

November 2021

“Testing Random Assignment To Peer Groups”

Koen Jochmans

TESTING RANDOM ASSIGNMENT TO PEER GROUPS

KOEN JOCHMANS*

TOULOUSE SCHOOL OF ECONOMICS, UNIVERSITY OF TOULOUSE CAPITOLE

This version: November 10, 2021

Abstract

Identification of peer effects is complicated by the fact that the individuals under study may self-select their peers. Random assignment to peer groups has proven useful to sidestep such a concern. In the absence of a formal randomization mechanism it needs to be argued that assignment is ‘as good as’ random. This paper introduces a simple yet powerful test to do so. We provide theoretical results for this test. As a by-product we equally obtain such results for an approach popularized by [Guryan, Kroft and Notowidigdo \(2009\)](#). These results help to explain why this approach suffers from low power, as has been observed elsewhere. Our approach can equally be used to test for the presence of peer effects in the linear-in-means model without modification.

Keywords: asymptotic power, bias, fixed effects, peer effects, random assignment, test.

JEL classification: C12, C21.

Introduction

A fundamental issue when trying to infer peer effects is the concern that the individuals under study, at least partially, self-select their reference group. Exploiting the random assignment of individuals to peer groups has proven to be a fruitful way forward. [Sacerdote](#)

*Address: Toulouse School of Economics, 1 esplanade de l’Université, 31080 Toulouse, France. E-mail: koen.jochmans@tse-fr.eu.

Support from the European Research Council through grant n° 715787 (MiMo), and from the French Government and the ANR under the Investissements d’Avenir program, grant ANR-17-EURE-0010 is gratefully acknowledged. Comments by Ying Zeng were much appreciated.

The Stata command `rassign` implements the test developed here and can be installed from within Stata by typing `ssc install rassign` in the command window. I am most grateful to Vincenzo Verardi for help in the development of this command.

(2001) and Zimmerman (2003) estimate peer effects in college achievement by making use of the (conditional) random assignment of students to roommates. Katz, Kling and Liebman (2001) and Duflo and Saez (2003) are other early examples that use such exogenous variation in other settings.

In many studies on peer effects there is no formal randomization mechanism. In others the randomization is done at a higher level than under the experimental ideal. Examples of the former situation are in the work of Bandiera, Barankay and Rasul (2009) and Mas and Moretti (2009), both of which concern workers being assigned to teams or shifts. An example of the latter is Project STAR, where students appear to have been randomly assigned only to classes of a certain size, and not to classrooms themselves; Sojourner (2013) provides a detailed discussion on this. In such settings more work is needed to convincingly argue that the assignment of peers is ‘as good as random’.

Sacerdote (2001) pioneered a regression-based approach to test for random assignment. Guryan, Kroft and Notowidigdo (2009) pointed out that this test favors alternatives where there is negative assortative matching between peers, and suggested a modification.¹ Their proposal has been used frequently—Carrell, Fullerton and West (2009), Sojourner (2013), and Lu and Anderson (2015) are examples—but it has not been subject to theoretical investigation. The limited simulation evidence available suggests that it is size correct but has low power (Stevenson, 2015). Thus, the test would have difficulty in detecting violations of the null of random assignment.

In this paper we propose an alternative adjustment to the test of Sacerdote (2001), and study its properties under the null and under various local alternatives. The approach is based on a bias calculation and is straightforward to implement (a Stata implementation is also available). It is related in spirit to calculations in Angrist (2014) and Caeyers and

¹The intuition given in Guryan, Kroft and Notowidigdo (2009) and repeated elsewhere in the literature (Caeyers and Fafchamps, 2020) is that individuals cannot be their own peers. While this argument explains why the test favors negative alternatives it does not explain the cause of the size distortion. In fact, minor modifications to the proof of (1.2) below show that size distortion would also be present when individuals can be their own peers. Furthermore, in such a case the test will tend to favor alternatives where assortative matching is positive. In all cases, the cause of the (asymptotic) size distortion is the presence of fixed effects.

[Fafchamps \(2020\)](#) in a specific case but formalizes, operationalizes, and extends it in various directions; more detail on this is given below. The test allows both peer groups and urns from which peers are drawn to be of the same or of different sizes, accommodates designs in which peer groups need not be mutually exclusive, and is robust to heteroskedasticity of arbitrary form. Because assignment is usually random only conditional on allocation to urns, our test procedure, like [Sacerdote's \(2001\)](#), controls for fixed effects at the urn level. A straightforward modification to the test that allows to control for additional covariates is also presented.

An important remark is that the null model underlying [Sacerdote's \(2001\)](#) approach is formally equivalent to a linear-in-means model of social interactions in which all coefficients involving peer effects are equal to zero. Consequently, our test can equally be applied to test for the presence of peer effects in this workhorse model, and this without modification. This is a useful observation because the test does not require the usual conditions for identification in such settings (see, e.g., [Bramoullé, Djebbari and Fortin 2019](#) for a recent overview of this literature). Furthermore, identification is much easier to establish once such effects can be ruled out.

The derivations underlying our test allow to establish formal results for the influential test of [Guryan, Kroft and Notowidigdo \(2009\)](#). First, we confirm that this test is indeed size correct. Furthermore, their proposal corresponds to an alternative (and implicit) way of performing the bias correction that is inherent in our own procedure, at least when either an urn-level homoskedasticity assumption is satisfied, or when peer groups are mutually exclusive. This alternative approach is only implementable when there is variation in urn size, however. Second, we provide an asymptotic representation that helps to explain the low power that has been observed for the test of [Guryan, Kroft and Notowidigdo \(2009\)](#). We illustrate the power loss through theoretical power calculations and show that the test can have trivial power against a wide range of alternatives. In all cases considered our test is more powerful than theirs, and considerably so. A comparison with other proposals available in the literature, including permutation tests and sample splitting, is also made below.

