

Making Effects Manifest in Randomized Experiments

Jake Bowers *

July 19, 2010

Experimentalists desire precise estimates of treatment effects and nearly always care about how treatment effects may differ across subgroups. After data collection, concern may focus on random imbalance between treatment groups on substantively important variables. Pursuit of these three goals — enhanced precision, understanding treatment effect heterogeneity, and imbalance adjustment — requires background information about experimental units. For example, one may group similar observations on the basis of such variables and then assign treatment within those blocks. Use of covariates after data have been collected raises extra concerns and requires special justification. For example standard regression tables only approximate the statistical inference that experimentalists desire. The standard linear model may also mislead via extrapolation. After providing some general background about how covariates may, in principle, enable pursuit of precision and statistical adjustment, this paper presents two alternative approaches to covariance adjustment: one using modern matching techniques and another using the linear model — both use randomization as the basis for statistical inference.

1 What is a manifest effect?

A manifest effect is one we can distinguish from zero. Of course, we cannot talk formally about the effects of an experimental treatment as manifest without referring to probability: a scientist asks, “Could this result have occurred merely through chance?” or “If the true effect were zero, what is the chance that we’d observe an effect as large as this?” More formally, for a frequentist, saying a treatment effect is manifest is saying that the statistic we observe casts a great deal of doubt on a hypothesis of no effects. We are most likely to say that some observed effect casts doubt on the null hypothesis of no effect when we have a large sample and/or when noise in the outcome that might otherwise drown out the signal in our study has been well controlled. Fisher reminds us that while randomization alone is sufficient for a valid test of the null hypothesis of no effect, specific features of a given design allow equally valid tests to differ in their ability to make a treatment effect

*Assistant Professor, Dept of Political Science, University of Illinois @ Urbana-Champaign
Corresponding Author Contact Information: 231 Computing Applications Building, 605 E Springfield Ave Champaign IL 61820 —217.333.3881 — jwbowers@illinois.edu. *Acknowledgements:* Many thanks to Jamie Druckman, Don Green, Ben Hansen, Jim Kuklinski, Thomas Leeper, Costas Panagopoulos, and Cara Wong. Parts of this work were funded by NSF Grants SES-0753168 and SES-0753164.

manifest:

With respect to the refinements of technique [uses of covariates in the planning of an experiment], we have seen above that these contribute nothing to the validity of the experiment, and of the test of significance by which we determine its result. They may, however, be important, and even essential, in permitting the phenomenon under test to manifest itself. (Fisher 1935, 24).

Of course, one would prefer a narrow confidence interval to a wide confidence interval even if both excluded the hypothesis of no effects. As a general rule, more information yields more precision of estimation. One may increase information in a design by gathering more observations and/or gathering more data about each observation. This paper considers covariates as a refinement of technique to make treatment effects manifest in randomized studies. I focus first on simple uses of covariates in design, and then offer some ideas about their use in post-treatment adjustment.

What is a covariate? How should we use them in experiments?

A covariate is a piece of background information about an experimental unit — a variable unchanged and unchangeable by the experimental manipulation. Such variables might record the groups across which treatment effects ought to differ according to the theory motivating and addressed by the experimental design (say, men and women ought to react differently to the treatment), or might provide information about the outcomes of the experiment (say, men and women might be expected to have somewhat different outcomes even if reactions to treatment are expected to be the same across both groups). Covariates may be used profitably in experiments either *before* treatment is assigned (by creating subgroups of units within which treatment will be randomized during the recruitment and sample design of a study) and/or *after* treatment has been assigned and administered and outcomes measured (by creating subgroups within which outcomes ought to be homogeneous or adjusting for covariates using linear models).

Common Uses (and Potential Abuses) of Covariates in the Workflow of Randomized Experimentation

Every textbook on the design of experiments is, in essence, a book about the use of covariates in the design and analysis of experiments. This chapter ought not, and cannot, substitute for such sources. For the newcomer to experiments, I here summarize in broad strokes, and with minimal citations, the uses to which covariates may be put in the design and analysis of randomized experiments. After this summary, I offer a perspective on the use of covariates in randomized experimentation which, in fundamental ways, is the same as that found in books such as: Fisher (1935, 1925), Cox (1958), Cochran and Cox (1957), and Cox and Reid (2000). I differ from those previous scholars in hewing more closely and explicitly to: 1) the now well-known potential outcomes framework for causal inference (Neyman 1990; Rubin 1974, 1990; Brady 2008; Sekhon 2008) and 2) randomization as the basis for statistical inference.

Covariates allow precision enhancement

Blocking on background variables before treatment assignment allows the experimenter to create sub-experiments within which the units are particularly similar in their outcomes; adjustment using covariates after data have been collected may also reduce non-treatment related variation in outcomes. In both cases, covariates can reduce noise that might otherwise obscure the effects of the treatment.

Of course, such precision enhancements arrive with some costs: Implementing a blocking plan may be difficult if background information on experimental units is not available before recruitment/arrival at the lab (but see Nickerson 2005); care must be taken to reflect the blocking in the estimation of treatment effects to avoid bias and to take advantage of the precision enhancements offered by the design; and, in principle, analysts can mislead themselves by performing many differently adjusted hypothesis tests until they reject the null of no effects even when the treatment has no effect.

Covariates enable subgroup analyses

When theory implies differences in treatment effects across subgroups, subgroup membership must be recorded and, if at all possible, the experiment ought to be designed to enhance the ability of the analyst to distinguish group-differences in treatment effects. Covariates on subgroups may also be quite useful for post-hoc exploratory analyses designed not to cast doubt on common knowledge but to suggest further avenues for theory.

Covariates allow adjustments for random imbalance

All experiments may display random imbalance. Such baseline differences can arise even if the randomization itself is not suspect: recall that one out of twenty unbiased hypothesis tests will reject the null of no difference at the predetermined error rate of $\alpha = .05$ merely due to chance. An omnibus balance assessment such as that proposed by Hansen and Bowers (2008) is immune from this problem, but any unbiased one-by-one balance assessment will show imbalance in $100\alpha\%$ of the covariates tested. Thus, I call this problem “random imbalance” to emphasize that the imbalance could easily be due to chance and need not cast doubt on the randomization or administration of a study (although discovery of extensive imbalance might suggest scrutiny of the randomization and administration is warranted).

Random imbalance in a well-randomized study on substantively important covariates still may confuse the reader. In the presence of random imbalance, comparisons of treated to controls will contain *both* the effect of the treatment *and* the differences due to the random imbalance. One may attempt to remove the effects of such covariates from the treatment effect by some form of adjustment. For example, one may use the linear regression model as a way to adjust for covariates or one may simply group together observations on the basis of the imbalanced covariate. Adjustment on one or a group of observed covariates may, however, produce now-non-random imbalance on unobserved covariates. And, adjustment raises concerns that estimates of treatment effects may come to depend more on the details of the adjustment method rather than on the randomization and design of the study. Thus, the quandary: either risk known confusions of comparisons or risk unknown confounding and bear the burden of defending and assessing an adjustment method. An

obvious strategy to counter concerns about cherry-picking results or modeling artifacts is to present both adjusted and unadjusted results and to specify adjustment strategies before randomization occurs.

Randomization is the Primary Basis for Statistical Inference in Experiments

A discussion of manifest effects is also a discussion of statistical inference: statistical tests quantify doubt against hypotheses and a manifest effect is evidence which casts great doubt on the null of no effects. On what basis can we justify statistical tests for experiments?

In national surveys we draw random samples. Statistical theory tells us that the mean in the sample is an unbiased estimator of the mean in the population as long as we correctly account for the process by which we drew the sample in our estimation. That is, in a national survey, often (but not always) our *target of inference* is the population from which the sample was drawn and we are *justified* in so inferring by the sampling design.

In other studies we may not know how a sample was drawn (either we have no well-defined population or no knowledge of the sampling process or both). But we may know how our observed outcome was generated: say we know that at the micro-level our outcome was created from discrete events occurring independently of each other in time. In that case, we would be justified in claiming that our population was created via a Poisson process: in essence we have a population generating machine, data generating process, or model of outcomes as the target of our inference.

Now, what about randomized experiments? Although we might want to infer to a model, or even to a population, the strength of experiments is inference to a counter-factual. The primary targets of inference in a randomized experiment are the experimental treatment groups: we infer from one to another. Randomization makes this inference meaningful. But, randomization also can justify the statistical inference as well: the mean in the treatment group is a good estimator for what we would expect to observe if all of the experimental units were treated: the treatment group in a randomized study is a random sample from the finite “population” of the experimental pool.¹

All of the standard textbooks note this fact, but they also point out that estimating causal effects using randomization-based theory can be mathematically inconvenient or computationally intensive and that, thus, using the large-sample sampling theory (and/or Normal distribution models) turns out to provide very good approximations to the randomization-based results most of the time. Since the 1930s, the computational apparatus of randomization-based inference has expanded, as has its theoretical basis and applied reach. In this paper, all of the statistical inference I present will be randomization-based even if most of it also uses large-sample theory: for example, it takes no

¹See Bowers and Panagopoulos (2009) and Rosenbaum (2002*b*, ch. 2) for accessible introductions to randomization inference; a mode of inference developed in different yet compatible ways by Neyman (1990) (as a method of estimating mean differences) and (Fisher 1935, ch. 2) (as a method of testing). In this paper, I follow the Fisher-style approach in which causal effects are inferred from testing hypotheses rather than estimated as points. Both methods (producing an plausible interval for causal effects using a point \pm a range or testing a sequence of hypotheses) often produce identical confidence intervals.

more time to execute a randomization-based test of the null hypothesis of no effect using mean differences than it does using the linear regression model-based approximation.²

Recently Freedman (2008*c,b,a*) reminded us that the target of inference in randomized experiments was the counterfactual and he noted that linear regression and logistic regression were not theoretically justified on this basis. Green (2009) and Schochet (2009) reminded us that often linear regression can be an excellent approximation to the randomization-based difference of means. The need for this exchange is obvious, even if it is echoed in the early textbooks: those early authors moved very quickly to the technicalities of the approximations rather than dwell on the then uncomputable but theoretically justifiable procedures. As experimentation explodes as a methodology in political science, we are seeing many more small experiments, designs where randomization is merely a (possibly weak) instrument, and theories implying very heterogeneous treatment effects. I expect we will find more of these along with many more large studies with fairly Normal looking outcomes where treatment plausibly just shifts the Normal curve of the treated away from the Normal curve of the controls. Rather than hope that the linear-model approximation works well, this paper presents analyses which do not require that approximation and thereby offers others the ability to check the approximation.

