

Differences in Differences

A. Colin Cameron
Univ. of California, Davis

These slides are part of the set of slides
A. Colin Cameron, Introduction to Causal Methods
<https://cameron.econ.ucdavis.edu/causal/>

March 2023

Introduction

- These slides give an introductory example of differences-in-differences (DID) estimation
 - ▶ DID is a method for causal inference
 - ▶ it is a general method for when an exogenous policy comes into being that affects one group more than another
 - ▶ it is often used with repeated cross-section data over time
 - ▶ but can also be used by comparing subgroups.
- DID relies crucially on an assumption called parallel trends
 - ▶ in the absence of treatment the trends for treated and untreated groups are equal.

- Separately the Stata file `dind.do` implements these methods
 - ▶ using dataset `AED_HEALTHACCESS.DTA`
- The data are from chapter 13.6 of A. Colin Cameron (2022) *Analysis of Economics Data: An Introduction to Econometrics* <https://cameron.econ.ucdavis.edu/>.
- Data are originally from Shinsuke Tanaka (2014), “Does Abolishing User Fees Lead to Improved Health Status? Evidence from Post-Apartheid South Africa”, *American Economics Journal: Economic Policy*, 6(3), pages 282-312.

Outline

- 1 Introduction
- 2 Differences in differences
- 3 Example: Access to health care and health outcomes
- 4 Results
- 5 Further Details
- 6 References

Differences in differences: Two Time Periods

- Consider a “natural” experiment where an exogenous policy (called a treatment) comes into being that effects one group more than another.
- Let y denote the outcome and d denote the treatment
 - ▶ with $d = 1$ if treated and $d = 0$ if not treated.
- 1. Method 1: Treatment-control comparison (at a point in time)
 - ▶ Treatment effect = (\bar{y} for treated) - (\bar{y} for not treated) = $\bar{y}_{d=1} - \bar{y}_{d=0}$.
 - ▶ Problem: This is misleading if the treated and untreated groups differ in their characteristics
 - ★ e.g. if the policy was targeted towards poor people.
- 2. Method 2: Before-after comparison over time for treated only
 - ▶ Treatment effect = (\bar{y} for treated after treatment) - (\bar{y} for treated before treatment)
 - ▶ Problem: Misleading if other things also effect the treated over time.
- 3. Differences-in-differences combines methods 1. and 2.
 - ▶ it uses change over time for the untreated to control for nontreatment changes over time (assuming both groups have the same time trend).

Differences in differences formula

- Introduce time before (pre) and after (post) the policy comes into effect
 - ▶ $t = 0$ is a time period before and $t = 1$ is a time period after.
- Then the difference in difference estimate of the effect of treatment is
 - ▶ $DinD = \Delta \bar{y}$ for those treated $- \Delta \bar{y}$ for those not treated
 - ▶ $= (\bar{y}_{d=1,post} - \bar{y}_{d=1,pre}) - (\bar{y}_{d=0,post} - \bar{y}_{d=0,pre})$.
- Equivalently we can use
 - ▶ $DinD = (\bar{y}_{d=1,post} - \bar{y}_{d=0,post}) - (\bar{y}_{d=1,pre} - \bar{y}_{d=0,pre})$
 - ▶ the post-period difference in the two groups less that in the pre-period.
- $DinD$ can be estimated by computing the four separate means and then computing the differences.

Regression computation

- The same difference-in-difference estimate can be obtained as the coefficient of $t \times d$ in the OLS regression

$$y_i = \beta_1 + \beta_2 t_i + \beta_3 d_i + \beta_4 t_i \times d_i + u_i.$$

- where $t_i = 1$ in the post-period and $t_i = 0$ in the pre-period
 - and $d_i = 1$ if treated and $d_i = 0$ if not treated
 - $t_i \times d_i = 1$ if treated and in the post-period and $= 0$ otherwise.
- Proof: The model implies that y equals the following

	Treated ($d = 1$)	Not Treated ($d = 0$)	Difference over treatment
Pre ($t = 0$)	$\beta_1 + \beta_3$	β_1	β_3
Post ($t = 1$)	$\beta_1 + \beta_2 + \beta_3 + \beta_4$	$\beta_1 + \beta_2$	$\beta_3 + \beta_4$
Change over time	$\beta_2 + \beta_4$	β_2	Diff in diff = β_4!

Differences in differences regression computation

- So suppose we have data on each individual, not just the means.
- The OLS regression is

$$y_i = \beta_1 + \beta_2 t_i + \beta_3 d_i + \beta_4 t_i \times d_i + u_i.$$

- This is often written as

$$y_i = \beta_1 + \beta_2 Post_i + \beta_3 Treat_i + \beta_4 Post_i \times Treat_i + u_i.$$

- The difference-in-differences estimate is β_4 .
- The advantages of using an OLS regression are
 - ▶ 1. A t -test of $H_0 : \beta_4 = 0$ is a test of statistical significance of the treatment
 - ▶ 2. We can add control variables as additional regressors.
 - ▶ 3. We can compute robust standard errors of $\hat{\beta}_4$.

