

STATA SKILL PORTFOLIO

THE REFERENCE

A Stata Reference for Empirical Economics

Methods, Code, and Verified Results

from data management to causal inference, one technique at a time

DUC V. LE

Lê Vĩnh Đức

PhD Candidate, Department of Economics · Georgetown University

Office ICC-580 · 3700 O St NW, Washington, DC 20057, USA

dv111@georgetown.edu

Built and verified on Stata 19.5 SE · Updated June 2026

Contents

Preface	3
Roadmap: choosing a method	4
I Setup and environment	7
1 Getting started	8
1.1 Running the code	8
1.2 Packages	8
II Empirical foundations	9
2 Stata basics	10
2.1 The code	10
2.2 The results	11
3 Data management	13
3.1 The code	13
3.2 The results	16
4 Linear regression	18
4.1 The method	18
4.2 The code	18
4.3 The results	19
5 Panel data, DiD, and event studies	22
5.1 The method	22
5.2 The code	22
5.3 The results	24
6 Time series	26
6.1 The method	26
6.2 The code	26
6.3 The results	27
7 Importing public data	29
7.1 The code	29
7.2 The results	30

III	The causal-inference toolkit	32
8	Instrumental variables (2SLS)	33
8.1	The method	33
8.2	The code	33
8.3	The results	35
9	Binary outcome models	37
9.1	The method	37
9.2	The code	37
9.3	The results	39
10	Staggered-adoption DiD	40
10.1	The method	40
10.2	The code	40
10.3	The results	42
11	Regression discontinuity	44
11.1	The method	44
11.2	The code	44
11.3	The results	47
12	Synthetic control	49
12.1	The method	49
12.2	The code	49
12.3	The results	51
13	Dynamic panel data (GMM)	53
13.1	The method	53
13.2	The code	53
13.3	The results	55
IV	Appendices	57
A	Python/R → Stata cheat sheet	58
B	The methods, explained	59
B.1	Ordinary least squares	59
B.2	Panel data, difference-in-differences, and event studies	60
B.3	Time series	61
B.4	Instrumental variables (2SLS)	62
B.5	Binary outcome models	63
B.6	Staggered-adoption difference-in-differences	63
B.7	Regression discontinuity	64
B.8	Synthetic control	65
B.9	Dynamic panel data (GMM)	66
	References	68

Preface

This reference collects, in one place, the twelve scripts of the Stata skill portfolio together with the method behind each, the code, and the results it produces. The portfolio's individual files (01_basics.do ... 11_synthetic_control.do) and the package documentation (packages/PACKAGES.md) are excellent as references, but they are separate. This document binds them into a single, coherent textbook so the material reads as one body of knowledge.

Code availability. The twelve .do scripts behind every figure, table, and number in this reference — and this document itself — are openly available at github.com/duc-v-le/stata-empirical-methods.

A companion volume, *Stata Line-by-Line: An Empirical Economics Guide*, walks the same eleven scripts command by command for readers new to Stata's syntax; this book assumes that fluency and concentrates on the method, the code, and the results.

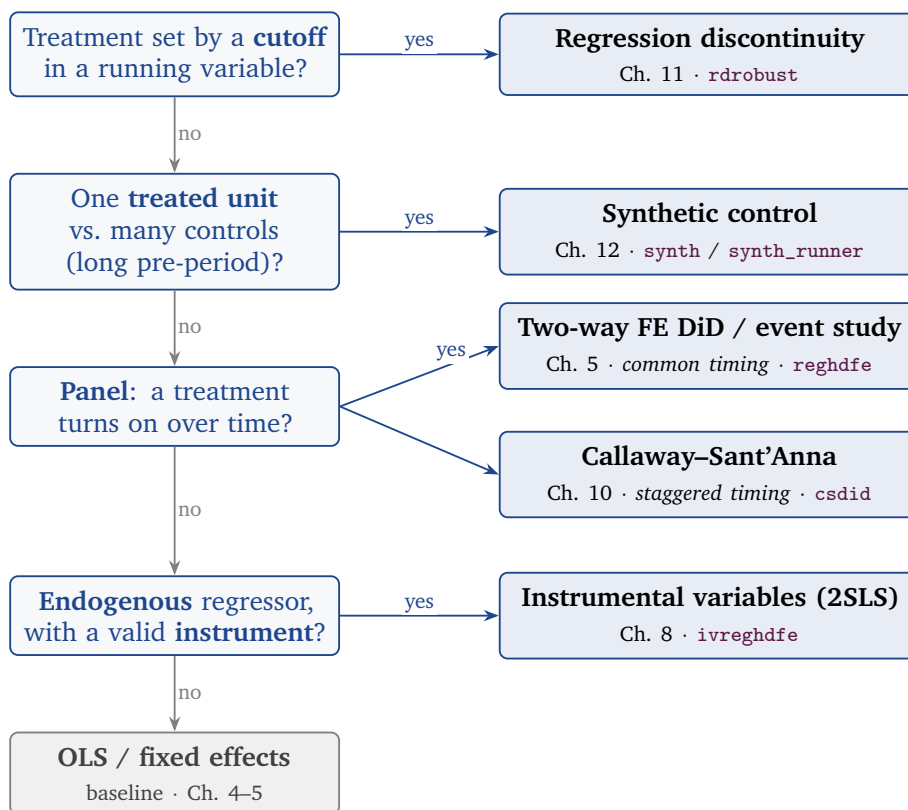
What this is — and is not. This is a *learning and demonstration* portfolio, not a research project. The data are either Stata's bundled example datasets, a real public series from FRED, or simulated with a fixed seed so that each estimator can be checked against a *known* true parameter. Every chapter ends by comparing the estimate to that truth.

How the book is organised. Part I sets up the environment and packages. Part II covers the empirical foundations — data management, regression, panel methods, time series, and importing public data. Part III is the modern causal-inference toolkit — instrumental variables, limited-dependent-variable models, staggered difference-in-differences, regression discontinuity, and synthetic control. The appendices give a Python/R→Stata translation table; a from-scratch explanation of every method ([Appendix B](#), with — for each estimator — its idea, symbols, formula, *why* it works, the assumption and what breaks it, a worked example, and a tie-back to that chapter's own run showing it recovers the planted answer); and the references.

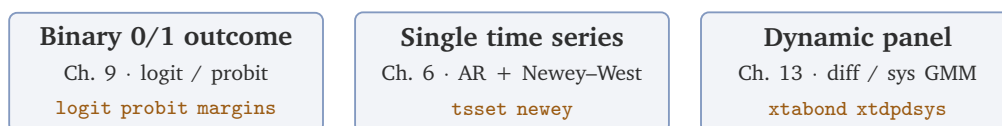
Reproducing everything. From the portfolio root, run `do 00_run_all.do` in Stata (after installing dependencies with `do packages/install_packages.do`). Each script writes its log to `logs/` and its figures and tables to `output/` — the same files embedded throughout this reference.

Roadmap: choosing a method

The questions lead to the identification strategy that fits the data, then to that chapter. Outcome- and data-type choices (binary outcomes, a single time series, dynamic panels) are *orthogonal* to the design and are boxed separately at the bottom. For an unfamiliar method, [Appendix B](#) (“The methods, explained”) builds each one up from scratch — intuition, formula, a worked example, and the result the chapter recovers.



Orthogonal to the design — pick by the outcome / data type:



Always clean and document the data first (Ch. 3); report robust or clustered standard errors.

How to read the roadmap

The diagram is a *decision tree*, read from the top-left box straight down, one yes/no question at a time about the data and how the treatment came about. The first “yes” follows the arrow to the method on the right — the identification strategy — and the box names the chapter and the Stata command; a “no” drops to the next question. The questions run from the most specific design to the most general, so the *first* “yes” is the answer; if every answer is “no,” the path ends at the plain-regression baseline at the bottom. The tree sorts by *identification strategy*, not by data shape: one panel dataset can land on regression discontinuity, difference-in-differences, instrumental variables, or the baseline, depending on where its causal leverage comes from.

The five questions, in plain words.

- **A cutoff in a running variable?** Did treatment switch on the instant some number crossed a threshold — a test score clearing a pass mark, an age passing 65, a vote share above 50%? Then units just below and just above the cutoff are otherwise alike, and the *jump* in the outcome at the threshold is the effect. ⇒ **regression discontinuity** (Ch. 11).
- **One treated unit, many controls?** Exactly one state (or firm, or country) was treated, with a long run of pre-treatment history? Build a weighted blend of the untreated units that mimics the treated one’s path *before* treatment, and read the effect as the gap that opens *after*. ⇒ **synthetic control** (Ch. 12).
- **A panel where treatment turns on over time?** Many units followed over time, with a treatment that switches on. If everyone switches on the *same* date, two-way fixed-effects difference-in-differences (Ch. 5); if they switch at *different* dates (a staggered roll-out), naive two-way FE is biased, so use Callaway–Sant’Anna (Ch. 10).
- **An endogenous regressor with a valid instrument?** The key right-hand-side variable is correlated with the error — reverse causality, or an omitted confounder — but an *instrument* exists that shifts that regressor without otherwise touching the outcome. This question is about the *regressor*, not the data shape: it applies to a panel just as well as a cross-section — `ivreghdfe` instruments *while* absorbing fixed effects. (It sits on the “no” branch only because the prior box asks whether the design is a treatment-switches-on *difference-in-differences*, not whether the data happen to be a panel.) ⇒ **instrumental variables / 2SLS** (Ch. 8).
- **None of the above?** With no special design, fall back to OLS / fixed effects (Ch. 4–5) — the honest baseline. Report the result as a *descriptive association*, not a causal effect.

The three boxes at the bottom answer a different question. The tree chooses the *design* — how a correlation is argued to be causal. The bottom boxes are about the *shape of the outcome or data*, and they are *orthogonal*: one layers on top of whatever design the tree picked, not instead of it. A **binary 0/1 outcome** calls for logit/probit and `margins` (Ch. 9); a **single time series** for an autoregression with Newey–West HAC errors (Ch. 6); a **panel with a lagged dependent variable** (persistence) for dynamic-panel GMM (Ch. 13). They mix freely with the tree — a difference-in-differences *on a binary outcome*, say, uses a design box *and* the logit box.

Two rules always apply. Whatever the path: **clean and document the data first** (Ch. 3), and **report robust or clustered standard errors** — clustered at the level treatment varies (state, firm, individual).

A worked walk-through. Consider 50 states, a minimum-wage increase that different states enacted in different years, and the goal of estimating its effect on teen employment. Walking the tree: no cutoff (skip RD); more than one treated state (skip synthetic control); *yes* — it is a panel with treatment turning on over time, and the timing is *staggered*, so the arrow points to **Callaway–Sant’Anna** (Ch. 10). Employment is continuous, so no bottom box applies; the standard errors would be clustered by state. Four questions, one decision.

Part I

Setup and environment

Chapter 1

Getting started

1.1 Running the code

Each script is a self-contained `.do` file. Run one interactively with `do 01_basics.do`, or in batch from a shell:

```
/Applications/StataNow19SE/StataSE19.app/Contents/MacOS/stata-se -b do 00_run_all.do
```

The master script `00_run_all.do` runs all twelve in order. Outputs land in `output/` (figures, \LaTeX /CSV tables, saved datasets) and logs in `logs/`.

1.2 Packages

Everything runs on base Stata except for a small set of community commands, all installed in one step with `do packages/install_packages.do`. Table 1.1 lists them and why each is needed; full provenance is in `packages/PACKAGES.md`. The two newest chapters add *no* dependencies: dynamic-panel GMM (Ch. 13) uses the built-in `xtabond/xtdpdsys`, and the efficient-GMM step (Ch. 8) the built-in `ivregress gmm` — Roodman’s `xtabond2` is the popular community alternative for dynamic panels, but is not required here.

Table 1.1: Community (SSC/GitHub) packages used in the portfolio

Package	Used in	Purpose
<code>reghdfe</code>	04, 09	Regression absorbing high-dimensional fixed effects
<code>ftools</code>	(dep)	Fast data operations; <code>reghdfe</code> dependency
<code>require</code>	(dep)	Version checker; <code>reghdfe</code> dependency
<code>estout</code>	03, 04, 07, 08	Publication tables (<code>esttab</code>) to \LaTeX /CSV
<code>coefplot</code>	04	Coefficient / event-study plots
<code>ivreghdfe</code>	07	IV/2SLS with high-dimensional fixed effects
<code>ivreg2, ranktest</code>	(dep)	IV engine and weak-IV tests
<code>csdid</code>	09	Callaway–Sant’Anna staggered-adoption DiD
<code>drdid</code>	(dep)	Doubly-robust DiD; <code>csdid</code> dependency
<code>rdrobust</code>	10	Local-polynomial regression discontinuity
<code>rddensity, lpdensity</code>	10	RD manipulation / density test
<code>synth</code>	11	Abadie synthetic control
<code>synth_runner</code>	11	Placebo inference for <code>synth</code> (from GitHub)
<code>distinct</code>	(dep)	<code>synth_runner</code> dependency

Part II

Empirical foundations

Chapter 2

Stata basics

The leap from Python or R is mostly syntactic: one dataset is held in memory at a time; `generate/replace` create and overwrite variables; `egen` handles group and window operations; and variable labels keep output self-documenting. The script below loads the bundled auto data, explores it, builds variables, and exports two figures.

2.1 The code

The workflow. Load a bundled dataset, inspect and summarise it, build and label variables with `generate/egen`, subset with `if`, draw a quick graph, and save a `.dta` file.

```
1  *=====
2  * 01_basics.do --- Stata fundamentals for applied economists
3  *-----
4  * Run: from the repository root, in Stata: do 01_basics.do
5  * Batch: stata-se -b do 01_basics.do
6  *=====
7  version 17 // write code against a fixed language version
8  clear all
9  set more off // don't pause for --more-- in long output
10 set linesize 90
11
12 capture log close
13 capture mkdir logs // create the log folder if absent (gitignored)
14 log using "logs/01_basics.log", replace text
15
16 *-----
17 * 1. Load data. sysuse loads datasets shipped with Stata (works offline).
18 *-----
19 sysuse auto, clear // 1978 automobile data (74 cars)
20 describe // structure: variables, types, labels (~ R str())
21 codebook mpg foreign, compact // quick audit of a few variables
22 list make price mpg in 1/5 // peek at rows (~ head())
23
24 *-----
25 * 2. Summarize / explore.
26 *-----
27 summarize // means, sd, min, max for all numeric vars
28 summarize price, detail // percentiles, skew/kurtosis
29 tabulate foreign // one-way frequencies
30 tabulate foreign rep78, row // two-way, with row percentages
31 bysort foreign: summarize mpg // grouped summary (~ groupby().describe())
32
33 *-----
34 * 3. Create & transform variables.
35 * generate = new var; replace = overwrite; egen = "extended" generate.
36 *-----
37 generate gpm = 1/mpg // gallons per mile
```

```

38 label variable gpm "Gallons per mile"
39 generate price_k = price/1000
40 generate byte expensive = price > 6000 // boolean -> 0/1 indicator
41 egen mpg_z = std(mpg) // standardized mpg
42 egen price_mean_by_for = mean(price), by(foreign) // group mean (window fn)
43
44 * Conditional logic: replace ... if
45 generate size_class = "small"
46 replace size_class = "large" if weight > 3500
47
48 * Missing values: Stata stores numeric missing as "." (sorts as +infinity).
49 codebook rep78 // note the 5 missing values
50 count if missing(rep78)
51
52 *-----
53 * 4. Labels make output self-documenting (a Stata habit worth keeping).
54 *-----
55 label define yesno 0 "No" 1 "Yes"
56 label values expensive yesno
57 tabulate expensive
58
59 *-----
60 * 5. A couple of graphs, exported to disk (no display needed in batch).
61 *-----
62 histogram price, frequency title("Distribution of price")
63 graph export "output/01_hist_price.png", replace width(1200)
64
65 twoway (scatter price weight) (lfit price weight), ///
66 title("Price vs. weight with linear fit") legend(off)
67 graph export "output/01_scatter_price_weight.png", replace width(1200)
68
69 *-----
70 * 6. Save a cleaned copy (Stata's native .dta format).
71 *-----
72 save "output/auto_clean.dta", replace
73
74 log close
75 display "01_basics.do finished OK"

```

2.2 The results

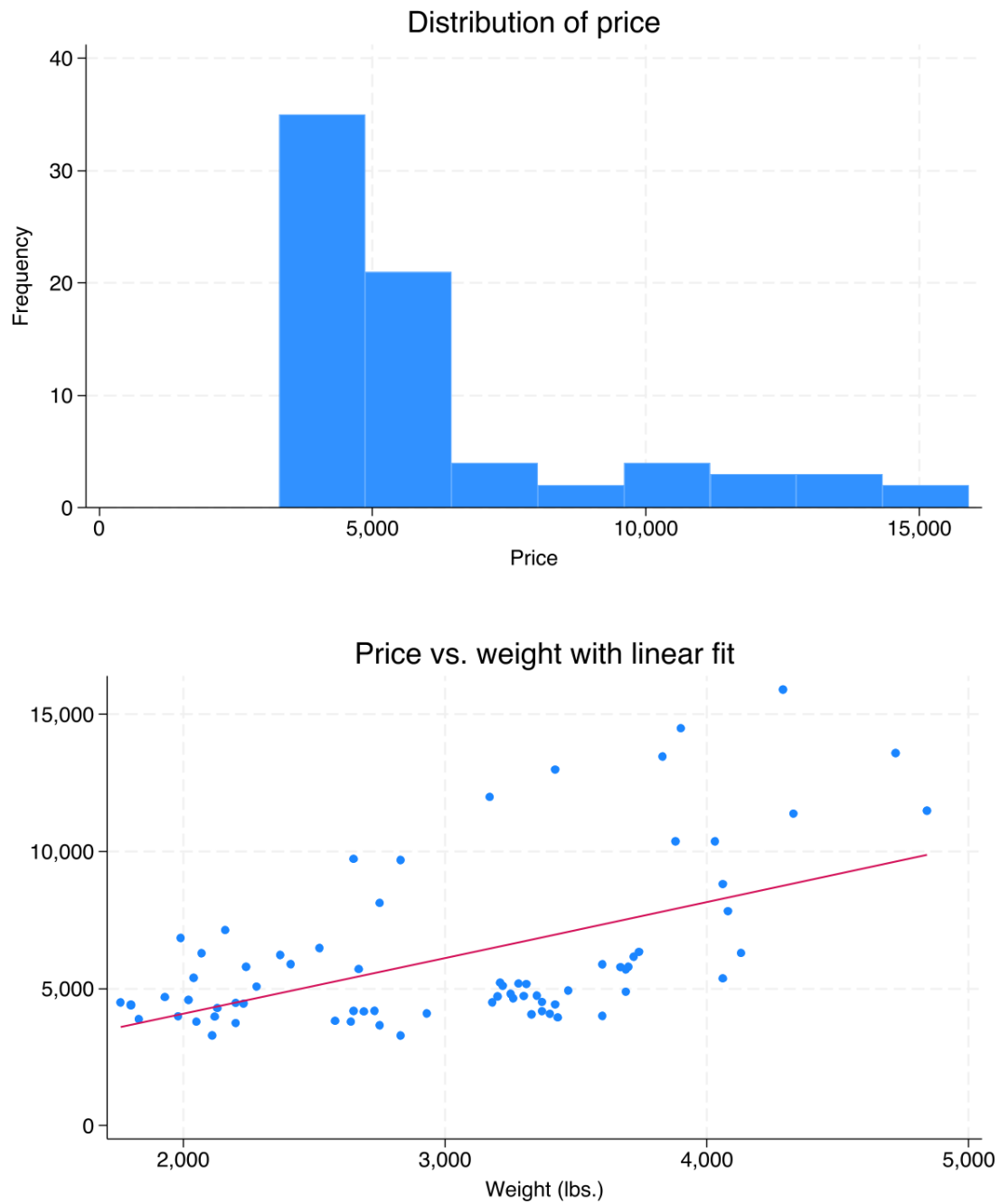


Figure 2.1: Distribution of price (top) and price vs. weight with a linear fit (bottom).

