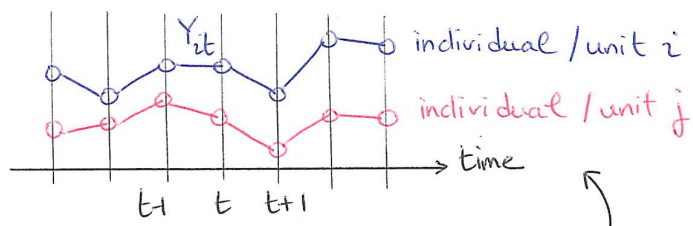


## CI = PANEL DATA METHODS

Panel data (aka longitudinal data or cross-sectional time series) is data that include observations of the same (groups of) units (such as individuals, firms, households, products) over multiple time periods.



[ $Y_{it}$  = observation for individual  $i$  at time  $t$ .]

↳ We want to use panel data to estimate the effect of an intervention that affects some units in some time periods.

(generalisation from previous chapters where iid observations were considered)

This chapter covers:

- ↳ the difference estimator with panel data in an RCT
- ↳ the difference-in-difference estimator in an observational study & its application to RCTs.
- ↳ generalisation to stratified designs: block by block analysis & IPW
- ↳ generalisation to clustered experiments.

## I. THE DIFFERENCE ESTIMATOR

2

Consider  $n$  units  $i = 1, \dots, n$ , each receiving a binary treatment assignment  $W_i \in \{0, 1\}$  completely at random. Let  $n_t$  denotes the number of treated units, and  $n_c = n - n_t$  the number of control units. In an RCT,  $W_i \perp \{Y_i(0), Y_i(1)\}$ , where  $Y_i(j)$  is the Potential Outcome of unit  $i$  receiving treatment  $j$  [ $(j=1)$  = treated unit ;  $(j=0)$  = control unit].

For each unit  $i$ , we observe the potential outcome corresponding to its treatment assignment  $Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$ . Throughout this chapter, we assume an infinite superpopulation model, so that  $\{Y_i(0), Y_i(1)\} \sim \mathbb{P}$ . We are interested in this section in the estimation of the  $ATE = \Delta^{\infty}$

$$= \mathbb{E}(Y_i(1) - Y_i(0)).$$

$\mathbb{E}(\cdot)$  under the joint probability  $\mathbb{P}$ .

No panel data yet: each unit is associated with a single observation  $Y_i$  [no time index]

In chapter CI: RANDOMISED CONTROL TRIALS, we introduced the unbiased & consistent difference estimator of  $\Delta^{\infty}$ ,

$$\hat{\Delta} = \frac{1}{n_t} \sum_{i=1}^n W_i Y_i - \frac{1}{n_c} \sum_{i=1}^n (1 - W_i) Y_i.$$

We re-express  $\hat{\Delta}$  as the OLS estimate of  $\Delta^0$  in the linear model (3)

$$Y_i = \beta_0 + \Delta^0 W_i + \varepsilon_i.$$

[DIFF|OLS]

In matrix notation,  $Y = X\beta + \varepsilon$  where  $\beta = \begin{pmatrix} \beta_0 \\ \Delta^0 \end{pmatrix}$ ,

$$X = \begin{pmatrix} 1 & W_1 \\ \vdots & \vdots \\ 1 & W_n \end{pmatrix}, \quad Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix},$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\Delta} \end{pmatrix} = (X^t X)^{-1} X^t Y; \quad \hat{\Delta} = \text{diff-estimator}$$

Assuming  $E\varepsilon = 0$  and writing  $\text{Cov}(\varepsilon) = E\varepsilon\varepsilon^t = \Sigma_\varepsilon$

the covariance matrix of  $\hat{\beta}$  is

$$\Sigma_{\hat{\beta}} = (X^t X)^{-1} (X^t \Sigma_\varepsilon X) (X^t X)^{-1}.$$

This linear representation together with the expression  $\Sigma_{\hat{\beta}}$  is particularly useful for the construction of confidence intervals for  $\Delta$  under various assumptions on the correlation structure of the residual error (this appears to be crucial when later on extending the present results to panel data).

