

Matching for Adjustment and Causal Inference

Escuela de Invierno en Métodos y Análisis de Datos UCU-DCSP

Jake Bowers

jwbowers@illinois.edu

Online: <http://jakebowers.org/>

14–17 julio 2014

Overview

This class is an introduction to statistical adjustment using matching and propensity scores in the style pioneered by Rubin and Rosenbaum and currently in rapid development by many other methodologists across the social science and statistical disciplines. An important motivation for matching is to approximate an experimental design. And, since such a motivation arises from a desire to make transparent and defensible statements about causal relations, we will introduce the counterfactual conception of causal inference and the potential outcome formalization of these ideas. We will also spend some time on statistical inference (hypothesis testing, confidence interval creation) after the creation of a matched design. Finally, we will grapple with some of the questions that are current research topics in this area: When and how one can claim to have adjusted “enough”? How can we engage with concerns about unobserved confounds even if we have adjusted for what we observe?

Since methods of matching are rapidly developing in the methodology literature, we will here focus on the simplest and oldest form: post-stratification. The general concepts and work-flow should be transportable to more sophisticated methods of matched adjustment.

Goals and Expectations

This course aims to help you think about statistical adjustment using stratification and matching as compared to statistical adjustment using the linear model directly (adjustment by “residualization”).

The course ought to give you opportunities to practice producing matched designs for your data and to ask questions that puzzle you as you do this work.

The point of the course is to position you to do the future learning that is at the core of your work as an academic analyzing data.

This course does not delve deeply into the theories of causal inference, statistical inference, or algorithms at the heart of these methods of statistical adjustment. Rather, through practice using tools, I hope that your curiosity is awakened and you begin to read more broadly and understand more deeply on your own.

Expectations I assume some previous engagement with high school mathematics, probability and statistical computing in the R statistical computing environment. If you have not used R, you are welcome to take the class, but I encourage you to get a little experience with R before the first class session. Feel free to email me to ask for advice about how to practice with R before the class begins.

Participation We will be doing hands-on work. I plan to lecture very little and instead will hope to pose problems of statistical theory, research design, and data for you to solve at your computers. I anticipate that you’ll work in small groups, asking me and/or the group questions as you proceed. I will break away to draw on the board or demonstrate on my own computer now and then to clarify points or help you around particularly difficult tasks.

Computing We will be using R in class so those of you with laptops available should bring them. Of course, I will not tolerate the use of computers for anything other than class related work during active class time. Please install R (<http://www.r-project.org>) on your computers before the first class session. You may prefer to use R in the context of the Rstudio IDE (<http://www.rstudio.com/>).

Computing is an essential part of modern statistical data analysis — both for turning data into information and for conveying that information persuasively (and thus transparently and reliably) to the scholarly community. In this course we will pay attention to computing, with special emphasis on understanding what is going on behind the scenes. You will be writing your own routines for a few simple and common procedures.

Books We will use the Rosenbaum book as our primary source. The other books are useful for further study.

Required: Rosenbaum, P. R. (2010). *Design of Observational Studies*. Springer (pdf free to download from some university ip addresses or via university library springerlink subscriptions: <http://www.springerlink.com/content/978-1-4419-1212-1/contents/>)

Recommended: Becker, H. S. (1986). *Writing for Social Scientists: How to Start and Finish Your Thesis, Book, or Article*. University of Chicago Press

Berk, R. (2004). *Regression Analysis: A Constructive Critique*. Sage

Gelman, A. and Hill, J. (2007). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press (particularly chapters 9,10 and 23 see <http://www.stat.columbia.edu/~gelman/arm/>).

Morgan, S. L. and Winship, C. (2007). *Counterfactuals and Causal Inference: Methods and Principles for Social Research (Analytical Methods for Social Research)*. Cambridge University Press (See <http://www.wjh.harvard.edu/~cwinship/cfa.html> for some links and background reading)

Rosenbaum, P. R. (2002b). *Observational Studies*. Springer-Verlag, second edition (see <http://www-stat.wharton.upenn.edu/~rosenbap/index.html> for lots of papers and presentations).

Rubin, D. B. (2006). *Matched sampling for causal effects*. Cambridge University Press, Cambridge; New York

Schedule

Note: This schedule is preliminary and subject to change. We will spend roughly 4 hours together for three of the days and 3 hours on one of the days. I anticipate mixing group discussions of your questions from the readings with in-class work using your own laptops or those provided by the school to get practice creating, analyzing, and assessing matched designs.