Chapter 3

Data management

The everyday plumbing of empirical work: joining datasets (`merge`), moving between wide and long layouts (`reshape`), aggregating (`collapse`), and handling dates, strings, and duplicates. These are the operations that consume most of the time in empirical work before any model runs.

The concept that is hardest to picture is `reshape`: the same data can sit in *long* form (one row per id–period) or *wide* form (one row per id, with the periods spread across columns). Figure 3.1 shows the two layouts and the commands that convert between them.

<table border="1"><thead><tr><th>id</th><th>year</th><th>income</th></tr></thead><tbody><tr><td>1</td><td>2010</td><td>100</td></tr><tr><td>1</td><td>2011</td><td>110</td></tr><tr><td>1</td><td>2012</td><td>130</td></tr><tr><td>2</td><td>2010</td><td>200</td></tr><tr><td>2</td><td>2011</td><td>220</td></tr><tr><td>2</td><td>2012</td><td>250</td></tr></tbody></table>	id	year	income	1	2010	100	1	2011	110	1	2012	130	2	2010	200	2	2011	220	2	2012	250	$\xrightarrow{\text{reshape wide}}$	<table border="1"><thead><tr><th>id</th><th>income2010</th><th>income2011</th><th>income2012</th></tr></thead><tbody><tr><td>1</td><td>100</td><td>110</td><td>130</td></tr><tr><td>2</td><td>200</td><td>220</td><td>250</td></tr></tbody></table>	id	income2010	income2011	income2012	1	100	110	130	2	200	220	250
id	year	income																																	
1	2010	100																																	
1	2011	110																																	
1	2012	130																																	
2	2010	200																																	
2	2011	220																																	
2	2012	250																																	
id	income2010	income2011	income2012																																
1	100	110	130																																
2	200	220	250																																
	$\xleftarrow{\text{reshape long}}$																																		

Figure 3.1: The same data in long form (left) and wide form (right); `reshape` converts between them.

In practice. Data cleaning is most of the work in any empirical project: real data arrive with inconsistent identifiers to harmonize, missing values to handle explicitly (Stata stores numeric missing as `.`, which sorts as $+\infty$), dates and strings to parse, and joins to verify — always check `_merge` after a `merge`.

3.1 The code

The workflow. Join two files with `merge`, reshape between wide and long with `reshape`, aggregate with `collapse`, and tidy dates, strings, and duplicate rows.

```
1  *-----  
2  * 02_data_management.do --- merge, reshape, collapse, dates, strings, dups  
3  *-----  
4  version 17  
5  clear all  
6  set more off  
7  set linesize 90  
8  capture log close  
9  capture mkdir logs // create the log folder if absent (gitignored)  
10 log using "logs/02_data_management.log", replace text  
11  
12 *-----  
13 * 1. MERGE (joins). Split auto into two files, then merge 1:1 on a key.
```

```

14 *-----
15 sysuse auto, clear
16 keep make price mpg
17 save "output/_left.dta", replace
18
19 sysuse auto, clear
20 keep make weight foreign
21 save "output/_right.dta", replace
22
23 use "output/_left.dta", clear
24 merge 1:1 make using "output/_right.dta"
25 tabulate _merge // 3 = matched both; 1 = master only; 2 = using only
26 assert _merge==3 // sanity check the join is clean
27 drop _merge
28
29 * m:1 (many-to-one): a country-year panel picks up one GDP figure per country.
30 clear
31 input str3 country gdp
32 "USA" 21000
33 "CAN" 1700
34 "MEX" 1100
35 end
36 save "output/_gdp.dta", replace
37 clear
38 input str3 country int year
39 "USA" 2020
40 "USA" 2021
41 "CAN" 2020
42 "CAN" 2021
43 "MEX" 2020
44 "MEX" 2021
45 end
46 merge m:1 country using "output/_gdp.dta" // each country's GDP fans out to its year rows
47 sort country year
48 list country year gdp _merge, sepby(country) noobs
49 drop _merge
50 erase "output/_gdp.dta"
51
52 * _merge values: an imperfect join shows 1 (master only), 2 (using only), 3 (matched).
53 clear
54 input str3 country pop
55 "USA" 331
56 "CAN" 38
57 "MEX" 126
58 end
59 save "output/_pop.dta", replace
60 clear
61 input str3 country gdp
62 "CAN" 1700
63 "MEX" 1100
64 "BRA" 1600
65 end
66 merge 1:1 country using "output/_pop.dta" // master {CAN,MEX,BRA} vs using {USA,CAN,MEX}
67 tabulate _merge
68 sort _merge country
69 list country pop gdp _merge, sepby(_merge) noobs
70 drop _merge
71 erase "output/_pop.dta"
72

```

```

73 *-----
74 * 2. RESHAPE (wide <-> long), the panel-builder's workhorse.
75 *-----
76 clear
77 input id year income
78 1 2010 100
79 1 2011 110
80 1 2012 130
81 2 2010 200
82 2 2011 220
83 2 2012 250
84 end
85 list, sepby(id)
86 reshape wide income, i(id) j(year) // long -> wide (income2010, income2011,...)
87 list
88 reshape long income, i(id) j(year) // wide -> long (back again)
89 list, sepby(id)
90
91 *-----
92 * 3. COLLAPSE (aggregate to a coarser level; ~ groupby().agg()).
93 *-----
94 sysuse auto, clear
95 collapse (mean) mean_price=price mean_mpg=mpg ///
96 (sd) sd_price=price (count) n=price, by(foreign)
97 list
98
99 * --- export the collapse summary as a table for the reference ---
100 capture file close cl
101 file open cl using "output/O2_collapse_table.tex", write replace
102 file write cl "\begin{table}[htbp]\centering" _n
103 file write cl "\caption{\cmd{collapse}: the 74-car sample reduced to one row per origin}" _n
104 file write cl "\begin{tabular}{lrrrr}" _n "\toprule" _n
105 file write cl "Origin & Mean price & Mean mpg & SD price & \\\\" _n "\midrule" _n
106 forvalues i = 1/'=N' {
107     local org = cond(foreign['i']==0, "Domestic", "Foreign")
108     file write cl "'org' & " %6.0f (mean_price['i']) " & " %5.1f (mean_mpg['i']) ///
109         " & " %6.0f (sd_price['i']) " & " %3.0f (n['i']) " \\" _n
110 }
111 file write cl "\bottomrule" _n "\end{tabular}" _n "\end{table}" _n
112 file close cl
113
114 *-----
115 * 4. DATES (Stata stores dates as integers; display via formats).
116 *-----
117 clear
118 input str10 raw
119 "2020-01-15"
120 "2020-06-30"
121 "2021-12-01"
122 end
123 generate edate = date(raw, "YMD") // string -> numeric daily date
124 list raw edate, noobs // edate is a number: days since 1960-01-01
125 format edate %td
126 list raw edate, noobs // same column, now shown as a calendar date
127 generate year = year(edate)
128 generate month = month(edate)
129 generate mdate = mofd(edate) // monthly date
130 format mdate %tm
131 list

```

```

132
133 *-----
134 * 5. STRING handling.
135 *-----
136 sysuse auto, clear
137 generate brand = word(make, 1)           // first token of make
138 generate make_up = upper(make)
139 generate is_amc = strpos(make, "AMC")>0
140 generate spc = strpos(make, " ")       // position of the first space (0 if none)
141 list make brand make_up spc in 1/6, noobs
142 tabulate brand if is_amc
143
144 *-----
145 * 6. DUPLICATES.
146 *-----
147 clear
148 input id v
149 1 10
150 1 10
151 2 20
152 3 30
153 3 30
154 end
155 duplicates report                       // how many dup rows
156 duplicates drop                         // drop exact duplicates
157 list
158
159 * tidy up scratch files
160 erase "output/_left.dta"
161 erase "output/_right.dta"
162
163 log close
164 display "02_data_management.do finished OK"

```

3.2 The results

A `collapse` reduces the 74-car auto sample to one row per origin — foreign cars are pricier on average and noticeably more fuel-efficient in this 1978 sample:

Table 3.1: `collapse`: the 74-car sample reduced to one row per origin

Origin	Mean price	Mean mpg	SD price	<i>N</i>
Domestic	6072	19.8	3097	52
Foreign	6385	24.8	2622	22

Join quality. After a `merge`, `_merge` labels every row — 1 master-only, 2 using-only, 3 matched — and an unmatched row receives missing (.) for the other file's variables. An imperfect join (master {CAN, MEX, BRA} against using {USA, CAN, MEX}):

```

+-----+
| country  pop   gdp           _merge |
+-----+
|      BRA   .  1600  Master only (1) |

```

```

+-----+
|   USA   331   .   Using only (2) |
+-----+
|   CAN   38  1700   Matched (3) |
|   MEX  126  1100   Matched (3) |
+-----+

```

A clean join is all 3, which `assert _merge==3` enforces after every merge in the script.

Dates and strings. `date("2020-01-15","YMD")` returns the integer 21929 (days since 1 Jan 1960); `format %td` only changes the display to 15jan2020. `String` helpers parse text in place — `word(make,1)` (first token), `upper(make)`, and `strpos(make," ")` (position of the first space):

```

+-----+
| make          brand      make_up  spc |
+-----+
| AMC Concord   AMC        AMC CONCORD  4 |
| AMC Pacer     AMC        AMC PACER    4 |
| AMC Spirit    AMC        AMC SPIRIT   4 |
| Buick Century Buick     BUICK CENTURY 6 |
| Buick Electra Buick     BUICK ELECTRA 6 |
| Buick LeSabre Buick     BUICK LESABRE 6 |
+-----+

```

Chapter 4

Linear regression

↔ the method, explained: [Appendix B.1](#)

4.1 The method

The workhorse model is

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + \varepsilon_i, \quad \hat{\beta} = (X'X)^{-1}X'y,$$

with [White \(1980\)](#) heteroskedasticity-robust (HC1) standard errors,

$$\widehat{\text{Var}}(\hat{\beta}) = (X'X)^{-1} \left(\sum_i \hat{\varepsilon}_i^2 x_i x_i' \right) (X'X)^{-1}.$$

Factor-variable notation (`i.`, `c.`, `##`) builds dummies and interactions; `margins` reports the average marginal effect

$$\overline{\text{ME}} = \frac{1}{n} \sum_i \frac{\partial \hat{E}[y | x_i]}{\partial x},$$

the interpretable quantity once interactions are present. Diagnostics: a joint *F* test, the [Breusch–Pagan](#) test for heteroskedasticity, and variance inflation factors.

In practice. OLS underlies hedonic price regressions (housing or car prices on characteristics), [Mincer \(1974\)](#) earnings equations (log wage on schooling and experience), and demand estimation — wherever the target is a conditional mean or a descriptive association.

4.2 The code

The workflow. Fit OLS with robust standard errors and factor-variable interactions, run post-estimation tests, and export a publication-ready table with `esttab`.

```
1  *=====
2  * 03_regression.do --- OLS, robust SE, factor variables, interactions,
3  *                               marginal effects, postestimation, publication tables
4  *=====
5  version 17
6  clear all
7  set more off
8  set linesize 90
9  capture log close
10 capture mkdir logs // create the log folder if absent (gitignored)
11 log using "logs/03_regression.log", replace text
12
13 sysuse auto, clear
```

```

14 generate weight_t = weight/1000           // weight in 1000s of lbs
15 label variable weight_t "Weight (1000 lbs)"
16 label variable mpg "Mileage (mpg)"
17 label variable foreign "Foreign"
18
19 *-----
20 * 1. OLS with heteroskedasticity-robust standard errors.
21 *   i.foreign = factor (dummy); c.var = continuous; ## = full interaction.
22 *-----
23 eststo clear
24 eststo m1: regress price mpg weight_t, vce(robust)
25 eststo m2: regress price mpg weight_t i.foreign, vce(robust)
26 eststo m3: regress price c.mpg##c.weight_t i.foreign, vce(robust)
27
28 * Print a comparison table to the log.
29 esttab m1 m2 m3, se star(* 0.10 ** 0.05 *** 0.01) b(%9.2f) ///
30     r2 label mtitles("Base" "+Foreign" "+Interaction") ///
31     title("Determinants of car price")
32
33 * Export the same table to LaTeX and CSV (publication-ready deliverables).
34 esttab m1 m2 m3 using "output/03_regression_table.tex", replace ///
35     se star(* 0.10 ** 0.05 *** 0.01) b(%9.2f) r2 label booktabs ///
36     mtitles("Base" "+Foreign" "+Interaction") ///
37     title("Determinants of car price\label{tab:price}")
38 esttab m1 m2 m3 using "output/03_regression_table.csv", replace ///
39     se star(* 0.10 ** 0.05 *** 0.01) b(%9.3f) r2 label plain
40
41 *-----
42 * 2. Marginal effects (margins) --- interpret models on the outcome scale.
43 *-----
44 quietly regress price c.mpg##c.weight_t i.foreign, vce(robust)
45 margins, dydx(mpg)           // average marginal effect of mpg
46 margins foreign             // predicted price by foreign status
47 margins, dydx(mpg) at(weight_t=(2 3 4)) // how the mpg slope varies with weight
48
49 *-----
50 * 3. Postestimation diagnostics.
51 *-----
52 quietly regress price mpg weight_t i.foreign
53 test mpg weight_t           // joint F-test (both = 0)
54 estat hettest               // Breusch-Pagan heteroskedasticity test
55 estat vif                   // variance inflation (multicollinearity)
56 predict yhat                // fitted values
57 predict ehat, residuals     // residuals
58 rvfplot, yline(0) title("Residual vs. fitted")
59 graph export "output/03_rvfplot.png", replace width(1200)
60
61 log close
62 display "03_regression.do finished OK"

```

4.3 The results

The three columns build the specification up step by step. In Stata's factor-variable notation `i.foreign` is the indicator $1[\text{foreign} = 1]$ and `c.mpg##c.weight_t` expands to *both* main effects

and their product, so the columns estimate

- (1) $\text{price}_i = \beta_0 + \beta_1 \text{mpg} + \beta_2 \text{weight_t} + \varepsilon_i,$
- (2) $\text{price}_i = \beta_0 + \beta_1 \text{mpg} + \beta_2 \text{weight_t} + \beta_3 \mathbf{1}[\text{foreign}] + \varepsilon_i,$
- (3) $\text{price}_i = \beta_0 + \beta_1 \text{mpg} + \beta_2 \text{weight_t} + \beta_3 (\text{mpg} \cdot \text{weight_t}) + \beta_4 \mathbf{1}[\text{foreign}] + \varepsilon_i,$

each with heteroskedasticity-robust (HC1) standard errors (`vce(robust)`). In this 1978 sample, weight and foreign manufacture are strong positive predictors of price and the mileage–weight interaction is small. These are *descriptive associations* in a teaching dataset, not causal effects—features the model omits (trim, brand, engine quality) plausibly drive both a car’s weight and its price, so a coefficient describes how price and that regressor co-move, not what would happen if it were changed.

Table 4.1: Determinants of car price

	(1) Base	(2) +Foreign	(3) +Interaction
Mileage (mpg)	-49.51 (95.81)	21.85 (80.75)	292.83* (153.90)
Weight (1000 lbs)	1746.56** (777.84)	3464.71*** (777.62)	5382.75*** (1084.04)
Domestic		0.00 (.)	0.00 (.)
Foreign		3673.06*** (664.94)	3369.81*** (757.93)
Mileage (mpg) × Weight (1000 lbs)			-118.91 (85.68)
Constant	1946.07 (4213.79)	-5853.70 (3873.72)	-10105.04*** (3212.97)
Observations	74	74	74
R^2	0.293	0.500	0.524

Standard errors in parentheses

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

The residual-versus-fitted plot fans out as the fitted value grows — visible heteroskedasticity, which is precisely why the standard errors are heteroskedasticity-robust.

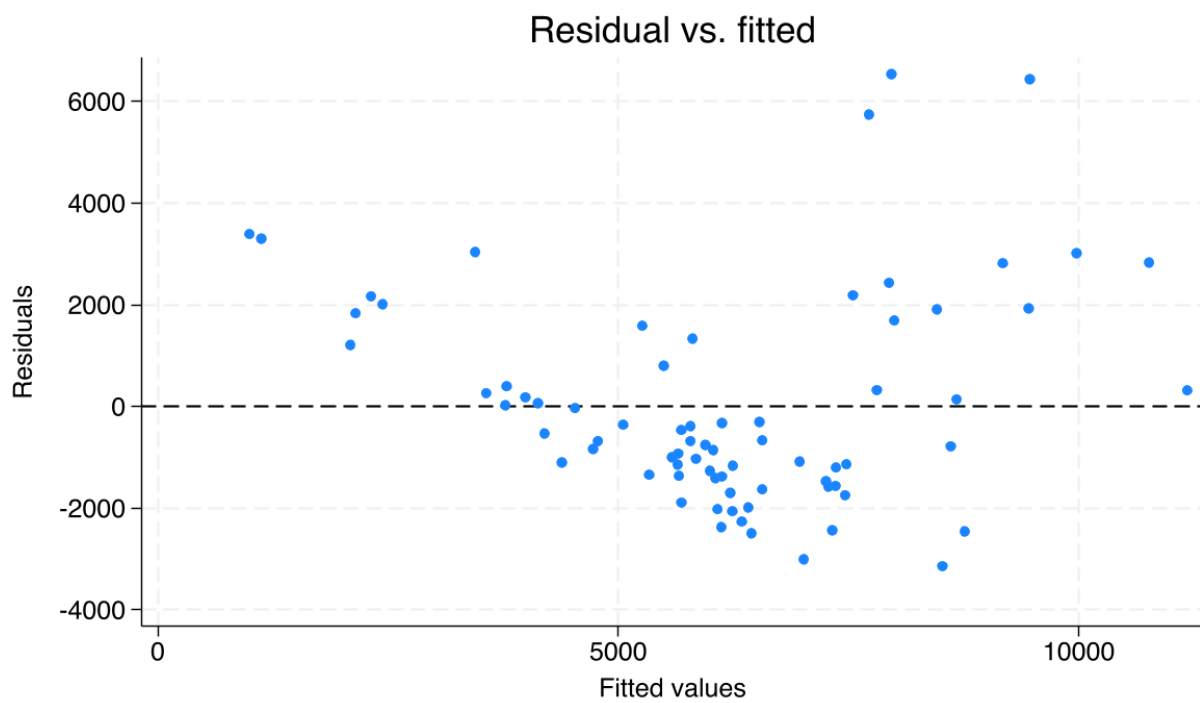


Figure 4.1: Residuals versus fitted values — the spread widens, indicating heteroskedasticity.