x Example 1: Homoskedastic errors:  $\Sigma_\varepsilon = \sigma^2 \underline{I}_n$ ,

so that  $\sigma^2 = \text{Var } \varepsilon_i$  for all  $i$ . In this case

$$\Sigma_{\hat{\beta}} = \sigma^2 (X^t X)^{-1} \quad (\text{plug-in estimator for } \sigma^2)$$

x Example 2: Heteroskedastic errors;  $\text{var } \varepsilon_i = \sigma_i^2$  (4)

$$\Sigma_\varepsilon = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{pmatrix} \leftarrow \text{independent errors.}$$

$$\text{Then } \Sigma_{\hat{\beta}} = (X^t X)^{-1} \left( \sum_{i=1}^n \sigma_i^2 X_i X_i^t \right) (X^t X)^{-1},$$

$$\text{where } X = \begin{pmatrix} -X_1^t - \\ \vdots \\ -X_n^t - \end{pmatrix}.$$

A consistent estimator of  $\Sigma_{\hat{\beta}}$  in this case is given by the Eicker-Huber-White (EHW) estimator:

$$\hat{\Sigma}_{\text{EHW}} = (X^t X)^{-1} \left( \sum_{i=1}^n \hat{\varepsilon}_i^2 X_i X_i^t \right) (X^t X)^{-1} \\ = (X^t X)^{-1} \left( X^t \begin{pmatrix} \hat{\varepsilon}_1^2 & & 0 \\ & \ddots & \\ 0 & & \hat{\varepsilon}_n^2 \end{pmatrix} X \right) (X^t X)^{-1},$$

where  $\hat{\varepsilon}_i = i$ -th residual error  $= Y_i - X_i^t \hat{\beta}$ . Standard errors based on the EHW estimator are commonly called ROBUST.

↳ The consistency of  $\hat{\Sigma}_{\text{EHW}}$  is derived in White (1980) under uniformly bounded assumptions of the error variances & covariance matrix of the regressors. Together with the asymptotic normality of  $\sqrt{n}(\hat{\beta} - \beta)$ , the EHW estimator can be used to construct appropriate confidence intervals with a desired nominal coverage.

x Example 3 = More generally, we may consider a block diagonal matrix

(5)

$$\Sigma_{\varepsilon} = \begin{pmatrix} \Sigma_1 & 0 & & 0 \\ 0 & \Sigma_2 & & \\ & & \ddots & \\ 0 & 0 & & \Sigma_B \end{pmatrix} \begin{matrix} \updownarrow n_1 \\ \updownarrow n_2 \\ \vdots \\ \updownarrow n_B \end{matrix} \quad [n=n_1+\dots+n_B]$$

$$= (\Sigma_{ij})_{\substack{i=1,\dots,n \\ j=1,\dots,n}} \quad \uparrow \text{B blocs}$$

Units are clustered and assumed uncorrelated provided they belong to two different clusters.

To start with, we assume further the homoskedastic structure

$$\Sigma_k = \begin{pmatrix} \sigma^2 & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 \end{pmatrix} \quad \forall k=1,\dots,B$$

constant correlation  $\rho$  for two units belonging to the same cluster.

In this case,

$$\Sigma_{\beta} = (X^t X)^{-1} \left( \sum_{k=1}^B X_k^t \Sigma_k X_k \right) (X^t X)^{-1}$$

with

$$X = \begin{pmatrix} X_1 \\ X_B \end{pmatrix} \begin{matrix} \updownarrow n_1 \\ \updownarrow n_B \end{matrix}$$

$\leftarrow 2 \rightarrow$

In an RCT, units clustered together typically receive

the same treatment status:  $W_i = W_j \quad \forall i, j \in C_k$ , (6)

where  $C_k$  denotes the  $k$ -th cluster of size  $n_k$ .

