard Errors for Clustering?,” Working Paper 24003, National Bureau of Economic Research November 2017.
- , —, —, **and** —, “Sampling-Based Versus Design-Based Uncertainty in Regression Analysis,” *Econometrica*, 2020, 88 (1), 265–296.
- Bertrand, Marianne, Esther Duflo, and Sendhil Mullainathan**, “How Much Should We Trust Differences-In-Differences Estimates?,” *The Quarterly Journal of Economics*, 02 2004, 119 (1), 249–275.
- Bracha, Anat, Uri Gneezy, and George Loewenstein**, “Relative Pay and Labor Supply,” *Journal of Labor Economics*, 2015, 33 (2), 297–315.
- Cameron, A. Colin, Jonah B. Gelbach, and Douglas L. Miller**, “Bootstrap-Based Improvements for Inference with Clustered Errors,” *The Review of Economics and Statistics*, 2008, 90 (3), 414–427.
- , —, **and** —, “Robust Inference With Multiway Clustering,” *Journal of Business & Economic Statistics*, 2011, 29 (2), 238–249.
- Carpenter, Jeffrey**, “The labor supply of fixed-wage workers: Estimates from a real effort experiment,” *European Economic Review*, 2016, 89, 85–95.
- Cipriani, Marco, Antonio Guarino, Giovanni Guazzarotti, Federico Tagliati, and Sven Fischer**, “Informational Contagion in the Laboratory,” *Review of Finance*, 06 2017, 22 (3), 877–904.
- Corgnet, Brice, Mark Desantis, and David Porter**, “What Makes a Good Trader? On the Role of Intuition and Reflection on Trader Performance,” *The Journal of Finance*, 2018, 73 (3), 1113–1137.
- de Chaisemartin, Clément and Jaime Ramirez-Cuellar**, “At What Level Should One Cluster Standard Errors in Paired and Small-Strata Experiments?,” *American Economic Journal: Applied Economics*, January 2024, 16 (1), 193–212.

- Duffy, John and Dietmar Fehr**, “Equilibrium selection in similar repeated games: experimental evidence on the role of precedents,” *Experimental Economics*, 2018, 21, 573–600.
- Embrey, Matthew, Guillaume R. Fréchette, and Sevgi Yuksel**, “Cooperation in the Finitely Repeated Prisoner’s Dilemma,” *The Quarterly Journal of Economics*, 08 2017, 133 (1), 509–551.
- Engelmann, Dirk and Guillaume Hollard**, “Reconsidering the Effect of Market Experience on the “Endowment Effect”,” *Econometrica*, 2010, 78 (6), 2005–2019.
- Ericson, Keith M.**, “Design Issues in Economics Lab Experiments: Randomization,” 2018.
- Fréchette, Guillaume R.**, “Session-effects in the laboratory,” *Experimental Economics*, Sep 2012, 15 (3), 485–498.
- Gallo, Edoardo and Chang Yan**, “The effects of reputational and social knowledge on cooperation,” *Proceedings of the National Academy of Sciences*, 2015, 112 (12), 3647–3652.
- Hortala-Vallve, Rafael, Aniol Llorente-Saguer, and Rosemarie Nagel**, “The role of information in different bargaining protocols,” *Experimental Economics*, 2013, 16, 88–113.
- Ibragimov, Rustam and Ulrich K. Müller**, “Inference with Few Heterogeneous Clusters,” *The Review of Economics and Statistics*, 03 2016, 98 (1), 83–96.
- Kézdi, Gábor**, “Robust Standard Error Estimation in Fixed-Effects Panel Models,” *Hungarian Statistical Review*, 2004, Special 9, 96–116.
- MacKinnon, James G., Morten Ørregaard Nielsen, and Matthew D. Webb**, “Testing for the appropriate level of clustering in linear regression models,” *Journal of Econometrics*, 2023, 235 (2), 2027–2056.
- Moffatt, Peter G.**, *Experiments: Econometrics for Experimental Economics*, London New York, NY: Macmillan International Higher Education, 2016.
- Müller, Ulrich K.**, “A More Robust t-Test,” Working Paper 2020.
- Petersen, Mitchell A.**, “Estimating Standard Errors in Finance Panel Data Sets: Comparing Approaches,” *The Review of Financial Studies*, 06 2008, 22 (1), 435–480.

Robbett, Andrea, “Local Institutions and the Dynamics of Community Sorting,” *American Economic Journal: Microeconomics*, August 2014, 6 (3), 136–156.

Roodman, David, Morten Ørregaard Nielsen, James G. MacKinnon, and Matthew D. Webb, “Fast and wild: Bootstrap inference in Stata using boottest,” *The Stata Journal*, 2019, 19 (1), 4–60.

Tversky, Amos and Daniel Kahneman, “Judgment under Uncertainty: Heuristics and Biases,” *Science*, 1974, 185 (4157), 1124–1131.

A More simulation results

(* The number of simulation iterations is 1,000, not 10,000 in this Appendix.)