Chapter 5

Panel data, DiD, and event studies

↔ the method, explained: [Appendix B.2](#)

5.1 The method

With repeated observations on units over time, two-way fixed effects difference-in-differences is

$$y_{it} = \alpha_i + \gamma_t + \beta (\text{treat}_i \times \text{post}_t) + \varepsilon_{it},$$

where α_i absorbs unit level differences, γ_t common shocks, and β is the DiD estimand. The dynamic version traces the effect by event time k :

$$y_{it} = \alpha_i + \gamma_t + \sum_{k \neq -1} \beta_k \mathbf{1}[\text{treat}_i] \mathbf{1}[t - t_i^* = k] + \varepsilon_{it}.$$

Here t_i^* is unit i 's treatment date, so $k = t - t_i^*$ is *event time* (periods since onset) and the product $\mathbf{1}[\text{treat}_i] \mathbf{1}[t - t_i^* = k]$ switches on only for a treated unit in the period k steps from its own treatment; β_k is the effect at that horizon. The sum skips $k = -1$ because one event-time coefficient is *not estimable* — the full set of lead/lag dummies is collinear with the fixed effects, so one period must serve as the dropped *baseline*. By convention that period is $k = -1$, the one *just before* treatment: dropping it fixes $\beta_{-1} \equiv 0$, and every other β_k is then read *relative to* the eve of treatment. So the pre-treatment coefficients ($k < -1$) should sit near zero — a visual parallel-trends check — while the post-treatment ones ($k \geq 0$) trace out the dynamic effect. Keeping the leads in the sum — writing $k \neq -1$ rather than only $k \geq 0$ — is the whole point: with no treatment yet they ought to be zero, so they serve as a placebo test, the built-in pre-trend check an event study provides over a single DiD coefficient. Standard errors cluster by unit.

In practice. The two-by-two DiD is the template for policy evaluation — [Card and Krueger's \(1994\)](#) New Jersey–Pennsylvania minimum-wage study is the canonical example — and the event-study version traces the dynamics of a policy that switches on at a known date.

5.2 The code

The workflow. Declare the panel, estimate difference-in-differences by two-way fixed effects with `reghdfe`, and trace the dynamic event study with `coefplot`.

```
1 *=====
2 * 04_panel_did_eventstudy.do --- panel FE, difference-in-differences,
3 *                               two-way FE with reghdfe, dynamic event study
4 *-----
5 * Builds a SYNTHETIC firm-by-year panel (fully reproducible, no external data),
6 * then runs the canonical causal-inference workflow of applied empirical work.
```

```

7  *-----
8  version 17
9  clear all
10 set more off
11 set linesize 90
12 set seed 90210
13 capture log close
14 capture mkdir logs // create the log folder if absent (gitignored)
15 log using "logs/04_panel_did_eventstudy.log", replace text
16
17 *-----
18 * 1. Simulate a panel: 150 firms x 12 years (2010-2021).
19 *   Treated firms (half) get a +1.5 effect starting in 2015 that grows over time.
20 *-----
21 set obs 150
22 generate firm = _n
23 generate fe_firm = rnormal(0, 2) // firm fixed effect
24 generate byte treat = mod(firm, 2)==0 // half the firms are treated
25 expand 12
26 bysort firm: generate year = 2009 + _n // 2010..2021
27 generate byte post = year >= 2015 // treatment turns on in 2015
28 generate rel = year - 2015 // event time relative to treatment
29
30 generate year_fe = 0.3*(year-2015) + 0.5*sin(year) // common time shocks
31 generate tau = 0
32 replace tau = 0.5*(rel+1) if treat==1 & year>=2015 // dynamic treatment effect
33 generate y = 3 + fe_firm + year_fe + tau + rnormal(0,1)
34 label variable y "Outcome"
35
36 xtset firm year // declare panel structure
37 xtdescribe
38
39 *-----
40 * 2. Pooled OLS vs. fixed-effects DiD (the 2x2 interaction).
41 *-----
42 eststo clear
43 eststo ols: regress y i.treat##i.post, vce(cluster firm)
44 eststo fe: xtreg y i.treat##i.post, fe vce(cluster firm)
45 * Under firm FE the time-invariant treat main effect drops; the DiD estimate
46 * is the 1.treat#1.post coefficient (true value = average post effect).
47
48 *-----
49 * 3. Two-way fixed effects with reghdfe (absorb firm AND year FE).
50 *-----
51 eststo twfe: reghdfe y 1.treat#1.post, absorb(firm year) vce(cluster firm)
52
53 esttab ols fe twfe, se star(* 0.10 ** 0.05 *** 0.01) b(%9.3f) ///
54 mtitles("Pooled OLS" "Firm FE" "Two-way FE") ///
55 keep(1.treat#1.post) label title("Difference-in-differences estimates")
56 esttab ols fe twfe using "output/04_did_table.tex", replace booktabs ///
57 se star(* 0.10 ** 0.05 *** 0.01) b(%9.3f) keep(1.treat#1.post) label ///
58 mtitles("Pooled OLS" "Firm FE" "Two-way FE") title("DiD estimates")
59
60 *-----
61 * 4. Dynamic event study: treatment-effect path by event time.
62 *   The full set of event-time dummies sums to the (time-invariant) treated
63 *   indicator, which the firm FE absorbs, so one dummy must be dropped. We omit
64 *   t = -1 (the eve of treatment) as the single reference period, so every other
65 *   coefficient is read relative to it. Leads should sit near zero (no pre-trend);

```

```

66 *   lags trace the dynamic effect.
67 *-----
68 forvalues k = -5/6 {
69     if 'k'==-1 continue           // omit t = -1: the single reference period
70     local j = 'k' + 6           // 1..12 index -> a legal variable name
71     generate evt'j' = (treat==1 & rel=='k')
72     label variable evt'j' "t = 'k'"
73 }
74 reghdfe y evt*, absorb(firm year) vce(cluster firm)
75
76 coefplot, keep(evt*) vertical yline(0) xline(4.5, lpattern(dash)) ///
77     coeflabels(evt1="-5" evt2="-4" evt3="-3" evt4="-2" evt6="0" evt7="1" evt8="2" ///
78         evt9="3" evt10="4" evt11="5" evt12="6") ///
79     title("Event-study estimates (ref: t = -1)") ///
80     xtitle("Event time (years relative to treatment)") ///
81     ytitle("Effect on outcome")
82 graph export "output/04_event_study.png", replace width(1400)
83
84 log close
85 display "04_panel_did_eventstudy.do finished OK"

```

5.3 The results

The three columns apply the same DiD interaction $\text{treat} \times \text{post}$ under progressively richer fixed effects. In Stata's factor-variable notation $i.\text{treat}\#i.\text{post}$ expands to both main effects *and* their product, while $1.\text{treat}\#1.\text{post}$ is the product alone, so the columns estimate

- (1) $y_{it} = \beta_0 + \beta_1 \text{treat}_i + \beta_2 \text{post}_t + \beta (\text{treat}_i \times \text{post}_t) + \varepsilon_{it}$,
- (2) $y_{it} = \alpha_i + \beta_2 \text{post}_t + \beta (\text{treat}_i \times \text{post}_t) + \varepsilon_{it}$,
- (3) $y_{it} = \alpha_i + \gamma_t + \beta (\text{treat}_i \times \text{post}_t) + \varepsilon_{it}$,

for (1) pooled OLS, (2) firm fixed effects (`xtreg, fe`), and (3) two-way firm-and-year fixed effects (`reghdfe`). Each richer specification absorbs a main effect: the firm effect α_i soaks up the time-invariant treat_i , so column 2 drops it, and the year effect γ_t soaks up post_t , so column 3 drops it too. That is why the code switches to the bare interaction $1.\text{treat}\#1.\text{post}$ in column 3 — the main effects are now collinear with the fixed effects and cannot be estimated separately. The DiD estimand β survives throughout, and the simulated post-treatment effect (about 1.5) is recovered by all three columns (Table). The event study shows a flat pre-trend and a rising dynamic effect (Figure).

Table 5.1: DiD estimates

	(1)	(2)	(3)
	Pooled OLS	Firm FE	Two-way FE
$\text{treat}=1 \times \text{post}=1$	1.973*** (0.104)	1.973*** (0.104)	1.973*** (0.104)
Observations	1800	1800	1800

Standard errors in parentheses

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

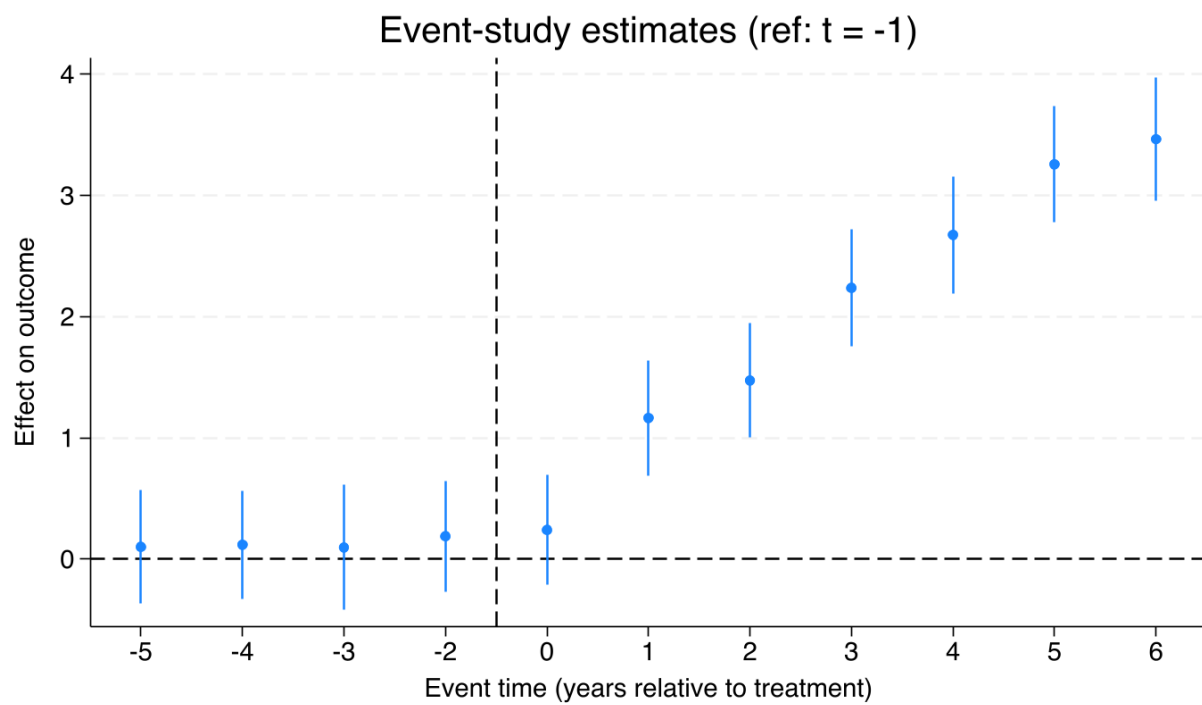


Figure 5.1: Event-study estimates: flat, insignificant leads; rising lags.

Chapter 6

Time series

↔ the method, explained: [Appendix B.3](#)

6.1 The method

A first-order autoregression with a level break,

$$y_t = c + \phi y_{t-1} + \delta \mathbf{1}[t \geq t^*] + \varepsilon_t,$$

declared with `tsset` so that the lag (L.), difference (D.), and lead (F.) operators work. [Newey and West \(1987\)](#) HAC standard errors are robust to serial correlation up to a chosen lag. The simulated series uses $\phi = 0.6$ with a +3 shift in 2015m1.

In practice. Persistence and HAC standard errors are everywhere in macro and finance — modelling inflation or output dynamics, volatility clustering in asset returns, and any regression whose residuals are serially correlated, where Newey–West is the default fix.

6.2 The code

The workflow. Declare the series with `tsset`, apply lag and difference operators, inspect the autocorrelation, and fit a lagged regression with Newey–West (HAC) standard errors.

```
1  *=====
2  * 05_timeseries.do --- tsset, lag/diff operators, autocorrelation,
3  *                      HAC (Newey-West) SE, and a simple event/break window
4  *-----
5  * Simulates a monthly AR(1) series with a level shift in 2015m1.
6  *=====
7  version 17
8  clear all
9  set more off
10 set linesize 90
11 set seed 4321
12 capture log close
13 capture mkdir logs // create the log folder if absent (gitignored)
14 log using "logs/05_timeseries.log", replace text
15
16 *-----
17 * 1. Build a monthly time series and declare it with tsset.
18 *-----
19 set obs 120
20 generate t = _n
21 generate mdate = tm(2010m1) + t - 1 // Stata monthly date
22 format mdate %tm
23 tsset mdate // declare time-series structure
```

```

24
25 generate eps = rnormal(0,1)
26 generate y = eps in 1
27 replace y = 0.6*L.y + eps in 2/L           // AR(1): recursion fills in order
28 replace y = y + 3 if mdate >= tm(2015m1) // structural level shift
29 label variable y "Simulated series"
30
31 *-----
32 * 2. Time-series operators: L. (lag), D. (difference), F. (lead).
33 *-----
34 generate dy = D.y                          // first difference
35 generate y_l1 = L.y                        // first lag
36 summarize y dy
37
38 *-----
39 * 3. Plots: the series and its autocorrelation function.
40 *-----
41 tsline y, title("Monthly series with a 2015 level shift") ///
42     tline(2015m1, lpattern(dash))
43 graph export "output/05_tsline.png", replace width(1400)
44
45 ac y, title("Autocorrelation function")
46 graph export "output/05_acf.png", replace width(1200)
47
48 *-----
49 * 4. Regression with a lag + HAC (Newey-West) standard errors.
50 *-----
51 newey y L.y, lag(3)                       // robust to serial correlation up to 3 lags
52
53 *-----
54 * 5. Test the level break with an indicator (a mini event study).
55 *-----
56 generate byte post = mdate >= tm(2015m1)
57 regress y L.y i.post, vce(robust)         // coefficient on post = estimated jump
58
59 log close
60 display "05_timeseries.do finished OK"

```

6.3 The results

The two regressions map directly onto the time-series operators that `tsset` unlocks:

$$\begin{aligned} \text{newey y L.y, lag(3)} : & y_t = c + \phi y_{t-1} + \varepsilon_t, \\ \text{regress y L.y i.post} : & y_t = c + \phi y_{t-1} + \delta \mathbf{1}[t \geq 2015m1] + \varepsilon_t, \end{aligned}$$

where the lag operator `L.y` is y_{t-1} and `i.post` is the break indicator $\mathbf{1}[t \geq 2015m1]$. The first fits the pure AR(1) with Newey–West HAC standard errors (robust to serial correlation up to three lags); the second adds the level-shift dummy, whose coefficient δ is the estimated jump. The simulated series is visibly persistent — its autocorrelation function decays slowly rather than dropping to zero — and the +3 level shift at 2015m1 is recovered by the break indicator ($\hat{\delta} \approx 3$).

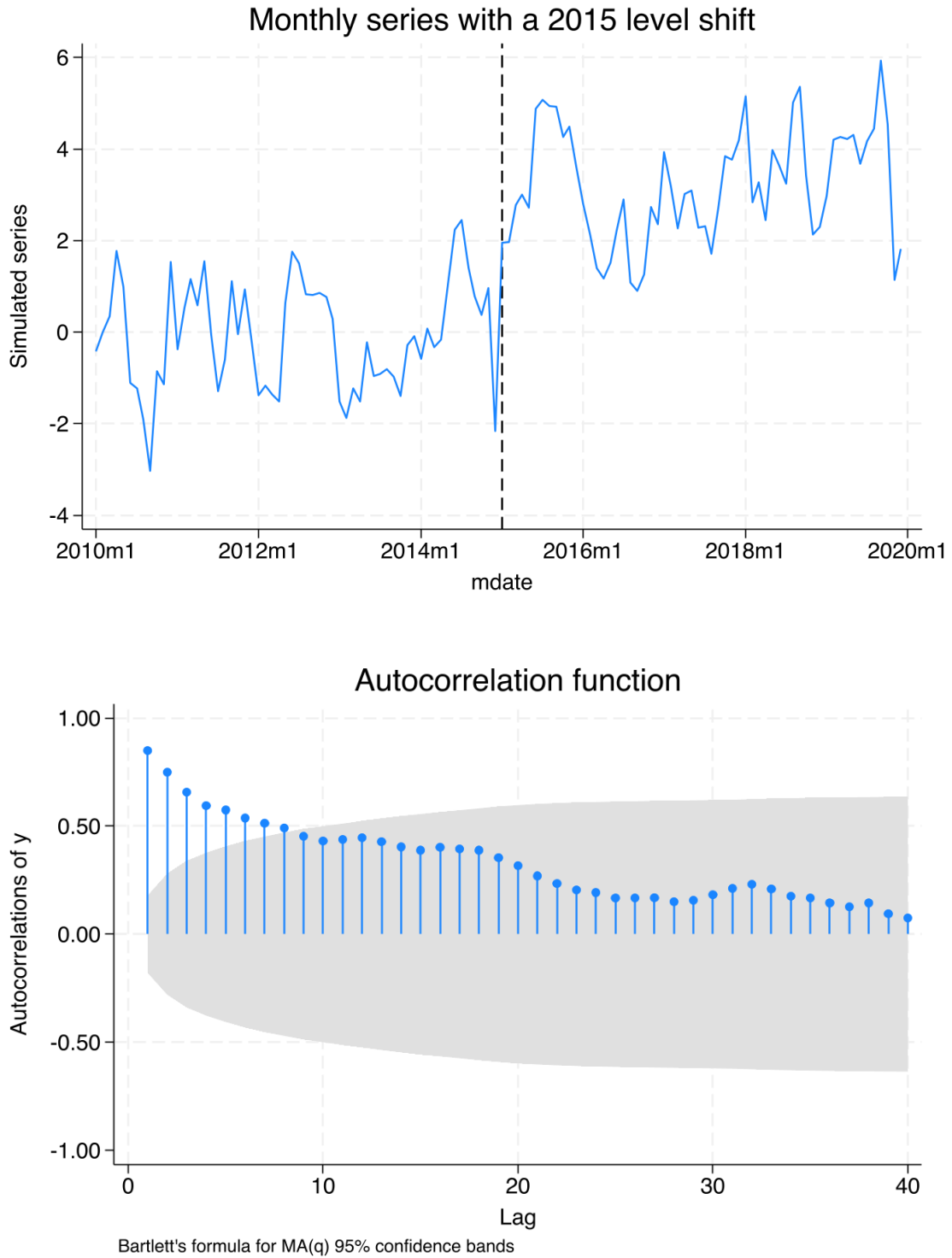


Figure 6.1: The simulated series with its 2015 level shift (top) and its autocorrelation function (bottom).