$$Y_{ik} = \beta_0 + \Delta^0 W_k + \varepsilon_{ik}$$

$i=1,\dots,n_k$   
 $k=1,\dots,B$

In this case,  $X = \begin{pmatrix} X_1 \\ \vdots \\ X_B \end{pmatrix}$  with

$$X_k = \begin{pmatrix} 1 & W_k \\ \vdots & \vdots \\ 1 & W_k \end{pmatrix} \begin{matrix} \updownarrow n_k \\ \leftarrow 2 \rightarrow \end{matrix} = \mathbb{1}_{n_k} \tilde{X}_k^t \quad ; \quad \begin{cases} \mathbb{1}_{n_k} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ \tilde{X}_k = \begin{pmatrix} 1 \\ W_k \end{pmatrix} \end{cases}$$

$$\Rightarrow X^t X = \sum_{k=1}^B n_k \tilde{X}_k \tilde{X}_k^t$$

$$\Rightarrow X^t \Sigma_{\varepsilon} X = \sum_{k=1}^B X_k^t \Sigma_k X_k = \sum_{k=1}^B \tilde{X}_k \mathbb{1}_{n_k}^t \Sigma_k \mathbb{1}_{n_k} \tilde{X}_k^t$$

with  $\Sigma_k = \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$

$$\bullet \Sigma_k \mathbb{1}_{n_k} = \left( 1 + (n_k - 1)\rho \right) \sigma^2$$

$$\bullet \mathbb{1}_{n_k}^t \Sigma_k \mathbb{1}_{n_k} = n_k \sigma^2 (1 + (n_k - 1)\rho)$$



$$X^t \Sigma_{\varepsilon} X = \sum_{k=1}^B \sigma^2 n_k (1 + (n_k - 1)\rho) \tilde{X}_k \tilde{X}_k^t \quad (7)$$

defining  $\tau_k = \sum_{k=1}^B \sigma^2 n_k \tau_k \tilde{X}_k \tilde{X}_k^t$

$$\tau_k = 1 + (n_k - 1)\rho$$

$$\Rightarrow \Sigma_{\beta}^{-1} = \sigma^2 \left( \sum_k n_k \tilde{X}_k \tilde{X}_k^t \right)^{-1} \left( \sum_k n_k \tau_k \tilde{X}_k \tilde{X}_k^t \right) \left( \sum_k n_k \tilde{X}_k \tilde{X}_k^t \right)$$

In the special case where all groups have equal size  $n_k = m \quad \forall k=1, \dots, B$ ,  $\tau_k = 1 + (m-1)\rho \equiv \tau$ ,

and

$$\Sigma_{\beta}^{-1} = \sigma^2 \left( \sum_{k=1}^B m \tilde{X}_k \tilde{X}_k^t \right)^{-1} \tau$$

- Covariance of  $\beta$  ignoring the correlation structure -

MOULTON FACTOR

$\Rightarrow$  The correlation structure of  $\varepsilon$  inflates the (co)variance by a factor  $\tau$ . The Moulton Factor is maximal when  $\rho = 1$ , in which case  $\tau = m$ : each cluster carries a single unit of information. See Moulton (1990).

\* Remark: The bloc-diagonal covariance structure can be justified using the representation

$$\varepsilon_{ik} = \underbrace{u_k}_{\text{zero mean}} + \underbrace{\eta_{ik}}_{\text{additive random effect model}}$$

$\&$  variance  $\sigma_u^2, \sigma_{\eta}^2, \text{ iid}$

Then  $\text{cov}(\varepsilon_{ik}, \varepsilon_{jk}) = \mathbb{E}(u_k + \eta_{ik})(u_k + \eta_{jk}) \quad (8)$

$$= \sigma_u^2 \quad i \neq j$$

$$\text{var}(\varepsilon_{ik}) = \sigma_u^2 + \sigma_{\eta}^2 \equiv \sigma_{\varepsilon}^2$$

$$\text{corr}(\varepsilon_{ik}, \varepsilon_{jk}) = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_{\eta}^2} \equiv \rho$$