2 Strategies for enhancing Precision Before Assignment

... we consider some ways of reducing the effect of uncontrolled variations on the error of the treatment comparisons. The general idea is the common sense one of grouping the units into sets, all the units in a set being as alike as possible, and then assigning the treatments so that each occurs once in each set. All comparisons are then made within sets of similar units. The success of the method in reducing error depends on using general knowledge of the experimental material to make an appropriate grouping of the units into sets (Cox 1958, 23).

We have long known that covariates enhance the precision of estimation to the extent that they predict outcomes. This section aims to make this intuition more concrete in the context of a randomized experiment.

An example by simulation

Imagine that we desire to calculate a difference of means. In this instance we are using a fixed covariate x and the purpose of this difference in means is to execute a placebo test or a balance test not to assess the causal effects of a treatment. Imagine two scenarios, one in which a binary treatment $Z_{ib} = 1$ is assigned to $m_b = 1$ subject within each of B pairs $b = 1, \dots, B; i = 1, \dots, n_b, B \leq n; n = \sum_{b=1}^B n_b$ (for pairs $n = 2B$ and thus $n_b = 2$), and another in which a binary treatment

²This paper is written in the mixture of R and \LaTeX known as Sweave (Leisch 2002, 2005) and, as such, all of the code required to produce all of the results, tables, and figures (as well as additional analyses not reported) are available for learning and exploration from <http://jakebowers.org>. Thus, I will spend relatively little time discussing the details of the different methods, assuming that those interested in learning more will download the source code of the paper and apply themselves to adapting it for their own purposes.

$Z_i = 1$ is assigned to $m = n - m = (n/2)$, subjects $i = 1, \dots, n$ without any blocking. Consider the test statistics

$$\begin{aligned} d_{\text{pairs}} &= \frac{1}{B} \sum_{b=1}^B \left(\sum_{i=1}^{n_b} Z_{ib} x_{ib} / m_b - \sum_{i=1}^{n_b} (1 - Z_{ib}) x_{ib} / (n_b - m_b) \right) \\ &= \frac{1}{B} \sum_{b=1}^B \left(\sum_{i=1}^2 (Z_{ib} x_{ib} - (1 - Z_{ib}) x_{ib}) \right) \end{aligned} \quad (1)$$

which reduces to the difference in means of x between treated and control units within pairs summed across pairs and

$$d_{\text{no pairs}} = \sum_{i=1}^n Z_i x_i / m - \sum_{i=1}^n (1 - Z_i) x_i / (n - m) \quad (2)$$

which sums across all units within control and treatment conditions. These two quantities are the same even if one is written as an average over B pairs: because pairs are blocks of equal size and therefore each block-specific quantity ought to contribute equally to the sum.

The theory of simple random sampling without replacement suggests that the variances of these statistics differ. First,

$$\begin{aligned} \text{Var}(d_{\text{no pairs}}) &= \frac{n}{m(n - m)} \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n - 1} \\ &= \left(\frac{4}{n} \right) \left(\frac{1}{n - 1} \right) \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned} \quad (3)$$

Second,

$$\begin{aligned} \text{Var}(d_{\text{pairs}}) &= \left(\frac{1}{B} \right)^2 \sum_{b=1}^B \frac{n_b}{m_b(n_b - m_b)} \sum_{i=1}^{n_b} \frac{(x_{ib} - \bar{x}_b)^2}{n_b - 1} \\ &= \left(\frac{2}{B^2} \right) \sum_{b=1}^B \sum_{i=1}^2 (x_{ib} - \bar{x}_b)^2. \end{aligned} \quad (4)$$

If pairs were created on the basis of similarity on x , then $\text{Var}(d_{\text{no pairs}}) > \text{Var}(d_{\text{pairs}})$ because $\sum_{i=1}^n (x_i - \bar{x})^2 > \sum_{b=1}^B \sum_{i=1}^2 (x_{ib} - \bar{x}_b)^2$. Any given x_i will be farther from the overall mean (\bar{x}) than it would be from the mean of its pair (\bar{x}_b). Note also that the constants multiplying the sums are $(4/n(n - 1))$ in the unpaired case and $8/n^2$ (since $B = (n/2)$) in the paired case. As long as

$n > 2$, $(4/(n^2 - n)) < (8/n^2)$ and this difference diminishes as n increases. Thus, benefits of pairing can diminish as the size of the experiment increases as long as within-pair homogeneity does not depend on sample size.

To dramatize the benefits possible from blocking, I include a simple simulation study based roughly on real data from a field experiment of newspaper advertisements and turnout in US cities (Panagopoulos 2006). The cities in this study differed in baseline turnout (with baseline turnout ranging from roughly 10 percent to roughly 50 percent). Panagopoulos paired the cities before randomly assigning treatment, with baseline turnout differences within pair ranging from one to seven percentage points in absolute value. The simulation presented here takes his original eight-city dataset and makes two fake versions, one with thirty two cities and another with 160 cities. The expanded datasets were created by adding small amounts of uniform random noise to copies of the original dataset. These new versions of the original dataset maintain the same general relationships between treatment and outcomes and covariates and pairs as the original, but allow us to examine the effects of increasing sample size. In this simulation the covariate values are created so that the average difference of means within pair is zero. For each of the 1000 iterations of the simulation and for each dataset ($n = 8, 32, 160$), the procedure was:

Create fake data based on original data Each pair receives a random draw from a Normal distribution with mean equal to the baseline outcome of its control group and standard deviation equal to the standard deviation of the pair: For the original dataset the within pair differences on baseline turnout of 1, 4, and 7 translate into standard deviations of roughly .7, 3, and 5. The “true” difference is required to be zero within every pair, but each pair may have a different baseline level of turnout and a different amount of variation — mirroring the actual field experiment.

Estimate variance of treatment effects under null Apply equations (3) and (4).³

Figure 1 shows that estimates of $\sqrt{\text{Var}(d_{\text{no pairs}})}$ were always larger than the estimates for $\sqrt{\text{Var}(d_{\text{pairs}})}$ although the ratio $\sqrt{\text{Var}(d_{\text{no pairs}})}/\sqrt{\text{Var}(d_{\text{pairs}})}$ diminished as the size of the experiment increased. null standard deviation across the simulations for the paired eight city example (1.8) is close to the average standard deviation for the non-paired tests in the 160 city example (2.3). That is, these simulations show that a paired experiment of eight units could be more precise than an unpaired experiment of 160 units (although of course a paired experiment of 160 units would be yet more precise (the average null sd in that case is .5)). Notice that, in this case, the advantages of the paired design diminish as the size of the experiment increases but do not disappear.

Pairing provides the largest opportunity for enhancement of precision in the comparison of two treatments — after all, it is hard to imagine a set more homogeneous than two subjects nearly identical on baseline covariates (Imai, King, and Nall 2009). And, even if it is possible for an unpaired design to produce differences of means with lower variance than a paired design, it is improbable given common political science measuring instruments.

What do we do once the experiment has run? One can enhance precision using covariates either via post-stratification (grouping units with similar values of covariates) or covariance adjustment

³The actual computations used the `xBalance` command found in the `RIttools` library for R (Bowers, Fredrickson, and Hansen 2009) as described in Hansen and Bowers (2008).

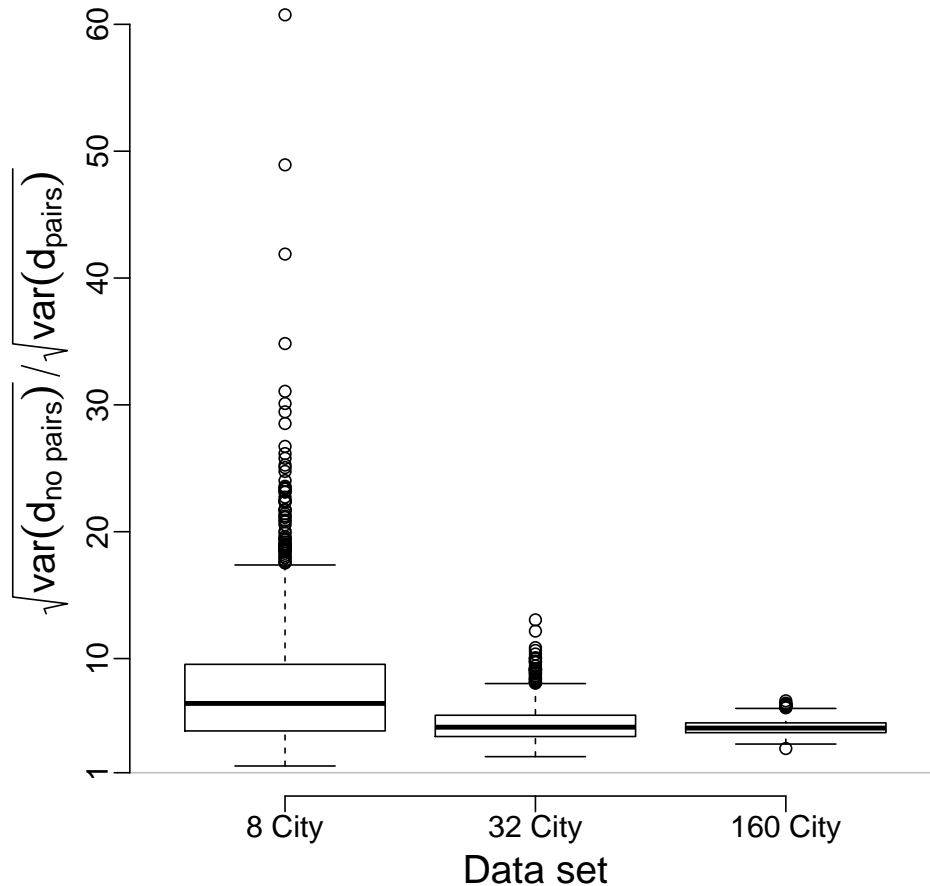


Figure 1: Paired designs were always more efficient than unpaired designs in the simulations described in Section 2. The unpaired design dominates the paired design less dramatically as the size of the experiment increases.