1—July 14—Experiments, Potential Outcomes, and Treatment Effects

Questions and Reading: What is the point of experiments? What are the key characteristics of experiments? Why are experiments special?

[Kinder and Palfrey, 1993](#)

[Gerber and Green, 2012](#), Chap 1

How can we bolster the interpretability of our comparisons if we do not have an experiment?

[Rosenbaum, 2010](#), Chap 1

What do we mean by “causal inference”? Can we convince ourselves that experiments have special advantages like “unbiased estimation of averages of potential outcomes”?

[Gerber and Green, 2012](#), Chap 2

Extra Reading: [Gelman and Hill, 2007](#), Chap 9.0 – 9.3 (On potential outcomes and causal inference)

[Angrist and Pischke, 2009](#), Chap 2

[Holland, 1986](#) (on the Counterfactual/Manipulationist conception of causality)

Brady, 2002 (for a survey of other major conceptions of causality from the perspective of applied social science)

2—July 15—Adjustment by Simple and Complex, Algorithmic Stratification

Given the problems that arise from the use of the linear model, how can we use research design to approximate the randomized experiment? How would we assess whether a matched design adjusted enough or not? What does it mean to adjust enough? What actual steps do we take in order to flexibly express our ideas about what it means for “like to be compared with like.”

Questions and Reading: Why not use the linear model for adjustment? How do we know when we have adjusted enough to make a strong case for clear, interpretable comparisons?

Rosenbaum, 2010, Chap 6

Gelman and Hill, 2007, Chap 9.5–9.6

What is the basic intuition behind post-stratified designs and modern optimal versions of post-stratification that we call “matched designs”? What is “balance”? How do we use R to create and evaluate matched designs?

Rosenbaum, 2010, Chap 7–9, 13 (Especially Chap 7, Chap 8.6, Chap 9)

Extra: Rosenbaum, 2010, Chap 3

Hansen (2004)

Hansen (2011) for an example walk-through of a matched analysis including a discussion of missing data on covariates.

Ho et al. (2007) [esp. their discussion of model sensitivity, for example their Fig 2]

Three different ideas about balance testing: (1) Imai et al. (2008); (2) Sekhon (2007a)¹; (3) Hansen and Bowers (2008) Hansen (2008) or for a less mathematical version of the same argument (Bowers, 2011, §3).

3—July 16—Statistical Inference for Matched/Post-stratified Designs

Questions and Reading: Given a matched design, how can we produce tests of substantively meaningful hypotheses about the unobserved comparisons of potential outcomes that we call “causal effects?”

Rosenbaum, 2010, Chap 2

How can we produce confidence intervals for an estimate of an average treatment effect?

Dunning, 2012, Chap 6.1 and Appendix 6.1

Imbens and Rubin, 2009, Chap 17

Extra: Berk, 2004, Chap 4 on general requirements for statistical inference (i.e. what does it mean to do statistical inference at all, what are we inferring to?)

Lin (2011) provides some useful proofs supporting the idea that linear regression with “robust” HC2 standard errors provides a useful large-sample way to do statistical inference about average treatment effects.

Miratrix et al. (2012) teach us about statistical inference for average treatment effects when matching after experimental outcomes have been collected.

Freedman (2008b,a, 2007, 2006) Suggesting that even the large sample statistical inference from using linear regression in randomized experiments is biased. Also arguing that the Huber-White standard errors are not a good idea.

¹ <http://sekhon.berkeley.edu/papers/SekhonBalanceMetrics.pdf>

[Rosenbaum \(2002a\)](#) and [Bowers and Panagopoulos \(2011\)](#) showing how covariance adjustment is compatible with Fisher’s randomization inference (and thus can be unproblematic after matching).

[Schochet \(2009\)](#); [Green \(2009\)](#) Suggesting that in large samples these biases worried about by Freedman ought not to worry us.

[Imbens and Rubin, 2009](#), Chap 6–8 Suggesting, similarly to Green and Schochet, that regression is fine for statistical inference in experiments (and further suggesting the use of the Huber-White robust standard errors).

[Abadie and Imbens, 2004](#) suggesting that the bootstrap is not a good approach with matched designs.

For advanced reading on the latest in statistical theory for statistical inference for “matching estimators” (which include but are not restricted to post-stratified studies) see:

[Hansen \(2009\)](#) for theory using randomization-inference.

[Abadie and Imbens \(2009\)](#) for a large-sample, Normal theory approach.

4—July 17— Sensitivity Analysis: An observational study is not a randomized experiment.

Questions and Reading: We addressed concerns about variables that we do observe using matching. Although randomization addresses concerns about all background covariates (observed and unobserved), any non-randomized study may be criticized on the basis that it does not adjust for variables that were not observed.