$\&$  thus  $\text{cov}(\varepsilon_{ik}, \varepsilon_{jk}) = \rho(\sigma_u^2 + \sigma_{\eta}^2) = \rho \sigma_{\varepsilon}^2$

Liang & Zeger (1986) relax Moulton's model of constant correlation and consider general covariance matrices  $\Sigma_k : \Sigma_{\beta}^{-1} = (X^t X)^{-1} \left( \sum_{k=1}^B X_k^t \Sigma_k X_k \right) (X^t X)^{-1}$

Estimated using

$$\hat{\Sigma}_{LZ} = (X^t X)^{-1} \sum_{k=1}^B \left( \sum_{i=1}^{n_k} (Y_{ik} - \hat{\beta}^t X_{ik}) X_{ik} \right) \left( \sum_{i=1}^{n_k} (Y_{ik} - \hat{\beta}^t X_{ik}) X_{ik} \right)^t \times (X^t X)^{-1}$$

- plugged-in residuals -  
since

$$\sum_{k=1}^B X_k^t \Sigma_k X_k = \mathbb{E} \left\{ \sum_{k=1}^B \underbrace{X_k^t}_{(2 \times n_k)} \underbrace{\Sigma_k}_{(n_k \times n_k)} \underbrace{X_k}_{(n_k \times 2)} \right\}$$

$$= \mathbb{E} \left\{ \sum_{k=1}^B \left( \sum_{i=1}^{n_k} \underbrace{\varepsilon_i}_{(2 \times 1)} X_{ik} \right) \left( \sum_{i=1}^{n_k} \underbrace{\varepsilon_i}_{(1 \times 2)} X_{ik}^t \right) \right\}$$

$$X_k = \begin{bmatrix} - & X_{ik}^t & - \\ \vdots & & \vdots \end{bmatrix}$$

$(n_k \times 2)$   $\updownarrow n_k$

$$X = \begin{pmatrix} \boxed{X_1} \\ \vdots \\ \boxed{X_B} \end{pmatrix}$$

$(n \times 2)$   $\updownarrow n_1$   $\updownarrow n_B$



Summary :  $\hat{\Delta}$  is the OLS estimate of  $\Delta^\infty$  in the linear model  $Y_i = \beta_0 + \Delta^\infty W_i + \varepsilon_i$ . We established that under (i) homoskedastic errors, (ii) heteroskedastic errors and (iii) general block-diagonal covariance structure for  $\varepsilon$ , we can consistently estimate the covariance matrix of  $\begin{pmatrix} \hat{\beta}_0 \\ \hat{\Delta} \end{pmatrix}$ , and thus construct (asymptotically) valid confidence intervals for the difference in means estimator of the ATE.

x Remark : Relative Scale

In the linear representation  $Y_i = \beta_0 + \Delta^\infty W_i + \varepsilon_i$ ,  
 $\beta_0 = E(Y_i | W_i=0) = E(Y_i(0) | W_i=0)$   
 $\Delta^\infty = E(Y_i | W_i=1) - E(Y_i | W_i=0)$   
 $= E(Y_i(1) | W_i=1) - E(Y_i(0) | W_i=0)$   
 $= E(Y_i(1) - Y_i(0))$

$\Rightarrow \delta^\infty := \frac{E(Y_i(1) - Y_i(0))}{E(Y_i(0))} = \frac{\Delta^\infty}{\beta_0} = \text{relative "lift"}$

A natural estimator for  $\delta^\infty$  is  $\hat{\delta}^\infty := \frac{\hat{\Delta}^\infty}{\hat{\beta}_0}$ ; whose variance can be easily derived from  $\Sigma_{\hat{\beta}}$  using the delta method