(fitting linear models with covariates predicting outcomes). In the first case analyses proceeds as if pre-stratified: estimates of treatment effects are made conditional on the chosen grouping. In the second case linear models “adjust for” covariates — increased precision can result from their inclusion in linear models under the same logic as that taught in any introduction to statistics classes.

3 Balance Assessment: Graphs and Tests

We expect that the bias reduction operating characteristics of random assignment would make the baseline outcomes in the control group comparable to the baseline outcomes in the treatment group. If the distribution of a covariate is similar between treated and control groups we say that this covariate is “balanced”, or that the experiment is balanced with respect to that covariate. Yet, it is always possible that any given randomization might make one or more observed (or unobserved) covariates imbalanced merely by chance. If the imbalanced covariates seem particularly relevant to substantive interpretations of the outcome (as would be the case with outcomes measured before the treatment was assigned), we would not want such differences to confuse treatment-versus-control

differences of post-treatment outcomes.

One graphical mode of assessment

Figure 2 provides both some reassurance and some worry about balance on the baseline outcome in the thirty two-city data example. The distributions of baseline turnout are nearly the same (by eye) in groups two and three, but differ (again, by eye) in groups one and four. Should these differences cause concern?

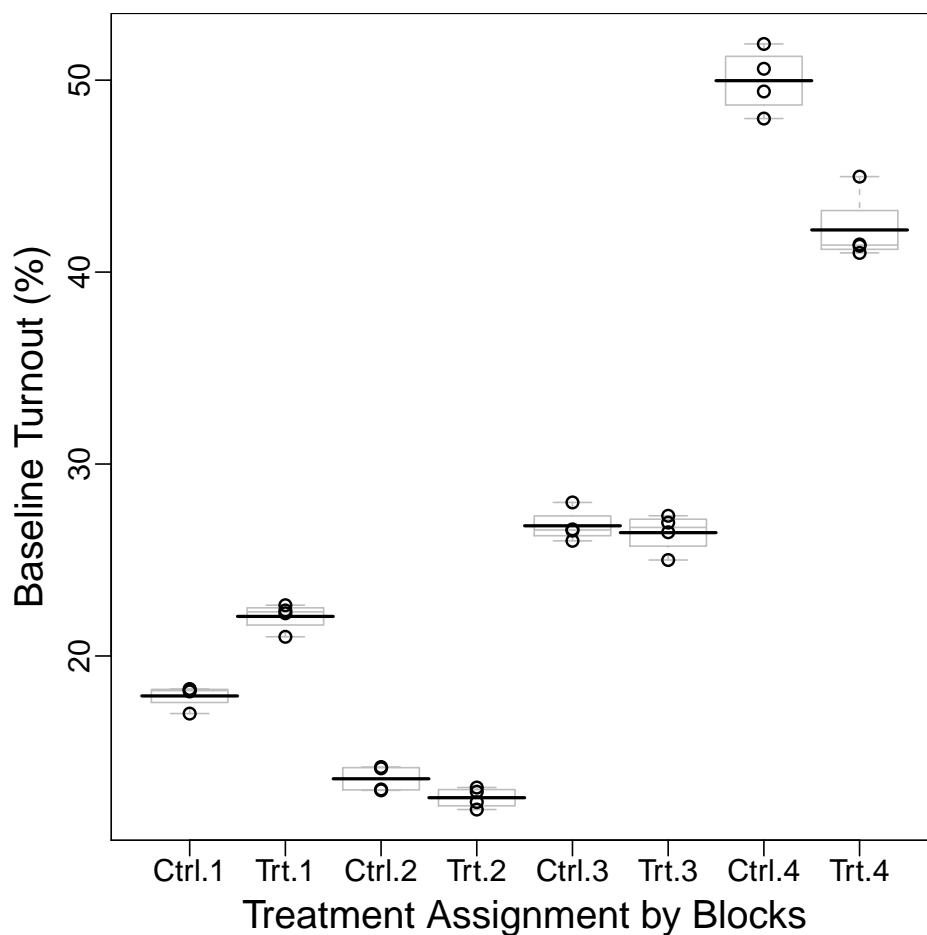


Figure 2: Within each of the four strata (1,2,3,4) the baseline outcomes (past turnout) for the thirty two-city data are plotted by assignment condition (Trt=Treatment, Ctrl=Control). Means are long, thick black lines. Boxplots in gray.

The top row of Table 1 provides one answer to the question about worries about random imbalance on baseline outcome. We would not be very surprised to see a mean-difference of $d_{\text{blocks}} = -1.2$ if there were no real difference given this design ($p = .2$). Of course, questions about balance of a study are not answerable by looking at only one variable. Table 1 thus shows a set of randomization-based balance tests to assess the null of no difference in means between treated and control groups on covariates one-by-one and also all together using an omnibus balance test. We easily reject the null of balance on the linear combination of these variables in this case ($p=0.00003$), although

| x | \bar{x}_{Ctrl} | \bar{x}_{Tt} | d_{blocks} | $\sqrt{\text{Var}(d_{\text{blocks}})}$ | Std. d_{blocks} | z | p |
|----------------------|-------------------------|-----------------------|---------------------|--|--------------------------|------|-----|
| Baseline Outcome | 27.1 | 25.8 | -1.2 | 0.9 | -0.1 | -1.4 | 0.2 |
| Percent Black | 17.3 | 2.6 | -14.7 | 4.2 | -1.2 | -3.5 | 0.0 |
| Median HH Income | 30.6 | 46.3 | 15.7 | 3.2 | 2.6 | 4.9 | 0.0 |
| Number of Candidates | 13.2 | 10.3 | -2.9 | 2.0 | -0.5 | -1.5 | 0.1 |

Table 1: One-by-one Balance Tests for Covariates adjusted for blocking in the blocked thirty two-city Study. Two-sided p -values are reported in the p column based on referring z to Normal distribution approximating the distribution of the mean-difference under null of no effects. An omnibus balance test casts doubt on the null hypothesis of balance on linear combinations of these covariates at $p=0.00003$. Test statistics (d_{blocks} , etc..) are generalizations of equations (1) and (4) developed in Hansen and Bowers (2008) and implemented in Bowers, Fredrickson, and Hansen (2009). Statistical inference here (z and p -values) is randomization-based but uses large-sample Normal approximations for convenience. Other difference-of-means tests without the large-sample approximation, and other tests such as Kolmogorov-Smirnov tests and Wilcoxon rank sum tests, provide the same qualitative interpretations. For example, the p -values on the tests of balance for baseline outcome (row 1 in the table) ranged from $p = 0.16$ and $p = .17$ for the simulated and asymptotic difference of means tests respectively, to $p = .33$ and $p = .72$ for exact and simulated Wilcoxon rank sum and Kolmogorov-Smirnov tests.

whether we worry about the evident differences on percent black or median income of the cities may depend somewhat on the extent to which we fear that such factors matter for our outcome — or whether these observed imbalances suggest more substantively worrisome imbalances in variables that we do not observe. Does this mean the experiment is broken? Not necessarily.⁴

Understanding p -values in balance tests.

Say we didn't observe thirty-two observations in sets of four but only eight observations in pairs. What would our balance tests report then? Table 2 shows the results for such tests, analogous to those shown in Table 1.

Notice that our p -values now quantify less doubt about the null. Is there something wrong with randomization-based tests, if by reducing the size of our experiment we would change our judgement about the operation of the randomization procedure? The answer here is no.⁵

The p -values reported from balance tests used here summarize the extent to which random imbalance

⁴Note that it would be strange to see such imbalances in a study run with thirty-two observations rather than an eight observation study expanded to thirty-two observations artificially as done here. These imbalances, however, dramatize and clarify benefits and dangers of post-treatment adjustment as explained throughout this paper.