1 Testing random assignment

Consider a setting where we observe stratified data on r independent urns containing, respectively, n_1, \dots, n_r individuals. Within each urn individuals are assigned to peer groups.

The assignment of peers in urn g is recorded in the $n_g \times n_g$ matrix

$$(\mathbf{A}_g)_{i,j} := \begin{cases} 1 & \text{if } i \text{ and } j \text{ are peers} \\ 0 & \text{if they are not} \end{cases};$$

as individuals cannot be their own peer matrix \mathbf{A}_g has only zeros on its main diagonal. The number of peers of individual i is $m_g(i) := \sum_{j=1}^{n_g} (\mathbf{A}_g)_{i,j}$. We assume that each individual has at least one peer but do not otherwise restrict peer groups; they may be of different sizes and are allowed to overlap. The goal is to test whether individuals are randomly assigned to their respective peer groups. Clearly, while we maintain the terminology of urns and peers for simplicity, this setup covers general adjacency matrices and, therefore, arbitrary network structures.²

Let $x_{g,i}$ be an observable characteristic of individual i in urn g . Sacerdote (2001) noted that, under random assignment, $x_{g,i}$ will be uncorrelated with $x_{g,j}$ for all $j \in [i]$, where $[i] := \{j : (\mathbf{A}_g)_{i,j} = 1\}$ is the set of i 's peers. Letting $\bar{x}_{g,[i]} := m_g(i)^{-1} \sum_{j=1}^{n_g} (\mathbf{A}_g)_{i,j} x_{g,j}$, the average value of the characteristic among i 's peers, he then proceeded by testing whether the slope coefficient in a within-group regression of $x_{g,i}$ on $\bar{x}_{g,[i]}$ is statistically different from zero. The within-group estimator controls for fixed effects at the urn level. This is important as, even if assignment is randomized within urns, individuals might be assigned to an urn based on other attributes. In the data of Sacerdote (2001), for example, students are randomly assigned to rooms conditionally on gender and their answers to a set of survey questions. If peer assignment within urns is presumed to only be random conditional on a set of additional covariates $\mathbf{w}_{g,i}$, say, they can equally be controlled for by including them as additional regressors.

²Everything to follow can be modified to deal with situations where the adjacency matrices $\mathbf{A}_1, \dots, \mathbf{A}_r$ are asymmetric, have non-binary entries, and have a non-zero main diagonal. To maintain focus we do not pursue the most general case here.

1.1 Bias calculation

As observed by [Guryan, Kroft and Notowidigdo \(2009\)](#), the test just described will typically not be size correct. To see the problem, and a path forward, we start by a bias calculation. For now we ignore any additional covariates $\mathbf{w}_{g,i}$ and thus consider a fixed-effect regression of $x_{g,i}$ on $\bar{x}_{g,[i]}$. The within-group estimator, $\hat{\rho}$, is defined as the solution to the normal equation

$$\sum_{g=1}^r \sum_{i=1}^{n_g} \bar{x}_{g,[i]} (\tilde{x}_{g,i} - \hat{\rho} \tilde{\bar{x}}_{g,[i]}) = 0, \quad (1.1)$$

where $\tilde{x}_{g,i}$ and $\tilde{\bar{x}}_{g,[i]}$ are deviations of, respectively, $x_{g,i}$ and $\bar{x}_{g,[i]}$ from their within-urn mean. A calculation given in the Appendix shows that the normal equation is biased. Moreover,

$$\mathbb{E}_0 \left(\sum_{g=1}^r \sum_{i=1}^{n_g} \bar{x}_{g,[i]} \tilde{x}_{g,i} \right) = - \sum_{g=1}^r \sigma_g^2, \quad (1.2)$$

where the subscript on the expectations operator indicates that the expectation is taken under the null of random assignment, and we have assumed that $\mathbb{E}_0((x_{g,i} - \mathbb{E}_0(x_{g,i}))^2) =: \sigma_g^2$ does not vary across individuals. This urn-level homoskedasticity assumption can be dispensed with and we do so below. Furthermore, it will turn out that, when peer groups are mutually exclusive, the test derived under this homoskedasticity assumption is, in fact, robust to heteroskedasticity.

Equation (1.2) implies that the within-group estimator is inconsistent under asymptotics where the number of urns grows large but their size is held fixed. In the Appendix we show that (under the null)

$$\text{plim}_{r \rightarrow \infty} \hat{\rho} = - \frac{\lim_{r \rightarrow \infty} \frac{1}{r} \sum_{g=1}^r \sigma_g^2}{\lim_{r \rightarrow \infty} \frac{1}{r} \sum_{g=1}^r \sigma_g^2 \mathbb{E}_0 \left(\sum_{i=1}^{n_g} \frac{1}{m_g(i)} - \frac{1}{n_g} \sum_{i=1}^{n_g} \sum_{j=1}^{n_g} \frac{m_g(i \cap j)}{m_g(i) m_g(j)} \right)}, \quad (1.3)$$

where $m_g(i \cap j) := \sum_{k=1}^{n_g} (\mathbf{A}_g)_{i,k} (\mathbf{A}_g)_{k,j}$ is the number of peers that individuals i and j have in common. The probability limit is always negative. All else equal its magnitude is decreasing in urn sizes and increasing in the degree of overlap between peer groups. When peer groups do not overlap it is also increasing in the size of the peer groups. Furthermore, in the special case where all urns are of size n and are partitioned into peer groups so that

each individual has m peers,

$$\text{plim}_{r \rightarrow \infty} \hat{\rho} = -\frac{m}{n-m},$$

which no longer depends on the urn variances. This last expression co-incides with the one reported in Proposition 1 of [Caeyers and Fafchamps \(2020\)](#).