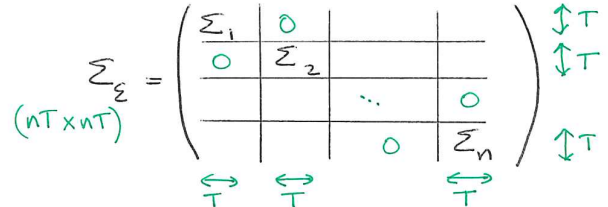
$$\text{var}\left(\frac{X}{Y}\right) \approx \frac{\mu_X^2}{\mu_Y^2} \left( \frac{\sigma_X^2}{\mu_X^2} - 2 \frac{\text{cov}(X, Y)}{\mu_X \mu_Y} + \frac{\sigma_Y^2}{\mu_Y^2} \right)$$

• Panel Data : Units  $i=1, \dots, n$  are observed over some period of time  $t=1, \dots, T$ . Each unit  $i$  receives a treatment status  $W_i \in \{0, 1\}$  that they keep throughout the experiment.  $\rightarrow$  completely at random

Assume a common baseline  $\gamma_0$  for all P.O.  $Y_{it}(0)$ , so that  $Y_{it}(0) = \gamma_0 + \eta_{it}$ ;  $E\eta_{it} = 0$ . In addition, let  $\Delta_{it}$  be the effect of the treatment on unit  $i$  at time  $t$ ,  $Y_{it}(1) = Y_{it}(0) + \Delta_{it}$ .  
 $\rightarrow$  may vary across units and time.  
 $\rightarrow$  assume a general bivariate distribution for  $\Delta_{it}$ ; over units & time.

x Estimation Method :  $Y_{it} = \beta_0 + \Delta W_i + \varepsilon_{it}$ ,  $E\varepsilon_{it} = 0$  where  $Y_{it} = W_i Y_{it}(1) + (1-W_i) Y_{it}(0)$ .

$\rightarrow$  Errors are clustered over time



The OLS estimator of  $\Delta$  is

$$\hat{\Delta} = \frac{1}{n_e T} \sum_{i|W_i=1} \sum_{t=1}^T Y_{it} - \frac{1}{n_c T} \sum_{i|W_i=0} \sum_{t=1}^T Y_{it}$$

$$\xrightarrow{n_e, n_c \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) | W_i=1) - \frac{1}{T} \sum_{t=1}^T E(Y_{it}(0) | W_i=0)$$
  

$$= \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) - Y_{it}(0)) = \text{averaged treatment effect over the test period.}$$

## II - DIFFERENCE-IN-DIFFERENCES

(11)

We relax the common baseline assumption of page 10 and put

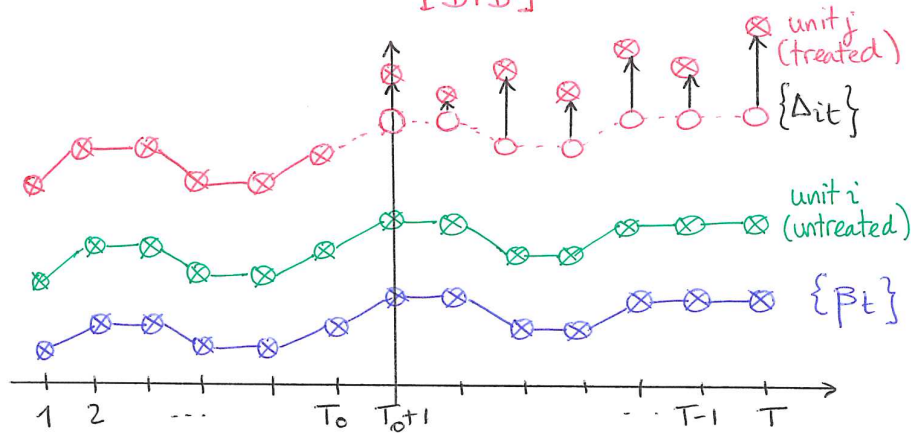
$$Y_{it}(0) = \alpha_i + \beta_t + \varepsilon_{it}$$

unit-specific effect  $\alpha_i$     time-specific effect  $\beta_t$

⇒ Units evolve in parallel following a baseline time series  $\{\beta_t\}$ , shifted by unit-specific quantities  $\alpha_i$ .