⁵Echoing Senn (1994) among others in the clinical trials world, Imai, King, and Stuart (2008) provide some arguments against hypothesis tests for balance in observational studies. Hansen (2008) answers these criticisms with a formal account of the intuition provided here pointing out that 1) randomization-based tests do not suffer the problems highlighted by those authors and 2) highlighting the general role that p -values play in randomization based inference.

| x | \bar{x}_{Ctrl} | \bar{x}_{Tt} | d_{blocks} | $\sqrt{\text{Var}(d_{\text{blocks}})}$ | Std. d_{blocks} | z | p |
|----------------------|-------------------------|-----------------------|---------------------|--|--------------------------|------|-----|
| Baseline Outcome | 26.0 | 24.8 | -1.2 | 2.0 | -0.1 | -0.6 | 0.5 |
| Percent Black | 16.8 | 2.4 | -14.4 | 11.0 | -1.1 | -1.3 | 0.2 |
| Median HH Income | 29.2 | 44.5 | 15.4 | 8.1 | 2.5 | 1.9 | 0.1 |
| Number of Candidates | 13.0 | 10.0 | -3.0 | 5.1 | -0.5 | -0.6 | 0.6 |

Table 2: One-by-one Balance Tests for Covariates in the blocked 8 City Study. An omnibus balance test casts doubt on the null hypothesis of balance on linear combinations of these covariates at $p=0.41$. Test statistics (d_{blocks} , etc..) are generalizations of equations (1, 4) developed in Hansen and Bowers (2008) and implemented in Bowers, Fredrickson, and Hansen (2009). Statistical inference here (z and p -values) is randomization-based but uses large-sample Normal approximations for convenience.

is worrisome. With a sample size of eight, the confidence interval for our treatment effect will be large — taking the pairing into account, a 95% interval will be on the order of ± 3.5 (as roughly calculated on the baseline outcomes in Section 2). For example, both Tables 2 and 1 show the block-adjusted mean-difference in percent black between treated and control groups to be roughly 14.5 percentage points. In the thirty-two-city example, this difference cast great doubt against the null of balance, while in the eight-city example this difference casts less doubt. Now, evaluating the difference between controls and treatment on actual outcome in the eight-city case gives $d_{\text{pairs}} = 1.5$ and under the null of no effects gives $\sqrt{\text{Var}(d_{\text{pairs}})} = 2.6$ — and inverting this test leads to an approximate 88% confidence interval of roughly $[-7, 7]$: The width of the confidence interval itself is about fourteen points of turnout. Even if percent black were a perfect predictor of turnout (which it is not, with a pair adjusted linear relationship of -0.07 in the eight-city case), the p -value of $.2$ indicates that the relationship with treatment assignment is weak enough, and the confidence intervals on the treatment effect itself would be wide enough, to make any random imbalance from percent black a small concern. That is, the p -values reported in Table 2 tell us that random imbalances of the sizes seen here will be small relative to the size of the confidence interval calculated on the treatment effect. With a large sample, a small random imbalance is proportionately more worrisome because it is large relative to the standard error of the estimated treatment effect. Given the large confidence intervals on a treatment effect estimated on eight units, the random imbalances shown here are less worrisome — and the p -values encode this worry just as they encode the plausibility of the null.

Thus, even though our eyes suggested we worry about the random imbalance on baseline turnout we saw in Figure 2, that amount of imbalance on baseline outcome is to be expected in both the thirty-two and eight-city cases — it is an amount of bias that would have little effect on the treatment effect were we to adjust for it. Of course, the omnibus test for imbalance on all four covariates simultaneously reported in Table 1 does cast doubt on the null of balance — and the tests using the d -statistics in the table suggest that the problem is with percent black and median household income rather than baseline outcomes or number of candidates.⁶

⁶If we had twenty covariates and rejected the null of balance with $p < .05$, we would expect to falsely see evidence of imbalance in one of twenty covariates. Thus, Hansen and Bowers (2008) urge the use of an omnibus test — a test which assesses balance across all linear combinations of the covariates in the table. Yet, the variable by variable display is useful in the same way that graphs such as Figure 2 are useful — in suggesting (not proving) the sources of imbalance.

4 Covariates allow adjustment for random imbalance

Even with carefully designed experiments there may be a need in the analysis to make some adjustment for bias. In some situations where randomization has been used, there may be some suggestion from the data that either by accident effective balance of important features has not been achieved or that possibly the implementation of the randomization has been ineffective (Cox and Reid 2000, 29).

A well-randomized experiment aiming to explain something about political participation showing manifest imbalance on education poses a quandry. If the analyst decides to adjust, she then may fall under suspicion: even given a true treatment effect of zero, one adjustment out of many tried will provide a p -value casting doubt on the null of no effect merely through chance. Without adjustment we know how to interpret p -values as expressions of doubt about a given hypothesis: low p -values cast more doubt than high p -values. Adjustment in and of itself does not invalidate this interpretation: a p -value is still a p -value. Concerns center rather on 1) whether an “adjusted treatment effect” is substantively meaningful and how it relates to different types of units experiencing the treatment in different ways — that is, the concerns center on the meaning of “adjustment” in the context of the adjustment method (a linear model or a post-stratification) and 2) whether some specification search was conducted with only the largest adjusted treatment effect reported representing a particularly rare or strange configuration of types of units. Such worries do not arise in the absence of adjustment. Yet, if the analyst declines to adjust, then she knows that part of the treatment effect in her political participation study is due to differences in education, thereby muddying the interpretation of her study.

One may answer such concerns by announcing in advance the variables for which random imbalance would be particularly worrisome and also provide a proposed adjustment and assessment plan a priori. Also, if one could separate adjustment from estimation of treatment effects, one may also avoid the problem of data snooping. For example, Bowers and Panagopoulos (2009) propose a power-analysis based method of choosing covariance adjustment specifications which can be executed independently of treatment effect estimation, and it is well-known that one may post-stratify and/or match without ever inspecting outcomes. Post-stratification may also relieve worries about whether comparisons adjusted using linear models are artifacts of the functional form (Gelman and Hill 2007, ch.9).

I have noted that there are two broad categories of statistical adjustment for random imbalance: adjustment by stratification and adjustment using models of outcomes. In both cases, adjustment amounts to choice of weights; and in both cases adjustment may be executed entirely to enhance precision even if there is no appreciable evidence of random imbalance. Notice that the “unadjusted” estimate may already be an appropriately weighted combination of block-specific treatment effects — and that to fail to weight (or “adjust”) for block-specific probabilities of treatment assignment will confound estimates of average treatment effects (if the probabilities of assignment differ across blocks) and decrease precision (if the variation in the outcomes is much more homogeneous within blocks than across blocks).

Post-stratification enables adjustment but must respect blocking and design.

Say, for example, that within blocks defined by town, the treated group on average contained too many men (and that although gender was important in the study, the researcher either could not or forgot to block on it within existing blocks). An obvious method of preventing “male” from unduly confusing with estimates of treatment effects is to only compare men to men, within block. Analysis then proceeds using the new set of blocks (which represent both the pre-treatment blocks and the new post-treatment strata within them) as before.

One may also use modern algorithmic matching techniques to construct strata. Keele, McConnaughy, and White (2008) argue in favor of matching over linear models for covariance adjustment and show simulations suggesting that such post-stratification can increase precision. Notice that matching to adjust experiments is different from matching in observational studies: matching here must be done without replacement in order to respect the assignment process of the experiment itself and matching must be full. That is, although common practice in matching in observational studies is to exclude certain observations as unmatchable or perhaps to reuse certain excellent control units, in a randomized experiment every observation must be retained and matched only once. This limits the precision enhancing features of matching (at least in theory) since homogeneity will be bounded first by the blocking structure before random assignment and then again by requiring that all observations be matched.