Chapter 7

Importing public data

Most empirical projects begin by getting public data into Stata. This script first round-trips a CSV — writing Stata’s built-in auto data out and reading it straight back, to demonstrate the `export/import` delimited commands — then imports a *real* FRED series, the U.S. unemployment rate (UNRATE), *live* through Stata’s native importer (`import fred`): it parses the date, declares a monthly time series, and plots it. The importer needs a free key, set once with `set fredkey` (which keeps it in the Stata config, not this script); a previously downloaded CSV is the offline fallback.

Setting up the FRED key. A one-time setup, so `import fred` works:

1. Register for a free API key at fredaccount.stlouisfed.org/apikeys.
2. In Stata, run `set fredkey YOUR_KEY`, permanently *once*. The key is stored in the Stata configuration — it persists across sessions and never appears in a do-file.
3. Thereafter `import fred UNRATE` (or any series ID, e.g. GDPC1, CPIAUCSL) pulls that series live; rerun step 2 to replace the key.

```
1 set fredkey YOUR_KEY, permanently // once: free key at fredaccount.stlouisfed.org
2 import fred UNRATE, clear // then pull any series live from FRED's API
```

In practice. The same pattern pulls in most public macro and finance data — FRED (Federal Reserve series), SEC EDGAR (filings), the World Bank and IMF, the BLS, and exchange or market feeds — which is where almost every empirical project begins.

7.1 The code

The workflow. Import and export delimited files, and pull a real public series (unemployment, from FRED) into Stata.

```
1 *****
2 * 06_import_public_data.do --- getting public data INTO Stata
3 *****
4 * Empirical work often requires scraping publicly available data and
5 * cleaning it. This shows the three routes used in practice.
6 *****
7 version 17
8 clear all
9 set more off
10 set linesize 90
11 capture log close
12 capture mkdir logs // create the log folder if absent (gitignored)
13 log using "logs/06_import_public_data.log", replace text
14
```

```

15 *-----
16 * (A) Round-trip a CSV to demonstrate import/export delimited (the CSV workhorses).
17 *   We write the built-in 'auto' data out and read it straight back --- auto is just
18 *   a throwaway here, the point is the file I/O, not the data. varnames(1)=header.
19 *-----
20 sysuse auto, clear
21 keep make price mpg foreign
22 export delimited using "output/_cars.csv", replace
23 clear
24 import delimited "output/_cars.csv", varnames(1) clear
25 describe
26 list in 1/3
27 erase "output/_cars.csv"
28
29 *-----
30 * (B) Import a REAL public series: the U.S. unemployment rate (FRED: UNRATE).
31 *   PRIMARY route --- Stata's native FRED importer, a live API call (the most
32 *   reproducible: no manual download step). It needs a free key, set ONCE with
33 *   set fredkey YOUR_KEY, permanently (free key: fredaccount.stlouisfed.org)
34 *   The key then lives in the Stata config, never in this script. With no key
35 *   set we fall back to a CSV previously downloaded from FRED, so this always runs.
36 *-----
37 capture import fred UNRATE, clear
38 if _rc == 0 {
39     display "Live FRED API pull succeeded: " _N " observations."
40     rename UNRATE unrate // the native importer keeps FRED's case
41     generate mdate = mofd(daten) // import fred already provides daten (%td)
42 }
43 else {
44     display "No FRED key configured (rc = " _rc "); using the bundled CSV instead."
45     import delimited "data/UNRATE.csv", varnames(1) clear
46     rename observation_date datestr
47     generate daten = date(datestr, "YMD") // parse the YYYY-MM-DD string ourselves
48     format daten %td
49     generate mdate = mofd(daten) // collapse to monthly
50 }
51 format mdate %tm
52 tsset mdate
53 summarize unrate
54 list datestr unrate in 1/3
55 list datestr unrate in -3/L // most recent observations
56
57 tsline unrate, title("U.S. unemployment rate (FRED: UNRATE)") ///
58     ytitle("Percent") xtitle("")
59 graph export "output/06_unrate.png", replace width(1400)
60 save "output/fred_unrate.dta", replace
61 display "FRED UNRATE imported: " _N " monthly observations"
62
63 * To pull several series at once over a date window:
64 * import fred UNRATE GDPC1, daterange(2000-01-01 .) aggregate(monthly)
65
66 log close
67 display "06_import_public_data.do finished OK"

```

7.2 The results

This is the one genuinely real series in the book — pulled *live* from FRED's API by `import fred`:

```
. import fred UNRATE, clear
Series ID      Nobs   Date range      Frequency
UNRATE        940    1948-01-01 to 2026-05-01  Monthly

Live FRED API pull succeeded: 941 observations.
```

The unemployment rate climbs in every postwar recession and spikes to nearly 15% in April 2020 at the onset of the COVID-19 pandemic, then falls back — a reminder that importing and cleaning data is the gateway to any empirical analysis.

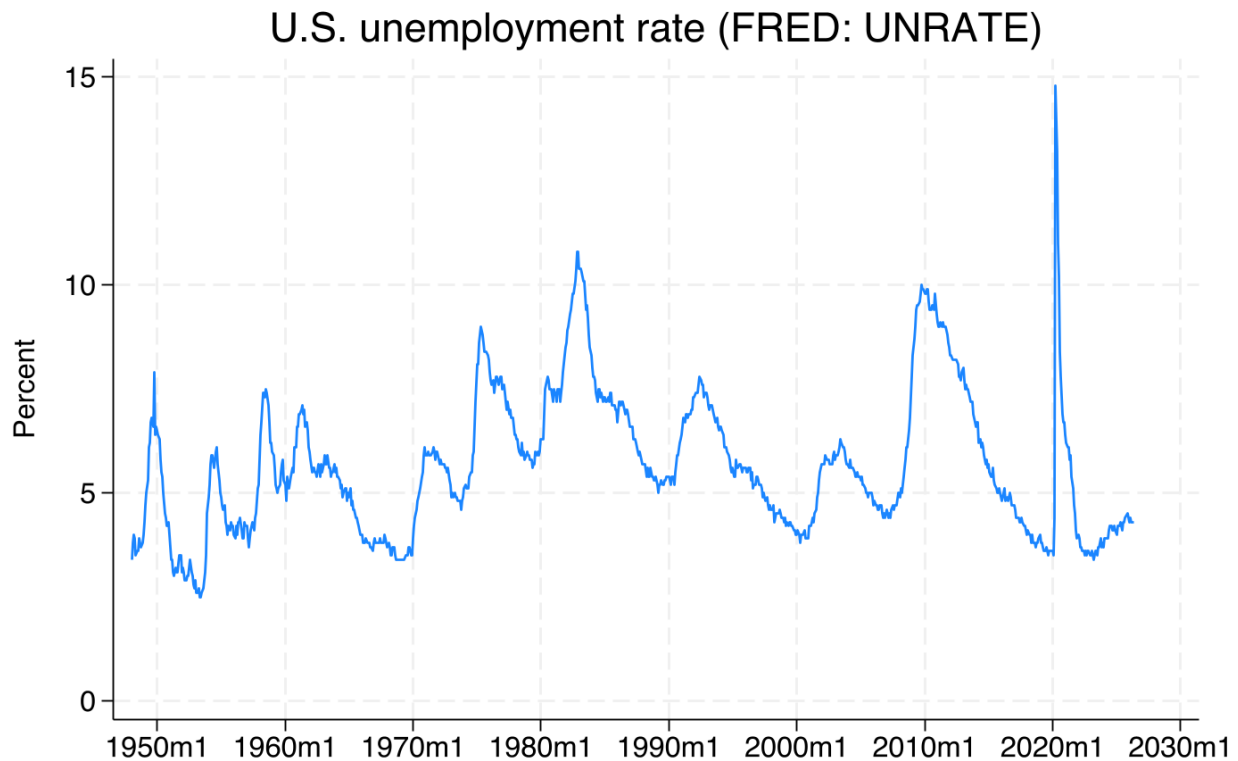


Figure 7.1: U.S. unemployment rate, 1948–2026 (FRED: UNRATE), imported and plotted in Stata.

Part III

The causal-inference toolkit

Chapter 8

Instrumental variables (2SLS)

↔ the method, explained: [Appendix B.4](#)

8.1 The method

When a regressor is correlated with the error, $\text{Cov}(x_i, \varepsilon_i) \neq 0$, OLS is inconsistent. A valid instrument z is relevant ($\text{Cov}(z, x) \neq 0$) and exogenous ($\text{Cov}(z, \varepsilon) = 0$) — think of a lottery that shifts who gets treated but is otherwise unrelated to the outcome. With one endogenous regressor and two instruments the model is

$$\underbrace{y = \beta_0 + \beta_1 x + \varepsilon,}_{\text{structural equation}} \quad \underbrace{x = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + v,}_{\text{first stage}}$$

written ($x = z_1 \ z_2$) in Stata. Two-stage least squares is

$$\hat{\beta}_{2SLS} = (X'P_Z X)^{-1} X'P_Z y, \quad P_Z = Z(Z'Z)^{-1} Z';$$

this is GMM with the weight matrix $(Z'Z)^{-1}$. Efficient two-step GMM (`ivregress gmm`) instead uses the optimal heteroskedasticity-robust weight matrix, gaining efficiency when the errors are heteroskedastic — and coinciding with 2SLS when they are not. `ivreghdfe` adds high-dimensional fixed effects. The standard checks are the **first-stage F** (weak instruments; [Staiger and Stock, 1997](#)), the **Durbin–Wu–Hausman** test (endogeneity), and the **Hansen J** (overidentification).

In practice. Famous instruments include quarter of birth for the returns to schooling ([Angrist and Krueger 1991](#)), the Vietnam draft lottery for veteran earnings ([Angrist 1990](#)), and cost shifters for demand estimation; modern work often uses judge- or examiner-leniency designs.

8.2 The code

The workflow. Contrast biased OLS with 2SLS, run the weak-instrument, endogeneity, and over-identification diagnostics, and refit with high-dimensional fixed effects via `ivreghdfe`.

```
1  *=====
2  * 07_iv_2sls.do --- instrumental variables / 2SLS, with diagnostics, and the
3  *                    high-dimensional-FE version (ivreghdfe)
4  *-----
5  * Simulates an endogenous regressor x correlated with an unobserved confounder
6  * u; two valid instruments z1, z2 (independent of u). True slope on x is 2.0,
7  * so OLS should be biased and 2SLS should recover ~2.0.
8  *=====
9  version 17
10 clear all
11 set more off
```

```

12 set linesize 90
13 set seed 1234
14 capture log close
15 capture mkdir logs // create the log folder if absent (gitignored)
16 log using "logs/07_iv_2sls.log", replace text
17
18 set obs 5000
19 generate id = _n
20 generate z1 = rnormal() // instrument 1
21 generate z2 = rnormal() // instrument 2
22 generate u = rnormal() // unobserved confounder
23 generate x = 0.6*z1 + 0.5*z2 + 0.6*u + rnormal() // endogenous regressor
24 generate e = 0.6*u + rnormal() // structural error (corr with x via u)
25 generate grp = mod(id,50) + 1 // 50 groups, for the FE demo
26 bysort grp (id): generate gfe = rnormal() if _n==1
27 bysort grp (id): replace gfe = gfe[1]
28 generate y = 1 + 2*x + e + gfe // TRUE slope on x = 2
29
30 *-----
31 * 1. OLS (biased) vs. 2SLS (base Stata: ivregress).
32 *-----
33 eststo clear
34 eststo ols: regress y x, vce(robust) // biased upward
35 eststo tsls: ivregress 2sls y (x = z1 z2), vce(robust) first // first-stage shown
36
37 *-----
38 * 2. IV diagnostics --- the standard validity checks for any IV design.
39 *-----
40 estat firststage // weak-instrument F (rule of thumb: > 10)
41 estat endogenous // Durbin-Wu-Hausman: is x actually endogenous?
42 estat overid // over-identification test (2 instruments, 1 endog var)
43
44 *-----
45 * 2b. Efficient (two-step) GMM. 2SLS is GMM with a fixed weight matrix; 'gmm' uses
46 * the optimal heteroskedasticity-robust weight matrix, so it is more efficient
47 * when errors are heteroskedastic. With homoskedastic errors (as here) it
48 * essentially coincides with 2SLS; the over-id test is Hansen's J.
49 *-----
50 ivregress gmm y (x = z1 z2), vce(robust)
51
52 *-----
53 * 3. IV with high-dimensional fixed effects (ivreghdfe = ivreg2 + reghdfe).
54 * Absorbs the group FE that base ivregress would need as dummies.
55 *-----
56 eststo ivfe: ivreghdfe y (x = z1 z2), absorb(grp) vce(robust)
57
58 *-----
59 * 4. Compare the slope on x across estimators.
60 *-----
61 esttab ols tsls ivfe, b(%9.3f) se star(* 0.10 ** 0.05 *** 0.01) ///
62 keep(x) mtitles("OLS (biased)" "2SLS" "IV + FE") ///
63 title("True slope on x = 2.0: OLS is biased, IV recovers it")
64 esttab ols tsls ivfe using "output/07_iv_table.tex", replace booktabs ///
65 b(%9.3f) se star(* 0.10 ** 0.05 *** 0.01) keep(x) ///
66 mtitles("OLS" "2SLS" "IV+FE") title("IV vs OLS")
67
68 *-----
69 * 5. Export an IV-diagnostics table for the reference (re-fit quietly first; the
70 * last active estimates above are the IV+FE model, not the plain 2SLS).

```

```

71 *-----
72 quietly regress x z1 z2, vce(robust)
73 test z1 z2
74 local fsF = r(F)
75 quietly ivregress 2sls y (x = z1 z2), vce(robust)
76 quietly estat endogenous
77 local dwh = r(r_score)
78 local dwhp = r(p_r_score)
79 quietly estat overid
80 local oj = r(score)
81 local ojpt = r(p_score)
82 local dwhpt = cond('dwhp' < 0.001, "\(<)0.001", string('dwhp', "%5.3f"))
83 local ojpt = string('ojpt', "%5.3f")
84
85 capture file close it
86 file open it using "output/07_iv_diag_table.tex", write replace
87 file write it "\begin{table}[htbp]\centering" _n
88 file write it "\caption{IV diagnostics (2SLS, heteroskedasticity-robust)}" _n
89 file write it "\begin{tabular}{lcc}" _n "\toprule" _n
90 file write it "Test & Statistic & \(\rho\) \\" _n "\midrule" _n
91 file write it "First-stage \(\rho\) (weak instruments) & " %6.1f ('fsF') " & --- \\" _n
92 file write it "Endogeneity: Durbin--Wu--Hausman (\(\chi^2_{1}\)) & " %6.1f ('dwh') " & 'dwhpt' \\"
93 _n
94 file write it "Overidentification: robust score (\(\chi^2_{1}\)) & " %5.2f ('oj') " & 'ojpt' \\" _n
95 file write it "\bottomrule" _n "\end{tabular}" _n "\end{table}" _n
96
97 log close
98 display "07_iv_2sls.do finished OK"

```

8.3 The results

The true slope is 2.0. The three columns are three estimators of the *same* structural equation, differing only in how the endogenous x enters:

- | | | |
|--------------|---|------------------------------------|
| (1) OLS: | $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ | (raw $x \rightarrow$ biased), |
| (2) 2SLS: | $y_i = \beta_0 + \beta_1 \hat{x}_i + \varepsilon_i$ | (\hat{x} from the first stage), |
| (3) 2SLS+FE: | $y_i = \alpha_{g(i)} + \beta_1 \hat{x}_i + \varepsilon_i$ | (α_g absorbed), |

fit by `regress` y x , `ivregress 2sls` y ($x = z1\ z2$), and `ivreghdfe` with `absorb(grp)`. The two IV estimators replace the contaminated x with the fitted value \hat{x} from the first-stage regression of x on the instruments,

$$\hat{x}_i = \hat{\pi}_0 + \hat{\pi}_1 z1_i + \hat{\pi}_2 z2_i.$$

OLS is biased upward because $\text{Cov}(x, \varepsilon) > 0$ through the unobserved confounder; both IV estimators are consistent and recover the truth. Efficient GMM (`ivregress gmm`) returns $\hat{\beta}_1 \approx 1.96$, indistinguishable from 2SLS here because the simulated errors are homoskedastic — the optimal weight matrix buys nothing when 2SLS is already efficient.

The diagnostics confirm the design: the first-stage F is far above the weak-instrument rule-of-thumb of 10; the **Durbin–Wu–Hausman** test strongly rejects exogeneity, so instrumenting is warranted; and the over-identification test does *not* reject, consistent with both instruments being valid.

Table 8.1: IV vs OLS

	(1) OLS	(2) 2SLS	(3) IV+FE
x	2.173*** (0.015)	1.960*** (0.027)	1.994*** (0.020)
N	5000	5000	5000

Standard errors in parentheses
 * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table 8.2: IV diagnostics (2SLS, heteroskedasticity-robust)

Test	Statistic	p
First-stage F (weak instruments)	1243.7	—
Endogeneity: Durbin–Wu–Hausman (χ^2_1)	95.6	<0.001
Overidentification: robust score (χ^2_1)	1.48	0.224

Chapter 9

Binary outcome models

↔ the method, explained: [Appendix B.5](#)

9.1 The method

For a 0/1 outcome,

$$P(y_i = 1 | x_i) = G(x_i'\beta),$$

with $G = \Lambda$ (logit) or $G = \Phi$ (probit); the linear probability model is OLS on the indicator. Raw coefficients are not comparable across these — the comparable quantity is the average marginal effect

$$\text{AME}_j = \frac{1}{n} \sum_i g(x_i'\beta) \beta_j,$$

which is nearly identical for logit and probit. Logit odds ratios are e^{β_j} . Fit is summarised by a classification table and the area under the ROC curve.

