

# Econometrics III

Štěpán Jurajda

First Draft 1997

February 18, 2023

## Contents

<b>I</b>	<b>Introduction</b>	<b>8</b>
1	Causal Parameters and Policy Analysis in Econometrics . . . . .	9
2	A Check-list for Empirical Work . . . . .	12
3	A Preview of Identification Approaches <sup>1</sup> . . . . .	16
3.1	Control for $X$ (or $P(X)$ ) . . . . .	17
3.2	Group-Level Variation and Identification . . . . .	18
3.3	Exogenous Variation (IV) . . . . .	20
3.4	Identification Through Heteroscedasticity . . . . .	21
3.5	‘Natural’ Experiments . . . . .	22
3.6	Experimental Setup and Solution . . . . .	23
3.7	Difference in Differences . . . . .	24
3.7.1	Fixed Effects . . . . .	26
3.7.2	IV DD . . . . .	27
3.8	Regression Discontinuity . . . . .	28
3.9	AI/Big Data . . . . .	29
4	Reminder . . . . .	30
4.1	Note on Properties of Joint Normal pdf . . . . .	30
4.2	Testing Issues . . . . .	31
5	Deviations from the Basic Linear Regression Model . . . . .	35

---

<sup>1</sup>See <https://www2.bc.edu/arthur-lewbel/identification-zoo-final.pdf>

<b>II</b>	<b>Panel Data Regression Analysis</b>	<b>38</b>
6	GLS with Panel Data . . . . .	38
6.1	SURE . . . . .	39
6.2	Random Coefficients Model . . . . .	39
6.3	Random Effects Model . . . . .	42
7	What to Do When $E[\varepsilon   x] \neq 0$ . . . . .	44
7.1	The Fixed Effect Model . . . . .	44
7.2	Errors in Variables . . . . .	49
8	Inference in Panel Data Analysis . . . . .	52
8.1	Inference with Clustered Data and in “Difference in Differences” . . . . .	52
8.2	Hausman test . . . . .	62
8.3	Using Minimum Distance Methods in Panel Data . . . . .	63
8.3.1	The Minimum Distance Method . . . . .	63
8.3.2	Arbitrary Error Structure . . . . .	64
8.3.3	Testing the Fixed Effects Model . . . . .	65
9	Simultaneous Equations and IV Estimation . . . . .	67
9.1	Testing for exclusion restrictions and IV validity . . . . .	71
9.2	Regression Discontinuity . . . . .	73
9.3	Dealing with Weak Instruments . . . . .	76
10	GMM and its Application in Panel Data . . . . .	80
11	Dynamic Panel Data . . . . .	83
12	Other Regression Applications . . . . .	84
12.1	Quantile Regression . . . . .	84
12.2	The Reflection Problem . . . . .	86
12.3	Oaxaca-Blinder Decompositions . . . . .	89
12.4	Meta Analysis . . . . .	92
12.5	Misc Topics . . . . .	93
13	Nonparametrics . . . . .	93
13.1	Multidimensional Extensions and Semiparametric Applications . . . . .	95
<b>III</b>	<b>Qualitative and Limited Dependent Variables</b>	<b>96</b>
14	Qualitative response models . . . . .	96
14.1	Binary Choice Models . . . . .	97
14.1.1	Linear Probability Model . . . . .	97
14.1.2	Logit and Probit MLE . . . . .	98
14.1.3	The WLS-MD for Multiple Observations . . . . .	100

14.1.4	Panel Data Applications of Binary Choice Models . . . . .	100
14.1.5	Relaxing the distributional assumptions . . . . .	103
15	Duration Analysis . . . . .	105
15.1	Hazard Function . . . . .	105
15.2	Estimation Issues . . . . .	107
15.2.1	Flexible Heterogeneity Approach . . . . .	110
15.2.2	Left Censored Spells . . . . .	114
15.2.3	Expected Duration Simulations . . . . .	114
15.2.4	Partial Likelihood . . . . .	115
16	Multinomial Choice Models . . . . .	115
16.0.5	Unordered Response Models . . . . .	116
16.0.6	Ordered Response Models . . . . .	122
16.0.7	Sequential Choice Models . . . . .	123
17	Models for Count Data . . . . .	124
17.1	Threshold Models . . . . .	125
18	Limited Dependent Variables . . . . .	126
18.1	Censored Models . . . . .	126
18.2	Truncated Models . . . . .	129
18.3	Semiparametric Truncated and Censored Estimators . . . . .	130
19	Sample Selection . . . . .	132
19.1	Sampling on Outcome . . . . .	133
19.1.1	Choice-based sampling . . . . .	133
19.1.2	Endogenous Stratified Sampling . . . . .	134
19.2	Models with Self-selectivity . . . . .	135
19.2.1	Roy's model . . . . .	136
19.2.2	Heckman's $\lambda$ . . . . .	137
19.2.3	Switching Regression . . . . .	139
19.2.4	Semiparametric Sample Selection . . . . .	141
20	Program and Policy Evaluation . . . . .	142
20.1	Setup of the Problem . . . . .	143
20.1.1	Parameters of Policy Interest . . . . .	144
20.1.2	Experimental Solution . . . . .	145
20.2	Matching . . . . .	148
20.3	Local IV . . . . .	155
20.3.1	Local Average Treatment Effect . . . . .	155
20.3.2	Marginal Treatment Effect . . . . .	159

## Preamble

These lecture notes were originally written (before the Wooldridge textbook became available) for a 2nd-year Ph.D. course in cross-sectional econometrics. The goal of the course is to introduce tools necessary to understand and implement empirical studies in economics focusing on other than time-series issues. The main emphasis of the course is twofold: (i) to extend regression models in the context of panel data analysis, (ii) to focus on situations where linear regression models are not appropriate and to study alternative methods. Due to time constraints, I am not able to cover dynamic panel data models. Examples from applied work will be used to illustrate the discussed methods. Note that the course covers much of the work of the Nobel prize laureates for 2000. The main **reference textbook** for the course is *Econometric Analysis of Cross Section and Panel Data*, [W], Jeffrey M. Wooldridge, MIT Press 2002. Other useful references are:

1. *Econometric Analysis*, [G], William H. Greene. 5th edition, Prentice Hall, 2005.
2. *Analysis of Panel Data*, [H], Cheng Hsiao, Cambridge U. Press, 1986.
3. *Limited-dependent and Qualitative Variables in Econometrics*, [M], G.S. Maddala, Cambridge U. Press, 1983.
4. *Structural Analysis of Discrete Data and Econometric Applications* [MF], Manski & McFadden <[elsa.berkeley.edu/users/mcfadden/discrete.html](http://elsa.berkeley.edu/users/mcfadden/discrete.html)>
5. *Panel Data Models: Some Recent Developments*, [AH] Manuel Arellano and Bo Honoré <[ftp://ftp.cemfi.es/wp/00/0016.pdf](http://ftp.cemfi.es/wp/00/0016.pdf)>

I provide suggestions for reading specific parts of these additional references throughout the lecture notes, but these suggestions are always additional to already having read the relevant part of the Wooldridge textbook.

Below find *selected* suggested readings:<sup>2</sup>

### I.1 Introduction, causal parameters and policy analysis in econometrics

- Heckman, J.J. (2000) “Causal parameters and policy analysis in econometrics: A twentieth century perspective” *QJE* February 2000.

---

<sup>2</sup>You can find most of the listed papers as well as additional readings at [\\ftp\LECTURES\YEAR\\_2\EconometricsIII](http://ftp\LECTURES\YEAR_2\EconometricsIII).

- Geweke, J.F., Horowitz, J.L., Pesaran, M.H. (2006) “Econometrics: A Bird’s Eye View,” IZA DP No. 2458.
- Chattopadhyay, R., & Duflo, E. (2004). Women as policy makers: Evidence from a randomized policy experiment in India. *Econometrica*, 72(5), 1409-1443.
- Athey and Imbens (2017) “The State of Applied Econometrics: Causality and Policy Evaluation,” *Journal fo Economics Perspectives*, 31(2): 3–32.

#### II.4 Cases where residuals are correlated

- GLS: Deaton A. (1997) *Analysis of Household Surveys*, Chapter 2.<sup>3</sup>
- Panel Data: Hsiao C. and M.H. Pesaran (2004) “Random Coefficient Panel Data Models,’ IZA Discussion Paper no. 1236.

#### II.5 Cases where residuals and regressors are correlated

- Fixed effects: Ashenfelter O. and A. Kruger (1994) “Estimates of the Economic Return to Schooling from a New Sample of Twins,” *American Economic Review* 84: 1157-1173.
- Errors in variables: Griliches Z. and J. Hausman (1986) “Errors in Variables in Panel Data,” *Journal of Econometrics* 31:93-118.
- Difference in Differences:
  - La Ferrara, E., Chong, A., & Duryea, S. Soap Operas and Fertility: Evidence from Brazil. *American Economic Journal: Applied Economics* 4(4):1-31
  - Bertrand M., Duflo E. & Mullainathan S. (2004) “How - How Much Should We Trust Difference-in-Differences Estimates,” *Quarterly Journal of Economics* 119: 249-275.
  - Donald, S.G. and Lang, K. (2007) “Inference with Difference-in-Differences and Other Panel Data,” *The Review of Economics and Statistics* 89: 221–233.
  - Cameron, A. Colin, and Douglas L. Miller (2015) “A Practitioner’s Guide to Cluster-Robust Inference” *Journal of Human Resources* 50 (2): 317–372.

---

<sup>3</sup><http://www.worldbank.com/lsms/tools/deaton/index.htm>

## II.7-8 Simultaneity

- Kling, J. (1999) “Interpreting IV Estimates of the Returns to Schooling,” Princeton University, Industrial Relations Section WP 415.
- Hahn J. and J. Hausman (2002) “A New Specification Test for the Validity of Instrumental Variables,” *Econometrica* 70(1)163-189.

## III.12 Sample Selection

- The Reflection Problem: Waldinger F (2012) “Peer Effects in Science - Evidence from the Dismissal of Scientists in Nazi Germany” *Review of Economic Studies* 79(2):838-861
- Regression discontinuity: Lee, D. (2008). “Randomized experiments from non-random selection in US House elections”. *Journal of Econometrics*, 142(2), 675-697. Agrawal A, Goldfarb A (2008) “Restructuring Research: Communication Costs and the Democratization of University Innovation” *American Economic Review* 98(4):1578–1590
- Censored models and self-selection models: Heckman, J.J. (1979) “Sample Selection Bias as a Specification Error,” *Econometrica* 47:153-161.<sup>4</sup>

## III.13 Program Evaluation

- Difference-in-Differences based on Matching: Blundell, R. and Costa-Dias, M. (2000) “Evaluation Methods for Non-Experimental Data,” *Fiscal Studies* 21: 427–468.
- A unified framework for IV and selection models: Manning, A. (2003) “Instrumental Variables for Binary Treatments with Heterogeneous Treatment Effects: A Simple Exposition,”<sup>5</sup>
- Program Evaluation:  
DiNardo, J. and Lee, D.S. (2010) “Program Evaluation and Research Designs,” NBER WP No. 16016.

---

<sup>4</sup>See also Lewbel A. (2004) “Simple Estimators For Hard Problems: Endogeneity in Discrete Choice Related Models”

<sup>5</sup><http://econ.lse.ac.uk/staff/amanning/work/econometrics.html>

Heckman, J.J. (2010) “Building Bridges Between Structural and Program Evaluation Approaches to Evaluating Policy,” *Journal of Economic Literature* 48: 356–398.

- LaLonde, R. J. (1986). Evaluating the econometric evaluations of training programs with experimental data. *American Economic Review*, 604-620

#### IV.16 Introduction to nonparametric and semiparametric methods

- Kernel estimation and Local Linear Regression: *Applied Nonparametric Regression*, Härdle, Cambridge U Press, 1989.
- Censored and sample-selection models: Chay & Powell (2001) “Semiparametric Censored Regression Models”
- Carneiro, Hansen, Heckman (2003)<sup>6</sup>
- Imbens, G.W. (2003) “Semiparametric Estimation of Average Treatment Effects under Exogeneity: A Review,” UC Berkeley and NBER

---

<sup>6</sup><http://www.nber.org/papers/w9546>

# Part I

## Introduction

What do economists mean by “ $x$  affects  $y$ ”?

We traditionally ask about effects using the conditional expectation function of  $y$  given  $x$ . Why? Because it is the best predictor of  $y$  in the sense of minimizing the mean of squared errors. We also like expectation as a concept that speaks to the typical values of  $y$  (given a particular  $x$ ).

We typically quantify effects using regression analysis, which, originally, was a descriptive statistical tool.<sup>7</sup> *The* regression function, i.e., a conditional mean  $E[y|x] = \int_{-\infty}^{\infty} y dF(y|x)$ , has no behavioral meaning.

We can comfortably speak of  $x$  as being a *causal determinant* of  $y$  when we have (i) a theoretical model suggesting  $x$  causes  $y$ , and (ii) exogenous variation in  $x$  identifying the model (with identification being part of the model, ideally). When we then regress  $y = x\beta + \varepsilon$  to estimate  $\hat{\beta}$ , we can measure how much  $x$  causes  $y$ .<sup>8</sup> A causal variable  $x$  typically captures some treatment (policy) and we wish to answer policy-relevant “what if” questions. With an empirically identified structural model, one can answer an array of policy-relevant counterfactual questions. The three steps of developing a causal model, asking identification questions, and estimation/inference typically go hand-in-hand.

Economists typically define causal treatment effects of interest within the potential-outcome Rubin Causal Model (Holland, 1986), which considers two potential outcomes, one with and one without treatment:  $y_{1i} = \mu_1 + u_{1i}$  and  $y_{0i} = \mu_0 + u_{0i}$  where  $T \in \{0, 1\}$  denotes treatment.<sup>9</sup> We may want to know  $E[y_{1i} - y_{0i}] = \mu_1 - \mu_0$ , i.e. the treatment effect under random assignment, the *average treatment effect* (ATE), but we typically want to find out about

---

<sup>7</sup>See, e.g., Deaton (1997) p. 63.

<sup>8</sup>Often, we focus on the effect of one *causal* variable (for which we have an exogenous source of variation) and use other regressors as *control* variables.

<sup>9</sup>Imbens (2019, <https://arxiv.org/pdf/1907.07271.pdf>) argues that the recent alternative Acyclic Graph Approach (of *The Book of Why*, Pearl and Mackenzie, 2018) to causal inference does not have realistic application examples and does not sufficiently engage with estimation. He also discusses recent interactions between econometrics and machine learning techniques.



$E[y_{1i} - y_{0i}|T_i = 1]$ , the *average effect of treatment on the treated* (ATT), which equals  $ATE + E[u_{1i} - u_{0i}|T_i = 1]$ . The fundamental problem with estimating the ATT is that the data only provides  $E[y_{1i}|T_i = 1]$  and  $E[y_{0i}|T_i = 0]$  but not the “what-if” counterfactual  $E[y_{0i}|T_i = 1]$ .

$$\begin{aligned} E[y_{1i}|T_i = 1] - E[y_{0i}|T_i = 0] &= ATT + \{E[y_{0i}|T_i = 1] - E[y_{0i}|T_i = 0]\} \\ &= \mu_1 - \mu_0 + E[u_{1i} - u_{0i}|T_i = 1] + \\ &\quad \{E[u_{0i}|T_i = 1] - E[u_{0i}|T_i = 0]\}. \end{aligned}$$

The sample selection bias (the term in curly brackets) comes from the fact that the treatments and controls may have a different outcome if neither got treated. This model suggests that there is non-trivial heterogeneity in treatment effects that should be at the center of estimated models.<sup>10</sup>

## 1. Causal Parameters and Policy Analysis in Econometrics

Econometrics<sup>11</sup> differs from statistics in defining the identification problem (in terms of structural versus reduced form equations). “Cross-sectional” econometrics (as opposed to time-series) operationalizes Alfred Marshall’s (1890) comparative-statics controlled-variation idea (*ceteris paribus*) into its main notion of causality (compare to time-series analysis and its statistical Granger causality definition).<sup>12</sup> The ultimate goal of econometrics is to provide policy evaluation.

In the classical paradigm of econometrics, economic models based on clearly stated axioms allow for a definition of well-defined structural “policy invariant” parameters. Recovery of the structural models allows for induction of causal parameters and policy-relevant analysis.

This paradigm was built within the work of the Cowles Commission starting in the 1930s. The Commission’s agenda concerned macroeconomic Simultaneous Equation Models and was considered an intellectual success, but empirical failure due to incredible identifying assumptions.

---

<sup>10</sup>See equation (20.1).

<sup>11</sup>This introductory class is based on a recent survey by J.J. Heckman (2000). By the end of the course, make sure to come back to this introduction. By then, you could also read A. Deaton’s 2008 Keynes Lecture (NBER WP no. 14690) and the reply by G. Imbens (NBER WP no. 14896).

<sup>12</sup>For problems with causality, see <https://www.samsi.info/wp-content/uploads/2016/03/Halpern-causalitysurveytk-2.pdf>

A number of responses to the empirical failure of SEM developed, including first VAR and structural estimation methodology and later calibration, non-parametrics (sensitivity analysis), and the “natural experiment” approach. Let us in brief survey the advantages (+) and disadvantages (–) of each approach:

- VAR is “innovation accounting” time-series econometrics, which is not rooted in theory.
  - (+) accurate data description
  - (–) black box; may also suffer from non-credible identifying restrictions (as macro SEM); most importantly, results hard to interpret in terms of models.
- Structural estimation is based on explicit parametrization of preferences and technology. Here we take the economic theory as the correct full description of the data. The arguments of utility or technology are expressed as functions of explanatory variables. Given these  $i$ –specific arguments and an initial value of structural parameters, the optimization within the economic model (e.g., a nonlinear dynamic optimization problem) is carried out for each decision unit (e.g., unemployed worker). The predicted behavior is then compared with the observed decisions which leads to an adjustment of the parameters. Iteration on this algorithm (e.g., within MLE framework) provides the final estimates.
  - (+) ambitious
  - (–) computer hungry; empirically questionable: based on many specific functional form and distributional assumptions, but little sensitivity analysis is carried out given the computational demands, so estimates are not credible.<sup>13</sup>
- Calibration: explicitly rely on theory, but reject “fit” as the desired main outcome, focus on general equilibrium issues.
  - (+) transparency in conditional nature of causal knowledge
  - (–) casual in use of micro estimates, poor fit.
- Non-parametrics (as an extension of sensitivity analysis): do not specify any functional form of the “regression” in fear of biasing the results by too much unjustified structure.

---

<sup>13</sup>This criticism applies less in recent years, see Example 20.11. See Andrews, Gentzkow, and Shapiro (2017, QJE) for a simple way to assess the sensitivity of the estimates.

(+) transparency: clarify the role of distributional and functional form assumptions.

(−) non-parametrics is very data hungry.

- Natural experiment: search for situations in the real world that remind us of an experimental setup. Use such experiments of nature (as instrumental variables) to identify causal effects; keep theory at an intuitive level. Ultimately, argue that randomized controlled trials (RCT) are the gold standard of science (as in evidence-based medicine<sup>14</sup>) and where possible abandon work with non-randomized data to generate credibility.<sup>15</sup>

(+) internal validity and transparency: clear and mostly credible identification solving the assignment (selection into treatment) problem.

(−) often we focus on *whether* treatment works but do not learn *why* it works (mechanism); causal parameters are relative to IV (LATE<sup>16</sup>); it may be hard to cumulate knowledge or generalize; often not clear there isn't an important selection problem of participation in experiment.

We return to the issue of economic-theory-based structural-model estimation, which allows for the estimation of well-defined policy-relevant parameters within an ex ante policy evaluation, versus IV- or randomization-based ‘theory-free’ evaluation in Section 20 (see Heckman, 2010 JEL, for a recent bridge-building exposition with many key references). Of course, the best work uses experimental identification to estimate structural models. To make the best choice, consider the *Marschak’s Maxim*: estimators should answer *well-posed economic* problems

---

<sup>14</sup>It is hard to get acceptance on causality without experiments. The overwhelming correlation between smoking and lung cancer has been accepted as evidence of causality in absence of direct experimental evidence, but it took a very very long time. In medicine, only about a third of distinct *treatments* are supported by RCT evidence (as beneficial or likely to be beneficial) although most of medical *activity* (treatments used) is likely to be evidence-based. See <http://clinicalevidence.bmj.com/ceweb/about/knowledge.jsp>

<sup>15</sup>The RCT approach may lead one to abandon the use of economic theory. Further, Deaton (2009) argues that in some cases “...instrumental variables have moved from being solutions to a well-defined problem of inference to being devices that induce quasi-randomization.” Heckman (2010, JEL) sees the program evaluation literature as *defining* parameters of interest as summaries of experimental interventions. Indeed, Rubin and Holland argue that causal effects are defined only if an experiment can be performed. See also Deaton and Cartwright (2016, NBER WP No. 22595) and Rothstein and von Wachter (2016, NBER WP No. 22585).

<sup>16</sup>See Section 20.

with *minimal assumptions*.<sup>17</sup> Often, policy invariant *combinations* of structural parameters, not necessarily the whole structural model, are needed to answer well-posed economic questions.

Also, the fundamental problem of econometric policy evaluation is that to predict the future, it must be like the past, but the goal is to predict effects of a new policy, i.e. to predict a future that will not be like the past. Here, Marschak (1953) argues that predicting effects of future policy may be possible by finding past variation related to variation induced by new policy. The relationship between past and future variation is made using an economic model. Using this approach we may not need to know the full structural model to evaluate a particular policy. See Ichimura and Taber (2000).<sup>18</sup>

In this course, we will mostly remain within the classical paradigm and discuss parametric reduced-form econometric models. We will also occasionally touch on non-parametric and natural-experiment research and return to discussing causal inference when introducing the program evaluation literature in Section 20.

## 2. A Check-list for Empirical Work

Here are some tips for how to do (and present) empirical research, based largely on John Cochrane.<sup>19</sup>

---

<sup>17</sup>Marschak's Maxim is an application of Occam's Razor to policy evaluation.

<sup>18</sup>Marschak (1953) considers a monopolistic firm that experiments with output levels and observes profit as the outcome (this reduced form can be tabulated without knowledge of any structural parameters). Now, consider the government changing the tax rate on the monopolist's output. This changes the reduced form relationship of profit and output so it would seem that a new round of experimentation (tabulation) is necessary under the new policy for the government to predict the effects on profits. But one can estimate the demand function using the "old" data and then use an economic model (demand only depends on taxes through price) to predict the new profits. By using some aspect of a behavioral model one can exploit other types of variation to mimic a policy change.

Ichimura and Taber (2000) use this insight to provide a framework to link variation from a "natural experiment" to policy-equivalent variation. This requires stronger assumptions than just IV, but weaker compared to a fully structural model. For a related discussion see Deaton (2009).

<sup>19</sup>[http://faculty.chicagosb.edu/john.cochrane/research/Papers/phd\\_paper\\_writing.pdf](http://faculty.chicagosb.edu/john.cochrane/research/Papers/phd_paper_writing.pdf)

He also provides many useful writing and presentation tips, such as: get to the main result of the paper as fast as possible (both in a presentation and in the text of the paper), move everything that's not essential to the main story line of the paper to an appendix, use active and present tense, simple words and sentences, etc.

Based on this course, you should understand that the three most important things for empirical work are Identification, Identification, Identification. So, in your written work, describe your identification strategy clearly. In particular:

1. Explain what *economic mechanism* caused the dispersion in your right hand variables.
2. Explain what *economic mechanism* constitutes the error term. What things other than your right hand variable cause variation in the left hand variable? (For example, is it prediction error?) This will help you explain in *economic terms* why you think the error term is uncorrelated with the right hand variables.
3. Explain the *economics* of why your instruments are correlated with the right-hand-side variable (sufficiently so that they are not weak) and not with the error term.<sup>20</sup> If your instrument has no economics to it, only quasi-randomized assignment, is it not only *external* (see Remark 72) but also *exogenous*? Will the variation in treatment within the group assigned to treatment be random or related to the size of the effect?
4. Are you sure causality doesn't run from  $y$  to  $x$ , or from  $z$  to  $y$  and  $x$  simultaneously? (Do interest rates cause changes in housing demand or vice versa (or does the overall state of the economy cause both to change)? Does the size of police force affect crime rates or vice versa?)
5. Describe the source of variation in the data that drives your estimates, for every single number you present. (Think of fixed effects vs. random effects.)
6. Are you sure you are looking at a demand curve, not a supply curve? As one way to clarify this question, ask "whose behavior are you modeling?" Example: Suppose you are interested in how interest rates affect housing demand, so you run the number of new loans on interest rates. But maybe when housing demand is large for other reasons, demand for mortgages drives interest rates up. Are you modeling the behavior of house purchasers or the behavior of savers (how savings responds to interest rates)?

---

<sup>20</sup>Also, do you understand the difference between an instrument and a control? In regressing  $y$  on  $x$ , when should  $z$  be used as an additional variable on the right hand side and when should it be an instrument for  $x$ ? How would you use an ability measure when running a wage regression with education on the right-hand side (i.e., when you are worried about unobservable ability bias)?

7. Consider carefully what controls should and should not be in the regression. You may not want to include all the “determinants” of  $y$  on the right hand side. High  $R^2$  is usually bad — it means you ran  $y$  on a version of  $y$  and some other  $x$ . Do you want to condition on group fixed effects in situations when you do not need to remove group unobservables? If you do include them, you are asking about the effect of  $x$  within these groups, not across. (Think of wage regressions with industry dummies.) Similarly in matching exercises: economic theory should tell us what  $x$  need to be controlled for.<sup>21</sup>
8. For every parametric or distributional assumption (a column in your table of results based on traditional techniques) provide a more flexible (semi-parametric) version to assess the sensitivity to arbitrary assumptions.
9. Explain the *economic* significance of your results. Explain the economic magnitude of the central estimates,<sup>22</sup> not just their statistical significance.<sup>23</sup> In cases where decision makers are faced with deciding whether to implement a new policy or not, confidence intervals can be a more useful way of communicating uncertainty of point estimates (than p values are).<sup>24</sup>

---

<sup>21</sup>For selecting covariates in RCT analysis, see Remark 187.

<sup>22</sup>If the variables in the regression have a natural scale, such as dollars, and if you report the mean of the  $Y$  variable, one can easily see how large  $\beta$  is. When there is no intuitive meaning to  $Y$  units, one reports effects in standard deviations. For example, no one would know what a change in GPA of 0.3 points meant but saying that as a result of treatment it moved two standard deviations tells us something that students really moved in rankings. But if a variable is very tightly defined, reporting standard deviations can artificially make results look more important than they are.

Standardizing  $x$ s (mean zero and SD of 1 by subtracting the mean and dividing by the SD) leads to  $\beta$ s having “effect sizes” of a one standard deviation change in the regressor. The % of  $y$  explained = standardized  $\beta$  / SD of  $y$ .

<sup>23</sup>Ziliak and McCloskey (2008, U of Michigan Press) offer a million examples of how statistical significance is misused in empirical work and argue that one should primarily focus on economic size of effects. Of course, a large imprecisely estimated effect can be more important for the real world than a tiny precisely estimated one. Especially if several studies obtain a similar estimate: large, but statistically imprecise within an individual sample. It’s also clear that sampling-based confidence intervals (having to do with noise in small samples) capture only one of many sources of error in estimation. Yet, statistical significance can help differentiate among competing models. For more on this issue, see Section 4.2. Also, seeing  $p = 0.03$  for men and  $p = 0.13$  for women does not mean that different associations were seen in males and females; instead, one needs a p-value for the difference in the gender-specific associations.

<sup>24</sup><https://pubs.aeaweb.org/doi/pdfplus/10.1257/jep.35.3.157>

You should be able to link these suggestions to the discussion of the search for exogenous variation and the sensitivity to parametric and distributional assumptions provided in Section 1.

Ad point 7: Often we have multiple, potentially overlapping  $x$  proxies for, say, economic conditions or financial development of a country or a region, with little theory to guide selection of  $X$ s. Using a single proxy (a strategy used in corporate finance for example) would lead to attenuation bias. Lubotsky and Wittenberg (2006) suggest a strategy that minimizes attenuation bias by weighting betas of all available proxies. The advantage of the method is that it leaves the covariances between the error terms of these multiple proxies unrestricted. In comparison, principle component analyses of IV strategies require these to have zero covariances. For an application of the approach, see Ruhm (2018, NBER WP No. 24188) who runs long changes in drug-related mortality on corresponding changes in economic conditions and control variables. Dippel, Gold, Heblich, and Pinto (2017) have multiple measures of labor market conditions and only one IV. They aggregate their labor market indicators to a single index using principle component analysis. If the effects are concentrated in one principle component with a clear interpretation, this approach may be appealing. For adding intercorrelated  $X$ s in a sequence, see Gelbach (2016, JoLE). His decomposition nests IV, Oaxaca-Blinder, and the Hausman test.

Ad point 9: Economists often estimate log-linear (Mincerian, Euler, production-function, gravity) equations and present elasticities that are easy to interpret in terms of their magnitude. However, see Santos Silva and Tenreyro (2006) or Blackburn (2007) for how misleading such elasticities estimated by OLS from log-linearized models can be in the presence of heteroskedasticity. To see why, start with a constant-elasticity model  $y_i = \exp(\alpha)x_{1i}^{\beta_1}x_{2i}^{\beta_2}\epsilon_i$  or simply  $y_i = \exp(\alpha)\exp(x_i'\beta)\epsilon_i$  where  $E[\epsilon_i|x_i] = 1$  and, assuming that  $y_i > 0$ , arrive at  $\ln(y_i) = \alpha + x_i\beta + \ln(\epsilon_i)$ . Now, the Jensen's inequality ( $E[\ln y] \neq \ln E[y]$ ) implies that  $E[\ln(\epsilon_i)|x_i] \neq 0$ . This residual expectation will be a constant (i.e., will affect only the constant  $\alpha$  of the model) only under very strong assumptions on the distribution of  $\epsilon$  (because of the non-linear transformation of the dependent variable, the conditional expectation of  $\epsilon$  depends on the *shape* of the conditional distribution of  $\epsilon$ ). Specifically, with non-negative  $y_i$  the conditional variance of  $y$  (and of  $\epsilon$ ) vanishes when  $y$  is near zero. With heteroskedasticity (with the variance of  $\epsilon$  depending on  $x$ ), the expectation  $E[\ln(\epsilon_i)|x_i]$  will be a function of  $x$ , thus rendering OLS of the log-linearized model inconsistent. Solution? WNLLS on levels, i.e., minimize squares of  $(y_i - \exp(x_i'\beta))$ , but one needs an assumption on the functional

form of  $V[y_i|x_i]$  to be able to weight efficiently (see exercise 14.3). Santos Silva and Tenreyro (2006) propose a Poisson pseudo-maximum likelihood estimator,<sup>25</sup> Blackburn (2007) also suggests quasi likelihoods.<sup>26</sup> Note that often the reason for applying the  $\ln y$  transformation is that  $y$  has a probability mass at zero, but such models have an extra parameter that is generally not determined by theory but whose values have enormous consequences for point estimates, see NBER WP no. 30735, which argues for the use of two-part models (see, for example, Section 17.1), OLS on the untransformed  $y$ , and Poisson regressions.<sup>27</sup>

**Remark 1.** *Manning and Mullahy (2001; also NBER Technical Working Paper No. 246) has been cited almost 2,000 times; it considers various alternatives to log models.*

**Remark 2.** *When using test score data (e.g., PISA, TIMSS, IALS, etc.), which are fundamentally ordinal, not cardinal, transformations of said data can alter estimated coefficients; also, comparing coefficients across alternative tests is tricky. Penney (REStat 2017, Economics of Education Review in press) cracks both problems.<sup>28</sup>*

**Remark 3.** *Look up some empirical fallacies, such as the ecological fallacy,<sup>29</sup> the hot hand fallacy,<sup>30</sup> etc.*

### 3. A Preview of Identification Approaches<sup>31</sup>

Understanding variation in  $X$ . The variation giving rise to coefficient estimates should be linked to the problem studied.

---

<sup>25</sup>The WLS FOCs are  $\sum_{i=1}^n (y_i - \exp(x_i'\hat{\beta})) \exp(x_i'\hat{\beta}) x_i = 0$  while the Poisson FOCs drop the  $\exp(x_i'\hat{\beta})$  weight.

<sup>26</sup>When dealing with 0 values on the LHS with the  $\ln$  transformation, one often adds the value 1 as in  $\ln(y + 1)$ . Alternatively, the transformation  $\ln\left(y + \sqrt{y^2 + 1}\right)$  can accommodate zero  $y$  values while maintaining the elasticity interpretation of the estimated RHS coefficients.

<sup>27</sup>The ideal Stata Poisson command is `ppmlhdfc` written by Sergio Correia and not the native `xtpoisson`. The former allows clustering on arbitrary IDs (i.e., not only on the differenced-out FEs), it supports multi-way clustering, and `y` supports Stata's factor/interaction notation.

<sup>28</sup>See Bond and Lang (JPE, 2019) for a related application to happiness scales.

<sup>29</sup><https://www.kellogg.northwestern.edu/faculty/spenkuch/research/ei-iv.pdf>

<sup>30</sup>Miller and Sanjurjo (Econometrica). Miller and Sanjurjo (Journal of Economic Perspectives, 2019).

<sup>31</sup>See <https://www2.bc.edu/arthur-lewbel/identification-zoo-final.pdf>



**Example 3.1.** Suppose that you compare crime rates and deterrence across Czech districts. Specifically, you regress district crime rates on district police force size per crime, after controlling for a district fixed effect. But this differences in deterrence may actually largely come from different trends in regional crime rates combined with even and fixed distribution of police force size inherited from communism. So it's unclear to what extent we measure one way causality here.

**Example 3.2.** You want to regress age at first marriage on wages using RLMS, but you only see young women getting married if they still live with parents, which is itself a function of labor income.

Different sources of variation lead to different interpretation of (different estimates of) the coefficients.

**Example 3.3.** See Bell et al. (2002) for an application relating wages and unemployment in several dimensions: (i) aggregate time series, (ii) cross-sectional compensating wage differentials as a no-moving equilibrium, (iii) regional equilibrium conditional on time and regional fixed effects.

**Example 3.4.** You motivate a paper by comparing living with parents and employment uncertainty for young workers across Italy and Sweden, but then you estimate the effect using within-Italy variation. Is it  $X$  or  $\beta$ ?

Effective causal estimation methods in nonexperimental settings often combine modeling of the conditional mean of outcomes using regressions with covariate balancing (e.g., through matching, see Bang and Robins, 2005).

### 3.1. Control for $X$ (or $P(X)$ )

**Example 3.5.** Returns to education, ability bias and IQ test scores.

When is controlling for  $X$  enough to identify a causal effect? I.e., when is *selection on observables* plausible? (When is it plausible that conditional on  $X$ , assignment to treatment is as good as random?) Using regressions or matching.<sup>32</sup>

**Example 3.6.** If applicants to a college are screened based on  $X$ , but conditional on passing the  $X$  test, they are accepted based on a first-come/first-serve basis.

---

<sup>32</sup>IPWRA

### 3.2. Group-Level Variation and Identification

Often variation of interest in  $x$  does not occur across individuals but across groups of individuals (firms, regions, occupations). When using individual-level data with group-level variation in the variable of interest, one needs to correct standard errors to admit the actual number of degrees of freedom (dimension of the variation of interest).

Also, sometimes you ask why individuals belonging to the same group act in a similar way (Manski (1995)).

**Example 3.7.** *Does the propensity to commit crimes depend on the average crime rate in the neighborhood or on the average attributes of people living in the neighborhood or on some exogenous characteristics of the neighborhood like the quality of schools etc.? Or think of high-school achievement as another example.*

(Groups of) **movers vs. stayers:** Consider the effect of a union dummy (0/1 variable) in levels and in first differences:

$$y_{it} = UNION_{it}\beta + \epsilon_{it}$$

$$y_{it} - y_{it-1} = (UNION_{it} - UNION_{it-1})\beta + \Delta\epsilon_{it}$$

and note that only those who switch status between  $t$  and  $t - 1$  are used in the ‘differenced’ specification. A similar argument can be made when using aggregate data.<sup>33</sup> This strategy is thought of as being closer to causal evidence; it relies on “movers” — but are they exogenous?

**Example 3.8.** *Gould and Paserman (2002) ask if women marry later when male wage inequality increases. They use variation across U.S. cities in male wage inequality and marriage behavior and allow for city-specific fixed effects and time trends to establish causality. To write a paper like this, start with graphs of levels and changes, then condition on other  $X$  variables, check if female wage inequality has any effect (it doesn’t), and conclude. It is not clear where changes in male wage inequality come from, but one would presumably not expect these changes to be driven by a factor that would also affect marriage behavior.*

---

<sup>33</sup>For example, if you want to study the effect of part-time work on fertility, you can hardly run fertility on part-time status of individuals and pretend part-time status is assigned exogenously. But perhaps, if there is variation across regions and time in *availability* of part-time jobs, one could estimate a relationship at the aggregate level.

**Example 3.9.** *Movers are also necessary for identification of (multiple FE) models, such as wage regressions with worker and employer fixed effects. To make these models work, one needs substantial amount of movers, see NBER WP No. 27368 for bias adjustments for limited mobility.*

**Identification of National Policy Effects:** In case of national policy changes, within-country identifying variation is hard to come by while cross-country variation is often plagued by country-level unobservables. Manning (2001) studies a national increase in minimum wages by relating changes in employment before and after the minimum wage introduction to the fraction of low paid workers in the pre-minimum wage period. See also a classic paper by Card (1992) who considers the imposition of a federal minimum wage: the “treatment effect” varies across states depending on the fraction of workers initially earning less than the new minimum. Ahern and Dittmar (2012) compares the effect of the introduction of gender quotas for boards of directors in Norway on company performance by comparing companies that were farther away from the quota target and those that were closer before the introduction of the policy. Similarly for SOX and independent directors in Paligorova (2007), etc. She first shows that those that did not have independent board do indeed show a stronger increase in independence in comparison to those firm that did have independent boards as of before SOX. This is step 0 in all of program evaluation: establish that there is a program!

Cross-country indirect strategies: It is usually hard to use country-wide before/after and cross-country comparisons to identify national policy effects. See, e.g., the discussion of identifying effects of institutions in Freeman (1998). But avoiding the main identification issue and focusing on interactions of the main causal variable can shed some light on the direction and mechanism of the causal effect. Rajan and Zingales (1998, JF) study the finance-growth nexus. One should isolate the part of the variation in financial development that is unrelated to current and future growth opportunities, which are inherently unobservable. Tackling this identification problem at the country level is very difficult. So, Rajan and Zingales (1998) give up on the big question and provide qualitative evidence on causality using industry-country comparisons. They come up with an industry-specific index of the need for tapping the financial system (using external finance) and regress industry growth from a sample of countries on country and global-industry fixed effects as well as on the *interaction* between U.S. industry external finance dependence (EFD) and country financial development. Such regression asks whether industries predicted to be in more need of external finance grow

faster in countries with more developed financial markets, conditional on all (potentially unobservable) country- and industry-specific factors driving growth.

### 3.3. Exogenous Variation (IV)

You want to estimate  $\beta$  in  $y = X\beta + \varepsilon$  but  $E[\varepsilon|X] \neq 0$  because of endogeneity or measurement error. A valid instrument  $Z$  is correlated with  $X$  but not with  $\varepsilon$ . The  $R^2$  of the first stage should not be too high or too low. Where do you get such a variable? One solution is to find a “natural” experiment (more correctly quasi-experiment) which generates such variation and then rely on this one source alone (read Angrist and Krueger, 2001, for a reader-friendly exposition).<sup>34</sup>

**Example 3.10.** *Card (1993) estimates returns to schooling, which may be affected by ability endogeneity bias, using proximity to college as an instrument for education. You may think of the distribution of student distance from college as providing a quasi experiment that the regression is using. Ideally, you want to drop students randomly from helicopter. Is this case close enough? Whose effect are we estimating?*

**Example 3.11.** *Changes in wage structure, which occur in a supply-demand framework: “Women, War and Wages” by Acemoglu, Autor and Lyle. First, establish that there is a treatment—variation in draft causes differences in female labor supply. Second, ask whether there is an effect—of female labor supply on wage dispersion.*

**Example 3.12.** *In the Moving to Opportunity (MTO) experiment in Boston, vouchers supporting re-location of families were assigned at random, but compliance with the randomized experiment was imperfect. One uses the random assignment as IV for actual treatment to zoom in on the ‘compliers’ and to avoid the ‘defiers’ (and this also omits the ‘always takers’ and the ‘never takers’).*

Ever since the Imbens and Angrist (1994) Local Average Treatment Effect (LATE) interpretation of IV (2SLS), our understanding of IV has changed fundamentally. We return to this insight in more detail in Section 20.3.1.

**Remark 4.** *Other than econometric tests for IV validity, see below, there are also intuitive tests in situations when identification comes from some quasi-experiment. For example, ask whether there is an association between the instrument and outcomes in samples where there should be none.*

---

<sup>34</sup>For deep neural network IV, see <http://proceedings.mlr.press/v70/hartford17a/hartford17a.pdf>

**Remark 5.** *IV for binary treatments in LimDep models is discussed in Angrist (JBES, 2001).*<sup>35</sup>

**Example 3.13.** *See Berg 2007, IZA DP No. 2585 for a discussion of IVs, which derive from the time interval between the moment the agent realizes that they may be exposed to the policy and the actual exposure. Berg presents an economic model in which agents with larger treatment status have a stronger incentive to find out about the value of the IV, which invalidates the IV. In other words, the exclusion restriction is likely to be violated if the outcome depends on the interaction between the agent’s effort and his treatment status.*

**Example 3.14.** *One popular IV is the Bartik (shift-share) treatment intensity used widely in the immigration impact literature and in the research on import penetration effects. Import exposure to China in a local labor market (and time) is defined as the product of the local pre-existing industrial structure with the economy-wide Chinese import penetration by industry, which is then replaced by the corresponding figures from another developed country (to focus on the effect driven by supply shocks in China).*

**Example 3.15.** *Sometimes, we have one IV to study the effect of a ‘mediated’ mechanism. For example, Dippel, Gold, Heblich, and Pinto (2017) ask how import penetration from China affects voting in Germany through its effects on German labor markets (it may have other channels of effect, of course). Their Bartik IV  $Z$  affects treatment  $T$ , which then affects  $Y$ , the ultimate outcome, but also  $M$ , the labor market situation, which also affects  $Y$ . If one assumes that unobservables that confound the relationship between  $T$  and  $M$  affect  $Y$  only through  $T$  and  $M$ , this identifies the whole channel of effects, i.e., it allows the authors to ask what portion of the total effect of  $T$  on  $Z$  operates through  $M$ . This so-called mediated effect corresponds to the product of the effect of  $T$  on  $M$  and the effect of  $M$  on  $Y$  (from a regression of  $Y$  on both  $M$  and  $T$ ).*

### 3.4. Identification Through Heteroscedasticity

A completely different approach to identification working off second moments: Hogan and Rigobon (2002). Estimate returns to education when education is

---

<sup>35</sup>In Stata, `ivprobit` requires a continuous endogenous regressor; with binary endogenous  $X$ , the command is `biprobit`. But Stata’s “margins” command is slow, so an alternative is distribution-free LMP, citing Angrist and Lewbel.

endogenous by splitting the sample into two groups based on different covariance matrices. They suggest this strategy is stronger when compared to IV because IVs are weak and there is a lot of variance in heteroskedasticity, so one can use it to solve measurement error, simultaneity and omitted variable biases in cross-sectional data.

As an illustration consider a model for wages  $w$  and schooling  $s$

$$\begin{aligned}w_i &= \beta s_i + X_i \mu_1 + \epsilon_i \\s_i &= \alpha w_i + X_i \mu_2 + \eta_i.\end{aligned}$$

The covariance of the reduced form, which we can estimate,

$$\Omega = \frac{1}{(1 - \alpha\beta)^2} \begin{bmatrix} \nu_\epsilon + \beta^2 \nu_\eta & \alpha \nu_\epsilon + \beta \nu_\eta \\ \cdot & \alpha^2 \nu_\epsilon + \nu_\eta \end{bmatrix},$$

consists of 3 equations in 4 unknowns ( $\nu_\epsilon, \nu_\eta, \alpha, \beta$ ). Now, suppose you split the sample into two parts, which have empirically different  $\Omega$ . If the regression coefficients are stable across groups, suddenly you have 6 equations in 6 unknowns.

The crucial condition for this approach to be credible is to find a situation where coefficients are stable across sub-populations with different variances. Unlike in the natural-experiment literature, here it is harder to explain the economic meaning behind the identification.

### 3.5. ‘Natural’ Experiments

Meyer (1995) and Angrist and Krueger (1999): “Natural experiment” examine outcome measures for observations in treatment groups and comparison (control) groups that are not randomly assigned. In absence of randomization, we look for sources of variation that resemble an experimental design.

**Example 3.16.** *For example, when studying the effect of unemployment benefits on labor supply, it is hard to differentiate the effect of the benefits from the effect of past labor supply and earnings. So a quasi-experimental design would use changes in benefits applying to some groups but not others (benefits such as maternity benefits, unemployment insurance, workers’ compensation, Medicaid, AFDC) to define the treatment and control groups.*

**Example 3.17.** *Other examples of quasi-experimental variation: Vietnam-era draft lottery,<sup>36</sup> state-level minimum wage laws changes, large influxes of immigrants, family size effect on family choice and the delivery of twins, the variation in number of children coming from the gender sequence of children (preference for a boy), returns to education and quarter of birth (compulsory schooling), differential distance in effect of medical treatment (heart disease). Think of these events as providing an IV.<sup>37</sup>*

### 3.6. Experimental Setup and Solution

Consider a study of the effect of a training program where workers are randomized into and out of treatment (training). The effect of the program:  $y_{1i}$  is earning with training,  $y_{0i}$  is earnings without training. We only look at the population of eligible workers. They first choose to apply for the training program or not. We observe  $y_{1i}$  only when  $D_i = 1$  (the person applied for and took training) and observe  $y_{0i}$  only when  $D_i = 0$  (these are the so called eligible non-participants, ENPs). We want to know  $E[y_{1i} - y_{0i} | D_i = 1]$ ; however, the data only provides  $E[y_{1i} | D_i = 1]$  and  $E[y_{0i} | D_i = 1]$  is not observed—it is the counterfactual. This problem is solved by randomization: take the  $D = 1$  group and randomize into *treatment* ( $R = 1$ ) and *control* ( $R = 0$ ) group. Then construct the experimental outcome:  $E[y_{1i}^* | D_i^* = 1, R_i = 1] - E[y_{0i}^* | D_i^* = 1, R_i = 0]$ .<sup>38</sup>

However, experiments are costly, often socially unacceptable (in Europe), and people may behave differently knowing they are in an experiment (think of expanding medical coverage).<sup>39</sup>

**Remark 6.** *Mediation analysis is an exercise where one asks by what channel did a randomly assigned manipulation work. Mediation is present if the estimated effect of a randomized  $X$  on outcome  $Y$  gets smaller when one additionally controls*

---

<sup>36</sup> Angrist, J. D. (1990). Lifetime earnings and the Vietnam era draft lottery: evidence from social security administrative records. *American Economic Review*, 313-336.

<sup>37</sup> Similarly, if there is a random assignment to treatment, but imperfect compliance, the assignment indicator is the right IV for the treatment dummy. See, e.g. the Moving to Opportunity papers.

<sup>38</sup> This can be used as a benchmark for the accuracy of sample selection techniques that we need when we have no experiment, see Section 8.

<sup>39</sup> For a practical guide to randomization, see <http://www.povertyactionlab.com/papers/Using%20Randomization%20in%20Development%20Economics.pdf>  
For experimental panel data design, see <https://www.nber.org/papers/w26250>

for  $M$  (so one controls for both  $X$  and  $M$ ). This exercise has three potential shortcomings: (1) The mediator of interest correlates with other omitted mediators. (2) Non-linearities and measurement error can produce spurious mediation, and there may be reverse causality with the 'dependent variable' causing the change in the mediator. (3) If  $M$  and  $Y$  are correlated outside of the experiment (i.e., independently of  $X$ ), then randomization no longer allows one to run the regression with both  $X$  and  $M$ . We can leave other variables out of  $Y = \blacksquare + cX$  due to randomization,  $\hat{c}$  is unbiased. But including  $M$  as in  $Y = \blacksquare + bM + cX$  makes  $\hat{b}$  and  $\hat{c}$  biased since  $b$  picks up the causal effect of  $M$  but it also picks up the confounded effect of all things that correlate with  $M$  outside the experiment, and  $\hat{c}$  is biased downwards.<sup>40</sup>

### 3.7. Difference in Differences

A simple research design, referred to as “Differences,” compares one group before and after the treatment (i.e., employment before and after minimum wage increase):  $y_{it} = \alpha + \beta d_t + \varepsilon_{it}$ , where  $d_{it} \in \{0, 1\}$  is the dummy for the treatment group. The crucial assumption is that without treatment,  $\beta$  would be 0 (no difference in means of  $y$  for treatment and control (before and after) groups). So estimate of beta is just mean of  $y$  after minus mean of  $y$  before. If there are changes in other conditioning variables, add  $x'_{it}\gamma$ . However, there are often underlying trends and/or other possible determinants (not captured by  $x$ ) affecting the outcome over time, making this identification strategy rather weak.

Therefore, a popular alternative is the “Difference in differences” design, that is a before/after design with an untreated comparison group. Here, we have a treatment ( $j = 1$ ) and a comparison ( $j = 0$ ) group for both the before ( $t = 0$ ) and after ( $t = 1$ ) time period:

$$\begin{aligned} y_{it}^j &= \alpha + \alpha_1 d_t + \alpha^j d^j + \beta d_t^j + \gamma' x_{it}^j + \varepsilon_{it}^j \\ \beta_{DD} &= \bar{y}_1^1 - \bar{y}_0^1 - (\bar{y}_1^0 - \bar{y}_0^0). \end{aligned}$$

In other words, you restrict the model so that

$$E[y_i^1 | i, t] = E[y_i^0 | i, t] + \beta.$$

The main threat to this method is the possibility of an interaction between group and time period (changes in state laws or macro conditions may not influence all groups in the same way). Note that we must enforce  $\gamma$  to be the same

<sup>40</sup>See <http://datacolada.org/103> for details.



across  $j$  and that we consider  $x$  as control variable, while  $d_t^j$  is the causal variable. One typically tests for parallel trends in the period prior to treatment, which is suggestive of counterfactual parallel trends, but parallel pre-trends are neither necessary nor sufficient for the parallel counterfactual trends condition to hold as argued in Kahn-Lang and Lang (NBER WP no. 24857), who also argue that any DiD paper should address why the original levels of the experimental and control groups differed, and why this would not impact trends, and that the parallel trends argument requires a justification of the chosen functional form.

**Example 3.18.** *Famous studies: Card and Krueger (1994) NJ-PA minimum wage study or Card (1990) Mariel Boatlift study. While in the NJ-PA study, the comparison group is obvious, in the immigration paper, Card must select cities that will approximate what would have happened to Miami were there no Boatlift (resulting in a 7% increase in Miami labor force in 4 months). These cities better have similar employment trends before the immigration influx. But note: each study is really only one observation, see 8.*

The most important simple update for the DiDs methodology is the synthetic control approach developed by Abadie, Diamond, and Hainmueller (2010, 2014). In terms of the Mariel Boatlift paper, instead of comparing Miami to Houston or the simple average of Houston and Atlanta, one compares Miami to a weighted average of controls that is more similar to Miami than any single city (or a simple average of cities) would be.

**Remark 7.** *The best situation for the DD method is when*

- *the comparison group both before and after has a distribution of outcomes similar to that of the treatment group before treatment. This is important for non-linear transformations of the dependent variable (marginals differ based on the base);*
- $\widehat{\alpha}_1$  *is not too large (otherwise there are frequent changes all the time).*

**Example 3.19.** *Studies where  $t$  is not the time dimension: Madrian job-lock paper: How does insurance coverage affect the probability of moving between jobs? Hypothesis: those with both current coverage and a greater demand for insurance (because spouse doesn't have coverage at work, or greater demand for health care) should be less likely to change jobs. Let  $t = 0$  be low demand for*

insurance,  $t = 1$  high demand, and let  $j = 0$  denote uncovered workers, and  $j = 1$  covered workers. It is harder to assess interactions between  $j = 1$  and  $t = 1$  if  $t$  is something more amorphous than time. Does greater insurance demand have the same quantitative effect on the mobility of those with and without their own coverage even if health insurance were not an influence?

**Example 3.20.** *Treatments that are higher-order interactions: Treatment applies to only certain demographic groups in a given state and time. Do not forget to include first-order interactions when testing for the presence of second-order interactions! Gruber (1994) mandated maternity benefits paper: Treatment group: women of certain ages ( $k = 1$ ) in  $d = 1$  and  $t = 1$ .*

### 3.7.1. Fixed Effects

The difference in differences (DD) design is the basis of panel-data estimation with fixed effects. One runs these regressions when policy changes occur in time as well as across regions (states) of the US, Russia, etc.