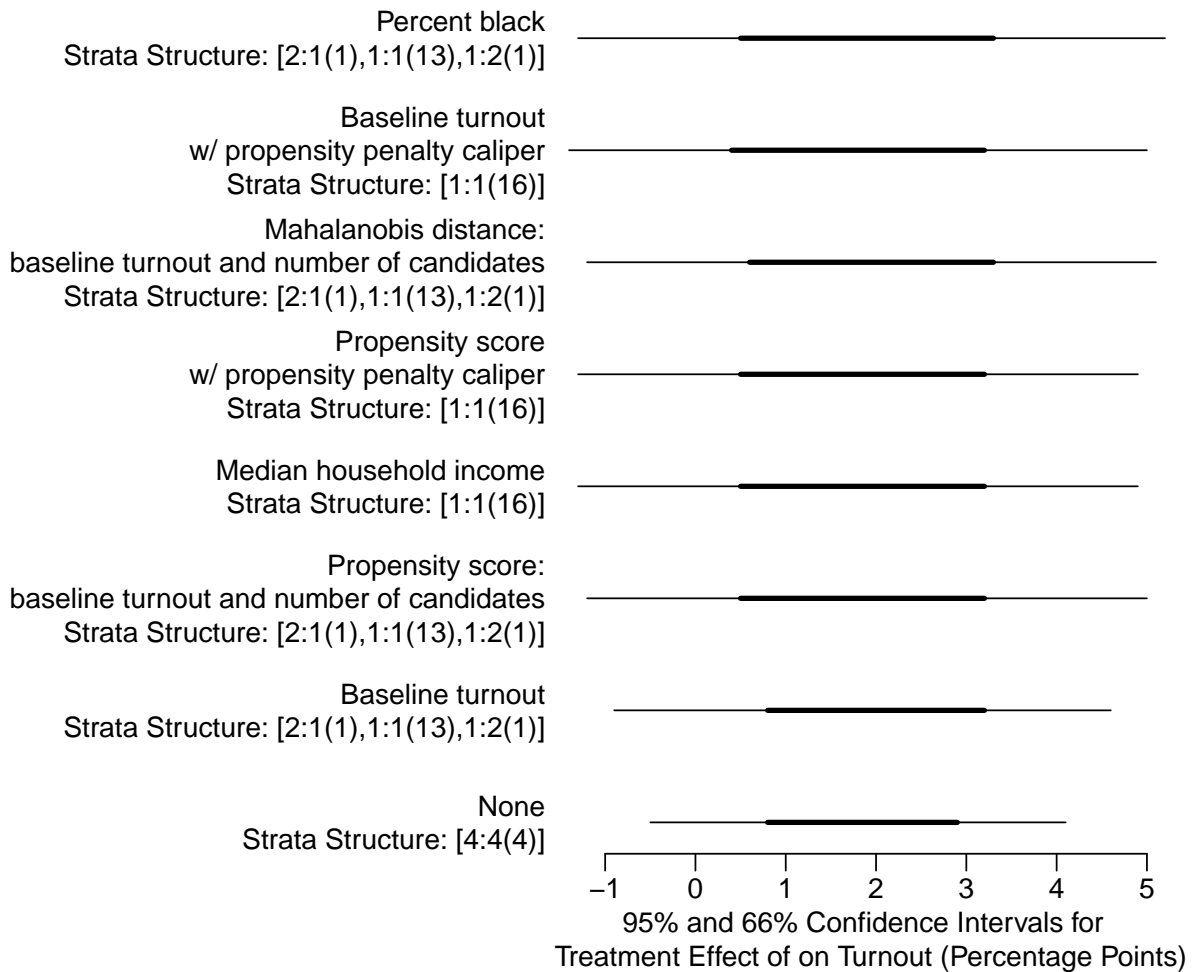


Figure 3: Confidence intervals for the difference in turnout between treated and control cities in the thirty two-city turnout experiment data ordered by width from narrowest at bottom to widest at top. Each line represents the interval for the treatment effect after different post-stratification adjustments have been applied within the existing 4 blocks of 8 cities. Thin lines show the 95% intervals. Thick lines show the 66% intervals. The propensity-score was calculated from a logistic regression of treatment assignment on baseline turnout, number of candidates running for office, percent black and median household income (both from the 2000 Census), and block-indicators. All post-stratification and interval estimation was done with full, optimal matching using the `optmatch` (Hansen and Fredrickson 2010) and `RITools` (Bowers, Fredrickson, and Hansen 2009) packages for R. Numbers below the post-stratification label show the structure of the stratification: for example, without any post-stratification the experiment had 4 blocks, each with 4 treated and 4 control cities. The matching on absolute distance on the propensity score with a propensity caliper penalty produced 16 pairs of treated and control cities (1:1(16)). The match on absolute distance on baseline turnout produced 1 set with 2 treated and 1 control (2:1(1)), 13 paired sets (1:1(13)) and 1 set with 1 treated and 2 controls (1:2(1)). Since no observation could be excluded, calipers implied penalties for the matching algorithm rather than excluding observations from the matching (Rosenbaum 2010, ch 8).

Figure 3 shows confidence intervals resulting from a variety of post-stratification adjustments made to the thirty two-city turnout data. In this particular experiment the within-set homogeneity increase resulting from post-stratification did not outweigh the decrease in degrees of freedom occurring from the need to account for strata: the shortest confidence interval was for the unadjusted data (shown at the bottom of the plot).

Did the post-stratification help with the balance problems with the Census variables evident in Table 1? Table 3 shows the strata-adjusted mean differences and p -value for balance tests now adjusting for the post-stratification in addition to the blocking. Balance on baseline turnout and number of candidates does improve somewhat with the matching but balance on percent black and median household income does not appreciably improve. Notice a benefit of post-stratification here: the post-adjustment balance test shows us that we have two covariates which we could not balance.

| x | Post-Hoc Full-Matching on: | | | | | |
|----------------------|----------------------------|-----|---------------------|-----|---------------------|-----|
| | Blocks | | Baseline Turnout | | Propensity Score | |
| | d_{strata} | p | d_{strata} | p | d_{strata} | p |
| Baseline Outcome | -1.2 | 0.2 | -1.2 | 0.3 | -1.2 | 0.3 |
| Percent black | -14.7 | 0.0 | -14.7 | 0.0 | -14.7 | 0.0 |
| Median HH Income | 15.7 | 0.0 | 15.8 | 0.0 | 15.7 | 0.0 |
| Number of candidates | -2.9 | 0.1 | -2.6 | 0.3 | -2.9 | 0.3 |

Table 3: One-by-one Balance Tests for Covariates adjusted for covariates by post-stratification in the blocked thirty two-city Study. Two-sided p -values assess evidence against the null of no effects. An omnibus balance test casts doubt on the null hypothesis of balance on linear combinations of these covariates adjusted for Blocks, and two kinds of Post-Hoc full-matching within blocks at $p=0.00003, 0.003, 0.004$. Strata adjusted mean-differences (d_{strata}) are generalizations of equations (1) developed in Hansen and Bowers (2008) and implemented in the `RIttools` package for R (Bowers, Fredrickson, and Hansen 2009). Statistical inference here (p -values) is randomization-based but uses large-sample Normal approximations for convenience. Post-hoc stratification results from optimal, full-matching (Hansen 2004) on either absolute distance on baseline turnout or absolute distance on a propensity score with propensity caliper penalty shown in Figure 3.

Discussion of the advantages of blocking in Section 2 is, in essence, a discussion about how to analyze blocked (pre- or post-stratified) experimental data. The rest of the paper is devoted to understanding what it is that we mean by “covariance adjusted treatment effects.”

Linear models enable adjustment but may mislead the unwary

Even with carefully designed experiments there may be a need in the analysis to make some adjustment for bias. In some situations where randomization has been used, there may be some suggestion from that data that either by accident effective balance of important features has not been achieved or that possibly the implementation of the randomization has been ineffective (Cox and Reid 2000, 29).

Cox and Reid’s discussion in their Section 2.3 entitled “Retrospective adjustment for bias” echoes Cox (1958, 51–2) and Fisher (1925). What they call “bias,” I think might more accurately be called “random imbalance”.

Although Fisher developed the analysis of covariance using an asymptotic F test which approximated the randomized-based results, others have since noted that the standard sampling-based infinite-population or Normal-model theory of linear models does not justify their use in randomized experiments. For example, in discussing regression standard errors, Cox and McCullagh (1982) note “It is known . . . that ‘exact’ second-order properties of analysis of covariance for precision improvement do not follow from randomization theory . . .” (547). In this section, I provide an overview of a randomization-based method for covariance adjustment which can use Normal approximations in the same way as those used for balance assessment and placebo-tests above. This method is not subject to the concerns of Freedman (2008*a,b,c*), and thus suggests itself as useful in circumstances where the linear model as an approximation may cause concern or perhaps as a check on such approximations.

First, though, let us get clear on what it means to “adjust for” random imbalance on a covariate.

What does it mean to say that an estimate has been “adjusted” for the “covariance” of other variables?

Let us look first at how covariance adjustment might work in the absence of blocking by looking only at the first block of eight units in the thirty two-city dataset. Figure 4 is inspired by similar figures in Cox (1958, ch. 4) with dark gray showing treated units and black showing control units. The unadjusted difference of means is 6.6 (the vertical distance between the open squares that are not vertically aligned on the gray vertical line). The thick diagonal lines are the predictions from a linear regression of the outcome on an indicator of treatment and baseline outcome. The adjusted difference of means is the vertical difference between the regression lines, here, five. If there had been no relationship between baseline outcomes and post-treatment outcomes, the regression lines would have been flat and the vertical distances between those lines would have been the same as the unadjusted difference of means (the thin dark gray and black horizontal dashed lines). As ought to be clear here, parallel effects is a required assumption for covariance adjustment to be meaningful. In this case, a regression allowing different slopes and intercepts between the treatment groups shows the treatment slope of .25 and the control group slope of .23, thus, the assumption is warranted.

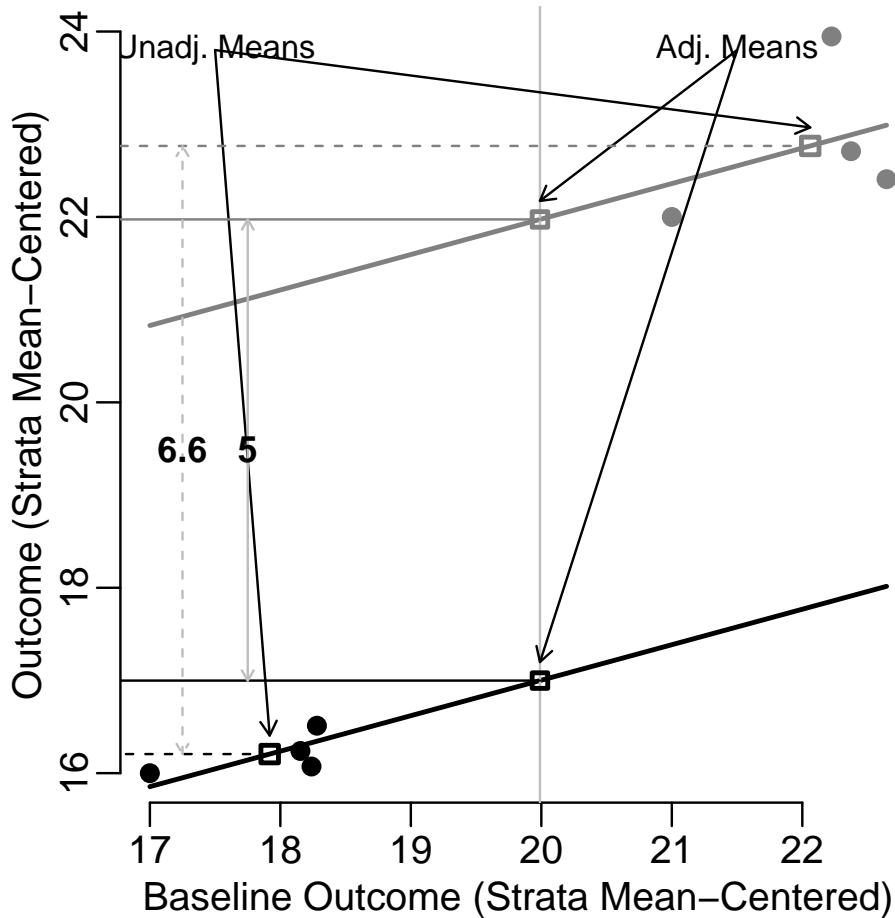


Figure 4: Covariance-adjustment in a simple random experiment. Dark gray and black circles show treated and control units respectively. The unadjusted difference of means is 6.6. The thick diagonal lines are: $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 Z_i + \hat{\beta}_2 y_{i,t-1}$ with $Z_i = \{0, 1\}$ and the adjusted average treatment effect is $(\hat{Y}_i|Z_i = 1, y_{i,t-1} = \bar{y}_{t_1}) - (\hat{Y}_i|Z_i = 0, y_{i,t-1} = \bar{y}_{t_1}) = 5$.