In practice. Binary-outcome models fit any yes/no decision — loan default and credit approval, labour-force participation, mortgage acceptance, or vote choice — reporting the average marginal effects rather than the raw coefficients.

9.2 The code

The workflow. Fit the linear-probability, logit, and probit models, read the odds ratios, compute average marginal effects with `margins`, and plot the predicted probabilities.

```
1  *=====
2  * 08_logit_probit.do --- binary-outcome models and how to interpret them
3  *-----
4  * Models the probability a car is foreign as a function of mpg, weight, price.
5  * The key lesson: logit/probit coefficients are NOT directly interpretable ---
6  * use margins for average marginal effects and predicted probabilities.
7  *=====
8  version 17
9  clear all
10 set more off
11 set linesize 90
12 capture log close
13 capture mkdir logs // create the log folder if absent (gitignored)
14 log using "logs/08_logit_probit.log", replace text
15
16 sysuse auto, clear
17 generate weight_t = weight/1000
18 generate price_k = price/1000
19 label variable weight_t "Weight (1000 lbs)"
```

```

20 label variable price_k "Price ($1000s)"
21
22 *-----
23 * 1. Three ways to model a 0/1 outcome: LPM, logit, probit.
24 *-----
25 eststo clear
26 eststo lpm: regress foreign mpg weight_t price_k, vce(robust) // linear prob.
27 eststo logit: logit foreign mpg weight_t price_k, vce(robust)
28 eststo probit: probit foreign mpg weight_t price_k, vce(robust)
29 * NOTE: raw coefficients are on different scales --- do NOT compare directly.
30 esttab lpm logit probit, b(%9.3f) se star(* 0.10 ** 0.05 *** 0.01) ///
31 mtitles("LPM" "Logit" "Probit") label ///
32 title("Coefficients differ by scale --- compare marginal effects, not these")
33
34 *-----
35 * 2. Odds ratios (logit) --- a common reporting convention.
36 *-----
37 logit foreign mpg weight_t price_k, or
38
39 *-----
40 * 3. Average marginal effects --- the comparable, interpretable quantities.
41 * Logit and probit AMEs are typically very close.
42 *-----
43 quietly logit foreign mpg weight_t price_k
44 margins, dydx(*) post
45 estimates store ame_logit
46 quietly probit foreign mpg weight_t price_k
47 margins, dydx(*) post
48 estimates store ame_probit
49 esttab ame_logit ame_probit, b(%9.4f) se ///
50 mtitles("Logit AME" "Probit AME") title("Average marginal effects (compare these)")
51
52 *-----
53 * 4. Predicted probabilities + a plot across the range of mpg.
54 *-----
55 quietly logit foreign mpg weight_t price_k
56 margins, at(mpg=(12(4)40))
57 marginsplot, title("Predicted P(foreign) across mpg") ///
58 ytitle("Predicted probability")
59 graph export "output/08_logit_pr.png", replace width(1300)
60
61 *-----
62 * 5. In-sample fit: classification table and ROC area.
63 *-----
64 quietly logit foreign mpg weight_t price_k
65 estat classification // sensitivity/specificity at p=0.5
66 lroc, nograph // area under ROC curve
67
68 log close
69 display "08_logit_probit.do finished OK"

```

9.3 The results

The three columns fit the same linear index $x'_i\beta = \beta_0 + \beta_1 \text{mpg} + \beta_2 \text{weight_t} + \beta_3 \text{price_k}$ through three different link functions:

- (1) LPM: $P(\text{foreign}_i = 1 | x_i) = x'_i\beta,$
- (2) logit: $P(\text{foreign}_i = 1 | x_i) = \Lambda(x'_i\beta),$
- (3) probit: $P(\text{foreign}_i = 1 | x_i) = \Phi(x'_i\beta),$

fit by `regress`, `logit`, and `probit` respectively, with Λ the logistic CDF and Φ the standard-normal CDF. Because the link functions differ, the raw coefficients are not comparable across columns — the comparable quantity is the average marginal effect (`margins`, `dydx(*)`). The logit and probit average marginal effects coincide to the third decimal — as expected, the two link functions imply nearly the same marginal effects — and in-sample fit is high (area under the ROC curve ≈ 0.96). The plot shows the predicted probability falling smoothly with mileage.

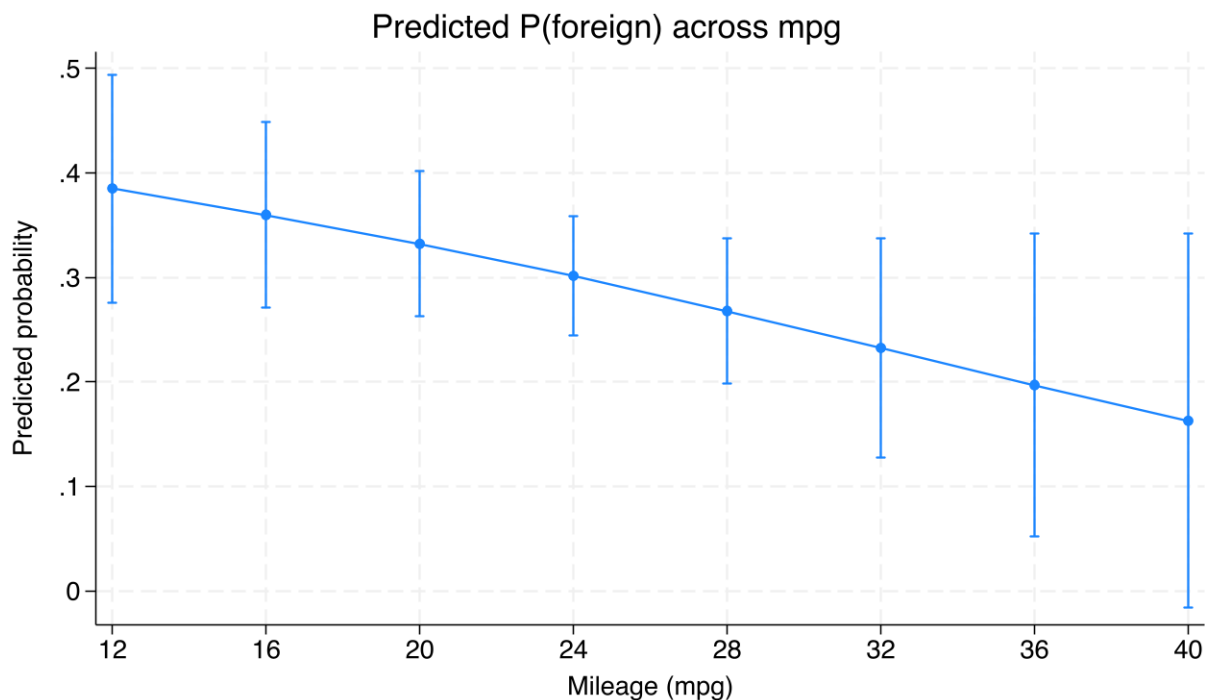


Figure 9.1: Predicted probability that a car is foreign, across mileage.

Chapter 10

Staggered-adoption DiD

↔ the method, explained: [Appendix B.6](#)

10.1 The method

When units are treated at different times (say, states adopting a policy in different years) and effects are dynamic, the two-way FE coefficient is a non-convex weighted average of effects and can be badly biased ([Goodman-Bacon, 2021](#)). [Callaway and Sant'Anna \(2021\)](#) instead estimate clean group-time average treatment effects,

$$ATT(g, t) = E[Y_t - Y_{g-1} | G = g] - E[Y_t - Y_{g-1} | C],$$

using not-yet- or never-treated units C as comparisons, via a doubly-robust estimator. These are then aggregated (overall, by event time, by cohort, by calendar period).

In practice. This is the tool for staggered policy adoption — minimum-wage changes, Medicaid expansion, or marijuana legalization rolling out across states in different years — the very setting where naive two-way FE goes wrong.

10.2 The code

The workflow. Build a staggered-adoption panel, estimate clean group-time effects with `csdid`, aggregate them, and contrast with the biased naive two-way fixed-effects estimate.

```
1  *=====
2  * 09_staggered_did_csdid.do --- staggered-adoption DiD (Callaway & Sant'Anna)
3  *-----
4  * When units are treated at DIFFERENT times, two-way FE DiD can be biased
5  * (the "forbidden comparison" / negative-weighting problem). Callaway &
6  * Sant'Anna (2021) estimate clean group-time ATTs and aggregate them.
7  *
8  * Simulates 3 cohorts: first-treated in 2014, in 2017, and never-treated.
9  * True effect is 0 pre-treatment and grows after each cohort's own start.
10 *-----
11 version 17
12 clear all
13 set more off
14 set linesize 90
15 set seed 555
16 capture log close
17 capture mkdir logs // create the log folder if absent (gitignored)
18 log using "logs/09_staggered_did_csdid.log", replace text
19
20 *-----
21 * 1. Build a staggered panel: 300 units x 12 years (2010-2021).
```

```

22 *   gvar = year of first treatment; 0 = never treated (csdid convention).
23 *-----
24 set obs 300
25 generate id = _n
26 generate g = 0 // never-treated by default
27 replace g = 2014 if id <= 100 // early cohort
28 replace g = 2017 if id > 100 & id <= 200 // late cohort
29 generate u_id = rnormal()*1.5 // unit fixed effect
30 expand 12
31 bysort id: generate year = 2009 + _n // 2010..2021
32
33 generate att = 0
34 replace att = 0.4*(year - g + 1) if g>0 & year>=g // dynamic, grows post-onset
35 generate yr_fe = 0.2*(year-2015) // common year trend
36 generate y = 3 + u_id + yr_fe + att + rnormal()
37 label variable y "Outcome"
38
39 tabulate g, missing // cohort sizes (in unit-years)
40
41 *-----
42 * 2. Callaway & Sant'Anna estimator (doubly robust). No covariates here.
43 *-----
44 csdid y, ivar(id) time(year) gvar(g) method(dripw)
45
46 *-----
47 * 3. Aggregations of the group-time ATTs.
48 *-----
49 estat simple // single overall ATT
50 matrix _Tsimple = r(table) // capture the aggregate ATT now, before
51 // estat group/event overwrite r(table)
52 estat group // ATT by treatment cohort
53 estat event // dynamic / event-study path (LAST, so
54 // csdid_plot below plots THIS)
55
56 *-----
57 * 4. Event-study plot of the dynamic aggregation.
58 * csdid_plot graphs the most recent estat aggregation -> run estat event last.
59 *-----
60 csdid_plot, title("Callaway-Sant'Anna event study") ///
61 ytitle("ATT") xtitle("Event time (years relative to treatment)")
62 graph export "output/09_csdid_event.png", replace width(1400)
63
64 *-----
65 * 5. For contrast: the naive two-way FE estimate on the same data.
66 * (Can be biased under staggered timing + dynamic effects --- that's the point.)
67 *-----
68 generate byte treated_now = (g>0 & year>=g)
69 reghdfe y treated_now, absorb(id year) vce(cluster id)
70
71 *-----
72 * 6. Export a table contrasting the Callaway-Sant'Anna ATT with the naive TWFE
73 * estimate (the bias this method removes).
74 *-----
75 local twfe = _b[treated_now]
76 local twse = _se[treated_now]
77 local csatt = _Tsimple[1,1]
78 local csse = _Tsimple[2,1]
79
80 capture file close ct

```

```

81 file open ct using "output/09_csdid_table.tex", write replace
82 file write ct "\begin{table}[htbp]\centering" _n
83 file write ct "\caption{Overall ATT: Callaway--Sant'Anna vs.\ naive two-way FE}" _n
84 file write ct "\begin{tabular}{lcc}" _n "\toprule" _n
85 file write ct "Estimator & ATT & Std.\ err. \\" _n "\midrule" _n
86 file write ct "Callaway--Sant'Anna (\pkg{csdid}) & " %5.3f ('csatt') " & " %5.3f ('csse') " \\"
      _n
87 file write ct "Naive two-way FE (\pkg{reghdfe}) & " %5.3f ('twfe') " & " %5.3f ('twse') " \\" _n
88 file write ct "\bottomrule" _n "\end{tabular}" _n "\end{table}" _n
89 file close ct
90
91 log close
92 display "09_staggered_did_csdid.do finished OK"

```

10.3 The results

The two rows of the table are two estimators on the *same* staggered panel. `csdid` — with `gvar(g)` passing each unit's treatment cohort (g = year of first treatment, 0 = never treated) — aggregates the clean group-time effects of §The method, whereas the naive two-way FE regresses the outcome on a single post-onset indicator:

$$\begin{aligned} \text{Callaway--Sant'Anna: } & \widehat{ATT} \text{ aggregated from } ATT(g, t), \\ \text{naive two-way FE: } & y_{it} = \alpha_i + \gamma_t + \beta^{\text{TWFE}} \text{treated_now}_{it} + \varepsilon_{it}. \end{aligned}$$

The single coefficient β^{TWFE} is the non-convex weighted average that misleads under staggered timing. The estimator returns a flat pre-trend (no anticipation) and a rising post-treatment path that tracks the planted dynamic effect. The overall ATT recovers the simulated average; the naive two-way FE estimate on the *same* data is substantially smaller — the negative-weighting bias this method is designed to remove.

Table 10.1: Overall ATT: Callaway--Sant'Anna vs. naive two-way FE

Estimator	ATT	Std. err.
Callaway--Sant'Anna (<code>csdid</code>)	1.495	0.112
Naive two-way FE (<code>reghdfe</code>)	1.009	0.073

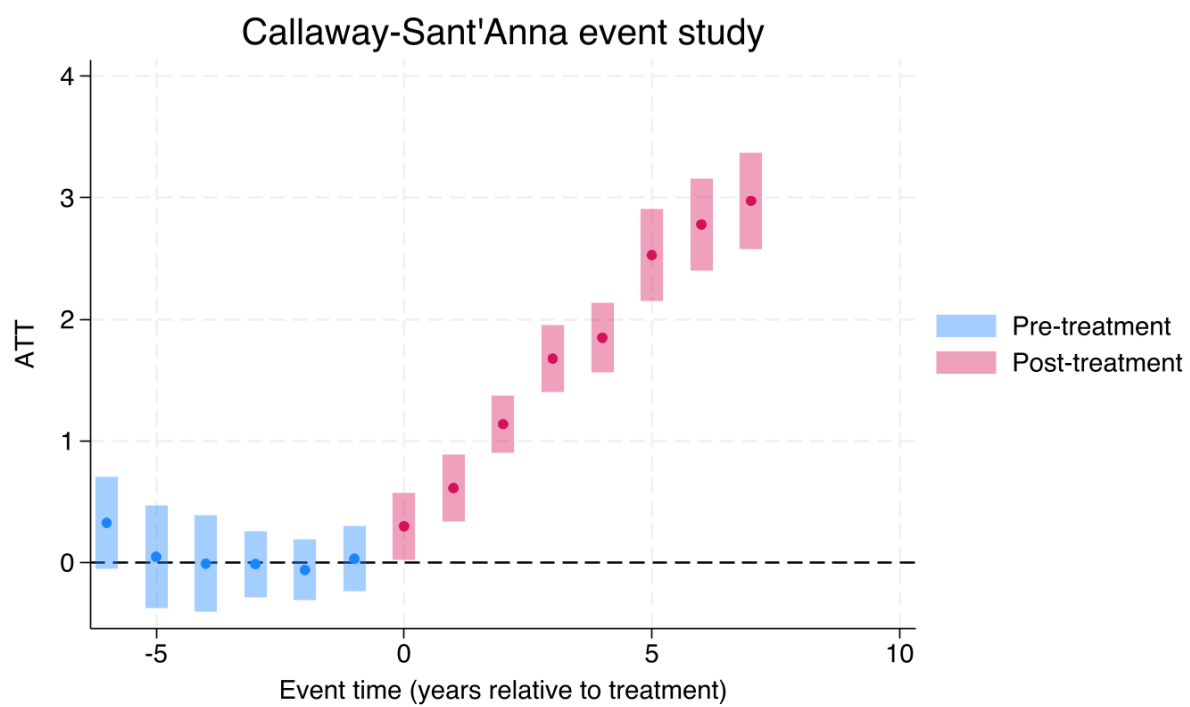


Figure 10.1: Callaway–Sant’Anna event study: pre-treatment (covering zero) vs. post-treatment.

Chapter 11

Regression discontinuity

↔ the method, explained: [Appendix B.7](#)

11.1 The method

When treatment switches at a cutoff of a running variable (say, a scholarship awarded once a test score clears a threshold), $d_i = \mathbf{1}[x_i \geq c]$, the sharp RD estimand is the jump in the conditional mean,

$$\tau = \underbrace{\lim_{x \downarrow c} E[Y | X = x]}_{\text{just ABOVE the cutoff}} - \underbrace{\lim_{x \uparrow c} E[Y | X = x]}_{\text{just BELOW the cutoff}},$$

estimated by local-linear regression each side of c within an MSE-optimal bandwidth, with [Calonico, Cattaneo and Titiunik \(2014\)](#) bias-corrected robust confidence intervals. A **manipulation test** checks that the density of X is continuous at c ([Cattaneo, Jansson and Ma, 2020](#)).

In practice. Sharp RD exploits threshold rules: scholarship and test-score cutoffs (its original use, [Thistlethwaite and Campbell 1960](#)), class-size caps (the “Maimonides rule” of [Angrist and Lavy 1999](#)), vote-share thresholds for incumbency in close elections ([Lee 2008](#)), and program-eligibility cutoffs.