**Example 3.21.** *Or consider the union status effect on wages discussed above: movers vs. stayers.*

**Example 3.22.** *Ashenfelter and Greenstone “Using Mandated Speed Limits to Measure the Value of a Statistical Life” In 1987 states were allowed to raise speed limits on rural interstate highways above 55 mph, 40 did (to 65 mph), 7 did not. You study the increase in speed (and time saved) and contrast this with the number of fatalities. Comparison groups are states that remained at 55 mph and other highways within states that went for 65 mph. They estimate*

$$\ln(\text{hours of travel})_{srt} = \beta \ln(\text{miles of travel})_{srt} + \gamma \ln(\text{fatalities})_{srt} + \alpha_{sr} + \eta_{rt} + \mu_{st} + \nu_{srt}$$

*but there is endogeneity problem in that people adjust travel speed to reduce fatalities when the weather is bad etc. So they use a dummy for having the 65 mph speed limit as an IV. In the end they get \$1.5m per life.*

**Remark 8.** *There is an alternative to using panel data with fixed effects that uses repeated observations on cohort averages instead of repeated data on individuals. See Deaton (1985) Journal of Econometrics.*

### 3.7.2. IV DD

Note that we often used the state-time changes as IV, instead of putting the  $d_{it}^j$  dummies on the RHS. Alternatively, in field experiments with imperfect compliance, we use the randomized assignment to IV for the actual treatment status.

**Example 3.23.** *State-time changes in laws generate exogenous variation in workers' compensation in Meyer et al. (AER) paper on injury duration. Instead of using  $d_{it}^j$  on the right-hand-side, include benefits as a regressor and instrument for it using the dummies  $d_{it}^j$ . This approach directly estimates the derivative of  $y$  w.r.t. the benefit amount.*

**Example 3.24.** *Unemployment Insurance effects on unemployment hazards (duration models). Meyer (1990) using state-time variation in benefits. Here we insert the benefits because who knows how to do IV in a nonlinear model.<sup>41</sup>*

**Example 3.25.** *Cutler and Gruber (1995) estimate the crowding out effect of public insurance in a large sample of individuals. They specify a model*

$$Coverage_i = \beta_1 Elig_i + X_i \beta_2 + \varepsilon_i$$

*As usual in U.S. research design, there is variation in state-time rules governing eligibility. Eligibility is potentially endogenous and also subject to measurement error. To instrument for  $Elig_i$  they select a national random sample and assign that sample to each state in each year to impute an average state level eligibility. This measure is not affected by state level demographic composition and serves as an IV since it is not correlated with individual demand for insurance or measurement error, but is correlated with individual eligibility.*

**Example 3.26.** *Angrist (1990) Vietnam era draft lottery—can't just use difference-in-differences in examining effect of veteran status on earnings (some people went anyway, and others avoided)—draft lottery numbers and military status are highly correlated, so use IV. Or quarter of birth study of Angrist and Krueger (1991).*

---

<sup>41</sup>But note that benefits tied to unemployment level, which is tied to duration! Juraйда and Tannery (2003) use within-state variation in unemployment levels to provide a stronger test of job search theory.

### 3.8. Regression Discontinuity

When assignment to treatment is (fully or partly) determined by the value of a covariate lying on either side of an (administrative) threshold, such assignment may be thought of as a natural experiment. Assume that the covariate has a *smooth* relationship with the outcome variable, which can be captured using parametric or semi-parametric models, and infer causal effects from discontinuity of the conditional expectation of the outcome variable related to assignment to treatment, which was determined by the ‘forcing’ variable being just below or just above the assignment threshold.<sup>42</sup>

**Example 3.27.** Angrist and Lave (1998) study of the class-size effect using the Maimonides rule: not more than 40 pupils per class. Class size is endogenous because of potential quality sorting etc. Assuming cohorts are divided into equally sized classes, the predicted class size is

$$z = \frac{e}{1 + \text{int}[(e - 1)/40]},$$

where  $e$  denotes the school enrollment. Note that in order for  $z$  to be a valid instrument for actual class size, one must control for the smooth effect of enrollment because class size increases with enrollment as do test scores.

**Example 3.28.** Matsudaira (JEcm) studies the effect of a school program that is mandatory for students who score on a test less than some cutoff level.

**Example 3.29.** Or think of election outcomes that were just below or just above 50%.

**Remark 9.** Clearly, there is some need for ‘local’ extrapolation (there is 0 common support), so one assumes that the conditional regression function is continuous.

**Remark 10.** Using, e.g., Local Linear Regressions, one estimates an ATT parameter, but only for those who are at the regression discontinuity and only for compliers.

---

<sup>42</sup>See the guide to practice of regression discontinuity by Imbens and Lemieux (2007). It is an NBER WP no. 13039 and also the introduction to a special issue of the *Journal of Econometrics* on regression discontinuity.

### 3.9. AI/Big Data

How does AI fit in all this? Several algorithms designed for use in highly dimensional data: Lasso (in Stata 16, each R package gives a different version), random forests, support vector machines,... What should be the role of machine learning in identification strategies used in economics? Machine learning provides predictions with good out-of-sample MSE. That is not what we typically care about in economics (we care about causal effects, not predictions). Also, highly dimensional data occur more likely in finance than in, say, labor economics. But AI can supply input to answering causal questions: choice of controls, finding relevant subsamples (Cengiz, Dube, Lindner and Zentler-Munro), or creating usable data from raw information (Buckles et al., 2020).

**Example 3.30.** *Instead of controlling for a highly dimensional  $X$  within a linear regression  $Y = \beta D + X'\gamma + \epsilon$  one can partial out  $X$  from both sides of the equation  $Y - \widehat{E[Y|X]} = \beta \left( D - \widehat{E[D|X]} \right) + \epsilon$ , where  $\widehat{E[.|X]}$  corresponds to a Lasso procedure. Lasso can select the  $X$  subset on each side of the equation, and the union of these two subsets can then be used as controls within a linear regression; see the post-double selection Lasso (Belloni, Chernozhukov, Hansen, 2014).*

## 4. Reminder

This section aims at reminding ourselves with some basic econometrics. We started the introduction with the conditional expectation function  $E[Y | X]$ . The law of iterated expectations decomposes a random variable  $Y$  into the conditional expectation function of  $Y$  given  $X$  and a residual that is mean independent of  $X$  (i.e.,  $E[\varepsilon|x] = 0$ ) and also uncorrelated with (orthogonal to) any function of  $X$ .

**Exercise 4.1.** *Prove or provide a counterexample for the following statements:*

- (a)  $Y \perp X \iff COV(X, Y) = 0$ . See also Exercise 4.2.
- (b)  $E[X | Y] = 0 \iff E[XY] = 0 \iff COV(X, Y) = 0$
- (c)  $E[X | Y] = 0 \implies E[Xg(Y)] = 0 \forall g(\cdot)$ . Is  $COV(X, Y) = 0$  ?
- (d)  $E[Y] = E_X[E_Y(Y | X)]$  and  $V[Y] = \underbrace{E_X[V_Y(Y | X)]}_{\text{residual variation}} + \underbrace{V_X[E(Y | X)]}_{\text{explained variation}}$ .

Why do we often run linear regressions (OLS)? Because when  $X$  and  $Y$  are jointly normal (see subsection 4.1) or when we work with a fully saturated model (with parameters for every combination of  $X$  values), then  $E[Y | X]$  is linear and the linear regression function is it. More importantly, OLS is also the best linear predictor (best approximation of  $E[Y | X]$  within the class of linear functions in terms of the minimum mean square error criterion). See also [W]1.

### 4.1. Note on Properties of Joint Normal pdf

In this note we show that the “true” regression function is linear if the variables we analyze are jointly Normal.<sup>43</sup> Let

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

**Exercise 4.2.** *Show that*

$$\Sigma_{12} = 0 \iff f(x | -) = f(x_1 | \mu_1, \Sigma_{11})f(x_2 | \mu_2, \Sigma_{22})$$

---

<sup>43</sup>Galton (1886), the fore-father of econometrics, studied height of parents and their children, two normally-distributed variables, and ran the first (linear) regression; he found “regression toward mediocrity in hereditary stature,” what we call today regression to the mean.

i.e., under normality, linear independence is equivalent to independence in probability.

**Theorem 4.1.**  $E[X_2 | X_1]$  is linear in  $X_1$ .

**Proof.** To get the conditional distribution of  $X_2 | X_1$  first find a linear transformation of  $X$  which block-diagonalizes  $\Sigma$  :

$$\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} I_1 & 0 \\ -\Sigma_{21}\Sigma_{11}^{-1} & I_2 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \\ \implies VAR \begin{pmatrix} X_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22.1} \end{pmatrix}$$

and  $X_1$  and  $Y_2$  are independent i.e.,  $Y_2 \equiv Y_2 | X_1 \sim N(\mu_2 - \Sigma_{21}\Sigma_{11}^{-1}\mu_1, \Sigma_{22.1})$ . Now note that  $X_2 = Y_2 + \Sigma_{21}\Sigma_{11}^{-1}X_1$  and conditioning on  $X_1$  the last term is a constant  $\implies X_2 | X_1 \sim N(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(X_1 - \mu_1), \Sigma_{22.1})$  or equivalently  $X_2 | X_1 \sim N(\mu_{2.1} + \Delta_{2.1}X_1, \Sigma_{22.1})$ . ■

**Remark 11.**  $\mu_{2.1} = \mu_2 - \Delta_{2.1}\mu_1$  is the intercept,  $\Delta_{2.1} = \Sigma_{21}\Sigma_{11}^{-1}$  is the regression coefficient, and  $\Sigma_{22.1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$  is the conditional covariance matrix which is constant i.e., does not depend on  $X_1$  (homoscedasticity).

**Remark 12.** OLS is attractive because it gives the minimum mean square error linear approximation to the conditional expectation function even when the linear model is misspecified.

## 4.2. Testing Issues

### Basic Principles<sup>44</sup>

---

<sup>44</sup>Does one need inferential tools (tests, confidence intervals) when working with large population data (i.e., the complete census)? Statistical analysis is about decision making (say about size of effects) under conditions of uncertainty. Consequently, one may want to use statistical procedures when using a complete census when the numbers are quite small in the 'table' being analysed. If the mortality rate of some small area is large, is it large in every year? Even when there is no sampling variation, there is 'natural variation'. Is the descriptive statistic (based on population) data a reliable measure for the underlying parameter (regional differences in mortality)? For arguments against the use of the inferential approach with population data, see Gorard (2013). Finally, the Bayesian approach has a very different perspective on inference and does not concern itself with repeated samples from hypothetical super-populations.

- Wald, Lagrange Multiplier, Likelihood Ratio. In class we provide a visualization of these in a graph. Note that they are asymptotically equivalent. So, obtaining different answers from each test principle may signal misspecification.<sup>45</sup>
- Specification tests: preview of Hansen and Hausman.
- Non-nested testing.

Generally, one needs to worry the about the false-positive rate and the power of a test.<sup>46</sup>

**Data Mining (Fishing, Overfitting) and Inference** The validity of empirical work is often questioned because researchers do not test their theory by running one regression, but they “data mine” or “p-hack”, even unknowingly.<sup>47</sup>

Today, the Cowless commission paradigm (Haavelmo, 1944; Popper, 1959) is abandoned in favor of more *interaction with data* (learning) so that the paradigm is merely used as a reporting style (Leamer, 1978).<sup>48</sup> Even our initial theory comes in part from previous empirical work. Because it’s clear that we all try different specifications, we report many alternative ones in our papers (sensitivity analysis) in order to convince the audience of the validity of our story.<sup>49</sup> Note

---

<sup>45</sup>Also, using the same level of significance for  $N = 100$  and  $N = 1,000,000$  is not right. With one million observations, you will reject any  $H_0$ . The Leamer’s (1978) rule for an F-test is to reject if  $F > \frac{N-k}{r} (N^{r/N} - 1)$ , where  $r$  is the number of restrictions and  $N - k$  is the number of degrees of freedom of the unrestricted error sum of squares. (See Kmenta, 2nd edition, p. 422). For example, consider  $k = 3$  and  $r = 2$ ; when  $N = 30$ , reject if  $F > 3.44$  (5% at 3.32) but with  $N = 1,000$ , reject if  $F > 6.93$  (5% at 3.00).

<sup>46</sup>For example, see <http://datacolada.org/62> for how the quadratic regression is not a diagnostic for u-shapedness. The proposed test consists of estimating two lines and testing for opposite slopes; the breakpoint is selected to increase the statistical power by strengthening the weaker line at the expense of the stronger line.

<sup>47</sup>[http://www.stat.columbia.edu/~gelman/research/unpublished/p\\_hacking.pdf](http://www.stat.columbia.edu/~gelman/research/unpublished/p_hacking.pdf)

<sup>48</sup>Let neither measurement without theory // Nor theory without measurement dominate // Your mind but rather contemplate // A two-way interaction between the two // Which will your thought processes stimulate // To attain syntheses beyond a rational expectation! Arnold Zellner [Zel96]

<sup>49</sup>But this will never be perfect. In Randomized controlled trials (RCT), one should post the research design before conducting the trial. Otherwise, you stop measuring when you reach a significant effect, or you look (data-mine) for the sub-group where treatment works, etc. If you measure the effect of  $x$  (treatment) on multiple outcomes (various  $ys$ ), one of the  $ys$  may seem affected by accident. Also, look up the Bonferroni’s correction (and Romano



that in econometrics we either test theory by means of estimation or use theory to identify our models (e.g., by invoking the Rational Expectations hypothesis in estimation of dynamic models in order to identify valid instruments).

A closely related problem of inference concerns *sequential testing*: While test properties are derived based on a one-shot reasoning, in practice we carry out a sequence of such tests, where the outcome of one test affects the next test (where both are based on the same data), invalidating the test properties. These concerns may be dealt with by setting aside a portion of the data before the start of the analysis and verifying the ‘final’ regression on this subset at the end of the analysis by means of a one-shot specification test. Another response is that you first have to make your model “fly” (i.e. achieve Durbin Watson =2) and only later can you go about testing it by, e.g., predicting out of sample. See Romano, Shaik, and Wolf (2010) on resampling. See also Gelbach (JoLE, 2016).

Another problem results from looking for mistakes in data and programs as long as we don’t like the results – as long as we do not confirm our expectation. There is no journal that would publish papers that fail to reject  $H_0$  (a paper with good motivation, good data and regressions, and insignificant estimates). There is clearly publication bias with regard to significance, sign, and size of estimated effects. Meta analysis drives this point home in most fields of science. Should results be a random sample? Ashenfelter, Harmon and Oosterbeek (1999) test for publication bias by running  $\widehat{\beta}_{IV}$  estimates from several studies on their standard error. In the medical RCT literature, people use a funnel plot of treatment effect against trial size for the same purpose.<sup>50</sup> Much work demonstrates the departure of published p-value distributions from uniform, i.e., spikes at 5% significance. So one is to be sceptical of results with  $p$  just below 0.05.<sup>51</sup> Andrews and Kasy (2019, AER) identify the conditional probability of publication as a function of a study’s result and propose solutions to selective publications. P-hacking (Elliott et al.

---

and Wolf, 2005, implemented in Stata) for multiple comparisons. And the dead salmon fMRI study. <https://blogs.scientificamerican.com/scicurious-brain/ignobel-prize-in-neuroscience-the-dead-salmon-study/> Together with clustering (see Remarks 23 and 24), this correction is brutal to statistical significance of treatment effects.

<sup>50</sup>For a related study, see IZA DP No. 7268. See also note n. 242. Also, see Ioannidis (August 2005 Issue of PLoS Medicine) for why most published research findings are wrong.

<sup>51</sup>A recent update from biology is here: Head et al. (2015) “The Extent and Consequences of P-Hacking in Science” PLoS Biol. See also Simmons, Nelson “False-Positive Economics: Simple Solutions to the Selective Reporting of Economics Research” Journal of Economic Perspectives. See also Brodeur et al. (2016, AEJ: Applied Econ). Finally, see also <http://datacolada.org/58> and “The power of bias in economics research” by Ioannidis, Stanley, and Doucouliagos.

2021, Ecm) may be more of a problem for IV and DiDs than for RCTs and RD designs (Brodeur, et al., 2020, AER).<sup>52</sup> Snyder and Zhuo (2018, NBER WP No. 25058) uncover huge drawer bias in robustness check tests.<sup>53</sup> Unfortunately, the null-result publication penalty is real.<sup>54</sup>

The growing problems with reproducibility of scientific results and the continuing tyranny of  $p < 0.05$  led the ASA to issue principles of p-value use in 2016 with the purpose of focusing more on size of effects, confidence intervals, and data choices:<sup>55</sup> p-values can indicate how incompatible the data are with a specified statistical model (similarly confidence intervals show effect sizes that are most compatible with the data under the given model), but they do not measure the probability that the studied hypothesis is true and  $p < 0.05$  should not be the only basis for reaching conclusions. (See note n. 242 for a discussion of why using statistical significance alone to claim conclusive research findings may be ill-founded. See also Abadie, NBER WP No. 24403.) By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.<sup>56</sup> Imbens<sup>57</sup> argues that in cases where decision makers are faced with deciding whether to implement a new policy or not, confidence intervals are a more useful way of communicating uncertainty of point estimates (than p values are). Also, statistical inference is not equivalent to scientific inference (Ziliak, 2019; Manski, 2019).

By 2019, *The American Statistician* had a special issue on p-value use and argued that we should abandon the use of the term “statistical significance”, but continue reporting continuous p-values (without dichotomization). They issued a number of recommendations including the reporting of p-values not just for difference from zero, but also for difference from a pre-specified alternative such as minimal important effect size. One is to report the probability that the hypothesis is true, one is to use second-generation p values (SGPV) and report false discovery rates (false positive risk—the probability the result occurred by chance<sup>58</sup>). Ziliak

---

<sup>52</sup>Brodeur, Cook, and Heyes (2020) find evidence for p-hacking and publication bias, but their conclusions are sensitive to adjusting for rounding errors (Kranz and Pütz).

<sup>53</sup>For recent extensions of standard sensitivity analyses see Andrews, Gentzkow and Shapiro, (2017) or Andrews and Oster (2019).

<sup>54</sup><https://www.cesifo.org/en/publications/2022/working-paper/null-result-penalty>

<sup>55</sup><http://retractionwatch.com/2016/03/07/were-using-a-common-statistical-test-all-wrong-statist>

<sup>56</sup>The difference between “significant” and “not significant” is not itself statistically significant (German and Stern, 2006).

<sup>57</sup><https://pubs.aeaweb.org/doi/pdfplus/10.1257/jep.35.3.157>

<sup>58</sup>Colquhoun (2019) suggests the use of a prior probability of effect being real of 0.5. <http://for-calc.ucl.ac.uk/>

and McCloskey (2008) offer several useful pieces of practical advice: (i) Adjust  $\alpha$  to sample size (see also Gannon, Pereira, and Polpo, 2019) and deal with the power of the test. (ii) Report coefficients in elasticity form or in a way that allows one to easily see the economic magnitude of the estimates. (iii) Do not drop from regression specifications statistically insignificant controls with economically large estimated effects. (iv) Present an independent simulation that would allow some perspective on whether the coefficient estimates are of reasonable magnitude.

**Remark 13.** *To analyze z-curves to look for mean power and signs of publication bias consider R packages `Z-curve.2.0` and `p-curve`. Different research designs (DinDs vs. experiments) may differ in how large a jump in  $p$  one can achieve for one unit of change in data/research design (<http://datacolada.org/91>).*

## 5. Deviations from the Basic Linear Regression Model

Here, we consider 3 main departures from the basic classical linear model: (a) when they occur, (b) what the consequences are, and (c) how to remedy them. This preview sets the stage for our subsequent work in panel-data and limited-dependent-variable (LIMDEP) estimation techniques.

- (i)  $V[\varepsilon_i | x_i] = \sigma_i^2 \neq \sigma_\varepsilon^2$ , i.e. the diagonal of the variance-covariance matrix is not full of 1s: (a) e.g., linear prediction vs.  $E[y | x]$  or heteroskedasticity,<sup>59</sup> (b) the inference problem of having underestimated standard errors and hence invalidating tests, (c) GLS based on assumed form of heteroskedasticity or the heteroskedasticity-consistent standard errors (White, 1980). The Huber-White idea is that you don't need to specify the usually unknown form of how  $V[\varepsilon_i | x_i]$  depends on  $x_i$ .<sup>60</sup> The method ingeniously avoids having to estimate  $N$  of  $\sigma_i^2(x_i)$  by pointing out that the  $k$  by  $k$  matrix  $\sum_{i=1}^N x_i x_i' \hat{\varepsilon}_i^2$ ,

---

<sup>59</sup>Arises all the time. For example when working with regional averages  $y_r = \frac{1}{N_r} \sum_{i=1}^{N_r} y_{ir}$  we have  $V(y_r) = \frac{1}{N_r} V(y_{ir})$ .

<sup>60</sup>Before White (1980), econometricians corrected for heteroscedasticity by WLS based on an explicit assumption about the form (source) of heteroscedasticity. Weighting observations can have a dramatic effect on estimated parameters under parameter heterogeneity, which is a more important issue than inference (see NBER WP no. 18859). The robust-standard-error approach of White (1980) takes the opposite approach: it does not move the estimates from the OLS benchmark and it corrects inference around this unchanged set of coefficients. The key issue is therefore whether weighting is based on some informative economic model assumption or whether it is driven by ad hoc unfounded heteroscedasticity assumptions.

where  $\widehat{\epsilon}_i$  is the OLS predicted residual<sup>61</sup>, converges to the true matrix with all of the  $V[\epsilon|x]$  so that

$$\widehat{V}(\widehat{\beta}_{OLS}) = \left( \sum_{i=1}^N x_i x_i' \right)^{-1} \sum_{i=1}^N x_i x_i' \widehat{\epsilon}_i^2 \left( \sum_{i=1}^N x_i x_i' \right)^{-1} .$$

(Here we also preview the Hausman test by comparing the OLS and Huber/White variance-covariance matrix. See [G]11.2,11.4, [W]4.2.3.

- (ii)  $COV[\epsilon_i, \epsilon_j | x_i, x_j] \neq 0$  : (a) time series or unobserved random effect (family effects), (b) possible inconsistency of  $\beta$  (for example when estimating  $y = \alpha + \epsilon$ , the asymptotic variance of  $\widehat{\alpha}$  does not converge to 0) , (c) GLS, Chamberlin's trick (see below).

**Example 5.1.** *GLS in spacial econometrics (see p.526 in Anselin, 1988) Here we present a way of parametrizing cross-regional correlation in  $\epsilon$ s (using analogy between time correlation coefficient and spacial correlation) and provide an example of how non-nested testing arises (e.g., with respect to how we specify the contiguity matrix summarizing prior beliefs about the spacial correlation) and what it means to concentrate the likelihood. Most importantly, we remind ourselves of how FGLS works in two steps. The first part of the panel data analysis (Section 6) will all be FGLS.*

- (iii)  $E[\epsilon_i | x_i] \neq 0$  : (a) Misspecification, Simultaneity, Lagged dependent variables and serial correlation in errors, Fixed effect model, Measurement error, Limited dependent variables; (b) inconsistency of  $\beta$ , (c) GMM/IV, non-parametrics, MLE.

In the first part of the course on regression techniques for panel data, we will first deal with (i) and (ii) by running various GLS estimators. Second we will also explore panel data strategies of dealing with (iii), chiefly the fixed-effect and the IV techniques. The second part of the course on LIMDEP models will all address (iii). Finally, within the program evaluation part, we will focus on estimating the causal effect of a binary treatment on outcomes—something that can be done using regressions, but we will cover more recent techniques. Similar to the (i) vs (ii) distinction (i.e., GLS vs. IV) in regression analysis, we will distinguish two

<sup>61</sup>Remember that with heteroscedasticity OLS still provides unbiased estimates of  $\beta$ s, so that  $\widehat{\epsilon} = y - x' \widehat{\beta}_{OLS}$  is also unbiased.

strands of methods depending on whether we require the conditional independence assumption to hold, i.e., the independence of potential outcomes and treatment assignment conditional on a set of observable covariates.

## Part II

# Panel Data Regression Analysis

Reading assignment: [H] 1.2, 2, 3.2 - 3.6, 3.8, 3.9.

## 6. GLS with Panel Data

So far we talked about cases when OLS fails to do its job and GLS fixes the problem, i.e. cases where the variance assumption is violated. Now, we are going to apply that reasoning in the panel data context.

The model we have in mind is

$$\begin{aligned}
 y_{it} &= x'_{it}\beta_i + \epsilon_{it} \text{ with } i = 1, \dots, N \text{ and } t = 1, \dots, T, \text{ or} & (6.1) \\
 y_i &= X_i \beta_i + \epsilon_i \text{ with } i = 1, \dots, N \text{ or} \\
 y_{T \times 1} &= \begin{matrix} T \times k & k \times 1 \end{matrix} \\
 y_{NT \times 1} &= \begin{matrix} \begin{bmatrix} X_1 & 0 & \cdots & 0 \\ 0 & X_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & X_N \end{bmatrix} & \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_N \end{pmatrix} & \begin{matrix} kN \times 1 \\ kN \times 1 \end{matrix} \end{matrix} + \epsilon \text{ or } y = \begin{matrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix} & \begin{matrix} NT \times k \\ NT \times k \end{matrix} \end{matrix} \beta + \epsilon
 \end{aligned}$$

where the covariance structure of  $\epsilon_{it}$  will again be of interest to us. In a panel model we can allow for much more flexible assumptions than in a time series or a cross-section.

**Remark 14.**  $N$  and  $T$  do not necessarily refer to number of individuals and time periods respectively. Other examples include families and family members, firms and industries, etc.

**Remark 15.** The number of time periods  $T$  may differ for each person. This is often referred to as unbalanced panel.

**Remark 16.**  $T$  is usually smaller than  $N$  and most asymptotic results rely on  $N \rightarrow \infty$  with  $T$  fixed.

The first question is whether we constrain  $\beta$  to be the same across either dimension. We cannot estimate  $\beta_{it}$  as there is only  $NT$  observations.

### 6.1. SURE

Suppose we assume  $\beta_{it} = \beta_i \forall t$ , that is for some economic reason we want to know how  $\beta$ s differ across cross-sectional units or F test rejects  $\beta_{it} = \beta \forall i, t$ .

If  $E[\varepsilon_{it} | x_{it}] = 0 \forall t$  and  $V[\varepsilon_{it} | x_{it}] = \sigma_{ii}^2$  and  $(x_{it}, \varepsilon_{it})$  is *iid*  $\forall t$  then we estimate  $\beta_i$  by running  $N$  separate OLS regressions. (Alternatively we can estimate  $y_{it} = x'_{it}\beta_t + \varepsilon_{it}$ .)

Now, if the covariance takes on a simple structure in that  $E(\varepsilon_{it}\varepsilon_{jt}) = \sigma_{ij}^2$  and  $E(\varepsilon_{it}\varepsilon_{js}) = 0$  there is cross-equation information available that we can use to improve the efficiency of our equation-specific  $\beta_i$ s. We have  $V[\varepsilon] = E[\varepsilon\varepsilon'] = \Sigma \otimes I_T \neq \sigma^2 I_{NT}$ , i.e. the  $\varepsilon$ 's are correlated across equations and we gain efficiency by running GLS (if  $X_i \neq X_j$ ) with  $\widehat{\sigma}_{ij}^2 = \frac{1}{T} \widehat{\varepsilon}_i' \widehat{\varepsilon}_j$  where the  $\widehat{\varepsilon}$  first comes from OLS as usual. Iterated FGLS results in MLE in asymptotic theory. In class we demonstrate the GLS formula for SURE and get used to having two dimensions in our data (formulas) and variance-covariance matrices.

### 6.2. Random Coefficients Model

What if we still want to allow *parameter* flexibility across cross-sectional units, but some of the  $\beta_i$ s are very uninformative. Then one solution may be to combine the estimate of  $\beta_i$  from each time series regression 6.2 with the 'composite' estimate of  $\beta$  from the pooled data in order to improve upon an imprecise  $\widehat{\beta}_i$  using information from other equations.<sup>62</sup> In constructing  $\beta$ , each  $\beta_i$  should then be given a weight depending on how informative it is.

To operationalize this idea, the RCM model allows the coefficients to have a random component (something typical for Bayesians, see [H 6.2.2]), i.e. we assume

$$y_i = X_i \beta_i + \varepsilon_i \quad (6.2)$$

$T \times 1$

where the error terms are well behaved, but

$$\beta_i = \underbrace{\beta}_{\text{nonstochastic}} + \nu_i \text{ with } E[\nu_i] = 0 \text{ and } E[\nu_i \nu_i'] = \Gamma.$$

OLS on 6.2 will produce  $\widehat{\beta}_i$  with  $V[\widehat{\beta}_i] = \sigma_i^2 (X_i' X_i)^{-1} + \Gamma = V_i + \Gamma$

**Exercise 6.1.** Show that the variance-covariance matrix of the residuals in the pooled data is  $\Pi = \text{diag}(\Pi_i)$ , where  $\Pi_i = \sigma_i^2 I + X_i' \Gamma X_i$ .

<sup>62</sup>Note that in a SURE system, each  $\widehat{\beta}_i$  is coming from equation by equation OLS.

$V_i$  tells us how much variance around  $\beta$  is in  $\widehat{\beta}_i$ . Large  $V_i$  means the estimate is imprecise.

Let  $\widehat{\beta} = \sum_{i=1}^N W_i \widehat{\beta}_i$ , where  $\sum_{i=1}^N W_i = I$ . The optimal choice of weights is

$$w_i = \left[ \sum_{j=1}^N (V_j + \Gamma)^{-1} \right]^{-1} (V_i + \Gamma)^{-1} \quad (6.3)$$

$\Gamma$  can be estimated from the sample variance in  $\widehat{\beta}_i$ 's ([G] p318). Note that  $\widehat{\beta}$  is really a matrix weighted average of OLS.

**Exercise 6.2.** Show that  $\widehat{\beta}$  is the GLS estimator in the pooled sample.

**Remark 17.** As usual we need asymptotics to analyze the behavior of  $\widehat{\beta}$  since weights are nonlinear. Also note that  $\widehat{\Gamma}$  is coming from the cross-sectional dimension, while  $\widehat{\beta}_i$  is estimated off time series variation.

Finally, we recombine  $\widehat{\beta}_i = A_i \widehat{\beta} + (I - A_i) \widehat{\beta}_i$  with optimal<sup>63</sup>  $A_i = (\Gamma^{-1} + V_i^{-1})^{-1} \Gamma^{-1}$ .

**Remark 18.** If  $E[\nu_i] = f(X_i) \implies E[\nu_i | X_i] \neq 0 \implies \widehat{\beta}_i$  is not consistent for  $\beta_i$ .

**Remark 19.** Today, the key (alternative) motivation for the use of Random Coefficients models comes from postulating the Rubin Causal Model (Holland, 1986) with alternative outcomes under and in absence of treatment  $T_i = 0, 1$ :

$$\begin{aligned} y_{i0} &= \alpha_0 + x_i' \theta_0 + u_{i0} \\ y_{i1} &= \alpha_1 + x_i' \theta_1 + u_{i1}, \end{aligned}$$

so that there is a distribution of causal treatment effects

$$\beta_i = y_{i1} - y_{i0} = \alpha_1 - \alpha_0 + x_i' (\theta_1 - \theta_0) + u_{i1} - u_{i0}. \quad (6.4)$$

We return to this topic in Section 20. One can also run the estimated  $\beta_i$  on person-specific external regressors.

<sup>63</sup>See [H] p.134 if you are interested in the optimality of  $A_i$ .



**Remark 20.** Arellano and Bonhomme (2012) study fixed-effect panel data models of the following type  $y_{it} = \alpha_i + \beta_i T_i + \delta' x_{it} + \epsilon_{it}$  with strictly exogenous regressors and exploit limited time dependence of time-varying errors for identification of variances of coefficients and for identification of their distributions (using non-parametric deconvolution techniques). They derive asymptotic properties of density estimators when covariates are discrete.<sup>64</sup>

**Remark 21.** As a digression, consider a situation when simple cross-sectional data are not representative across sampling strata, but weights are available to re-establish population moments.<sup>65</sup> First consider calculating the expectation of  $y$  (weighted mean). Then consider weighting in a regression. Under the assumption that regression coefficients are identical across strata, both OLS and WLS (weighted least squares) estimators are consistent, and OLS is efficient.<sup>66</sup> If the parameter vectors differ for each sampling strata  $s = 1, \dots, S$  so that  $\beta_s \neq \beta$ , a regression slope estimator analogous to the mean estimator is a weighted average of strata-specific regression estimates:

$$\widehat{\beta} = \sum_{s=1}^S w_s \widehat{\beta}_s, \quad \widehat{V}(\widehat{\beta}) = \sum_{s=1}^S w_s^2 \widehat{V}(\widehat{\beta}_s), \quad (6.5)$$

where  $w_s$  are scalar strata-specific weights, and where  $\widehat{\beta}_s$  is an OLS estimate based on observations from stratum  $s$ . In contrast, the WLS procedure applied to pooled

<sup>64</sup>They illustrate their approach by estimating the effect of smoking on birth weight; programs available at [http://www.cemfi.es/bonhomme/Random\\_codes.zip](http://www.cemfi.es/bonhomme/Random_codes.zip).

<sup>65</sup>For source see Deaton's *Analysis of Household Surveys* (1997, pp. 67-72).

<sup>66</sup>Although, see the Imbens and Hellerstein (1993/1999) study mentioned at the end of Section 10. And see NBER WP No. 23826 for weighting and external validity of RCTs. Also, when data are grouped, running the average  $Y$  for each value of  $X$  on  $X$  will replicate the microdata regression of  $Y$  on the grouped  $X$  when weighting by the number of  $Y$  observations for each grouped  $X$ . On the other hand, weighting using heteroscedasticity weights (including those used in the linear probability model of Section 14.1.1) is also questionable since the conditional variance model may be poorly estimated, thus messing up the hoped-for efficiency improvements. Unweighted regressions will always be the minimum mean square error approximation of the population conditional expectation, BLUE under homoscedasticity. Also weighting can reduce efficiency under some circumstances (Solon, Haider, and Wooldridge, 2015). For an opposing view, see Romano and Wolf (2017, JEcm) who argue for efficiency gains through the use of WLS (but with Huber-White standard errors, which allow for correct inference even if the variance-covariance structure is misspecified) in presence of conditional heteroscedasticity (i.e., when OLS is unbiased and consistent, but no longer BLUE). They show MC evidence this helps in small samples and asymptotical theory there can be valid inference in WLS even when the model for reweighting the data is misspecified.

data from all strata results in an estimator  $\widehat{\beta}_{WLS}$ ,

$$\widehat{\beta}_{WLS} = \left( \sum_{s=1}^S w_s X'_s X_s \right)^{-1} \sum_{s=1}^S w_s X'_s y_s = \left( \sum_{s=1}^S w_s X'_s X_s \right)^{-1} \sum_{s=1}^S w_s X'_s X_s \widehat{\beta}_s,$$

which is in general not consistent for the weighted average of  $\beta_s$ .<sup>67</sup>

### 6.3. Random Effects Model

Assuming  $\beta_{it} = \beta \forall i, t$  in Equation 6.1 one can impose a covariance structure on  $\epsilon$  and apply the usual GLS approach. The random effects model (REM) specifies a particularly simple form of the residual covariance structure, namely  $\epsilon_{it} = \alpha_i + u_{it}$  with  $E[\alpha_i \alpha_j] = \sigma_\alpha^2$  if  $i = j$  and is 0 otherwise. Other than that the only covariance is between  $u_{it}$  and  $u_{it}$  which is  $\sigma_u^2$ . We could also add a time random effect  $\lambda_t$  to  $\epsilon_{it}$ .

Given this structure  $V \equiv V \begin{pmatrix} \epsilon_i \\ \epsilon_t \end{pmatrix}_{T \times 1} = \sigma_u^2 I_T + \sigma_\alpha^2 e_T e_T'$ , where  $e_T$  is a  $T \times 1$  column of numbers 1. We write down  $E[\epsilon \epsilon']$  using  $V$  and invert  $V$  using the partitioned inverse formula to write down the GLS formula:

$$\widehat{\beta}_{GLS} = \left( \sum_{i=1}^N X'_i V^{-1} X_i \right)^{-1} \sum_{i=1}^N X'_i V^{-1} y_i \quad (6.6)$$

The GLS random effects estimator has an interpretation as a weighted average of a “within” and “across” estimator. We show this in class by first skipping to the fixed effect model to describe the within estimator. Then we return to the above GLS formula, re-parametrize  $V^{-1}$  using the matrix  $Q = I_T - \frac{1}{T} e_T e_T'$ , which takes things in deviation from time mean, and gain intuition by observing the two types of elements inside the GLS formula: (i) the “within” estimator based

---

<sup>67</sup>Neither is OLS, of course. The WLS estimator is consistent for  $\beta$  if the parameter variation across strata is independent of the moment matrices and if the number of strata is large (see, e.g., Deaton, 1997, p. 70). Neglecting coefficient heterogeneity can have serious consequences; for example it can result in significant estimates of incorrectly included regressors and bias other parameters even if the erroneously included variables are orthogonal to the true regressors (Pesaran et al., 2000). See NBER WP 18859 for an overview of why we (should not) weight in regression. One should contrast coefficients and standard errors from both weighted and unweighted LS, but this is rarely done. If OLS and WLS do lead to dramatically different coefficient estimates, it signals either parameter heterogeneity (misspecification) or endogenous sampling (see Section 19.1.1).

on deviations from mean  $x_{it} - \bar{x}_i$  and (ii) the “across” estimator working off the time averages of the cross-sectional units, i.e.  $\bar{x}_i - \bar{x}$ . Treating  $\alpha_i$  as random (and uncorrelated with  $x$ ) provides us with an intermediate solution between treating  $\alpha_i$  as being the same ( $\sigma_\alpha^2 = 0$ ) and as being different ( $\sigma_\alpha^2 \rightarrow \infty$ ). We combine both sources of variation: (i) over time within  $i$  units and (ii) over cross-sectional units.

As usual, the random effects GLS estimator is carried out as FGLS (need to get  $\widehat{\sigma}_u^2$  and  $\widehat{\sigma}_\alpha^2$  from OLS on within and across dimensions).

**Remark 22.** *Of course, one does not have to impose so much structure as in REM: (i) one can estimate the person specific residual covariance structure, see the next Remark, and (ii) one can use minimum distance methods and leave the structure of error terms very flexible (see section 8.3.2).*

**Remark 23.** *Cross-sections with group-level variables can be thought of as panel data. If you are interested in the effect of the group-level variable, you need to admit that it does not vary independently across individual observations.<sup>68</sup> This is done by adjusting standard errors by clustering: Apply the White (1980) unconditional heteroskedasticity idea while allowing for both unconditional heteroskedasticity as well as for correlation over time within a cross-sectional unit (or a group of individuals, a cluster  $g$  subscript replacing the individual  $i$  subscript) within the standard “sandwich matrix” (Liang and Zeger, 1986).<sup>69</sup>*

$$V_{clu}(\widehat{\beta}_{OLS}) = \left( \sum_{i=1}^N X_i' X_i \right)^{-1} \sum_{i=1}^N X_i' \widehat{\varepsilon}_i \widehat{\varepsilon}_i' X_i \left( \sum_{i=1}^N X_i' X_i \right)^{-1}. \quad (6.7)$$

*Note that the averaging across clusters (individuals in this case) that makes the middle  $k \times k$  matrix accurate requires the number of clusters to be large (see Section 8). Also note that conditioning on cluster-level fixed effects (the topic of the next section) does not control for all of the within-cluster correlation.*

---

<sup>68</sup>A haiku by Keisuke Hirano:  
*T-stat looks too good  
 Try clustered standard errors—  
 Significance gone*

<sup>69</sup>This is implemented in Stata regressions using the `vce(cluster clustvar)` option. See 10.7 and <http://www.stata.com/support/faqs/stat/cluster.html>. Also see Wooldridge (2003, AER).

**Exercise 6.3.** Does  $V_{clu}(\widehat{\beta}_{OLS})$  have to be positive definite? What is its rank? Note that its rank equals that of the middle matrix in the “sandwich”, which can be expressed as  $V'V$  where  $V$  is a  $k \times G$  matrix (where  $G$  is the number of clusters). Write out  $V$  to conclude.

**Remark 24.** Along similar lines, one needs to adjust inference for multi-stage survey design (such as selecting randomly villages and then randomly households within villages—unobservables will be related within a village). There is an extensive set of commands for such adjustments available in *Stata*: see `svy` commands.

**Remark 25.** Of course, group-level variables need not be just exogenous  $x$  controls. Group-level IVs are discussed in Remark 74. Sometimes we also ask why individuals belonging to the same group act in a similar way and whether they reflect each other’s behavior—in that case, the group-level right-hand-side variable is the average of  $y$  for the group. This makes for some thorny identification issues; see Manski (1995), Durlauf (2002) and Brock and Durlauf (2001).

## 7. What to Do When $E[\varepsilon | x] \neq 0$

### 7.1. The Fixed Effect Model

One of the (two) most important potential sources of bias in cross-sectional econometrics is the so called heterogeneity bias arising from unobserved heterogeneity related to both  $y$  and  $x$ .

**Example 7.1.** Estimation of the effect of fertilizer on farm production in the presence of unobserved land quality; an earnings function and schooling when ability is not observed, or a production function when managerial capacity is not in the data, imply possibility of heterogeneity bias.

If we have valid IVs (exclusion restriction), we can estimate our model by 2SLS. If we have panel data, however, we can achieve consistency even when we do not have IVs available. If we assume that the unobservable element correlated with  $x$  does not change over time, we can get rid of this source of bias by running the fixed effect model (FEM). This model allows for an individual specific constant, which will capture all time-constant (unobserved) characteristics:

$$y_{it} = \alpha_i + x'_{it}\beta + \epsilon_{it} \quad (7.1)$$

When  $T \geq 2$  the fixed effects  $\alpha_i$  are estimable, but if  $N$  is large, they become nuisance parameters and we tend to get rid of them: by estimating the model on data taken in deviation from the time mean or by time differencing.

To summarize, the FEM is appropriate when the unobservable element  $\alpha$  does not vary over time and when  $COV[\alpha_i, X_i] \neq 0$ . This nonzero covariance makes the  $\widehat{\beta}_{OLS}$  and  $\widehat{\beta}_{GLS}$  inconsistent. We'll come to the testing issue in section 8.

Suppose  $x'_{it} = (w_{it}, z_i)$  and partition  $\beta$  appropriately into  $\beta^w$  and  $\beta^z$ . In this case note that we cannot separately identify  $\beta^z$  from  $\alpha_i$ . This shows that when we run the fixed effect model,  $\widehat{\beta}$  is identified from individual variation in  $X_i$  around the individual mean, i.e.  $\widehat{\beta}$  is estimated off those who switch (change  $x$  over time).  $\widehat{\alpha}_i$ 's are unbiased, but inconsistent if  $T$  is fixed. Despite the increasing number of parameters as  $N \rightarrow \infty$ , OLS applied to 7.1 yields consistent  $\widehat{\beta}^w$  because it does not depend on  $\widehat{\alpha}_i$ . To see this solve the following exercise.

**Exercise 7.1.** Let  $M_D = I_{NT} - D(D'D)^{-1}D'$ , where

$$D = \begin{bmatrix} e_T & 0 & \dots & 0 \\ 0 & e_T & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & e_T \end{bmatrix} \quad \text{and} \quad e_T = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}.$$

Using the definition of  $M_D$  show  $\widehat{\beta}^w$  is estimated by a regression of  $y_{it} - \bar{y}_i$  on  $w_{it} - \bar{w}_i$ , where  $\bar{w}_i = \frac{1}{T} \sum_{t=1}^T w_{it}$ .<sup>70</sup>

**Remark 26.** While the LSDV (the Least Square Dummy Variable estimator—the FEM OLS with all of the fixed effects estimated) and the “within” FEM estimator give the same slope  $\beta$  coefficients, they lead to different standard errors in small samples due to different degrees-of-freedom correction. They are based on different assumptions:  $E[\epsilon_{it} - \bar{\epsilon}_i \mid w_{it} - \bar{w}_i] = 0$  versus  $E[\epsilon_{it} - \epsilon_{it-1} \mid w_{it} - w_{it-1}] = 0$ . For small  $T$  the average  $\bar{w}_i$  is not a constant, but a r.v.; that is why  $E[\epsilon_{it} \mid w_{it}] = 0$  is no longer enough and we need  $E[\epsilon_{it} - \bar{\epsilon}_i \mid w_i] = 0$ .

**Remark 27.** Of course, we may also include time dummies, i.e. time fixed effects. We may also run out of degrees of freedom.

<sup>70</sup>This is an application for the Frisch-Waugh-Lovell theorem.

**Remark 28.** *There is an alternative to using panel data with fixed effects that uses repeated observations on cohort averages instead of repeated data on individuals. See Deaton (1985) Journal of Econometrics.*

**Remark 29.** *While effects of time-constant variables are generally not identified in fixed effects models,<sup>71</sup> there are several tricks. One can estimate the change in the effect of these variables (Angrist, 1995 AER). Greiner and Rubin (2011, REStat) use changing perception of immutable characteristics. For fixed effects strategies allowing one to estimate the effect of time-constant variables in fixed-effect non-linear panel data models, see Honore and Kesina (2017). In the presence of a noisy proxy for a key time-constant variable, which is absorbed in the FE, panel data can still be used to provide useful information on the effect of the time-constant variable (M. Běln, CERGE-EI WP No. 624).*

**Remark 30.** *Two-way fixed effects: There could be several fixed effects at different data dimensions, e.g., worker as well as firm time-constant unobservables as well as match-specific fixed effects for each combination of worker and firm. The estimation was thought to be technically difficult<sup>72</sup> until Nikolas Mittag came up with a simple algorithm to implement the inversion of the required matrices for such a multiple-dimension fixed-effect specification, to be installed by running `ssc install twfe` in Stata.<sup>73</sup> See NBER WP no. 27368 for biases (and bias-correction methods) that arise because of limited mobility of workers across firms. For a two-way FEs application of recently developed heterogeneity-robust difference-in-differences methods, see IZA DP No. 14682.*

**Remark 31.** *The two-way FE model can be biased in presence of heterogeneous and dynamic treatment effects; it can fail to estimate averages (even weighted averages) of unit-specific events (unless the unit-specific coefficients are statistically independent of all the other terms in the model, which is unlikely) and can fall outside of the range of the unit-specific true coefficients. However, if there is a unit (state) that is not affected at all, then the TWFE estimator is still problematic, but it is possible to construct a DiDs estimator that is centered around the*

---

<sup>71</sup>I.e., we cannot estimate the effect of gender (except for sex change operations, which should make clear how important selection of movers is for this strategy).

<sup>72</sup>See Abowd, Kramarz, and Margolis (*Econometrica*) who discuss the estimation of a three-way error-component models. Andrews, Schank and Upward (2006, *Stata Journal* 6 (4)) support the implementation of some of these methods into Stata (see `st0112`).

<sup>73</sup>Also downloadable from <http://home.cerge-ei.cz/mittag/programs.html#fe>

average of the true unit-specific effects by using this one unit as the comparison unit for each (differently) treated units.<sup>74</sup> De Chaisemartin and D’Haultfoeulle (2018, 2020) propose an estimator, which corrects for diverging trends due to differential exposure (Stata packages `fuzzydid` and `did_multipldedgt`). Relatedly, in staggered adoption settings, where different units (for example, US states) adopt a policy (for example, unilateral divorce) at different times, policy effects may differ over time or across units based on when they adopt the policy, leading to biases. Goodman-Bacon (2021) (Stata package `bacondecomp`) offers diagnostics. In particular, apparent pre-trends can arise solely from treatment effects heterogeneity, and Sun and Abraham (2021, JEcm) avoid the issue by taking advantage of the presence of never-treated units in the sample (see Stata package `eventstudyinteract`).<sup>75</sup> Callaway and Sant’Anna (2021) propose a similar estimator that uses not-yet-treated units as control and can efficiently adjust for covariates using approaches developed in Sant’Anna and Zhao (2020). The Stata package `csdid` implements this estimator. Athey and Imbens (2022) consider the interpretation and variability of the difference-in-differences estimator in situations in which a unit’s date of adoption is randomly assigned. See also Borusyak et al. (CEPR DP No. DP17247). For up-to-date surveys of DiDs methods that deal with biases arising when the policy’s effect is heterogeneous between groups or over time, see de Chaisemartin and D’Haultfoeulle (2022, NBER WP no. 29691) and Roth et al. (2022).

**Remark 32.** Verdier (2020, REStat) shows that with few observations per fixed effect, inference in two-way fixed-effect models is affected by the noise in fixed effect estimation (one-way clustering will not do) and offers a solution.

**Remark 33.** The Dif-in-Difs FE approach rests on the identifying assumption of similar pre-treatment trends, i.e., it rules out the possibility that states experiencing or (correctly) expecting a boom increase minimum wages. A standard DiDs application plots pre-trends for treatment and control groups and perform formal tests of whether treatment has an effect on the outcome before it occurs. As an alternative to tests for pre-trends, Freyaldenhoven, Hansen, and Shapiro (2019 AER) consider a linear panel data model with possible endogeneity and show how to exploit a covariate related to the confound but unaffected by the

---

<sup>74</sup>See <https://pubs.aeaweb.org/doi/pdfplus/10.1257/jep.36.4.193>

<sup>75</sup>For an example of application of the Sun and Abraham (2021) method, see IZA DP No. 15843 (<https://docs.iza.org/dp15843.pdf>)

policy of interest to perform causal inference.<sup>76</sup> Consider a setup where a covariate  $\tilde{x}$  (adult employment) is affected by the unobservable (labor demand  $\eta$ ), which affects treatment allocation (minimum wage increases  $x$ ), but is not affected by treatment  $x$  itself. They propose to study the dynamics of  $\tilde{x}$  around the timing of treatment and a 2SLS model to identify the treatment effect by regressing  $y$  (youth employment) on  $\tilde{x}$  and  $x$  where  $x_{i,t+1}$  serve as instruments based on assuming that the dynamic relationship of  $\tilde{x}$  to  $x$  mimics that of  $\eta$  to  $x$ .

**Remark 34.** Note that there may not be variation in the treatment  $x$  inside each group corresponding to fixed effects. With heterogenous treatment effects, this has consequences (Miller et al., NBER Working Paper No. 26174).

**Remark 35.** Today, one typically implements the DiD design based on matching. See Section 20.2.

**Remark 36.** One typically motivates the DiDs estimator by arguing that the transformed (e.g., first-differenced) outcome is independent of unobservables. Alternatively, one searches for controls that make the unobservables independent of the treatment assignments. Combining the two approaches leads to more robust estimates of treatment effects according to Arkhangelsky and Imbens (2019). They generalize the fixed effect estimator that compares treated and control units at the same time within the set of units with the same fraction of treated periods.

**Remark 37.** The most important methodological update of the DiDs methodology is the synthetic control (SC) approach (Abadie, Diamond, and Hainmueller, 2010, 2014; Arkhangelsky, et al., 2019 NBER WP No. 25532, AER 2021). For each treated group one uses a control composed of an average of controls weighted to resemble the treated group (the implementation relies on the minimum distance method, see below). Arkhangelsky, et al. (2019) link the SC method to WLS DiDs regressions, discuss robustness and inference.<sup>77</sup> More specifically, DiDs relies on the parallel trends assumption, and one solution to this challenge has been the application of synthetic control (SC) methods, where one generates a single synthetic control from a unique convex weighting of underlying control units (Think of Miami counterfactual in the Mariel boat-lift paper.). The Synthetic DiDs (SDID) estimator of Arkhangelsky et al. (2021, AER) uses optimally chosen time weights

<sup>76</sup><https://www.brown.edu/Research/Shapiro/pdfs/prerends.pdf>

<sup>77</sup>Be wary of conditioning on all lagged outcome values: [http://www.gregor-pfeifer.net/files/SCM\\_Predictors.pdf](http://www.gregor-pfeifer.net/files/SCM_Predictors.pdf)



and separate unit weights to match treated and control units on pre-treatment trends (not necessarily on both pre-treatment trends and levels as the Abadie et al. (2010) SC did). The unit weights ensure the parallel pre-treatment trends, the time weights draw more weight from pre-treatment periods which are more similar to post-treatment periods.<sup>78</sup>

**Remark 38.** As usual, one can test the SC DiD approach using a placebo test based on data from before the intervention. Inference can be based on a permutation test (Abadie, Diamond, & Hainmueller, 2015; Ando and Ando, 2015) or bootstrap (Gobillon and Magnac, 2016; Firpo and Possebom, 2017). As the weights in the SCM approach are themselves estimated, this introduces measurement error in the key RHS variable. See, e.g., Ferrari and Galindo, for how to deal with this.

**Remark 39.** Kaul et al. (2021, JBES) argue against controlling for the entire pretreatment path of the outcome variable in applications of the synthetic control method.

## 7.2. Errors in Variables

([H] 3.9) One particular form of endogeneity of RHS variables was of concern in the previous section. We used the fixed effect model to capture time constant person-specific characteristics. The second most important potential source of bias is measurement error. Its effects are typically opposite to those of a typical unobserved fixed effect. Consider the model 7.1, where  $x$  is measured with error, i.e., we only observe  $\tilde{x}$  such that

$$\tilde{x}_i = x_i + \nu_i \quad (7.2)$$

In the case of classical measurement error, when  $E[\nu x] = E[\nu \epsilon] = 0$ , OLS is inconsistent and biased towards 0. For a univariate  $x_{it}$  we show in class that

$$\hat{\beta}_{OLS} \xrightarrow{p} \frac{\sigma_x^2}{\sigma_x^2 + \sigma_\nu^2} \beta \quad (7.3)$$

Note that what matters is the ratio of the ‘signal’  $\sigma_x^2$  to ‘noise’  $\sigma_\nu^2$  (the reliability ratio  $\sigma_x^2 / (\sigma_x^2 + \sigma_\nu^2)$ ). The standard error (inference) is also underestimated in proportion to the reliability ratio. Also note that adding additional regressors will

---

<sup>78</sup>The SDID estimator is now Stata (`sdid`) for balanced panels, thanks to Pailanir and Clarke (2002), see <https://arxiv.org/pdf/2301.11859.pdf>.

typically exacerbate the measurement error bias because the additional regressors absorb some of the signal in  $\tilde{x}$ .