What about with blocked data? Figure 5 shows the covariance adjustment for three different covariates: a) baseline outcomes (on the left), b) percent black (in the middle), and c) median household income (in \$1000s, on the right). In each case the data are *aligned* within each block by subtracting the block-mean from the observed outcome (i.e., block centered). A linear regression of the block-mean centered outcome on the block-mean centered covariate plus the treatment indicator is equivalent to a linear regression of the observed outcome on the covariate plus treatment indicator plus indicator variables recording membership in blocks (i.e., “fixed effects” for blocks).

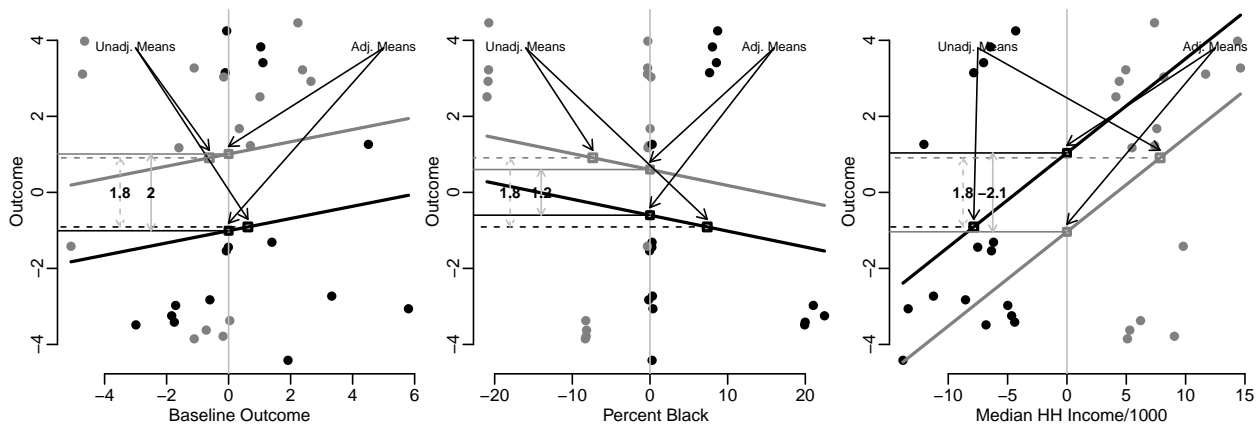


Figure 5: Covariance-adjustment in a Blocked Random Experiment (4 Blocks). Dark gray and black circles show treated and control units respectively. All variables are block-mean centered.

In the leftmost plot, we see that adjustment for baseline outcomes in addition to the blocking structure does very little to change the treatment effect: in fact, these blocks were chosen with reference to baseline outcomes, and so adjusting for blocks is roughly equivalent to adjusting for baseline outcomes (only it does not require what is clearly a dubious parallel lines assumption). If the parallel lines assumption held in this case, however, we might talk meaningfully about an adjusted effect. The average treatment effect in the first block is 6.6, but the treated units were more likely to participate, on average, than the control units even at baseline (a mean difference of 4.1). Some of the 6.6 points of turnout may well be due to baseline differences (no more than 4.1 we assume, and probably less so, since it would only matter to the extent that baseline turnout is also related by chance in a given sample to treatment assignment). In this case, the block-specific relationship between baseline outcomes and treatment assignment is vanishingly small (difference of means is 1.9), so only about two points of the average treatment effect is due to the baseline treatment effect (where “due to” is in a very specific linear smoothed conditional means sense). The adjusted effect is actually closer to 5 than 4.6 because the intuitions here provided with differences of means is not identical to what is happening with an analysis of covariance (although it is close and provides helpful intuition in this case).

The middle and right hand plots show two instances in which the intuitions using differences of means become more difficult to believe. In strata 1, 2, and 4, every control unit has a higher percent black than every treated unit. The unadjusted average treatment effect is 1.8, but after adjustment for percent black the effect is 1.2. The middle plot, however, shows that the assumption of parallel effects is hard to sustain and that there is little overlap in the distributions of percent black between the treatment and control groups. In this case, however, the adjustment makes little difference in the qualitative interpretation of the treatment effect.

The right hand figure is a more extreme case of the middle figure. This time there is no overlap at all between the distributions of median income between the controls (in black, and all on the left side of the plot) and the treated units (in dark gray, and all on the right side of the plot). The adjustment causes the treatment effect to change sign — from 1.8 to -2.1 percentage points of turnout. Is -2.1 a

better estimate of the treatment effect than 1.8? Clearly median income has a strong relationship with outcomes and also, via random imbalance, with treatment assignment (recall the test casting doubt on the null of balance in Table 1).

What is the problem with covariance adjustment in this way? First, as noted earlier, the assumption of parallel lines is not correct. Second, we begin to notice another problem not mentioned in textbooks like Cox and Reid (2000) or Cox (1958) — random assignment will, in large samples, ensure balance in the distributions of covariates but will not ensure such balance in small samples. This means that the distributions of the covariates in theory, ought to be quite similar between the two groups. However, the theory does not exclude the possibility of random imbalance on one of many covariates, and it is well known that random imbalance can and does appear in practice. Adjustments for such imbalance can be done in such a way that adjusted mean differences are still meaningful representations of the treatment effect (as shown by the adjustment for baseline outcomes in the left plot of Figure 5). But, as the distributions of covariates become more unbalanced, the covariance adjustment can mislead. It is hard to claim, based on these data, that adjusting for median household income is a meaningful operation — one just cannot imagine (in these data) finding groups of treated and control units with the same median household income.⁷

Thus, we can see where some criticisms of covariance adjustment might easily come from: covariance adjustment done with multiple regression without additional diagnostics poses a real problem for diagnosing whether the imbalance is so severe as to provide no “common support” in the distributions of the covariates. In such cases, one really cannot “adjust for” the imbalance and one must be resigned to the fact that treatment versus control comparisons, in the data, reflect something other than treatment assignment even if they would not in a larger sample or across repeated experiments. Of course, as sample size grows, given a finite set of covariates and a treatment with finite variance (i.e., a treatment which only has a few values that does not gain values as sample size grows), we would expect the problem of common support to disappear in randomized studies. Luckily, in many studies, one can assess such problems before treatment is administered.

Randomization alone can justify statistical inference about covariance-adjusted quantities without requiring the common assumptions required to justify standard linear regression tables

A predominant use for covariance adjustment is not to ameliorate random imbalance but to enhance statistical precision. To the extent that a covariate predicts outcomes, one may use it to reduce the noise in the outcome unrelated to treatment assignment, and thus help make treatment effects manifest. Covariance adjustment (whether for precision or for random imbalance) means linear regression. In theory, counter-factual statistical inference using the linear regression model for covariance adjustment estimator is biased (Freedman 2008*b,c,a*) but, in practice, it is often an excellent approximation (Green 2009; Schochet 2009). What should we do when we worry about the approximation: for example when the experiment is small, or there is great heterogeneity in effects and/or variance of effects across blocks, or great heterogeneity or discreteness in the outcome (such that the central limit theorem takes longer to kick in than one would prefer)? Rosenbaum

⁷Notice that this problem also occurred in Figure 4 but was not mentioned in order not to detract from the pedagogical task of describing the mechanism of covariance adjustment.

(2002a) presents a simple argument which builds on the basics of Fisher’s randomization-based inference. Here I provide some very brief intuition to guide study of that paper. This method of randomization-justified covariance adjustment does not rely on the linear model for statistical inference but does “adjust” using the linear model.

Say an outcome is measured with noise caused in part by covariates. When we randomly assign treatment we are attempting to isolate the part of the variation in the outcome due to the treatment from that due to other factors. Say we are still interested in the difference in mean outcomes between treatment and control groups as assigned. The standard deviations of those means may be large (making the treatment effect hard to detect) or small (making the treatment effect more easily manifest). If part of the noise in the outcome is due to covariates, then the residual from regressing the outcome on the covariates represents a less noisy version of the outcome — the outcome without noise from linear relationships with covariates. This residual e_{ib} (for unit i in block b) is measured in units of the outcome (i.e., “percent turning out to vote” in our running fake data example). The potential outcomes to treatment and control for units i in blocks b , y_{Tib} and y_{Cib} , are fixed, Y_{ib} is random by virtue of its relationship with random assignment Z because $Y_{ib} = Z_{ib}y_{Tib} + (1 - Z_{ib})y_{Cib}$. A null hypothesis tells us what function of Y_i and Z_i would recover y_{Ci} : that is, if the null hypothesis is correct, then removing the effect (say, τ_{ib}) from the treated units $Y_{ib,Z=1}$ would tell us how the treated units would behave under control. Under a constant, additive model of effects, $y_{Tib} = y_{Cib} + \tau$ and so $Y_{ib} - Z_{ib}\tau_{ib} = y_{Cib}$.⁸ So, considering our null hypothesis for the sake of argument, $H_0 : \tau_0 = \tau_{ib}$, regressing $(Y_{ib} - Z_{ib}\tau_{ib})$ on x_i is regressing a fixed quantity (i.e., y_{Cib}) on another fixed quantity, x_{ib} and so the residuals from that regression are a fixed quantity.⁹ Thus one may substitute e for x in (1) and (4). Fisher’s style of inference begins with a test of a null hypothesis and inverts the hypothesis for a confidence interval: thus, the method allows for us to infer consistent estimates of the causal effect by testing a sequence of causal hypotheses τ_0 . Very loosely speaking, the point estimate is the causal effect hypothesised by the best-supported hypothesis tested.¹⁰ Note that this is a method of hypothesis testing, not of

⁸If this model is incorrect, then randomization-based inferences will be conservative but the coverage of the confidence intervals will still be correct as noted independently by (Gadbury 2001) and Robins (2002, § 2.1), and of course, other substantively meaningful models of effects are available Rosenbaum (ch. 5 2002b, see) and Rosenbaum (§ 6 2002a, and also see) or Rosenbaum (2010, ch. 2). For example, as Rosenbaum (2002c, 323) notes, if the treatment effect varies by a binary covariate x coding 0 for group 1 and 1 for group 2 (such that the parallel lines assumption is incorrect), we would then specify the potential responses to control as $Y_{ib} - Z_{ib}(\tau_{x=1}x_{ib} + \tau_{x=0}(1 - x_{ib}))$ for treatment effects that differ by group. I use the constant additive effects model in this paper to map most closely onto the causal quantities implied by the choice of a linear regression model for covariance adjustment: indeed, for this very reason I can show how both styles of covariance adjustment can produce identical quantities in Figure 6. Interested readers might find the discussion in Rosenbaum (2002c, § 3–6) useful for thinking about the equivalence of estimating an average treatment effect and testing a sequence of hypotheses about individual causal effects.