11.2 The code

The workflow. Simulate a sharp-RD dataset, estimate the jump with `rdrobust`, select the bandwidth `rdbwselect`, draw the annotated plot `rdplot`, test for manipulation with `rddensity`, and export the results table.

```
1  *=====
2  * 10_rdd_regression_discontinuity.do --- sharp regression discontinuity (RD)
3  *-----
4  * Units with running variable x >= 0 get treated. The outcome is a smooth
5  * function of x with a JUMP of 2.0 at the cutoff = the treatment effect.
6  * Uses Calonico-Cattaneo-Titiunik local-polynomial estimation (rdrobust).
7  *=====
8  version 17
9  clear all
10 set more off
11 set linesize 90
12 set seed 1010
13 capture log close
14 capture mkdir logs // create the log folder if absent (gitignored)
15 log using "logs/10_rdd_regression_discontinuity.log", replace text
16
17 set obs 3000
```

```

18 generate x = runiform(-1,1)           // running variable; cutoff at 0
19 generate byte d = (x >= 0)           // sharp assignment
20 generate y = 3 + 1.5*x + 0.5*x^2 + 2.0*d + rnormal(0,0.5) // TRUE jump = 2.0
21 label variable x "Running variable"
22 label variable y "Outcome"
23
24 *-----
25 * 1. RD point estimate: local-linear fit each side of the cutoff, MSE-optimal
26 *   bandwidth, bias-corrected robust confidence interval.
27 *-----
28 rdrobust y x, c(0)
29
30 *-----
31 * 2. Data-driven bandwidth selection (what rdrobust chose, and alternatives).
32 *-----
33 rdbwselect y x, c(0) all
34
35 *-----
36 * 3. The annotated RD picture. Binned means (rdplot genvars) + a local fit each
37 *   side, then three annotations that make the estimate legible:
38 *   - open dots at the two one-sided limits at the cutoff;
39 *   - a grey double arrow, offset just right of the cutoff, spanning the
40 *     estimated jump and labelled with rdrobust's tau-hat (so the picture
41 *     reports the SAME number as the results table); short dashed leaders
42 *     link the two limit dots to the arrow's ends;
43 *   - the control fit extrapolated across c as the "no-treatment"
44 *     counterfactual (dashed) -- what the treated units would have done.
45 *-----
46 rdplot y x, c(0) genvars hide
47 egen _bin = tag(rdplot_id)
48
49 * local fit each side (quadratic, matching the smooth DGP) + the two limits at c
50 quietly regress y c.x##c.x if x<0
51 predict double _fit_l if x<0           // control fit
52 predict double _cf if inrange(x,0,0.45) // control fit extended past c
53 local yL = _b[_cons]                   // control limit lim_{x up c}
54 quietly regress y c.x##c.x if x>=0
55 predict double _fit_r if x>=0           // treated fit
56 local yR = _b[_cons]                   // treated limit lim_{x down c}
57
58 * headline estimate (rdrobust local-linear) drives the arrow + its label
59 quietly rdrobust y x, c(0)
60 local tau = e(tau_cl)
61 local yT = 'yL' + 'tau'                 // treated limit implied by tau-hat
62 local taus : display %4.3f 'tau'
63 local ymid = ('yL' + 'yT')/2
64
65 twoway ///
66   (scatter rdplot_mean_y rdplot_mean_x if rdplot_mean_x<0 & _bin, mcolor(midblue%35) msze(
67     small)) ///
68   (scatter rdplot_mean_y rdplot_mean_x if rdplot_mean_x>=0 & _bin, mcolor(cranberry%35) msze(
69     small)) ///
70   (line _fit_l x if x<0, sort lcolor(midblue) lwidth(medthick)) ///
71   (line _fit_r x if x>=0, sort lcolor(cranberry) lwidth(medthick)) ///
72   (line _cf x if inrange(x,0,0.45), sort lcolor(midblue%55) lpattern(dash) lwidth(medium))
73   ///
74   (pci 'yT' 0 'yT' 0.24, lcolor(gs9) lpattern(dash) lwidth(thin)) ///
75   (pci 'yL' 0 'yL' 0.24, lcolor(gs9) lpattern(dash) lwidth(thin)) ///
76   (pcarrowi 'yL' 0.24 'yT' 0.24, lcolor(gs9) lwidth(medthick) mcolor(gs9) barbsize(2.2)) ///

```

```

74     (pcarrowi 'yT' 0.24 'yL' 0.24, lcolor(gs9) lwidth(medthick) mcolor(gs9) barbsize(2.2)) ///
75     (scatteri 'yL' 0 'yT' 0, mcolor(white) mlcolor(black) mlwidth(thin) msymbol(0) msize(medium))
    ///
76     , ///
77     xline(0, lpattern(dash) lcolor(gs10)) ///
78     text('ymid' 0.27 "{&tau} = 'taus'", color(gs9) place(e) size(medium)) ///
79     legend(order(3 "Control fit (x < c)" 4 "Treated fit (x {&ge} c)" ///
80             5 "Counterfactual (no treatment)") ///
81            cols(1) position(11) ring(0) size(small) region(lstyle(none) fcolor(none))) ///
82     title("Sharp RD: the outcome jumps at the cutoff") ///
83     subtitle("binned means + local fit each side; grey arrow = rdrobust estimate") ///
84     xtitle("Running variable X (sharp cutoff at c)") ytitle("Outcome Y") ///
85     xlabel(-1 -.5 0 "c = 0" .5 1)
86 graph export "output/10_rdplot.png", replace width(1600)
87 drop _fit_l _fit_r _cf _bin
88
89 *-----
90 * 4. Manipulation test: is the density of x continuous at the cutoff?
91 * (A jump would suggest units sorted around the threshold --- an RD red flag.)
92 * Here x is uniform, so the test should NOT reject.
93 *-----
94 rddensity x, c(0)
95
96 *-----
97 * 5. Export a compact results table for the reference. rdplot's genvars call
98 * overwrote e(), so re-fit quietly to recover the rdrobust/rddensity scalars.
99 *-----
100 quietly rdrobust y x, c(0)
101 local tau = e(tau_cl)
102 local cil = e(ci_l_rb)
103 local cir = e(ci_r_rb)
104 local pv = e(pv_rb)
105 quietly rddensity x, c(0)
106 local mpv = e(pv_q)
107 local ptxt = cond('pv' < 0.001, "\(<)0.001", string('pv', "%5.3f"))
108
109 capture file close rt
110 file open rt using "output/10_rd_table.tex", write replace
111 file write rt "\begin{table}[htbp]\centering" _n
112 file write rt "\caption{Sharp RD estimate vs. \ the true jump (\(\hat{\tau}=2.0\))}" _n
113 file write rt "\begin{tabular}{lc}" _n "\toprule" _n
114 file write rt " & Estimate \\" _n "\midrule" _n
115 file write rt "RD effect (\(\hat{\tau}\)) & " %5.3f ('tau') " \\" _n
116 file write rt "Robust 95% CI & [" %5.3f ('cil') ", " %5.3f ('cir') "]" _n
117 file write rt "Robust \(\rho\)-value & 'ptxt' \\" _n
118 file write rt "\addlinespace" _n
119 file write rt "Manipulation test \(\rho\) (\pkg{rddensity}) & " %5.3f ('mpv') " \\" _n
120 file write rt "\bottomrule" _n "\end{tabular}" _n "\end{table}" _n
121 file close rt
122
123 log close
124 display "10_rdd_regression_discontinuity.do finished OK"

```

11.3 The results

`rdrobust` y x , $c(0)$ fits a separate local-linear regression on each side of the cutoff $c = 0$, inside an MSE-optimal bandwidth h :

$$\text{below } (x < 0) : E[y | x] = \alpha_- + \beta_- x,$$

$$\text{above } (x \geq 0) : E[y | x] = \alpha_+ + \beta_+ x,$$

and the sharp-RD estimate is the gap at the cutoff, $\hat{\tau} = \alpha_+ - \alpha_-$ — the $\lim_{x \downarrow c} - \lim_{x \uparrow c}$ of §The method. The estimate recovers the true jump of 2.0 within a tight robust confidence interval, so the local-linear RD is approximately unbiased here. The **manipulation test** does not reject, so the running variable shows no sign of sorting at the cutoff — the design’s key validity check passes.

Table 11.1: Sharp RD estimate vs. the true jump ($\tau = 2.0$)

	Estimate
RD effect ($\hat{\tau}$)	1.938
Robust 95% CI	[1.770, 2.082]
Robust p -value	<0.001
Manipulation test p (<code>rddensity</code>)	0.953

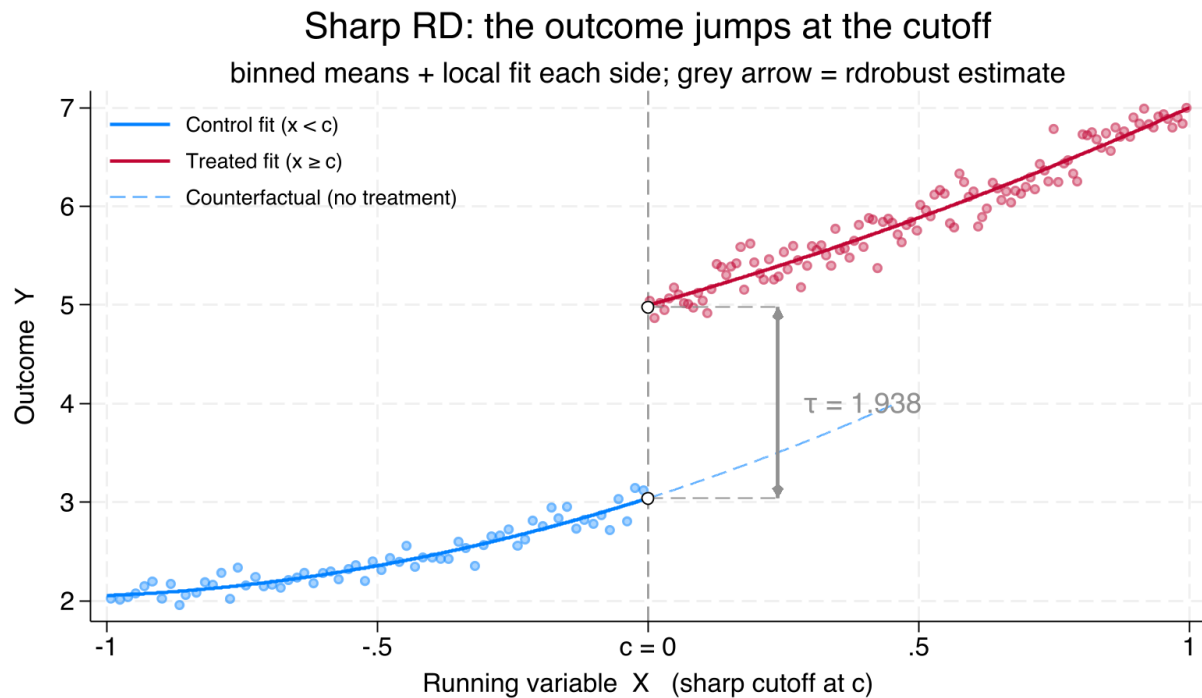


Figure 11.1: Sharp RD: binned means and a local fit each side of the cutoff. The grey double arrow (just right of the cutoff) marks the estimated jump (`rdrobust` $\hat{\tau} = 1.938$, the same figure as Table 11.1), with open dots at the two one-sided limits; the dashed line extrapolates the control fit across c as the no-treatment counterfactual.

Chapter 12

Synthetic control

↔ the method, explained: [Appendix B.8](#)

12.1 The method

The synthetic control method (Abadie, Diamond and Hainmueller, 2010) takes one treated unit and many donors (say, one state that adopts a policy, set against a weighted blend of states that did not) and chooses non-negative weights w_j summing to one that minimise the pre-treatment distance

$$(X_1 - X_0W)'V(X_1 - X_0W).$$

The synthetic control is $\sum_j w_j Y_{jt}$ and the estimated effect is the post-treatment gap $Y_{1t} - \sum_j w_j Y_{jt}$. Inference is by placebo: re-estimate treating each donor as treated and rank the treated unit's gap against the placebo distribution.

In practice. Synthetic control suits single-unit comparative cases: the effect of California's Proposition 99 tobacco tax (Abadie, Diamond and Hainmueller), terrorism and Basque GDP (Abadie and Gardeazabal 2003), and German reunification.

12.2 The code

The workflow. Build a donor panel, fit the `synth` weighted counterfactual for the treated unit, plot treated-vs-synthetic paths, and run placebo inference with `synth_runner`.

```
1  *=====
2  * 11_synthetic_control.do --- synthetic control method (Abadie et al.)
3  *-----
4  * One treated unit, many controls. We build a "synthetic" version of the treated
5  * unit as a weighted average of controls that matches its PRE-treatment path,
6  * then read the treatment effect as the post-treatment gap between the two.
7  *
8  * Simulates 20 units x 30 periods; unit 1 is treated from t = 21 with a true
9  * effect of +3. A factor structure makes a good synthetic control feasible.
10 *-----
11 version 17
12 clear all
13 set more off
14 set linesize 90
15 set seed 2020
16 capture log close
17 capture mkdir logs // create the log folder if absent (gitignored)
18 log using "logs/11_synthetic_control.log", replace text
19
20 *-----
21 * 1. Build the panel (factor model: shared time factors x unit-specific loadings).
```

```

22  *-----
23  set obs 20
24  generate id = _n
25  generate u1 = runiform()           // loading on factor 1 (fixed within unit)
26  generate u2 = runiform()           // loading on factor 2
27  expand 30
28  bysort id: generate t = _n         // 1..30
29  generate f1 = sin(t/3)             // common factors (functions of time)
30  generate f2 = t/30
31  generate y = 5 + 3*u1*f1 + 4*u2*f2 + rnormal(0,0.3)
32  replace y = y + 3 if id==1 & t>20  // TRUE post-treatment effect = +3 on unit 1
33  label variable y "Outcome"
34
35  xtset id t                          // declare the panel
36
37  *-----
38  * 2. Fit the synthetic control for unit 1 (treatment starts t = 21). We save the
39  *   treated/synthetic paths with keep(), then draw a COLORED version of them
40  *   (synth's built-in 'fig' is monochrome).
41  *-----
42  synth y y(5) y(10) y(15) y(20) u1 u2, ///
43      trunit(1) trperiod(21) keep("output/synth_results.dta", replace)
44
45  preserve
46      use "output/synth_results.dta", clear
47      keep if !missing(_time)
48      twoway (line _Y_treated _time, lcolor(navy)          lwidth(medthick)) ///
49             (line _Y_synthetic _time, lcolor(cranberry) lwidth(medthick) lpattern(dash)), ///
50             xline(20.5, lpattern(dash) lcolor(gs10)) ///
51             legend(order(1 "Treated unit" 2 "Synthetic control") rows(1) position(6)) ///
52             title("Synthetic control: treated unit vs. its synthetic") ///
53             xtitle("t") ytitle("Outcome")
54      graph export "output/11_synth_path.png", replace width(1400)
55  restore
56
57  *-----
58  * 3. Inference by placebo: re-assign treatment to each control unit in turn and
59  *   compare the treated unit's gap to the placebo distribution (synth_runner).
60  *-----
61  synth_runner y y(5) y(10) y(15) y(20) u1 u2, ///
62      trunit(1) trperiod(21) gen_vars
63
64  * Colored placebo "gap" plot: each control gap in light gray, the treated gap bold.
65  * synth_runner's gen_vars stores each unit's gap in 'effect'.
66  levelsof id if id!=1, local(donors)
67  local plots ""
68  foreach d of local donors {
69      local plots `plots' (line effect t if id==`d', lcolor(gs13) lwidth(vthin))
70  }
71  twoway `plots' (line effect t if id==1, lcolor(cranberry) lwidth(thick)), ///
72          yline(0, lcolor(gs9)) xline(20.5, lpattern(dash) lcolor(gs10)) legend(off) ///
73          title("Placebo gaps: treated unit (bold) vs. control units") ///
74          xtitle("t") ytitle("Gap (treated - synthetic)")
75  graph export "output/11_synth_placebo.png", replace width(1400)
76
77  log close
78  display "11_synthetic_control.do finished OK"

```

12.3 The results

The true post-treatment effect is +3. `synth y y(5) y(10) y(15) y(20) u1 u2, trunit(1) trperiod(21)` chooses non-negative donor weights w_j (summing to one) so that treated unit 1's pre-period predictors — the listed outcome lags and covariates — best match the weighted donors. The synthetic control is then $\sum_j w_j y_{jt}$, and the estimated effect is the post-period gap

$$\hat{\tau}_t = y_{1t} - \sum_j w_j y_{jt}, \quad t \geq 21.$$

The synthetic control tracks the treated unit closely *before* treatment — a good pre-fit is what licenses the comparison — and diverges by about 3 *after*. In the placebo distribution the treated unit's gap exceeds every control's, so the permutation p -value is near zero: the effect is unlikely to be an artefact of chance.

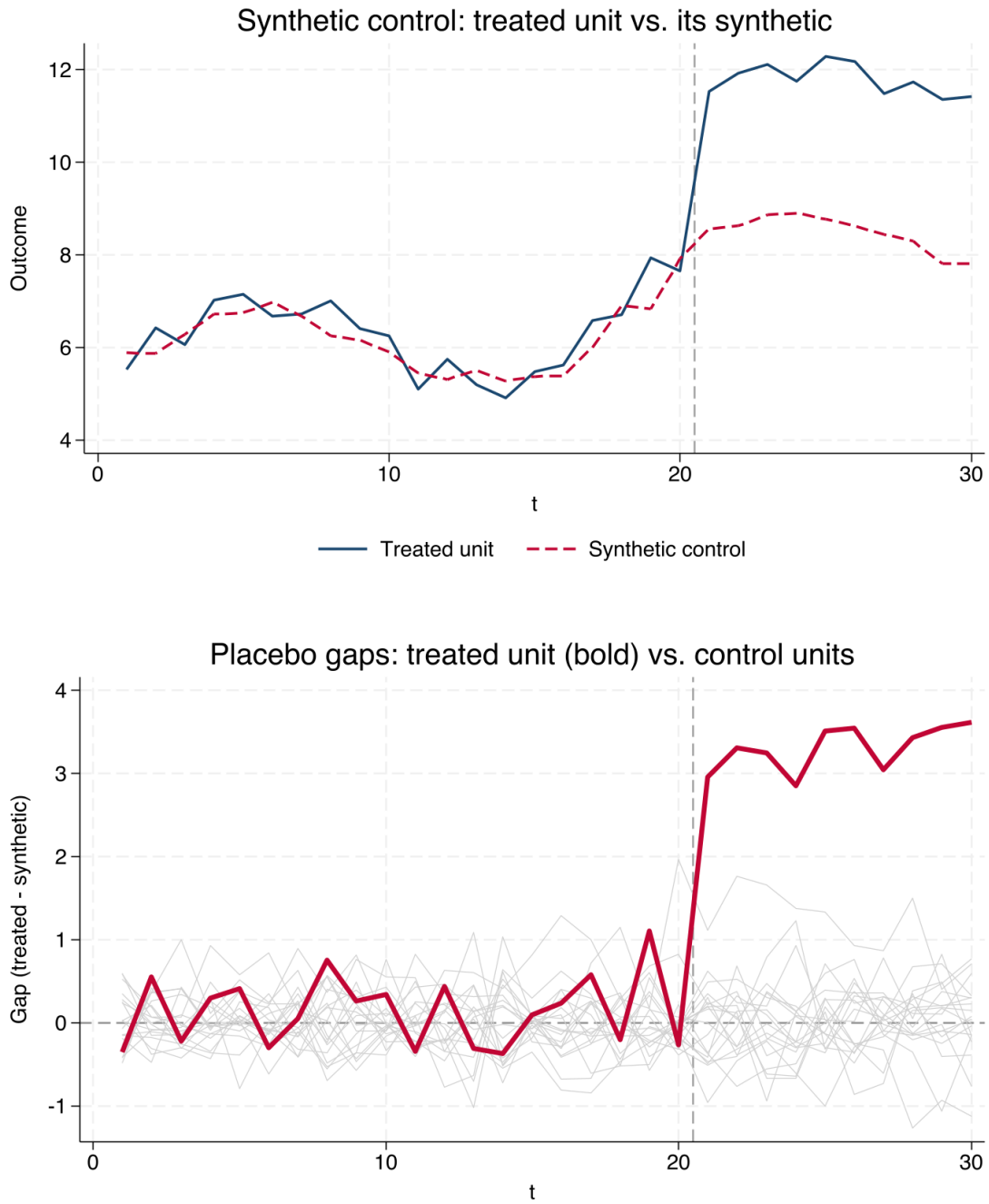


Figure 12.1: Treated vs. synthetic outcome path (top) and the treated gap against placebo gaps (bottom).