The implication of the inconsistency is that the regression-based test will be biased toward negative alternatives and that its size will tend to one as the number of urns grows large.

1.2 A corrected test

The bias calculated in (1.2) is surprisingly simple and suggests a natural adjustment to the proposal of [Sacerdote \(2001\)](#). Observe that an unbiased estimator of σ_g^2 (under the null) is

$$\frac{1}{n_g - 1} \sum_{i=1}^{n_g} x_{g,i} \tilde{x}_{g,i}.$$

Therefore, the re-centered covariance

$$q_r^{\text{HO}} := \sum_{g=1}^r \sum_{i=1}^{n_g} \bar{x}_{g,[i]} \tilde{x}_{g,i} + \sum_{g=1}^r \frac{1}{n_g - 1} \sum_{i=1}^{n_g} x_{g,i} \tilde{x}_{g,i} = \sum_{g=1}^r \sum_{i=1}^{n_g} \tilde{x}_{g,i} \left(\bar{x}_{g,[i]} + \frac{x_{g,i}}{n_g - 1} \right)$$

will be exactly unbiased under random assignment. An estimator of the standard deviation of q_r^{HO} is a conventional standard error that clusters observations at the urn level. It equals

$$s_r^{\text{HO}} := \sqrt{\sum_{g=1}^r \left(\sum_{i=1}^{n_g} \tilde{x}_{g,i} \left(\bar{x}_{g,[i]} + \frac{x_{g,i}}{n_g - 1} \right) \right)^2}.$$

Hence, an adjusted test statistic is $t_r^{\text{HO}} := q_r^{\text{HO}}/s_r^{\text{HO}}$. Note that the entire construction of this statistic is based on calculations under the null. As such it is in the spirit of a Lagrange-multiplier test.

Theorem 1 states the asymptotic behavior of the statistic t_r^{HO} under the null and under alternatives where $\mathbb{E}(q_r^{\text{HO}}) = b_r$ for a sequence of constants $b_r = O(\sqrt{r})$. In the theorem, we let v_g denote the variance of $\sum_{i=1}^{n_g} \tilde{x}_{g,i} (\bar{x}_{g,[i]} + x_{g,i}/n_g - 1)$.

Theorem 1. *Let $\mathbb{P}(n_g > 2) > 0$. If $\max_{g,i} \mathbb{E}(x_{g,i}^8) = O(1)$ and $\max_{g,i} (\text{var}(x_{g,i}^2))^{-1} = O(1)$ and $\mathbb{P}(v_g > 0) > 0$, then*

$$t_r^{\text{HO}} - \frac{b_r}{s_r^{\text{HO}}} \xrightarrow{d} N(0, 1),$$

as $r \rightarrow \infty$.

It is easy to verify that urns of size two would not contribute to the test statistic and so can be dropped. Hence the need for the first condition in the theorem. The second condition contains standard moment requirements. The third condition, finally, ensures non-degeneracy of the test statistic. The prime case where this requirement fails is the case where, in each urn, all individuals are peers of each other, i.e, in the complete-network setting. Of course, in the context of testing random assignment, such a design is of little interest.

An implication of the theorem is that, for any $\alpha \in (0, 1)$,

$$\lim_{r \rightarrow \infty} \mathbb{P}_0(t_r^{\text{HO}} > z_{1-\alpha}) = \alpha,$$

where z_α is the α -quantile of the standard-normal distribution. One-sided and two-sided tests then follow in the usual manner. The theorem also implies that the test is consistent against any alternative for which b_r does not grow slower than \sqrt{r} . Several such deviations, along with asymptotic power calculations, are considered in the Appendix.

The probability limit in (1.3) is smaller (in magnitude) for urns of larger size. This may suggest that in settings where peers are drawn from large urns, ignoring the bias issue in the test of [Sacerdote \(2001\)](#) is inconsequential ([Guryan, Kroft and Notowidigdo, 2009](#)). Such reasoning ignores the fact that the standard deviation of the within-group estimator, too, is decreasing in urn sizes. The conclusion, then, in line with results in the panel data literature (e.g., [Hahn and Kuersteiner 2002](#)), is that the bias will only be ignorable for testing purposes when the size of the urns is substantially larger than the number of urns.

2 Connections to the literature

[Guryan, Kroft and Notowidigdo \(2009\)](#) [Guryan, Kroft and Notowidigdo \(2009\)](#)

proposed to augment the within-group regression of [Sacerdote \(2001\)](#) by including the leave-one-out average

$$\frac{1}{n_g - 1} \sum_{j \neq i} x_{g,j} = \frac{n_g}{n_g - 1} \left(\frac{1}{n_g} \sum_{j=1}^{n_g} x_{g,j} - \frac{x_{g,i}}{n_g} \right) = \frac{n_g}{n_g - 1} \left(\bar{x}_g - \frac{x_{g,i}}{n_g} \right)$$

as an additional regressor. The within-group transformation sweeps out all terms that do not vary within urns, and so the approach is equivalent to a within-group regression of $x_{g,i}$ on $\bar{x}_{g,[i]}$ and $x_{g,i}/(n_g - 1)$. This highlights why variation in urn size is required for this approach to be implementable. When n_g does not vary across urns this regression will yield a perfect fit that satisfies the null whether or not peer assignment is random. [Guryan, Kroft and Notowidigdo \(2009\)](#) offer an intuition of why their strategy yields size control and provide supporting simulations. However, a theoretical analysis of the test is, to our knowledge, not available.