⁹For a single covariate, x and a regression fit $(Y_{ib} - Z_{ib}\tau_{ib}) = \hat{\beta}_0 + \hat{\beta}_1x_{ib}$, $e_{ib} = (Y_{ib} - Z_{ib}\tau_{ib}) - (\hat{\beta}_0 + \hat{\beta}_1x_{ib})$. The residual is written e , not \hat{e} , because the regression fit is not an estimate of an unknown quantity but merely calculating a function of fixed features of the existing data.

¹⁰See discussion of the Hodges-Lehmann point estimate in Rosenbaum (2002a) and Rosenbaum

estimation. It would be quite incorrect to interpret the difference means of residuals as an estimate of a treatment effect, because the residuals have already have specific causal hypotheses built into them as just described.

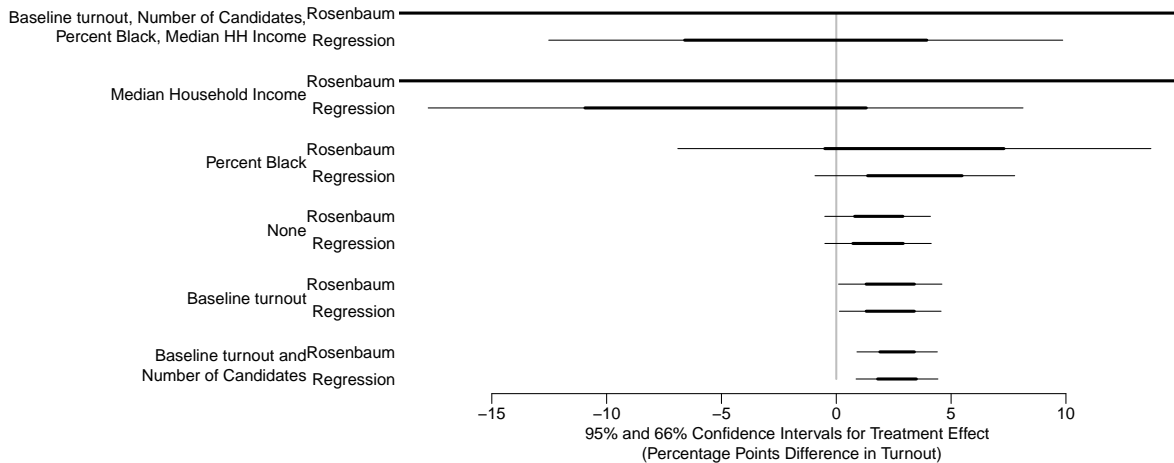


Figure 6: Confidence intervals for the difference in turnout between treated and control cities in the thirty-two-city turnout experiment data by type of covariance adjustment using least squares regression ordered by width from narrowest at bottom to widest at top. “Regression”-style adjustment regressed turnout on treatment indicator, block indicators, and covariates and referred to the standard iid+ t -based reference distribution for confidence intervals. “Rosenbaum”-style adjustment regressed turnout on block indicators, and covariates, and then used the residuals as the basis for tests of the null hypotheses implied by these confidence intervals. The reference distributions for the Rosenbaum-style are large-sample approximations to the randomization distribution implied by the design of the experiment using the `RITools` (Bowers, Fredrickson, and Hansen 2009) package for R. Thin lines show the 95% intervals. Thick lines show the 66% intervals. The randomization based confidence intervals for outcomes after adjustment by functions including median household income are essentially infinite.

Figure 6 shows that the Rosenbaum covariance adjustment in this particular data is well approximated by the direct regression-style covariance adjustment in the unadjusted case or when the adjustment is made for baseline turnout and number of candidates — and the version adjusted for baseline turnout and number of candidates (just) excludes zero from its 95% confidence interval.

(2002*b*, ch. 2) for more formal discussion of randomization-justified point-estimates of causal effects. In the context of a large, cluster randomized, field experiment with binary outcomes and non-random non-compliance, Hansen and Bowers (2009) show how one may approximate the results of testing sequences of hypotheses with simple calculations of means and associated randomization-based variances, including a method for randomization-based covariance adjustment. Because the Hansen and Bowers (2009) method approximates the results of testing sequences of hypotheses with simple means and variances their method requires an asymptotic justification. That paper (and related reproduction materials) (Bowers, Hansen, and Fredrickson 2008) also provide some tools for assessing the asymptotic justification.

The two approaches differ when the difference of means is adjusted for the Census variables. The most notable difference here is for median household income where the direct adjustment method is based entirely on the linear extrapolation between the groups while the Rosenbaum approach correctly captures the sense in which there is no reasonable way to adjust the treatment effect for this covariate. Since adjustment for percent black also requires much linear extrapolation, the randomization-based confidence interval again reflects the fact that the design itself has little information about what it means to adjust for this variable.

The advantages of this style of covariance adjustment are: 1) that it side-steps Freedman’s critiques of covariance adjustment for experiments,¹¹ 2) although we used large-sample normal approximations to evaluate the differences of means here, neither differences of means nor large-sample normal approximations are necessary (and the large-sample approximations are checkable),¹² 3) unmodeled heteroskedasticity or incorrect the functional form do not invalidate the statistical inference based on this method as they do for the standard approach. For example, the parallel lines assumption is no longer relevant for this approach since the only job of the linear model is to reduce the variance in the outcome. Finally, this particular data example allows us to notice another side benefit of the statistical property of correct coverage: when there is no (or little) information available in the design, the randomization-based approach will not overstate the certainty of conclusions in the same way as the model-based approach.¹³

¹¹In particular, Freedman (2008*b*, 189) notes that the Fisher-style covariance adjustment is valid “If $T_i = C_i$ for all i (the “strict null hypothesis”), then $\beta \equiv 0$ and adjustment will help—unless $\overline{\alpha Z} = 0$, i.e., the remaining variation (in C_i) is orthogonal to the covariate.” Another method, elaborated most recently in Hansen and Bowers (2009) also does not directly model the relationship between treatment and outcomes and so also avoids this critique.

¹²Rosenbaum (2002*a*) focuses on normal approximations to the Wilcoxon rank sum statistic as his preferred summary of treatment effects (and the normal approximations there are not necessary either, but merely convenient and often correct in large-enough samples with continuous outcomes).

¹³The disadvantages of this mode are not on display here (although they will be obvious to those who use the code contained in this paper for their own work). First, remember, this approach produces confidence intervals by testing sequences of hypotheses. It does not “estimate” causal effects as a standard regression estimator would but rather assesses the plausibility of causal effects using tests. Of course, such assessments of plausibility are also implicit in confidence intervals for standard regression estimators. However, the mechanics of the two methods of covariance adjustment are quite different. The randomization-based approach as implemented here builds a confidence interval by direct inversion of hypothesis tests: in this case, we tested hypotheses about the treatment effect from $\tau_0 = -20$ to $\tau_0 = 20$ by .1. This can be computationally burdensome if the number of hypotheses to test is large or if we eschew large-sample approximations. Second, we did not work to justify our choice of mean difference (rather than rank, or other summary of observed outcomes and treatment assignment). The standard linear regression estimator requires attention to mean-differences as the quantity of interest whereas any test-statistic may be used in the randomization-based method of adjustment shown here.

Best Practices for Regression-Based Adjustment

Adjusted treatment effect estimates always invite suspicion of data snooping or modeling artifacts. None of the techniques discussed here entirely prevents such criticism. Of course, the easy way to avoid such criticism is to announce in advance what kinds of random imbalance are most worrisome, and announce a plan for adjustment (including a plan for assessing the assumptions of the adjustment method chosen). Covariance adjustment using the standard linear regression model requires that one believe the assumptions of that model. For example, this model as implemented in most statistical software requires a correct model of the relationship between outcomes and covariates among treatment and control units (i.e. a correct functional form), that the heteroskedasticity induced by the experimental manipulation is slight, and that the sample size is large enough to overcome the problems highlighted by Freedman (2008*a,b,c*). As with any use of linear regression, one may assess many of these assumptions. If one or more of these assumptions appear tenuous, however, this paper has shown that one may still use the linear model for adjustment but do so in a way that avoids the need to make such commitments. Readers interested in the Rosenbaum (2002*a*) style of covariance-adjustment should closely study that paper. The code contained in the reproduction archive for this paper may also help advance understanding of that method.

5 The more you know, the more you know

A randomized study which allows “the phenomenon under test to manifest itself” provides particularly clear information and thus enhances theory assessment, theory creation, and policy implementation. Thus, researchers should attend to those elements of the design and analysis that would increase the precision of their results. This paper points to only a very small part of the enormous body of methodological work on the design of experiments.