Chapter 13

Dynamic panel data (GMM)

↔ the method, explained: [Appendix B.9](#)

13.1 The method

A *dynamic panel* carries a lagged dependent variable,

$$y_{it} = \rho y_{i,t-1} + \beta \mathbf{x}_{it} + \alpha_i + \varepsilon_{it},$$

and that lag is the problem: $y_{i,t-1}$ contains α_i , so it is correlated with the unit effect. Pooled OLS therefore *overstates* ρ , while the within (fixed-effects) transform correlates the demeaned lag with the demeaned error — the [Nickell \(1981\)](#) bias — so FE *understates* it, badly when T is small; the truth is bracketed between the two. [Arellano and Bond \(1991\)](#) first-difference the equation to sweep out α_i and instrument the differenced lag $\Delta y_{i,t-1}$ with *deeper lagged levels* $y_{i,t-2}, y_{i,t-3}, \dots$ (valid, since they are uncorrelated with $\Delta \varepsilon_{it}$) — this is *difference GMM*. [Blundell and Bond \(1998\)](#) add moment conditions in levels (lagged differences instrumenting the levels), giving the more efficient *system GMM*, far better behaved when ρ is near one.

The key validity check is the [Arellano–Bond](#) test for serial correlation in the differenced residuals: [AR\(1\)](#) should reject (differencing mechanically induces it) but [AR\(2\)](#) should *not* — a significant AR(2) means the deeper lags are themselves endogenous and the design fails. Overidentification is the [Sargan/Hansen \$J\$](#) .

In practice. Dynamic panels are everywhere in corporate finance and macro — the persistence of firm leverage, of investment, of employment, or of national output — wherever a lagged dependent variable sits beside unit effects. Roodman’s community `xtabond2` is the most popular implementation; the built-in `xtabond` and `xtdpdsys` used here need no install.

13.2 The code

The workflow. Simulate a dynamic panel with a known ρ , show that pooled OLS and within FE bracket the truth from above and below, then recover it with Arellano–Bond difference GMM (`xtabond`) and Blundell–Bond system GMM (`xtdpdsys`), checking the AR(2) moment condition.

```
1 *=====
2 * 12_dynamic_panel_gmm.do --- dynamic panel data: Arellano-Bond &
3 *                               Blundell-Bond (system) GMM
4 *-----
5 * A dynamic panel  $y_{it} = \rho y_{i,t-1} + \beta x_{it} + \alpha_i + e_{it}$  has a
6 * built-in endogeneity: the lagged dependent variable is mechanically correlated
7 * with the unit effect  $\alpha_i$ . Pooled OLS therefore OVERstates  $\rho$ , and the
8 * within (FE) estimator UNDERstates it (Nickell bias, severe when  $T$  is small).
9 * Difference and system GMM instrument the lagged level with deeper lags and
```

```

10 * recover the truth. TRUE rho = 0.6, beta = 1.
11 *-----
12 version 17
13 clear all
14 set more off
15 set linesize 90
16 set seed 90210
17 capture log close
18 capture mkdir logs // create the log folder if absent (gitignored)
19 log using "logs/12_dynamic_panel_gmm.log", replace text
20
21 *-----
22 * 1. Simulate a dynamic panel: 500 firms; keep 12 periods after a burn-in so the
23 * process has reached its stationary distribution.
24 *-----
25 local N = 500
26 local Tgen = 20
27 local Tkeep = 12
28 local rho = 0.6
29 local beta = 1
30 set obs `N'
31 generate id = _n
32 generate alpha = rnormal() // unit (firm) fixed effect
33 expand `Tgen'
34 bysort id: generate t = _n
35 xtset id t
36 generate x = 0.5*alpha + rnormal() // regressor, correlated with alpha
37 generate eps = rnormal()
38 generate y = .
39 replace y = alpha + `beta'*x + eps if t==1 // starting value
40 forvalues s = 2/`Tgen' { // iterate the dynamics
41     replace y = `rho'*L.y + `beta'*x + alpha + eps if t==`s'
42 }
43 drop if t <= `Tgen' - `Tkeep' // burn-in: discard the early periods
44 bysort id (t): replace t = _n // re-index time to 1..12
45 xtset id t
46
47 *-----
48 * 2. The bias: pooled OLS (rho too high) vs. within/FE (rho too low).
49 *-----
50 regress y L.y x, vce(cluster id) // pooled OLS: rho biased UP
51 local rho_ols = _b[L.y]
52 local se_ols = _se[L.y]
53 xtreg y L.y x, fe vce(cluster id) // within: rho biased DOWN (Nickell)
54 local rho_fe = _b[L.y]
55 local se_fe = _se[L.y]
56
57 *-----
58 * 3. Arellano-Bond difference GMM and Blundell-Bond system GMM.
59 * estat abond: AR(1) in the differenced errors SHOULD reject; AR(2) should NOT
60 * -- the key check that the lagged-level moment conditions are valid.
61 *-----
62 xtabond y x, lags(1) vce(robust) // difference GMM (Arellano-Bond)
63 local rho_ab = _b[L.y]
64 local se_ab = _se[L.y]
65 estat abond
66
67 xtdpdsys y x, lags(1) vce(robust) // system GMM (Blundell-Bond)
68 local rho_bb = _b[L.y]

```

```

69 local se_bb = _se[L.y]
70 estat abond
71
72 *-----
73 * 4. Persistence rho across estimators (true rho = 0.6).
74 *-----
75 display "rho: OLS=" %5.3f 'rho_ols' " FE=" %5.3f 'rho_fe' ///
76        " DiffGMM=" %5.3f 'rho_ab' " SysGMM=" %5.3f 'rho_bb'
77
78 capture file close dp
79 file open dp using "output/12_dynpanel_table.tex", write replace
80 file write dp "\begin{table}[htbp]\centering" _n
81 file write dp "\caption{Dynamic panel: the persistence \(\rho\) (true \(\rho=0.6\))}" _n
82 file write dp "\begin{tabular}{lcccc}" _n "\toprule" _n
83 file write dp "& Pooled OLS & Within FE & Diff GMM & Sys GMM \\" _n "\midrule" _n
84 file write dp "\(\hat{\rho}\) (coef. on \(\hat{y}_{t-1}\)) & " ///
85        %5.3f ('rho_ols') " & " %5.3f ('rho_fe') " & " %5.3f ('rho_ab') " & " %5.3f ('rho_bb') " \\"
86        _n
87 file write dp "(std. err.) & (" %5.3f ('se_ols') ") & (" %5.3f ('se_fe') ") & (" ///
88        %5.3f ('se_ab') ") & (" %5.3f ('se_bb') ") \\" _n
89 file write dp "\bottomrule" _n "\end{tabular}" _n "\end{table}" _n
90 file close dp
91
92 log close
93 display "12_dynamic_panel_gmm.do finished OK"

```

13.3 The results

With a true persistence of $\rho = 0.6$, the four columns differ only in how each removes (or fails to remove) the unit effect α_i :

- | | | |
|-----------------|--|----------------------------|
| (1) Pooled OLS: | $y_{it} = \rho y_{i,t-1} + \beta \mathbf{x}_{it} + (\alpha_i + \varepsilon_{it})$ | (ρ biased up), |
| (2) Within FE: | $\tilde{y}_{it} = \rho \tilde{y}_{i,t-1} + \beta \tilde{\mathbf{x}}_{it} + \tilde{\varepsilon}_{it}$ | (Nickell: down), |
| (3) Diff GMM: | $\Delta y_{it} = \rho \Delta y_{i,t-1} + \beta \Delta \mathbf{x}_{it} + \Delta \varepsilon_{it}$ | (IV: $y_{i,t-2}, \dots$), |
| (4) Sys GMM: | (3) plus the untransformed levels equation | (IV: $\Delta y_{i,t-1}$), |

with $\tilde{z}_{it} = z_{it} - \bar{z}_i$ the within (demeaned) transform and $\Delta z_{it} = z_{it} - z_{i,t-1}$ the first difference. Pooled OLS leaves α_i in the error, where it correlates with the lag and biases ρ up; within FE demeans it away but the demeaned lag correlates with the demeaned error (**Nickell bias**), biasing ρ down; both GMM estimators first-difference α_i out and instrument the lag with the panel's own history, so only they recover the truth.

Table 13.1: Dynamic panel: the persistence ρ (true $\rho = 0.6$)

	Pooled OLS	Within FE	Diff GMM	Sys GMM
$\hat{\rho}$ (coef. on y_{t-1})	0.803	0.524	0.576	0.577
(std. err.)	(0.004)	(0.009)	(0.017)	(0.016)

The **Arellano–Bond** test confirms the design — **AR(1)** rejects while **AR(2)** does not, so the lagged-level instruments are valid:

```
. estat abond
Arellano-Bond test for zero autocorrelation in first-differenced errors
Order      z          Prob > z
   1    -17.12     0.0000    <- AR(1): expected (differencing induces it)
   2     -0.94     0.3490    <- AR(2): does NOT reject -> instruments valid
```

A significant AR(2) would be the red flag; here $p \approx 0.35$, so the design passes.

Part IV
Appendices

Appendix A

Python/R → Stata cheat sheet

Table A.1: Common tasks across pandas / R / Stata

Task	pandas / R	Stata
Load CSV	<code>pd.read_csv / read_csv</code>	<code>import delimited "f.csv", varnames(1) clear</code>
Inspect	<code>df.head() / head()</code>	<code>list in 1/5, describe, codebook</code>
New column	<code>df["x"]=... / mutate</code>	<code>generate x = ... (replace to overwrite)</code>
Group aggregate	<code>groupby().agg</code>	<code>collapse (mean)..., by(g) / egen ..., by(g)</code>
Reshape	<code>pivot/melt</code>	<code>reshape wide / reshape long</code>
Join	<code>merge</code>	<code>merge 1:1 key using "f.dta"</code>
OLS robust	<code>cov_type="HC1" / feols</code>	<code>regress y x, vce(robust)</code>
Fixed effects	<code>PanelOLS / feols(. fe)</code>	<code>xtreg ..., fe / reghdfe ..., absorb()</code>
Marginal effects	<code>get_margeff() / margins</code>	<code>margins, dydx(x)</code>
Lag	<code>.shift() / lag()</code>	<code>L.x (after <code>tsset</code>)</code>

Appendix B

The methods, explained

This appendix is the explanatory companion to the body. The chapters state each method crisply assuming familiarity with the econometrics; here we build the nine estimators up from scratch. Each entry runs through the same beats: the *idea*; the *symbols*; what the defining *formula* means; *why it works* (why the estimator recovers the effect); *the catch* (the identifying assumption and what breaks it); a *worked example* to follow with pen and paper; and *in the chapter* — how that chapter’s own run recovers the planted answer. (Stata basics, data management, and data import are tooling rather than estimators, so they are not repeated here.)

B.1 Ordinary least squares

↪ *worked Stata chapter: Chapter 4*

The idea. OLS draws the straight line (or hyperplane) through a cloud of points that makes the total *squared* vertical distance to the points as small as possible. Each slope β_j answers a *ceteris paribus* question: “if x_j rises by one unit and the other regressors are held fixed, how much does the average of y move?”

Symbols. y the outcome (a car’s price); x_1, \dots, x_k the regressors (mileage, weight); β_0 the intercept and β_1, \dots, β_k the slopes; ε the error — everything else that moves y ; n the number of observations; X the $n \times (k+1)$ matrix stacking a column of ones and the regressors.

The formula. Choose β to minimise the sum of squared residuals $\sum_i (y_i - x_i' \beta)^2$. Differentiating and setting the gradient to zero gives the *normal equations* $X'(y - X\beta) = 0$, whose solution is

$$\hat{\beta} = (X'X)^{-1}X'y.$$

Why it works. Geometrically $\hat{\beta}$ *projects* y onto the space spanned by the regressors: it makes the residual orthogonal to every x_j , so no linear function of the regressors can predict what is left over — that orthogonality is precisely the content of the normal equations. When the error variance is non-constant the point estimate is unchanged but its standard error must use the heteroskedasticity-robust “sandwich” $(X'X)^{-1}(\sum_i \hat{\varepsilon}_i^2 x_i x_i')(X'X)^{-1}$.

The catch. $\hat{\beta}$ recovers the *causal* effect only if the regressors are *exogenous*, $E[\varepsilon | X] = 0$ — no omitted variable correlated with x , no reverse causality. The orthogonality OLS imposes *in sample* (residual $\perp x$) is mechanical; the assumption is that it also holds in the *population*. When it fails — unobserved ability in a wage–schooling regression — the slope is biased, which is exactly what instrumental variables repairs.

Worked example. Three points $(x, y) = (1, 2), (2, 2), (3, 5)$, so $\bar{x} = 2$ and $\bar{y} = 3$:

$$\text{slope} = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sum(x - \bar{x})^2} = \frac{(-1)(-1) + 0 + (1)(2)}{1 + 0 + 1} = \frac{3}{2} = 1.5,$$

$$\text{intercept} = \bar{y} - 1.5\bar{x} = 3 - 3 = 0.$$

The fitted line is therefore $\hat{y} = 1.5x$.

In the chapter. Chapter 4 fits car price on mileage, weight, and origin in Stata’s 74-car auto data with robust standard errors. Adding the foreign indicator and an interaction lifts the fit from $R^2 = 0.29$ to 0.52, and weight enters strongly ($\approx 1,750$ per 1,000 lbs). Because these are real data there is no planted coefficient to recover — the chapter reads the table mechanically, not as economics.

B.2 Panel data, difference-in-differences, and event studies

↔ worked Stata chapter: [Chapter 5](#)

The idea. The goal is the effect of a treatment that some units receive at a known date. Comparing treated to untreated confounds pre-existing differences between them; comparing each group’s before to after confounds whatever else changed over time. Difference-in-differences removes both at once: take each group’s *change* across the date, then difference those two changes.

Symbols. y_{it} outcome for unit i in period t ; α_i a unit fixed effect (time-invariant traits, differenced away); γ_t a time fixed effect (shocks common to all units); treat_i marks the eventually-treated group, post_t the post period; β the DiD effect.

The formula. The 2×2 DiD is the difference of two differences,

$$\beta = (\bar{y}_{\text{tr,post}} - \bar{y}_{\text{tr,pre}}) - (\bar{y}_{\text{co,post}} - \bar{y}_{\text{co,pre}}),$$

which is identical to the coefficient on the interaction $\text{treat}_i \times \text{post}_t$ in the two-way fixed-effects regression $y_{it} = \alpha_i + \gamma_t + \beta(\text{treat}_i \times \text{post}_t) + \varepsilon_{it}$.

Why it works. The control group’s change estimates the common trend — what *would* have happened to the treated group anyway. Subtracting it strips that trend out of the treated group’s change, leaving only the treatment. The unit fixed effects α_i absorb every fixed difference between units, so the comparison is made within unit, over time.

The catch. *Parallel trends*: absent treatment, the treated group’s average would have moved in step with the control group’s. It cannot be tested directly (the counterfactual is unobserved), but the *event study* — the path of effects by period relative to treatment — gives indirect evidence: if the pre-treatment “leads” sit flat near zero, the groups were moving together beforehand, which is reassuring.

Reading the event study. The dynamic version swaps the single β for one coefficient β_k per *event time* $k = t - t_i^*$ (periods relative to a unit's own treatment date t_i^*), fitted for every k in the data except $k = -1$:

$$k : \underbrace{-5 \quad -4 \quad -3 \quad -2}_{\text{leads: } \beta_k \approx 0} \quad \bigg| \quad \underbrace{-1}_{\text{baseline}} \quad \bigg| \quad \underbrace{0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6}_{\text{lags: the effects}}$$

One indicator must be dropped — the full lead/lag set is collinear with the fixed effects — and $k = -1$, the eve of treatment, is the conventional baseline ($\beta_{-1} \equiv 0$), so every β_k is read relative to it. The *lags* ($k \geq 0$) are the effects of interest: the treated-minus-control gap 0, 1, 2, ... periods after onset. The *leads* ($k \leq -2$) are kept on purpose — they measure that same gap *before* treatment, where it ought to be zero, so they act as a placebo test. Flat, insignificant leads are the visual parallel-trends check; a large lead would mean the groups were already diverging, leaving the post-treatment estimates untrustworthy. Summing over $k \neq -1$ rather than only $k \geq 0$ is precisely what builds that test in — the advantage an event study has over a single DiD number.

Worked example. Treated mean 10 → 16 (change +6); control mean 10 → 12 (change +2). DiD = 6 − 2 = 4: of the treated group's six-point rise, two points are the common trend and four are the effect.

In the chapter. Chapter 5 simulates 150 firms over 2010–2021 with treatment switched on in 2015 whose effect grows 0.5 per year — a true *average* post-treatment effect of 2.0. Pooled OLS, firm FE, and two-way FE all return 1.973 (SE 0.104), and the event-study leads sit flat at zero while the lags climb: the estimator recovers the planted truth and shows no pre-trend.

B.3 Time series

↔ worked Stata chapter: [Chapter 6](#)

The idea. In a series where today depends on yesterday, the observations are not independent draws — they carry information forward. Two consequences follow: we can *model* the persistence (an autoregression), and when we don't, we must *correct* ordinary standard errors, which otherwise treat correlated observations as independent and so overstate precision.

Symbols. y_t the value at time t ; ϕ the persistence (AR coefficient); c an intercept; ε_t the period- t shock; L the lag operator, $Ly_t = y_{t-1}$.