**Exercise 7.2.** Show that the reverse regression of  $\tilde{x}$  on  $y$ , which suffers from endogeneity, helps to bracket the true coefficient.

**Exercise 7.3.** Suppose there are two variables in  $x_{it}$ , only one of which is measured with error. Show when the coefficient estimator for the other variable is affected as well.

**Remark 40.** In the case of misclassification of a binary variable  $E[\nu\varepsilon] = 0$  cannot hold. This still biases the coefficient towards 0 (Aigner, 1973). However, the bias can go either way in other cases of non-classical measurement error.

**Remark 41.** Measurement error in a binary endogenous variable is a headache. See DiTralia and Garcia-Jimeno (JEcm, 2019) for the identification of a misclassified binary endogenous regressor when a discrete-valued instrumental variable is available.

If  $\tilde{x}$  is the dependent variable, this will only increase standard errors under the classical measurement error assumptions. However, if the covariance of  $x$  and  $\nu$  is not zero, then we can express  $\nu = \delta x + \varsigma$ , where  $\delta$  is the attenuation bias, which will affect regressions with  $\tilde{x}$  on the left hand side. Bound et al. (1994) show that the  $\hat{\beta}$  in the regression  $\tilde{x} = x + \nu = x(1 + \delta) + \varsigma = z\beta + \epsilon$  will be biased since  $E[\hat{\beta}/\beta] = (1 + \delta)$ . They also derive the reliability ratio for this case. Meyer and Mittag (2016) address the non-classical situation for the case when the dependent variable  $\tilde{x}$  is binary and derive a closed form solution for the bias in the linear probability model.

Within estimators (differencing) will typically make the measurement error bias worse. The signal-to-noise ratio will depend on  $\sigma_x^2$  and on  $\sigma_x^2 + \sigma_\nu^2[(1 - \tau)/(1 - \rho)]$  where  $\tau$  is the first-order serial correlation in the measurement error and  $\rho$  is the first-order serial correlation in  $x$ . Again, the intuition is that differencing kills some of the signal in  $\tilde{x}$  because  $x$  is serially correlated, while the measurement error can occur in either period.

**Exercise 7.4.** Derive the above-stated result.

**Exercise 7.5.** Explain how we could use a second measurement of  $x_{it}$  to consistently estimate  $\beta$ . Would this strategy work also in the case of a binary  $\tilde{x}$ ?<sup>79</sup>

<sup>79</sup>Look up Kane, Rouse, and Staiger (1999) after you're done with the exercise.

**Example 7.2.** In estimating the labor supply equation off PSID data the measure of wages is created as earnings over hours. If there is a measurement error in hours, the measurement error in wages will be negatively correlated with the error term in the hours equation.

**Remark 42.** When you don't have an IV, use reliability measures (separate research gives you these). See Abowd and Stinson (2013, REStat) and Bingley and Martinello (2017, JoLE) for an approach where neither of your two measures is assumed to be the truth.

**Remark 43.** IV estimation method for errors in variables does not generalize to general nonlinear regression models. If the model is polynomial of finite order it does: see Hausman et al. (1991). See Schennach for use of Fourier transformation to derive a general repeated-measurement estimator for non-linear models with measurement error.

**Exercise 7.6.** Assume a simple non-linear regression model  $y_i = \beta f(x_i) + \varepsilon_i$  with one regressor  $x_i$  measured with error as in Equation 7.2. Use Taylor series expansion around  $\tilde{x}$  to illustrate why normal IV fails here.

Griliches and Hausman (1986): “Within” estimators are often unsatisfactory, which was blamed on measurement error. Their point: we may not need extraneous information. If  $T > 2$  differencing of different lengths and the deviations-from-mean estimator will eliminate fixed effects and have a different effect on potential bias caused by measurement error. Therefore, differencing may suggest if measurement error is present, can be used to test if errors are correlated, and derive a consistent estimator in some cases. Note that here again (as with the fixed effect model) panel data allows us to deal with estimation problems that would not be possible to solve in simple cross-section data in absence of valid instruments.

**Example 7.3.** Aydemir and Borjas (2006) estimate the impact of the share of immigrants on a local labor markets on wages. The regional share of immigrants  $\pi_r$  is estimated from random samples of the population, which introduces sampling error in the key independent variable. The variance of the sampling error is binomial ( $\pi_r(1 - \pi_r)/n_r$ ) for very small sampling rates ( $\tau$ ). It can be shown that the inconsistency takes the form

$$\widehat{\beta}_{OLS} \xrightarrow{p} \left(1 - \frac{(1 - \tau)E_r[\pi_r(1 - \pi_r)/n_r]}{(1 - R^2)\sigma_p^2}\right) \beta, \quad (7.4)$$

where  $\sigma_p^2$  is the variance of the observed  $\pi_r$  across the labor markets and where the  $R^2$  comes from an auxiliary regression of  $\pi_r$  on all right-hand-side variables in the model. The authors show that the use of regional fixed effects, typical in this literature, makes the auxiliary-regression  $R^2$  close to 1 and the bias can be very large.<sup>80</sup> Equation (7.4) can be used to predict the correct coefficient or run IV using two measures. Note that the same problem is present when regressing outcomes on group-specific shares of unemployed, women, foreign firms, etc. as long as random samples are used to estimate the group-specific shares.<sup>81</sup>

**Remark 44.** Measurement error is important for field experiments, see Leeat, Gillen, and Snowberg (*JPE*, in press).

## 8. Inference in Panel Data Analysis

Tests like *Breusch-Pagan* tell us whether to run OLS or random effects (GLS). What we really want to know is whether we should run fixed effects or random effects, i.e., is  $COV[\alpha_i, X_i] \neq 0$ ? Next, we wonder how to test for specific departures from the fixed effect model. But first, we discuss how to conduct inference when regressors vary across groups (clusters) of observations as in industry-specific variables in firm-level regressions or as in difference-in-differences (fixed-effect) regressions.<sup>82</sup>

### 8.1. Inference with Clustered Data and in “Difference in Differences”

One of the most popular (panel-data) identification strategies, pioneered by Ashenfelter and Card (1985), is to estimate the difference in before/after time differences in the mean outcome variable across (two) groups, where some (one) of the groups

<sup>80</sup>A back-of-the-envelope percent bias is  $(1 - \tau) (\bar{\pi}(1 - \bar{\pi})/\bar{n}) [(1 - R^2)\sigma_p^2]^{-1}$ .

<sup>81</sup>There is another problem related to ‘shares’ variables: Gerdes (IZA DP No. 5171) uses earlier results to highlight potential inconsistency of the fixed-effect specification linear in shares, which implicitly weights the differences (changes) in shares by the denominator variable of the shares (e.g., population size). To the extent that both the dependent variable and the denominator variable can be determined by common factors, this leads to inconsistency. The proposal is to use the  $\ln(\text{share})$  transformation, which eliminates this source of inconsistency (see also Ashenfelter and Greenstone, 2004).

<sup>82</sup>Cameron and Miller (JHR, 2015) is a key practical reference for clustering. For a useful overview of more recent advances in this area, see Michler and Josephson (2021, arXiv:2107.09736). See also note n. 67 and Remark 7.3 for examples of the fragility of the fixed effect model.

experience a treatment (policy change). For example, think of using repeated cross-sections of workers across U.S. states, where some of the states change the statutory minimum wage at some year, and estimate the following fixed effect model of the effect of this treatment  $T$  on  $y$ :

$$y_{ist} = \alpha_s + \delta_t + \gamma x_{ist} + \beta T_{st} + (\varepsilon_{ist} + \eta_{st}), \quad (8.1)$$

where  $i$  denotes workers,  $s$  denotes states (groups) and  $t$  denotes time (year).<sup>83</sup> So, the typical approach is to measure the differences in the before/after control/treatment averages.<sup>84</sup> With  $t = 0, 1$  and  $s = 0, 1$  and no  $x$ s and no  $\eta$ s,  $\widehat{\beta}_{DD} = \beta + (\bar{\varepsilon}_{11} - \bar{\varepsilon}_{10}) - (\bar{\varepsilon}_{01} - \bar{\varepsilon}_{00}) \longrightarrow \beta$  thanks to averaging of  $\varepsilon$ s inside  $st$  groups. When the random effects  $\eta_{st}$  are present,  $\widehat{\beta}_{DD} \longrightarrow \beta + (\eta_{11} - \eta_{10}) - (\eta_{01} - \eta_{00})$  and so is not consistent for  $\beta$  with  $t = 0, 1$  and  $s = 0, 1$ .<sup>85</sup>

When forming inference based on the Diff-in-Diffs approach, or when using group-level variation (such as a class-level variable not varying across individual students; Moulton, 1986, 1990), as well as when using stratified samples (sample villages, then households within villages), we need to acknowledge that  $T$  does not vary at the  $i$  level. (See Remark 23 and Example 10.7.) The downward bias to standard errors is the larger the more correlated  $T$  is within clusters (for example, correlated over time when using state-time policy changes) and the more correlated residuals are within clusters.<sup>86</sup> So far, the practical upshot of the recent work is that when there are many treatment and control groups (more than 50

<sup>83</sup>Of course, one can go further and focus on the effect on a sub-group within  $ist$ , defining a difference-in-difference-in-differences estimate.

<sup>84</sup>DinDs only identifies the ATT, which we introduce in Section 20.1 (it does not say anything about the effect of staying unionized or staying out of the union). Duleep (2012, IZA DP no. 6682) points out that the variance of the difference in the treatment and control group averages ( $\beta$ ) is larger than the variance of the average of the individual differences, such that instead of running the regression 8.1, one should just manually calculate the before/after difference for each individual within the control and the treatment groups, take the averages of these individual differences and then difference them out. This assumes no  $x$ , of course, and will become useful once we learn how to match in Section 20.2.

<sup>85</sup>In the famous Card and Krueger 1994 AER minimum wage study, there is only one treatment group (NJ) and only one control group (PA). No matter how many restaurants are covered in each state, the treatment effect may not be consistently estimated.

<sup>86</sup>And the bias is worse with fewer clusters, skewed  $X$  and with unbalanced clusters. The bias could reverse sign in the unusual case when within-cluster regressor and error correlations are of the opposite sign. See Carter, Schnepel & Steigerwald (2017, REStat) for the effective contribution of each cluster depending on its properties.

clusters), it is OK to just cluster  $\epsilon$  at group (state) level.<sup>87</sup> If the number of clusters is small, however, there is a serious inference problem.

More specifically, Bertrand et al. (2004) suggest that OLS applied to (8.1) gives wrong standard errors because it relies on long time series, the dependent variables are typically highly positively serially correlated, and the treatment dummy  $T_{st}$  itself changes little over time. In their paper, placebo laws generate significant effects 45% of the time, as opposed to 5% (they repeatedly generate ‘law’ changes at random and estimate the model and they count the number of estimates suggesting a significant effect).<sup>88</sup> If the number of groups is small, they propose randomization-based inference: using the distribution of estimated placebo law effects to form the test statistic (to recover the distribution of the treatment effect parameter conditional on state and time fixed effects). See the discussion on randomized inference below.

DinDs regressions are just one example of research designs where treatment is correlated within clusters. Clustering is also expected when one works with a stratified random sample.

**Remark 45.** *An example of when clustering makes sense is the use of the Bartik (1991) (shift-share) IVs, where, for example, a demand shock in a geographic location is constructed by combining the local composition of industries in the pre-period with the national growth rates of industries or national trade liberalization by industry, etc. (e.g., Autor, Dorn, Hansen, 2013; Acemoglu and Restrepo, 2017). Identification comes from exogenous industry-level shocks at the national level and/or exogenous local industrial pre-period structure. For instance, Autor et al. (2013) study the local labor market effects of Chinese import penetration by industry. To avoid reverse causality, they IV for US national-level import penetration using European penetration rates by industry. Borusyak and Jaravel (2017) show that consistency of the Bartik approach requires that there are sufficiently many industries with a substantial employment share and they also*

---

<sup>87</sup>Again, see Remark 23. When we worry about serial correlation of state-level shocks, we cluster at state level, not at state-time level as state-level clustering allows for any time-series process within each state. Even when one conditions on fixed effects at the cluster level, clustering may be important (Arellano, 1987).

<sup>88</sup>As a solution they propose to aggregate up the time series dimension into pre- and post-treatment observations or allow for arbitrary covariance over time and within each state. These solutions work fine if the number of groups is sufficiently large—but then we simply cluster. See Hansen (2007) for a proposed alternative solution that is based on (correctly) estimating the AR shock process at the group (state) level of variation.

highlight the importance of clustering to allow for the correlation of residuals in locations with similar industrial composition.<sup>89</sup> Adao, Kolesar, and Morales (QJE 2019) conduct a placebo exercise where tests based on commonly used standard errors with 5% nominal significance level reject the null of no effect in up to 55% of the placebo samples. Overrejection problem arises because regression residuals are correlated across regions with similar sectoral shares, independently of their geographic location. They derive useful inference methods.

Goldsmith-Pinkham, Sorkin, and Swift (2018) show that the Bartik instruments are numerically equivalent to using the initial shares (interacted with time fixed effects when multiple periods are used) as instruments in a weighted GMM estimation—the shifts only provide the weights and affect the instrument relevance, but not the endogeneity. So one needs to argue why the initial shares are exogenous.<sup>90</sup> Next, they show that in practice, a small number of industries often account for a large share of the identifying variation, which makes clear which shares in particular one needs to argue are exogenous. See <https://github.com/paulgp/bartik-weight>. Jaeger, Ruist, and Stuhler (2018) point out that if it takes time for markets to adjust to shocks, then the Bartik instrument will conflate the short-term response (e.g. a fall in wages when new immigrants enter) and the long-term response (e.g. a positive move back as capital has time to adjust). They propose to add lagged immigrant flows (and IV for them using the corresponding Bartik IV). However, for this to work, one needs independent variation in the two periods in where migrants are coming from. Borusyak, Hull, and Jaravel (2018, NBER WP No. 24997) derive orthogonality condition for these instruments to identify causal effects and suggest a new type of regressions to help visualize identifying shock variation and to correct standard errors. For inference issues with shift-share IVs, see also Remark 45.

**Remark 46.** When cluster sizes are small, FEM inference should be based on the “within” estimator, not on the LSDV estimator because they use different degrees-

---

<sup>89</sup>Alternatively, one may view the endogeneity problem as stemming from the possibility that import penetration at the national level depends on local supply shocks, in which case leave-one-out estimation would be the remedy.

<sup>90</sup>To be more precise, with a fixed T, a fixed number of industries, and the number of countries going to infinity, shares provide all of identifying variation. With a large number of growth shocks (number of industries going to infinity or large T), the growth rates can provide a consistent estimate of the second-stage parameter. As proved in Borusyak, Hull, and Jaravel (2018), exogenous independent shocks to many industries leads [sic] the Bartik estimator to be consistent, even when the shares are not exogenous.

of-freedom corrections.<sup>91</sup> See Cameron and Miller (2015, *JHR*) and Carter and Steigerwald (2013).

**Remark 47.** Carter et al. (2017, *REStat*) provide a proof of asymptotic normality of  $t$ -statistics when the clusters are not of equal size<sup>92</sup> and offer a degree-of-freedom (cluster-heterogeneity) correction that approximates the effective number of clusters once you account for their unequal sizes and for covariates' variation over clusters.<sup>93</sup> Related to the issue of unequal cluster sizes, Brewer, Crossley and Joyce (2017, *JEcmMethods*) provide Monte Carlo evidence on low power with few clusters and suggest GLS with robust inference, an unusual thing in the literature.

**Remark 48.** The larger clusters we define, the safer we are in terms of the cluster bias of inference, but the more imprecise the middle matrix in the “sandwich” variance estimator (given in Equation (6.7)) will be. This reflects the general bias-variance trade-off.<sup>94</sup> But there are other problems with over-clustering, see the last remark in this section.

**Remark 49.** Traditional inference only considers uncertainty entering through sampling error in the estimation of the group (before/after) means of  $y$ . Recent work, however, argues that one should allow for additional sources of uncertainty when forming inference, in particular, uncertainty in the quality of the control groups (Abadie, Diamond, and Hainmueller, 2007).

When the number of clusters is small, i.e., below 50 (where this number comes from placebo simulation studies), the “sandwich” estimator will lead to too many rejections of  $H_0$ , especially when clusters are not balanced. As suggested by papers discussed below, this is probably happening on a large scale. There are three approaches to fix this issue: bias-corrected variance, bootstrap with asymptotic refinement, and the use of the  $T$  distribution with adjusted degrees of freedom

---

<sup>91</sup>Hence, in Stata, use `xtreg y x, fe vce(robust)`, not `regress y x i.id_clu` or `regress y x, absorb(id_clu) vce(cluster id_clu)`. With multi-way clustering, use `cmgreg` or `ivreg2`.

<sup>92</sup>Consider a sample of 20 observations divided into two clusters and the number of non-zero elements of the  $\widehat{\varepsilon}_i \widehat{\varepsilon}_i'$  matrix as one cluster size increases to 19.

<sup>93</sup>They also show that one cannot conduct inference on cluster-level fixed effects using the cluster-robust  $t$ -statistic.

<sup>94</sup>Of course, multi-way clustering is available for all Stata commands that allow clustering: `vcemway`.



(Cameron and Miller, 2015). Recently, randomized inference is also becoming popular.

Bootstrapped inference is based on re-sampling new samples out of the data in a way that preserves key features of the original sampling scheme. A widely-used workable bootstrap- $t$  solution for standard errors when the number of clusters is as low as six is provided in Cameron et al. (2008).<sup>95</sup>

The use of the  $T$  distribution is the simplest approach.<sup>96</sup> The use of the  $T$  distribution makes a difference: with  $G = 10$ , the two-sided  $\alpha = 0.05$  critical value for  $T(9)$  is 2.26 as opposed to the 1.96 based on  $N$ . Donald and Lang (2007) further suggest subtracting from  $G$  the number of regressors that are invariant within clusters.

Donald and Lang (2007, REStat) focus on the case where the number of groups (both treatment and control) is small. They propose a *two-step approach* where one first estimates an individual-level regression with fixed effects corresponding to the group-level variation; in the second stage one runs these fixed effects on the group-level RHS variable.<sup>97</sup> There are two possible cases for inference here: First, if the number of individuals within each group is large, this two-step estimator is efficient and its  $t$ -statistics are distributed  $T$  if the underlying group random effects (errors) are normally distributed.<sup>98</sup> If one does not want to assume that the

---

<sup>95</sup>Only their ‘wild bootstrap’ method works with such low number of clusters. For details of implementation, see their Appendix B and Stata code from one of the authors at <http://www.econ.ucdavis.edu/faculty/dlmiller/statafiles/> and also Bansil Malde’s ado file at <http://tinyurl.com/c8vz3br>. For an introduction to Stata’s `boottest`, see [http://qed.econ.queensu.ca/working\\_papers/papers/qed\\_wp\\_1406.pdf](http://qed.econ.queensu.ca/working_papers/papers/qed_wp_1406.pdf) For an extension to IV methods, see Finlay and Magnusson (2010). See also Barrios et al. (NBER WP No. 15760) and <http://ftp.iza.org/dp7742.pdf>

<sup>96</sup>Note that the number of parameters (for small sample degree-of-freedom corrections is  $k + G$ , where  $G$  is the number of groups, it is not just  $k$ . Stata correctly uses the  $T(G - 1)$  critical values after `reg y x, vce(cluster)` but not after `xtreg y x, vce(robust)`.

<sup>97</sup>One can think of the second stage as a Minimum Distance problem (see Section 8.3.1) where one ought to weight with the inverse of the variance of the estimated fixed effects. One may also need to think about the different implied weighting in the individual- and group-level regressions, and about omitted variable bias (see, e.g., Baker and Fortin, 2001).

<sup>98</sup>The error term in the 2nd stage regression has two components, one of which is the prediction error from the 1st stage, which will be approximately normal if samples within clusters are large. If one then also assumes that cluster-level shocks (the other component of the 2nd stage residual) are also normally distributed, then the result follows because  $t$  statistics are known to have the  $T$  distribution with normal errors in small samples.

More generally, they think of the problem as follows: a small number of groups (NJ and PA, see note n. 85) are drawn from the population of potential groups. Given the groups, one gets

cluster-level shocks are normally distributed, we're back to the standard "cluster-robust" generalization of White's heteroskedasticity-robust standard errors, where one only assumes that the cluster-level shocks are uncorrelated across clusters, i.e., that clusters are defined correctly. If, however, the number of members of each group is small, one needs special circumstances to argue that the two-step estimator of Donald and Lang (2007) is efficient with  $t$ -statistics having the  $T$  distribution. They even recommend running the first stage separately for each group. Ibragimov and Muller (2007) also propose estimating the regression within each group and averaging the estimates, dividing by their std. deviation and using  $T$  critical values.<sup>99</sup>

**Remark 50.** See IZA DP 12369 for an example of where the Donald and Lang CIs make a difference.

**Remark 51.** A useful way to test for the restrictions involved in these different approaches is to frame the issue as a minimum distance problem (see Section 8.3.1 and Loeb and Bound, 1996).

**Remark 52.** If the group-time shocks are not random (so that the Donald and Lang procedure is not applicable), Conley and Taber (REStat 2011) propose a more general inference method for cases where the number of treatment groups is small, but the number of control groups is large.<sup>100</sup> They propose to use estimates obtained from the many controls to characterize the small-sample distribution of the treatment parameter. Think of comparing the (change in the) outcome in the treatment state to the distribution of the corresponding outcomes from all other states—is it an outlier?  $\beta$  is obviously not consistent with only few treatment groups, but one can test that  $\beta$  equals a particular value and to get the corresponding confidence intervals based on the empirical distribution of residuals from

---

the group averages based on random samples taken within groups. The NJ-PA dif-in-difs is based on 4 means, each of which is obtained from a large sample—so they are approximately normal.

<sup>99</sup>See also Wooldridge (2003, AER) and an update of the AER paper, Wooldridge (2006) at <https://www.msu.edu/~ec/faculty/wooldridge/current%20research/clus1aea.pdf>, and Hansen (2007), who offers some simulations. One of the upshots is, again, that one should not cluster with few groups and large group sizes. See also Athey and Imbens (2006) for a non-parametric difference-in-differences.

<sup>100</sup>In terms of the Mariel boatlift paper, compare Miami to many other cities. Note that up until now, we discussed the asymptotics in terms of the number of groups without distinguishing the 'treatment' and 'control' groups.

the (many) control groups.<sup>101</sup> Their method is a generalization of the placebo laws experiment of Bertrand et al. (2004). In terms of equation (8.1) the inconsistency of  $\widehat{\beta}$  depends on  $T_{st} - \bar{T}_s$ , which is observable, and on  $\eta_{st} - \bar{\eta}_s$ . The distribution of the latter term can be estimated from the ‘controls’ as their number increases (when considering the DinDs method, we have assumed at the start that  $\eta_{st} - \bar{\eta}_s$  is unrelated to the imposition of treatment). Finally, Ferman and Pinto (2019, RE-Stat) derive an inference method that works in differences-in-differences settings with few treated and many control groups in the presence of heteroskedasticity.

**Remark 53.** Chandar et al. (2020, NBER WP No. 26389) consider field experiments conducted with the village, etc. as the unit of randomization. In short panel data, unit-level estimation with unit fixed effects and cluster-level estimation weighted by the number of units per cluster tend to be robust.

**Remark 54.** See Fafchamps and Gubert (2007, JDE) for clustering in a dyadic regressions (`ngreg` in STATA), where each observation expresses a relationship between pairs within a group. When studying network/peer effects on multiple outcomes, the question of whether to apply a Bonferoni correction (see note n. 49) vs. its opposite (due to strategic complementarities) is discussed in Casey, Glennerster, and Miguel (2012). See IZA DP No. 12584 for a procedure to cluster in data where one has measures of spatial/network linkages. To the extent that clustering is asked to reflect spatial relationships through unobservables (outcomes), data-driven procedures for cluster formation are also suggested here: <https://arxiv.org/pdf/2107.14677.pdf>. Note that such methods can also be applied within the first stage to the extent that clustering is motivated by correlation of treatment assignment, not in outcomes. For DinDs clustering in presence of spatial correlation, see also <https://arxiv.org/abs/1909.01782>.

**Remark 55.** Imbens and Kolesar (2016, REStat) provide references to the literature that documents the poor behavior of robust and clustered standard errors in small samples, particularly if the distribution of the covariates is skewed. They advertize the simple-to-implement improvement to clustered inference suggested by Bell and McCaffrey (2002) and provide R codes as well.

**Remark 56.** Abadie, Athey, Imbens, and Wooldridge (2017, NBER WP no. 24003, now in the QJE) argue in a cross-sectional setting that clustering is jus-

---

<sup>101</sup><http://www.ssc.wisc.edu/~ctaber/DD/notes.pdf>

tified in a research (experimental) design where assignment (not residuals or regressors<sup>102</sup>) is correlated within clusters.<sup>103</sup> Clustered inference corresponds to the assumption that there are clusters in the population beyond those in the sample. Standard robust inference corresponds to a random sample from the population. Abadie et al. propose a new variance estimator for the case when all (population) clusters are in the sample and there is variation in treatment assignment within clusters. Without cluster-level fixed effects, one should cluster if (a) there is clustering in the sampling and heterogeneity in the treatment effect, or (b) there is clustering in assignment. In contrast, with cluster-level fixed effects, they argue that clustering is necessary only if there is heterogeneity in the treatment effect (and either clustering in assignment or sampling).

**Remark 57.** *How does one determine the optimal level of clustering? MacKinnon, Nielsen, and Webb (2023) provide a testing procedure to assess the appropriate level of clustering in linear regression models (ArXiv:2301.04522), which draws on the logic of the Hausman specification test.*

Athey and Imbens (2017) advocate the use of randomization-based inference (Fisher, 1935) as opposed to sampling-based inference when analyzing experimental data.<sup>104</sup> Similar to bootstrap, the randomization inference (RI) thought experiment is based on repeated sampling from the underlying population (data). But in RI uncertainty in estimates arises from the random assignment of treatments, rather than from hypothesized sampling from a large population. While the RI p-values (based on the Fisherian reasoning in fixed experimental samples) are constructed similarly to bootstrap-based p-values (based on Neyman-type population sampling approach), there is a key difference: bootstrap captures uncertainty over the specific sample drawn from the population, while RI reflects uncertainty over which units in the sample are assigned to the treatment. In bootstrap, one re-samples observations from the actual sample with replacement, to reflect sampling variation. In RI, the randomness is in which units in the sample are treated: Could it be that even in absence of a treatment effect, the specific occurred selection of units into treatment generates a treatment-control gap? So, one re-assigns

---

<sup>102</sup>Clustering affects inference even if both residuals and regressors are uncorrelated within clusters as long as there is within-cluster correlation of their product. Does this mean one should use clustering?

<sup>103</sup>Clustering is also justified in the presence of a stratified sampling design, but most applications of clustering are not based on this.

<sup>104</sup><https://www.mattblackwell.org/files/teaching/s05-fisher.pdf>

“treatment” at random (keeps the share treated constant, unlike across bootstrap samples), to get a distribution of placebo “treatment effects” under the null hypothesis that the treatment has no effect, which in turn gives the right probability of rejecting the null. Note that even in small samples there is a large number of perturbations of treatment assignment.

**Remark 58.** *RI can make a big difference in experiments: Young (2016) relies on exact RI in to highlight the over-rejection of clustered inference: Leaving out one cluster (jackknife) kills  $p$  values of 50 experimental AER/AEJ papers. He shows that the asymptotics underlying clustered inference depends on maximal leverage in the regression going to zero, while it actually is close to 1 thanks to a long list of controls added to improve precision or as robustness checks.<sup>105</sup> But this conclusion may be driven by a bad Stata default setting for robust standard errors.<sup>106</sup>*

**Remark 59.** *It is natural to use robust inference in controlled experiments (because treatment affects the variance of  $y$ , not just its mean). Should we cluster by treatment level?<sup>107</sup>*

In the DiDs setting of this section, RI is an alternative to clustering.<sup>108</sup> Barrios et al. (2012) argue that researchers should care about the data structure beyond state-level clustering. When treatment status does not depend on just one variable, but on several factors, RI corresponds to permuting all treatment-determining variables randomly and assigning the treatment status thereafter, according to such pseudo characteristics (Pfeifer, Reutter, and Strohmaier, JHR).

---

<sup>105</sup>Leverage corresponds to the influence of data points on their predicted values (i.e., elements of the  $H = X(X'X)^{-1}X'$  matrix). Maximum leverage may stay high as  $n$  increases for cluster dummies for example. When conventional and randomization-based tests are exact and have accurate size, they have identical power as well. See, Jakiela and Ozier, Gendered Language paper for another application.

<sup>106</sup><http://datacolada.org/99> suggests using R or in Stata running `reg y x, vce(hc3)`.

<sup>107</sup>Robinson (2021) distills the discussion in Abadie et al. (2017, 2020) into a practical guide for clustering in experimental data (Table 3 in <https://ts-robinson.com/publication/robinson-whenshouldwe-2020/robinson-whenshouldwe-2020.pdf>)

<sup>108</sup>RI can also be applied in RDD designs Cattaneo et al. (2015) or the potential outcome framework (Ho and Imai, 2006)

## 8.2. Hausman test

- Basic idea is to compare two estimators: one consistent under both null hypothesis (no misspecification) and under the alternative (with misspecification), the other consistent only under the null. If the two estimates are significantly different, we reject the null.

$$\begin{array}{l}
 H_0 : COV[\alpha_i, X_i] = 0 \\
 H_A : COV[\alpha_i, X_i] \neq 0
 \end{array}
 \begin{array}{c}
 \hat{\beta}_{LSDV} \text{ fixed effects} \\
 \text{consistent, inefficient} \\
 \text{consistent}
 \end{array}
 \left|
 \begin{array}{c}
 \hat{\beta}_{GLS} \text{ random effects} \\
 \text{consistent, efficient} \\
 \text{inconsistent}
 \end{array}
 \right.$$

- The mechanics of the test:

**Theorem 8.1.** Under  $H_0$  assume  $\sqrt{n}(\hat{\beta}_j - \beta) \xrightarrow{D} N(0, V(\hat{\beta}_j))$ ,  $j \in \{LSDV, GLS\}$  and  $V(\hat{\beta}_{LSDV}) \geq V(\hat{\beta}_{GLS})$  and define  $\sqrt{n} \hat{q} = \sqrt{n}(\hat{\beta}_{LSDV} - \hat{\beta}_{GLS}) \xrightarrow{D} N(0, V(\hat{q}))$  where

$$V_q \equiv V(\hat{q}) = V(\hat{\beta}_{LSDV}) + V(\hat{\beta}_{GLS}) - COV(\hat{\beta}_{LSDV}, \hat{\beta}'_{GLS}) - COV(\hat{\beta}_{GLS}, \hat{\beta}'_{LSDV}).$$

then

$$COV(\hat{\beta}_{LSDV}, \hat{\beta}'_{GLS}) = COV(\hat{\beta}_{GLS}, \hat{\beta}'_{LSDV}) = V(\hat{\beta}_{GLS})$$

so that we can easily evaluate the test statistic  $\hat{q}' V_q^{-1} \hat{q} \rightarrow \chi^2(k)$ .

We prove the theorem in class using the fact that under  $H_0$  the  $\hat{\beta}_{GLS}$  achieves the Rao-Cramer lower bound.

**Remark 60.** The Hausman test asks if the impact of  $X$  on  $y$  within a person is the same as the impact identified from both within and cross-sectional variation.

**Remark 61.** Mundlak's formulation connects random and fixed effects by parametrizing  $\alpha_i$  (see [H] 3).

**Remark 62.** Similar to the Hansen test (see Section 10), Hausman is an all-encompassing misspecification test, which does not point only to  $COV[\alpha_i, X_i] \neq 0$ , but may indicate misspecification. Of course, tests against specific alternatives will have more power.

**Remark 63.** The power of the Hausman test might be low if there is little variation for each cross-sectional unit. The fixed effect  $\hat{\beta}$  is then imprecise and the test will not reject even when the  $\beta$ s are different.

**Remark 64.** *There is also a typical sequential testing issue. What if I suspect both individual and time fixed effects: which should I first run Hausman on. Since  $T$  is usually fixed, it seems safe to run Hausman on the individual effects, with time dummies included. But then we may run out of degrees of freedom.*

**Remark 65.** *The  $V_q$  given above assumes, unnecessarily, that REM is fully efficient under  $H_0$ , i.e. that there is no heteroskedasticity. When clustering, use the Stata command `xtoverid` to implement the Wooldridge (2010) cluster-robust version of the Hausman test.*

### 8.3. Using Minimum Distance Methods in Panel Data

Hausman test might reject  $COV[\alpha_i, X_i] = 0$  and one may then use the fixed effect model. But the fixed effect model is fairly restrictive and eats up a lot of variation for  $\alpha_i$ s. When  $T$  is small we can test the validity of those restrictions using the MD methods. The same technique allows for estimation of  $\beta$  with a minimal structure imposed on  $\alpha$ , allowing for correlation between the unobservable  $\alpha$  and the regressors  $x$ . We will first understand the MD method and then apply it to panel data problems.

#### 8.3.1. The Minimum Distance Method

Suppose we have a model that implies restrictions on parameters that are hard to implement in the MLE framework. When estimation of an unconstrained version of our model is easy (OLS) and consistent, the MD method offers a way to impose the restrictions and regain efficiency and also to test the validity of the restrictions ([H] 3A).

Denote the unconstrained estimator as  $\hat{\pi}_N$ , where  $N$  is the sample size in the unconstrained estimation problem, and denote the constrained parameter of interest as  $\theta$ . Next, maintain the assumption that at the true value of  $\theta$  the restrictions  $\pi = f(\theta)$  are valid. The objective is to find  $\hat{\theta}$  such that the distance between  $\hat{\pi}$  and  $f(\hat{\theta})$  is minimized:<sup>109</sup>

$$\hat{\theta}_N = \arg \min \{S_N\} \text{ where } S_N = N[\hat{\pi}_N - f(\theta)]' A_N [\hat{\pi}_N - f(\theta)], \quad (8.2)$$

and where  $A_N \xrightarrow{p} A$  is a weighting matrix and  $\sqrt{N}[\hat{\pi}_N - f(\theta)] \xrightarrow{D} N(0, \Delta)$ .<sup>110</sup>

<sup>109</sup>Find the minimum distance between the unconstrained estimator and the hyperplane of constraints. If restrictions are valid, asymptotically the projection will prove to be unnecessary.

<sup>110</sup>See Breusch-Godfrey 1981 test in Godfrey, L. (1988).

The minimization problem 8.2 is of considerably smaller dimension than any constrained estimation with the  $N$  data points!

**Theorem 8.2.** *Under the above assumptions and if  $f$  is 2nd order differentiable and  $\frac{\partial f}{\partial \theta'}$  has full column rank then a)  $\sqrt{N}[\hat{\theta}_N - \theta] \xrightarrow{D} N(0, V(A))$ , b) the optimal  $A = \Delta^{-1}$ , and c)  $\widehat{S}_N \xrightarrow{D} \chi^2(r)$  where  $r = \dim(\pi) - \dim(\theta)$  is the number of overidentifying restrictions.*

We provide the proof in class. To show a) simply take a FOC and use Taylor series expansion to relate the distribution of  $\hat{\theta}_N$  to that of  $\hat{\pi}_N$ .

**Remark 66.** *Note that the Minimum Distance Method is applicable in Simultaneous Equation Models to test for exclusion restrictions.*

$$\Gamma y_t + Bx_t = u_t \Rightarrow y_t = \Pi x_t + v_t \text{ where } \Pi = -\Gamma^{-1}B$$

and we can test zero restrictions in  $\Gamma$  and  $B$ .

**Remark 67.** *MD is efficient only among the class of estimators which do not impose a priori restrictions on the error structure.*

### 8.3.2. Arbitrary Error Structure

When we estimate random effects,  $COV[\alpha, x]$  must be 0; further, the variance-covariance structure in the random effect model is quite restrictive. At the other extreme, when we estimate fixed effects, we lose a lot of variation and face multicollinearity between  $\alpha_i$  and time constant  $x$  variables.

However, when  $T$  is fixed and  $N \rightarrow \infty$ ,<sup>111</sup> one can allow  $\alpha$  to have a general expectations structure given  $x$  and estimate this structure together with our main parameter of interest:  $\beta$  (Chamberlain 1982, [H] 3.8). That is we will not eliminate  $\alpha_i$  (and its correlation with  $x$ ) by first differencing. Instead, we will control for (absorb) the correlation between  $\alpha$  and  $x$  by explicitly parametrizing and estimating it. This parametrization can be rich: In particular, serial correlation and heteroskedasticity can be allowed for without imposing a particular structure on the variance-covariance matrix. In sum, we will estimate  $\beta$  with as little structure on the omitted latent random variable  $\alpha$  as possible.<sup>112</sup> The technique of estimation will be the MD method.

<sup>111</sup>So that  $(N - T^2K)$  is large.

<sup>112</sup>The omitted variable has to be either time-invariant or individual-invariant.



Assume the usual fixed effect model with only  $E[\varepsilon_{it} | x_{it}, \alpha_i^*] = 0$

$$y_i = e_T \alpha_i^* + \underset{T \times K}{X_i} \beta + \varepsilon_i \quad (8.3)$$

and let  $x_i = \text{vec}(X_i')$ .<sup>113</sup> To allow for possible correlation between  $\alpha_i$  and  $X_i$ , assume  $E[\alpha_i^* | X_i] = \mu + \lambda' x_i = \sum_{t=1}^T \lambda_t' x_{it}$  (note  $\mu$  and  $\lambda$  do not vary over  $i$ ) and plug back into 8.3 to obtain

$$y_i = e_T \mu + (I_T \otimes \beta' + e_T \lambda') x_i + [y_i - E[y_i | x_i]] = e_T \mu + \underset{T \times KT}{\Pi} x_i + v_i \quad (8.4)$$

We can obtain  $\widehat{\Pi}$  by gigantic OLS and impose the restrictions on  $\Pi$  using MD.<sup>114</sup> We do not need to assume  $E[\alpha_i | X_i]$  is linear, but can treat  $\mu + \lambda X_i$  as a projection, so that the error term  $v_i$  is heteroscedastic.

**Exercise 8.1.** Note how having two data dimensions is the key. In particular, try to implement this approach in cross-section data.

**Remark 68.** Hsiao's formulae (3.8.9.) and (3.8.10.) do not follow the treatment in (3.8.8.), but use time varying intercepts.

### 8.3.3. Testing the Fixed Effects Model

Jakubson (1991): In estimating the effect of unions on wages we face the potential bias from unionized firms selecting workers with higher productivity. Jakubson uses the fixed effect model and tests its validity. We can use the MD framework to test for the restrictions implied by the typical fixed effect model. The MD test is an omnibus, all-encompassing test and Jakubson (1991) offers narrower tests of the fixed effect model as well:

- The MD test: Assume

$$y_{it} = \beta_t x_{it} + \varepsilon_{it} \text{ with } \varepsilon_{it} = \gamma_t \alpha_i + u_{it}$$

where  $\alpha_i$  is potentially correlated with  $x_i \in \{0, 1\}$ <sup>115</sup>. Hence specify  $\alpha_i = \underset{T \times k}{\lambda'} x_i + \xi_i$ . Now, if we estimate

$$y_i = \underset{T \times T}{\Pi} x_i + \nu_i$$

<sup>113</sup>Here,  $\text{vec}$  is the vector operator stacking columns of matrices on top of each other into one long vector. We provide the definition and some basic algebra of the  $\text{vec}$  operator in class.

<sup>114</sup>How many underlying parameters are there in  $\Pi$ ? Only  $K + KT$ .

<sup>115</sup>If  $\alpha_i$  is correlated with  $x_{it}$  then it is also correlated with  $x_{is} \forall s$ .

the above model implies the non-linear restrictions  $\Pi = \text{diag}(\beta_1, \dots, \beta_T) + \gamma\lambda'$  which we can test using MD. If  $H_0$  is not rejected, we can further test for the fixed effect model, where  $\beta_t = \beta \forall t$  and  $\gamma_t = 1 \forall t$ .

- Test against particular departures:<sup>116</sup>

- Is differencing valid? Substitute for  $\alpha_i$  to get

$$y_{it} = \beta_t x_{it} + \left(\frac{\gamma_t}{\gamma_{t-1}}\right) y_{it-1} - \left(\beta_{t-1} \frac{\gamma_t}{\gamma_{t-1}}\right) x_{it-1} + [u_{it} - \left(\frac{\gamma_t}{\gamma_{t-1}}\right) u_{it-1}]$$

Estimate overparametrized model by 3SLS with  $x$  as an IV for lagged  $y$ , test exclusion restrictions (see Remark 68), test  $\left(\frac{\gamma_t}{\gamma_{t-1}}\right) = 1$  (does it make sense to use  $\Delta y_{it}$  on the left-hand side?), if valid test  $\beta_t = \beta \forall t$ .

- Is the effect “symmetric”?<sup>117</sup>

$$\Delta y_{it} = \delta_{1t} ENTER_{it} + \delta_{2t} LEAVE + \delta_{3t} STAY + \Delta \mu_{it}$$

- Does the effect vary with other  $X$ s?

**Remark 69.** *In the fixed effect model we rely on changing  $x_{it}$  over time. Note the implicit assumption that union status changes are random.*

- In “fuzzy” DiDs, i.e., when treatment rate only increases more in the treatment group, a popular estimator of the treatment effect is the DID of the outcome divided by the DID of the treatment (see the Wald IV above). Chaisemartin and D’Haultfoeuille (2017, REStud) study identification in presence of parameter heterogeneity and show that this ratio identifies a local average treatment effect only if the effect of the treatment is stable over time. They propose several alternative estimators that allow for treatment effect heterogeneity and change over time when there is a control sub-group whose exposure to the treatment does not change over time.<sup>118</sup>
- See Hoderlein and WeWl for a Binary Choice Difference-in-Differences Model with Heterogeneous Treatment Effects. They show identification of the ATT and of the joint distribution of the actual and counterfactual latent outcome variable in the treatment group.

<sup>116</sup>These tests are more powerful than the omnibus MD test. Further, when MD test rejects  $H_0$  then the test against particular departure can be used to point to the *source* of misspecification.

<sup>117</sup>See Gilraine (2020, JOLE) for how entry into vs. exit from treatment can identify distinct treatment effects within the RDD setup.

<sup>118</sup>See <https://sites.google.com/site/clementdechaisemartin/> for the **Stata** package.

## 9. Simultaneous Equations and IV Estimation

Simultaneous Equations are unique to social science. They occur when more than one equation links the same observed variables. Reverse causality. Identification.

Solution: IV/GMM: find variation in the  $X$  that suffers from simultaneity bias which is not related to the variation in the  $\varepsilon$ s, i.e., use  $\widehat{X}$  instead—the projection of  $X$  on  $Z$ —the part of variation in  $X$  that is generated by an instrument.<sup>119</sup> Theory or intuition is often used to find an “exclusion restriction” postulating that a certain variable (an instrument) does not belong to the equation in question. We can also try to use restrictions on the variance-covariance matrix of the structural system errors to identify parameters which are not identified by exclusion restrictions.

**Example 9.1.** Consider the demand and supply system from *Econometrics I* (and of Haavelmo):

$$\begin{aligned}q_D &= \alpha_0 + \alpha_1 p + \alpha_2 y + \varepsilon_D \\q_S &= \beta_0 + \beta_1 p + \varepsilon_S \\q_D &= q_S\end{aligned}$$

where  $S$  stands for supply,  $D$  stands for demand and  $p$  is price and  $y$  is income. We solve for the reduced form

$$\begin{aligned}p &= \pi_1 y + v_p \\q &= \pi_2 y + v_q\end{aligned}$$

and note that one can identify  $\beta_1$  by instrumenting for  $p$  using  $y$  which is excluded from the supply equation. Here we note that in exactly identified models like this the IV estimate  $\widehat{\beta}_1 = \frac{\widehat{\pi}_2}{\widehat{\pi}_1}$  (show this as an exercise); this is called indirect least squares and demasks IV. To identify  $\alpha_1$  estimate  $\Omega$ , the variance-covariance matrix of the reduced form, relate the structural and reduced form covariance matrices and assume  $COV(\varepsilon_D, \varepsilon_S) = 0$  to express  $\alpha_1$  as a function of  $\beta_1$ .

<sup>119</sup>Another intuitive approach is to model the unobservables directly (as the residuals in the first stage regression) and include them as an explanatory variable into the main equation. This is the so called *control function* approach (CF). In cases where the endogenous variable enters linearly, the CF approach and 2SLS are equivalent. However, CF is advantageous for models that are non-linear in endogenous variables (even if linear in parameters) and especially for models that are non-linear in parameters. See also Section 19.2.2.

**Remark 70.** Here, we think of 2SLS as a simultaneous equation system, where both equations come from one economic model. A very different approach is when we have one structural equation with an endogeneity problem and we find a “natural” or controlled experiment” (see Sections 1 and 20) to use only the exogenous portion of the variation in the variable of interest (Deaton, 2009, NBER WP no. 14690).

Along the same divide, one can think of identification as either corresponding to our ability to go from reduced-form to structural parameters (within an economic model) or to the (potential) presence of an experiment that could answer the question posed. For some questions, there are no counterfactuals, no available interpretation of the results of an experiment.<sup>120</sup>

**Remark 71.** Of course, next to simultaneity (and reverse causality), the other two sources of bias that are dealt with using IV strategies are measurement error and omitted variable bias.

A good instrument  $Z$  must be correlated with the endogenous part of  $X$  (in the first-stage regression controlling for all exogenous explanatory variables!)<sup>121</sup> and it must be valid, i.e., not correlated with  $\varepsilon$ . The next two sections discuss testing of each desired IV property.<sup>122</sup>

Note that the lack of IV- $\varepsilon$  correlation corresponds to two underlying statements: (a) the IV is as good as randomly assigned (conditional independence assumption, traditionally referred to as *exogeneity*), and (b) the IV operates through a single known causal channel (the IV is correctly excluded from the main equation with  $y$  on the left-hand side, *exclusion restriction*). (Deaton, 2009, prefers to call IVs that satisfy the conditions (a) and (b) external and exogenous, respectively.)

---

<sup>120</sup>For example, answering whether starting school at seven is better than doing so at six is impossible when using elementary school tests as the outcome variable because it is not possible to disentangle the effect of maturation from the effect of schooling. One can answer the question with data on adults, where age no longer has an effect on ability (Angrist and Pischke, 2009, p.6).

<sup>121</sup>See [W] exercise 5.11. The key is that the prediction error from the first stage, which appears in the residual of the second stage, is not made orthogonal to the other exogenous variables by means of the first stage projection.

<sup>122</sup>See Murray (2006, J of Economic Perspectives) for an easy-to-read overview of the material covered in this section; he provides several interesting and useful examples from applied work that are not used here.

**Remark 72.** Often, IVs are derived from “natural experiments”.<sup>123</sup> These “natural” IVs are clearly not affected by the  $y$  variables and so are exogenous (external in Deaton’s terms, there is no simultaneity), which is, however, only a necessary condition for being correctly excluded (exogenous in Deaton’s terms)! The question (assumption) is whether these IVs have any direct effect on  $y$  (other than through their effect on the endogenous  $x$ ).

**Remark 73.** In 2SLS, consistency of the second-stage estimates does not depend on getting the first-stage functional form right (Kelejian, JASA 1971; Heckman, Econometrica 1978). In other words, a simple linear probability model for an endogenous dummy is sufficient in the first stage to obtain consistent estimates in the second-stage regression in 2SLS.<sup>124</sup>

**Exercise 9.1.** If the endogenous variable enters in a quadratic form, one needs two instruments (such as the original instrument and its square, if it’s not a dummy) and two first-stage regressions, one for  $x$  and the other one for  $x^2$ . Show why using only one first stage and plugging in the square of the predicted endogenous variable could be a problem. As usual, consider that the difference between the original and the predicted values is part of the residual of the main equation.<sup>125</sup>

**Remark 74.** IVs are often group specific. These work by averaging (eliminating) within-group heterogeneity (using the law of large numbers, similar to using group averaging to reduce measurement error bias) and using only between group variation. Obviously, the cost is the loss of information within groups. This strategy stands in contrast to within-data estimation (fixed effect model) used in panel data settings, where we discard between group heterogeneity. In particular, with only

<sup>123</sup>Draft lottery predicts participation in the military, which drives earnings, Angrist, 1990; distance from the equator predicts per capita GDP, which drives religiosity, McCleary and Barro, 2006; rivers instrument for the number of school districts in explaining educational outcomes, Hoxby, 2000; month of birth as an instrument for years of schooling in an earnings regression, Angrist and Krueger, 1991; or rainfall as an instrument for economic growth in explaining civil war, Miguel, Satyanath, and Sergenti, 2004.

<sup>124</sup>See Newey and Powell (2003, Ecma) and Chen and Pouzo (2012, Ecma) for follow-ups.

<sup>125</sup>Similarly, consider dealing with measurement error by using one noisy measure (proxy) of a true (latent)  $x$  as instrument for another proxy, see Exercise 7.5. In a model with interactions in two latent variables on the RHS, such as  $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \epsilon$ , we need three IVs—for example  $z_1, z_2$ , and  $z_1 z_2$ . If only  $x_2$  is measured with error, then we need only two IVs— $z_2$  and  $z_2 x_1$ . An unrelated note on interacted variables in OLS: if you want the (interpretation of the) coefficients on the main (uninteracted) variables to not be affected by the inclusion of the interaction term, you need to enter the interaction as  $(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)$ .

one IV,  $\widehat{\beta}_{IV} = COV_N(\widehat{X}, Y)/VAR(\widehat{X}) = COV_N(Z, Y)/COV(Z, X)$  and when  $Z$  is binary, one can further write this as

$$\widehat{\beta}_{IV} = \frac{E_N[Y|Z = 1] - E_N[Y|Z = 0]}{E_N[X|Z = 1] - E_N[X|Z = 0]},$$

where  $E_N$  is a shorthand for sample mean.<sup>126</sup> This is also known as the Wald IV estimator.

When there are many (exclusive) groups, the IV estimator can be expressed as consisting of all pairwise (Wald IV) group comparisons, i.e. consisting of all IV comparisons between each pair of groups (binary Wald IV estimators) and these pairwise Wald-IV estimators are efficiently (GLS) weighted in the grand 2SLS (see Angrist, 1991, Remark 21, Section 6.2, and Section 20.3).

**Remark 75.** Angrist and Imbens (1995) show that when a continuous treatment (with variation across both extensive and intensive margins) is coarsened into a dummy, and when the IV affects both margins (with the same sign), the IV estimate will be biased upwards as it will correspond to the sum of the extensive-margin effect and of the intensive margin effect. 2SLS is shown to be a weighted average of per-unit average causal effects along the length of an appropriately defined causal response function.

**Remark 76.** What if there is some correlation between the instrument and the error term  $\epsilon$ , i.e., what if the IV is imperfect? Nevo and Rosen (2012, REStat) assume that (a) the IV- $\epsilon$  correlation is of the same sign as the correlation between the endogenous  $x$  and  $\epsilon$ , and (b) that the IV is less correlated with  $\epsilon$  than the  $x$  is. Based on these assumptions, they derive analytic bounds for the parameters. There are other similar results (originating with Manski's work) on partial (point vs. interval) identification in 2SLS with imperfect instruments (Conley, Hansen, and Rossi, 2012; Ho and Rosen, 2016), i.e. imperfectly exogenous IVs, IVs that are at least less endogenous than the endogenous  $X$ .

---

<sup>126</sup>This is based on  $COV(Z, Y) = E[Z Y] - E[Z]E[Y] = (E[Z Y]/P(Z = 1) - E[Y])P(Z = 1) = \{E[Y|Z = 1] - (E[Y|Z = 1]P[Z = 1] + E[Y|Z = 0]P[Z = 0])\}P(Z = 1) = (E[Y|Z = 1] - E[Y|Z = 0])P[Z = 1](1 - P[Z = 1]) = (E[Y|Z = 1] - E[Y|Z = 0])VAR[Z]$

### 9.1. Testing for exclusion restrictions and IV validity

When we have more instruments than endogenous variables, we can test for their validity using the Sargan's test.<sup>127</sup> This test asks if any of the IVs are invalid based on maintaining the validity of an exactly-identifying sub-set of the IVs. But which IVs belong to the valid subset? The test is problematic when all IVs share a common rationale (see Example 9.2 below). The test can never test whether *all* of the IVs are valid. In particular, not rejecting H0 in a test of over-identifying IVs (restrictions) is consistent with all instruments (restrictions) being invalid while rejecting H0 is consistent with a subset of IVs being correct!

**Remark 77.** *For a test of the validity of over-identifying exclusion restrictions, which we have already covered, see remark 66.*

**Example 9.2.** *For a simple test of an exclusion restriction, see Card (1993) who estimates returns to schooling using proximity to college as an instrument for education and tests for exclusion of college proximity from the wage equation. To do this he assumes that college proximity times poverty status is a valid instrument and enters college proximity into the main wage equation. Notice that you have to maintain just identification to test overidentification. Also note that, unfortunately, the instrument, which is maintained as correctly excluded, in order to allow for this test, is based on the same economic argument as the first instrument, which is being tested.*

Aside from econometric tests for IV validity (overidentification), one can also conduct intuitive tests when the exogenous variation (IV) comes from some quasi-experiment. Chiefly, this consists of asking whether there is a direct association between the instrument and the outcome in samples where there was never any treatment.

**Example 9.3.** *For example, Angrist in the 1990 Vietnam-era draft lottery paper asks if earnings vary with draft-eligibility status for the 1953 cohort, which had a lottery, but was never drafted.*

**Remark 78.** *Das and Polachek (2019, IZA DP no. 12766) propose identification in situations where exogenous factors make the treatment ineffective for a subset*

---

<sup>127</sup>See `stata ivreg2`. See the next section on GMM for a general way of testing for overidentifying restrictions—the Hansen test.

of the treated population and so for these observations any observed correlation between the outcome and treatment must be due to the confounding endogeneity bias, which then can be used to identify the causal effects.