Random assignment has three main scientific aims: 1) it is designed to balance distributions of covariates (observed and unobserved) such that, across repeated randomizations, assignment and covariates should be independent, 2) it is designed to allow assessment of the uncertainties of estimated treatment effects without requiring populations or models of outcomes (or linearity assumptions), and 3) it is a method of manipulating putatively causal variables in a way that is impersonal — and thus enhances the credibility of causal claims. Political scientists have recently become excited about experiments primarily for the first and third aims but have ignored the second aim. This article has discussed some of the benefits (and pitfalls) of the use of covariates in randomized experiments while maintaining a focus on randomization as the basis for inference.

While randomization allows statistical inference in experiments to match the causal inference, covariate imbalance can and does occur in experiments. Balance tests are designed to detect worrisome imbalances. One ought to worry about random imbalances when they are 1) large enough (and relevant enough to the outcome) that they should make large changes in estimates of treatment effects and 2) large relative to their baseline values such that interpretation of the treatment effect could be confused.

Small studies provide little information to help detect either treatment effects or imbalance. The null randomization distribution for a balance test in a small study ought to have larger variance than said distribution in a large study. The same observed imbalance will cast more doubt on the null of balance in a large study than it will in a small study: the observed value will be farther into the tail of the distribution characterizing the hypothesis for the large study than it will in the small study. The same relationship between small and large studies holds when the test focuses on the treatment effect itself. Thus, a p -value larger than some acceptance threshold for the null hypothesis of balance tells us that the imbalance is not big enough to cause detectable changes in assessments of treatment effects. A p -value smaller than some acceptance threshold tells us that the imbalance is big enough to cause detectable changes when we gauge the effects of treatment.

Given random imbalance, what should one do? Adjustment can help, but adjustment can also hurt. This paper showed (using a fake dataset built to follow a real dataset) a case in which adjustment can help and seems meaningful and a case in which adjustment does not seem meaningful, as well as an intermediate case. One point to take away from these demonstrations is that some imbalance can be so severe that real adjustment is impossible. Just as is the case in observational studies, merely using a linear model without inspecting the data can easily lead an experimenter to mislead herself — and problems could multiply when more than one covariate is adjusted for at a time. Rosenbaum (2002a)’s proposal for a randomization-based use of linear regression models is attractive in that covariance-adjusted confidence intervals for the treatment effect do not depend on a correct functional form for the regression model. In this paper all of the adjustment for median household income depended on a functional form assumption so the randomization-based

confidence interval was essentially infinite (signaling that the design of the study had no information available for such adjustment) while the model-based regression confidence interval, while much wider than the unadjusted interval, was bounded. Modern matching techniques may also help with this problem (see Keele, McConnaughey, and White 2008). In this paper precision was not enhanced by matching within blocks but matchings including median household income did not radically change confidence intervals for the treatment effect—and balance tests before and after matching readily indicated that the matchings did not balance median household income.

This paper has not engaged with some of the other circumstances in which covariate information is important for randomized studies. In particular, if outcomes are missing, then prognostic covariates become ever more important in experimental studies given their ability to help analysts build models of missingness and models of outcomes (Barnard et al. 2003; Horiuchi, Imai, and Taniguchi 2007). Thus, I have understated the value of collecting more information about the units one has. The trade-offs between collecting more information about units versus including more units in a study ought to be understood from the perspectives long articulated in the many textbooks on experimental design: simple random assignment of units to two treatments (treatment versus control) can be a particularly inefficient research design.

References

- Barnard, John, Constantine E. Frangakis, Jennifer L. Hill, and Donald B. Rubin. 2003. “Principal Stratification Approach to Broken Randomized Experiments: A Case Study of School Choice Vouchers in New York City.” *Journal of the American Statistical Association* 98(462): 299–324.
- Bowers, Jake, and Costas Panagopoulos. 2009. ““Probability of What?”: A Randomization-based Method for Hypothesis Tests and Confidence Intervals about Treatment Effects.”.
- Bowers, Jake, Ben B. Hansen, and Mark Fredrickson. 2008. Reproduction archive for: ‘Attributing Effects to A Cluster Randomized Get-Out-The-Vote Campaign’. Technical report University of Illinois and University of Michigan.
- Bowers, Jake, Mark Fredrickson, and Ben Hansen. 2009. *RIttools: Randomization Inference Tools*.
- Brady, Henry E. 2008. “Causation and Explanation in Social Science.” *Oxford handbook of political methodology* pp. 217–270.
- Cochran, William G., and Gertrude M. Cox. 1957. *Experimental Designs*. New York: John Wiley & Sons.
- Cox, David R. 1958. *The Planning of Experiments*. John Wiley.
- Cox, David R., and Nancy Reid. 2000. *The Theory of the Design of Experiments*. Chapman & Hall/CRC.
- Cox, David R., and Peter McCullagh. 1982. “Some Aspects of Analysis of Covariance (with discussion).” *Biometrics* 38: 541–561.
- Fisher, Ronald A. 1925. *Statistical Methods for Research Workers*. Edinburgh Oliver & Boyd.

- Fisher, Ronald A. 1935. *The Design of Experiments*. Edinburgh: Oliver and Boyd.
- Freedman, David A. 2008a. “On Regression Adjustments in Experiments with Several Treatments.” *Annals of Applied Statistics* 2(1): 176 – 196.
- Freedman, David A. 2008b. “On Regression Adjustments to Experimental Data.” *Advances in Applied Mathematics* 40(2): 180–193.
- Freedman, David A. 2008c. “Randomization Does Not Justify Logistic Regression.” *Statistical Science* 23(2): 237–249.
- Gadbury, G.L. 2001. “Randomization Inference and Bias of Standard Errors.” *The American Statistician* 55(4): 310–313.
- Gelman, Andrew, and Jennifer Hill. 2007. *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Cambridge University Press.
- Green, Donald P. 2009. “Regression Adjustments to Experimental Data: Do David Freedman’s Concerns Apply to Political Science?”
- Hansen, Ben, and Mark Fredrickson. 2010. *Optmatch: Functions for Optimal Matching, Including Full Matching*. R package 0.6-1 ed.
- Hansen, Ben B. 2004. “Full Matching in an Observational Study of Coaching for the SAT.” *Journal of the American Statistical Association* 99(September): 609 – 618.
- Hansen, Ben B. 2008. “Comment: The Essential Role of Balance Tests in Propensity-matched Observational Studies.” *Statistics in Medicine* 27(May 30): 2050 – 2054.
- Hansen, Ben B., and Jake Bowers. 2008. “Covariate Balance in Simple, Stratified and Clustered Comparative Studies.” *Statistical Science* 23(2): 219–236.
- Hansen, Ben B., and Jake Bowers. 2009. “Attributing Effects to A Cluster Randomized Get-Out-The-Vote Campaign.” *Journal of the American Statistical Association* 104(Sep): 873—885.
- Horiuchi, Yusaku, Kosuke Imai, and Naoko Taniguchi. 2007. “Designing and Analyzing Randomized Experiments: Application to a Japanese Election Survey Experiment.” *American Journal of Political Science* 51(3): 669–687.
- Imai, Kosuke, Gary King, and Clayton Nall. 2009. “The Essential Role of Pair Matching in Cluster-randomized Experiments, with Application to the Mexican Universal Health Insurance Evaluation (with discussion).” *Statistical Science* 24(1): 29–72.
- Imai, Kosuke, Gary King, and Elizabeth Stuart. 2008. “Misunderstandings among Experimentalists and Observationalists about Causal Inference.” *Journal of the Royal Statistical Society, Series A* 171(2): 1–22.
- Keele, Luke J., Corrine McConnaughy, and Ismail K. White. 2008. Adjusting Experimental Data. In *Experiments in Political Science Conference*.

- Leisch, Friedrich. 2002. Dynamic Generation of Statistical Reports Using Literate Data Analysis. In *Compstat 2002 - Proceedings in Computational Statistics*, ed. W. Haerdle, and B. Roenz. Heidelberg, Germany: Physika Verlag pp. 575–580.
- Leisch, Friedrich. 2005. *Sweave User Manual*.
- Neyman, Jerzy. 1990. “On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9 (1923).” *Statistical Science* 5: 463–480.
- Nickerson, David W. 2005. “Scalable Protocols Offer Efficient Design for Field Experiments.” *Political Analysis* 13(3): 233.
- Panagopoulos, Costas. 2006. “The Impact of Newspaper Advertising on Voter Turnout: Evidence from a Field Experiment.”
- Robins, James M. 2002. “[Covariance Adjustment in Randomized Experiments and Observational Studies]: Comment.” *Statistical Science* 17(3): 309–321.
- Rosenbaum, Paul R. 2002a. “Covariance Adjustment in Randomized Experiments and Observational Studies.” *Statistical Science* 17(3): 286–327.
- Rosenbaum, Paul R. 2002b. *Observational Studies*. Second ed. Springer-Verlag.
- Rosenbaum, Paul R. 2002c. “Rejoinder.” *Statistical Science* 17(August): 321–327.
- Rosenbaum, Paul R. 2010. *Design of Observational Studies*. New York: Springer.
- Rubin, Donald B. 1974. “Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies.” *Journal of Educational Psychology* 66(5): 688–701.
- Rubin, Donald B. 1990. “[On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9.] Comment: Neyman (1923) and Causal Inference in Experiments and Observational Studies.” *Statistical Science* 5(4): 472–480.
- Schochet, Peter. 2009. “Is Regression Adjustment Supported by the Neyman Model for Causal Inference.” *Journal of Statistical Planning and Inference* .
- Sekhon, Jasjeet S. 2008. “The Neyman-Rubin Model of Causal Inference and Estimation via Matching Methods.” *Oxford handbook of political methodology* pp. 271–.
- Senn, Stephen J. 1994. “Testing for Baseline Balance in Clinical Trials.” *Statistics in Medicine* 13(17): 1715–26.