The formula. A first-order autoregression is $y_t = c + \phi y_{t-1} + \varepsilon_t$. Iterating it, a shock feeds into the next period at rate ϕ , the one after at ϕ^2 , and so on, so its effect k periods out is ϕ^k . When residuals stay serially correlated, the Newey–West (HAC) variance reweights the autocovariances up to a chosen lag L :

$$\hat{\Omega} = \hat{\Gamma}_0 + \sum_{\ell=1}^L \left(1 - \frac{\ell}{L+1}\right) (\hat{\Gamma}_\ell + \hat{\Gamma}'_\ell).$$

Why it works. The geometric decay ϕ^k means distant shocks barely register, so the cumulative effect of a one-unit shock is the finite sum $1 + \phi + \phi^2 + \dots = 1/(1 - \phi)$. Newey–West keeps the usual point estimate but adds back the cross-period covariances that the independent-errors formula ignores, with the triangular weights $(1 - \ell/(L+1))$ guaranteeing a positive variance.

The catch. *Stationarity* — for the AR(1), $|\phi| < 1$, so the series reverts rather than wandering off and its mean and variance are stable. A nearby trap is a *level break*: an unmodelled shift makes the series look more persistent than it is and biases $\hat{\phi}$ upward — control for the break and the bias disappears.

Worked example. With $\phi = 0.6$, a one-unit shock raises y by 1 now, 0.6 next period, 0.36 after that, ...; the cumulative effect is $1/(1 - 0.6) = 2.5$.

In the chapter. Chapter 6 simulates a monthly AR(1) with $\phi = 0.6$ plus a +3 level shift in 2015. Regressing y on its lag *without* the break gives $\hat{\phi} = 0.85$ — inflated, the level-break trap exactly; adding the break indicator recovers $\hat{\phi} = 0.57$ (≈ 0.6), and Newey–West (lag 3) widens the standard errors for the residual serial correlation.

B.4 Instrumental variables (2SLS)

↪ worked Stata chapter: [Chapter 8](#)

The idea. If x is correlated with the error — schooling and unobserved ability both raise wages — OLS cannot separate the effect of x from the confounder, and its slope is biased. An *instrument* z is a variable that moves x but reaches y *only through* x . Keeping only the variation in x that z explains discards the contaminated part and leaves a clean comparison.

Symbols. y outcome; x the endogenous regressor; z the instrument; β the causal slope; π the first-stage slope; ε the structural error.

The formula. *First stage*: regress x on z , $x = \pi z + v$, and form the fitted \hat{x} (the part of x driven by z). *Second stage*: regress y on \hat{x} . With a single instrument the two stages collapse to the Wald ratio

$$\hat{\beta}_{\text{IV}} = \frac{\widehat{\text{Cov}}(z, y)}{\widehat{\text{Cov}}(z, x)}.$$

Why it works. Because z touches y only through x , the numerator $\text{Cov}(z, y)$ equals β times $\text{Cov}(z, x)$; dividing by the denominator returns β . The confounder, being unrelated to z , drops out — that is the whole trick. The price is precision: only the z -driven slice of x is used, so IV standard errors exceed OLS's.

The catch. Two requirements. *Relevance*, $\text{Cov}(z, x) \neq 0$ — the instrument must actually move x , or the denominator is near zero and the estimate explodes (the “weak instrument” problem, watched with the first-stage F , rule of thumb > 10). *Exclusion*, $\text{Cov}(z, \varepsilon) = 0$ — no back-door from z to y . Exclusion is not testable; it is argued from how z arises.

Worked example. A lottery z raises program take-up x by 0.4 and raises the outcome y by 0.8. Then $\hat{\beta}_{\text{IV}} = 0.8/0.4 = 2$: each unit of participation raises y by 2.

In the chapter. Chapter 8 simulates an endogenous x with a true slope of 2.0 and a confounder that pushes OLS upward. OLS returns 2.173 (biased high); 2SLS with two instruments returns 1.960 and the high-dimensional-FE version 1.994 — both land on 2.0 once the confounded variation is purged, and the first-stage F sits far above the weak-instrument bar.

B.5 Binary outcome models

↔ worked Stata chapter: [Chapter 9](#)

The idea. When y is 0/1, fitting a straight line (the linear probability model) can predict probabilities below 0 or above 1 and forces a constant marginal effect. Logit and probit instead push a linear index $x'\beta$ through an S-shaped curve pinned to $[0, 1]$ that flattens at both ends.

Symbols. $y \in \{0, 1\}$; $x'\beta$ the linear index; Λ the logistic CDF (logit); Φ the standard-normal CDF (probit); g the matching density (Λ' or ϕ).

The formula. $P(y = 1 | x) = G(x'\beta)$ with $G = \Lambda$ or Φ . Because G is nonlinear, a coefficient β_j is *not* the effect on the probability; the comparable summary is the average marginal effect, which averages the curve's slope over the sample,

$$\text{AME}_j = \frac{1}{n} \sum_i g(x'_i \beta) \beta_j.$$

For logit, e^{β_j} is the odds ratio.

Why it works. The slope of the probability curve, $g(x'\beta)\beta_j$, is steepest in the middle and shallow at the extremes — a one-unit change in x moves the probability a lot for a borderline case and little for an almost-certain one. Averaging that slope across the sample yields a single number on the same footing as an OLS coefficient, which is why the AME, not β_j , is the quantity to report.

The catch. The index is assumed linear and the link (logistic or normal) correctly specified. In practice the choice barely matters: logit and probit deliver nearly identical marginal effects, differing mainly in the scale of the raw coefficients, not the AMEs.

Worked example. A logit index of 0 gives $P = \Lambda(0) = 0.5$; an index of 1 gives $P = \Lambda(1) = e/(1 + e) \approx 0.73$. A coefficient of 0.7 multiplies the odds by $e^{0.7} \approx 2$.

In the chapter. Chapter 9 models whether a car is foreign from its mileage, weight, and price. The logit coefficient on mileage is -0.121 — not a probability effect; the average marginal effect is -0.009 , essentially identical to the probit AME (-0.009) and close to the LPM slope (-0.015): the textbook “logit \approx probit” result. The odds ratio on mileage is 0.886.

B.6 Staggered-adoption difference-in-differences

↔ worked Stata chapter: [Chapter 10](#)

The idea. When units adopt treatment in *different* years, an ordinary two-way fixed-effects regression quietly uses already-treated units as if they were clean controls for later-treated ones. Because those “controls” are still responding to their own treatment, the comparison is contaminated — a “forbidden comparison” whose negative weights can shrink or even flip the estimate. Callaway–Sant’Anna rebuild the estimator from clean 2×2 blocks that only ever use not-yet-treated or never-treated units as controls, then average the blocks.

Symbols. $G = g$ the cohort (the first period a unit is treated); t calendar time; C a clean comparison group (never-treated or not-yet-treated); $ATT(g, t)$ the effect for cohort g at time t ; $e = t - g$ event time.

The formula. The building block compares a cohort’s change since just before its own treatment to a clean control group’s change over the same window,

$$ATT(g, t) = \underbrace{E[Y_t - Y_{g-1} \mid G = g]}_{\text{cohort } g\text{'s change}} - \underbrace{E[Y_t - Y_{g-1} \mid C]}_{\text{clean controls' change}},$$

and these $ATT(g, t)$ are then averaged — overall, by event time e , or by cohort.

Why it works. Each block is an honest DiD: a single cohort against controls that have *not* yet been treated, so no already-treated unit ever enters the control group. Keeping the cohorts separate until the final averaging step is what removes the negative weighting the pooled regression introduces.

The catch. Parallel trends must hold relative to the clean comparison group, and there must be *no anticipation* — units do not change behaviour before g in expectation of treatment. When effects are dynamic and timing is staggered, it is the naive pooled estimate to distrust, not this one.

Worked example. A cohort treated in 2015 rises by 5 from 2014 to 2016; the not-yet-treated units rise by 2 over the same span. Then $ATT(2015, 2016) = 5 - 2 = 3$.

In the chapter. Chapter 10 simulates two cohorts (treated in 2014 and 2017) plus a never-treated group, with an effect that grows 0.4 per year after each cohort’s own onset (a true average near 1.5). Callaway–Sant’Anna recovers 1.50; the naive two-way FE on the *same* data returns only 1.01 — the staggered-timing bias, made visible.

B.7 Regression discontinuity

↔ worked Stata chapter: [Chapter 11](#)

The idea. Imagine a treatment that turns on the exact moment a continuous score hits a specific line. The individuals landing a fraction below and a fraction above that line are essentially twins; they share the same background, talent, and luck. The only difference is that one group got the treatment and the other missed it. By comparing these two identical groups right at the border, any sudden jump in their outcomes reveals the exact impact of the treatment. The threshold naturally creates a natural, cost-free randomized experiment for anyone nearby.

Symbols. Y is the outcome (i.e. subsequent earnings, enrollment); X is the running variable (the continuous score determining treatment, e.g., a test score); c is the threshold or cutoff; τ (tau) is the treatment effect, measured by the vertical jump at the cutoff.

The formula. $E[Y | X = x]$ traces the average outcome as a function of the score; τ is its jump at c ,

$$\tau = \underbrace{\lim_{x \downarrow c} E[Y | X = x]}_{\text{just above (treated)}} - \underbrace{\lim_{x \uparrow c} E[Y | X = x]}_{\text{just below (control)}}.$$

The arrows are one-sided limits: $x \downarrow c$ approaches c from *above* (the treated side), $x \uparrow c$ from *below* (the control side).

Why it works. A student scoring 60 (just qualifies) and one scoring 59.9 (just misses) are alike in every way except the treatment, so nothing but the treatment can explain a jump in their outcomes *exactly* at c . A treated and an untreated unit cannot be observed at the *same* score, so the fit approaches c from each side — a line on each side — and reads off the gap between them. Away from c the curve is smooth; only at c does treatment create a step.

The catch. *Continuity:* absent treatment, average outcomes would pass smoothly through c , so any jump is the treatment's doing. It fails if units can *manipulate* the score to land on the favourable side (a grader nudging a 59 up to a 60); that sorting puts a jump at c that is not the treatment. A density/manipulation test checks whether units pile up just past the cutoff.

Worked example. Just below the cutoff the fitted line reaches 3.0; just above it reaches 5.0. Then $\tau = 5.0 - 3.0 = 2.0$.

In the chapter. Chapter 11 plants a jump of 2.0 at the cutoff. `rdrobust` (local-linear, bias-corrected) estimates $\hat{\tau} = 1.938$ with a robust 95% CI of [1.770, 2.082], and the `rddensity` manipulation test returns $p = 0.95$ — no bunching, because the p -value is high, we fail to reject the null hypothesis of a smooth density. This proves there is no manipulation of scores to "bunch" just above the threshold; confirming the design is clean, valid quasi-experiment, and the estimate lands on the truth.

B.8 Synthetic control

↔ *worked Stata chapter:* [Chapter 12](#)

The idea. For a single treated unit — one state, one country — there is no obvious comparison group, so one is *built*: a weighted average of untreated “donor” units, with the weights chosen so the blend tracks the treated unit’s path *before* treatment. Afterwards, the gap between the real unit and this synthetic twin is the estimated effect.

Symbols. Unit 1 is treated, units $j \geq 2$ are donors; $w_j \geq 0$ are weights summing to one; Y_{1t} is the treated outcome, Y_{jt} the donor outcomes; X_1, X_0 are the pre-treatment characteristics of the treated unit and of the donors.

The formula. Choose the weights to match the pre-treatment period as closely as possible,

$$\min_W (X_1 - X_0W)'V(X_1 - X_0W), \quad w_j \geq 0, \quad \sum_j w_j = 1,$$

then form the synthetic control $\sum_j w_j Y_{jt}$ and read the effect as the post-treatment gap $Y_{1t} - \sum_j w_j Y_{jt}$.

Why it works. If the weighted donors reproduce the treated unit's outcome for many pre-treatment periods, they plausibly share its underlying drivers (its exposure to the common shocks), so they stand in for what the treated unit *would* have done. The non-negativity and sum-to-one constraints keep the synthetic unit inside the donors' range — no extrapolation beyond the data.

The catch. A credible *pre-treatment fit* (a poor match means a poor counterfactual), *no interference* (the treatment must not spill onto the donors), and a comparable donor pool. With only one treated unit, inference is by *placebo*: pretend each donor was treated, compute its gap, and ask whether the real unit's gap is unusually large against that distribution.

Worked example. Two donors with weights 0.6 and 0.4 and post-treatment outcomes 8 and 13 give a synthetic value $0.6(8) + 0.4(13) = 10$. If the treated unit is 13, the effect is $13 - 10 = 3$.

In the chapter. Chapter 12 plants a +3 effect on the treated unit from period 21. The fitted synthetic control tracks it closely beforehand, and the post-treatment gaps average about 3.0 (e.g. 2.96, 3.31, 3.25, . . .); every placebo *p*-value is essentially zero, so the real unit's gap stands out sharply from the donor placebos.

B.9 Dynamic panel data (GMM)

↪ worked Stata chapter: [Chapter 13](#)

The idea. When today's outcome depends on yesterday's — leverage, investment, output all have momentum — the model needs a lagged dependent variable. But alongside unit fixed effects that lag is endogenous (it contains the fixed effect), and the standard estimators are biased in *opposite* directions. The fix is to instrument the lag with the panel's own deeper history.

Symbols. y_{it} outcome; ρ persistence (the coefficient on $y_{i,t-1}$); α_i the unit effect; Δ the first difference.

The formula. First-difference to kill α_i ,

$$\Delta y_{it} = \rho \Delta y_{i,t-1} + \beta \Delta x_{it} + \Delta \varepsilon_{it},$$

then instrument $\Delta y_{i,t-1}$ with lagged *levels* $y_{i,t-2}, y_{i,t-3}, \dots$ (difference GMM). System GMM adds the untransformed levels equation, using lagged *differences* as its instruments.

Why it works. A level two or more periods back is correlated with $\Delta y_{i,t-1}$ but, after differencing, uncorrelated with $\Delta \varepsilon_{it}$ — exactly the relevance and exogeneity an instrument needs. GMM stacks the many such moment conditions and weights them efficiently.

The catch. The lagged-level instruments are valid only if ε_{it} is itself serially uncorrelated — precisely what the **AR(2)** test checks. The instrument count also grows with T^2 , so deep lags are often “collapsed”; and system GMM is preferred when ρ is near one, where deep lags become weak instruments in differences.

Worked example. Pooled OLS reads $\hat{\rho} = 0.80$ and within FE reads 0.52; the truth (0.6) is bracketed between them. Instrumenting the lag with its own history, GMM returns 0.58 — inside the bracket and close to the truth.

In the chapter. Chapter 13 simulates 500 firms over 12 periods with $\rho = 0.6$. Pooled OLS returns 0.803, within FE 0.524, and both difference and system GMM about 0.58; the AR(2) test does not reject ($p \approx 0.35$), validating the lagged-level instruments.

References

Methods.

- **White, Halbert**, 1980, A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity, *Econometrica* 48, 817–838. [10.2307/1912934](https://doi.org/10.2307/1912934)
- **Newey, Whitney K., and Kenneth D. West**, 1987, A simple, positive semi-definite, heteroskedasticity and autocorrelation consistent covariance matrix, *Econometrica* 55, 703–708. [10.2307/1913610](https://doi.org/10.2307/1913610)
- **Staiger, Douglas, and James H. Stock**, 1997, Instrumental variables regression with weak instruments, *Econometrica* 65, 557–586. [10.2307/2171753](https://doi.org/10.2307/2171753)
- **Goodman-Bacon, Andrew**, 2021, Difference-in-differences with variation in treatment timing, *Journal of Econometrics* 225, 254–277. [10.1016/j.jeconom.2021.03.014](https://doi.org/10.1016/j.jeconom.2021.03.014)
- **Callaway, Brantly, and Pedro H. C. Sant’Anna**, 2021, Difference-in-differences with multiple time periods, *Journal of Econometrics* 225, 200–230. [10.1016/j.jeconom.2020.12.001](https://doi.org/10.1016/j.jeconom.2020.12.001)
- **Calonico, Sebastian, Matias D. Cattaneo, and Rocío Titiunik**, 2014, Robust nonparametric confidence intervals for regression-discontinuity designs, *Econometrica* 82, 2295–2326. [10.3982/ECTA11757](https://doi.org/10.3982/ECTA11757)
- **Cattaneo, Matias D., Michael Jansson, and Xinwei Ma**, 2020, Simple local polynomial density estimators, *Journal of the American Statistical Association* 115, 1449–1455. [10.1080/01621459.2019.1635480](https://doi.org/10.1080/01621459.2019.1635480)
- **Abadie, Alberto, Alexis Diamond, and Jens Hainmueller**, 2010, Synthetic control methods for comparative case studies, *Journal of the American Statistical Association* 105, 493–505. [10.1198/jasa.2009.ap08746](https://doi.org/10.1198/jasa.2009.ap08746)
- **Nickell, Stephen**, 1981, Biases in dynamic models with fixed effects, *Econometrica* 49, 1417–1426. [10.2307/1911408](https://doi.org/10.2307/1911408)
- **Arellano, Manuel, and Stephen Bond**, 1991, Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations, *Review of Economic Studies* 58, 277–297. [10.2307/2297968](https://doi.org/10.2307/2297968)
- **Blundell, Richard, and Stephen Bond**, 1998, Initial conditions and moment restrictions in dynamic panel data models, *Journal of Econometrics* 87, 115–143. [10.1016/S0304-4076\(98\)00009-8](https://doi.org/10.1016/S0304-4076(98)00009-8)

Canonical applications (cited in the “In practice” notes).

- **Abadie, Alberto, and Javier Gardeazabal**, 2003, The economic costs of conflict: a case study of the Basque Country, *American Economic Review* 93, 113–132. [10.1257/000282803321455188](https://doi.org/10.1257/000282803321455188)

- **Angrist, Joshua D.**, 1990, Lifetime earnings and the Vietnam era draft lottery: evidence from social security administrative records, *American Economic Review* 80, 313–336.
- **Angrist, Joshua D., and Alan B. Krueger**, 1991, Does compulsory school attendance affect schooling and earnings?, *Quarterly Journal of Economics* 106, 979–1014. [10.2307/2937954](https://doi.org/10.2307/2937954)
- **Angrist, Joshua D., and Victor Lavy**, 1999, Using Maimonides' rule to estimate the effect of class size on scholastic achievement, *Quarterly Journal of Economics* 114, 533–575. [10.1162/003355399556061](https://doi.org/10.1162/003355399556061)
- **Card, David, 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, 772–793.
- **Lee, David S.**, 2008, Randomized experiments from non-random selection in U.S. House elections, *Journal of Econometrics* 142, 675–697. [10.1016/j.jeconom.2007.05.004](https://doi.org/10.1016/j.jeconom.2007.05.004)
- **Mincer, Jacob**, 1974, *Schooling, Experience, and Earnings* (Columbia University Press for the NBER, New York).
- **Thistlethwaite, Donald L., and Donald T. Campbell**, 1960, Regression-discontinuity analysis: an alternative to the ex post facto experiment, *Journal of Educational Psychology* 51, 309–317. [10.1037/h0044319](https://doi.org/10.1037/h0044319)