Calculations summarized in the Appendix reveal that the approach of [Guryan, Kroft and Notowidigdo \(2009\)](#) tests whether

$$\sum_{g=1}^r \sum_{i=1}^{n_g} \tilde{x}_{g,i} \left(\bar{x}_{g,[i]} + \frac{x_{g,i}}{n_g - 1} \right) \left(1 - \frac{\delta}{n_g - 1} \right) + o_p(\sqrt{r}), \quad (2.4)$$

is statistically different from zero. Here,

$$\delta := \frac{\lim_{r \rightarrow \infty} \frac{1}{r} \sum_{g=1}^r \sigma_g^2}{\lim_{r \rightarrow \infty} \frac{1}{r} \sum_{g=1}^r \sigma_g^2 \mathbb{E}_0 \left(\frac{1}{n_g - 1} \right)},$$

is the probability limit of the slope coefficient of a within-group regression of $x_{g,i}$ on $x_{g,i}/(n_g - 1)$, under the null. The summand in the leading term in (2.4) is equal to the summand in q_r^{HO} , up to a scale factor that varies at the urn level. This factor is bounded and so, by virtue of [Theorem 1](#), we conclude that the test will indeed exhibit correct size in large samples.

The limited simulation evidence available suggests that the test of [Guryan, Kroft and Notowidigdo \(2009\)](#) may suffer from low power; see [Stevenson \(2015\)](#) and also the extended version of her analysis in the Appendix. Because the approach requires variation in urn sizes one may expect the test to be particularly underpowered when such variation is

limited (Stevenson 2015, Caeyers and Fafchamps 2020). While this is true, as evidenced by (2.4), low power also arises from a different source. Equation (2.4) is again useful here. The weights $1 - \delta/(n_g - 1)$ have mean zero, implying that they take on both positive and negative values. Hence, bias terms will tend to cancel each other out.

To see this it suffices to consider a design where urns are of size \bar{n}_1 with probability $(1 - p_n)$ and of size \bar{n}_2 with probability p_n , where $\bar{n}_1 < \bar{n}_2$. The non-centrality parameter in the limit distribution of the test statistic of Guryan, Kroft and Notowidigdo (2009) can be shown to equal

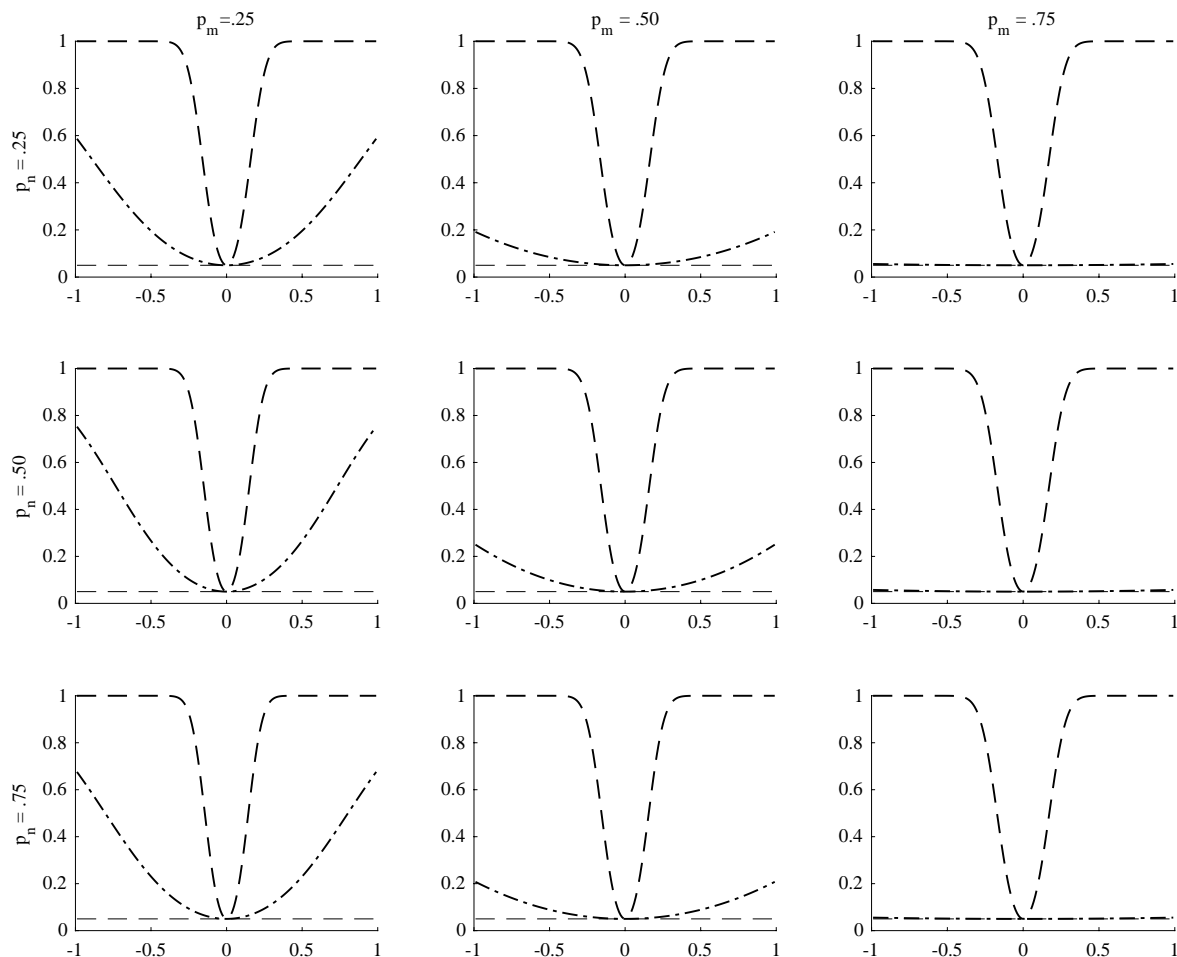
$$\mu^* := \sqrt{p_n(1 - p_n)} \frac{b(\bar{n}_2) - b(\bar{n}_1)}{\sqrt{v(\bar{n}_1)p_n + v(\bar{n}_2)(1 - p_n)}}, \quad (2.5)$$