**Example 9.4.** Altonji, Elder and Taber (2005, JHR) provide several evaluations of IVs. They study the effect of studying in a Catholic (private) high school,  $CH$ , (as opposed to a public one) on several education outcomes  $y$ . In this literature, people use the following IVs ( $Z_i$ ): being Catholic ( $C$ ), distance to a Catholic school (proximity,  $D$ ), and their interaction ( $C * D$ ). So run  $y_i = \alpha CH_i + X_i' \gamma + \epsilon_i$  and IV for  $CH$ . Let's focus on evaluating the first IV:  $C$ :

1. Ask what the direct effect of  $C$  is on  $y$  for those who attend public eighth grades, because this group of students never attends Catholic high schools regardless of their religion or proximity.<sup>128</sup> They find some direct effect, suggesting the IV is not valid.
2. Use an approach that guesses about the degree of selection on unobservables from the measured degree of selection on observables. We usually believe that unobservables and observables are correlated. Indeed, in their data, there are large gaps in observables in favor of Catholic students. They find that all of the IVs lead to substantial upward biases in 2SLS. The 2SLS estimates are implausibly large. See the generalization of this approach by Emily Oster (2016) "Unobservable Selection and Coefficient Stability".<sup>129</sup>
3. Finally, they not only ask about whether exclusion restrictions are correct but also compare power of the IV to that of non-linearity in identifying the coefficient of interest in non-linear models. Specifically, they compare 2SLS to bivariate Probits and find that non-linearities, rather than the IVs (exclusion restrictions), are the main sources of identification, which is why

---

<sup>128</sup>Let  $\text{proj}$  refer to the least squares projection:  $\text{proj}(C_i|X_i) = X_i' \pi$  and  $\text{proj}(CH_i|X_i, C_i) = X_i' \beta + \lambda C_i$ . Define  $\tilde{C}_i \equiv C_i - X_i' \pi$ . Then  $\hat{\alpha}_{IV} \xrightarrow{p} \alpha + \frac{\text{COV}(\tilde{C}_i, \epsilon_i)}{\lambda \text{VAR}(\tilde{C}_i)}$ . Now, run a regression of  $y$  on  $C$  and  $X$  for those who have  $\Pr(CH_i = 1) = 0$ . The coefficient on  $C_i$  will converge to  $\frac{\text{COV}(\tilde{C}_i, \epsilon_i)}{\text{VAR}(\tilde{C}_i)}$ .

So, obtain an estimate of the bias  $\psi$  by running  $y_i = X_i' \delta + [\tilde{C}_i \hat{\lambda}] \psi + \omega_i$ .

<sup>129</sup>In linear models, the omitted variable bias formula allows one to get a sense of the sensitivity of estimated parameters. Gelbach (2016, JoLE) uses the formula to consider adding intercorrelated  $X$ s in a sequence. Andrews, Gentzkow, and Shapiro (2017, QJE) offer a (local) analogue of the formula for non-linear models based on a local linearization of moment conditions.



*the bivariate Probits give more plausible estimates.*<sup>130</sup> *The lesson is to avoid non-linearities, since otherwise you may get an estimate seemingly based on an IV while all you're getting is actually based on the arbitrary distributional assumption.*

**Remark 79.** *Oster (2016) argues that omitted variable bias is proportional to coefficient movements only if such movements are scaled by the change in the R-squared when controls are included.*

**Remark 80.** *Conley, Hansen, and Rossi (2012, REStat) illustrate applications where instruments can yield informative results even under appreciable deviations from an exact exclusion restriction.*

## 9.2. Regression Discontinuity

A separate class of IV models, referred to as **regression discontinuity** design (RDD), arises when the endogenous variable, such as assignment to some treatment, is (fully or partly) determined by the value of a (forcing) covariate lying on either side of an (administrative) threshold. Such assignment may be thought of as a natural experiment. Assume that the ‘assignment covariate’ has a *smooth* relationship with the outcome variable, which can be captured using parametric or semi-parametric models, and infer causal effects from discontinuity of the conditional expectation of the outcome variable related to assignment to treatment, which was determined by the ‘forcing’ variable being just below or just above the assignment threshold.<sup>131</sup> The position of the forcing covariate just above/below the threshold plays the role of the IV here.<sup>132</sup>

---

<sup>130</sup>First, compare the Probits with and without IVs. (Strictly speaking, the bivariate Probit assumptions (linearity and normality) are enough to identify  $\alpha$  even in absence of exclusion restrictions.) Second, enter in the second stage not  $\Phi(X_i'\beta + Z_i'\lambda)$  but two separate predicted probabilities for  $CH_i$  holding either  $X_i$  or  $Z_i$  constant at their sample means. This is pretty informal, but it helps to focus on where the identification works: off the IV or not.

<sup>131</sup>Clearly, there is some need for ‘local’ extrapolation (there is 0 common support in terms of the terminology of Section 20.2), so one assumes that the conditional regression function is continuous.

<sup>132</sup>When it predicts perfectly, we talk of *sharp* regression discontinuity. When the first-stage R2 is below unity, we talk of *fuzzy* regression discontinuity. See Hahn, Todd and VanderKlaauw (2001).

Start by  $Y_i = \alpha T_i + \epsilon_i$  with  $g(x) = E[\epsilon_i | X_i = x]$  and continuity around  $x = x^*$  and  $T_i = 0$  if  $X_i < x^*$  and equals 1 otherwise. As in the Wald estimator

$$\alpha = \frac{\lim_{x \uparrow x^*} E[Y_i | X_i = x] - \lim_{x \downarrow x^*} E[Y_i | X_i = x]}{\lim_{x \uparrow x^*} E[T_i | X_i = x] - \lim_{x \downarrow x^*} E[T_i | X_i = x]}.$$

In practice, this can be implemented using a Kernel regression (with  $h$  being the bandwidth) as

$$E[Y|X = x] \approx \frac{\sum_{i=1}^N K\left(\frac{X_i - x}{h}\right) Y_i}{\sum_{i=1}^N K\left(\frac{X_i - x}{h}\right)}$$

or we run a (local linear or polynomial) regression on  $E[Y_i | X_i, T_i] = \alpha T_i + g(X_i)$ .

**Example 9.5.** Matsudaira (2009, *JEcm*) studies the effect of a school program that is mandatory for students who score on a test less than some cutoff level.

**Example 9.6.** Or think of election outcomes that were just below or just above 50%.

**Example 9.7.** Angrist and Lave (1998) study of the class-size effect using the Maimonides rule: not more than 40 pupils per class. Class size is endogenous because of potential quality sorting etc. Assuming cohorts are divided into equally sized classes, the predicted class size is

$$z = \frac{e}{1 + \text{int}[(e - 1)/40]},$$

where  $e$  denotes the school enrollment. Note that in order for  $z$  to be a valid instrument for actual class size, one must control for the smooth effect of enrollment because class size increases with enrollment as do test scores.

**Remark 81.** Lee and Lemieux (2009) argue that regression discontinuity designs are best viewed as a close “cousin” of randomized experiments (see Section 20.1), not as a special case of IV methods (Angrist and Krueger, 1999) or matching<sup>133</sup> (Heckman et al. 1998, see Section 20.2). DiNardo and Lee (2010) show how the presence of “local random assignment” around the threshold is not just a maintained assumption, but a consequence of a structural assumption (with testable predictions) about the extent to which agents can manipulate their assignment.

<sup>133</sup>For an introduction, see <http://wol.iza.org/articles/matching-as-regression-estimator.pdf>

**Remark 82.** See Gilraine (2020, *JOLE*) for how entry into vs. exit from treatment can identify distinct treatment effects within the RDD setup.

For estimation and specification testing of regression discontinuity models, see the monograph and SW packages by Calonico, Cattaneo and coauthors.<sup>134</sup> For example, the `rdwselect` package in Stata guides the choice of optimal RD bandwidth for assignment to treatment based on Calonico et al (2014 *Ecm*). For RDD estimation conditioning on (pre-determined) covariates, see Calonico et al. (2019, *REStat*) and `rdrobust`. Cattaneo, Jansson and Ma (2018, *Stata Journal*) provide a test based on the continuity of the pushing variable around the threshold.<sup>135</sup>

See also Canay and Kamat (2017, *REStud*) for a new test of the credibility of the RDD design. See Kolesar and Rothe (2018, *AER*) for inference in RD with a discrete running variable. See Bertanha and Imbens (2019, *JBES*) on external validity of RDD, Gelman and Imbens (2018, *JBES*) on why not to use high-order polynomials in RDD. For earlier contributions see Imbens and Lemieux (2007, 2008),<sup>136</sup> Van der Klaauw (2008), Lee and Lemieux (*JEL* 2010). Gelman and Imbens (2014, NBER WP. No. 20405) argue that local smoothing should be used to control for the forcing variable, not global polynomials.

**Example 9.8.** Palguta and Pertold (2017, *AEJ: Policy*) provide an example of a situation where the observations on either side of a threshold are certainly not randomly assigned.

**Remark 83.** Hausman and Rapson (2021, NBER WP no. 23602) discuss when the RDD framework is (not) applicable in time (*RDiT*). Along similar lines, see Cellini, Ferreira, and Rothstein (2010, *QJE*) for a ‘dynamic RDD’ estimator applicable in panel data where treatment may (almost) happen in multiple periods. They retain all of the data by absorbing variation coming from non-close observations using flexible controls for the ‘assignment covariate’, allow for lags (dynamic structure) in the treatment effect, and control for cross-sectional as well as time fixed effects. In effect they run a least squares conventional panel data regression with a polynomial in the ‘assignment covariate’ and a treatment effect with lag-structure. In event studies, we’d look only shortly before and after a treatment change.

<sup>134</sup><https://sites.google.com/site/rdpackages/replication/cit-2019-cup>

<sup>135</sup>See <https://rdpackages.github.io/rdrobust/>

<sup>136</sup>It is an NBER WP no. 13039 and also the introduction to a special issue of the *Journal of Econometrics* on regression discontinuity. You can wait with reading this one until after we cover Section 20.

**Remark 84.** Within the RD framework, Grembi, Nannicini, and Troiano (2016) propose a ‘difference in discontinuities’ method to identify the effect of a new policy when a previous policy was established using the same factor variable and cut-off of the new one. Millán-Quijano (2020) extends this method to a fuzzy RD situation.

**Remark 85.** Armstrong and Kolesár (2018, Ecm) study CIs for the RDD, in particular whether data-dependent tuning parameters can improve upon minmax CIs.

### 9.3. Dealing with Weak Instruments

IV is an asymptotic estimator, unlike OLS.  $\hat{\beta}_{IV} = (Z'X)^{-1}Z'Y = \beta_0 + (Z'X)^{-1}Z'\epsilon$  and so  $E[\hat{\beta}_{IV}|X, Z] = \beta_0 + (Z'X)^{-1}Z'E[\epsilon|X, Z]$  and since  $E[\epsilon|X] \neq 0$ ,  $E[\epsilon|X, Z] \neq 0$  for some  $Z$ . Hence, 2SLS is biased in small samples and it needs large samples to invoke consistency.

Other than testing for  $COV(Z, \epsilon) = 0$ , one needs to consider the *weak instrument* problem, that is make sure that  $COV(X, Z) \neq 0$ . Even a small omitted variable bias ( $COV(Z, \epsilon) \neq 0$ ) can go a long way in biasing  $\hat{\beta}$  even in very large samples if  $COV(X, Z)$  is small because  $p \lim \hat{\beta} = \beta_0 + COV(Z, \epsilon)/COV(X, Z)$ . See [W]5.2.6.<sup>137</sup>

When IVs are weak (and many), 2SLS is biased (even in really big samples).<sup>138</sup> The finite sample bias of 2SLS is larger when instruments are weaker and when there are more instruments (high over-identification), leading to a trade-off between efficiency and bias (Hahn and Hausman, 2002).<sup>139</sup> Furthermore, weak IVs also bias inference in 2SLS, making standard errors too small.

**Remark 86.** With grouped-data IV, the IV will be weak if the groups are small so there is not enough averaging out of the unobservables.

<sup>137</sup>The ratio may not be well approximated by the Normal distribution even in large samples when  $COV(X, Z) \simeq 0$ , but by Cauchy (the ratio of two normals centered close to zero).

<sup>138</sup>The bias of over-identified 2SLS with weak IVs goes towards the probability limit of the corresponding OLS (inconsistent itself thanks to endogeneity) as the first stage gets closer to 0. Just-identified 2SLS is approximately unbiased, but that’s of little help since with a weak IV 2SLS will be very noisy.

<sup>139</sup>According to Bound et al. (1995), the weak IV bias is proportional to  $(k - 2)\sigma_{\epsilon\nu}/NR^2\sigma_x^2$ , where  $\nu$  is the first-stage residual,  $k$  is the number of IVs, and  $R^2$  is the unadjusted first-stage R-squared. However, this bias formula is based on approximations which may not work well with some weak IVs (Hahn and Hausmann, 2001; Cruz and Moreira, 2005).

Two responses have been developed: (a) test for the presence of weak IVs so you don't use them, (b) devise estimation methods that work with weak IVs.

(a) Bound et al. (1995), Staiger and Stock (1997) suggest the use of F tests in the first stage to detect weak IVs. Stock and Yogo (2005) formalize/develop their original F rule of thumb test.<sup>140</sup> Montiel Olea and Pfueger (2013, 2014) improve upon the Stock and Yogo (2005) F test by allowing for heteroskedasticity (clustering).<sup>141</sup>

**Remark 87.** *A simple pre-testing approach is to estimate the so-called the reduced form regression of  $y$  on the exogenous  $x$  and on the instruments (while excluding the endogenous  $x$ ). Whenever the reduced form estimates are not significantly different from zero, one presumes that “the effect of interest is either absent or the instruments are too weak to detect it” (Angrist and Krueger 2001; Chernozhukov and Hansen, 2008). This is so because the IV (2SLS) estimator can be written as the ratio of the reduced-form and first-stage coefficients.*

(b) When considering estimation methods, we need to distinguish two situations: having few or many (weak) IVs:

(i) Let us start with having few weak IVs: Even when IVs are weak, there is hope for correct inference (around wrong point estimates). Moreira (2003) develops a conditional likelihood ratio test, which overcomes the distortion in standard errors: he changes critical values of the test using an estimate of the first-stage-regression coefficients on instruments. He then calculates confidence intervals (around the 2SLS point estimates) as the set of coefficient values that would not be rejected at a given level of statistical significance. These need not be symmetric. This test statistic is robust to weak instruments.<sup>142</sup> The argument is that

---

<sup>140</sup>The test is `ivreg2` in `Stata`.

See [http://mayoral.iae-csic.org/IV\\_2015/IVGot\\_lecture3.pdf](http://mayoral.iae-csic.org/IV_2015/IVGot_lecture3.pdf)

<sup>141</sup>See `weakivtest` in `Stata`. The test works only for a single endogenous regressor. Otherwise use the Kleibergen-Paap Wald statistic in `ivreg2`. Also note that pretesting for weak IVs changes the distribution of 2SLS (see our discussion of sequential testing in Section 4.2) as Andrews and Stock (2005) suggest. Also, see Cruz and Moreira (2005) for an argument that the ability to use standard inference in the second stage depends not only on the F from the first stage, but also on the degree of endogeneity of  $X$ . This is their motivation for developing the method described in point (b).

<sup>142</sup>The critical values are random. The conditional test statistic is available in `Stata` (get the latest version). See `condivreg`, Moreira and Poi (Stata Journal, 2003), Mikusheva and Poi (Stata Journal, 2006), Cruz and Moreira (JHR, 2005), and Moreira, Andrews and Stock in the Journal of Econometrics.

even with low first-stage  $F$ , it may be possible to obtain informative confidence intervals.<sup>143</sup> Of course, this method is mainly useful in the just-identified case with one endogenous variable<sup>144</sup> (when you pick your best IV and report a specification that's not over-identified) or with low degrees of over-identification, as in these cases, 2SLS will not be very misleading in most cases. Angrist and Kolesar (2021, NBER 29417) argue that in standard applications the just-identified IV estimator can typically be treated as all but unbiased (in terms of weak-IV bias).

(ii) When the number of instruments grows with sample size (but its ratio to sample size converges), the literature employs alternative asymptotics (Bekker, 1994 Ecm), which imply that 2SLS are inconsistent even with strong (boundedly informative) instruments. Even when there are many instruments, the Limited information maximum likelihood (LIML)<sup>145</sup> is consistent under both strong and weak IVs (even with many exogenous regressors), asymptotically normal (Hansen, Hausman, and Newey, 2008 JBES) and approximately unbiased for over-identified models.<sup>146</sup> However, with heteroskedasticity, LIML becomes inconsistent (Bekker and van der Ploeg, 2005). But a jackknife estimator of Angrist, Imbens, and Krueger (1999) is still consistent (see Hausman, Newey, Woutersen, Chao, and Swanson, 2012, for the practical upshot). For estimators usable with weak IVs, see Andrews and Stock (2005) and the Fuller estimators 2 and 4 in `ivreg2`. Over-identification tests also need to be modified in the many-IV case (Lee and Okui, 2012 JEcm; Anatolyev and Gospodinov, 2011 ET).

**Example 9.9.** *Returning to Angrist and Krueger (1991) IV estimation of returns to education, where the number of IVs is large, but does not grow with sample size, Cruz and Moreira (2005) show that in some specifications with relatively low first-stage  $F$  statistics, the confidence intervals around the wrong 2SLS estimate is not too far off from the original range given by Angrist and Krueger (A-K). In*

<sup>143</sup>However, the Moreira method is not asymptotically valid if one first decides to check the  $F$  of the first stage and proceeds with the Moreira method only when the  $F$  is not totally tragic. And one always does this. See Chioda and Jansson (2006) for a solution.

<sup>144</sup>For a generalization to multiple endogenous regressors see Kleibergen (2007, Journal of Econometrics).

<sup>145</sup>See Anderson (2005, JEcm) for the history of LIML, which was introduced by Anderson and Rubin in 1949. It is equivalent to 2SLS in the exactly identified case. It minimizes 
$$N \frac{(y - X\beta)' Z (Z' Z)^{-1} Z' (y - X\beta)}{(y - X\beta)' (y - X\beta)}.$$

<sup>146</sup>See Chao and Swanson (2005, Ecm). For a somewhat outdated survey of available solutions in highly over-identified cases, see Section 13 of the lecture notes to the Imbens/Wooldridge NBER Summer 07 course and the Flores-Lagunes (2007) simulation study.

the most interesting specification, they estimate a confidence interval that has the same lower bound as that given by A-K, but a much higher upper bound. They say this makes the correct  $\beta$  “likely to be larger” compared to that given by A-K.

**Remark 88.** Andrews, Stock, and Sun<sup>147</sup> provide a recent review of the weak IV literature.

**Remark 89.** Young (2017) provides bootstrap-based criticism of the weak IV pre-tests and weak-IV-robust methods.<sup>148</sup>

**Remark 90.** Of course, at the end, you test whether there is any endogeneity affecting  $\beta$ . This is typically done by Hausman, which simplifies to an auxiliary regression test where we include the first stage predicted residual in the second stage and test its significance ([W] p.118)<sup>149</sup>

**Remark 91.** It is still possible to estimate IV if only the instrument  $Z$  and  $y$  are in one data set and the instrument and the endogenous  $X$  are available in another data set (Angrist and Krueger, 1992). In this case, Inoue and Solon (2010, REStat) suggest the use of the two-sample 2SLS estimator  $\hat{\beta}_{TS2SLS} = (\hat{X}'_1 \hat{X}_1)^{-1} \hat{X}'_1 y_1$ , where  $\hat{X}_1 = Z_1 (Z'_2 Z_2)^{-1} Z'_2 X_2$  instead of the two-sample IV estimator  $\hat{\beta}_{TSIV} = (Z'_2 X_2 / n_2)^{-1} Z'_2 y_1 / n_1$ .<sup>150</sup> See also Arellano and Meghir (1991) for a Minimum-Distance version.

**Remark 92.** Actually, this logic applies to OLS just as well. We can replace any of the sample averages in the OLS formula by other averages that have the right plim. For example, we can ignore missing observations on  $y$  when calculating the  $X'X$  plim, use other data, “pseudo panels”, etc.

**Remark 93.** A simple and probably powerful test of joint IV validity and non-weakness is proposed by Hahn and Hausman (2002, Econometrica). It consists

<sup>147</sup><https://scholar.harvard.edu/iandrews/publications/weak-instruments-iv-regression-theory-and-practice>

<sup>148</sup><http://personal.lse.ac.uk/YoungA/ConsistencyWithoutInference.pdf>

<sup>149</sup>See also <http://www.stata.com/statalist/archive/2002-11/msg00109.html> .

<sup>150</sup>With exact identification, TS-2SLS is asymptotically equivalent to TS-IV, but with over-identification TS-2SLS is more efficient because it implicitly corrects for differences in the empirical distributions of the instrumental variables between the two samples.

of comparing the forward two stage estimate of the endogenous-variable coefficient to the inverse of the estimate from the reverse regression (the right hand side endogenous variables becomes dependent variable and the dependent variable from the forward regression becomes the new right-hand side variable) using the same instruments. It uses a more robust 2nd order asymptotic justification and compares two estimates of a structural parameter such that one can assess the economic magnitude of the departure from  $H_0$ .

**Remark 94.** When there are multiple endogenous variables (and multiple IVs), use the weak IV test by Sanderson & Windmeijer (2014, *Journal of Econometrics*).

**Remark 95.** Consider a model with many instruments that are weak, but valid (e.g., lags of stock returns). It is not desirable to pick just a few of them (e.g. the most recent lags) and run 2SLS. Instead, run the Continuously-Updated Estimator (CUE) with many weak instruments, which is a heteroskedasticity-autocorrelation-robust version of LIML. The structural coefficients are consistent and asymptotically normal. Using few weak instruments produces the asymmetric Moreira CI's or the discontinuous Andrews-Armstrong CI's.

## 10. GMM and its Application in Panel Data

Read at least one of the two handouts on GMM which are available in the reference folder for this course in the library. The shorter is also easier to read.

Theory (model) gives us population orthogonality conditions, which link the data to parameters, i.e.,  $E[m(X, Y, \theta)] = 0$ . The GMM idea: to find the population moments use their sample analogues (averages)  $\sum_{i=1}^N m(X_i, Y_i, \theta) = q_N$  and find  $\hat{\theta}$  to get sample analogue close to 0.

If there are more orthogonality conditions than parameters (e.g. more IV's than endogenous variables) we cannot satisfy all conditions exactly so we have to weight the distance just like in the MD method, and the resulting minimized value of the objective function is again  $\chi^2$  with the degrees of freedom equal to the number of overidentifying conditions. This is the so called **Hansen test** or J test or GMM test of overidentifying restrictions:

$$\hat{\theta}_N^{GMM} = \arg \min \{q_N(\theta)' W_N q_N(\theta)\} \quad (10.1)$$

To reach  $\chi^2$  distribution, one must use the optimal weighting matrix,  $\widehat{V(m)}^{-1}$ , so that those moment conditions that are better estimated are forced to hold more



closely (see Section 8.3.1 for similar intuition). A feasible procedure is to first run GMM with the identity matrix, which provides consistent  $\widehat{\theta}$  and use the resulting  $\widehat{\varepsilon}$ s to form the optimal weighting matrix.

**Remark 96.** *GMM nests most other estimators we use and is helpful in comparing them and/or pooling different estimation methods.*

**Example 10.1.** *OLS:  $y = X\beta + \varepsilon$ , where  $E[\varepsilon|X] = 0 \implies E[X'\varepsilon] = 0$  so solve  $X'(y - X\widehat{\beta}) = 0$ .*

**Example 10.2.** *IV:  $E[X'\varepsilon] \neq 0$  but  $E[Z'\varepsilon] = 0$  so set  $Z'(y - X\widehat{\beta}) = 0$  if  $\dim(Z) = \dim(X)$ . If  $\dim(Z) > \dim(X)$  solve 10.1 to verify that here  $= \widehat{\beta}_{TSLS}$ .*

**Example 10.3.** *Non-linear IV:  $y = f(X)\beta + \varepsilon$ , but still  $E[Z'\varepsilon] = 0$  so set  $Z'(y - f(X)\widehat{\beta}) = 0$ .<sup>151</sup>*

**Example 10.4.** *Euler equations:  $E_t[u'(c_{t+1})] = \gamma u'(c_t) \implies E_t[u'(c_{t+1}) - \gamma u'(c_t)] = 0$ . Use rational expectations to find instruments:  $Z_t$  containing information dates  $t$  and before. So  $E_t[Z_t'(u'(c_{t+1}) - \gamma u'(c_t))]$  is the orthogonality condition. Note that here  $\varepsilon$  is the forecast error that will average out to 0 over time for each individual but not for each year over people so we need large  $T$ .*

**Example 10.5.** *Maximum likelihood: Binary choice: Logit: Score defines  $\widehat{\beta}$ :*

$$\log L(\beta) = \log \prod_{n=1}^N \frac{\left(\exp(\beta' x_i)\right)^{y_i} 1^{1-y_i}}{1 + \exp(\beta' x_i)} = \beta' \sum_{n=1}^N x_i y_i - \sum_{n=1}^N \log \left(1 + \exp(\beta' x_i)\right)$$

$$\frac{\partial \log L(\beta)}{\partial \beta} = \sum_{n=1}^N x_i \left[ y_i - \frac{\exp(\beta' x_i)}{1 + \exp(\beta' x_i)} \right] = \sum_{n=1}^N x_i \varepsilon_i = 0$$

**Example 10.6.** *One can use GMM to jointly estimate models that have a link and so neatly improve efficiency by imposing the cross-equation moment conditions. For example, Engberg (1992) jointly estimates an unemployment hazard model (MLE) and an accepted wage equation (LS), which are linked together by a selection correction, using the GMM estimator.*

<sup>151</sup>See Kelejian (1971, JASA) for 2SLS with equations linear in parameters but non-linear in endogenous variables (such as a quadratic in X).

**Remark 97.** GMM does not require strong distributional assumptions on  $\varepsilon$  like MLE. Further, when  $\varepsilon$ s are not independent, the MLE will not piece out nicely, but GMM will still provide consistent estimates.

**Remark 98.** GMM is consistent, but biased in general as  $E[\widehat{\beta}_{IV}|X, Z] - \beta_0 = (Z'X)^{-1}E[\varepsilon|X, Z] \neq 0$  because  $E[\varepsilon|X, Z] \neq 0$  for some  $Z$  because  $E[\varepsilon|X] \neq 0$ . GMM is a large sample estimator. In small samples it is often biased downwards (Altonji and Segal 1994).

**Remark 99.** GMM allows us to compute variance estimators in situations when we are not using the exact likelihood or the exact  $E[y | x]$  but only their approximations.<sup>152</sup>

**Remark 100.** The GMM test of over-identifying restrictions can often fail to reject in IV situations because the IV estimates are often imprecise. When IV estimates are precise, on the other hand, rejection may correspond not to identification failure, but, rather, to treatment effect heterogeneity (see Section 20.3).

**Example 10.7.** The GMM analogue to 2SLS with general form of heteroskedasticity is

$$\widehat{\beta}_{GMM} = (X'Z\widehat{\Omega}^{-1}Z'X)^{-1}X'Z\widehat{\Omega}^{-1}Z'Y \quad (10.2)$$

and with panel data we can apply the White (1980) idea to estimate  $\widehat{\Omega}$  while allowing for any unconditional heteroskedasticity and for correlation over time within a cross-sectional unit:

$$\widehat{\Omega} = \sum_{i=1}^N Z_i' \widehat{\varepsilon}_i \widehat{\varepsilon}_i' Z_i$$

where the  $\widehat{\varepsilon}_i$  comes from a consistent estimator such as homoscedastic TSLS.<sup>153</sup>

**Remark 101.** See Cho and Phillips (2018, JEcm) for a simple-to-implement test of covariate matrix equality.

**Exercise 10.1.** Show that even with heteroscedastic errors, the GMM estimator is equivalent to TSLS when the model is exactly identified.

<sup>152</sup>See section 5 of the GMM handout by George Jakubson in the library.

<sup>153</sup>For differences between the traditional 3SLS and GMM, see [W]8.3.5.

**Exercise 10.2.** Compare the way we allow for flexible assumptions on the error terms in the estimator 10.2 to the strategy proposed in section 8.3.2.

**Example 10.8.** Nonlinear system of simultaneous equations. Euler equations.

McFadden (1989) and Pakes (1989) allow the moments to be simulated: SMM (see Remark 125). Imbens and Hellerstein (1993) propose a method to utilize exact knowledge of some *population* moments while estimating  $\theta$  from the *sample* moments: reweight the data so that the transformed sample moments would equal the population moments.

## 11. Dynamic Panel Data

Up to now, we relied mainly on the assumption of time-constant unobservables ( $\alpha_i$ ) to produce difference-in-differences estimates. But in many applications, this notion is not attractive.

For example, when evaluating labor-market interventions, there is a famous temporary drop in  $y$  right before the treatment takes place (the Ashenfelter's dip, Dehejia and Wahba, JASA 1999). Abbring and van den Berg (2003) point out that when assignment to treatment (policy change) is anticipated, individuals may react in an attempt to explore the policy rules and such anticipatory behavior would render 'before' individuals with similar characteristics incomparable, invalidating the D-in-Ds design.

One usually builds dynamics into the model by including  $y_{t-1,i}$  on the right-hand side, i.e.,  $y_{it} = \alpha_i + \rho y_{it-1} + \beta x_{it} + \epsilon_{it}$ . (This is often motivated with a partial adjustment model where agents adjust  $y$  over time to reach a latent target.) In a typical application, if we want to do both fixed effects and lagged  $y$ , we're probably asking too much of the data.<sup>154</sup> The problem is that after we get rid of  $\alpha_i$  by first differencing,  $\Delta\epsilon_{it}$  is correlated with  $\Delta y_{i,t-1}$  because they both depend on  $\epsilon_{i,t-1}$ . Solving this problem requires instrumenting with longer lags (having a long panel), but if there is too much serial correlation in  $\epsilon_{i,t}$  (a realistic situation) there may be no consistent estimator.<sup>155</sup>

<sup>154</sup>See Guryan (AER, 2004) for an argument that comparing the FEM (with no  $y_{i,t-1}$ ) to a regression with no  $\alpha_i$  but with  $y_{i,t-1}$  provides a useful interval of estimates.

<sup>155</sup>See this survey of production function estimation for detailed discussion of some of these issues: <http://www.economics.ox.ac.uk/Research/wp/pdf/paper513.pdf>

See the Arellano and Bond (1991) method (in `Stata`) or Blundell and Bond (1998).<sup>156</sup> The idea is to use as IVs not only  $y$  levels, but also differences, more of them (longer lags) as they become available with time, and also  $x$ s depending on whether one assume exogenous or only predetermined  $x$ s.<sup>157</sup> The consensus is that Arellano and Bond works well in small samples when  $|\rho| < 0.8$  but not for more strongly auto-regressive  $y$  processes. Instruments get weaker as  $|\rho|$  is closer to 1 and as the  $VAR(\alpha_i)/VAR(\epsilon_{it})$  increases. The solution of Blundell and Bond is to make the strong assumptions of stationarity and of  $E[(y_{it} - y_{is})\alpha_i] = 0$  and  $E[(x_{it} - x_{is})\alpha_i] = 0$ . The way out of this is offered by Hahn, Hausman and Kuersteiner (2007, JEcm), who propose to use only the longest available difference (and to thus turn the data into a cross-section) as this is the most informative moment to use when  $|\rho|$  is closer to 1. Aquaro and Čížek generalize this idea to multiple pairwise differences to improve the MSE.

## 12. Other Regression Applications

### 12.1. Quantile Regression

A regression method robust to outliers and censoring. Let's start with the 50th quantile. Sometimes we want to fit the regression line through medians of  $y|x$ , not the means, as OLS does. The LAD estimator

$$\min_{\beta} \sum_{i=1}^n |y_i - x_i' \beta|$$

corresponds to such median regression. One goal of this regression method is to reduce the influence of outliers (the LAD regression is not affected by taking a value of one  $y$  that is above the median and multiplying it by 1000). Quantile regressions only use ordinal information; see, e.g. the use of CLAD in censored data in Section 18.3.

However, note that OLS and LAD will only give the same answer when the distribution of  $\epsilon|x$  is symmetric around zero. A comparison of LAD and OLS

---

<sup>156</sup>Also, Arellano and Bover (1995) propose a system estimator that jointly estimates the regression in levels and in differences in order to re-incorporate levels variation and reduce the likelihood of weak-instruments bias. See also the Arellano and Honore chapter in the Handbook of Econometrics on new IV methods for dynamic panel models.

<sup>157</sup>For `Stata` commands see `xtabond` and `xtdpdsys`.

when  $y$  is a measure of wealth (a skewed  $y$ ) therefore does not necessarily speak about the importance of outliers.

Now consider estimating regression parameters  $\beta(q)$  corresponding to any quantile  $q$  you are interested in:

$$\min_{\beta(q)} \sum_{i=1}^n \left( y_i - x_i' \beta \right) \left( q - I \left[ y < x_i' \beta \right] \right).$$

The objective function is piece-wise linear and continuous. For intuition on this objective, think first of  $q = 0.5$ . In this case, the second bracket is either  $+0.5$  or  $-0.5$ .<sup>158</sup> When you are above (below) the regression line, you get  $0.5$  ( $-0.5$ ) multiplying a positive (negative)  $\epsilon$ . So, the objective function is just trying to minimize the sum of absolute values of  $\epsilon$ s (LAD), i.e., fit the line through (conditional) medians. Now think of  $q = 0.1$ . When you are above (below) the quantile regression line, you multiple  $|\epsilon|$  by  $0.1$  ( $0.9$ ). Why? You are minimizing and you want to punish for predicting below the line, because you want 90% of the points to be above the line.

You need to bootstrap for standard errors. There are now IV methods for quantile regressions as well as panel data quantile regressions.<sup>159</sup> See Abadie, Angrist and Imbens (2002, *Econometrica*) or Powel (1983) or Chernozhukov and Hansen (2006, *Econometrica*) for quantile IV regressions.<sup>160</sup>

See Section 14 of the Imbens/Wooldridge NBER Summer '07 course for panel data applications. In particular, the correlated random-effects quantile regression of Abrevaya and Dahl (2007) and the penalized fixed effects estimator of Koenker (2004). See Remark 178 for a reference to a sample selection correction for quantile regressions. B. Melly provides codes for quantile treatment effect estimation in Stata or quantile regression decomposition of differences in distribution.<sup>161</sup> Similarly, Machado and Mata (2005, *JAppliedEcm*) provide a Oaxaca-Blinder style gap decomposition for quantile regressions.

**Remark 102.** *Firpo, Fortin, and Lemieux (2009, Ecm) propose unconditional quantile regressions to measure the effects of treatment on the unconditional  $y$  distribution. A coefficient from the standard Koenker and Bassett (1978) (conditional) quantile regression does not deliver this object.*

<sup>158</sup>You can think of this the *sgn* function. See the MRE (Section 14.1.5).

<sup>159</sup>Of course, when the IVs vary at group level and we care about  $y$  distributions within groups, we can just run 2SLS at the group level with the group-specific quantiles on the LHS.

<sup>160</sup>Liu (2014, *J Popul Econ*) is a nice example using the last paper's methodology.

<sup>161</sup>See [http://www.econ.brown.edu/fac/Blaise\\_Melly/codes.html](http://www.econ.brown.edu/fac/Blaise_Melly/codes.html)

**Remark 103.** *Quantile regressions tell us about  $y$  distributions (for example, wage inequality) within  $x$  groups. Quantile regression coefficients speak about effects on distributions, not on individuals. For example, if  $T_i$  has a positive effect on a lower decile of the income distribution, this does not mean that someone who would have been poor will be less poor with treatment  $T_i$ , but that those who are poor with treatment are less poor than they would have been without treatment.*

**Remark 104.** *For treatment of classical measurement error in the dependent variable in a quantile regression, see Hausman et al. (2019, NBER WP No. 25819).*

**Remark 105.** *Abrevaya, Hsu, and Lieli (2015, JBES) consider the conditional average treatment effect (CATE) consisting of local treatment regressions and designed to capture the heterogeneity of a treatment effect across subpopulations when the unconfoundedness assumption applies. In contrast to quantile regressions, the subpopulations of interest are defined in terms of the possible values of a set of continuous covariates rather than the quantiles of the potential outcome distributions.*

## 12.2. The Reflection Problem

See Manski (1995). Sometimes you ask why individuals belonging to the same group act in a similar way. Economists usually think this may be because their  $x_i$  are similar or because the group has a common characteristic  $z_g$  (for example ethnic identity):

$$y_{ig} = \alpha + \beta' x_{ig} + \gamma' z_g + \varepsilon_{ig}.$$

Sociologists add that an individual may act in some way because other individuals within the group act that way, that is because of  $E[y|z]$  (herd behavior, contagion, norm effects; there is a social multiplier, an endogenous social effect) or because the individual outcome varies with the mean of the exogenous variables in the reference group  $E[x|z]$  (exogenous social effect).

**Example 12.1.** *Does the propensity to commit crimes depend on the average crime rate in the neighborhood or on the average attributes of people living in the neighborhood or on some exogenous characteristics of the neighborhood like the quality of schools etc.? Or think of high-school achievement as another example.*

Note that the regression

$$E[y|z, x] = \alpha + \beta'x + \gamma'z + \delta'E[y|z] + \lambda'E[x|z] \quad (12.1)$$

has a social equilibrium: taking an expectation of both sides w.r.t.  $z$  we get

$$E[y|z] = \alpha + \beta'E[x|z] + \gamma'z + \delta'E[y|z] + \lambda'E[x|z],$$

which you can solve for  $E[y|z]$  and plug back into equation 12.1 to show that in the reduced form equation, in which  $E[y|z, x]$  depends on  $(1, E[x|z], z)$ , the structural parameters are not identified. We cannot separately identify the so-called correlated, endogenous, and contextual effects. Under some conditions, we can say if at least one of the social effects is present, but we cannot determine which one.

See also Durlauf (2002) who trashes the recent social-capital (SC) literature by applying the same identification logic. He first considers an individual-level regression of the following type

$$y_{ig} = \alpha + \beta'x_{ig} + \gamma'z_g + \delta E(y_g|F_g) + \theta E(SC_g|F_g) + \varepsilon_{ig}, \quad (12.2)$$

where  $F_g$  is the collection of information one has for each group  $g$ . In this notation,  $z_g$  corresponds to the contextual effects (variables measured at group level predetermined at the the time of choice, including averages of individual characteristics). If  $SC_g$  is predetermined ( $E(SC_g|F_g) = SC_g$ ), it is simply another contextual effect and identification requires the presence of at least one individual-level variable whose group level average does not causally affect individuals.

If  $SC_g$  is an endogenous outcome of decisions that are contemporary to the behavioral choice  $y_{ig}$ , i.e., when we need a separate theory as to what determines social capital, then one needs two elements of  $x_{ig}$  not to be elements of  $z_g$  so as to provide instruments for  $E(y_g|F_g)$  and  $E(SC_g|F_g)$ . That is one needs two individual characteristics that affect individual behavior yet whose group analogues are excluded from the behavioral equation 12.2. This results is based on considering the two simultaneous equations: one determining  $y$  (equation 12.2), the other for  $SC$  :

$$SC_{ig} = \tilde{\alpha} + \tilde{\beta}'x_{ig} + \tilde{\gamma}'z_g + \tilde{\delta}E(y_g|F_g) + \tilde{\theta}E(SC_g|F_g) + \eta_{ig}. \quad (12.3)$$

Finally, Durlauf (2002) considers the case of having only aggregate (group-level) data. The system then boils down to:

$$y_g = \alpha + \gamma'z_g + \delta E(y_g|F_g) + \theta E(SC_g|F_g) + \varepsilon_g \quad (12.4)$$

$$SC_g = \tilde{\alpha} + \tilde{\gamma}' z_g + \tilde{\delta} E(y_g | F_g) + \tilde{\theta} E(SC_g | F_g) + \eta_g. \quad (12.5)$$

and identification is a textbook case asking whether one has instruments for the two endogenous variables  $E(y_g | F_g)$  and  $E(SC_g | F_g)$ : Are there variables in this world that would affect social capital formation but not other behavior (like GDP growth in Knack and Keefer, QJE 1997)?

Brock and Durlauf (2001) show that the lack of identification between endogenous and contextual effects does not occur in binary and multinomial choice models, essentially because of their non-linearity. Zanella (2007, JEEA) applies the nested logit structure (see Section 16.0.5) to a random utility framework in order to build a model with social interactions *and* endogenous choice of the group membership (neighborhood); the micro-founded model then suggests econometric identification strategies. For econometrics of (consumer) binary choice with social interactions (spillover-effects), for example subsidies for health-product adoption and vouchers for attending a high-achieving school, see NBER WP No. 25947.

One solution to the reflection problem, see Borjas (1992), is to use assumptions about the dynamics of social processes and to run

$$E_t[y|z, x] = \alpha + \beta x_t + \gamma z_t + \delta E_{t-1}[y|z].$$

This has been the strategy of several recent peer-effects papers where one uses pre-determined  $y$  (ability, smoking), i.e.,  $y_{-i}$  of others before they joined the group, to study the effect on one's current  $y_i$ . However, even this is problematic even with perfect randomization into groups (to deal with selection) as argued recently by Benjamin Elsnar as long as the  $y$  is continuous (ability, not being black). Randomizing entire groups (with continuous  $y$ ) assigns the whole distribution of (pre-determined)  $y$  (by group), not just the mean  $y_{-i}$ , but also other moments, the individual ranks, the distance to the first or last, etc. so one has to control for everything, not just the mean  $y$ . See the Gelbach JoLE paper on Oaxaca Blinder to deal with the order in which one introduces the various moments. For comparison, adding one marginal person to a group (e.g., using RD designs on thresholds of acceptance to selective education programs) allows one to run a clear reduced form for the effect of the group on the individual (plus the peer effect of the individual on the group).

Network data on the structure of interactions can really help. Bramoullé, Djebbari, and Fortin (2009) show that the reflection problem is solved in the presence of second order peers linked to first order peers but not to the individual, under the exogenous network formation assumption. See Johnsson and Moon (2015) for a related recently developed estimator. See Zacchia and Fernández



(2019) for identification of peer effects in networks without the assumption of the model's error term being conditionally independent of  $X$ s and of the structure of interactions.

### 12.3. Oaxaca-Blinder Decompositions

Often, you want to use Least Squares regressions to explain (account) the sources of the difference in the outcomes across two groups of workers, countries, etc. (Think of regression as a conditional expectation.) For example, a vast and ultimately unsuccessful literature aimed at measuring the extent of wage discrimination has followed Oaxaca (1973) and Blinder (1973) in decomposing the overall mean wage difference between the advantaged (men) and disadvantaged (women) into two parts: the first reflecting the difference in average productive endowments of individuals in each group and the second part due to the differences in coefficients. Following this approach, one first estimates logarithmic wage regressions separately for each gender, controlling for explanatory variables. The decomposition technique relies on the fact that the fitted regressions pass through the sample means<sup>162</sup> as follows:

$$\overline{\ln w_g} = \widehat{\beta}_g' \overline{X_g}, \quad g \in \{f, m\}, \quad (12.6)$$

where  $f$  denotes females and  $m$  denotes males,  $\overline{\ln w_g}$  is the gender-specific mean of the natural logarithm of hourly wage, and where  $\overline{X_g}$  represents the respective vectors of mean values of explanatory variables for men and women. Finally,  $\widehat{\beta}_m$  and  $\widehat{\beta}_f$  are the corresponding vectors of estimated coefficients. A general form of the mean wage decomposition is as follows:

$$\overline{\ln w_m} - \overline{\ln w_f} = (\overline{X_m} - \overline{X_f})' \widetilde{\beta} + [\overline{X_m}' (\widehat{\beta}_m - \widetilde{\beta}) + \overline{X_f}' (\widetilde{\beta} - \widehat{\beta}_f)], \quad (12.7)$$

where  $\widetilde{\beta}$  represents a counter-factual non-discriminatory wage structure. The first term on the right hand side of equation 12.7 represents that part of the total logarithmic wage difference which stems from the difference in average productive characteristics across gender. The second term originates in the differences in gender-specific coefficients from the non-discriminatory wage structure and is often interpreted as reflecting wage discrimination.<sup>163</sup>

<sup>162</sup>This idea does not work in quantile regressions. See Machado and Mata for the method applicable in median regressions.

<sup>163</sup>There have been objections to this decomposition approach. First, by focusing on the mean gap, it ignores meaningful differences in gender-specific wage distributions. Second, if character-

**Remark 106.** Using  $\beta_m$  or  $\beta_f$  for  $\tilde{\beta}$  corresponds to estimating the ATU or ATT, respectively (when being a female is the “treatment”). Sloczynski (2020, ILLR) argues that the interpretation of Oaxaca-Blinder-type decompositions depends on the relative sizes of subpopulations under study.<sup>164</sup>

**Remark 107.** Nopo (2004) and Black et al. (2005) and others now point out to matching as a preferred alternative to parametric methods when support is not perfectly overlapping. Lechner and Strittmatter (2019, *Econometric Reviews*) consider practical procedures dealing with the common support problems in matching estimators.

**Remark 108.** For the importance of careful interpretation of decompositions of this type in terms of discrimination, see the 2020 response of Durlauf and Heckman to Fryer’s (2019) paper on police shootings of black Americans (all in *JPE*).

**Remark 109.** <https://www.journals.uchicago.edu/doi/full/10.1086/710976>

**Remark 110.** Gelbach (*JoLE*, 2016) shows how the specific sequence of adding regressors affects the Oaxaca-Blinder-style decompositions in Neal and Johnson (1996). He develops a conditional decomposition: Start with  $Y = X_1\beta_1 + X_2\beta_2 + \epsilon$ , where we care about  $\beta_1$  (which captures the gap between two groups such as genders conditional on  $X_2$ ) and consider

$$\hat{\beta}_1^{base} = \left(X_1'X_1\right)^{-1} X_1'Y = \left(X_1'X_1\right)^{-1} X_1'X_1\beta_1 + \left(X_1'X_1\right)^{-1} X_1'X_2\beta_2 + \left(X_1'X_1\right)^{-1} X_1'\epsilon,$$

where  $\hat{\beta}_1^{base}$  is the regression coefficient from projecting  $Y$  on  $X_1$  alone. Now,  $E[\hat{\beta}_1^{base}] = \beta_1 + \Gamma\beta_2$ , where  $\Gamma$  is the matrix of coefficients from projecting columns of  $X_2$  on columns of  $X_1$ . The bias in the effect of the  $j$ -th element of  $X_1$  from excluding the  $k$ -th element of  $X_2$  from the equation (decomposition) is given by  $\sum_{k=1}^{K_2} \Gamma[j, k]\beta_2[k]$ . Testing the significance of this bias in a type of Hausman test.

There are a number of variants of this method depending on how one simulates the non-discriminatory wage structure  $\tilde{\beta}$ . Neumark (1988) and Oaxaca and

---

istics which might differ between males and females are omitted in the vector of regressors, the contribution of these characteristics will be captured by the constant term and will erroneously appear in the measure of discrimination.

<sup>164</sup>He also provides treatment-effects reinterpretations of the Reimers, Cotton, and Fortin decompositions.

Ransom (1994) suggest the use of regression coefficients based on pooled data including both men and women, arguing that they provide a good estimate of a competitive non-discriminatory norm.<sup>165</sup> Alternatively, one can use a similar approximation based on weighting the male and female coefficients with sample proportions of each sex (Macpherson and Hirsh, 1995).

It is not always clear how you apply the method in non-linear models (see Ham et al., 1998, AER). Recently, the decomposition has been extended to quantile (median) regressions by Machado and Mata (2000).<sup>166</sup> There is a versions of this decomposition for Probit (Myeong-Su Yun, 2004). Choe et al. (IZA DP No. 10530) study the serious consequences of the normalization on the error variance in Logit/Probit for sample mean probability decompositions. In a recent paper, he also adds standard errors for this decomposition.<sup>167</sup> Finally, there is an invariance problem that has to do with the choice of the base category (affecting the constant and hence the unexplained part).

There are important extensions taking the idea beyond first moments and into decomposing whole distributions. See DiNardo, Fortin, and Lemieux (1996, *Econometrica*) and Bourguignon, Ferreira, and Leite “Beyond Oaxaca-Blinder: Accounting for Differences in Household Income Distributions”. The DiNardo et al. decomposition has been programmed into Stata.<sup>168</sup> Matching-based weighting methods are to be found in Fortin and Firpo (2011). See also the chapter *Decomposition Methods in Economics in the Handbook of Labor Economics*.

**Remark 111.** *Maasoumi and Wang (2019, JPE) argue that instead of evaluating the gender wage gap at each quantile, which assumes rank invariance (male and female wages being ordered the same way according to skill levels), an assumption that is rejected by data, one should first characterize each gender-specific wage distribution and then evaluate their difference. They characterize distributions using entropy functions to assign more weight to lower wage levels (to express inequality aversion), provide tests of stochastic dominance rankings, and deal with selection into work (participation) by applying the Arellano and Bonhomme (2017) quantile-copula approach to jointly estimate wages and participation decisions in order to replace missing wages of those who do not work by their reservation wages.*

<sup>165</sup>Neumark (1988) provides a theoretical justification for this approach using a model of discrimination with many types of labor where employers care about the proportion of women they employ.

<sup>166</sup>See, e.g., Albrecht JOLE for an application of quantile regressions to gender wage gaps.

<sup>167</sup>See also Fairlie (2005) *Journal of Economic and Social Measurement*.

<sup>168</sup>For a discrete simple version of the DFL decomposition, see Biewen (2001, REStat).

*As a result, there is less gender wage convergence than one would think based on the traditional decompositions.*

#### 12.4. Meta Analysis

In medicine, and since the 1990s also in economics (especially since Stanley, 2008 OBES), we often combine estimates of parameters across studies to ask whether they converge to a meta-average that could be considered the true parameter. (Assuming no parameter heterogeneity in the simplest form.) Further, one can try to explain differences in estimates across studies using explanatory variables, either those related to the country where a given parameter was estimated, or those related to the estimation techniques used in various studies of the same topic. For example, there may be multiple studies exploring a given question (in detail) using one-country data. (In keeping with the example from the previous subsection, researchers often estimate the unexplained portion of the gender wage gap in one country.) Next, the question is how we cumulate knowledge across such studies. When you want to learn about the impact of institutions or policies on the unexplained portion of the wage gap you may collect data that consists of the estimates of other studies, which you then regress on explanatory variables capturing the country-time specific variables.<sup>169</sup>

There is another use of Meta analysis: When scientists report their results, they are naturally driven to report important useful findings, that is those that reject the null hypothesis of no effect. One can analyze the set of existing results to see if there is “reporting” “drawer” bias. That is, one can estimate a regression using the results from other studies, asking about the effect on the published results of the method of estimation used, type of data, etc. and the size of the standard error. Consider for example the estimation of returns to education. IV studies typically have larger standard errors and typically report larger (significant) returns. See Ashenfelter, Harmon and Oosterbeek (1999 Labour Economics).<sup>170</sup> If there is no bias in reporting, the estimates should not be correlated with their standard error. If, however, researchers are more likely to report higher estimates when standard errors increase (IV), this will result in sample selection (non-representative sample of all estimated results).

---

<sup>169</sup>See work by Winter-Ebmer and others explaining the gender wage gap across countries. It is important to know that there are a number of econometrics problems with this approach.

<sup>170</sup>For another application of meta-analysis see Card and Krueger “Myth and Measurement” book on minimum wages.

The simplest presentation of publication bias is in the form of a funnel plot where one plots precision (1/standard error) against the estimated parameter values. The plot should be symmetric, but often isn't. One can also regress the estimated parameters on their standard errors (this is the funnel asymmetry test, FAT), which is often estimated with precision as weights, clustered at study level, and on variables such as having a sponsor or the journal impact factor.

### 12.5. Misc Topics

Expectations: Manski (2004, Ecm). Regression to the Mean: Krueger and Mueller (2002). Test Scores on the LHS: "A Self-Reference Problem in Test Score Normalization," *Ec of Edu Review* and "Test Score Measurement and the Black-White Test Score Gap," *REStud* 2017.

Rank Inference: Rankings are often based on estimates, generating rank uncertainty. Mogstad et al. (2010, NBER WP 26883) construct confidence sets for the rank of each population.

## 13. Nonparametrics

Can we estimate a regression (a density) without any functional form or distributional assumptions?<sup>171</sup>

**Kernel estimation** A typical OLS regression will use information from the whole range of  $x \in [\underline{x}, \bar{x}]$  to estimate  $E[y_i | x = x_i] = \beta'x_i$ . Here, we will estimate a conditional expectation function  $E[y | x] = m(x)$  using 'local' information from an area  $A(x)$  'close' to  $x$ :

$$\widehat{E[y | x]} = \widehat{m(x)} = \frac{\sum_{i=1}^n I\{i \in A(x)\}y_i}{\sum_{i=1}^n I\{i \in A(x)\}} = \sum_{i=1}^n w_i(x)y_i.$$

Two questions: (i) how to define  $A(x)$ , (ii) are the weights  $w_i(x)$  from above optimal. Instead of the indicator function  $I\{\cdot\}$  let us use a bounded, symmetric *Kernel* function  $K(\cdot)$  such that  $\int K(u)du = 1$ . For asymptotic theory on choosing the optimal Kernel and bandwidth<sup>172</sup>, see [N] and Silverman (1986).