Example: Access to health care and health outcomes

- Does better access to health care lead to better health outcomes?
- Dataset AED_HEALTHACCESS has data on 1,071 South African children aged 1 to 4 years in 54 communities.
- In 1993 26 of 54 communities had access to a health care clinic.
- In 1998 all 54 communities had access to a health care clinic.
- Outcome y is waz is a weight-for-age z -score
- Treatment $d = 1$ if have access to a health care clinic.
- Time $t = 0$ in 1993 (pre-period) and $t = 1$ in 1998 (post-period).

Example (continued)

- Summary statistics for key variables

Variable name	Storage type	Display format	Value label	Variable label
waz	double	%6.2f		Weight for age z Score
hightreat	float	%9.0g		= 1 if community has clinic in 1993
post	float	%9.0g		= 1 if year==98 and =0 if year==93
postXhigh	float	%9.0g		= post times hightreat
waz	double	%6.2f		Weight for age z Score
whz	double	%6.2f		Weight for height z Score

```
. summarize waz hightreat post postXhigh waz whz
```

Variable	Obs	Mean	Std. dev.	Min	Max
waz	1,071	-.205873	1.587432	-5.88	4.94
hightreat	1,071	.4276377	.4949671	0	1
post	1,071	.4668534	.4991332	0	1
postXhigh	1,071	.1979458	.3986373	0	1
waz	1,071	-.205873	1.587432	-5.88	4.94
whz	1,071	.6390009	2.199942	-9.89	9.99

Results: Manual computation

- The following table gives the mean values of *waz*
 - ▶ for the high treated and low treated children
 - ▶ before and after the expansion in free health care.

	High treated	Low treated
Before (1993)	-0.545 ($n = 246$)	-0.414 ($n = 325$)
After (1998)	0.321 ($n = 212$)	-0.069 ($n = 288$)
Change over time	0.867	0.345
Difference in differences		0.521

- High treated: *waz* increased by 0.867, from -0.545 to 0.321.
- Low treated: *waz* increased by 0.345, from -0.414 to -0.069.
- DID estimate is $0.867 - 0.345 = 0.521$.
- This is a very substantial effect
 - ▶ a third of a standard deviation change in *waz* for this sample.

Results: Regression computation

- Again greater access to health clinics increased waz by 0.521
- Since the treatment was at the community level, use cluster-robust standard errors with clustering on community
 - ▶ the standard error is 0.236 whereas heteroskedastic-robust s.e. is 0.194.

```
. * Diff-in-diff - no controls and cluster-robust standard errors
. reg waz postXhigh post hightreat, vce(cluster idcommunity) noheader
      (Std. err. adjusted for 54 clusters in idcommunity)
```

waz	Coefficient	Robust std. err.	t	P> t	[95% conf. interval]	
postXhigh	.5216188	.2352991	2.22	0.031	.0496685	.993569
post	.3450874	.1371018	2.52	0.015	.070096	.6200788
hightreat	-.1310593	.1968084	-0.67	0.508	-.525807	.2636884
_cons	-.4141846	.1151423	-3.60	0.001	-.6451308	-.1832384

Further analysis

- A richer and better model
 - ▶ controls for community by adding fixed effects for each community
 - ▶ controls for each individual by adding regressors such as parental education and household income
- For child i in community c
 - ▶ $y_{ic} = \beta_1 + \beta_2 t_i + \beta_3 d_i + \beta_4 t_i \times d_i + \gamma_c + \beta_5 x_{ic} + \dots + u_i.$

```
. * D in D with fixed effects for community and individual controls
. reg waz postXhigh post hightreat i.idcommunity ///
> fedu medu hhsized lntotmnc immuniz nonclinic, ///
> vce(cluster idcommunity) noheader
note: 242.idcommunity omitted because of collinearity.
      (Std. err. adjusted for 54 clusters in idcommunity)
```

waz	Coefficient	Robust std. err.	t	P> t	[95% conf. interval]	
postXhigh	.6428807	.2710993	2.37	0.021	.0991243	1.186637
post	-.6807024	.3487963	-1.95	0.056	-1.380299	.0188944
hightreat	-.2911247	.2360665	-1.23	0.223	-.7646142	.1823648

Further Details

- Differences-in-differences analysis is not restricted to one with time.
 - ▶ e.g. we might have a policy that affected only single women.
 - ▶ then compare the difference between married and single women with the difference between married and single men
 - ★ assuming that without the policy the change from married to single would be the same for men and women.