where $b(n)$ and $v(n)$ are the bias and variance of $\sum_{i=1}^{n_g} \tilde{x}_{g,i} (\bar{x}_{g,[i]} + x_{g,i}/(n_g - 1))$ conditional on $n_g = n$. This equation confirms that $\mu^* \rightarrow 0$ as $p_n(1 - p_n) \rightarrow 0$ and formalizes the notion that the test will tend to have low power when variation in urn sizes is small. The formula also shows that the test will have trivial asymptotic power when $b(\bar{n}_1) - b(\bar{n}_2) = 0$, i.e., in designs where the bias contributions coming from the different urn sizes cancel each other out.

We confirm these findings in Figure 1; additional power calculations are provided in the Appendix. We consider designs where each of 25 urns contains six individuals with probability p_n and four individuals with probability $1 - p_n$. Within urns of size four, each individual is assigned one peer at random while in the larger urn peer groups are of size three with probability p_m and of size two with probability $1 - p_m$. Figure 1 plots (theoretical) power (as a function of ρ) against alternatives of the form $x_{g,i} = \rho \bar{x}_{g,[i]} + \varepsilon_{g,i}$ with $\varepsilon_{g,i} \sim \text{independent}(\alpha_g, \sigma^2)$. The plots in Figure 1 are arranged so that p_n increases when going down rows and p_m increases when moving through columns. Dashed curves refer to our test. Dashed-dotted curves represent the test of Guryan, Kroft and Notowidigdo (2009). Both tests are two-sided at the 5% level, and we included a dashed horizontal line in the plots to mark the size.

Figure 1 shows high power for our test across all designs. The test of Guryan, Kroft and Notowidigdo (2009) is less powerful against all alternatives, and substantially so. There

Figure 1: Power analysis



Power for our test (dashed line) and for the test of [Guryan, Kroft and Notowidigdo \(2009\)](#) (dashed-dotted line) in a design with two possible urns sizes (4 and 6) and two possible peer-group sizes (2 and 3). $p_n := \mathbb{P}(n_g = 6)$ and $p_m := \mathbb{P}(m_g(i) = 2 | n_g = 6)$. A horizontal dashed line indicates the size of the test. Plots are based on theoretical calculations and are for 25 urns.

is a reduction in its power when p_n moves away from .50 (i.e., across rows); recall that power converges to size as p_n approaches either zero or one. For the values considered here, this effect is small relative to the impact of changing p_m , with power initially going down considerably when p_m moves from .25 to .50, and afterwards essentially flattening out completely when $p_m = .75$. This is a reflection of the numerator in μ^* getting close to zero; the bias in urns of size four cancels out with the bias in urns of size six. As μ^* is

multiplicative in ρ these changes are uniform on $(-1, 1)$.

Guryan, Kroft and Notowidigdo (2009) also describe an alternative randomization test (see, e.g., Lehmann and Romano 2006, Chapter 15, for a general treatment of such tests) that is based on the sampling distribution of the (uncorrected) within-group estimator obtained from randomly re-assigning individuals to peer groups within each urn. In general, randomization tests have many attractive properties. However, in the current context, the proposed test will fail (even in large samples) when errors are heteroskedastic, for example. This is so because, as Equation (3.6) in the next section implies, the probability limit of the within-group estimator (under the null) is not invariant to random re-assignment of individuals to peer groups in that case.

It may be of interest to note that our approach of testing whether $q_r^{\text{HO}} = 0$ can be cast into a permutation test by appealing to the developments in Hemerik, Goeman and Finos (2020). Moreover, it follows from their Corollary 2 that a standard sign-flipping test applied across the urns will be asymptotically size correct. When $x_{g,i}$ is continuous and the distribution of $\sum_{i=1}^{n_g} \tilde{x}_{g,i}(\bar{x}_{g,[i]} + x_{g,i}/(n_g - 1))$ is symmetric around zero, the test will further be size correct for a fixed number of urns, as per Proposition 1 in Hemerik, Goeman and Finos (2020).

Caeyers and Fafchamps (2020) In designs where all urns are of the same size, n , and each urn is partitioned into equally-sized groups of size $m + 1$ (so that each individual has m peers) we had that

$$\text{plim}_{r \rightarrow \infty} \hat{\rho} = -\frac{m}{n - m} =: \rho_0.$$

This expression was equally obtained by Caeyers and Fafchamps (2020), albeit under a full homoskedasticity assumption (see also Booser and Cacciola 2001 and Angrist 2014 for analogous calculations for the setting without fixed effects). To test for random assignment Caeyers and Fafchamps (2020) suggest to perform the test of Sacerdote (2001), only with the dependent variable $x_{g,i}$ replaced by $x_{g,i} - \rho_0 \bar{x}_{g,[i]}$. Of course, this strategy is numerically identical to constructing the bias-corrected estimator $\hat{\rho} - \rho_0$, and performing a standard two-sided t-test on it. Whichever implementation is chosen, it is important to stress that

this approach is valid because, in this simple case, the probability limit, ρ_0 , does not depend on any unknown parameters and so need not be estimated from the data.

It is immediately clear that this idea can be generalized to our setup, making use of the expression for the probability limit we obtained in (1.3).³ However, this probability limit now has to be estimated. This implies that the usual standard errors reported with a within-group regression are invalid and have to be adjusted in order for the test so constructed to be size correct.

The chief difference between our proposal and the alternative possibility just described is that we base our test directly on a recentered normal equation of the within-group estimator and not on the within-group estimator itself. While the bias in the normal equation is very simple and independent of the design, the probability limit of the within-urn estimator depends in a complicated manner on the size and overlap between the different peer groups.