<sup>171</sup>We will come back to the use of semi-parametric methods in the estimation of discrete choice models (section 14.1.5) and apply non-parametric approaches when matching on unobservables (as in the selection bias model of Section 19.2.4) as well as observables (see section 20.2).

<sup>172</sup>The bandwidth can be also data-dependent.

**K-th Nearest Neighbor** Define  $J(x) = \{i : x_i \text{ is one of the } K \text{ nearest neighbors}\}$  and use  $w_i(x) = \frac{1}{K}$  if  $i \in J(x)$ . Kernel estimation lets precision vary and keeps bias constant. KNN does the opposite.

**Local Linear Regression** See Fan and Gijbels (1996). Kernel estimation has problems at the boundary of the space of  $x$  which LLR is able to remedy.

$$\widehat{m}(x_0) = \widehat{\alpha}, \text{ where } \widehat{\alpha} = \arg \min_{\alpha, \beta} \sum_{i=1}^n \{y_i - \alpha - \beta(x_i - x_0)\}^2 K\left(\frac{x_i - x_0}{a_n}\right)$$

Fan (1992, 1993) demonstrates advantages of LLR (lowess) over standard kernel estimators. Standard errors are bootstrapped.

The kernel  $K$  and  $a_n$  are chosen to optimize the asymptotic MSE.<sup>173</sup>

Kernels used in practice are:

- Epanechnikov:  $K(u) = \frac{3}{4}(1 - u^2)I\{|u| \leq 1\}$  (optimal  $K$  in both LLR and Kernel estimation, optimal  $a_n$  differ)
- Quartic:  $K(u) = \frac{15}{16}(1 - u^2)^2I\{|u| \leq 1\}$
- Triangle:  $K(u) = (1 - |u|)I\{|u| \leq 1\}$

The choice of  $a_n$  can be made using

- a point-wise plug-in method which relies on an initial estimate,
- a cross-validation method which chooses *global*  $a_n$  to minimize the MSE  $\sum_{i=1}^n (y_i - \widehat{m}_i(x_i))^2$ .
- a fishing expedition: increase  $a_n$  as long as linearity is not rejected.

**Remark 112.** STATA has a kernel smoother, kernel density estimation (*kdensity*), and does local linear regression. Advanced programs are available on the www for S-PLUS (<http://lib.stat.cmu.edu>).

**Remark 113.** Kernel estimation is basically a LLR on just the constant term.

<sup>173</sup>The bias in Kernel estimation depends on the distribution of regressors and on the slope of the regression function. The LLR bias only depends on the second derivative of the regression function. The asymptotic variance of the two methods is close unless data is sparse or  $m$  is changing rapidly around the  $x_0$  data point.

**Remark 114.** *LLR are used in the regression discontinuity design (see Section 9.1).*

**Remark 115.** *There are also extensions of the localization idea to the MLE framework, see Tibshirani and Hastie (1987).*

**Remark 116.** *There are other local regressions. For example, see the LOWESS procedure in S-PLUS. See Fan and Gijbels (1996) book (p.201) for local versions of quantile regressions.*

### 13.1. Multidimensional Extensions and Semiparametric Applications

The curse of dimensionality is severe. To have a reasonable speed of convergence we need very large samples. There are a few ways how to proceed:

- Regression trees: recursively split  $x$  to estimate step functions; derive a stopping rule to minimize mean square error.
- Impose additive separability or Projection pursuit regression:

$$m(x) = g_1(x\beta_1) + g_2(x\beta_2) + \dots$$

- Partial Linear Model: For a model  $y = z\beta + f(x) + \varepsilon$ , where both  $z$  and  $x$  are scalars, estimators of  $\beta$  can be constructed which are asymptotically normal at the  $\sqrt{n}$  speed of convergence. See Yatchew, A. (1998).
- Average derivative estimation: If I am interested in  $\theta = E \left\{ \frac{\partial m(x_i)}{\partial x_i} \right\}$  then  $\theta$  can be estimated with  $\sqrt{n}$  speed of convergence. Example: binary choice or Tobit models.
- Index sufficiency. See the semi-parametric Heckman's  $\lambda$  application by Powell (1989) in Section 19.2.4 or matching on propensity score in Section 20.2.
- See Athey and Imbens (Econometrica, 2006)<sup>174</sup> for a generalization of the difference-in-differences method (see Section 8.1). They relax linearity and allow treatment effects to vary across individuals and average treatment effects to vary across groups (such as states that do and don't adopt some policy).

---

<sup>174</sup><http://kuznets.fas.harvard.edu/%7Eathey/CIC.pdf>

## Part III

# Qualitative and Limited Dependent Variables

When using linear models, we spent much of our time dealing with endogeneity issues by applying instrumental variable approaches. But many economic decisions are qualitative or involve corner solutions making linear models inappropriate. Unfortunately, it is very difficult to deal with endogeneity in limited dependent variables models (esp. in dynamic panel data ones) unless strong assumptions are made on the exact relationship between the endogenous regressors and the instruments—something we do not want to do generally.

## 14. Qualitative response models

Our usual regression methods are designed for a continuous dependent variable. In practice, we very often analyze a qualitative response - a discrete dependent variable. For example: decide to buy a car, quit a job, retire, move, work; choose among many alternatives such as how to commute to work; choose sequentially the level of education; influence the number of injuries in a plant, etc. While it was entirely plausible to assume that  $\varepsilon$  in our usual regression model with a continuous  $y$  had a continuous pdf, this assumption is not valid here. The usual  $E[y | x]$  no longer does the job in those situations.

Most traditional models are estimated by MLE which allows us to write down even very complicated models.<sup>175</sup> As a consequence, IV is not easily possible and panel data analysis is difficult. Further, heteroskedasticity or omission of an explanatory variable orthogonal to the included regressors cause bias unlike in the linear regression analysis!<sup>176</sup> Since MLE crucially hinges on distributional assumptions, recent literature focuses on estimation methods not requiring specification of any parametric distribution.

---

<sup>175</sup>Also testing using the LR principle is very convenient.

<sup>176</sup>For example, Arabmazar and Schmidt (1981), give some examples of the asymptotic biases for the Tobit model, see section 18.1.



## 14.1. Binary Choice Models

### 14.1.1. Linear Probability Model

In the Linear Probability Model we assume our usual linear regression even though  $y_i \in \{0, 1\}$ . As a consequence the interpretation of  $E[y_i | x_i] = \beta'x_i$  being the probability the event occurs breaks down when  $\widehat{\beta}'x_i \notin [0, 1]$ .

**Exercise 14.1.** Show that given  $E[\varepsilon_i] = 0$ , the residuals  $\varepsilon_i$  which can take on only two values are heteroscedastic.

The main advantage of the LPM is its ability to handle IV estimation easily (2SLS, see Remark 73).<sup>177</sup> Applying the LPM is going to be close to ok asymptotically when the empirical  $\widehat{y}_i$ s are not close to 0 or 1 and in most applications, it is the preferred econometric model.<sup>178</sup> One should also allow the  $x$ s to enter as finite order polynomials, allowing for non-linearity, which will help with predicting out of admissible range.

**Example 14.1.** Cutler and Gruber (1995) estimate the crowding out effect of public insurance in a large sample of individuals. They specify a LPM:

$$\text{Coverage}_i = \beta_1 \text{Elig}_i + X_i \beta_2 + \varepsilon_i$$

*Eligibility is potentially endogenous and also subject to measurement error. To instrument for  $\text{Elig}_i$  they select a national random sample and assign that sample to each state in each year to impute an average state level eligibility. This measure*

<sup>177</sup>There are 2SLS procedures available for non-linear models (Achen 1986), but these require strong distributional assumptions (Angrist and Krueger, 2001; [W]). A method sometimes used here is a Probit with a *control function* approach (see note n. 119) based on joint normality of residuals in the first and second stage. See Rivers and Vuong (1988, JEconometrics) for Probit, Smith and Blundell (1986, Econometrica) for Tobit, and Blundell and Powell (2004, REStud) for semiparametric models with a control function approach to endogeneity. See also Chesher and Rosen (AERpp 2013) for a discussion of 2SLS with binary outcome and binary endogenous variable: they agree that the IV LATE interpretation may be useful here, but offer nonparametric IV as means of getting at ATE and other parameters.

<sup>178</sup>The LPM is unlikely to be consistent (Horrace and Oaxaca, 2006) and missclassification on the LHS is a particularly serious issue (Hausman, Abrevaya, and Scott-Morton, 1998), which Meyer and Mittag (2017) deal with successfully. In practice the arbitrary-assumed non-linearity-based model is also likely to be ‘wrong’, i.e., almost always inconsistent. See Angrist and Pischke (2009) for a vigorous defense of linear regressions in any situation, including qualitative choice, limited-dependent-variable models or IV with heterogenous treatment effects.

is not affected by state level demographic composition and serves as an IV since it is not correlated with individual demand for insurance or measurement error, but is correlated with individual eligibility.

#### 14.1.2. Logit and Probit MLE

The MLE methods transform the discrete dependent variable into a continuous domain using cumulative distribution functions. This is a natural choice as any  $F(\cdot) \in [0, 1]$ .

Assume existence of a continuous latent variable  $y_i^* = \beta' x_i + u_i$  where we only observe  $y_i = 1$  if  $y_i^* > 0$  and  $y_i = 0$  otherwise. Then for a symmetric  $F(\cdot)$  we have

$$P[y_i = 1 | x_i] = P[u_i > -\beta' x_i] = 1 - F(-\beta' x_i) = F(\beta' x_i). \quad (14.1)$$

Two common choices for  $F(\beta' x_i)$  are  $\Lambda(\beta' x_i) = \frac{\exp(\beta' x_i)}{1 + \exp(\beta' x_i)}$  (Logit) and  $\Phi(\beta' x_i)$  (Probit). The sample likelihood is then built under random sampling.<sup>179</sup>

**Remark 117.** MLE maximizes the log-likelihood  $\mathcal{L}(\theta) = \sum_{i=1}^N \log f(x_i, \theta)$ , where  $f(x_i, \theta)$  is the individual likelihood contribution, for computational convenience. It is a natural thing to do since

$$E \{ \mathcal{L}(\theta) - \mathcal{L}(\theta_0) \} \stackrel{iid}{=} n E \left\{ \log \left[ \frac{f(x_i, \theta)}{f(x_i, \theta_0)} \right] \right\} \stackrel{Jensen}{\leq} n \log \left[ E \left\{ \frac{f(x_i, \theta)}{f(x_i, \theta_0)} \right\} \right] = 0.$$

Therefore we construct a sample analogue to  $E \log f(x_i, \theta)$  and maximize w.r.t.  $\theta$ . Random sampling guarantees that  $\frac{1}{N} \sum_{i=1}^N \log l(x_i, \theta)$  converges to  $E \log f(x_i, \theta)$ . Hence, lack of independence will not be a problem if the marginals do not shift around, even though  $\mathcal{L}(\theta)$  is no longer the right likelihood. Similar convergence property underlies the GMM.

**Remark 118.** Both models are suitable for non-linear optimization using the Newton-Raphson methods as the Hessian is always n.d.

---

<sup>179</sup>When one rejects equality of coefficients, one cannot just interact all right-hand side variables with group dummies in the case of limited dependent variable models (as one would in the case of OLS) as this would lead to incorrect inference in presence of group-level unobservables. Williams (2009), "Using heterogeneous choice models to compare logit and probit coefficients across groups", *Sociological Methods and Research*, 37:4, 531-59.

**Remark 119.**  $\hat{\beta}$ s from Logit and Probit are not directly comparable ( $\hat{\beta}_{\text{Logit}} \simeq 1.6\hat{\beta}_{\text{Probit}}$ , see [M p.23]). More importantly, while  $\hat{\beta}_{\text{OLS}} = \frac{\partial E[y_i|x_i]}{\partial x_i}$  we need to find the probability derivatives for logits and probits, e.g.  $\hat{\beta}_{\text{Logit}} \neq \frac{\partial P[y_i=1|x_i]}{\partial x_i} = \Lambda(-\beta'x_i)[1 - \Lambda(-\beta'x_i)]\hat{\beta}_{\text{Logit}}$ . These derivatives will only most rarely be any different from the LPM coefficients, even when the mean probability is different from 0.5. (And this is probably true for Tobit marginal effects in many applications.) Note that in effect, there is treatment effect heterogeneity built into these non-linear models with the shape of the heterogeneity corresponding to distributional assumptions.<sup>180</sup>

**Remark 120.** Parametric methods (e.g., Probit and Logit) assume strict monotonicity and homoscedasticity.<sup>181</sup>

**Remark 121.** There are bivariate extensions in the SURE spirit ([G] 21.6).<sup>182</sup>

**Exercise 14.2.** Show that in Probit, one can only estimate  $\beta/\sigma$ .

**Exercise 14.3.** Estimates from binary response models are essentially WLS estimates: Find the corresponding GMM/IV interpretation for logit model using the FOC's of the MLE. Compare it to the corresponding probit expression and find the WNLLS interpretation for probit. Will they give the same answer as the MLE in small samples? Think of the intuition behind the size of the weight in the variance-covariance matrix as a function of  $x'\beta$ .

**Remark 122.** See Davidson and MacKinnon (1993) textbook, chapter 15.4 for a useful auxiliary regression connected to qualitative response models.

**Remark 123.** When you estimate 2SLS with an endogenous binary variable, it may be tempting to use Logit or Probit in the first stage and to plug in the first-stage predicted values. But this will only generate consistent estimates if the first-stage assumptions are exactly valid, which is unlikely. So just use LPM in the first stage. Only OLS makes absolutely sure that first-stage residuals are orthogonal to both fitted value and covariates (as the first-stage prediction error will, of course, show up in the residual of the main equation). See remark 73.

<sup>180</sup>For discrete  $X$ , one simply averages predicted probabilities for  $X = 0$  and for  $X = 1$  for all scenarios. See <https://www3.nd.edu/~rwilliam/stats3/Margins01.pdf>

<sup>181</sup>There is a heterogeneity test for probit ([G]p.680).

<sup>182</sup>Also see section 5.2. of the GMM handout by George Jakubson in the library for an example with correlated probits and their univariate approximation.

### 14.1.3. The WLS-MD for Multiple Observations

([M] 2.8, [G] 21.4.6) Suppose we have  $n_i$  observations corresponding to  $x_i$  and that for  $m_i$  of them the event occurred. Then assume  $\hat{p}_i = \frac{m_i}{n_i} = p_i + u_i = \beta' x_i + u_i$  and correct for heteroskedasticity. For non-linear models we invert the *cdf* and we need a Taylor series expansion to find the form of heteroskedasticity.

**Example 14.2.** For the logit model  $p_i = \Lambda(\beta' x_i)$  and we get  $\Lambda^{-1}(p_i) = \ln \frac{p_i}{1-p_i} = \beta' x_i + u_i$ .  $\square$

**Exercise 14.4.** Show the WLS is a genuine MD.

See Papke and Wooldridge (1996, 2008) for a fractional response logit quasi likelihood (glm in Stata).

### 14.1.4. Panel Data Applications of Binary Choice Models

The usual suspects: Random and Fixed Effects. See [H] 7. We cover both below, but it is not clear why one should not use the linear probability model for the usual reasons. The maximum likelihood estimator for a panel binary response model with fixed effects can be severely biased if  $N$  is large and  $T$  is small as a consequence of the incidental parameters problem. This has led to the development of conditional maximum likelihood estimators.<sup>183</sup>

**Remark 124.** To build dynamics into these models, one either includes the lagged outcome as an explanatory variable (see Section 11) or one conditions on the duration in the current state as a covariate (see Section 15).

**Random Effect Probit** Probit does not allow the fixed effect treatment at all. Random effects model is feasible but has been difficult because of multidimensional integration. To prevent contamination of  $\beta$ 's, we need to integrate the random effects  $\alpha$  out. For MLE we must assume a particular distribution for  $\alpha$ , say  $g(\alpha)$  depending on parameters  $\delta$ . Then allowing for correlation of  $\alpha$  over time for the same person we can maximize the following with respect to both  $\beta$  and  $\delta$  :

$$\mathcal{L}(\theta) = \sum_{i=1}^N \log \int \prod_{t=1}^T \Phi(\beta' x_{it} + \alpha)^{y_{it}} [1 - \Phi(\beta' x_{it} + \alpha)]^{1-y_{it}} dG(\alpha|\delta) \quad (14.2)$$

<sup>183</sup>See, e.g., IZA DP No. 11182, for one example of a recent update on this work.

Notice the multidimensional integral (each  $\Phi$  is an integral inside the  $T$ -dimensional integral over  $\alpha$ ). We can simplify matters by assuming that the correlation of  $\alpha$  between any two time periods is the same. Then we can look at each  $y_{it}$  and  $\int P[y_{it} = 1 \mid x_{it}, \alpha_i] g(\alpha) d\alpha = P[y_{it} = 1 \mid x_{it}]$ . For each  $y_{it}$  we then have a double integral.

**Remark 125.** When we allow for general structure of  $g(\alpha)$  we need the simulated method of moments (SMM) to evaluate the integrals fast (McFadden, 1988): When computing the  $P[y_i = 1 \mid X_i] = P_i$  presents a formidable computational problem one solution is to use their unbiased estimates. To illustrate this method return back to a cross-section and consider the GMM interpretation of probit where  $P_i = \Phi(\beta' x_i)$  (see exercise 14.3):

$$0 = \sum_{i=1}^N (y_i - P_i) \frac{X_i \phi(x_i' \beta)}{P_i(1 - P_i)} = \sum_{i=1}^N (\varepsilon_i) w_i.$$

Suppose  $P_i$  is hard to calculate. Solution? Use  $w_i = X_i$ , which will deliver inefficient but consistent estimates. You still need to evaluate the  $P_i$  inside the  $\varepsilon_i$ . To do this, let  $I(\cdot)$  be the indicator function and consider

$$\Phi(\beta' x_i) = \int_{-\infty}^{\beta' x_i} \phi(s) ds = \int_{-\infty}^{\infty} \phi(s) I(s < \beta' x_i) ds = E_s[I(s < \beta' x_i)].$$

To simulate the integral generate  $R$  values  $s_r \sim N(0, 1)$  and evaluate  $\frac{1}{R} \sum_{r=1}^R I(s_r < \beta' x_i)$  to obtain an unbiased estimate of  $P_i$ . (It's not consistent as long as  $R$  is finite so can't use it in the  $w_i$ ). To conclude, drive the simulated moment condition to 0.

**Remark 126.** To allow (flexible) correlation between  $x_i$  and  $\alpha_i$  we may follow Chamberlain (1980), but we now need the true regression function (see section 8.3.2) and a distributional assumption on the  $\alpha$  equation error term.

**Remark 127.** There is a specific counterpart to random effects usable with the logit model: NP-MLE (Non-Parametric Maximum Likelihood, see Heckman and Singer, 1984, for such duration models). Simply approximate the  $g(\alpha)$  with a discrete distribution. Estimate the points of support and the respective probabilities as part of your likelihood maximization. See Section 15.2.1 for an application of this approach.

**Conditional Fixed Effect Logit** The motivation for a fixed effect model is similar as in panel data linear regression. In MLE the  $\alpha_i$ s are again consistent only with  $T \rightarrow \infty$ . Since  $T$  is usually fixed and since MLE relies on consistency, the  $\alpha_i$ s must be swept out. But how do you “difference out” an additive element from a non-linear function?

Logit does allow for such a trick. Consider the  $T$  observations on  $y_{it}$  as one  $T$ -variate observation  $y_i$ . The suggestion of Chamberlain (1980) is to maximize the conditional likelihood (see section 16.0.5) of  $y_i$  given  $\sum_{t=1}^T y_{it}$  which turns out to remove the heterogeneity. Conditional on  $\alpha_i$ s we have independence over both  $i$  and  $t$ .

**Exercise 14.5.** To verify this, write down the conditional likelihood contribution of  $y_i' = (0, 1)$  when  $T = 2$ .

**Remark 128.** Again, use Hausman test to compare the fixed effect model with the  $\alpha_i = \alpha$  simple pooled-data logit.

**Remark 129.** The conditional fixed effect logit is computationally cumbersome for  $T \geq 10$ .

**Exercise 14.6.** Explain why  $y_i$ s with no change over time are not used in the estimation and show that observations with time constant  $x$ s are not used either.

**Remark 130.** In dynamic discrete choice panel data models, the dynamics is usually handled by including the lagged outcome as an explanatory variable. Frederiksen, Honoré and Hu (2007) propose an alternative model in which the dynamics is handled by using the duration in the current state as a covariate. They allow for group-specific effect (such as a firm fixed effect) in parametric and semiparametric versions of the model. A Stata program is available (*fhh*).

**Remark 131.** It is not trivial to calculate marginal effects after FE Logit since the fixed effects are generally not consistently estimated.<sup>184</sup>

<sup>184</sup>[http://repec.org/usug2016/santos\\_uksug16.pdf](http://repec.org/usug2016/santos_uksug16.pdf)

Yoshitsugu Kitazawa, 2011. “Average elasticity in the framework of the fixed effects logit model,” Discussion Papers 49, Kyushu Sangyo University, Faculty of Economics.

### 14.1.5. Relaxing the distributional assumptions

Parametric models of choice (like logit or probit) are inconsistent if the distribution of the error term is misspecified, including the presence of heteroskedasticity.

One can go fully non-parametric. Matzkin (1992): Let  $E[y_i | x_i] = m(x) = F(h(x))$  and study the identification of  $h$  from  $F$ . This is the most general and least operational we can go.

**Index models** Cosslett (1981):  $\max \mathcal{L}(\theta)$  w.r.t. both  $\beta$  and  $F(g(\beta, x_i))$ , where  $g(\cdot)$  is assumed parametric. Only consistency derived, but no asymptotic distribution. Further research on index models includes Ichimura (1993) with a  $\sqrt{n}$  estimator and Klein and Spady (1993, in Stata as `sml`). All of these require  $\varepsilon$  and  $x$  to be independent. See also the average derivative estimator of Powell, Stock, and Stoker (Ecm, 1989, not implemented in standard packages).

**Maximum rank estimators** Manski's Maximum Score Estimator (1975, 1985) maximizes the number of correct predictions, is  $n^{-1/3}$  consistent, and is in LIMDEP. The idea is based on  $E[y_i | x_i] = F(\beta_0' x_i)$ . Assume  $F(s) = .5$  iff  $s = 0$ .<sup>185</sup> Then  $\beta_0' x_i \geq (\leq) 0$  iff  $E[y_i | x_i] \geq (\leq) .5$  and we use  $\text{sgn}(\beta' x_i) - \text{sgn}(E[y_i | x_i] - .5) = 0$  as a moment condition.<sup>186</sup> Then

$$\hat{\beta}_{MRE} = \arg \max_{s.t. \beta' \beta = 1} \frac{1}{n} \sum_{i=1}^n \left[ (2y_i - 1) \text{sgn}(\beta' x_i) \right] \quad (14.3)$$

Functionally related regressors are excluded by identification assumptions and  $\hat{\beta}_{MRE}$  is identified up to a scaling factor. Asymptotic distribution is not normal and not easy to use since variance is not the right measure of variation so we need to bootstrap. The method allows for conditional heteroskedasticity and generalizes to multinomial setting.

Smoothed MSE by Horowitz (1992) can be made arbitrarily close to  $\sqrt{n}$  convergence. The idea is to smooth the score function to make it continuous and differentiable by using *cdf* in place of *sgn*.

Another method of maximizing correct predictions is based on the Powell's idea of comparing pairs of people.<sup>187</sup> Assume the model  $y_i = d_i = 1\{x_i\beta + \epsilon_i > 0\}$  and assume  $\epsilon_i$  independent of  $x_i$  (no heteroskedasticity), then  $E[d_i - d_j | x_i, x_j] =$

<sup>185</sup>The choice of the median can be generalized to any quantile.

<sup>186</sup>Note that  $\text{sgn}(\cdot)$  is not invertible.

<sup>187</sup>See the discussion on selection the Powell's way below in Section 19.2.4

$E[d_i|x_i] - E[d_j|x_j] = F_\epsilon(x_i\beta) - F_\epsilon(x_j\beta) > 0$  iff  $(x_i - x_j)\beta > 0$  so estimate  $\beta$  by maximum rank estimator such as

$$\max_{\beta} \sum_{i < j} \text{sign}(d_i - d_j) \text{sign}((x_i - x_j)\beta) \quad (14.4)$$

This, of course gets rid of the intercept, so Heckman (1990) proposed that in presence of exclusion restriction, one can get the intercept off those who have  $p(d_i = 1)$  almost equal to one.

Finally, Sognian Chen (1999) uses the additional assumption of symmetry of the distribution of  $\epsilon$  to allow for  $\sqrt{n}$  estimation of the constant term. (All other semiparametric methods make for a slower rate for the constant even if they deliver  $\sqrt{n}$  for the slope.) (You still need to normalize the scale.) He also allows for heteroskedasticity of a particular form:  $f(\epsilon|x) = f(\epsilon|\tilde{x})$  where  $\tilde{x}$  is a subset of  $x$ . Assuming  $f$  is symmetric implies that  $E[d_i + d_j|x_i, x_j] = F_\epsilon(x_i\beta) + F_\epsilon(x_j\beta) > 1$  iff  $(x_i + x_j)\beta > 0$  (draw a picture of  $f$  symmetric around 0 to see why). Note that sum does not get rid of the intercept. So, estimate something like

$$\max_{\beta} \sum_{i < j} \text{sign}(d_i - d_j - 1) \text{sign}((x_i + x_j)\beta). \quad (14.5)$$

**Remark 132.** Two semiparametric approaches are based on approximating the unknown distributions in a flexible fashion. The *Stata* commands `snp` and `sml` implement the polynomial-approximation pseudo-likelihood approach of Gallant and Nychka (1987) and the Klein and Spady (1993) kernel-based (local-approximation) likelihood approach, respectively. These are  $\sqrt{n}$  consistent and asymptotically normal estimators that allow one to estimate univariate, bivariate and sample-selection binary-outcome models semiparametrically (see De Luca, 2008 *Stata Journal*). The `sml` routines are robust to heteroskedasticity. See also the `dfbr` routine, which implements NLLS-based semiparametric binary-outcome estimators under a conditional median restriction that allow for general forms of heteroskedasticity.

**Remark 133.** Differencing a fixed effect out of a non-linear function is difficult, which affects discrete-choice and limited-dependent-variables models. Assuming index structure and additivity, Gayle (2013) or Chernozhukov et al. (2013) estimate marginal effects in nonlinear panel models. Hoderlein and White (2012) derive a generalized version of differencing that identifies local average responses in a general nonseparable model without the single-index structure. They identify



effects for the subpopulation of individuals who have not experienced a change in covariates between the two time periods. Čížek and Lei (2018, JEcm) build on Hoderlein and White (2012) and propose an identification strategy where the distribution of individual effects depends on the explanatory variables only by means of their time averages. Their approach allows for lagged dependent and discrete explanatory variables.

## 15. Duration Analysis

[W] 20. Here we continue in reduced-form distribution-based maximum likelihood modelling, which is designed to fit the processes that result in variation in duration (length).<sup>188</sup>

**Example 15.1.** *Length of a job, duration of a marriage, how long a business lasts, when a worker retires, duration of a strike, length of an unemployment spell, length of a stay in a hospital depending on the type of insurance, spacing of births, time spent off drugs while fighting addiction, etc.*

The advantage of duration models is in their ability to handle time changing  $x$ s (both with respect to calendar and duration time), duration dependence, and right censoring. The models can also handle multiple exits and multiple states. Read Kiefer (1988), [G]22.5, [L].

### 15.1. Hazard Function

Duration models build upon the concept of a hazard function  $\lambda(t)$ , which is defined as the probability of leaving a given state at duration  $t$  *conditional* upon staying there up to that point. Using this definition one can build a likelihood function for the observed durations and estimate it using standard methods (MLE, GMM). For example, if the hazard does not depend on either  $x$ s or duration  $t$ , then we can express the unconditional probability of observing a spell of duration  $t$ , denoted  $f(t)$  as  $f(t) = \lambda(1 - \lambda)^{t-1}$ . The probability that a spell lasts at least  $T$  periods is called survival  $S(T) = \Pr[t \geq T] = 1 - F(t) = (1 - \lambda)^{T-1}$ . This type of spell, where we do not observe the end of the spell, is called *right censored*. A *left censored* spell occurs when we do not observe the first part of the spell, but do observe

---

<sup>188</sup>In this regard, the approach is similar to how we build a model for count data from the Poisson distribution in Section 17.

when it ended. What makes a tremendous difference is whether we know when a left censored spell started or not. Of course  $\lambda(t) = \frac{f(t)}{S(T)}$ .

**Exercise 15.1.** Suppose the hazard depends on  $t$  and write down the likelihood contribution for a completed spell and for a right censored spell. Next assume that there is no duration dependence and write down the likelihood contribution of a left censored spell. Finally, how would your last answer differ in presence of duration dependence, depending on whether you know when a left censored spell started.

**Remark 134.** A first approximation to the hazard, ignoring both observed and unobserved differences is the so called Kaplan-Meier statistic (also called empirical hazard).<sup>189</sup>

$$\lambda(t) = \frac{\#[exit(t)]}{\#[risk(t)]} \text{ with } \sigma_\lambda(t) = \sqrt{\frac{\lambda(t)(1 - \lambda(t))}{\#[risk(t)]}}. \quad (15.1)$$

**Exercise 15.2.** Verify the formula for  $\sigma_\lambda$ . Also, think of how you would estimate the empirical hazard in a case of competing risks, i.e., when there are two or more ways how to leave a given state.

One can use either **discrete time or continuous time** hazard models. In a discrete time model, the transition can occur at most once in a given time period, i.e., these models depend on the unit of the time interval. In a continuous time model

$$\lambda(t) = \lim_{h \rightarrow 0} \frac{1}{h} \Pr(t \leq t^* < t + h \mid t^* \geq t) \quad (15.2)$$

A widely used continuous time model is the *proportional hazard* model (relative risk, Cox 1972),  $\lambda_i(t) = \exp(h(t)) \exp(x'_i \beta) = \lambda_0(t) \exp(x'_i \beta)$ , where  $\lambda_0(t)$  is the so called baseline hazard. Here,  $x' \beta$  shifts the baseline hazard (the shape of the duration dependence) up or down by a constant.<sup>190</sup>

**Remark 135.** Note that in continuous time, the hazard equals

$$-\frac{d \ln S(t)}{dt} = -\frac{d \ln[1 - F(t)]}{dt} = \frac{f(t)}{1 - F(t)} = \lambda(t),$$

<sup>189</sup>See work by Jiun-Hua Su on counterfactual Kaplan Meier evaluation.

<sup>190</sup>A major alternative is the accelerated failure model (log linear in  $t$ )  $\lambda(t) = \lambda_0(t\theta)\theta$ , where  $\theta = \exp(x'_i \beta)$  so that the effect of  $x' \beta$  is to accelerate or decelerate the degree of duration dependence over the course of the spell. Consequently,  $f(t) = \theta f_o(\theta t)$  and  $\log(T) = -\log(\theta) + \log(T\theta)$ .

which implies that

$$S(t) = \exp - \int_0^t \lambda(\tau) d\tau, \text{ and } f(t) = \lambda(t) \exp - \int_0^t \lambda(\tau) d\tau.$$

**Example 15.2.** One possible choice of a discrete time hazard is the logit specification:

$$\lambda_j(t, x_t | \theta_k^j) = \frac{1}{1 + e^{-h_j(t, x_t | \theta_k^j)}}$$

where  $h_j(t, x_t | \theta_k^j) = \beta_j' x_t + g_j(t, \gamma_j) + \theta_k^j$ . Here,  $g_j(t, \gamma_j)$  is a function capturing the duration dependence.<sup>191</sup>

**Exercise 15.3.** Can the logit model be interpreted as an approximation to a proportional hazard model?

**Remark 136.** One can trick LIMDEP or other software to estimate the logit duration model.

**Remark 137.** Kalbfleish and Prentice (p. 47) argue that the complementary log-log model (`cloglog` in Stata) provides a “uniquely appropriate” discrete-time approximation of the continuous Cox model:  $\Pr(t^* = t \mid t^* \geq t - 1) = 1 - (1 - \lambda(t))^{\exp(x_i' \beta)}$ , where  $\lambda(t) = \exp[\int_{t-1}^t \lambda_0(u) du]$ .

## 15.2. Estimation Issues

First, there is a possibility of the so called *length-biased (stock) sampling*: correct sampling is from inflow during a certain time window (sampling frame). Sampling from stock oversamples long spells (somebody starting a quarter ago with a short spell will not be in today’s stock).

Second, *left censored spells* with an unknown date of start create a difficult estimation problem (see Exercise 15.1 and below).<sup>192</sup>

<sup>191</sup>For proportional hazard models, Elbers and Ridder show that the heterogeneity distribution and the duration dependence can be separately identified.

<sup>192</sup>We can fix things if we know start of spell unless there are unobservables, which would lead to dynamic distortion of the distribution of unobservables by selection on who of the left-censored makes it into the sample.

Third, it is well known that in the presence of *unobserved person-specific characteristics* affecting the probability of exit, all of the estimated coefficients will be biased.<sup>193</sup>

One of the widely used methods of controlling for unobserved factors is the flexible semi-parametric heterogeneity MLE estimator proposed by Heckman and Singer (1984) (also called NP-MLE as in non-parametric MLE). They show that if there is a parametric continuous distribution of unobservables, the estimated distribution has to be that of a discrete mixing distribution with a step function nature. (Think of random effect probit.) Using simulations, a small number of points of support has been shown to remove the bias in  $\beta$ s. There is no known way of correctly constructing the asymptotic standard errors, since the dimension of the parameter space depends on the sample size. So assume the number of points of support is fixed to invoke standard asymptotics and determine that actual number of points of support of the distribution of unobservables from the sample likelihood.<sup>194</sup>

**Remark 138.** *The heterogeneity bias in duration dependence coefficients has been shown to be negative. To see why, think of two flat hazards  $\lambda_{M/S}(t)$  of married and single women and construct the empirical hazard in absence of the marital status info.*

**Remark 139.** *Note that if there is no variation in the  $x$ s independent of duration, identification will be difficult.*

**Remark 140.** *Competing risk models are generally unidentified in the sense that for every dependent distribution of time to exit by cause, one can find an inde-*

---

<sup>193</sup>Here, we are concerned with the effects of unobserved heterogeneity in duration models. For an example of similar methods in other settings, see Berry, Carnall and Spiller (1995), where they explicitly allow for two types of airline customers (businessmen vs. tourists), unobserved by the econometrician.

<sup>194</sup>A simulation study by Xianghong Li and Barry Smith (2015, *Journal of Econometrics*) provides guidance on the performance of different types of optimization algorithms with respect to finding the global optimum (depending on starting values) in single-spell models. They advocate the use of a step function for duration dependence and suggest the use of the simulated annealing (SA) algorithm in place of derivative-based optimization techniques. They also suggest a bootstrap procedure for choosing the number of support points and argue that likelihood ratio tests may still be appropriate for this purpose. Effectively, they contradict some of the well known results of Baker and Melino (2000) who suggest using Schwarz or Akaike criterion for picking the number of points of support and they also contradict the conclusion of Gaure et al. (2007) on the SA algorithm.

pendent one that is observationally equivalent. Typically, people assume independence of causes conditional on  $X$  (random censoring) or they assume a parametric model.

**Example 15.3.** *Honore and Lleras-Muney (Econometrica, 2006) look at the competing risk of dying of cardiovascular diseases (CVD) and cancer. Age-adjusted mortality rate from cancer has been flat for 30 years and some view this as evidence on little progress on cancer. However, if the causes of CDV and cancer are dependent, this interpretation is wrong. Honore and Lleras-Muney make no assumptions of the joint distribution of the underlying durations and estimate bounds on the marginal distributions of each of the competing duration variables.*

**Example 15.4.** *A timing-of-events approach (Abbring and van den Berg, 2003, Econometrica) based on the mixed proportional hazard (MPH) model is used when estimating the effect of a treatment (assigned without randomization) on exit from a given state, when we know the exact time of assignment to treatment, which occurs during an on-going spell in a given state. Identification is established of the causal relationship between the two durations (outcome duration, duration to treatment) in absence of IVs or CIA (based on conditionin on observables only), based on assuming (i) MPH structure, (ii) no anticipation, and (iii) conditional independence of the two hazards when conditioning on both observables and unobservables, which are estimated using the NP-MLE approach. The idea is to use variation (conditional on  $X$ s and the other assumptions) in the timing of treatment to generate comparisons in the next period between those who were already treated and those who were note (but could have been). The timing is the “IV” here. As usual, the method has to deal with dynamic sample selectino on unobservables over spell durations, and multiple-spell data make identification easier.<sup>195</sup>*

**Remark 141.** *As usual in a non-linear setting, IV typically requires specifying a fully parametric ‘first stage’, which is unattractive because (as opposed to the linear model), this first stage must be exactly true to avoid inconsistencies. An alternative is to treat durations as an outcome variable and set-up a non-parametric*

---

<sup>195</sup>In a Monte Carlo study Gaure, Røed and Zhang (2007) argue that separating causality (of some treatment and/or duration dependence) from selection (on unobservables) within non-experimental duration data by means of estimating the mixing discrete distributions of unobservables (i.e., by means of the non-parametric maximum likelihood estimator, NPMLE) is a piece of cake (and restrictions on the number of support points proposed by Baker and Melino (2000) in a single-spell framework may cause more harm than good).

*LATE analysis.* Froelich (2007, *JEcm*) covers this approach when covariates are needed. See also Abbring and Van den Berg (2004) and Aakvik, Helge, and Kjerstad (2012) who follow Heckman and Vytlacil (2005). van den Berg, Bonev and Mammen (2016) develop a nonparametric instrumental variable approach for the estimation of average treatment effects on hazard rates.

### 15.2.1. Flexible Heterogeneity Approach

Let us concentrate on a discrete time logit hazard model. We need to allow the likelihood to pick up the presence of unobservable person-specific heterogeneity. To use the “random effects” approach, estimate a discrete mixing distribution  $p(\theta)$  of an unobserved heterogeneity term  $\theta$  as a part of the optimization problem. In doing so, one can approximate any underlying distribution function of unobservables.

More specifically, let  $\lambda_j(t, x_t | \theta_k^j)$  be the conditional probability (hazard) of leaving a given state at time (duration)  $t$  for someone with person specific characteristics  $x_t$ , conditional upon this person having the unobserved factor  $\theta_k^j$ ,  $k = 1, 2, \dots, N_\theta^j$ . The  $j$  subscript stands for the different ways of leaving a given state and serves, therefore, as a state subscript as well. For example one can leave unemployment for a new job or for a recall, in which case  $j \in \{r, n\}$ , or one can leave employment through a quit or through a layoff, in which case  $j \in \{q, l\}$ . This is often referred to as a *competing risk model*. Below, we will use the example of quit, layoff, recall and new job. See also the discussion in [P]6.5.1.

To give an example of how the sample likelihood is evaluated using the concept of a hazard function, assume away any complications arising from the competing risks for now. Let  $\lambda$  denote the overall hazard out of a given state. In the absence of any unobserved heterogeneity, the likelihood function contribution of a single employment spell which ended at duration  $t$  would be

$$L_e(t) = \lambda(t, x_t) \prod_{v=1}^{t-1} [1 - \lambda(v, x_v)]. \quad (15.3)$$

In a competing risks specification with layoff and quit hazards (not allowing for unobserved factors), the unconditional probability of someone leaving employment through a quit at duration  $t$  would become

$$L_e^q(t) = \lambda_q(t, x_t) \prod_{v=1}^{t-1} [1 - \lambda_q(v, x_v)][1 - \lambda_l(v, x_v)], \quad (15.4)$$

where  $\lambda_q$  and  $\lambda_l$  denote the quit and layoff hazards respectively. Similarly, for someone who gets laid off in week  $t$  of an employment spell, the likelihood contribution becomes

$$L_e^l(t) = \lambda_l(t, x_t) \prod_{v=1}^{t-1} [1 - \lambda_q(v, x_v)][1 - \lambda_l(v, x_v)]. \quad (15.5)$$

Hazard models are natural candidates for dealing with the problem of right-censoring. For an employment spell which is still in progress at the end of our sampling frame (i.e., no transition out of employment has been observed), one enters the survival probability

$$S_e(T) = \prod_{v=1}^T [1 - \lambda_q(v, x_v)][1 - \lambda_l(v, x_v)]. \quad (15.6)$$

Here,  $T$  denotes the highest duration at which we observe the spell in progress and  $S_e(T)$  gives the probability of a given spell lasting at least  $T$  periods. The sample likelihood then equals the product of individual likelihood contributions. Now, if we introduce the unobserved heterogeneity, the likelihood function contribution for someone leaving unemployment at duration  $t$  for a new job would be

$$L_u^n(t) = \sum_{k=1}^{N_\theta^n} \sum_{m=1}^{N_\theta^r} p(\theta_k^n, \theta_m^r) L_u^n(t|\theta_k^n, \theta_m^r), \quad (15.7)$$

where  $p(\theta_k^n, \theta_m^r)$  is the probability of having the unobserved components  $\theta_k^n$  and  $\theta_m^r$  in the new job and recall hazards respectively, and where

$$L_u^n(t|\theta_k^n, \theta_m^r) = \lambda_n(t, x_t|\theta_k^n) \prod_{v=1}^{t-1} [1 - \lambda_n(v, x_v|\theta_k^n)] [1 - \lambda_r(v, x_v|\theta_m^r)]. \quad (15.8)$$

The likelihood of leaving an employment spell in week  $s$ , denoted  $L_e(s)$ , is specified in a similar fashion (with quit and layoff being the different reasons for exit here).

The previous discussion focuses on examples with a single spell of each type. Equation 15.9 gives the likelihood contribution of a person with two completed spells of employment. The first spell starts in week  $t + 1$  and ends with a layoff in week  $s$  (at duration  $s - t$ ); the second spell starts in week  $r + 1$  and ends with a quit in week  $w$  (at duration  $w - r - s - t$ ).

$$L(s, w) = \sum_{k=1}^{N_\theta^q} \sum_{m=1}^{N_\theta^l} p(\theta_k^q, \theta_m^l) L_e^l(s|\theta_k^q, \theta_m^l) L_e^q(w|\theta_k^q, \theta_m^l) \quad (15.9)$$

Here  $\theta^q$  and  $\theta^l$  denote the unobserved terms entering quit and layoff hazards respectively and

$$L_e^l(s|\theta_k^q, \theta_m^l) = \lambda_l(s, x_s|\theta_m^l) \prod_{v=t+1}^{s-1} [1 - \lambda_q(v, x_v|\theta_k^q)] [1 - \lambda_l(v, x_v|\theta_m^l)] , \quad (15.10)$$

$$L_e^q(w|\theta_k^q, \theta_m^l) = \lambda_q(w, x_w|\theta_m^l) \prod_{v=r+1}^{w-1} [1 - \lambda_q(v, x_v|\theta_k^q)] [1 - \lambda_l(v, x_v|\theta_m^l)] .$$

Using multiple spell data provides greater variation and improves identification of the unobserved heterogeneity distribution (need to separate duration dependence from unobserved heterogeneity). However, use of this type of data raises the possibility of *selection bias*; i.e., the individuals with more than one spell of either type may be a non-random sample. To control for this problem, one can estimate the whole duration history of all states jointly while allowing the unobserved heterogeneity to be correlated across these spells. To continue in the example we used up to now, the unemployment and employment hazard have to be estimated *jointly* in order to control for selection bias into multiple spells. One has to take into account the joint density of the unobservables across the two hazards, denoted by  $p(\theta^u, \theta^e)$ . Suppose we want to estimate a competing risks specification for quits and layoffs jointly with an overall hazard for unemployment. The likelihood contribution of someone leaving the first unemployment spell after  $t$  weeks, then getting laid off after  $s - t$  weeks on a job and staying in the second unemployment spell till the date of the interview, say at  $T - s - t$  weeks into the last spell, then becomes

$$L^{u,l,u}(t, s, T) = \sum_{k=1}^{N_\theta^u} \sum_{m=1}^{N_\theta^q} \sum_{n=1}^{N_\theta^l} p(\theta_k^u, \theta_m^q, \theta_n^l) L_u(t|\theta_k^u) L_e^l(s|\theta_m^q, \theta_n^l) S_u(T|\theta_k^u), \quad (15.11)$$

where

$$L_u(t|\theta_k^u) = \lambda_u(t, x_t|\theta_k^u) \prod_{v=1}^{t-1} [1 - \lambda_u(v, x_v|\theta_k^u)] .$$

The employment contribution,  $L_e^l$  is defined in equation 15.10 . Finally

$$S_u(T|\theta_k^u) = \prod_{v=s+1}^T [1 - \lambda_u(v, x_v|\theta_k^u)]$$



is the survivor function expressing the probability of a given spell lasting at least  $T$  periods.

One can compute individual contributions to the sample likelihood for other labor market histories in a similar way. The number of points of support of the distribution of unobservables ( $N_\theta^u$ ,  $N_\theta^q$  and  $N_\theta^l$ ) is determined from the sample likelihood.<sup>196</sup> Note the assumption of  $\theta^u$ ,  $\theta^q$  and  $\theta^l$  staying the same across multiple unemployment and employment spells respectively. There are many possible choices for the distribution of unobservables:

### Heterogeneity Distributions

1. Independent Heterogeneity:  $p(\theta^u, \theta^e) = p_u(\theta^u)p_e(\theta^e)$

2. Bivariate Heterogeneity Distribution:

	$\theta_1^l$	$\theta_2^l$	...	$\theta_N^l$
$\theta_1^q$	$p(\theta_1^q, \theta_1^l)$	$p(\theta_1^q, \theta_2^l)$	...	$p(\theta_1^q, \theta_N^l)$
$\theta_2^q$	$p(\theta_2^q, \theta_1^l)$	$p(\theta_2^q, \theta_2^l)$	...	$p(\theta_2^q, \theta_N^l)$
...	...	...	...	...
$\theta_M^q$	$p(\theta_M^q, \theta_1^l)$	$p(\theta_M^q, \theta_2^l)$	...	$p(\theta_M^q, \theta_N^l)$

3. One factor loading: pairs of  $\{\theta^l, \theta^q\}$  such that

$p(\Theta_1)$	$\Theta_1 = \{\theta_1^l, c\theta_1^l\}$
$p(\Theta_2)$	$\Theta_2 = \{\theta_2^l, c\theta_2^l\}$
...	...
$p(\Theta_N)$	$\Theta_N = \{\theta_N^l, c\theta_N^l\}$

4. Heterogeneity distribution with 3-tuples (corresponding to one way of leaving unemployment and 2 ways of leaving employment.)

$p(\Theta_1)$	$\Theta_1 = \{\theta_1^u, \theta_1^l, \theta_1^q\}$
$p(\Theta_2)$	$\Theta_2 = \{\theta_2^u, \theta_2^l, \theta_2^q\}$
...	...
$p(\Theta_N)$	$\Theta_N = \{\theta_N^u, \theta_N^l, \theta_N^q\}$

5. ‘Stayer’ heterogeneity: Suppose that we want to allow for the possibility of never leaving employment through a quit (or for the possibility of never

<sup>196</sup>Simulation provide important guidance. See Baker and Melino (2000) and more importantly Xianghong and Smith (2009).

returning to a prison.) Assume, for now, that the only way to transit out of employment is to quit. Furthermore, assume that there is no unobserved heterogeneity. A typical stayer model would then parametrize an individual's contribution to the likelihood as follows:

$$L(t) = p_s + (1 - p_s) \{ \lambda_q(t, x_t) \prod_{v=1}^{t-1} [1 - \lambda_q(v, x_v)] \},$$

where  $p_s$  is the probability of never leaving employment and  $\lambda_q$  is a quit hazard. See Jurajda (2002) for details on estimation.

6. Continuous parametric distributions of heterogeneity, for example Weibull.

### 15.2.2. Left Censored Spells

We need to know when they started. In presence of unobserved heterogeneity, dropping left censored spells will cause bias. See Ham and Lalonde (1997) for an example where the bias matters. Heckman and Singer (1984) suggest to model the interrupted spells with a separate hazard, i.e., a new hazard with a different  $\beta$  from the fresh spells. See also exercise 15.1.

### 15.2.3. Expected Duration Simulations

How to evaluate the magnitude of coefficients? Use the unconditional probability of leaving a given state to compute the expected durations under different values of  $x_s$ . Interpret the difference between expected durations as the magnitude of the particular  $\beta$ . The expected duration is computed as

$$E(t|X) = \sum_{i=1}^I \frac{\sum_{t=1}^{\infty} t f_i(t)}{I}, \quad (15.12)$$

where  $I$  is the number of spells in the sample,  $x_{it}$  is the vector of all explanatory variables for a spell  $i$  at duration  $t$ , and  $X$  represents the collection of all  $x_{it}$  vectors.<sup>197</sup> Finally, using the example of a *recall* and *new* job hazard out of unemployment, the unconditional probability of leaving unemployment at duration  $t$  denoted  $f_i(t)$  is computed as follows:

---

<sup>197</sup>A simpler (biased) approach is to evaluate the expected duration at a mean individual  $\bar{x}$ .

$$\begin{aligned}
f_i(t) &= \sum_{k=1}^N p(\theta_k^r, \theta_k^n) f_i(t|\theta_k^r, \theta_k^n), \text{ where} \\
f_i(t|\theta_k^r, \theta_k^n) &= \{\lambda_r(t, x_{it}|\theta_k^r) + \lambda_n(t, x_{it}|\theta_k^n) - \lambda_r(t, x_{it}|\theta_k^r)\lambda_n(t, x_{it}|\theta_k^n)\} \times \\
&\quad \prod_{v=1}^{t-1} [1 - \lambda_r(v, x_v|\theta_k^r)] [1 - \lambda_n(v, x_v|\theta_k^n)].
\end{aligned}$$

**Remark 142.** *In multiple-state multiple-spell data, single-spell duration simulations do not provide a full picture. See, e.g., Jurajda (2002). Garcia-Suaza (2015) provides Oaxaca-Blinder decompositions for duration models. So does IZA DP No. 9909.*

#### 15.2.4. Partial Likelihood

Cox (1972, 1975): estimate  $\beta$  in the proportional hazard model  $\lambda_i(t) = \lambda_0(t) \exp(x_i'\beta)$ , without specifying the form of the baseline hazard  $\lambda_0(t)$ . Order the completed durations  $t_i$  into  $t_{(i)}$ . The conditional probability that individual 1 concludes a spell at time  $t_{(1)}$ , given that all other individuals could have completed their spells at that duration is

$$\frac{\lambda(t_{(1)}, x_{(1)})}{\sum_{i=1}^n \lambda(t_{(1)}, x_{(i)})} = \frac{\exp(x'_{(1)}\beta)}{\sum_{i=1}^n \exp(x'_{(i)}\beta)}. \quad (15.13)$$

In the absence of information about the form of duration dependence, only the information about the *order* of the spell durations is used.

**Remark 143.** *This method alone does not allow the expected duration simulations. It is possible, though, to construct a nonparametric estimate of the baseline hazard using the estimated  $\exp(x_i'\hat{\beta})$ . See [P].*

## 16. Multinomial Choice Models

See McFadden (1984) for a summary of pioneering research and Pudney (1989), *Modelling Individual Choice*, [P], chapter 3, for discussion of the material in view of the underlying economic theory. ([M]2,3, [A]9, [G]21)<sup>198</sup>

<sup>198</sup>The LIMDEP v.7 (1995) manual discusses several extensions and examples of application.

### 16.0.5. Unordered Response Models

So far we talked about 0/1 decisions. What if there are more choices?

**Example 16.1.** *Choice of commuting to work, choice among occupations, purchasing one of many product brands, etc.*

We want to analyze *simultaneous* choice among  $m$  alternatives. The idea is to look at pairwise comparisons to some reference outcome:

$$\frac{p_j}{p_j + p_m} = F(\beta'_j x) \Rightarrow \frac{p_j}{p_m} = \frac{F(\beta'_j x)}{1 - F(\beta'_j x)} = G(\beta'_j x) \Rightarrow p_j = \frac{G(\beta'_j x)}{1 + \sum_{k=1}^{m-1} G(\beta'_k x)} \quad (16.1)$$

**Remark 144.** *Note that in our binary choice example ( $m = 2$ ), we also started by defining  $p_j(p_j + p_m)^{-1} = p/(p + 1 - p) = p = F(\beta'_j x)$ .*

**Multinomial Logit (MNL)** If  $F(\beta'_j x) = \Lambda(\beta'_j x)$  then  $G(\beta'_j x) = \exp(\beta'_j x)$  and the estimation does not require any integration. Simply define  $y_{ij} = 1$  if person  $i$  chooses the  $j$ -th choice and,  $y_{ij} = 0$  otherwise, and

$$\max_{\beta_1, \beta_2, \dots, \beta_{m-1}} \log L = \sum_{i=1}^N \sum_{j=1}^m y_{ij} \log p_{ij}, \text{ where} \quad (16.2)$$

$$p_{ij} = \frac{\exp(\beta'_j x_i)}{1 + \sum_{l=1}^{m-1} \exp(\beta'_l x_i)} \text{ for } j = 1, \dots, m-1 \text{ and } p_{im} = \frac{1}{1 + \sum_{l=1}^{m-1} \exp(\beta'_l x_i)}.$$

**Remark 145.** *The FOCs again have the familiar GMM interpretation*

$$\sum_{i=1}^N (y_{ij} - \hat{p}_{ij}) x_i = 0 \text{ for } j = 1, \dots, m-1$$

and again imply that if  $x$  consists only of a constant, the model predicts the actual frequencies (see exercise 14.3). This can be used to define a measure of fit based on comparing our log-likelihood with the benchmark that one would obtain by merely regressing the outcome on constants  $\alpha_j$ .