Since an observational study (no matter how well matched) cannot adjust for unobserved confounders, how can we address concerns about such unobserved variables?

[Cornfield et al. \(1959\)](#)

[Hosman et al. \(2010\)](#)

Extra: [Rosenbaum, 2002b](#), Chap 4
[Rosenbaum, 2010](#), Chap 3, 14
[Imbens \(2003\)](#)

*—Other Topics

If we move quickly, or if the class has some organized preferences, we could change the syllabus to discuss some of these ideas.

Non-bipartite matching – Advances in Multivariate Matching: Beyond Binary Treatment [Rosenbaum, 2010](#), Chap 11; [Lu et al. \(2011\)](#); [Imai and van Dyk \(2004\)](#)

Longitudinal Matching – Advances in Multivariate Matching: Matching with Longitudinal Data [Rosenbaum, 2010](#), Chap 12

Stuff that was painfully left out but which is important

Here are just a few extra citations to launch self-study of aspects of matching which we did not cover in our class.

The class elected to focus on matching for longitudinal problems for the last class. We thus are unable to cover other approaches to matching that have been developed by political methodologists such as Genetic Matching ([Diamond and Sekhon, 2006](#); [Sekhon, 2007b](#)) or Coarsened Exact Matching ([Iacus et al., 2009, 2011](#)) or Balance Optimization Subset Selection ([Nikolaev et al., 2012](#)) or more fine-tuned versions of the optimal post-stratification that we consider in this class (for example, ([Zubizarreta, 2012](#))). Nor did we have time to engage with many other applied and theoretical topics in causal inference for observational studies such as the work establishing causal interpretation of the propensity score (cited in the Rosenbaum textbook), or the alternative approaches to causal inference based on weighting by

functions of the propensity score such as those arising from work by Jamie Robins (Glynn and Quinn, 2010), let alone alternative conceptualizations of causal relations such as those developed by Judea Pearl (Pearl, 2000) or the work on estimation by Heckman or bounding causal inferences by Manski.

References

- Abadie, A. and Imbens, G. (2004). On the Failure of the Bootstrap for Matching Estimators. *NBER, Unpublished Manuscript*.
- Abadie, A. and Imbens, G. (2009). Matching on the estimated propensity score.
- Angrist, J. and Pischke, J. (2009). *Mostly harmless econometrics: an empiricist's companion*. Princeton Univ Pr.
- Becker, H. S. (1986). *Writing for Social Scientists: How to Start and Finish Your Thesis, Book, or Article*. University of Chicago Press.
- Berk, R. (2004). *Regression Analysis: A Constructive Critique*. Sage.
- Bowers, J. (2011). Making effects manifest in randomized experiments. In Druckman, J. N., Green, D. P., Kuklinski, J. H., and Lupia, A., editors, *Cambridge Handbook of Experimental Political Science*, chapter 32. Cambridge University Press, New York, NY.
- Bowers, J. and Panagopoulos, C. (2011). Fisher's randomization mode of statistical inference, then and now. Unpublished manuscript.
- Brady, H. (2002). Models of causal inference: Going beyond the neyman-rubin-holland theory.
- Cornfield, J., Haenszel, W., Hammond, E., Lilienfeld, A., Shimkin, M., and Wynder, E. (1959). Smoking and lung cancer: Recent evidence and a discussion of some questions. *Journal of the National Cancer Institute*.
- Diamond, A. and Sekhon, J. (2006). Genetic matching for estimating causal effects: A general multivariate matching method for achieving balance in observational studies.
- Dunning, T. (2012). *Natural experiments in the social sciences: a design-based approach*. Cambridge University Press.
- Freedman, D. A. (2006). On the So-called "Huber Sandwich Estimator" and "Robust Standard Errors". *The American Statistician*, 60(4):299–302.
- Freedman, D. A. (2007). On regression adjustments in experiments with several treatments. *Annals of Applied Statistics (To Appear)*.
- Freedman, D. A. (2008a). On regression adjustments to experimental data. *Advances in Applied Mathematics*, 40(2):180–193.
- Freedman, D. A. (2008b). Randomization does not justify logistic regression. *Statistical Science*, 23(2):237–249.
- Gelman, A. and Hill, J. (2007). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press.
- Gerber, A. and Green, D. (2012). *Field experiments: Design, analysis, and interpretation*. WW Norton.
- Glynn, A. and Quinn, K. (2010). An introduction to the augmented inverse propensity weighted estimator. *Political Analysis*, 18(1):36.
- Green, D. P. (2009). Regression adjustments to experimental data: Do david freedman's concerns apply to political science? Unpublished Manuscript.
- Hansen, B. (2008). Comment: The essential role of balance tests in propensity-matched observational studies. *Statistics in Medicine*, 27(12).