Stevenson (2015) Next we mention the suggestion of [Stevenson \(2015\)](#), which is based on data splitting. Although its properties have not been established theoretically, the subsampling scheme she proposed circumvents bias under the null, at least when peer groups are mutually exclusive, and so should lead to size correct inference in this case (under regularity conditions). Of course, the scheme is also computationally substantially more demanding than the bias-adjustment proposal made here.

³After the current paper was circulated, [Caeyers and Fafchamps \(2020\)](#) included an extension of their probability-limit calculation to allow for the size of urns and peer groups to be heterogenous. This is a special case of (1.3). Moreover, under random sampling of urns and homoskedasticity our formula simplifies to

$$\frac{1}{\mathbb{E}_0 \left(\sum_{i=1}^{n_g} \frac{1}{m_g(i)} - 1 \right)}$$

in this case. The statement in [Caeyers and Fafchamps \(2020, Proposition 2\)](#), in contrast, is considerably more complicated and seems to implicitly assume that the frequency with which each combination of urn and peer-group size appears in the sample is fixed and independent of the sample size. It is not obvious that their expression can be reduced to the above equation even in this case. Their proof discusses only the example where urns can be of one of two sizes. A proof of the general statement in their Proposition 2 is not provided.

An alternative approach Finally, we remark that an alternative procedure to testing for random assignment that has been employed (see, e.g., Wang 2009 and Chetty, Friedman, Hilger, Saez, Schanzenbach, and Yagan 2011) is to first regress $x_{g,i}$ on a set of urn dummies and a set of peer-group dummies, and next test that the latter set are all zero by means of a standard F-test. Such a test, however, does not yield size control, even in the absence of urn fixed effects (in which case, recall, an adjustment to the approach of Sacerdote (2001) is not needed to yield asymptotic size control). Indeed, this testing problem is a generalization of the one-way analysis of variance problem, which has a long history in statistics; see Akritas and Papadatos (2004) for example.

3 Extensions

3.1 Heteroskedasticity

So far we have worked under an assumption of urn-level homoskedasticity. We now drop this restriction and allow that $\sigma_{g,i}^2 := \mathbb{E}_0((x_{g,i} - \mathbb{E}_0(x_{g,i}))^2)$ varies both between and within urns in an arbitrary way.

First, calculations analogous to those that gave rise to (1.2) show that, now,

$$\mathbb{E}_0 \left(\sum_{g=1}^r \sum_{i=1}^{n_g} \bar{x}_{g,[i]} \tilde{x}_{g,i} \right) = - \sum_{g=1}^r \mathbb{E}_0 \left(\frac{1}{n_g} \sum_{i=1}^{n_g} \frac{1}{m_g(i)} \sum_{j=1}^{n_g} (\mathbf{A}_g)_{i,j} \sigma_{g,j}^2 \right). \quad (3.6)$$

Hence, the contribution of each urn to the bias equals (minus) the expected within-urn mean of peer-group averaged variances.

Appealing to a result of Hartley, Rao and Kiefer (1969), we show in the Appendix that an unbiased estimator of the bias in (3.6) is

$$- \sum_{g=1}^r \sum_{i=1}^{n_g} \omega_{g,i} x_{g,i} \tilde{x}_{g,i}, \quad \omega_{g,i} := \frac{1}{n_g - 2} \left(\sum_{i' \in [i]} \frac{1}{m_g(i')} - \frac{1}{n_g - 1} \right),$$

which is again well-defined for all urns of size $n_g > 2$. Hence, a modification of q_r^{HO} that is

robust to heteroskedasticity of arbitrary form is given by

$$q_r^{\text{HC}} := \sum_{g=1}^r \sum_{i=1}^{n_g} \tilde{x}_{g,i} (\bar{x}_{g,[i]} + \omega_{g,i} x_{g,i}), \quad (3.7)$$

which satisfies $\mathbb{E}_0(q_r^{\text{HC}}) = 0$. It differs from q_r^{HO} only in that the weight $(n_g - 1)^{-1}$ is replaced by $\omega_{g,i}$, which varies at the individual level. Construction of $\omega_{g,i}$ is nonetheless immediate from \mathbf{A}_g .

Observe that, in the important special case where peer groups do not overlap we have $m_g(i') = m_g(i)$ for all $i' \in [i]$, and so

$$\omega_{g,i} = \frac{1}{n_g - 1}.$$

This is the weight we used to construct our test statistic under homoskedasticity. It thus follows that t_r^{HO} is robust to heteroskedasticity in this case.

The standard deviation of q_r^{HC} can be estimated by

$$s_r^{\text{HC}} := \sqrt{\sum_{g=1}^r \left(\sum_{i=1}^{n_g} \tilde{x}_{g,i} (\bar{x}_{g,[i]} + \omega_{g,i} x_{g,i}) \right)^2}.$$

A modified version of our test statistic that remains size correct under heteroskedasticity of arbitrary form also when peer groups overlap is $t_r^{\text{HC}} := q_r^{\text{HC}}/s_r^{\text{HC}}$. This statistic is asymptotically normal under the same conditions as before. In the following theorem, $b_r := \mathbb{E}(q_r^{\text{HC}}) = O(\sqrt{r})$ and $v_g := \mathbb{E}((\sum_{i=1}^{n_g} \tilde{x}_{g,i} (\bar{x}_{g,[i]} + \omega_{g,i} x_{g,i}))^2)$.