**Exercise 16.1.** Verify that the benchmark likelihood equals  $\prod_{j=1}^m \left(\frac{N_j}{N}\right)^{N_j}$  where  $N_j = \sum_{i=1}^N y_{ij}$ .

**Exercise 16.2.** What happens in the commuting choice example when all males choose to drive?

**Remark 146.** To interpret the estimates  $\hat{\beta}$  we need the derivatives w.r.t.  $x_k$  ( $k$ -th element of  $x$ ) even more as  $\hat{\beta}_{jk}$  shows up in  $p_l$   $l = 1, \dots, m$ .  $\frac{\partial p_j}{\partial x_k} = p_j[\beta_{jk} - \sum_{s=1}^{m-1} p_s \beta_{sk}]$ .

**Remark 147.** There is a utility maximization model of individual choice which leads to the multinomial logit, assuming additivity of disturbances  $y_{ij}^* = V_j(x_i) + \epsilon_{ij}$  with  $y_{ik} = 1$  iff  $\forall j \neq k$   $y_{ik}^* > y_{ij}^*$  and assuming  $\epsilon_j$ 's are iid type I extreme value distribution.<sup>199</sup>  $P[y_k = 1|x_i]$  corresponds to the joint occurrence of  $V_k(x_i) + \epsilon_k > V_j(x_i) + \epsilon_j$   $\forall j \neq k$ , that is

$$P[y_k = 1|x_i] = \int \prod_{j \neq k} F(\epsilon_k + V_k(x_i) - V_j(x_i)) f(\epsilon_k) d\epsilon_k.$$

In class we show that this equals  $\exp(V_k(x_i)) / \sum_{j=1}^m \exp(V_j(x_i))$ . Set  $V_k(x_i) = x_i' \beta_k$  to conclude.

**Remark 148.** Matejka and McKay (AER) provide a more general underpinning for the Logit model—one where individuals make choices with imperfect information about payoffs and, before choosing, they can study the payoffs, which is costly.

**Exercise 16.3.** Verify the derivatives in [M] p.36 and show  $\frac{p_j}{p_k} = \exp[(\beta_j - \beta_k)' x]$ .

**McFadden's Conditional Logit** So far we focused on the question of how individual characteristics influence the choice. Next, answer the question of how often will individuals choose a new alternative, i.e., express the probability of choice as a function of the *characteristics of the choice*  $k$  (as perceived by individual  $i$ ), say  $z_{ik}$ , not necessarily the characteristics of the individual  $x_i$ .

$$P[y_i = k] = \frac{\exp(\beta' z_{ik})}{\sum_{s=1}^m \exp(\beta' z_{is})} \quad (16.3)$$

<sup>199</sup>Because the difference of two random variables following type I extreme value actually follows logistic distribution. Of course, this is much simpler with Normal distribution, where a difference is again Normal. See Multinomial Probit below.

**Remark 149.** *Individual characteristics which do not change with the choice drop out unless we combine the two models, i.e., allow for both choice and personal characteristics.*

**Exercise 16.4.** Show  $\frac{p_j}{p_k} = \exp[(z_j - z_k)' \beta]$ .

**Remark 150.** *The elimination by aspect model ([M]3.4) represents another way how to account for similarities between alternative choices.*

In the multinomial logit (MNL) model, the joint distribution of extreme value  $\varepsilon$ s does not involve any unknown parameters and is therefore not capable of approximating a wide range of stochastic structures. Furthermore, the MNL assumes that disturbances are *independent* (see Remark 147). When there is correlation, consistency suffers. Consider for example the choice between a blue bus, a red bus and a train. Hence, multinomial logit conforms to the *IIA hypothesis* (independence from irrelevant alternatives). See exercise 16.3 which shows that  $\frac{p_j}{p_k}$  does not depend on characteristics or even the existence of choices other than  $j$  or  $k$ . Hence an introduction of a new alternative means that all of the existing probabilities are reduced by the same amount, irrespective of the new choice degree of similarity to any of the existing ones. The model restricts the choice probabilities to share a uniform set of cross-elasticities.<sup>200</sup>

Inclusion of some potentially correlated alternatives can be tested with a typical Hausman test (Hausman and McFadden, 1984). Under  $H_0 : \text{IIA}$ , one can estimate a subset of the  $\beta_j$  parameters consistently but inefficiently by dropping the individuals who choose the potentially correlated alternatives. These  $\hat{\beta}_j$ s can then be compared to those estimated off the whole data set with all options. Of course, if IIA is violated, the latter will be inconsistent.

**Remark 151.** *In absence of some natural grouping of the alternatives, the choice of the subset to leave out is arbitrary and, hence, so is the test. See also IZA DP No. 5826 for technicalities of the test.*

**Remark 152.** *In study of social interactions (social capital, peer effects, herding behavior), one may want to condition on the group- or region-specific propensity of the particular choice. For example, will I be more likely to go to college when many of my schoolmates do; or will I be more likely to become self-employed when*

---

<sup>200</sup>  $\frac{\partial \log P[y_i=j|z]}{\partial \log z_k} = -P[y_i=k|z] \frac{\partial \beta'_k z_k}{\partial \log z_k}$ . See Pudney, p. 118.

many of my neighbors do so. One then conditions on the predicted probability of a given choice in the group/neighborhood  $E(y_g)$  (endogenous effect) together with group-specific  $z_g$  variables (contextual effects) and individual  $x_{ig}$ . In linear regressions, there is a formidable identification problem with estimating these social equilibria (see Manski's, 1995, "reflection problem" and Durlauf, 2002);<sup>201</sup> however, in non-linear models identification is easier (Brock and Durlauf, 2001). These group-specific propensities are correlated across choices and may absorb the correlation of  $\epsilon$  across choices.<sup>202</sup> See also Brock and Durlauf (in press, JEcm) and Remark 54 for inference issues.

**Fixed Effect MNL** `Femlogit` (in `Stata`) implements a fixed effects (panel-data) MNL model a-la Chamberlain (1980) where choices depend on  $\alpha_{ij} + \beta'_j x_{it}$ . The model assumes no autocorrelation of residuals and strictly exogenous  $x$  conditional on  $\alpha$ s. Using the Chamberlain's trick (see Section 8.3.2),  $\alpha_{ij}$  disappears from the likelihood.<sup>203</sup> For interpretation of estimated effects, see the discussion in Pffor (2014, *Stata Journal*).

**Multinomial Probit and Nested Logit** Given how unattractive the IIA assumption of MNL is in this setting, there are two responses: Multinomial Probit and Nested Logit.<sup>204</sup>

**Multinomial Probit (MNP)** Unlike MNL, the MNP model allows for a full correlation structure with  $\epsilon \sim N(0, \Sigma)$  and requires  $m - 1$  dimensional numerical integration. One has to impose normalization on the  $m(m - 1)$  free elements  $\sigma$  of the  $m \times m$  matrix  $\Sigma$ .

With  $m = 3$ , the choice of the first alternative  $P[y_i = 1|x_i]$  corresponds to the joint occurrence of  $\eta_{12} \equiv \epsilon_1 - \epsilon_2 > V_2(x_i) - V_1(x_i)$  and  $\eta_{13} \equiv \epsilon_1 - \epsilon_3 > V_3(x_i) - V_1(x_i)$ .

<sup>201</sup>Note that the regression  $E[y|z, x] = \alpha + \beta'x + \gamma'z + \delta'E[y|z] + \lambda'E[x|z]$  has a social equilibrium:  $E[y|z] = \alpha + \beta'E[x|z] + \gamma'z + \delta'E[y|z] + \lambda'E[x|z]$ , which you can solve for  $E[y|z]$  and plug back to show that the structural parameters are not identified.

<sup>202</sup>There is still selection bias from the choice of group membership.

<sup>203</sup>In the special case of a binary logit with fixed effects, things simplify to `clogit`.

<sup>204</sup>It is also possible to avoid the IIA assumption by applying the Random Coefficient Model of Section 6.2. Simply allow the  $\beta$ s inside the  $V_j(x_i)$  of Remark 147 to vary across people as individual tastes imply that coefficients are correlated across choices with different characteristics. See McFadden and Train (2000) and the use of the Gibbs sampler for estimation. Another possibility is to follow the Heckman and Singer (1984) approach that we first mentioned in Remark 127 and that we cover in detail in Section 15.2.1.

One can then derive the variance-covariance of the joint normal *pdf* of  $\eta_{12}$  and  $\eta_{13}$ , the 2x2 matrix  $\tilde{\Sigma}$ , from the original  $\sigma$  elements. Finally,

$$P[y_i = 1|x_i] = \int_{-\infty}^{V_2-V_1} \int_{-\infty}^{V_3-V_1} \frac{1}{2\pi\sqrt{|\tilde{\Sigma}|}} \exp\left[-\frac{1}{2}(\eta_{12}, \eta_{13})' \tilde{\Sigma}^{-1}(\eta_{12}, \eta_{13})\right] d\eta_{12}d\eta_{13}.$$

The likelihood requires  $m-1$  dimensional numerical integration, numerical 1st and 2nd derivatives and is therefore potentially messy with  $m > 3$ . The Simulated Method of Moments of Remark 125 is used here; see, e.g. Geweke, Keane and Runkle (1994).

**Nested Logit** Alternatively, the independence assumption of MNL can be relaxed using the generalized extreme value (GEV) models ([M]3.7). The GEV distribution generalizes the independent univariate extreme value *cdfs* to allow for  $\varepsilon$  correlation across choices:

$$F(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_m) = \exp[-G(\exp(-\varepsilon_1), \dots, \exp(-\varepsilon_m))], \quad (16.4)$$

where the function  $G$  is such that  $F$  follows properties of (multinomial) *cdf*. The GEV approach has been widely used in the context of the nested multinomial logit model (see section 16.0.5).

**Example 16.2.** With  $G(a_1, a_2, \dots, a_m) = \sum a_m$  we obtain the simple MNL model. With  $m = 2$  and

$$G(a_1, a_2) = \left[ a_1^{\frac{1}{1-\sigma}} + a_2^{\frac{1}{1-\sigma}} \right]^{1-\sigma}$$

we can interpret the  $\sigma$  parameter as correlation. In this case

$$P[y_i = j|x_i] = \frac{\exp(\frac{V_j(x_i)}{1-\sigma})}{\exp(\frac{V_1(x_i)}{1-\sigma}) + \exp(\frac{V_2(x_i)}{1-\sigma})}$$

where  $V_j$  is the valuation of choice  $j$  (see Remark 147).

Our goal is to (a) study the use of the multinomial logit model in tree structures, and (b) use GEV to allow for departure from IIA *within* groups of alternatives, whilst assuming separability between groups.



**Example 16.3.** *Choice of house: choose the neighborhood and select a specific house within a chosen neighborhood. Choose to travel by plane, then choose among the airlines. Allow for unobservables to be correlated within neighborhoods.*

(a) In presence of a nested structure of the decision problem we assume the utility from house  $j$  in neighborhood  $i$  looks as follows:  $V_{ij} = \beta'x_{ij} + \alpha'z_i$ , where  $z_i$  are characteristics of neighborhoods and  $x_{ij}$  are house-specific characteristics. To facilitate estimation when the number of choices is very large but the decision problem has a tree structure, we use  $p_{ij} = p_i p_{j|i}$ ,<sup>205</sup> where as it turns out  $p_{j|i}$  only involves  $\beta$  but not  $\alpha$ :

$$p_{j|i} = \frac{\exp(\beta'x_{ij} + \alpha'z_i)}{\sum_{n=1}^{N_i} \exp(\beta'x_{in} + \alpha'z_i)} = \frac{\exp(\beta'x_{ij})}{\sum_{n=1}^{N_i} \exp(\beta'x_{in})}. \quad (16.5)$$

Similarly,

$$p_i = \sum_{n=1}^{N_i} p_{ij} = \frac{\exp(I_i + \alpha'z_i)}{\sum_{m=1}^C \exp(I_m + \alpha'z_m)}, \text{ where } I_i = \log \left[ \sum_{n=1}^{N_i} \exp(\beta'x_{in}) \right] \quad (16.6)$$

is the so-called inclusive value (the total contribution of each house in a neighborhood). One can therefore first estimate  $\beta$  off the choice within neighborhoods (based on  $p_{j|i}$ ) and then use the  $\hat{\beta}$  to impute  $\hat{I}_i$  and estimate  $\alpha$  by maximizing a likelihood consisting of  $p_i$ . This sequential estimation provides consistent estimates, but MLE iteration based on these starting values can be used to improve efficiency. If MLE gives different results it suggests misspecification.<sup>206</sup>

**Remark 153.** *The assumed forms of utility functions can differ across branches and decisions.*

**Remark 154.** *The NMNL gives identical fits to data as the hierarchical elimination by aspect model.*

(b) Next, we use generalized extreme value distribution to allow for correlation of the disturbances. Start off by assuming stochastic utility maximization along

<sup>205</sup>Of course  $p_{ijk} = p_i p_{j|i} p_{k|i,j}$ .

<sup>206</sup>The likelihood is no longer globally concave in all parameters. For estimation methods see [MF]p.1426.

the lines of Example 147 but assume GEV instead of type I extreme value. This will lead to a generalization of the NMNL model that actually nests independence (and generalizes to multivariate setting):

$$p_i = \frac{\exp[(1 - \sigma)I_i + \alpha'z_i]}{\sum_{m=1}^C \exp[(1 - \sigma)I_m + \alpha'z_m]} \quad (16.7)$$

Here one can test for within-neighborhood correlation by asking whether  $\hat{\sigma} = 0$ .

**Remark 155.** Zanella (2007, JEEA) extends the social interactions literature (see Remark 152) by applying the nested logit structure to a random utility framework in order to build a model with social interactions and endogenous choice of the group membership (neighborhood). The micro-founded theory then suggests econometric identification strategies.

#### 16.0.6. Ordered Response Models

**Example 16.4.** Ratings, opinion surveys, attained education level.  $0 < 1 < 2$  but  $1 - 0 \neq 2 - 1$ .

Use threshold “constants” to split the range of  $\epsilon$ s. A common  $\beta$  affects the decision among many alternatives.

**Example 16.5.** With 3 ordered choices assume that the latent  $y_i^* = -\beta'x_i + u_i$ . Then (i)  $y_i = 1$  if  $y_i^* < 0 \Leftrightarrow u_i < \beta'x_i$ , (ii)  $y_i = 2$  if  $y_i^* \in (0, c) \Leftrightarrow \beta'x_i < u_i < \beta'x_i + c$ , (iii)  $y_i = 3$  if  $y_i^* > c \Leftrightarrow \beta'x_i + c < u_i$ , where  $c$  is another parameter to be estimated. Following the usual logic the likelihood is based on a product of individual  $i$  contributions, which depend on choice:  $P[y_i = 1|x_i] = F(\beta'x_i)$  while  $P[y_i = 2|x_i] = F(\beta'x_i + c) - F(\beta'x_i)$  and  $P[y_i = 3|x_i] = 1 - F(\beta'x_i + c)$ .

The model generalizes to multinomial settings. Interpreting the coefficients based on their sign (!) is *not* obvious in the ordered response model (see [G] p.738).

**Remark 156.** See Blundell, Pistaferri, and Preston (2008, AER) for a novel two-step estimation procedure that allows applying instrumental variable techniques with ordinal data.

**Remark 157.** See Muris (2017, *REStat*) for a new method for estimating the fixed-effects ordered logit model, allowing one to estimate bounds on the marginal effects because it estimates the cut-off points together with the fixed effects.

**Remark 158.** When comparing (means of) two groups using a parametric ordered response model such as ordered Probit, we cardinalize the reported ordinal responses by assuming that both groups use the same reporting function to translate the latent unobserved opinion, happiness, etc. into the observable discrete responses and that both groups have the same distribution of happiness including the same variance. See Bond and Lang (*JPE* in press) for how this is rejected in happiness research. See Penney (2017) for a variance normalization approach in test score data applications. Note that *Stata* imposes that the ordered Probit variance is 1 and the constant is 0. If one estimates the ordered probit separately for each group, one obtains different cutpoints, implying a different reporting function but the same mean and variance for all groups. One can alternatively normalize the first two cutpoints to 0 and 1 and estimate separate means and variances for each group.

### 16.0.7. Sequential Choice Models

These models have a much richer set of coefficients than the ordered response models. They arise naturally when decisions take place at different points in time (e.g. choice of education level).

In the simplest case assume independence of disturbances and estimate the model using a sequence of independent binary choice models. (In doing so, one places severe restrictions on the underlying preferences and opportunity sets.) On the other hand, they have been used to lower the computational burden of simultaneous choice among  $m$  alternatives with correlated disturbances.

**Example 16.6.** First, choose to graduate from high school or not (this occurs with probability  $1 - F(\beta'_H x_i)$ ); if you do then choose to go to college ( $F(\beta'_H x_i)F(\beta'_C x_i)$ ) or not ( $F(\beta'_H x_i)[1 - F(\beta'_C x_i)]$ ). Note that the likelihood can be optimized separately with respect to  $\beta_H$  and  $\beta_C$  – we can run two separate logit/probit likelihoods, one over the choice of high school, the other over the choice of college (for those who did graduate from high school).

In the most advanced case of modelling intertemporal choice under uncertainty, it is more satisfactory to use dynamic programming techniques.

## 17. Models for Count Data

**Example 17.1.** *Number of accidents in a given plant. Number of visits to a doctor.*

The essential limiting form of binomial processes is **Poisson** distribution:  $P[y = r] = \exp(-\lambda)\lambda^r(r!)^{-1}$ . Assume the number of accidents in each plant follows Poisson with plant-specific parameter  $\lambda_i$  and that these processes are independent across plants. To bring in  $x'\beta$  assume  $\ln \lambda_i = \beta' x_i$  and maximize the likelihood:

$$\max_{\beta} L = \prod_i \exp(-\lambda_i) \lambda_i^{y_i} (y_i!)^{-1}. \quad (17.1)$$

However, Poisson is restrictive in many ways: First, the model assumes independence of number of occurrences in two successive periods. Second, the probability of occurrence will depend on time length of interval. Third, the model assumes the equality of mean and variance:

$$E[y_i | x_i] = V[y_i | x_i] = \lambda_i = \exp(\beta' x_i) \implies \frac{\partial E[y_i | x_i]}{\partial x_i} = \lambda_i \beta \quad (17.2)$$

The last assumption is relaxed by the **Negative Binomial** extension of Poisson, which allows for overdispersion: Let  $\ln \lambda_i = \beta' x_i + \epsilon$  where  $\epsilon \sim \Gamma(1, \alpha)$ . Integrate the  $\epsilon$  out of likelihood before maximization (as in Random Effect Probit) and maximize w.r.t. both  $\beta$  and  $\alpha$ , the overdispersion parameter.<sup>207</sup>

See, e.g., LIMDEP manual, section 26.2, for extensions of both models to panel data, censoring and truncation, the zero-inflated probability (see immediately below) and sample selection (see section 19).

**Remark 159.** *For example, the quasi-maximum likelihood conditional fixed-effects Poisson model (Poisson QML; see Hausman et al., 1984) has several desirable properties, including consistency of the coefficient estimates independently of any assumption on the conditional variance as long as the mean is correctly specified (Woolridge, 1997) and consistency of the standard errors even if the data generating process is not Poisson. This estimator can also be used for fractional and non-negative variables (Santos Silva and Tenreiro, 2006). It is implemented in Stata as `xtqmlp`<sup>208</sup> and its coefficient estimates are interpreted as  $(\exp(\beta) - 1) * 100$  percentage change or approximately as  $\beta * 100$  percentage change.*

<sup>207</sup> $V[y_i | x_i] = E[y_i | x_i](1 + \alpha E[y_i | x_i])$

<sup>208</sup><http://people.bu.edu/tsimcoe/code/xtqml.txt>

### 17.1. Threshold Models

Combine a binary choice model with other likelihoods. In case of the count data estimation this approach has been coined as the zero inflated probability model:

Consider the example of accident counts in plants. The zero-inflated version allows for the possibility that there could not be any accidents in plant  $i$ : When we observe  $y_i = 0$ , it can either correspond to our usual Poisson data generating process where out of luck, there were no accidents in the given time period or it can correspond to a plant where the probability of having an accident is zero (in that event  $Z_i = 1$ ): that is we can express  $P[y_i = 0|x_i]$  as

$$\begin{aligned} P[0|x_i] &= P[Z_i = 1|x_i] + P[Z_i = 0|x_i]P[y_i^* = 0|x_i] = \\ &= F(\gamma' x_i) + (1 - F(\gamma' x_i)) \exp(-\lambda_i) \lambda_i^{y_i} (y_i!)^{-1}. \end{aligned} \quad (17.3)$$

Ideally, there is at least one variable affecting  $Z$ , but not  $y_i^*$  to aid identification.

## 18. Limited Dependent Variables

See [W] 16, [M]6, [G]22, [P]4. Let's combine qualitative choice with continuous variation.

**Example 18.1.** *Wage data censored from above at the maximum social security contribution level.*

**Example 18.2.** *Zero expenditure and corner solutions: labor force participation or R&D expenditure. A change in the  $x$ s affects both the usual intensive margin and the extensive margin of the corner solution.*

Note the fundamental difference between the two examples.

The difference between censoring and truncation, which both have to do with thresholds on observable  $y$ s, is in observing the  $x$ s for the censored values of  $y$ s.

### Technical Preliminaries

(a) Means of truncated distributions:

$$E[\varepsilon \mid \varepsilon \geq c] = \int_c^\infty \frac{\varepsilon f(\varepsilon)}{1 - F(c)} d\varepsilon \quad (18.1)$$

**Exercise 18.1.** *Show that if  $\varepsilon \sim N(\mu, \sigma^2)$  then  $E[\varepsilon \mid \varepsilon \geq c] = \mu + \sigma \lambda(\frac{c-\mu}{\sigma})$ , where  $\lambda(\cdot) = \frac{\varphi(\cdot)}{1-\Phi(\cdot)}$  is the so called inverse of the Mills' ratio.<sup>209</sup> Also find  $V[\varepsilon \mid \varepsilon \geq c]$ .*

(b) Means of censored distributions, where  $\varepsilon^c = \max\{c, \varepsilon\}$  :

$$E[\varepsilon^c] = F(c)c + [1 - F(c)]E[\varepsilon \mid \varepsilon \geq c] \quad (18.2)$$

### 18.1. Censored Models

**Example 18.3.** *When actual income is above \$ 100 000, the reported income is \$ 100 000.*

The structural model is using the concept of an underlying latent variable  $y_i^*$  :

$$\begin{aligned} y_i^* &= \beta' x_i + u_i \text{ with } u_i \sim N(\mu, \sigma^2) \\ y_i &= y_i^* \text{ iff } y_i^* > c \\ y_i &= c \text{ iff } y_i^* \leq c \end{aligned} \quad (18.3)$$

<sup>209</sup>Using the fact that  $\int \varepsilon \phi(\varepsilon) d\varepsilon = -\int d\phi(\varepsilon) = -\phi(\varepsilon)$

**Tobit Model** When the data are censored, variation in the observed variable will understate the effect of the regressors on the “true” dependent variable.<sup>210</sup> As a result, OLS will typically result in coefficients biased towards zero.

WLOG<sup>211</sup> suppose the threshold occurs at  $c = 0$ . OLS is inconsistent no matter whether we include or exclude the zero observations because  $E[\widehat{\beta}_{OLS} - \beta]$  depends on the truncated expectation of  $u_i$  in either case.

**Exercise 18.2.** Characterize the bias of the OLS estimator when applied only to the nonzero  $y$  observations and show that OLS estimator when applied to all  $y$  observations is inconsistent (we are essentially omitting the explanatory variable  $\lambda(x'_i\beta/\sigma)$ ).

The Tobit likelihood is

$$L = \prod_{y_i^* > c} \frac{1}{\sigma} \varphi\left(\frac{y_i - x'_i\beta}{\sigma}\right) \prod_{y_i^* \leq c} \Phi\left(\frac{c - x'_i\beta}{\sigma}\right), \quad (18.4)$$

which has a single maximum, but two step procedures have been devised by Heckman ([M]8.2) and Amemiya ([M]6.5).

**Remark 160.** The two step procedure of Heckman starts with a Probit on  $y_i > 0$  or not. This delivers consistent  $\widehat{\beta}/\sigma$ . In the second step, bring in the continuous information and consider

$$\begin{aligned} E[y_i|x_i] &= P[y_i^* > 0]E[y_i|y_i^* > 0] + P[y_i = 0]E[y_i|y_i = 0] = \\ &= \Phi\left(\frac{x'_i\beta}{\sigma}\right) x'_i\beta + \sigma\varphi\left(\frac{x'_i\beta}{\sigma}\right) + 0 = \Phi_i x'_i\beta + \sigma\varphi_i. \end{aligned}$$

Use the first-step  $\widehat{\beta}/\sigma$  to predict  $\widehat{\Phi}_i$  and  $\widehat{\varphi}_i$  and estimate  $y_i = \widehat{\Phi}_i x'_i\beta + \sigma \widehat{\varphi}_i$  for a new set of  $\widehat{\beta}$  and  $\widehat{\sigma}$ . As usual, drastic differences between first- and second-step estimates signal misspecification.

**Remark 161.** The likelihood is only piece-wise continuous.

<sup>210</sup>Plot a graph with censored data to see that you get no change in  $y$  for changes in  $x$  for the censored values of  $y$ .

<sup>211</sup>Assuming there is a constant in  $x$  ([M]p.159).

**Remark 162.** *The estimator is biased in the presence of heteroskedasticity. See Arabmazar and Schmidt (1981) for the potential magnitude of the bias. See Koenker and Bassett (1982) for quantile regression tests for heteroskedasticity. Pagan and Vella (1989) propose a test for heteroskedasticity when the dependent variable is censored. Need zero-expected-value residuals to construct the test. These can be obtained by a trimmed LS estimator (Powell 1986). See section 18.3 for recent heteroskedasticity-robust alternatives to Tobit such as CLAD.*

Of course, the model can be easily extended to censoring from above *and* below. The model also has extensions allowing for multiple and variable thresholds  $c_i$ , endogenous thresholds, heteroskedasticity,<sup>212</sup> panel data random effects, sample selection (see section 19), SEM, nested structures, and non-normality (see LIMDEP manual, Ch. 27).<sup>213</sup>

How do we interpret the coefficients? Depends on the type of data generating mechanism / question we ask. There are 3 types of predictions we can consider, using the definitions of  $E[y^* | x]$  (for data censoring problems) and  $E[y | x]$  and  $E[y | x, y^* > 0]$  (for corner solution problems). But then we would not want to use a Tobit for a corner solution problem. Or other latent variable models for that matter, since what is the meaning of negative  $y^*$  latent health care expenditure? See also [W] 16.1, 16.5.<sup>214</sup>

**Exercise 18.3.** *Find the expressions for these 3 conditional mean functions and their derivatives w.r.t.  $x$ .  $E[y | x]$  has a particularly simple form. What is the intuition? Marginal variation in  $x$  does not lead to anyone crossing the extensive margin. In the case of evaluating the effect of a dummy variable, the formula makes clear that the effect is much smaller than the corresponding parameter, which is intuitive given that the latent outcome always changes when the dummy switches, while the observed outcome stays the same (at 0) for many individuals.*

**Remark 163.** *There is little theoretical justification for Tobit in rational choice models (see [P]p.141).*

---

<sup>212</sup>The available methods use a parametric assumption on the form of heteroscedasticity. Semi-parametric estimators are a focus of much current research, see section 18.3.

<sup>213</sup>For a survey of Tobit specification tests see [P]4.1.5. For further reading see special issues of the *Journal of Econometrics* (84-1,86-1,87-1). One strand of tests is based on conditional moment restrictions, see [G]22.3.4d.

<sup>214</sup>Bauer and Sinning (2005 IZA DP; 2008 AStA) provide a Oaxaca-Blinder decomposition for Tobit and, more generally, for non-linear models.



**Remark 164.** Under the strong assumption of joint normality of error terms, one can instrument in a Tobit.

### Grouped Data

**Example 18.4.** Wages reported only in ranges, i.e.  $w_i \in [\$10000, \$20000)$ , i.e.  $w_i \in [c_{j-1}, c_j)$   $j = 1, \dots, J$

The difference between this model and the ordered choice models is that the threshold values are known here. For  $c_j = H$  and  $c_{j-1} = L$  the likelihood contribution of observations with  $y_i$ s in those ranges is

$$\ln L_{HL} = \sum_{i=1}^N \left\{ \ln[\Phi(\eta H - x'_i \gamma) - \Phi(\eta L - x'_i \gamma)] \right\}, \quad (18.5)$$

where  $\gamma = \frac{\beta}{\sigma}$  and  $\eta = \frac{1}{\sigma}$  (similar reparametrization of the likelihood is used in the estimation of the Tobit mode, see Olsen 1978). Use

$$E[y_i^* | x_i] = x'_i \beta + \sigma \frac{\varphi_{iL} - \varphi_{iH}}{\Phi_{iH} - \Phi_{iL}} \quad (18.6)$$

for prediction. Again, the model can be extended to allow for sample selection.

### 18.2. Truncated Models

**Example 18.5.** Only have data on low income households when studying the impact of variable  $x$  on income  $y$ .

$$L = \prod_{y_i^* > c} \frac{1}{\sigma} \varphi \left( \frac{y_i - x'_i \beta}{\sigma} \right) \left[ 1 - \Phi \left( \frac{c - x'_i \beta}{\sigma} \right) \right]^{-1} \quad (18.7)$$

Tobit-type model is not feasible here as we do not observe  $x$ s for the  $y = c$  observations. To evaluate the impact of  $x$  on  $y$ , in a simple truncation, use

$$E[y_i | x_i, y_i < c] = x'_i \beta + \sigma \frac{\varphi(c/\sigma)}{\Phi(c/\sigma)}.$$

In a double truncation region, use Equation 18.6 for  $E[y_i | c_{iL} \leq y_i \leq c_{iH}]$ .

Finally, it is an opportune time, to note that the Tobit model is restrictive in constraining the coefficients and the  $x$ s affecting the extensive and intensive margins to be the same ([G]p.770).

**Example 18.6.** Consider studying the impact of the age of a building on the cost from a fire in that building. In some buildings there is no fire and cost is zero, in other buildings you observe a fire and the associated cost. It is likely that older buildings are more likely to experience fire, while the cost of fire, conditional on having one, is likely to be higher in a newer building.

We can relax the Tobit likelihood and split it into two (independent) parts: (i) 0/1 probit for whether there is a fire or not, and (ii) a truncated normal regression of the cost of fire estimated on those buildings where there was a fire. Further, we can allow different explanatory variables to enter each of the separate two likelihoods.

**Remark 165.** Assuming the  $x$ s affecting both margins (equations) are the same, note that under the equality of coefficients, the relaxed two-part model boils down to the restricted Tobit model. Hence, the equality of coefficients is testable using an LR test:

$$LR = -2\{\ln L_{PROB} + \ln L_{TRUNC} - \ln L_{TOBIT}\} \sim \chi^2(k) \text{ where } k = \dim(\beta).$$

**Remark 166.** But the disturbances from the two separate equations are likely dependent, which is why we need a sample selection model!

**Remark 167.** Note that (without any covariates) a linear regression of an “expenditure” (corner solution)  $Y$  on a binary treatment would give the (unconditional) average treatment effect for expenditures, but a truncated regression (conditional on positive) would not have a causal interpretation in a randomized trial because the experiment changes the composition of the group with positive expenditures.

### 18.3. Semiparametric Truncated and Censored Estimators

If the residual in a censored model is subject to heteroskedasticity of an unknown form or if we do not know the distribution of the  $\varepsilon$  for sure, then standard MLE will be inconsistent. Also, maximum likelihood estimation of censored panel-data fixed-effect models will be generally inconsistent even when we have the correct parametric form of the conditional error distribution (Honoré, 1992).

Below, we will continue to specify the regression function parametrically, but will try to do without assuming parametric distributions for  $\varepsilon$ . The estimators will alternate between additional “recensoring,” which will compensate for the original censoring in the data, and a “regression” step using only the “trimmed” data part. For simplicity, consider only censoring or truncation from below at 0.

**Symmetrically Trimmed Least Squares** How can we estimate truncated or censored models without relying on particular distributional assumptions? Consider truncation from below at 0 in a model  $y_i^* = x_i'\beta + \epsilon_i$ . The idea of the estimator is to trim (truncate) the dependent variable *additionally* from above to make it symmetrically distributed. The new dependent variable will be symmetrically distributed around the regression function so we can apply least squares. But where do you trim from above? Depends on  $\beta$ . Assume that  $f_\epsilon(s|x)$  is symmetric around zero and unimodal. Then for  $x_i'\beta > 0$ , the  $\epsilon$  is truncated at  $0 - x_i'\beta$  so a symmetric truncation of  $\epsilon$  is at  $x_i'\beta - 0$ . This corresponds to truncating  $y$  at  $2x_i'\beta - 0$  (plot a distribution graph to see this point).

Powell's (1986) Symmetrically Trimmed LS is consistent and asymptotically normal for a wide class of symmetric error distributions with heteroskedasticity of unknown form. With data truncated from below, the estimator minimizes

$$\sum I\{x_i'\beta > 0\}I\{y_i < 2x_i'\beta\} [y - x_i'\beta]^2. \quad (18.8)$$

Alternatively, with censoring from below, apply the same idea (Symmetrically Censored LS) to minimize

$$\sum I\{x_i'\beta > 0\} [\min(y_i, 2x_i'\beta) - x_i'\beta]^2. \quad (18.9)$$

**Censored Least Absolute Deviation** [W] 16.6.4. Powell's (1984) CLAD is again based on *additional* censoring of  $y$ . The main idea is to look at median as opposed to mean, because median is not affected by censoring. (The main assumption of the estimator is therefore zero median of  $F_\epsilon(s|x)$ .) Median is not affected as long as we are in the uncensored part of the data. If we are below the censoring point, then the median does not depend on  $x'\beta$ . So:  $median(y_i^*|x_i) = \max\{x_i'\beta, 0\}$ . (Again, we work only with variation in the  $x_i'\beta > 0$  area.) We note that medians are estimated using LAD.<sup>215</sup> The CLAD estimator is then found by minimizing

$$\sum |y_i - \max\{0, x_i'\beta\}|. \quad (18.10)$$

---

<sup>215</sup>Estimate  $\delta$ , the median of  $z_i$ , by  $\min \sum_i |z_i - \delta|$ . LAD is not a least-squares, but a median (quantile) regression (these are in general more robust to outliers, see Section 12.1). Quantile regression fits a linear model for conditional quantiles, just as OLS fits a linear model for conditional means.

It alternates between deleting observations with estimates of  $x_i'\beta$  that are outside of the uncensored region and estimating the median regression based on the remaining observations.

It can be used in a data censoring (Tobit) or corner solution (Heckman's  $\lambda$  of Section 19.2.2) setup and in presence of heteroskedasticity and it is even more robust than STLS. CLAD is programmed into Stata.<sup>216</sup> CLAD has a small sample bias opposite to the OLS bias; for a two-step remedy, see Khan and Powell (2001, JEcm). Honore, Khan and Powell (2002) allow the censoring thresholds to be not always observed. Blundell and Powell (2007, JEcm) allow for IVs.<sup>217</sup> Also see Newey and Powell (1990). More accessible treatment of related topics can be found in a book on *Methods of Moments and Semiparametric Econometrics for Limited Dependent Variable Models* by Myoung-jae Lee.<sup>218</sup>

## 19. Sample Selection

A focus of an enormous volume of empirical and theoretical literature. It involves features of both truncated and censored models. The treatment of the data where sampling depends on outcomes is different in cases where the variables determining selection are observable and when they are not (for an introductory discussion see [P]2.5).

We first consider situations where the econometrician (data collection agency) chooses to base sampling on  $y$ .

**Example 19.1.** *Consider a sample of families and estimate the impact of  $x$  on family income. However, the sample is such that low-income families are over-sampled.*

Second, we consider situations where the individual behavior results in sample selection: situations where people/firms/etc. select themselves into different states based on potentially unobserved characteristics.

**Example 19.2.** *You can only measure the impact of  $x$  on wages ( $y$ ) for those women who work (selection on  $y$ ). Whether or not a woman works, depends on the wage she could get when working.*

<sup>216</sup>By D. Jolliffe, a former CERGE-EI faculty member, and by two of your co-students.

<sup>217</sup><http://www.ucl.ac.uk/~uctp39a/BlundellPowellCRQDec05.pdf>

<sup>218</sup>Stata ado files for these semiparametric models can be downloaded from <http://emlab.berkeley.edu/users/kenchay> .

**Example 19.3.** *Average wage over the business cycle: seems flatter due to selective drop out of work.*

### 19.1. Sampling on Outcome

We first think about sampling based on qualitative choice and then on a continuous outcome variable.

#### 19.1.1. Choice-based sampling

Consider a binary choice problem. To analyze a rare event when population probabilities  $p(y_i) = P[y_i = 1]$  are tiny (training treatment, violent crime), we can decide to save money and instead of collecting a large random sample of the population (in order to have a decent number of  $y_i = 1$  group members), we sample randomly within each  $y$  group to obtain  $f(x_i | y_i)$ . Then we note  $f(x_i, y_i) = p(y_i | x_i)f(x_i) = f(x_i | y_i)p(y_i)$  and write the likelihood function for the two samples

$$L(\cdot) = \prod_{i \in S_1} f(x_i | y_i = 1) \prod_{i \in S_2} f(x_i | y_i = 0) \quad (19.1)$$

in terms of  $p(y_i | x_i) = F(\beta' x_i)$  (in the bivariate example).<sup>219</sup>

**Remark 168.**  $P[y_i = 1], P[y_i = 0]$  usually come from a different data set (are known), but can be estimated as part of the problem.

Next, consider a multinomial choice problem. Manski and McFadden (1981) set up an intuitive conditional maximum likelihood estimator using the formula for the conditional probability of  $i$  given  $x$  in the sample. For  $j = 1, \dots, M$  choices:

$$L(\theta) = \sum_{i=1}^N \ln \frac{p(y = y_i | x_i, \theta) H_{y_i} / p(y = y_i)}{\sum_{j=1}^M p(y_i = j | x_i, \theta) H_j / p(y_i = j)}, \quad (19.2)$$

where  $H_j$  is the probability of sampling from a strata  $j$  (can be unknown to the researcher).<sup>220</sup>

<sup>219</sup>See also Pudney (1989), chapter 3.2.3.

<sup>220</sup>However, as noted in Cosslett (1981) paper, this estimator is not efficient. Partly because the sample distribution of  $x$  actually depends on  $\theta$ :  $g(x) = \sum_{s=1}^S H_s / p(s) \sum_{j \in I(s)} p(j|x, \theta) f(x)$ . So we should use this information to help us estimate  $\theta$  better. But then we do not get rid of  $f(x)$ , which was the beauty of 19.2.

**Remark 169.** Replacing  $H$  with  $\hat{H}$  actually improves efficiency. Counterintuitive. Only works with an inefficient estimator.<sup>221</sup>

To repeat: selecting “balanced” data with the same number of 0 and 1 outcomes and running a simple Logit is wrong, even though it is regularly done in practice. See Scott and Wild (1997) for a bias correction. There is much practical research on economically optimal sampling design, i.e., saving money on data collection and delivering efficiency. Some of the sampling designs combine sampling on the key right-hand-side variable with sampling on the outcome.<sup>222</sup>

### 19.1.2. Endogenous Stratified Sampling

It occurs when the probability that an individual is observed in the sample depends on the (continuous)  $y$ .<sup>223</sup>

**Example 19.4.** Suppose that you wish to estimate returns to schooling using the PSID, a survey that over-samples low-income households. This leads to inconsistent estimation, which can be corrected using weighting, because the sampling criterion (family income) is related to the error term in the regression for log earnings.

**Remark 170.** Stratification (over/undersampling) based on  $x$  variables presents no problem for OLS, as long as there is no parameter heterogeneity across strata (see Remark 21).

---

<sup>221</sup>Cosslett (1981) devised a pseudo-likelihood estimator, replacing the  $f(x)$  with a set of discrete densities  $f_n$ . Counterintuitively, even though the number of parameters climbs with  $n$ , this estimator is efficient. (Think of estimating a mean this way.) However, it is not practical. So, Imbens (1992) comes up with a reparametrization of Cosslett moment conditions which is implementable. It is based on the intuition that to devise a moment condition based on  $x$  with many points of support, I do not need to know the points of support themselves (see the example below). He uses change in variables in the FOC (moment conditions) of the Cosslett estimator between a subset of the points of support of  $x$  and the population marginal densities  $p$  to come up with nice moment conditions.

Example: Suppose you want to estimate  $\delta = \Pr(z > 0)$ . If  $z$  is discrete with  $\{z^1, z^2, \dots, z^L\}$  points of support and unknown probabilities  $\{\pi_1, \pi_2, \dots, \pi_L\}$  one could efficiently estimate  $\delta$  on the basis of  $i = 1, \dots, N$  independent observations of  $z_i$  by ML as  $\hat{\delta} = \sum_{l|z^l > 0} \hat{\pi}_l = \frac{1}{N} \sum_{n=1}^N I[z_n > 0]$  where the last representation of the estimator does not depend on the points of support. It can also be used when  $\delta$  does not have a discrete distribution.

<sup>222</sup>This is a detailed Monte Carlo study of alternative sampling designs for Logit: <http://www.occ.treas.gov/ftp/workpaper/wp2001-3.pdf>

<sup>223</sup>Simple estimation using data split into subsamples based on the level of the dependent variable is a no-no thing in econometrics.

Assume the final sample is obtained by repeated random drawings, with each draw being made from stratum  $i$  with probability  $p_j = \frac{n_j}{N_j}$  which is independent of the  $x$ s. Here  $n_j$  is the number of observations in data from strata  $j$  and  $N_j$  is the population size of strata  $j$ . Let  $y_{ij}$  denote the value of  $y$  for person  $i$  from strata  $j$ . The typical solution in practice is WLS (Wooldridge, 1999):

$$\min \sum_{i,j} \frac{1}{p_j} (y_{ij} - x'_{ij}\beta)^2,$$

which works asymptotically. In small samples it will be biased.

A potentially better solution is MLE. Consider an example of endogenous stratification (think oversampling or undersampling) with known threshold  $L$  and with 2 strata ( $j = 1, 2$ ) of the level of  $y$  ([M]6.10.). Assume Normality and maximize a likelihood based on<sup>224,225</sup>

$$\begin{aligned} L(y_i|x_i) &= L_i^{-1} p_1 \phi((y_i - x'_{ij}\beta)/\sigma) \text{ if } y_i < L \text{ and} & (19.3) \\ L(y_i|x_i) &= L_i^{-1} p_2 \phi((y_i - x'_{ij}\beta)/\sigma) \text{ if } y_i > L \text{ where} \\ L_i &= p_1 \Phi[(L - x'_{ij}\beta)/\sigma] + p_2 (1 - \Phi[(L - x'_{ij}\beta)/\sigma]). \end{aligned}$$

In the next subsection, we will consider cases when truncation or censoring occurs with stochastic or unobservable thresholds and where individuals make the sample-inclusion choice, not the data collection agency.

## 19.2. Models with Self-selectivity

Now it is the individuals that we study who makes the sampling decision.

**Example 19.5.** *Fishing and hunting: the Roy's model ([M]9.1); workers choose their union status based on the wage "in" and on the wage "out"; labor force participation; returns to education; migration and income; effect of training programs, evaluation of social policy.*

There are two main types of models: first, when we do not observe the  $y$  under one choice and observe it under the other (labor force participation, Heckman's  $\lambda$ ), second, when we observe  $y$  under all *chosen* alternatives (union wages, switching regression).

<sup>224</sup>The formulas in [M]6.10 are conditional on  $y_{ij}$  actually being drawn from a given strata  $j$ .

<sup>225</sup>See [Mp.173] for the asymptotic justification of WLS based on this MLE.

### 19.2.1. Roy's model

First consider a classical theory (paradigm) on the topic. A worker  $i$  chooses to either hunt or fish, depending on which of corresponding outputs  $y_{iH}$  and  $y_{iF}$  is larger. Note that we never observe both  $y_{iH}$  and  $y_{iF}$  for each worker, but only one of the two outcomes.<sup>226</sup>

Assuming that

$$\begin{pmatrix} y_{iH} \\ y_{iF} \end{pmatrix} \sim N \left( \begin{pmatrix} \mu_H \\ \mu_F \end{pmatrix}, \begin{pmatrix} \sigma_H^2 & \sigma_{HF} \\ \sigma_{HF} & \sigma_F^2 \end{pmatrix} \right),$$

one can show that

$$E[y_{iH}|y_{iH} > y_{iF}] = \mu_H + \frac{COV(y_{iH}, y_{iH} - y_{iF})}{\sqrt{V(y_{iH} - y_{iF})}} \frac{\phi(z)}{\Phi(z)}, \text{ where } z = \frac{\mu_H - \mu_F}{\sqrt{V(y_{iH} - y_{iF})}}.$$

In short,  $E[y_{iH}|y_{iH} > y_{iF}] = \mu_H + \frac{\sigma_H^2 - \sigma_{HF}}{\sigma} \frac{\phi(z)}{\Phi(z)}$  and similarly for  $E[y_{iF}|y_{iH} > y_{iF}]$ . There are 3 possible cases of a Roy's economy:

1. If  $\sigma_H^2 - \sigma_{HF} > 0$  and  $\sigma_F^2 - \sigma_{HF} > 0$ , those who hunt are better off than an average hunter (similarly for the fishermen). This is the case of absolute advantage.
2. When  $\sigma_H^2 - \sigma_{HF} > 0$  and  $\sigma_F^2 - \sigma_{HF} < 0$ , those who hunt are better than average in both occupations, but they are better in hunting (comparative advantage).
3. Reverse of 2.

**Remark 171.** Notice that individuals with better skills choose the occupation with higher variance of earnings. Also notice the importance of  $\sigma_{HF} \neq 0$  (see Remark 166).

**Remark 172.** Note that  $\sigma_H^2 - \sigma_{HF} < 0$  and  $\sigma_F^2 - \sigma_{HF} < 0$  cannot happen due to Cauchy-Schwartz inequality.

**Remark 173.** Think of real-life applications when you really want to know the population-wide  $\mu_H$ .

---

<sup>226</sup>For purposes of policy evaluation we will need to deal with estimating the counterfactual. See subsection 20.



**19.2.2. Heckman's  $\lambda$** 

See Heckman (1980), [M]6.11, 8.4. Consider a two-equation behavioral model:

$$\begin{aligned} y_{i1} &= x'_{i1}\beta_1 + u_{i1} \\ y_{i2} &= x'_{i2}\beta_2 + u_{i2}, \end{aligned} \quad (19.4)$$

where  $y_{i1}$  is observed only when  $y_{i2} > 0$ .

**Example 19.6.** Observe wages ( $y_{i1}$ ) only for women who work ( $y_{i2} > 0$ ).

Note that the expectation of data on  $y_{i1}$  you observe depends on the selection rule which determines that  $y_{i1}$  is observable:

$$\begin{aligned} E[y_{i1}|x_i, y_{i2} > 0] &= x'_{i1}\beta_1 + E[u_{i1}|selection\ rule] = \\ &= x'_{i1}\beta_1 + E[u_{i1}|y_{i2} > 0] = x'_{i1}\beta_1 + E[u_{i1}|u_{i2} > -x'_{i2}\beta_2]. \end{aligned} \quad (19.5)$$

We have an omitted variable problem:  $x_{i2}$  enters the  $y_{i1}$  equation. Of course  $E[u_{i1}|u_{i2} > -x'_{i2}\beta_2] = 0$  if  $u_{i1}$  and  $u_{i2}$  are independent (again, think of Remark 166 and  $\sigma_{HF}$  in the Roy's model).

If we assume that  $u_{i1}$  and  $u_{i2}$  are jointly normal with covariance  $\sigma_{12}$  and variances  $\sigma_1^2$  and  $\sigma_2^2$  respectively, we know what  $E[u_{i1}|u_{i2} > -x'_{i2}\beta_2]$  looks like: It is the usual inverse of the Mills' ratio, which we will call here Heckman's lambda:

$$E[y_{i1}|x_i, y_{i2} > 0] = x'_{i1}\beta_1 + \frac{\sigma_{12}}{\sigma_2} \frac{\phi(x'_{i2}\beta_2/\sigma_2)}{\Phi(x'_{i2}\beta_2/\sigma_2)} = x'_{i1}\beta_1 + \sigma_\lambda \lambda(x'_{i2}\beta_2). \quad (19.6)$$

While we can numerically identify  $\sigma_\lambda$  from  $\beta_1$  even when  $x_{i2} = x_{i1}$  because  $\lambda$  is a non-linear function, there is need for exclusion restrictions (variables in  $x_{i2}$  not included in  $x_{i1}$ ) in order to avoid identification by functional form (i.e. by distributional assumption implying nonlinearity in  $xs$ ).

The model can be estimated by FIML or in two stages. The two-stage estimation starts with a probit on  $y_{i2} > 0$  which delivers  $\widehat{\beta}_2$  which can be used to calculate  $\widehat{\lambda}_i = \lambda(x'_{i2}\widehat{\beta}_2)$ . In the second stage  $y_{i1}$  is run on  $x_{i1}$  and  $\widehat{\lambda}_i$  to estimate  $\widehat{\beta}_1$  and  $\widehat{\sigma}_\lambda$ . Of course, if  $\widehat{\sigma}_\lambda = 0$ , selection is not important.

The joint normality implies a particular form of heteroskedasticity at the second step regression (GLS matrix  $\Gamma$ ). Further, we have to make another GLS correction for the fact that we're not using  $\lambda_i(z)$  but only  $\widehat{\lambda}_i(z)$  so that the error term contains the following:  $\sigma_\lambda(\lambda_i - \widehat{\lambda}_i) \cong \frac{\partial \lambda_i}{\partial z}(\beta_2 - \widehat{\beta}_2)x_{i2}$  evaluated at  $z = x'_{i2}\beta_2$

(this approach of using the first order Taylor series approximation is often called the Delta method—the point is general: whenever you use predicted regressors in a non-linear function, you need to correct your standard errors!). Hence, the variance-covariance of the error term in the second-step regression is composed of  $\Gamma$  plus  $(\frac{\partial \lambda}{\partial z})' \text{Var}(\widehat{\beta}_2) (\frac{\partial \lambda}{\partial z})$ .

Recent semiparametric literature is relaxing the assumption of joint normality of disturbances (see section 19.2.4 below).

**Example 19.7.** First run probit on labor force participation and obtain  $\widehat{\lambda}$ , then run the wage regression to get the effect of education on wages  $\widehat{\beta}$  (and  $\widehat{\sigma}$ ).

**Example 19.8.** Consider the hours labor-supply regression with wages on the RHS. First, you need to correct the hours equation for sample selection into labor force (only observe  $h$  for those who work). This correction comes from a comparison of behavior equations governing reservation wages  $w_i^R$  and market wages  $w_i$  which leads to a 0/1 participation estimation depending on  $Z_i' \gamma$ , where  $Z$  is the collection of RHS variables from both  $w_i^R$  and  $w_i$  equations. Second, you need to instrument for  $w_i$  which is likely endogenous. The first stage regression where you predict  $\widehat{w}_i$  also needs to have a selection correction in it. Finally, you can estimate

$$h_i = \delta \widehat{w}_i + x_i' \beta + \sigma \lambda(Z_i' \widehat{\gamma}) + \varepsilon_i.$$

There is serious need for exclusion restrictions: you need an exclusion restriction for running IV for  $w_i$  (that is a variable predicting wages but not hours) and you need another exclusion restriction to identify the selection correction in the first-stage wage equation (that is you need a variable affecting participation, but not wages).

**Remark 174.** Asymptotic distribution: two stage methods are efficient in one iteration.

**Remark 175.** If the unobservable selection threshold is time constant we can use a fixed effect panel data model to deal with it.

**Example 19.9.** The  $\lambda$  method is applicable in unbalanced panel data, see Lafontaine and Shaw (1995) for an example. Franchisees that go out of business have shorter  $T_i$ . Their fixed effect model appears to eliminate most of the selection bias suggesting that within-firm variation in selection has little effect. In other words,

the FEM approach will work when the reason why firms exit is the time constant unobservable. To apply the  $\lambda$  method, one needs (time changing) predictors of survival excluded from the  $y$  equation. Finally, in cases when it is reasonable that those who disappear from the data have  $y$  that is always below the median of the (surviving) sample, one can insert such a small value of  $y$  for the firms that went belly up and use them in an estimation of a median (LAD) regression, where their particular  $y$  value is of no importance, as long as it is below the median.

**Remark 176.** What about sample selection in panel data when we're also instrumenting? Semykina and Wooldridge (2010, *JofEcm*) analyze pooled 2SLS and fixed effects-2SLS estimators and propose tests for selection bias as well as a parametric and a semiparametric estimation procedure that correct for selection in the presence of endogenous regressors.<sup>227</sup>

**Remark 177.** Estimating parametric Limdep models with a  $\lambda$  on the RHS is a big problem, especially with heteroskedasticity, which kills consistency. The  $\lambda$  problem is that selection affects the whole distribution and  $\lambda$  only fixes the expectation (centering).

**Remark 178.** Buchinsky (1998, *JAppliedEcm*) provides a (rarely used) method of estimating quantile (median) regressions with a sample selection correction. See Section 12.1 for an introduction to quantile regressions.

### 19.2.3. Switching Regression

In this case we observe  $y$  (and  $x$ ) under all *chosen* alternatives.

**Example 19.10.** The union-nonunion or migrants-stayers wage model. The owner-rental housing demand. The privatized-state profit function.