- Hansen, B. (2009). Propensity score matching to recover latent experiments: diagnostics and asymptotics. Technical Report 486, Statistics Department, University of Michigan.
- Hansen, B. and Bowers, J. (2008). Covariate balance in simple, stratified and clustered comparative studies. *Statistical Science*, 23:219.
- Hansen, B. B. (2004). Full matching in an observational study of coaching for the sat. *Journal of the American Statistical Association*, 99:609.
- Hansen, B. B. (2011). Propensity score matching to extract latent experiments from nonexperimental data: A case study.
- Ho, D., Imai, K., King, G., and Stuart, E. (2007). Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Political Analysis*, 15:199–236.
- Holland, P. W. (1986). Statistics and causal inference (with discussion). *Journal of the American Statistical Association*, 81:945–970.
- Hosman, C. A., Hansen, B. B., and Holland, P. W. (2010). The sensitivity of linear regression coefficients' confidence limits to the omission of a confounder. *The Annals of Applied Statistics*, 4(2):849–870.
- Iacus, S., King, G., and Porro, G. (2009). Causal inference without balance checking: Coarsened exact matching. Retrieved September, 13:2010.
- Iacus, S., King, G., and Porro, G. (2011). Multivariate matching methods that are monotonic imbalance bounding. *Journal of the American Statistical Association*, 106(493):345–361.
- Imai, K., King, G., and Stuart, E. (2008). Misunderstandings among experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society, Series A*, 171(2):1–22.
- Imai, K. and van Dyk, D. A. (2004). Causal inference with generalized treatment regimes: Generalizing the propensity score. *Journal of the American Statistical Association*, 99(467):854–866.
- Imbens, G. (2003). Sensitivity to Exogeneity Assumptions in Program Evaluation. *The American Economic Review*, 93(2):126–132.
- Imbens, G. and Rubin, D. (2009). Causal inference in statistics. Unpublished book manuscript. Forthcoming at Cambridge University Press.
- Kinder, D. and Palfrey, T. (1993). On behalf of an experimental political science. *Experimental foundations of political science*, pages 1–39.
- Lin, W. (2011). Agnostic notes on regression adjustments to experimental data: reexamining freedman's critique. Unpublished manuscript.
- Lu, B., Greevy, R., Xu, X., and Beck, C. (2011). Optimal nonbipartite matching and its statistical applications. *The American Statistician*, 65(1):21–30.
- Miratrix, L., Sekhon, J., and Yu, B. (2012). Adjusting treatment effect estimates by post-stratification in randomized experiments. *JR Stat. Soc. Ser. B. Stat. Methodol.* To appear.
- Morgan, S. L. and Winship, C. (2007). *Counterfactuals and Causal Inference: Methods and Principles for Social Research (Analytical Methods for Social Research)*. Cambridge University Press.
- Nikolaev, A. G., Jacobson, S. H., Cho, W. K. T., Sauppe, J. J., and Sewell, E. C. (2012). Balance optimization subset selection (boss): An alternative approach for causal inference with observational data. Technical report, University of Buffalo.
- Pearl, J. (2000). *Causality : Models, Reasoning, and Inference*. Cambridge University Press, reprinted with corrections edition.
- Rosenbaum, P. R. (2002a). Covariance adjustment in randomized experiments and observational studies. *Statistical Science*, 17(3):286–327.

- Rosenbaum, P. R. (2002b). *Observational Studies*. Springer-Verlag, second edition.
- Rosenbaum, P. R. (2010). *Design of Observational Studies*. Springer.
- Rubin, D. B. (2006). *Matched sampling for causal effects*. Cambridge University Press, Cambridge; New York.
- Schochet, P. (2009). Is regression adjustment supported by the neyman model for causal inference. *Journal of Statistical Planning and Inference*.
- Sekhon, J. (2007a). Alternative balance metrics for bias reduction in matching methods for causal inference. *Survey Research Center, University of California, Berkeley*.
- Sekhon, J. (2007b). Multivariate and propensity score matching software with automated balance optimization: The matching package for r. *Journal of Statistical Software*, 10(2):1–51.
- Zubizarreta, J. (2012). Using mixed integer programming for matching in an observational study of kidney failure after surgery. *Journal of the American Statistical Association*, *Forthcoming*.