Theorem 2. *Let $\mathbb{P}(n_g > 2) > 0$. If $\max_{g,i} \mathbb{E}(x_{g,i}^8) = O(1)$ and $\max_{g,i} (\text{var}(x_{g,i}))^{-1} = O(1)$ and $\mathbb{P}(v_g > 0) > 0$, then*

$$t_r^{\text{HC}} - \frac{b_r}{s_r^{\text{HC}}} \xrightarrow{d} N(0, 1),$$

as $r \rightarrow \infty$.

3.2 Controlling for covariates

There may be situations where, in addition to urn fixed effects, it is desirable to control for other variables that vary at the individual level, $\mathbf{w}_{g,i}$. This would be needed when

randomization is assumed to take place within urns only conditional on these variables. A intuitive regression-based solution would be to first partial-out $\mathbf{w}_{g,i}$ from $x_{g,i}$ and $\bar{x}_{g,[i]}$ and then proceed in constructing our test statistic as before. We next show that, under regularity conditions, this approach is justified.

Let $\dot{x}_{g,i}$ denote the residual from an ordinary least-squares regression of $x_{g,i}$ on urn dummies and the vector of covariates $\mathbf{w}_{g,i}$. Then the modified test statistic takes the form

$$\hat{t}_r^{\text{HO}} := \frac{\hat{q}_r^{\text{HO}}}{\hat{s}_r^{\text{HO}}}$$

for

$$\hat{q}_r^{\text{HO}} := \sum_{g=1}^r \sum_{i=1}^{n_g} \dot{x}_{g,i} \left(\bar{x}_{g,[i]} + \frac{x_{g,i}}{n_g - 1} \right), \quad \hat{s}_r^{\text{HO}} := \sqrt{\sum_{g=1}^r \left(\sum_{i=1}^{n_g} \dot{x}_{g,i} \left(\bar{x}_{g,[i]} + \frac{x_{g,i}}{n_g - 1} \right) \right)^2}.$$

The statistic t_r^{HC} can be modified in the same way.

To state conditions under which Theorem 1 generalizes to partialling-out covariates we need

$$\check{x}_{g,i} := x_{g,i} - \mathbf{w}'_{g,i} \left(\sum_{g=1}^r \sum_{i'=1}^{n_g} \mathbb{E}(\mathbf{w}_{g,i'} \mathbf{w}'_{g,i'}) \right)^{-1} \left(\sum_{g=1}^r \sum_{i'=1}^{n_g} \mathbb{E}(\mathbf{w}_{g,i'} x_{g,i'}) \right).$$

This is the deviation of $x_{g,i}$ from its population linear projection on $\mathbf{w}_{g,i}$ (and no fixed effects).

The following theorem provides the result. Here, $\|\cdot\|$ refers to the Euclidean norm, b_r is once more suitably re-defined to be the bias in \hat{q}_r^{HO} under Pitman drifts towards the null hypothesis, and v_g again denotes the variance of the term that urn g adds to the test statistic.

Theorem 3. *Let $\mathbb{P}(n_g > 2) > 0$. If $\max_{g,i} \mathbb{E}(\check{x}_{g,i}^8) = O(1)$ and $\max_{g,i} (\text{var}(\check{x}_{g,i}^2))^{-1} = O(1)$ and $\mathbb{P}(v_g > 0) > 0$, then*

$$\hat{t}_r^{\text{HO}} - \frac{b_r}{\hat{s}_r^{\text{HO}}} \xrightarrow{d} N(0, 1),$$

as $r \rightarrow \infty$, provided that $\mathbb{E}(\check{x}_{g,i} | \mathbf{w}_{g,1}, \dots, \mathbf{w}_{g,n_g}) = \alpha_g$ for urn-specific constants $\alpha_1, \dots, \alpha_r$, that $\max_{g,i} \mathbb{E}(\|\mathbf{w}_{g,i}\|^4) = O(1)$ and that the matrix $\lim_{r \rightarrow \infty} r^{-1} \sum_{g=1}^r \mathbb{E}(\tilde{\mathbf{w}}_{g,i} \tilde{\mathbf{w}}'_{g,i})$ has maximal rank.

The conditions in this result are intuitive. First, the moment conditions on $x_{g,i}$ in Theorem 1 are replaced by corresponding conditions on $\check{x}_{g,i}$. Next, the mean-independence assumption is a requirement of strict exogeneity on $\mathbf{w}_{g,i}$. Finally, the conditions on the covariates are needed to ensure that the residuals from the auxiliary least-squares regression converge to their population counterparts.

4 Empirical illustration

Guryan, Kroft and Notowidigdo (2009) used the random assignment of golf players to playing partners in tournaments to estimate peer effects. Their data span the 2002, 2005, and 2006 seasons of the Professional Golfer’s Association (PGA) and cover 81 tournaments. We refer to Guryan, Kroft and Notowidigdo (2009) for a detailed description of the data. Here we only note the facts that are of direct relevance to our analysis. Players in the PGA are, at any point in time, assigned to one of four categories (cat 1, cat 1a, cat 2, and cat 3). At the start of each tournament, within these four categories, playing partners are assigned to groups of three golfers. These (mutually exclusive) peer groups play together for the first two rounds of the tournament. The analysis is limited to the first round. Conditional on the set of players who enter a tournament, the assignment is random within categories. Random assignment is tested by looking at the (corrected) within-group correlation between a measure of a golfer’s ability and the average ability of his playing partners.

The chief measure of ability used to do this is an estimate of the number of strokes more than 72 (i.e., above par) that a golfer typically takes in a round, on an average course, that is used for PGA tournaments. The more negative this number the better the player. Table 1 contains descriptive statistics for this variable, stratified by the four player categories. It shows that, broadly, average ability is higher in lower numbered categories, and that there remains substantial variation in this measure even conditional on category. To get a sense of urn sizes in these data the table also provides descriptive statistics of the number of players by tournament-by-category. These are based on a total of 8,791 observations in stead of the total of 8,801 observations as 10 observations concern urns of a size less

than three; recall that such urns do not contain any information for our purposes. We also included the same descriptive statistics for the weights $(n_g - 1)^{-1}$.