A first approximation in case of two choices is a restrictive **constant effect model** which pools data into one regression:

$$y_i = x_i' \beta + \alpha D_i + \varepsilon_i, \quad (19.7)$$

---

<sup>227</sup>The full working paper version is at <http://mailer.fsu.edu/~asemykina/selectionPS.pdf>  
Stata do-files (parametric and semiparametric, respectively) are at  
[http://mailer.fsu.edu/~asemykina/two\\_step\\_se\\_parametric.do](http://mailer.fsu.edu/~asemykina/two_step_se_parametric.do)  
<http://mailer.fsu.edu/~asemykina/semiparametric.zip>

which is estimated by IV under the assumption that  $y_{1i} - y_{0i} = \alpha$ , where 0 and 1 denote the two different states (union/nonunion, treatment/control). The first stage is based on  $P[D_i = 1 | z_i] = P[z_i' \gamma + \nu_i \geq 0 | z_i]$  so that  $\widehat{D}_i = F_\nu(z_i' \widehat{\gamma})$  for a symmetric  $F_\nu(\cdot)$ .

**Remark 179.** *The consistency of the main equation is contingent on correct specification of the error distribution  $F_\nu(\cdot)$ . See Remarks 73 and 123.*

**Remark 180.** *Again, the standard errors of  $\beta$  need to be corrected for the estimation of  $\gamma$ , see Lee (1981), and the discussion of Delta method above for  $\beta_2$ . With multiple choices, estimate a MNL model in the first stage. See a recent paper by Bourguignon, Fournier, and Gurgand, which corrects the formulas from Lee (1983).<sup>228</sup>*

A more general and widely used approach called **switching regression** assumes there are two (or more) regression functions and a discrete choice model determining which one applies. The typical estimation is similar to Heckman's  $\lambda$ .

**Example 19.11.** *Munich, Svejnar and Terrell study wages of Czech workers during 1991-1996. The workers either stay in post-communist firms or enter newly started private enterprises (de novo jobs, DNJ). Munich et al. obtain a positive coefficient on  $\lambda$  for the movers and a negative coefficient for the stayers.*

**Example 19.12.** *Engberg and Kim (1995) study the intra-metropolitan earnings variation: is it caused by person or place effects? Person  $i$  chooses his/her location  $j$  (inner city/poor suburb/rich suburb) based on his/her (latent) wage in each location  $w^*$  (think Roy model) and the location's amenities:*

$$U_{ij}^* = w_{ij}^* \gamma_j + x_i' \alpha_j + \nu_{ij}, \text{ where } w_{ij}^* = x_i' \beta_j + \varepsilon_{ij}.$$

$U^*$  is the latent utility of each location. Assuming that  $\varepsilon_{ij} \gamma_j + \nu_{ij}$  is iid logit, they look up the appropriate  $\lambda$  sample-selection formula and proceed to run switching wage regressions. The place effect is measured as  $\bar{x}'(\beta_{Suburb} - \beta_{City})$ . Actually, they present 6 different results based on what kind of control method they choose, starting with unconditional means. Assuming MNL iid error terms for the choice equation is not appropriate given their maintained assumption of no location

<sup>228</sup>This is implemented in the `selmllog` command in Stata. See also the `ssm` 'wrapper' command by Miranda and Rabe-Hesketh (2005).

effects for the most highly educated workers. They have no credible exclusion restrictions. Identifying the model off functional form blows it up. So they use non-linearities for identification: these come from non-parametric estimation of the selection equation. Do you find this credible? Finally, they run a semi-parametric selection (of location) model:

#### 19.2.4. Semiparametric Sample Selection

See Kyriazidou (1997) and Powell (1989). Assume  $d_i = 1\{x_i\gamma + v_{1i} > 0\}$ ,  $y_i = y_{i2} * d_i$ ,  $y_{i2} = x_i\beta + v_{2i}$  and assume that  $f(v_1, v_2)$  is independent of  $x_i$ . Then if I do not want to assume a particular form for the selection term (i.e., I am not willing to assume a particular distribution  $f$ ), follow Powell and choose person  $i$  and  $j$  such that  $x_i\gamma = x_j\gamma$  so that  $y_{i2} - y_{j2} = (x_i - x_j)\beta + \lambda(x_i\gamma) - \lambda(x_j\gamma) = (x_i - x_j)\beta$ . In practice do a Kernel on those pairs which are close, i.e., use an estimator such as

$$\left[ \sum K\left(\frac{(x_i - x_j)\hat{\gamma}}{n}\right) (x_i - x_j)(x_i - x_j)' \right]^{-1} \left[ \sum K\left(\frac{(x_i - x_j)\hat{\gamma}}{n}\right) (x_i - x_j)(y_i - y_j) \right] \quad (19.8)$$

**Example 19.13.** Return to Engberg and Kim and note the use of their maintained assumption as both a measuring stick for their different methods and as an identifying assumption for the semi-parametric sample selection estimation, where the constant is differenced out: Using  $\bar{x}'(\beta_{Suburb} - \beta_{City}) = 0$  for highly educated white males identifies the constant difference for other group.

Notice that what we are doing is essentially matching pairs of observations with the same probability of selection/participation, i.e. we are *matching on propensity score*.<sup>229</sup> Ahn and Powell further suggest that there is no need for any  $\gamma$  here and all can be done non-parametrically.

**Remark 181.** Of course, these models typically assume away heteroskedasticity, which is most likely to exist in large micro-data. Songian Chen uses symmetry assumption on  $f(v_1, v_2)$  to deal with heterogeneity of a particular form at both stages:  $f(|x) = f(v_1, v_2|\tilde{x})$  where  $\tilde{x}$  is a subset of  $x$ .

<sup>229</sup>See Rosenbaum and Rubin for early theoretical work on matching using the propensity score. We return to this issue in Section 20.2.

**Remark 182.** *An extension of the idea that the mean of the error term in the outcome equation for the selected sample is an invertible function of the selection probability from single-index to multiple-index models (i.e., with multiple choice in the ‘first stage’) is Dahl (2002, *Econometrica*). He provides a simple semiparametric correction for polychotomous selection models, where the outcome equation can be written as a partially linear model that depends on observed (first) choice probabilities.<sup>230</sup>*

The lesson from the example is that one should not attempt a problem of this type without an instrument (exclusion restriction). If there is an instrument, then a situation where we observe the outcome under both hypotheses allows for either selection-model-style switching-regression estimation or for simpler IV solutions. In the next section we will discuss the advantages and disadvantages of these two approaches. But note that the IV strategy is not available when we observe the outcome only for a subset of the individuals (one of the chosen alternatives).

**Remark 183.** *In absence of instrument, large support regressors have been proposed to estimate endogenous sample selection models. See IZA DP No. 12486 for references on this approach and for the Stata command `eqregsel`, which implements an estimation procedure that relies on neither exclusion restrictions nor large support regressors. Instead, the technique is based on the idea that if selection is endogenous, one can expect the effect of the outcome on selection to dominate those of the covariates for large values of the outcome (D’Haultfoeille et al., 2018). The method implements a series of quantile regressions in the tails of the outcome distribution (extremal quantile regressions).*

## 20. Program and Policy Evaluation

Consider estimating the impact of a treatment (binary variable): a medical procedure, a training program, etc. Evaluation of social programs is what much of true micro-econometrics is all about. We ask how to estimate the effect of a social policy (participation in a training program, change in college tuition) in absence of controlled experiments ([M]9.2.). How can we create counterfactual outcomes (such as what would have happened in absence of treatment)?<sup>231</sup>

<sup>230</sup><http://econweb.ucsd.edu/~gdahl/papers/multimarket-roy-model.pdf>

<sup>231</sup>Return to the first class of the course for a broad introduction to *causal inference*.

Methodological advances in program evaluation are important for all of cross-sectional econometrics.<sup>232</sup> First, the Heckman’s bivariate normal selection models dominated the field, but they were somewhat replaced by difference-in-differences models, which assumed that selection was based on time-constant unobservables. Recently, as data become richer, the matching method, introduced below, became popular; it focuses on controlling for observables, in contrast to previous approaches, which worried about selection on unobservables.

One of the main lessons of the recent literature, which we introduce in this section, is that the impacts of programs differ across individuals. Up to the discussion of the Roy model in Section 19.2.1, we estimated models where a given program (e.g., privatization) affected all participants in the same way. In fact, treatment effects often vary across individuals and these differences can be known to program participants or administrators such that they can act upon them, which makes those with the largest gains from the program treatment most likely to participate in the program.

At a fundamental level, we need to differentiate two types of problems: (1) *treatment effect problem*: What is the effect of a program in place on participants and nonparticipants compared to no program at all; (2) *structural problem*: What is the likely effect of a new program or an old program applied to a new setting. The latter problem is perhaps too ambitious and definitely requires more heroic assumptions. So focus on (1).

Crucially, we ask about *partial equilibrium* effects here; no answers given on across-board policy evaluation (such as making every student go to college) – no general equilibrium effects are taken into account!<sup>233</sup>

### 20.1. Setup of the Problem

Consider the effect of a (training) program within the Rubin Causal Model (Holland, 1986), which formulates causal questions as comparisons of potential outcomes:  $y_{1i} = \mu_1 + u_{1i}$  are earnings with training, and  $y_{0i} = \mu_0 + u_{0i}$  are earnings without training (go back to Roy’s model in Section 19.2.1).

Consider the population of eligible workers. They first choose to apply for the training program or not. We observe  $y_{1i}$  only when  $D_i = 1$  (the person applied

---

<sup>232</sup>See Imbens and Wooldridge (in press) “Recent Developments in the Econometrics of Program Evaluation”.

<sup>233</sup>See Ferracci, Jolivet, and van den Berg (2014, REStat) for an estimation method that identifies both individual treatment effects and effects of market-level distribution of assigned treatments (with share of treated leading to spillovers at market level).

for and took training) and observe  $y_{0i}$  only when  $D_i = 0$  (these are the so called eligible non-participants, ENPs).

### 20.1.1. Parameters of Policy Interest

We may want to know  $E[y_{1i} - y_{0i}] = \mu_1 - \mu_0$ , i.e. the treatment effect under random assignment, the *average treatment effect* (ATE), a.k.a. “the true causal effect.”<sup>234</sup> We may also want to know  $E[y_{1i} - y_{0i} | D_i = 1]$ , the *average effect of treatment on the treated* (ATT), which allows us to ask whether the current program, given the voluntary participation and all, is worth its costs or not. Finally, we may want to know the average effect on the untreated (ATU):  $E[y_{1i} - y_{0i} | D_i = 0]$ , which allows us to see whether we should expand the current program.<sup>235</sup>

Note that  $ATT = ATE + E[u_{1i} - u_{0i} | D_i = 1]$ , where the second term is the gain that the treatments obtain from treatment on top of the gain obtained by a random person. Of course, if the effect of treatment is the same for all individuals then  $ATE=ATT=ATU$ .

The fundamental problem: Consider estimation of ATT: The data only provides  $E[y_{1i} | D_i = 1]$  and  $E[y_{0i} | D_i = 0]$  but  $E[y_{0i} | D_i = 1]$  is the “what-if” counterfactual.

$$\begin{aligned} E[y_{1i} | D_i = 1] - E[y_{0i} | D_i = 0] &= ATT + \{E[y_{0i} | D_i = 1] - E[y_{0i} | D_i = 0]\} \\ &= \mu_1 - \mu_0 + E[u_{1i} - u_{0i} | D_i = 1] + \\ &\quad \{E[u_{0i} | D_i = 1] - E[u_{0i} | D_i = 0]\}. \end{aligned}$$

The sample selection bias (the term in curly brackets) comes from the fact that the treatments and controls may have a different outcome if neither got treated.

**Remark 184.** *We will mainly focus on the ATT and ATE, but one may be interested in other population parameters, for example in the voting criterion parameter, which asks about the fraction of the population that ex ante perceives a benefit from treatment (see Heckman, 2010 JEL; Cameron and Taber, 2004).*

<sup>234</sup>This is what biostatistics is after. Economists care about ATE when evaluating a mandatory program that affects for example all unemployed.

<sup>235</sup>Note that we may want to know the answer of the effect of expanding the program by 10% or by 30%. These two effects may differ if those who benefit more from the program are more likely to participate, all else equal.



### 20.1.2. Experimental Solution

Randomized Controlled Trials (RCTs), the golden standard of scientific evidence, provide an almost ideal (see end of Section) tool for measuring causal effects (think of double blind medical trials using placebo). Randomization ensures that on average the sample selection bias is zero:  $E[y_{0i}|D_i = 1] = E[y_{0i}|D_i = 0]$ .

Basic structure: Take the  $D = 1$  group and randomize into treatment ( $R = 1$ ) and control ( $R = 0$ ) group. Then construct the experimental outcome:  $E[y_{1i}|D_i = 1, R_i = 1] - E[y_{0i}|D_i = 1, R_i = 0]$ . This can be used as a benchmark for the accuracy of sample selection techniques that we need when we have no experiment.<sup>236</sup>

Results obtained from experiments are less controversial for general public than those based on non-experimental methodology and it may be harder to cheat when running an experiment compared to conducting a non-experimental study. However, experiments are (sometimes) costly and often socially unacceptable (in Europe). Further, randomization may disrupt the operation of the program in question and people may behave differently knowing they are in an experiment (think of temporarily expanding medical coverage). Some important treatment questions cannot be subject to experimental research design (family income while young, etc.).

**Remark 185.** *The ATT is an average! It can be positive when most of the treatment group gets hurt by treatment, but a few participants receive very large benefits.*<sup>237</sup> *One has to be careful with generalizations from RCTs as “it works on average” does not imply “it works for all”. More generally, randomized controlled trials only inform about the mean of the treatment effect and not about the distribution.*<sup>238</sup>

**Remark 186.** *Inference in experiments is not trivial. If we run  $y_i = \beta_0 + \beta_1 T_i + \epsilon_i$ ,*

<sup>236</sup>The LaLonde’s (1986) famous study, which showed that non-experimental methods are unable to replicate experimental results, motivated both further development of non-experimental methods and more use of social experiments. See McKenzie, Gibson, and Stillman (JEEA, 2010) for a recent example of such comparison.

<sup>237</sup>This is particularly a concern with small sample size when a few outliers can suggest a generalize-able positive effect.

<sup>238</sup>Note that both randomization and linearity are required for the ‘magic’ where we learn something about hypothetical outcomes without imposing any structure. In particular, randomization ensures that  $E[y_{0i}|R_i = 1] - E[y_{0i}|R_i = 0] = 0$  such that  $E[y_{1i}|R_i = 1] - E[y_{0i}|R_i = 0] = E[y_{1i}|R_i = 1] - E[y_{0i}|R_i = 1] + 0$  and thanks to expectation being a *linear* operator, this equals  $E[y_{1i} - y_{0i}|R_i = 1]$ , the ATT. We can’t find out about the median impact as the median of the difference is not the difference in medians.

using the randomized treatment dummy  $T$ , we can use the standard error only if we make a heteroskedasticity correction because treatment affects the variance of  $y$ . The difference in means divided by the standard error (approximately) has the Student's  $t$  distribution only when the number of treatments and controls is the same. See Duflo, Glennerster and Kremer (2008) for a discussion, which is useful also for determining the size of the experimental sample that should allow for detection of at least an effect that would justify the costs of running the program.<sup>239</sup>

**Remark 187.** *Selecting what covariates to condition on is not trivial in RCTs either. See the LOOP estimator by Wu and Gagnon-Bartsch and Jann Spiess (2018) job market paper.*

**Remark 188.** *There are other inference issues. Asking ex post, after conducting an experiment, if one group benefits more than another smells of data mining (there is no way to ex post adjust inference).<sup>240</sup> Looking at multiple outcomes and finding the one that “works” is another example (although there is a way to adjust inference—see the Bonferroni's correction for multiple comparisons<sup>241</sup>). Successful replication of “what works” experiments is more likely if we understand the mechanism of the treatment effect. Controlling for other  $X$ s may introduce finite sample bias (because of the correlation with treatment that will be there in small samples even if the population correlation is zero) but may improve precision if these covariates predict  $y$  for both treatments and controls.<sup>242</sup>*

<sup>239</sup>See IZA DP no. 8583 for a Stata simulation program helping one to determine the size of experimental samples depending on the research design. See also IZA DP No. 9908 for more advice.

<sup>240</sup>But see Abadie, et al. REStat (2018).

<sup>241</sup>See Romano and Wolf (2005, JASA) for a bootstrapped test with multiple outcome measures. The (Stata implementation of the) Romano-Wolf correction is (asymptotically) more powerful than earlier multiple testing procedures such as the Bonferroni and Holm corrections since it takes into account the dependence structure of the test statistics by resampling from the original data (see IZA DP No. 12845).

<sup>242</sup>There are larger issues as well. Some argue that most published research findings (in all sciences) are false as a result of the combination of the likely extent of publication bias (see note n. 4.2), power of tests, number of independent studies on a topic, and the 5% statistical significance level used to claim a research discovery (<http://www.edwardtufte.com/files/Study1.pdf>). This seems to be confirmed by the practice of pharmaceutical companies halting most target-validation projects where they are not able to confirm in their in-house experiments the results of published studies (<http://marginalrevolution.tumblr.com/post/9960121369/how-good-is-published-academic-research>).

**Remark 189.** *Even with experimental data, there are often problems of selectivity. Double blinding is not feasible in social sciences. First, worry about (obviously very selective) non-compliers: those with  $R_i = 1$  who do not show up for training and those with  $R_i = 0$  who get a similar treatment outside of the program. (Those who volunteer for the program but are randomized out of treatment may get treated elsewhere or may refuse to provide data.)<sup>243</sup> Second, dynamic sample selection may arise even in presence of initial randomization: see Ham and LaLonde (1997) for a duration study.*

**Remark 190.** *Even more fundamentally, RCTs solve bias and selection issues for the sub-population (of villages, schools, patients) on which the experiment is administered. They may suffer from selection issues in the process of selection of the experimental panel, either through willingness to participate, etc., or through small samples (100 villages is actually a small sample) and underestimated standard errors. In contrast, observational (or IV) studies often use almost entire populations, such that these methods can be more than just ‘RCT wannabes’ and are not automatically worse than RCTs. In general, RCTs should be subject to the same scrutiny as observational studies.<sup>244</sup>*

If we do not have an experiment, how can we deal with the fundamental estimation problem? Two answers are proposed, one when we have the key control variables (and do not have/need an instrument for treatment, 20.2), the other when we do have an exclusion restriction (which is needed because we do not have the key control variables, 20.3). The distinction between the two strands of methods is a familiar one: in regression analysis we either ran GLS or IV. In general, which of the two strands of methods one uses depends on whether the conditional independence assumption holds, i.e., whether potential outcomes are independent of treatment assignment conditional on a set of observable covariates.

---

<sup>243</sup>In those cases, one can exploit the random assignment and estimate the “intention-to-treat” (ITT) parameter. However, if we use the original treatment assignment as an IV for actual treatment, we’re back in the LATE IV world (see below), where one needs to model the behavior (non-compliance), i.e., the heterogeneity in the response to assignment.

<sup>244</sup>See Angus Deaton’s presentation on RCTs (with a comparison to Trial and Error as a method for figuring out what works) at <http://nyudri.org/events/past-events/annual-conference-2012-debates-in-development/>

## 20.2. Matching

This is a non-experimental solution based on controlling for enough observables such that selection on unobservables is (assumed) irrelevant.<sup>245</sup> Consider estimation of the ATT. When can we use  $E[y_0|D = 0]$  as a surrogate for the counterfactual  $E[y_0|D = 1]$ ? If  $E[y_0|D = 1, X] = E[y_0|D = 0, X]$  then there is hope. This condition (assumption) says that selection occurs based solely on observables, i.e. that “conditioning” on  $X$  (within each  $X$  cell) assignment to treatment is random. This is similar to standard exogeneity assumptions in regressions.<sup>246</sup>

**Example 20.1.** *An example of selection on observables is the following process of admission to university: Each applicant is ranked based on the number of “points” derived from observable characteristics. Admitted students are a random sample from those who pass the admission threshold (Angrist and Krueger, 1999).*

The assumption of unconfoundedness given  $X$  is that  $(y_{0i}, y_{1i}) \perp D_i | X_i$ . We assume that the treatment status is conditionally mean independent from the potential outcomes. This is also called the *Conditional Independence Assumption* (CIA). See Cochran and Rubin (1973) or Rosenbaum and Rubin (1983).<sup>247</sup>

To estimate only ATT (as opposed to ATE), one needs only a weaker version of the unconfoundedness assumption, namely that  $y_{i0} \perp D_i | X_i$ .<sup>248</sup> This assumption allows us to get at the counterfactual because

$$E[y_0|D = 1] = E\{E[y_0|D = 1, X]|D = 1\} = E\{E[y_0|D = 0, X]|D = 1\}$$

so we estimate the last term by

$$\frac{1}{N_1} \sum_{i \in \{D=1\}} \hat{E}[y_0|D = 0, X = x_i].$$

---

<sup>245</sup>The assumption is that a sufficiently detailed (non-parametric) conditioning on observables makes the sample selection problem go away.

<sup>246</sup>A recent review of matching techniques is available here:

<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2943670/pdf/nihms200640.pdf>

<sup>247</sup>Here are some sources for various names of this assumption: “unconfoundedness” or “ignorable treatment assignment” (Rosenbaum & Rubin, 1983), “selection on observables” (Barnow, Cain, & Goldberger, 1980), “conditional independence” (Lechner 1999, 2002), and “exogeneity” (Imbens, 2004).

<sup>248</sup>Note that the identifying assumptions motivating the panel-data fixed effect model of Section 7.1 are  $E[y_{0it}|\alpha_i, D_{it}, X_{it}] = E[y_{0it}|\alpha_i, X_{it}]$  together with a linear model for  $E[y_{0it}|\alpha_i, X_{it}]$ .

We can estimate a regression using  $D = 0$  data but predict outcome using  $D = 1$ . Alternatively, proceed non-parametrically and use matching to estimate the ATT: (i) for each value of  $X$  (combinations of  $x_k$  values) compute the difference in  $y$  between treatments and controls, (ii) average these differences across treatments'  $X$ . This is similar to a non-parametric (kernel) regression (see Section 13). The idea is to resemble an experiment—to make the distribution of observed determinants balanced across the two groups that we compare. Matching is a process of ‘re-building’ an experimental data set.

But what if the  $X$  support of  $E[y_1|D = 1, X]$  and  $E[y_0|D = 0, X]$  does not coincide? We cannot predict out of sample in terms of  $X$  for  $D = 0$  non-parametrically. Lack of *common support* will occur frequently if the dimension of  $X$  is high and is natural when comparing treatments and controls who are not randomly assigned.

**Remark 191.** *For the matching estimation of ATT, we need the presence of an ‘analogue’ control for each treated; it does not matter if there are controls with no ‘matched’ treated individuals. We drop those treated who have no ‘matched’ controls. So if impacts vary across people, this means that estimates of ATT by methods that do not drop any observations (i.e., regressions) will have different population analogues.*

**Remark 192.** *DiNardo and Lee (2010) argue that adding more controls in selection-on-observables matching exercises carries a risk of exacerbating selection biases.*

**Remark 193.** *See Duleep (2012, IZA DP no. 6682) for an efficiency argument supporting the use of matching when analyzing natural experiments (specifically, matching before and after observations within the treatment and control groups).*

**Matching on Propensity Score (Index Sufficiency)** However, ‘exact’ matching on the combination of all  $X$  values is not practical with high-dimensional  $X$  because there will be little data in each ‘data cell’ corresponding to a particular  $X$  combination.<sup>249</sup> We deal with the curse of dimensionality (of matching based on a high-dimensional  $X$ ) by conditioning only on a scalar: the propensity score

<sup>249</sup>Also see Abadie and Imbens (2002) on the bias of simple matching estimators when  $X$  dimension is high. Angrist and Hahn (2004) show that even though there is no asymptotic gain in efficiency from using the propensity score as opposed to matching on covariates, there will likely be a gain in efficiency in finite samples. Using the propensity score essentially corresponds to applying prior knowledge to reduce dimensionality and this will improve precision in small samples.

$p(X) = E[D|X] = \Pr[D = 1|X]$ . We match on  $p(X)$  over the common support – compare the outcome for those individuals (ENPs compared to treatments) with similar probability of participation in the program.<sup>250</sup>

Rosenbaum and Rubin (1983) show that this works, i.e., that if  $y_1, y_0 \perp D|X$  then  $y_1, y_0 \perp D|p(X)$  using the fact that  $D \perp X|p(X)$  (that is  $\Pr[D = 1|X, p(X)] = \Pr[D = 1|p(X)] = p(X)$ ).

In practice, matching on  $P(X)$  is not exact so we stratify the propensity score or do a Kernel (see Section 13) or Nearest Neighbor (see Section 13) non-parametric regression.<sup>251</sup> The choice of the matching method can make a difference in small samples (see Heckman, Ichimura and Todd, 1997). Standard errors are typically bootstrapped, although this has been shown to be invalid.<sup>252</sup> One of the most widely used methods today is dif-in-difs for matched data.<sup>253</sup>

Matching on  $P(X)$  shifts attention from estimating  $E[Y|X, D]$  to estimating  $P(X) = E[D|X]$ , which may be attractive when we have a better motivation or a model for the selection into treatment than for how the treatment operates. We often make parametric assumptions when estimating the Pscore (otherwise, it is not clear that the curse of dimensionality would be reduced).

How much matching reduces the bias is an empirical question.<sup>254</sup> It depends

---

<sup>250</sup>So always start by plotting the pscore for the treatment and control group in the same graph. See Caliendo and Kopeinig (2008) for a survey of practical issues with implementing matching strategies. Crump, Hotz, Imbens, and Mitnik (2009, *Biometrika*) suggest that for a wide range of distributions a good approximation to the optimal rule is provided by the simple selection rule to drop all units with estimated propensity scores outside the range  $[0.1, 0.9]$ .

<sup>251</sup>Stata programs to estimate treatment effects are available from Becker and Ichino (`att*`, 2002), Leuven and Sianesi (`psmatch2`, 2003) and Abadie et al. (`nnmatch`, 2004). See also notes by Andrea Ichino at <http://www.iue.it/Personal/Ichino/> Guido Imbens provides examples of matching programs for his re-evaluation of the LaLonde AER data at <http://emlab.berkeley.edu/users/imbens/estimators.shtml>

<sup>252</sup>Abadie and Imbens (2009, NBER WP No. 15301) provide the large sample distribution of propensity score matching estimators and a variance adjustment that corrects for the fact that the p-score itself is first estimated.

<sup>253</sup>In these exercises, one often wishes to match exactly on a few discrete variables (say industry and year) and within those groups matching is based on pscore. For a trick on how to do this in `psmatch2`, see <http://www.stata.com/statalist/archive/2010-09/msg00073.html>

<sup>254</sup>Matching is the *estimator du jour* in the program evaluation literature (Smith and Todd, 2004). The meta-analysis of active labor market policy program evaluations by Card, Kluve and Weber (2009, IZA DP no. 4002) suggests that “experimental and non-experimental studies have similar fractions of significant negative and significant positive impact estimates, suggesting that the research designs used in recent non-experimental evaluations are unbiased.”

on whether matching on observed  $X$ s balances unobserved determinants or not. See Heckman, Ichimura, Smith, and Todd (1998), Angrist (1995) or Dehejia and Wahba (1998). One of the lessons from this literature on the effects of training programs on the labor market is that one should match on sufficiently lagged pre-treatment performance (Ashenfelter's dip). To the extent that past performance and current  $X$ s do not capture individual's motivation or skills, the bias remains.

**Remark 194.** *One cannot directly test (reject) the unconfoundedness assumption, but in cases that there are two distinct control groups (as in Heckman, Ichimura and Todd, 1997), ideally two groups with different potential sources of bias, one can test whether the matched 'causal' effect estimate of being in one group as opposed to the other is zero as it should be.*

**Remark 195.** *Another related test is to estimate the 'causal effect' of treatment on lagged outcome (for example wages before training). If this 'effect' is estimated to be close to zero, the unconfoundedness assumption is more plausible. The logic of the Heckman and Hotz (1989) pre-program test is to look at  $F$  tests of the joint significance of dummies for treated units in the years prior to treatment (Ashenfelter's dip).*

**Remark 196.** *One can also quantify the sensitivity of the ATT estimate to the CIA assumption. Rosenbaum (2002) provides bounds on the degree of the departure from the unconfoundedness assumption that would make the estimated ATT insignificant.<sup>255</sup> A critical value of this test of 1.15 means that the confidence interval for the estimated effect would include zero if an unobserved variable caused the odds ratio of treatment assignment to differ between the treatment and comparison groups by 1.15. For another strand of analysis of sensitivity of matching to the Conditional Independence Assumption, see Ichino et. al. (2006) or Nannicini (2006).*

**Remark 197.** *The balancing property holds given the true propensity score (Rosenbaum and Rubin, 1983), but no such result exists for the estimated Pscore. Hence, to check that matching does its job, we need to use  $t$  tests to make sure that there are no significant differences in covariates means for both groups after matching, i.e., to test for balancing given Pscore. See Rosenbaum and Rubin (1985).*

---

<sup>255</sup>In Stata, `rbounds` gives the bounds for continuous variables and `mhbounds` works for binary outcome variables.

**Remark 198.** Another test for matching quality is in Smith and Todd (2005): a kernel weighted regression of the probit independent variables on a quartic of P-scores (interacted with treatment dummies) to show F tests of joint significance of the interaction terms. Also, see the Hotelling T2 test of the joint null of equal means of all the probit independent variables. In the ideal scenario, matching and re-weighting reduces differences across treatment/control, the tests (t, F, T2, H-H) are not significant.

**Remark 199.** One can implement the p-score matching by weighting observations using  $P(X)$  to create a balance between treatment and control groups.<sup>256</sup>

**Remark 200.** When the p-score cannot be consistently estimated because of choice-based sampling designs with unknown sampling weights, Heckman and Todd (2009, IZA DP No. 4304) show that the selection and matching procedures can be implemented using propensity scores fit on choice-based samples with misspecified weights, because the odds ratio of the propensity score fit on the choice-based sample is monotonically related to the odds ratio of the true propensity scores.

**Remark 201.** See IZA DP. No. 4841 for p-score-based estimation of distributional effects of interventions (i.e., not just of average effects).

**Remark 202.** See the coarsened exact matching method for a practical version of matching with several desirable features.<sup>257</sup>

**Remark 203.** IZA Discussion Paper No. 5140 is a Monte Carlo study that evaluates the relative performance of matching and other estimators based on the CIA assumption in presence of various types of measurement error. Kernel matching with a relatively large bandwidth appears more robust than other estimators.

**Matching with Multiple Treatments** What if the treatment is not binary? Multiple-treatment evaluation has been developed by Imbens (2000) and Lechner (2001). Let  $Y_i^0, Y_i^1, \dots, Y_i^K$  correspond to outcomes with  $K$  different intensities of

<sup>256</sup>See <http://www-personal.umich.edu/~jdinardo/bztalk5.pdf> and Nichols (2008, Stata Journal). Busso, DiNardo, and McCrary (2014, REStat) argue that reweighting is competitive with matching when the common support problem is not large.

<sup>257</sup>See <http://gking.harvard.edu/files/cem-plus.pdf> and the relevant software at <http://gking.harvard.edu/cem/>



treatment or no treatment ( $k = 0$ ). Let  $T_i = k$  denote the actual occurrence. Now, define ATT of treatment  $k$  using all the  $K + 1$  pairwise comparisons

$$E[Y^k - Y^{k'} | T = k] = E[Y^k | T = k] - E[Y^{k'} | T = k] \text{ for } k \in \{0, 1, \dots, K\}, k \neq k'$$

and estimate the second term (counterfactual) by  $E_X[E(Y^{k'} | T = k', X) | T = k]$ . Note that there is a lot of “common support” conditioning here (find a match for all  $k'$ ) so one really needs propensity score matching, which is estimated separately for all combinations of  $k$  and  $k'$  (ignoring the other groups).<sup>258</sup> The conditional independence assumption must hold for each of the  $k - k'$  comparison.

When treatment is continuous, one can apply the generalized Pscore technique introduced by Hirano and Imbens (2004) and Imai and van Dyk (2004). Follow Du and Girma (2009) implementation details. It all starts with estimating the fractional logit model (Papke and Wooldridge, 1996).

**Matching and DiDs** The estimator *du jour* corresponds to an approach suggested by Heckman, Ichimura, and Todd (1997) that makes the CIA assumption more likely to hold: they combine matching with diff-in-diffs, i.e., compare the *change* in the outcome variable for the  $k$ -treated groups with the *change* in the outcome variable for all non- $k$  groups, thus purging time invariant unobservables from the specification based on the familiar common trend assumption. Of course, this is easily and often done for binary treatment:

$$\frac{1}{N_1} \sum_{i \in I_1 \cap S_p} \left\{ (y_{1it} - y_{1it-1}) - \sum_{j \in I_0 \cap S_{pi}} W(i, j) (y_{0jt} - y_{0jt-1}) \right\},$$

where  $I_1$  is the set of treated individuals (with  $D_i = 1$ ) and where  $S_p$  represents the overall common support and  $S_{pi}$  is the set of non-treated individuals matched to a treated  $i$ . Finally,  $W(i, j)$  are either kernel weights corresponding to  $|P_i - P_j|$  or LLR (lowess) weights (see Section 13). P-score conditioning is still based on pre-treatment  $X$ s as post-treatment data would violate the CIA. See also the discussion in Blundell and Costa Dias (2000). This approach combines careful conditioning on observables through matching on pre-treatment performance (trends) with before/after differencing that eliminates time-constant unobservables. However, bootstrap is invalid for matching so inference could be based on

<sup>258</sup>Sianesi (2001) provides the Stata 8 code.

asymptotic theory for matched DiDs derived in Heckman, Smith, and Clements (1997, REStud).<sup>259</sup>

**Regression or Matching?** What if, instead of matching, we run a regression of  $y$  on  $D$  controlling for  $X$ ? What is the difference? First, the regression approach imposes functional form (linearity) over the common support area while matching is non-parametric. Regressions also use their functional form to work off the common support, i.e., use areas of  $X$  with non-overlapping support of treated and controls, which can be highly misleading.

But what if we include a full set of interactions among  $x$ s to approximate a non-parametric solution (i.e., use a fully saturated model)? This regression still differs from the matching ATT estimator in the implicit weighting scheme: Matching gives more weight to  $X$  cells with high probability of treatment  $P(X)$  (cells with high share of treated) and 0 weight to cells where  $p(X) = 0$ . Regression gives more weight to  $X$  cells where proportion of treated and untreated is similar, i.e., where the conditional variance of treatment is largest. (See Remark 21 on implicit weighting in regressions and Angrist, 1998). If those who are most likely to be selected (into training, the military, etc.) benefit less from the treatment (because they have the best earnings potential), then matching will give smaller effects than regression.

In the end, what matters most is not whether we run a regression or a matching exercise, but that we inspect the data on being balanced. Pscore estimation is useful for this. Recently, Crump, Hotz, Imbens, and Mitnik (2009) suggest that one select data on common support using propensity score and then run regressions on such balanced data. Further, inference in matching may be less standardized than in regressions. This approach is becoming popular.

**Example 20.2.** *Brown and Earle (2010, Growth Effects of Small Loan Programs) estimate the Pscore, then they manually match on Pscore within exactly matched groups and thus select their common support. Next, they run a panel regression (with bootstrapped std. errors and whilst weighting by the number of treatments in the group to get at ATT). They argue the regression technique adds efficiency.*

In short, one can equivalently either run a firm-fixed-effect & firm-time-trend panel-data regression with ATT-like weighting on the common support or match

---

<sup>259</sup>The `diff` package in `Stata` implements the HIT PS-matched DiDs (Villa, 2016, *Stata Journal*), but it assumes that all units are treated at the same time. Running a parametric DiDs on the matched sample is a simple alternative.

on pre-treatment trends in  $y$  and study the post-treatment change in  $y$ ; the two approaches will differ if much of the firm-specific time trend is estimated off post-treatment panel data.

**Remark 204.** *Arkhangelsky and Imbens (2019) first link the linear removal of average treatment and average covariate values across groups in a linear FE to weighting by the inverse of the propensity score and then generalize the propensity-score FE approach to non-linear removal of covariate values and for other group characteristics.*

### 20.3. Local IV

Now start worrying about unobservables again. Suppose, that we have instruments (affecting choice, but excluded from the  $y$  equation). There is an important contrast between IV methods and sample selection methods:

What are the *policy parameters of interest*? What do we want to know? ATE, ATT, ATU, or a treatment effect related to a specific new policy – affecting a specific subset of the population? If the effect of a given policy differs across parts of population (parameter heterogeneity) we need to be able to estimate the effect on each part. There are fundamental reasons why the treatment effect may differ across individuals.

#### 20.3.1. Local Average Treatment Effect

Imbens and Angrist (1994) prove the LATE theorem, which provides interpretation for IV estimates when treatment effects are heterogenous in the population.<sup>260</sup> If the effect of  $x$  on  $y$  varies in the population, then it can be shown that IV estimates are weighted averages of these group-specific effects where higher weight is given to those groups whose  $x$  is better explained (predicted) by the instrument (see Remark 21 and 74). So the IV estimate is the treatment effect on specific groups—it is a “local” effect. Specifically, a local average effect of the treatment on

---

<sup>260</sup>The assumptions underlying the theorem are IV exogeneity and exclusion restriction, as usual, and also “monotonicity”, which requires that anyone who would get treated if not induced so by the IV would also get treated if the IV is “pushing” for treatment. Alternatively called the “uniformity” assumption, it says that if the IV makes some agents enter treatment, it must not induce any agents to leave treatment. The assumption applies across people, not within a person across different values of the instrument (see Heckman, Urzua, and Vytlačil, 2006). For an introduction to the topic, see <http://www.irs.princeton.edu/pubs/pdfs/415.pdf>

those who change state (treatment status) in response to a change in the instrument (the *compliers*).<sup>261</sup> It is not informative about the effect for the never-takers or always-takers, i.e., those for whom the IV value does not help to predict treatment status, or for the defiers—those who were assigned to treatment (should get treated given their IV value and the way the first stage generally works), but did not get treated.

**Example 20.3.** Angrist and Krueger (1991) use quarter of birth and compulsory schooling laws requiring children to enrol at age 6 and remain in school until their 16th birthday to estimate returns to education. This approach uses only a small part of the overall variation in schooling; in particular, the variation comes from those who are unlikely to have higher education.

**Example 20.4.** Similarly, one may think of the Angrist (1990) estimate of the effect of military service as corresponding to the effect of the service on those drafted using the Vietnam-era lottery, but not those (majority) soldiers who volunteered.

**Remark 205.** Given that LATE identifies the effect on compliers, it is important to know as much as possible about the compliers, who, however, cannot be individually identified. One can, however, to describe the distribution of compliers' characteristics using variation in the first stage across covariates groups (p.171 of Angrist and Pischke, 2009). For example, are college graduates ( $x_{1i} = 1$ ) more likely to be compliers? Just look at the ratio of the first stage for them relative to the overall first stage:

$$\frac{E[D_i|Z_1 = 1, x_{1i} = 1] - E[D_i|Z_1 = 0, x_{1i} = 1]}{E[D_i|Z_1 = 1] - E[D_i|Z_1 = 0]}.$$

**Remark 206.** Note, that this is a general problem of all estimation. The only difference is that IV selects a specific part of variation (we know who identifies the effect) whereas OLS can be thought of as weighted average of many sources of variation, some potentially endogenous. In particular, we started the course by focusing on linear regressions, which provide the best linear predictor conditional on  $X$  when  $E[y|X]$  is not linear (such that there is misspecification error in the linear model residual), but in this case OLS coefficients (with robust standard errors) depend on the distribution of  $X$ , which makes the approach less attractive.

<sup>261</sup>DiNardo and Lee (2010) clarify that the “probabilistic monotonicity” assumption allows the LATE individual-specific weights to be interpreted as corresponding to the effect of the instrument on the *probability* of treatment.

**Remark 207.** *Blandhol et al. (NBER WP 29709) ask whether the LATE interpretation actually applies to the types of TSLS specifications that are used in practice. They advocate LATE specifications that are “saturated”, i.e., control for covariates non-parametrically.*

There is a strong analogy between IV, which we assume to be exogenous, but which does not give an  $R^2$  of 1 in the first stage, and a randomized experiment with non-perfect compliance with assigned status is important. In situations with no defiers, ATT is a weighted average of the effect of treatment on the compliers and on the always-takers, but LATE IV only identifies the effect on compliers. LATE will give ATT only when there are (almost) only compliers. IV will only give ATE under special conditions; for example, when something that was a matter of choice becomes legally binding such that everybody must comply (by law) and there are no never- or always-takers.<sup>262</sup> However, ATE can be viewed as an extrapolation that has LATE as its “leading term” (DiNardo and Lee, 2010).

**Remark 208.** *Estimating the so-called the reduced form (see Remark 87) of the outcome variable on the randomly assigned offer of treatment (the IV) gives the so-called intention to treat (ITT) parameter. LATE then equals ITT divided by the compliance rate (the first stage associated with the IV) based on the indirect least squares argument (see Example 9.1).*

**Remark 209.** *Note that because different instruments estimate different parameters, overidentification tests discussed in Section 9.1 are out the window. LATE overturns the logic of the Durbin-Wu-Hausman test for overidentification (see also note n. 271). Within the LATE framework, the relationship between OLS and IV estimates is completely uninformative about the existence of selection. One thus usually corrects for selection with no clear evidence to demonstrate selection into treatment (endogeneity) exists (and then one ends up with noisy 2SLS effects). Black et al. (IZA DP No. 9346) offers a solution: a simple test for the presence of selection bias (actually a necessary condition for the absence of selection): Conditional on covariates, they compare the mean outcomes of non-treated compliers to never-takers. They also compare outcomes for treated compliers to always-takers. I.e. they run regressions within treated and control groups of outcomes on covariates  $X$  and instruments  $Z$  (probability of selection). For the treated, they*

---

<sup>262</sup>See Manski (1990, 1996, 2003) for an approach that leads to bounds for ATE given assumptions one is willing to make.

run  $E[Y_{1i}|X_i, Z_i, D_i = 1] = g_1(X_i) + \alpha_1 Z_i$ . Here,  $\alpha_1 = 0$  implies not rejecting  $H^0$  of  $Y_{1i} \perp D_i | X_i$ .<sup>263</sup> Actually, finding that either  $\alpha_1 \neq 0$  or  $\alpha_0 \neq 0$  is evidence of either selection or violation of the exclusion restriction or both. They also provide a formula for the magnitude of the selection effect.

**Remark 210.** For treatment evaluation problems, in which there is nonrandom sample selection/attrition (that contaminates the original experiment or IV manipulation), Lee (2009, REStud) proposes a trimming procedure for bounding ATEs in the presence of sample selection. Confidence intervals for these bounds can be constructed following the procedure described in Imbens and Manski (2004, Econometrica).

LATE is a treatment effect at the *margin of participation* (in treatment) *relative to the instrument*. Suppose that your instrument must have an economically small effect, such as the presence of a small fee for training materials affecting the choice to participate in a training program. Then the LATE corresponds to people who are just about indifferent between participation or not.<sup>264</sup>

In *regression discontinuity* designs (see Section 9.1), one estimates a treatment parameter, but also only a local one—for those who are at the regression discontinuity and only for compliers.

**Example 20.5.** Costa Dias, Ichimura and van den Berg (2008, IZA DP no. 3280) propose a simple matching-and-IV method based on sharp discontinuity in treatment probability. They match on  $X$  and then apply a correction based on an IV  $Z$  that shifts  $P(X)$  to zero for values of  $Z$  above a certain threshold. In their case, the discontinuity arises thanks to eligibility rules for treatment: offer no training to unemployed workers above some age level. Just re-do matching (with kernel weights) but apply additional weights based on the distance from age eligibility cutoff.

The correction is based on the following argument: Assume that  $y_0 \perp Z | X$  and  $P[D = 0 | X, Z = z^*] = 1$ . Then  $E[y_0 | X] = E[y_0 | X, Z]$  and

$$E[y_0 | X, Z] = E[y_0 | X, Z, D = 0]P[D = 0 | X, Z] + E[y_0 | X, Z, D = 1]P[D = 1 | X, Z]$$

<sup>263</sup>The test looks for evidence of a non-constant control function.

<sup>264</sup>Vytlacil (2002) shows that LATE is equivalent to a non-parametric generalized Roy model, see Remark 217 below.

and at  $Z = z^*$  this implies that  $E[y_0|X] = E[y_0|X, Z = z^*, D = 0]$ . Next,

$$\begin{aligned} E[y_0|X, D = 1] &= \frac{E[y_0|X] - E[y_0|X, D = 0]P[D = 0|X]}{P[D = 1|X]} \\ &= \frac{E[y_0|X, Z = z^*, D = 0] - E[y_0|X, D = 0]P[D = 0|X]}{P[D = 1|X]} \\ &= E[y_0|X, D = 0] + \frac{E[y_0|X, Z = z^*, D = 0] - E[y_0|X, D = 0]}{P[D = 1|X]}, \end{aligned}$$

where the second term will equal zero if the conditional independence assumption that underpins standard matching holds (this can be tested).

**Example 20.6.** For an application of a wide array of matching and non-parametric IV estimators, see Ham, Li and Reagan (2009) migration study, which applies, among others, the Frölich (2007) non-parametric LATE estimator with covariates.

**Remark 211.** Mogstad et al. (2021, AER) inspect the monotonicity condition of the Imbens and Angrist (1994) LATE theorem and show that with multiple IVs the condition can only be satisfied if choice behavior is effectively homogeneous and can be at odds with data. So they consider a weaker partial monotonicity condition.

### 20.3.2. Marginal Treatment Effect

A central goal of economic analysis is to quantify marginal returns to a policy. In the Roy model of Section 19.2.1, the margin corresponds to people who are indifferent between hunting and fishing (treatment or no treatment). Recent work by Heckman and coauthors (summarized in Heckman, 2010 JEL) uses the concept of the marginal treatment effect (MTE), introduced by Bjorklund and Moffitt (1987), to provide a unified perspective on the IV and sample selection literature.<sup>265</sup>

To start from scratch: there are two possible states of the world:

(a) There is one “true” effect of  $D$  on  $y$ , namely  $\beta$ . There may be correlation between  $D$  and unobservables so we run IV for  $D$  in estimating:<sup>266</sup>

$$Y_i = \alpha + X_i\theta + \beta D_i + U_i.$$

<sup>265</sup>See [http://jenni.uchicago.edu/underiv/userguide\\_march\\_22\\_2006.pdf](http://jenni.uchicago.edu/underiv/userguide_march_22_2006.pdf)

<sup>266</sup>Alternatively, write the common effect model as  $Y_0 = \alpha + U$ ,  $Y_1 = \alpha + \beta + U$ .

Note that in this world, the mean marginal and average effect is the same for all people with the same  $X$ .

(b) Return to Roy and Section 20.1. Individuals know (estimate) their *varying* returns from choosing  $D$  and *act upon the size of the return*:

$$\begin{aligned} Y_{i0} &= \alpha_0 + X_i\theta_0 + U_{i0} \\ Y_{i1} &= \alpha_1 + X_i\theta_1 + U_{i1}, \end{aligned}$$

so that the causal effect

$$\beta_i = Y_{i1} - Y_{i0} = \alpha_1 - \alpha_0 + X_i(\theta_1 - \theta_0) + U_{i1} - U_{i0}. \quad (20.1)$$

There is, therefore, a *distribution* of returns (correlated random coefficients (CRC), ex post causal effects) that cannot be summarized by one number  $\beta$ .<sup>267</sup> Individuals are likely to act on their knowledge of the size of their own gain (effect).<sup>268</sup>

To complete the model (similar to Heckman's  $\lambda$  switching regressions), introduce an "instrument"  $Z$  and postulate the selection equation:

$$\begin{aligned} D_i^* &= \theta_D Z_i + U_{iD} \\ D_i &= 1 \text{ if } D_i^* \geq 0. \end{aligned}$$

Next, we define a new parameter: the marginal treatment effect  $MTE(x, u_D) = E[\beta|X = x, U_D = u_D]$  as the effect of  $D$  on a person with values  $(x, u_D)$  that is just indifferent between choosing  $D = 1$  and 0.<sup>269</sup> It is a willingness to pay measure for those at the margin of indifference given  $X$  and  $U_D$ .

The typical range of policy parameters (ATT, ATE, etc.) can be expressed as differentially weighted averages (integrals over population) of the MTE. IV and OLS can also be expressed as weighted averages of MTE, but the weights are not those of ATT, ATE, etc. IV weights are related to the type of instrument (LATE

<sup>267</sup>In the standard IV literature (a), we worry about (i) unobservable heterogeneity bias ( $COV(D, U) \neq 0$ ), (ii) downward measurement error bias ( $COV(D, U) \neq 0$ ), and (iii) the weak instrument bias, where  $COV(U, D)/COV(IV, D)$  is large because  $COV(IV, D)$  is small. Here, under (b), there are more econometric problems:  $COV(D, U_0) \neq 0$  as before, but also  $COV(\beta, U_0) \neq 0$  and crucially  $COV(\beta, D) \neq 0$ .

<sup>268</sup>For a test of whether subjects act on the knowledge of gains, see NBER Working Paper No. 15463.

<sup>269</sup>Start with writing down the LATE for the case of a binary IV (values  $z$  and  $z'$ ) and let  $z'$  go to  $z$ . This is MTE. The people who are affected by such a small change in IV are indifferent between the two choices.



interpretation—LATE estimates the mean return at the margin defined by the instrument manipulation). Heckman et al. conclude that while IV estimation may be more statistically robust compared to sample selection methods, IV may often not answer any economically interesting questions.

**Remark 212.** *CRC models where adoption of technology or investment in education are correlated with the return to such actions are typically identified using an IV (Heckman and Vytlacil, 1998; Wooldridge, 2003). Suri (2011, Econometrica) instead develops an identified structural model, which generalizes the Chamberlain (1984) fixed effect approach and relies on the method of minimum distance (see Section 8.3.2).<sup>270</sup>*

**Example 20.7.** *Consider an application of this approach to estimating the wage returns to education (Carneiro, Heckman and Vytlacil, 2002) where  $D \equiv S \in \{0, 1\}$  is college or high school. Either (a) human capital is homogenous (Griliches, 1977) and there is one  $\beta$  or (b) human capital is heterogeneous and people know their returns from  $S$  when making schooling decisions (Willis and Rosen, 1979).*

To deal with ability and measurement error biases, Card uses college proximity as an IV (see example 9.2). Typically  $\beta_{IV} > \beta_{OLS}$ . Now think of the LATE IV interpretation:  $\beta_{IV}$  is the effect of college on wages for those people whose college participation is affected by whether or not they grow up near college – these are students from low income families. Card therefore interprets  $\beta_{IV} > \beta_{OLS}$  as saying that students from low income families have high  $\beta$ , but don't attend college because of credit constraints. The instruments affect those with high MTE values disproportionately. However, these cost constraints must be really strong, because if we think of OLS as the average effect and IV as the marginal effect, one would expect the marginal student (affected by instrument, typically cost-related IV) to have lower returns than the average (typical) student who always goes to college (because benefits always exceed costs). Heckman et al. say Card's interpretation is wrong, because OLS is not the average effect. They say that there is a large positive sorting gain (comparative advantage in Roy model), Willis and Rosen prevail (marginal is below average), and true  $\beta > IV > OLS$ .

So what is MTE? First, run a probit to get  $\hat{D} = E[D|Z] = P(Z)$ . Second, using Equation 20.1, the observed outcome is

$$Y_i = \alpha_0 + X_i\theta_0 + D_i[\alpha_1 - \alpha_0 + X_i(\theta_1 - \theta_0)] + \{U_{i0} + D_i(U_{i1} - U_{i0})\}.$$

<sup>270</sup>See `randcoef` in Stata and Cabanillas et al. (2018) Stata Journal.

Now,

$$E[Y|Z, X] = \alpha_0 + X\theta_0 + E(D|Z)[\alpha_1 - \alpha_0 + X(\theta_1 - \theta_0) + E(U_1 - U_0|D = 1, Z)], \quad (20.2)$$

where the term in square brackets is the ATT. Next, invoke index sufficiency and condition on  $P(Z)$  instead on  $Z$  (this is like instrumenting for  $D$ ) and find the marginal treatment effect  $MTE(x, P(z))$  as  $\frac{\partial E[Y|X, P(Z)]}{\partial P(z)}$ .

Note that the non-linearity of Equation 20.2 implies heterogeneity in the treatment effects that is correlated with treatment status. This highlights the limitation of linear IV methods: the relationship between  $E[Y|X, P(Z)]$  and  $P(Z)$  is non-linear, so linear IV is misspecified and IV estimates depend on the instrument used (LATE). *Outcomes are non-linear function of participation probabilities.* The MTE is then simply the slope of this function at a given  $P(z)$ .<sup>271</sup>

The ATE equals the slope of the line that connects the  $E[Y|P(Z)]$  at  $P(Z) = 0$  and  $P(Z) = 1$ . Similarly, the  $TT(z)$  connects  $E[Y|P(Z)]$  at  $P(Z) = 1$  and  $P(z)$ . However, in practice, we do not have (treatments') data at  $P(Z) = 0$ . The LATE IV in the case of a binary instrument ( $Z$  taking only two values) just connects  $E[Y|P(Z)]$  at the two values  $z$  and  $z'$ .

Heckman (1990) shows that the sample-selection model is not identified without distributional assumptions (non-parametrically) unless we observe  $P(Z)$  at 0 and at 1 for some  $Z$  (we need the IV to be unbounded for this “identification at infinity”, which does not happen in practice). So how does one estimate the MTE? There are several ways. Heckman et al. use a **Local IV** method.<sup>272</sup> Think of a switching regression (with non-parametric  $\lambda$ s) applied at specific ranges of  $U_D$  (their rich instrument affects behavior in many parts of the ability distribution). See also Angrist (2004).