ρ_{ind}	ρ_{ss}	Pr(p-value<0.05)			Pr(p-value<0.01)			Pr(p-value<0.001)		
		SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}
0.8	-0.2	0.181	0.085	0.010	0.086	0.033	0.000	0.027	0.014	0.000
0.8	-0.1	0.265	0.092	0.026	0.134	0.031	0.004	0.065	0.012	0.000
0.8	0	0.352	0.103	0.051	0.216	0.033	0.006	0.114	0.014	0.000
0.8	0.1	0.426	0.099	0.086	0.298	0.035	0.021	0.182	0.013	0.005
0.8	0.2	0.493	0.100	0.130	0.381	0.042	0.051	0.254	0.017	0.009
0.5	-0.2	0.100	0.088	0.011	0.030	0.032	0.002	0.006	0.010	0.000
0.5	-0.1	0.137	0.086	0.027	0.050	0.033	0.006	0.011	0.009	0.000
0.5	0	0.197	0.081	0.042	0.083	0.039	0.008	0.028	0.012	0.000
0.5	0.1	0.283	0.088	0.074	0.139	0.035	0.017	0.069	0.013	0.003
0.5	0.2	0.373	0.090	0.109	0.219	0.033	0.039	0.116	0.012	0.006
0.2	-0.2	0.033	0.087	0.018	0.008	0.029	0.003	0.000	0.013	0.000
0.2	-0.1	0.060	0.090	0.026	0.011	0.032	0.006	0.002	0.010	0.000
0.2	0	0.084	0.088	0.043	0.017	0.032	0.008	0.004	0.009	0.000
0.2	0.1	0.118	0.084	0.058	0.042	0.032	0.011	0.007	0.009	0.001
0.2	0.2	0.158	0.082	0.086	0.067	0.035	0.022	0.019	0.008	0.006
0	0	0.034	0.093	0.037	0.008	0.030	0.008	0.000	0.012	0.000

Table 5: simulation results, with different p -values

ρ_{ind}	ρ_{ss}	4 sessions per treatment Pr(p-value<0.05)			6 sessions per treatment Pr(p-value<0.05)			8 sessions per treatment Pr(p-value<0.05)		
		SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}
0.8	-0.2	0.181	0.085	0.010	0.186	0.069	0.011	0.187	0.063	0.010
0.8	-0.1	0.265	0.092	0.026	0.281	0.064	0.021	0.257	0.063	0.028
0.8	0	0.352	0.103	0.051	0.352	0.068	0.037	0.322	0.064	0.061
0.8	0.1	0.426	0.099	0.086	0.437	0.068	0.076	0.406	0.069	0.096
0.8	0.2	0.493	0.100	0.130	0.493	0.070	0.124	0.475	0.068	0.131
0.5	-0.2	0.100	0.088	0.011	0.099	0.068	0.021	0.102	0.057	0.021
0.5	-0.1	0.137	0.086	0.027	0.146	0.068	0.032	0.135	0.053	0.032
0.5	0	0.197	0.081	0.042	0.200	0.073	0.047	0.194	0.062	0.051
0.5	0.1	0.283	0.088	0.074	0.279	0.068	0.072	0.258	0.065	0.084
0.5	0.2	0.373	0.090	0.109	0.358	0.065	0.116	0.340	0.061	0.111
0.2	-0.2	0.033	0.087	0.018	0.043	0.070	0.027	0.045	0.069	0.027
0.2	-0.1	0.060	0.090	0.026	0.070	0.070	0.032	0.066	0.065	0.035
0.2	0	0.084	0.088	0.043	0.091	0.068	0.047	0.095	0.062	0.050
0.2	0.1	0.118	0.084	0.058	0.129	0.066	0.067	0.119	0.056	0.072
0.2	0.2	0.158	0.082	0.086	0.177	0.071	0.084	0.165	0.053	0.093
0	0	0.034	0.093	0.037	0.047	0.068	0.046	0.045	0.068	0.051

Table 6: simulation results, with different numbers of sessions

ρ_{ind}	ρ_{ss}	18 subjects per session Pr(p-value<0.05)			27 subjects per treatment Pr(p-value<0.05)			36 subjects per session Pr(p-value<0.05)		
		SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}	SE_{het}	SE_{clu}^{ss}	SE_{clu}^{ind}
0.8	-0.2	0.181	0.085	0.010	0.174	0.076	0.008	0.189	0.094	0.016
0.8	-0.1	0.265	0.092	0.026	0.236	0.081	0.016	0.255	0.092	0.022
0.8	0	0.352	0.103	0.051	0.325	0.085	0.034	0.331	0.092	0.044
0.8	0.1	0.426	0.099	0.086	0.413	0.077	0.065	0.416	0.101	0.075
0.8	0.2	0.493	0.100	0.130	0.472	0.080	0.105	0.482	0.102	0.116
0.5	-0.2	0.100	0.088	0.011	0.089	0.093	0.017	0.099	0.111	0.020
0.5	-0.1	0.137	0.086	0.027	0.130	0.090	0.024	0.142	0.102	0.027
0.5	0	0.197	0.081	0.042	0.185	0.088	0.037	0.204	0.102	0.042
0.5	0.1	0.283	0.088	0.074	0.249	0.078	0.059	0.273	0.096	0.070
0.5	0.2	0.373	0.090	0.109	0.329	0.081	0.099	0.352	0.090	0.111
0.2	-0.2	0.033	0.087	0.018	0.041	0.094	0.023	0.047	0.103	0.031
0.2	-0.1	0.060	0.090	0.026	0.065	0.090	0.027	0.066	0.104	0.037
0.2	0	0.084	0.088	0.043	0.082	0.091	0.045	0.093	0.106	0.046
0.2	0.1	0.118	0.084	0.058	0.114	0.093	0.055	0.124	0.108	0.068
0.2	0.2	0.158	0.082	0.086	0.160	0.088	0.074	0.173	0.105	0.088
0	0	0.034	0.093	0.037	0.046	0.093	0.043	0.048	0.106	0.043

Table 7: simulation results, with different subjects per session