Table 1: The PGA data

	n obs	mean	std	min	max
	ability $(x_{g,i})$				
cat 1	3,205	-3.138	0.769	-5.159	1.440
cat 1a	3,436	-2.808	0.740	-4.326	6.732
cat 2	1,503	-2.857	0.894	-4.776	3.275
cat 3	657	-1.662	1.470	-4.776	6.315
	peer ability $(\bar{x}_{g,[i]})$				
cat 1	3,205	-3.132	0.599	-5.081	0.672
cat 1a	3,436	-2.811	0.591	-4.530	3.275
cat 2	1,503	-2.850	0.744	-4.776	3.275
cat 3	657	-1.690	1.270	-4.776	6.315
	urn size (n_g)				
tourn by cat	8,791	39.292	16.869	3	83
	weight $((n_g - 1)^{-1})$				
tourn by cat	8,791	0.037	0.040	0.012	0.500

The test statistics for the default (i.e., uncorrected) regression-based test, our corrected version, and the test where leave-me-out urn means are controlled for are collected in Table 2. The numbers in square brackets below are corresponding (two-sided) p -values. When fully stratifying the data by tournament and category we observe that the default test rejects the null of random assignment and would suggest there to be negative assortative matching between players. The other two tests have large p -values, finding little evidence to contradict the null.

We conclude this illustration by highlighting a caveat to the analysis of these data. Most, if not all, professional golf players participate to multiple tournaments per year and are also active for multiple years. Consequently, many players will appear in multiple urns, albeit with a different value for their ability measure, as this is updated over time. This, of course, induces dependence across urns which is in violation with our working assumption

Table 2: Results for the PGA data (test statistic [p-value])

stratification	default	corrected	control
tourn by cat	-3.957	-0.852	-1.209
	[0.000]	[0.394]	[0.227]

that urns are independent.

References

- Akritas, M. G. and N. Papadatos (2004). Heteroscedastic one-way ANOVA and lack-of-fit tests. *Journal of the American Statistical Association* 99, 368–382.
- Angrist, J. D. (2014). The perils of peer effects. *Labour Economics* 30, 98–108.
- Bandiera, O., I. Barankay, and I. Rasul (2009). Social connections and incentives in the workplace: Evidence from personnel data. *Econometrica* 77, 1047–1094.
- Boozer, M. and S. E. Cacciola (2001). Inside the ‘black box’ of Project STAR: Estimation of peer effects using experimental data. Yale Economic Growth Center Discussion Paper 832.
- Bramoullé, Y., H. Djebbari, and B. Fortin (2019). Peer effects in networks: A survey. Forthcoming in *Annual Review of Economics*.
- Caeyers, B. and M. Fafchamps (2020). Exclusion bias in the estimation of peer effects. NBER Working Paper No. 22565.
- Calhoun, G. (2011). Hypothesis testing in linear regression when k/n is large. *Journal of Econometrics* 165, 163–174.
- Carrell, S. E., R. L. Fullerton, and J. E. West (2009). Does your cohort matter? Measuring peer effects in college achievement. *Journal of Labor Economics* 27, 439–464.
- Chetty, R., J. N. Friedman, N. Hilger, E. Saez, D. W. Schanzenbach, and D. Yagan (2011). How does your kindergarten classroom affect your earnings? Evidence from Project STAR. *Quarterly Journal of Economics* 126, 1593–1660.
- Duflo, E. and E. Saez (2003). The role of information and social interactions in retirement plan decisions: Evidence from a randomized experiment. *Quarterly Journal of Economics* 118, 815–842.

- Guryan, J., D. Kroft, and N. J. Notowidigdo (2009). Peer effects in the workplace: Evidence from random groupings in professional golf tournaments. *American Economic Journal: Applied Economics* 44, 289–302.
- Hahn, J. and G. Kuersteiner (2002). Asymptotically unbiased inference for a dynamic panel model with fixed effects when both n and T are large. *Econometrica* 70, 1639–1657.
- Hartley, H. O., J. N. K. Rao, and G. Kiefer (1969). Variance estimation with one unit per stratum. *Journal of the American Statistical Association* 64, 841–851.
- Hemerik, J., J. J. Goeman, and L. Finos (2020). Robust testing in generalized linear models by sign flipping score contributions. Forthcoming in *Journal of the Royal Statistical Society, Series B*.
- Katz, L., J. Kling, and J. Liebman (2001). Moving to opportunity in Boston: Early results of a randomized mobility study. *Quarterly Journal of Economics* 116, 607–654.
- Lehmann, E. L. and J. P. Romano (2006). *Testing Statistical Hypotheses*. Springer.
- Lu, F. and M. Anderson (2015). Peer effects in microenvironments: The benefits of homogeneous classroom groups. *Journal of Labor Economics* 33, 91–122.
- Mas, A. and E. Moretti (2009). Peers at work. *American Economic Review* 99, 112–145.
- Sacerdote, B. (2001). Peer effects with random assignment: Results for Dartmouth roommates. *Quarterly Journal of Economics* 116, 681–704.
- Sojourner, A. (2013). Identification of peer effects with missing peer data: Evidence from Project STAR. *Economic Journal* 123, 574–605.
- Stevenson, M. (2015). Tests of random assignment to peers in the face of mechanical negative correlation: An evaluation of four techniques. Mimeo.
- Wang, L. C. (2009). Peer effects in the classroom: Evidence from a natural experiment in Malaysia. Mimeo.
- Zimmerman, D. (2003). Peer effects in academic outcomes: Evidence from a natural experiment. *Review of Economics and Statistics* 85, 9–23.