A simpler approach is proposed by Moffitt (2007, NBER WP no. 13534) who uses simple series estimation (splines in or polynomials of p-score, in his case  $\Phi(z'\theta_D)$ ) to non-parametrically trace out the shape of the relationship between participation probability and the outcome. Specifically, he runs a linear regression with regressors non-linear in  $\Phi(Z)$  to approximate  $Y_i = \alpha_0 + X\theta_0 + P(Z)g(P(Z)) +$

<sup>271</sup>Overidentification tests, which ask whether different IVs estimate different parameters, actually provide a test of this non-linearity (whether  $MTE$  depends on  $U_D$ ).

<sup>272</sup>They use non-parametric estimation by breaking the regression down into steps. See Section 13 for the definition of Local Linear Regression.  $\frac{\partial PE(U_1 - U_0|P)}{\partial P}$  is approximated using discrete differences. The local IV method has now been extended to unordered multiple choice models: <http://ftp.iza.org/dp3565.pdf>.

$\epsilon_i$ . The  $g$  function here is the TT effect (see equations 20.1 and 20.2 to see this) such that  $MTE = g(P(Z)) + P(Z)g'(P(Z))$ .

**Example 20.8.** *Moffitt (2007) can extrapolate out of the range of observed  $P(Z)$ , but one would have little trust in such results. He uses 3 IVs in order to operate in different portions of the  $P(Z)$  distribution. He first shows a histogram of  $\widehat{P}(Z)$ , which is thin towards  $P(Z) = 1$ , but he also explains that what matters is not only where the data is, but also where the instruments have incremental effects (where in the range of  $P(Z)$  they are not weak). Do see Figure 2 of the paper on this point. In his case, both conditions are satisfied in the region from 0.3 to 0.6.*

**Remark 213.** *Return to the issue of why  $\beta_{IV} > \beta_{OLS}$  in the returns to education estimation with constant returns discussed in example 20.7. For  $MTE > (\text{local}) OLS$  (as in the Card explanation of  $\beta_{IV} > \beta_{OLS}$ ), it is necessary that  $MTE > TT$  (outside of the neighborhood of  $P(Z) = 0$  or 1). Moffitt (2007) goes through a proof of this argument and rejects this explanation of why  $\beta_{IV} > \beta_{OLS}$  using his UK data.*

**Remark 214.** *An easy-to-read summary (including a discussion of identification of Equation 20.2) can be found in Manning (2003).<sup>273</sup> A 2016 CEPR DP No. 11390 by Cornelissen, Dustmann, Raute and Schoenberg is a more recent non-technical introduction to MTE estimation.*

**Remark 215.** *For MTE with a binary instrument  $Z$ , see Brinch, Mogstad, and Wiswall (2017, JPE) and Heckman and Vytlacil (2005, 2007) show that conditional independence of  $Y_0$  and  $Y_1$  implies constant MTEs.*

**Remark 216.** *Zhou and Xie (2019, JPE) define MTE as the expected treatment effect conditional on the propensity score and a latent variable representing unobserved resistance to treatment.*

**Remark 217.** *One can obtain MTE in the standard Heckman's  $\lambda$  approach using formulas in Heckman, Tobias, and Vytlacil (2000) NBER WP No. 7950. Assuming*

<sup>273</sup> Available at <http://econ.lse.ac.uk/staff/amanning/work/econometrics.html>

a parametric choice equation  $(\theta_D Z)$ ,<sup>274</sup> the parameters of interest are:<sup>275</sup>

$$\begin{aligned}ATE(x) &= x(\theta_1 - \theta_0) \\MTE(x, u_D) &= x(\theta_1 - \theta_0) + E(U_1 - U_0 | U_D = u_D) \\TT(x, z, D(z) = 1) &= x(\theta_1 - \theta_0) + E(U_1 - U_0 | U_D \geq -z\theta_D) \\LATE(x, D(z) = 0, D(z') = 1) &= x(\theta_1 - \theta_0) + E(U_1 - U_0 | -z'\theta_D \leq U_D \leq -z\theta_D)\end{aligned}$$

Assuming joint normality of  $(U_D, U_1, U_0)$  and normalizing variance in choice Probit, we have

$$\begin{aligned}TT(x, z, D(z) = 1) &= x(\theta_1 - \theta_0) + (\sigma_1 - \sigma_0) \frac{\varphi(z\theta_D)}{\Phi(z\theta_D)} \\LATE(x, D(z) = 0, D(z') = 1) &= x(\theta_1 - \theta_0) + (\sigma_1 - \sigma_0) \frac{\varphi(z'\theta_D) - \varphi(z\theta_D)}{\Phi(z'\theta_D) - \Phi(z\theta_D)} \\MTE(x, u_D) &= x(\theta_1 - \theta_0) + (\sigma_1 - \sigma_0)u_D.\end{aligned}$$

Now take integrals over sub-populations; e.g., for ATT average of the  $D = 1$  population. We assume the functional form (shape) of the non-linear relationship between  $P(Z)$  and  $E[Y|Z]$  so even a binary instrument (just observing two points) is enough to know the whole curve and the whole range of policy relevant parameters.

Heckman (2010, JEL) contrasts (structural) policy evaluation with randomization-defined program evaluation and then builds a bridge between the two. He applies the Marschak's Maxim (see Section 1) to provide a link between the (structural) Roy model and the (reduced-form) LATE framework. The formulas provided in this Remark are the bridge (see Vytlacil, 2002).

**Remark 218.** So far, we covered a binary treatment variable and so we discussed “local average treatment effects”. With a continuous  $X$ , the 2SLS estimates “local average partial effects”, it has a “average derivative” interpretation. Heckman, Urzua and Vytlacil (2006) and Heckman and Vytlacil (2007) extend the Imbens-Angrist model to ordered and unordered choice models, respectively, while Heckman (IZA DP no. 3980) extends the MTE approach to unordered choice.

**Remark 219.** See de Chaisemartin and D'Haultfoeuille NBER Working Paper No. 25904 for the ATE weights within two-way fixed effect linear models with heterogenous effects.

<sup>274</sup>We now hide  $\alpha_1 - \alpha_0$  inside the  $x(\theta_1 - \theta_0)$  to save space.

<sup>275</sup>Go back to example 19.12.: we estimated the ATE.

**Remark 220.** See Zhou and Xie (JPE) “Marginal Treatment Effects from a Propensity Score Perspective,”

In general, the lesson is that when responses are heterogeneous, there is no guarantee that 2SLS (with a valid instrument) is going to identify parameters of economic (policy) interest (any more than OLS). Especially when responses to choices vary among individuals and this variation affects the choices taken. When an IV is correlated with heterogeneous responses, 2SLS will not reveal the average partial effect. Heterogeneity is not a technical problem, but more often than not, it will correspond to the mechanics of the treatment effect. On the other hand, if the variation in the IV closely corresponds to the policy you have in mind (changing college tuition by the amount that corresponds to cost difference between those that live near a college or not), then LATE 2SLS is of interest and has high internal validity. Heckman (2010 JEL) argues that economists should use IV identification, but that they should use ranges of support of  $P(Z)$  (MTE weights) that correspond to actual policy changes.

Many of the “natural experiment” LATE 2SLS applications effectively estimate one structural equation whilst using an ad hoc external instrument (see Remarks 70 and 72 and Section 2 of Deaton 2009, NBER WP no. 14690), instead of estimating a simultaneous equation *model* (with exclusion restrictions) corresponding to economic theory. In this literature, the “questions [we answer]... are *defined* as probability limits of estimators and not by well-formulated economic problems” (Heckman, 2009, IZA DP no. 3980). The choice of the instrument determines the question we answer. We estimate unspecified “effects” instead of questions we care about.<sup>276</sup> Deaton (2009) adds that by using the LATE IV approach, we do not learn why and how the treatment works, only whether it works, which is a major issue for generalizing and applying the findings. Along similar lines, Heckman (2010 JEL) argues that one needs to separate two distinct tasks of policy analysis: (a) defining counterfactuals and (b) identifying causal models from data,<sup>277</sup> and that program evaluation ‘conflates’ the two. Deaton suggests to organize the experimental variation around the effect *mechanism*.

**Example 20.9.** Consider estimating the effect of railroads construction on poverty in a linear regression with an IV for railroad construction, such as whether the

<sup>276</sup>Heckman and Vytlacil (2005) then provide functions of IVs that answer well-posed questions (see above).

<sup>277</sup>Return again to the Introduction, to the last paragraph before Section 1, for a simpler restatement of causality being defined in terms of both a theoretical model and its empirical identification.

Government of China designates the given area as belonging to an “infrastructure development area” (Deaton 2009). The heterogeneity in the treatment parameter corresponds to different ways (context) of how railroads may alleviate poverty. The variation in the parameter is about the mechanisms that ought to be at the main objects of the enquiry. The deviation in the parameter from its average will be in the residual and will not be orthogonal to the IV if, within the group of cities designated as infrastructure zones, those that build railway stations are those that benefit the most in terms of poverty reduction, which one hopes to be the case. The Heckman’s local IV approach asks about how cities respond to their designation—asks about the mechanisms.

**Example 20.10.** Similarly, in the Maimonides rule regression-discontinuity example 9.7, the children that are shifted to classrooms of different sizes are not a random sample of all children in these classes. We need to know how children end up in different classes. See Urquiola and Verhoogen (2008).

Therefore, researchers recently *combine experimental variation with structural models*, as a way of testing the structure or to identify it. Structural models are then used to deliver parameters of policy interest (see Example 20.11 below). *Experiments ought to be theory driven* to deliver information on the mechanisms of treatment effects, to test predictions of theories that are generalizable, rather than to just test “what works”. For further reading, see “Giving Credit Where it is Due” by Banerjee and Duflo and the ‘theory-based impact evaluation’ calls by *3ie*.<sup>278</sup>

**Example 20.11.** See, e.g., papers by Attanasio, Meghir & Santiago and by Todd and Wolpin (AER) on the Progressa experiment in Mexico or the discussion of the use of structural estimation for evaluating labor-supply effects of earned-income-tax-credit policies by Blundell (2005, *Labour Economics*)<sup>279</sup> Todd and Wolpin: The Progressa experiment offered one level of subsidies to parents who send their children to school in Mexican villages. Todd and Wolpin estimate a model where parents make sequential decisions about sending their children to school or to

<sup>278</sup>At this point, read Blundell and Costa Dias (2008, IZA DP No. 3800). They present six related evaluation approaches in empirical microeconomics that we have covered in this course: (i) randomized (social) experiment, (ii) natural (quasi) experiment (i.e., difference in differences and when the method will identify the ATT), (iii) discontinuity design methods, (iv) matching, and (v) IV and control function methods.

<sup>279</sup><http://www.ucl.ac.uk/~uctp39a/Blundell%20-%20Adam%20Smith%20Lecture.pdf>

work, as well as about the timing and spacing of births. The model is estimated off the control group and it predicts well the response of the treatment group to the subsidy; it can therefore be used to infer the optimal size of the subsidy. Because in absence of treatment, there is no direct cost of schooling, they use observed child wages (the opportunity cost of attending school) to identify the model. The model is solved numerically, integrals are simulated, likelihood contributions are calculated by a smoothed frequency simulator.

**Remark 221.** Recently, Kline and Walters revisit the usual contrast between structural methods criticized for being sensitive to functional form assumptions and IV estimation. They study parametric LATE estimators based on Heckit and establish equivalence with IV to suggest that differences between structural and IV estimates often stem from disagreements about the target parameter rather than from functional form assumptions per se.

**Remark 222.** Henderson and Souto (2018, IZA DP no. 11914) overview IV estimation in a non-parametric regression setting.

## References

- Abadie, Alberto, Alexis Diamond, and Jens Hainmueller (2007) "Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California's Tobacco Control Program," *National Bureau of Economic Research Working Paper*: 12831, January 2007, 51 p.
- Abadie, Alberto, Joshua Angrist, and Guido Imbens (2002) "Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings," *Econometrica*: 70(1), Jan. 2002, pp. 91–117.
- Abadie, Alberto, and Guido Imbens (2002) "Simple and Bias-Corrected Matching Estimators for Average Treatment Effects," *National Bureau of Economic Research Technical Paper*: 283, October 2002, 57 p.
- Abowd, John M., Francis Kramarz, and David N. Margolis (1999) "High wage workers and high wage firms," *Econometrica* 67(2), March 1999: pp. 251–333.
- Abrevaya, J. and Dahl, C.M. (2008) "The effects of birth inputs on birthweight," *Journal of Business and Economic Statistics*: 26(4), pp. 379–397.
- Achen, Christopher H. (1986) *The statistical analysis of quasi-experiments*. University of California Press, 1986, 186 p.
- Ahn, H., L.J., Powell (1993) "Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism," *Journal of Econometrics*; 58(1-2), July 1993, pp. 3-29.
- Aigner, Dennis J. (1973) "Regression with a binary independent variable subject to errors of observation," *Journal of Econometrics*: 1(1), March 1973, pp. 49-59.
- Altonji, Joseph G., Todd E. Elder, and Christopher R. Taber (2005) "Selection on Observed and Unobserved Variables: Assessing the Effectiveness of Catholic Schools," *Journal of Political Economy*: 113(1), February 2005, pp. 151–184.
- Altonji, Joseph G. and Lewis M. Segal (1994) "Small Sample Bias in GMM Estimation of Covariance Structures", *National Bureau of Economic Research Technical Paper*: 156, June 1994.



- Amemiya T. (1984) "Tobit Models: A Survey," *Journal of Econometrics* 24(1-2), January-February 1984, pg: 3-61.
- Amemiya T. (1985): *Advanced Econometrics*, Harvard U. Press, 1985.
- Andrews, Schank and Upward (2006) "Practical fixed-effects estimation methods for the three-way error-components model," *Stata Journal* 6 (4), December 2006, pp. 461-481.
- Andrews, D. and J. Stock (2005) "Inference with Weak Instruments," *Advances in Economics and Econometrics*, Vol III, Blundel,, Newey and Persson (eds.), 122-173.
- Angrist, Joshua D. (1990) "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records," *American Economic Review*; 80(3), June 1990, pages 313-36.
- Angrist, Joshua D. (1991) "Grouped-data estimation and testing in simple labor-supply models," *Journal of Econometrics*: 47(2-3), February 1991, pp. 243-266.
- Angrist, Joshua D. (1995) "Conditioning on the Probability of Selection to Control Selection Bias", *National Bureau of Economic Research Technical Paper*: 181, June 1995.
- Angrist, Joshua D. (1998) "Estimating the Labor Market Impact of Voluntary Military Service Using Social Security Data on Military Applicants," *Econometrica* 66(2), March 1998, pages 249-88
- Angrist, Joshua D. (2004) "Treatment effect heterogeneity in theory and practice," *Economic Journal*: 114(494), March 2004, pp. C52-C83.
- Angrist, Joshua D. and Hahn, J. (2004) "When to control for covariates? Panel asymptotics for estimates of treatment effects," *The Review of Economics and Statistics*: 86(1), pp. 58-72.
- Angrist, Joshua D. and Alan B. Krueger (1992), "The Effect of Age at School Entry on Educational Attainment: An Application of Instrumental Variables with Moments from Two Samples," *Journal of the American Statistical Association*, 87, 328-336.
- Angrist, Joshua D. and Alan B. Krueger (1999) "Empirical Strategies in Labor Economics," Ashenfelter, Orley; Card, David, eds. *Handbook of Labor Economics*. Volume 3A. Handbooks in Economics, vol. 5. Amsterdam; New York and Oxford: Elsevier Science, North-Holland, 1999, pages 1277-1366.

- Angrist, Joshua D. and Alan B. Krueger. (2001) "Instrumental Variables and the Search for Identification: From Supply and Demand to Natural Experiments," *Journal of Economic Perspectives*; 15(4), Fall 2001, pages 69-85.
- Angrist, Joshua D. and V. Lavy (1999) "Using Maimonides' Rule to Estimate the Effect of Class Size on Scholastic Achievement", *Quarterly Journal of Economics* 114, pp. 533-575.
- Angrist, J.D. and Pischke, J. (2009): *Mostly Harmless Econometrics: An Empiricist's Companion*, Princeton University Press, 2009.
- Anselin (1988): *Spatial Econometrics: Methods and Models*, Kluwer Academic Press, 1988.
- Arabmazar A., P. Schmidt (1981) "Further Evidence on the Robustness of the Tobit Estimator to Heteroskedasticity", *Journal of Econometrics*, 17(2), Nov. 1981, pg: 253-58.
- Arellano, Manuel and Stephen Bond (1991) "Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations," *The Review of Economic Studies*: 58(2), April 1991, pp. 277-297
- Arellano, M., and Bonhomme, S. (2012) "Identifying Distributional Characteristics in Random Coefficients Panel Data Models," *Review of Economic Studies* 79, 987-1020..
- Arellano, Manuel and O. Bover (1995) "Another look at the instrumental variable estimation of error-components models," *Journal of Econometrics*: 68(1), July 1995, pp. 29-51.
- Arellano Manuel, Bo Honoré: *Panel Data Models: Some Recent Developments*, <<ftp://ftp.cemfi.es/wp/00/0016.pdf>>
- Arellano Manuel, C. Meghir (1992) "Female Labor Supply & On-the-Job Search: An Empirical Model Estimated Using Complementarity Data Sets," *Review of Economic Studies*, 59, 1992, 537-559.
- Ashenfelter, Orley and David Card (1985) "Using the Longitudinal Structure of Earnings to Estimate the Effect of Training Programs," *The Review of Economics and Statistics*: 67(4), Nov. 1985, pp. 648-660.
- Ashenfelter, O., Harmon, C. and Oosterbeek, H. (1999) "A review of estimates of the schooling/earnings relationship, with tests for publication bias," *Labour Economics*, 6(4),pp. 453-470.
- Ashenfelter, O. and A. Kruger (1994) "Estimates of the Economic Return to Schooling from a New Sample of Twins," *American Economic Review* 84: 1157-1173.

- Athey, Susan and Guido W. Imbens (2006) "Identification and Inference in Nonlinear Difference-in-Differences Models," *Econometrica*: 74(2), March 2006, pp. 431-497.
- Attanasio, Meghir & Santiago (in progress) "Education Choices in Mexico: Using a Structural Model and a Randomized Experiment to evaluate Progresas"
- Aydemir, Abdurrahman and George J. Borjas (2006) "A Comparative Analysis of the Labor Market Impact of International Migration: Canada, Mexico, and the United States," *National Bureau of Economic Research* (NBER) WP No. 12327, June 2006, 66 p.
- Baker, M., A. Melino (2000) "Duration Dependence and Nonparametric Heterogeneity: A Monte Carlo Study," *Journal of Econometrics*, 96(2), June, pp.357-393
- Baker, Michael and Nicole M. Fortin (2001). "Occupational Gender Composition and Wages in Canada, 1987-1988." *Canadian Journal of Economics*, 34(2), 345-376.
- Bauer, T. K. and Sinning, M. (2005) "Blinder-Oaxaca Decomposition for Tobit models," *IZA Discussion Papers*: 1795.
- Bauer, T. K. and Sinning, M. (2008) "An extension of the Blinder-Oaxaca Decomposition to nonlinear models," *AStA Advances in Statistical Analysis*: 92(2), May 2008, pp. 197-206.
- Becker, S.O. and Ichino, A. (2002) "Estimation of average treatment effects based on propensity scores," *Stata Journal*, StataCorp LP, 2(4), November 2002, pp.358-377.
- Berg, Van den, J. Gerard, G. Ridder (1993) "On the Estimation of Equilibrium Search Models from Panel Data," van Ours, Jan C., Pfann, Gerard A., Ridder, Geert, eds. Labor demand and equilibrium wage formation. *Contributions to Economic Analysis*, vol. 213. Amsterdam; London and Tokyo: North-Holland, pages 227-45.
- Berry, S., M. Carnall, T.P. Spiller, (1996) "Airline Hubs: Costs, Markups and the Implications of Customer Heterogeneity," *National Bureau of Economic Research* WP 5561, May 1996, pp. 20.
- Bertrand, Marianne, Esther Duflo, and Sendhil Mullainathan (2002) "How Much Should We Trust Differences-in-Differences Estimates?" *NBER Working Paper* No. 8841, March 2002, 39 p.
- Björklund, Anders and Robert Moffitt (1987) "The Estimation of Wage Gains and Welfare Gains in Self-Selection Models," *The Review of Economics and Statistics*: 69(1), Feb. 1987, pp. 42-49.

- Blundell, Richard W. and Richard J. Smith (1986) "An Exogeneity Test for a Simultaneous Equation Tobit Model with an Application to Labor Supply," *Econometrica*: 54(3), May 1986, pp. 679-685.
- Blundell, Richard W., Pierre-André Chiappori, Thierry Magnac and Costas Meghir (2005) "Collective Labour Supply: Heterogeneity and Nonparticipation," *IZA Discussion Paper* No. 1785, September 2005, 58 p.
- Blundell, Richard W. and James L. Powell (2004) "Endogeneity in Semiparametric Binary Response Models," *Review of Economic Studies*: 71(3), July 2004, pp. 655-679.
- Blundell, Richard W. and James L. Powell (2007) "Censored regression quantiles with endogenous regressors," *Journal of Econometrics*: 141(1), November 2007, pp. 65-83.
- Bound, J., D. Jaeger, and R. Baker (1995) "Problems with Instrumental Variables Estimation when the Correlation between Instruments and the Endogenous Explanatory Variable is Weak", *Journal of the American Statistical Association*, 90.
- Bourguignon, F., Fournier, M. and Gurgand, M. (2007) "Selection bias corrections based on the multinomial logit model: Monte Carlo comparisons, " *Journal of Economic Surveys*: 21(1), Blackwell Publishing, Feb. 2007, pp. 174-205.
- Brock, W.A. and Durlauf, S.N. (2001) "Discrete Choice with Social Interactions " *Review of Economic Studies*: 68(2), pp. 235-260.
- Brock, W.A. and S.N. Durlauf (2007) "Identification of binary choice models with social interactions," *Journal of Econometrics*: 140(1), September 2007, pp. 52-75.
- Buchinsky, M. (1998) "The dynamics of changes in the female wage distribution in the USA: a quantile regression approach," *Journal of Applied Econometrics*: 13(1), pp. 1-30.
- Cameron, A. Colin, Jonah B. Gelbach and Douglas L. Miller (2008) "Bootstrap-Based Improvements for Inference with Clustered Errors," *Review of Economics and Statistics*, 90(3), 414-427.
- Card, D. (1993) "Using Geographic Variation in College Proximity to Estimate the Return to Schooling," *National Bureau of Economic Research Working Paper*: 4483, October 1993, pg. 26.
- Card, D. and Alan B. Krueger (1994) "Minimum Wages and Employment: A Case Study of the Fast-Food Industry in New Jersey and Pennsylvania" *American Economic Review*: 84(4), Sep., 1994, pp. 772-793.

- Carneiro, Hansen, and Heckman (2003) "Estimating Distributions of Treatment Effects with an Application to the Returns to Schooling and Measurement of the Effects of Uncertainty on College," *National Bureau of Economic Research Working Paper*: 9546, March 2003, pg. 62.
- Carneiro P., E. Vytlacil, J.J. Heckman (2002) "Estimating The Return to Education When It Varies Among Individuals": [http://www.cerge.cuni.cz/events/seminars/sem-fall01/011108\\_a.asp](http://www.cerge.cuni.cz/events/seminars/sem-fall01/011108_a.asp)
- Chamberlain G. (1984) "Panel Data," in *Handbook of Econometrics* vol. II, pp. 1247-1318. Amsterdam, North-Holland.
- Chamberlain, G. (1980) "Analysis of Covariance with Qualitative Data," *Review of Economic Studies*; 47(1), Jan. 1980, pages 225-38.
- Chamberlain, G. (1982) "Multivariate Regression Models for Panel Data," *Journal of Econometrics* 18(1), Jan. 1982, pages 5-46.
- Cameron, A. Colin, Gelbach, Jonah B., and Douglas L. Miller (2008) "Bootstrap-based improvements for inference with clustered errors," *The Review of Economics and Statistics*, August 2008, 90(3): 414-427.
- Chay, Kenneth Y. and James L. Powell (2001) "Semiparametric Censored Regression Models," *The Journal of Economic Perspectives*: 15(4), Autumn 2001, pp. 29-42.
- Chen, S. (1999) "Distribution-free estimation of the random coefficient dummy endogenous variable model," *Journal of Econometrics*: 91(1), July 1999, pp. 171-199.
- Chernozhukov, V. and Hansen, C. (2008) "The reduced form: a simple approach to inference with weak instruments" *Economics Letters*: 100(1), pp. 68-71.
- Chioda, L. and M. Jansson (2006) "Optimal Conditional Inference for Instrumental Variables Regression," unpublished manuscript, department of economics, UC Berkeley.
- Cochran, W. G. and Rubin, D. B. (1973) "Controlling Bias in Observational Studies: A Review" *Sankhya*: Ser. A 35, pp. 417-446.
- Conley and Taber (2005) "Inference with "Difference in Differences" with a Small Number of Policy Changes," *NBER Technical Working Paper No. 312*, July 2005, 51 p.
- Cosslett S. (1981) "Maximum Likelihood Estimator for Choice-Based Samples", *Econometrica*, 49(5), Sept. 1981, pages 1289-1316.

- Costa Dias, Monica, Hidehiko Ichimura and Gerard J. van den Berg (2008) "The Matching Method for Treatment Evaluation with Selective Participation and Ineligibles" *IZA DP No. 3280*, January 2008, 36 p.
- Cox, D.R. (1972) "Regression models and life-tables," *Journal of the Royal Statistical Society. Series B (Methodological)*, 34(2), 1972, pp. 187-220.
- Cox, D.R. (1975) "Partial Likelihood," *Biometrika*: 62(2), 1975, p.269.
- Crump, R.K., Hotz, V.J., Imbens, G.W. and Mitnik, O.A. (2009) "Dealing with limited overlap in estimation of average treatment effects " *Biometrika*: Oxford University Press, 96(1), pp. 187-199.
- Cruz, L.M. and M.J. Moreira (2005) "On the Validity of Econometric Techniques with Weak Instruments: Inference on Returns to Education Using Compulsory School Attendance Laws," *Journal of Human Resources*, 40(2), 393-410.
- Cutler D., J. Gruber (1995) "Does Public Insurance Crowd Out Private Insurance?" *National Bureau of Economic Research*, Working Paper: 5082, April 1995, pages 33
- Davidson, R. and J.G. MacKinnon (1993): *Estimation and Inference in Econometrics*, Oxford University Press, 1993.
- Deaton, A. (1985) "Panel data from time series of cross-sections," *Journal of Econometrics*: 30(1-2), pp. 109-126.
- Deaton, A. (1997) *The analysis of household surveys: A microeconomic approach to development policy*, Baltimore and London: Johns Hopkins University Press for the World Bank, 1997, pg: 67-72.
- Deaton, A. (2009) "Instruments of development: Randomization in the tropics, and the search for the elusive keys to economic development, " *NBER WP No. 14690*.
- Dehejia, Rajeev H. and Sadek Wahba (1998) "Casual effects in non-experimental studies: re-evaluating the evaluation of training programs," *NBER WP No. 6586*, June 1998, 24 p.
- Donald, Stephen G. and Kevin Lang (2007) "Inference with Difference in Differences and Other Panel Data," *The Review of Economics and Statistics*, May 2007, 89(2): 221-233.
- Durlauf, S.N. (2002) "On the empirics of social capital," *Economic Journal*: 112(483), pp. F459-F479.

- Eckstein Z. and K. Wolpin (1989) "The Specification and Estimation of Dynamic Stochastic Discrete Choice Models: A Survey", *Journal of Human Resources*, 24(4), Fall 1989, pages 562-98.
- Elbers, Chris and Geert Ridder (1982) "True and Spurious Duration Dependence: The Identifiability of the Proportional Hazard Model," *The Review of Economic Studies*: 49(3), July 1982, pp. 403-409.
- Engberg, J. (1990) "Structural Estimation of the Impact of Unemployment Benefits on Job Search," *University of Wisconsin*, Ph.D. 1990
- Engberg J., S.L. Kim (1995) "Efficient Simplification of Bisimulation Formulas," in *TACAS*, pages 58-73.
- Epple D., H. Sieg (1999) "Estimating Equilibrium Models of Local Jurisdictions," *Journal of Political Economy* 107(4), August 1999, pages 645-81.
- Fan J., I. Gijbels (1996): *Local Polynomial Modelling and its Applications*, New York: Chapman & Hall, 1996.
- Flores-Lagunes, A. (2007) "Finite sample evidence of IV estimators under weak instruments," *Journal of Applied Econometrics*: 22(3), pp. 677-694.
- Frölich, M. (2007), "Nonparametric IV Estimation of Local Average Treatment Effects with Covariates," *Journal of Econometrics*, 139, 35-75.
- Galton, F. (1886) "Regression towards mediocrity in hereditary stature," *Journal of the Anthropological Institute*: 15, pp. 246-263.
- Geman, S., D. Geman (1984) "Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images," Polson,-Nicholas; Tiao,-George-C., eds. *Bayesian inference*. Volume 2. Elgar Reference Collection. International Library of Critical Writings in Econometrics, vol. 7. Aldershot, U.K.: Elgar; distributed in the U.S. by Ashgate, Brookfield, Vt., 1995, pages 149-69. Previously published: [1984].
- Geweke, John, Michael Keane and David Runkle (1994) "Alternative Computational Approaches to Inference in the Multinomial Probit Model," *The Review of Economics and Statistics*: 76(4), Nov. 1994, pp. 609-632.
- Godfrey, L.G. (1988) "Misspecification tests in econometrics: The Lagrange multiplier principle and other approaches," *Econometric Society Monographs series*, no. 16, Cambridge; New York and Melbourne: Cambridge University Press, 1988, pages xii, 252.

- Greene, William H. (2005): *Econometric Analysis*, fifth edition, Prentice Hall.
- Griliches Z. (1977) "Estimating the Returns to Schooling: Some Econometric Problems," *Econometrica* 45(1), Jan. 1977, pages 1-22.
- Griliches Z., J. Hausman (1986) "Errors in Variables in Panel Data," *Journal of Econometrics* 31: pp. 93-118.
- Haavelmo, T. (1994) "Statistical Testing of Business-Cycle Theories," Poirier, D.J., ed. *The methodology of econometrics*. Volume 1. Elgar Reference Collection. International Library of Critical Writings in Econometrics, vol. 6. Aldershot, U.K.: Elgar; distributed in the U.S. by Ashgate, Brookfield, Vt., 1994, pages 75-80. Previously published: [1943].
- Hahn J. and J. Hausman (2002a) "Notes on bias in estimators for simultaneous equation models," *Economics Letters* 75, pp. 237-41.
- Hahn J. and J. Hausman (2002b) "A New Specification Test for the Validity of Instrumental Variables," *Econometrica* 70 (1): 163-189.
- Hall, Peter (1992) *The Bootstrap and Edgeworth Expansion*, New York, Springer, 1992.
- Ham and LaLonde (1996) "The Effect of Sample Selection and Initial Conditions in Duration Models: Evidence from Experimental Data on Training", *Econometrica* 64(1), January 1996, pages 175-205.
- Hansen, C.B. (2007) "Asymptotic properties of a robust variance matrix estimator for panel data when T is large," *Journal of Econometrics*: 141(2), December 2007, pp. 597-620.
- Härdle, W., (1989): *Applied Nonparametric Regression*, Cambridge University Press, 1989.
- Hausman, J.A. (1978) "Specification Tests in Econometrics," *Econometrica* 46: 1251-1272.
- Hausman, J.A. et al. (1991) "Identification and Estimation of Polynomial Errors-in-Variables Models", *Journal of Econometrics*, 50(3), December 1991, pages 273-95.
- Hausman, J.A., D. McFadden (1984) "Specification Tests for the Multinomial Logit Model", *Econometrica* 52(5), September 1984, pages 1219-40.
- Heckman, J.J. (1979) "Sample Selection Bias as a Specification Error," *Econometrica* 47: 153-161.



- Heckman, J. J. (1980) "Addendum to sample selection bias as a specification error," E. Stromsdorfer, E. and Farkas, G. (Eds.) *Evaluation Studies Review Annua*. Volume 5. Beverly Hills: Sage Publications.
- Heckman, J.J. (1990) "A Nonparametric Method of Moments Estimator for the Mixture of Geometrics Model," Hartog, J.; Ridder,G.; Theeuwes, J., eds. *Panel data and labor market studies*. Contributions to Economic Analysis, vol. 192, Amsterdam; Oxford and Tokyo: North-Holland; distributed in the U.S. and Canada by Elsevier Science, New York, 1990, pages 69-79.
- Heckman, J.J. (2000) "Causal Parameters and Policy Analysis in Econometrics: A Twentieth Centurtury Perspective," *Quarterly Journal of Econometrics*, Feb. 2000.
- Heckman, J.J., H. Ichimura, J. Smith, P. Todd (1995) "Nonparametric characterization of selection Bias Using Experimental Data: A Study of Adult Males in JTPA," Presented at the *Latin American Econometric Society Meeting*, Caracas, Venezuela, 1994.
- Heckman, J., Ichimura, H., Smith, J. and Todd, P. (1998) "Characterizing selection bias using experimental data," *Econometrica*: 66(5), pp. 1017-1098.
- Heckman, J.J., H. Ichimura and P. Todd (1997) "Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Programme," *The Review of Economic Studies*: 64(4), Special Issue: Evaluation of Training and Other Social Programmes, Oct. 1997, pp. 605-654.
- Heckman J.J., B. Singer (1984) "A Method for Minimizing the Impact of Distributional Assumptions in Econometric Models for Duration Data," *Econometrica* 52(2), March 1984, pages 271-320.
- Heckman, J.J., Tobias, J.L. and Vytlacil, E. (2000) "Simple Estimators for Treatment Parameters in a Latent Variable Framework with an Application to Estimating the Returns to Schooling," *NBER WP No. 7950*, October 2000, 36 p.
- Heckman,J.J. and Urzua, S. (2009) "Comparing IV with Structural Models: What Simple IV Can and Cannot Identify," *IZADiscussion Paper no. 3980*. IZA, 2009.
- Heckman J.J., E. Vytlacil (1998) "Instrumental Variables Methods for the Correlated Random Coefficient Model: Estimating the Average Rate of Return to Schooling When the Return Is Correlated with Schooling", *Journal of Human Resources* 33(4), Fall 1998, pages 974-87.

- Heckman J.J., E. Vytlacil (2000) "The Relationship between Treatment Parameters within a Latent Variable Framework," *Economics Letters* 66(1), January 2000, pages 33-39.
- Heckman, J.J. and Vytlacil, E. (2005) "Structural equations, treatment effects and econometric policy evaluation," *NBER WP No. 11259*.
- Hellerstein, J.K., G.W. Imbens, (1999) "Imposing Moment Restrictions from Auxiliary Data by Weighting," *Review of Economics and Statistics* 81(1), February 1999, pages 1-14.
- Holland P. (1986) "Statistics and Causal Inference", *Journal of the American Statistical Association* 81(396), December 1986, pages 945-60.
- Honoré, B. (1992) "Trimmed Lad and Least Squares Estimation of Truncated and Censored Regression Models with Fixed Effects," *Econometrica*: 60(3), May 1992, pp. 533-565.
- Honoré, B., S. Khan and J.L. Powell (2002) "Quantile regression under random censoring," *Journal of Econometrics*: 109(1), July 2002, pp. 67-105.
- Honoré, B. and A. Lleras-Muney (2006) "Bounds in competing risks models and the war on cancer," *Econometrica*: 74(6), November 2006, pp. 1675-1698.
- Horowitz J. (1992) "A Smoothed Maximum Score Estimator for the Binary Response Model," *Econometrica*, 60(3), May 1992, pages 505-31.
- Hotz V., R. Miller (1993) "Conditional Choice Probabilities and the Estimation of Dynamic Models," *Review of Economic Studies* 60(3), July 1993, pages 497-529.
- Hoxby, C. (2000) "Does competition among public schools benefit students and taxpayers?," *American Economic Review*: 90(5), pp. 1209-1238.
- Hsiao, C. (1986) *Analysis of Panel Data*, Cambridge U. Press, 1986.
- Hsiao, C. and M.H. Pesaran (2004) "Random Coefficient Panel Data Models," *IZA Discussion Paper no. 1236*. IZA, 2004.
- Ichimura H. (1993) "Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single-Index Models," *Journal of Econometrics*, 58(1-2), July 1993, pages 71-120.
- Ichimura, H. and C. Taber (2000) "Direct Estimation of Policy Impacts," *NBER Technical WP No. 254*, 41 p.

- Imbens, Guido W. (1992) "An Efficient Method of Moments Estimator for Discrete Choice Models with Choice-Based Sampling," *Econometrica*, 60(5), September 1992, pages 1187-214.
- Imbens, Guido W. (2004) "Semiparametric Estimation of Average Treatment Effects under Exogeneity: A Review," *The Review of Economics and Statistics*, February 2004, 86(1): pp. 4-29.
- Imbens, Guido W. and J.K. Hellerstein (1993) "Raking and Regression," *Harvard Institute of Economic Research Working Paper No. 1657*, 15 p.
- Imbens, Guido W. and Thomas Lemieux (2007) "Regression Discontinuity Designs: A Guide to Practice," *National Bureau of Economic Research Working Paper: 13039*, April 2007, pg. 37.
- Imbens, Guido W. and Jeffrey M. Wooldridge (2007) "What's New in Econometrics," Unpublished Manuscript, National Bureau of Economic Research, Summer 2007.
- Inoue, Atsushi and Gary Solon (2005) "Two-Sample Instrumental Variables Estimators," *NBER Technical Working Paper 311*, June 2005, 19 p.
- Jakubson, G. (1991) "Estimation and Testing of the Union Wage Effect Using Panel Data," *Review of Economic Studies*, 58: 971-991.
- Jurajda, S. (2002) "Estimating the Effect of Unemployment Insurance Compensation on the Labor Market Histories of Displaced Workers," *Journal of Econometrics*, June 2002, Vol. 108, No. 2.
- Kelejian, Harry H. (1971) "Two-Stage Least Squares and Econometric Systems Linear in Parameters but Nonlinear in the Endogenous Variables," *Journal of the American Statistical Association*: 66(334), June 1971, pp. 373-374.
- Khan S., S. Chen (2000) "Estimating Censored Regression Models in the Presence of Nonparametric Multiplicative Heteroskedasticity," *Journal of Econometrics* 98, 283-316.
- Khan S., S. Chen (2001) "Semiparametric Estimation of a Partially Linear Censored Regression Model," *Econometric Theory* 17, pg. 567-590.
- Khan and Powell (2001) "Two-step estimation of semiparametric censored regression models," *Journal of Econometrics*: 103(1-2), July 2001, pp. 73-110.

- Kiefer N. (1988) "Economic Duration Data and Hazard Function," *Journal of Economic Literature* 26(2), June 1988, pg: 646-79.
- Kiefer, D. (1988) "Interstate Wars in the Third World: A Markov Approach," *Conflict Management and Peace Science*; 10(1), Spring 1988, pages 20-36.
- Klein R., R. Spady (1993) "An Efficient Semiparametric Estimator for Binary Response Models," *Econometrica*, 61(2), March 1993, pages 387-421.
- Kling, J. (1999) "Interpreting IV Estimates of the Returns to Schooling," *Princeton University Industrial Relations Section WP* 415, 1999.
- Kmenta, J. (1986): *Elements of Econometrics*, 2nd Ed., New York: Macmillan, 1986.
- Koenker, R. (2004) "Quantile Regression for Longitudinal Data," *Journal of multivariate analysis*: 91(1), pp. 74-89.
- Koenker, R., and G. Bassett (1982) "Robust Tests for Heteroscedasticity Based on Regression Quantiles," *Econometrica*, 50(1), Jan. 1982, pages 43-61.
- Kyriazidou, Ekaterini (1997) "Estimation of a Panel Data Sample Selection Model," *Econometrica*: 65(6), Nov. 1997, pp. 1335-1364.
- Lafontaine, F., J.K. Shaw (1995) "Firm-Specific Effects in Franchise Contracting: Sources and Implications", mimeo.
- Lafontaine, F., J.K. Shaw (1999) "The Dynamics of Franchise Contracting: Evidence from Panel Data," *Journal of Political Economy*, 107(5), October 1999, pages 1041-80.
- LaLonde, Robert J. (1986) "Evaluating the Econometric Evaluations of Training Programs with Experimental Data," *The American Economic Review*: 76(4), Sep. 1986, pp. 604-620.
- Lancaster T., (1990): *The Econometric Analysis of Transition Data*, Cambridge U. Press, 1990.
- Leamer, E. E. (1978) "Least-Squares versus Instrumental Variables Estimation in a Simple Errors in Variables Model," *Econometrica*, 46(4), July 1978, pages 961-68.
- Lee, L.F. (1981) "Fully Recursive Probability Models and Multivariate Log-Linear Probability Models for the Analysis of Qualitative Data," *Journal of Econometrics*, 16(1), May 1981, pages 51-69.

- Lee, L.F. (1983) "Generalized econometric models with selectivity," *Econometrica*: 51, pp. 507-512.
- Lee, Myoung-jae (1996) *Methods of Moments and Semiparametric Econometrics for Limited Dependent Variable Models*. Springer, 1996, 296 p.
- Lee, D. and Lemieux, T. (2009) "Regression Discontinuity Designs in Economics," *NBER Working Paper*: 14723.
- Lechner, M. (2001) "The empirical analysis of East Germany fertility after unification: An update," *European Journal of Population*: 17, pp. 61-74.
- Lewbel A. (2004) "Simple Estimators For Hard Problems: Endogeneity in Discrete Choice Related Models," *Boston College Working Papers in Economics* No. 604, 34 p.
- Loeb, Susanna and John Bound (1996) "The Effect of Measured School Inputs on Academic Achievement: Evidence from the 1920s, 1930s and 1940s Birth Cohorts," *The Review of Economics and Statistics*: 78(4), Nov. 1996, pp. 653-664.
- Machado, José A.F. and José Mata (2005) "Counterfactual decomposition of changes in wage distributions using quantile regression," *Journal of Applied Econometrics*: 20(4), pp. 445-465.
- Maddala G.S., (1983): *Limited-dependent and Qualitative Variables in Econometrics*, Cambridge U. Press, 1983.
- Manning, A. (2004) "Instrumental Variables for Binary Treatments with Heterogeneous Treatment Effects: A Simple Exposition," *Contributions to Economic Analysis & Policy*: 3(1), Article 9.
- Manski, Charles F. (1975) "Maximum Score Estimation of the Stochastic Utility Model of Choice," *Journal of Econometrics*, 3(3), Aug. 1975, pages 205-28.
- Manski, Charles F. (1985) "Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator," *Journal of Econometrics*, 27(3), March 1985, pages 313-33.
- Manski, Charles F. (1995) *Identification Problems in the Social Sciences*, Harvard University Press, 194 p.

- Manski, Charles F. and Daniel McFadden (1981) *Structural Analysis of Discrete Data and Econometric Applications*, <elsa.berkley.edu/users/mcfadden/discrete.html> Cambridge: The MIT Press, 1981.
- Marschak, J. (1953) "Economic Measurements for Policy and Prediction," Hood, W., Koopmans, T. and Hendry, D. F., eds. *Studies in Econometric Method*. Wiley, New York, 1953.
- Marschak, J. (1995) "Economic Interdependence and Statistical Analysis," Hendry, D. F.; Morgan, M. S., eds. *The foundations of econometric analysis*. Cambridge; New York and Melbourne: Cambridge University Press, 1995, pages 427-39. Previously published: [1942].
- Matsudaira, J.D. (2008) "Mandatory summer school and student achievement," *Journal of Econometrics*: 142(2), February 2008, pp. 829-850.
- Matzkin, R. (1992) "Nonparametric and Distribution-Free Estimation of the Binary Threshold Crossing and the Binary Choice Models," *Econometrica*, 60(2), March 1992, pages 239-70.
- McCleary, R. and Barro, R. (2006) "Religion and Economy," *Journal of Economic Perspectives*: 2(2), Spring 2006, pp. 49-72.
- McCulloch, R. E., P.E. Rossi (1994) "An Exact Likelihood Analysis of the Multinomial Probit Model," *Journal of Econometrics*; 64(1-2), Sept.-Oct. 1994, pages 207-40.
- McFadden, Daniel (1984) "Econometric Analysis of Qualitative Response Models" in *Handbook of Econometrics*. Volume II, Griliches, Z., Intriligator M. ed. Handbooks in Economics series, book 2. Amsterdam; New York and Oxford: North-Holland; distributed in the U.S. and Canada by Elsevier Science, New York, 1984.
- McFadden, Daniel (1989) "A Method of Simulated Moments for Estimation of Discrete Response Models without Numerical Integration", *Econometrica*, 57(5), September 1989, pages 995-1026.
- McFadden, Daniel, and Kenneth Train (2000) "Mixed MNL models for discrete response," *Journal of Applied Econometrics*: 15(5), pp. 447-470.
- Miguel, E., Satyanath, S. and Sergenti, E. (2004) "Economic Shocks and Civil Conflict: An Instrumental Variables Approach," *Journal of Political Economy* : 112(41), pp. 725-753.

- Miller, R. A. (1984) "Job Matching and Occupational Choice," *Journal of Political Economy*, 92(6), December 1984, pages 1086-120.
- Moffitt (2007) "Estimating Marginal Returns to Higher Education in the UK," *NBER WP no.* 13534, October 2007, 41 p.
- Moreira, Marcelo J. (2003) "A Conditional Likelihood Ratio Test for Structural Models," *Econometrica*: 71(4), July 2003, pp. 1027-1048.
- Moreira, Marcelo J. and Brian P. Poi (2003) "Implementing Conditional Tests with Correct Size in the Simultaneous Equations Model," *Stata Journal*: 1(1), pp. 1-15.
- Münich, Daniel, Jan Svejnar and Katherine Terrell (2005) "Returns to Human Capital under the Communist Wage Grid and during the Transition to a Market Economy," *The Review of Economics and Statistics*, February 2005, 87(1): 100–123.
- Murray, Michael P. (2006) "Avoiding Invalid Instruments and Coping with Weak Instruments" *Journal of Economic Perspectives*: 20(4), Fall 2006, pp. 111-132.
- Nevo, A. and Rosen, A. (2008) "Identification with Imperfect Instruments," *NBER Working Paper*: 14434.
- Newey, W. (1985) "Generalized Method of Moments Specification Tests," *Journal of Econometrics*, 29: 229-238.
- Newey, W. and J. Powell (1990) "Efficient Estimation of Linear and Type I Censored Regression Models under Conditional Quantile Restrictions," *Econometric Theory*: 6(3), Sep. 1990, pp. 295–317.
- Olsen, R. J. (1987) "Note on the Uniqueness of the Maximum Likelihood Estimator for the Tobit Model," *Econometrica*, 46(5), Sept. 1978, pages 1211-15.
- Pagan, A., F. Vella (1989) "Diagnostic Tests for Models Based on Individual Data: A Survey," *Journal of Applied Econometrics*, 4(0), Supplement, December 1989, pages S29-59.
- Pakes, A.S. (1986) "Patents as Options: Some Estimates of the Value of Holding European Patent Stocks," *Econometrica*, 54(4), July 1986, pages 755-84.
- Pakes, A.S. and D. Pollard (1989) "Simulation and the Asymptotics of Optimization Estimators," *Econometrica*, 57(5), pp. 1027-1057.

- Pesaran, M.H., Y. Shin, R.J. Smith (2000) "Structural Analysis of Vector Error Correction Models with Exogenous I(1) Variables," *Journal of Econometrics*, 97(2), August 2000, pages 293-343..
- Popper, K. (1959) *The logic of scientific discovery*. London: Hutchinson.
- Powell, J.L. (1983) "The Asymptotic Normality of Two-Stage Least Absolute Deviations Estimators," *Econometrica*: 51(5), Sep. 1983, pp. 1569-1575.
- Powell, J.L. (1984) "Least Absolute Deviation Estimation for the Censored Regression Model," *Journal of Econometrics*, 25(3), July 1984, pages 303-25.
- Powell, J.L. (1986) "Symmetrically Trimmed Least Squares Estimation for Tobit Models," *Econometrica*, 54(6), November 1986.
- Pudney, S. (1989) *Modelling Individual Choice*, Basil Blackwell, 1989.
- Rivers, Douglas and Quang H. Vuong (1988) "Limited information estimators and exogeneity tests for simultaneous probit models" *Journal of Econometrics*: 39, November 1988, pp. 347-366.
- Robinson, P.M. (1988) "Using Gaussian Estimators Robustly," *Oxford Bulletin of Economics and Statistics*; 50(1), February 1988, pages 97-106.
- Rosenbaum, P.R. (2002) "Covariance Adjustment in Randomized Experiments and Observational Studies," *Statistical Science*: 17(3), 2002, pp. 286-327.
- Rosenbaum, P.R. and D.B. Rubin (1983) "Assessing Sensitivity to an Unobserved Binary Covariate in an Observational Study with Binary Outcome," *Journal of the Royal Statistical Society. Series B (Methodological)*: 45(2), 1983, pp. 212-218.
- Rosenbaum P.R. and D.B Rubin (1984) "Estimating the Effects Caused by Treatments: Comment [On the Nature and Discovery of Structure]", *Journal of the American Statistical Association*, 79(385), March 1984, pages 26-28.
- Roy P.N. (1970) "On the Problem of Measurement of Capital," *Economics Affairs*, 15(6-8 156), Aug. 1970, pages 269-76.
- Rust, J. (1987) "A Dynamic Programming Model of Retirement Behavior," *National Bureau of Economic Research Working Paper*: 2470, December 1987, pages 64.



- Rust, J. (1992) "Do People Behave According to Bellman's Principle of Optimality?," *Hoover Institute Working Papers in Economics*: E-92-10, May 1992, pages 76.
- Rust J. (1994) "Structural Estimation of Markov Decision Processes," Engle,-Robert-F.; McFadden,-Daniel-L., eds. *Handbook of econometrics*. Volume 4. Handbooks in Economics, vol. 2. Amsterdam; London and New York: Elsevier, North-Holland, 1994, pages 3081-3143.
- Schennach, S (2004) "Nonparametric regression in the presence of measurement error," *Econometric Theory*: 20(6), December, 2004, pp. 1046-1093.
- Scott, A.J. and Wild, C.J. (1997) "Fitting regression models to case-control data by maximum likelihood," *Biometrika*: 84 (1), pp. 57-71.
- Sianesi, B. (2001) "Differential effects of Swedish active labour market programs for unemployed adults during the 1990s," *IFS Working Papers W01/25*, Institute for Fiscal Studies, 2001.
- Silverman (1986): *Density Estimation for Statistics and Data Analysis*, London, Chapman & Hall, 1986.
- Staiger, Douglas and James H. Stock (1997) "Instrumental Variables Regression with Weak Instruments," *Econometrica*: 65(3), May 1997, pp. 557-586.
- Suri, T. (2011). Selection and Comparative Advantage in Technology Adoption. *Econometrica*, 79(1), 159-209..
- Stock, J. H., and Yogo, M. (2005) "Testing for Weak Instruments in Linear IV Regression" in *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, ed. D.W. Andrews, D.W. ed. and Stock, J.H. Cambridge University Press. 2005.
- Taber, Chr.R. (2000) "Semiparametric Identification and Heterogeneity in Discrete Choice Dynamic Programming Models," *Journal of Econometrics*, 96(2), June 2000, pages 201-29.
- Tanner, M. A (1993) *Tools for Statistical Inference*, 2nd. edition. New York: Springer-Verlag.
- Tibshirani, Robert and Trevor Hastie (1987) "Local Likelihood Estimation," *Journal of the American Statistical Association*: 82(398), June 1987, pp. 559-567.

- Todd, P. E., and Wolpin, K.I. (2006) "Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Validate a Dynamic Behavioral Model of Child Schooling and Fertility," *American Economic Review*: 96(5), pp. 1384-1417.
- Urquiola, M. and Verhoogen, E. (2007). "Class size and sorting in market equilibrium: theory and evidence," *NBER Working Papers*: 13303.
- Van der Klaauw (2008) "Regression-Discontinuity Analysis: A Survey of Recent Developments in Economics" (forthcoming).
- White, H. (1980) "A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity," *Econometrica*, 48(4), May 1980, pages 817-38.
- Willis R. J., S. Rosen (1979) "Education and Self-Selection," *Journal of Political Economy*, 87(5), Part 2, Oct. 1979, pages S7-36.
- Wooldridge, Jeffrey M. (2002) *Econometric Analysis of Cross Section and Panel Data*, MIT Press, 2002.
- Wooldridge, Jeffrey M. (2003) "Cluster-Sample Methods in Applied Econometrics," *The American Economic Review*, Vol. 93, No. 2, Papers and Proceedings of the One Hundred Fifteenth Annual Meeting of the American Economic Association, Washington, DC, January 3-5, 2003 (May, 2003), pp. 133-138
- Wolpin, K. I. (1984) "An Estimable Dynamic Stochastic Model of Fertility and Child Mortality," *Journal of Political Economy*, 92(5), October 1984, pages 852-74.
- Yatchew A., (1998) "Nonparametric Regression Techniques," *Journal of Economic Literature* 1998.
- Zanella, Giulio (2007) "Discrete Choice with Social Interactions and Endogenous Memberships," *Journal of the European Economic Association*, March 2007, 5(1): 122-153.
- Zellner, Arnold (1996) *An Introduction to Bayesian Inference in Econometrics*. Wiley-Interscience, 448 p.
- Ziliak, S.T. and McCloskey, D.N.(2008) *The Cult of Statistical Significance: How the Standard Error Costs Us Jobs, Justice, and Lives*, Ann Arbor, MI: University of Michigan Press, 320 p.