Differences in Differences: Multiple Time Periods

- A more general model has several periods of data.
- Suppose we have individual i in state s in year t , and the treatment of interest d_{st} occurs at the state-year level.
- Then we estimate the two-way fixed effects model
 - ▶ $y_{ist} = \phi_s + \gamma_t + \alpha d_{st} + \beta_1 x_{1ist} + \dots + u_{ist}$
 - ▶ here ϕ_s and γ_t are state-specific and time-specific fixed effects.
- The key assumption is that of “parallel trends”
 - ▶ the time trend each period is the same for each state
 - ★ γ_t is the same for each state
 - ▶ this is partly testable in some applications using pretreatment data.
- Inference is based on standard errors clustered at the state (s) level
 - ▶ this leads to the “few clusters” problem if there are few clusters.
- While estimation by OLS is straightforward, interpretation is difficult if the treatment is binary ($d_{st} = 1$ or 0) and treatment is staggered, occurring at different times for different states.
 - ▶ this is an area of current academic research.

Example: From Stata Documentation

- Example
 - ▶ y outcome is `satis` (Patient satisfaction score)
 - ▶ d treatment is `procedure = 1`
 - ▶ s group is hospital (there are 46)
 - ▶ t time is month (there are 7 months: January to July)
 - ▶ i is individual
- Treatment begins in April at 18 of the 46 hospitals

- Summary statistics

```
. summarize
```

Variable	Obs	Mean	Std. dev.	Min	Max
hospital	7,368	22.83822	13.57186	1	46
frequency	7,368	2.473398	1.163957	1	4
month	7,368	3.625	2.117778	1	7
procedure	7,368	.2079262	.4058512	0	1
satis	7,368	3.619074	1.05576	.5467862	9.712885

- Stata didregress command: Treatment effect is 0.84789

```
. didregress (satis)(procedure), group(hospital) time(month)
```

Treatment and time information

Time variable: month

Control: procedure = 0

Treatment: procedure = 1

	Control	Treatment
Group		
hospital	28	18
Time		
Minimum	1	4
Maximum	1	4

Difference-in-differences regression
Data type: Repeated cross-sectional

Number of obs = 7,368

(Std. err. adjusted for 46 clusters in hospital)

	satis	Robust Coefficient	std. err.	t	P> t	[95% conf. interval]
ATET						
procedure						
(New vs Old)		.8479879	.0321121	26.41	0.000	.7833108 .912665

Note: ATET estimate adjusted for group effects and time effects.

- Following gives the same estimate using regress

```
. * The following gives the same results as didregress
. regress satis procedure i.hospital i.month, vce(cluster hospital)
```

```
Linear regression                Number of obs   =       7,368
                                F(6, 45)        =           .
                                Prob > F            =           .
                                R-squared           =       0.5333
                                Root MSE       =       .72384
```

(Std. err. adjusted for 46 clusters in hospital)

	Coefficient	Robust std. err.	t	P> t	[95% conf. interval]	
satis						
procedure	.8479879	.0321121	26.41	0.000	.7833108	.912665

- Visual test of parallel trends assumption
 - ▶ `estat trendplots`
- Formal test of parallel trends assumption
 - ▶ `estat ptrends`
- For more read the Stata pdf documentation.

References for DID

- These books are given in approximate order of increasing difficulty.
- A. Colin Cameron (2022), Analysis of Economics Data: An Introduction to Econometrics, chapter 13.6.
- Joshua D. Angrist and Jörn-Steffen Pischke (2015), Mastering Metrics, ch. 5.
- Cunningham, Scott (2021), Causal Inference: The MixTape, Yale UP, chapter 9.
- A. Colin Cameron and Pravin K. Trivedi (2022), Microeconometrics using Stata: Volumes 1 and 2, Second Edition, Stata Press, chapter 25.6.
- Joshua D. Angrist and Jörn-Steffen Pischke (2009), Mostly Harmless Econometrics: An Empiricist's Companion, Princeton University Press, chapter 5.
- A. Colin Cameron and Pravin K. Trivedi (2005), Microeconometrics: Methods and Applications, Cambridge University Press, chapter 22.6.
- Jeffrey M. Wooldridge, (2010), Econometric Analysis of Cross Section and Panel Data, Second Edition, MIT Press, chapter 6.5.

References for DID (continued)

- These books by non-economists are similar to *Mastering Metrics* in accessibility.
- Stephen L. Morgan and Christopher Winship (2015), *Counterfactuals and Causal Inference: Methods and Principles for Social Research*, Second edition, Cambridge University Press, chapter 11.3.
- Andrew Gelman, Jennifer Hill and Aki Vehtari (2022), *Regression and Other Stories*, Cambridge University Press, especially chapters 21.4.
- These are current more advanced econometrics articles
- B. Callaway and P.H.C. Sant'Anna (2021), "Difference-in-Differences with multiple time periods," *Journal of Econometrics*, 225, pages 200–230.
- Jeffrey M. Wooldridge (2021), "Two-way fixed effects, the two-way Mundlak regression, and difference-in-differences estimators,"
<http://doi.org/10.2139/ssrn.3906345>.