The diff-in-diff model is  $Y_{it} = \alpha_i + \beta_t + \Delta_{it}W_{it} + \varepsilon_{it}$

[DID]



x Estimation Method: TWFE

$$Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}$$

Let  $\hat{\Delta}$  denote the OLS estimator of  $\Delta$  in the TWFE model

Remark: Important to dissociate the TWFE model with the [DID] assumption. In section II.3 page 18 we introduce another linear model for estimation of a causal parameter within the [DID] framework.

12  
 $t = 1, \dots, T_0, T_0+1, \dots, T$  where indexes 1 to  $T_0$  denote observations before the treatment starts.  
 $t = T_0+1, \dots, T$  denotes the test period.

For treated units, we assume that the treatment starts at time  $t = T_0+1$  for all of them. In a subsequent section, we discuss the case of a staggered rollout.

[For a treated unit  $i$ ,  $W_{it} = 1 \forall t = T_0+1, \dots, T$ ]

[For all units,  $W_{it} = 0 \forall t = 1, \dots, T_0$ ]

The treatment is not necessarily randomized here (compared with section I & previous chapters). The TWFE model is usually considered in an observational study, where some units are more likely to receive the treatment than others.

In Appendix A, we show that the OLS estimator of  $\Delta$  is

$$\hat{\Delta} = \left\{ \frac{1}{n_t(T-T_0)} \sum_{i \in \text{trt}} \sum_{t \geq T_0+1} Y_{it} - \frac{1}{n_t T_0} \sum_{i \in \text{trt}} \sum_{t \leq T_0} Y_{it} \right\} - \left\{ \frac{1}{n_c(T-T_0)} \sum_{i \in \text{ctrl}} \sum_{t \geq T_0+1} Y_{it} - \frac{1}{n_c T_0} \sum_{i \in \text{ctrl}} \sum_{t \leq T_0} Y_{it} \right\}$$

↗ mean diff in the trt group  
↖ mean diff in the ctrl group

↑ The difference-in-differences estimator

Q: Which quantity does  $\hat{\Delta}$  identify?

## II.1. Identification with two time periods

(13)

In this section, we assume that  $T_0 = 1, T = 2$ : for each unit, we observe a single value  $Y_{i1}$  in pre-test, and  $Y_{i2}$  in test period. Then

$$\hat{\Delta} = \frac{1}{n_t} \sum_{i \in \text{trt}} (Y_{i2} - Y_{i1}) - \frac{1}{n_c} \sum_{i \in \text{ctrl}} (Y_{i2} - Y_{i1})$$

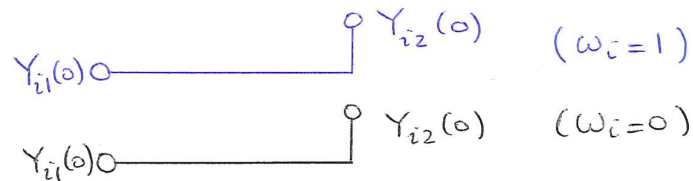
$$\xrightarrow{n_t, n_c \rightarrow \infty} \mathbb{E}(Y_{i2} - Y_{i1} | W_i = 1) - \mathbb{E}(Y_{i2} - Y_{i1} | W_i = 0)$$

Identification involves P.O. only.  
To get rid of observed values  $Y_{it}$ , we make the following two assumptions.

### (I) Parallel trends

$$\mathbb{E}(Y_{i2}(0) - Y_{i1}(0) | W_i = 1) = \mathbb{E}(Y_{i2}(0) - Y_{i1}(0) | W_i = 0)$$

[In absence of treatment, the two groups evolve by the same amount (on average)]



(II) No anticipation:  $Y_{i1}(0) = Y_{i1}(1)$ ,  $i$  treated.  
Individuals in the treatment group do not anticipate the upcoming treatment in pre-test period

Result: Under (I) and (II), the OLS estimator  $\hat{\Delta}$  of  $\Delta$  in the TWFE model  $Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}$  identifies the ATT =  $\Delta = \mathbb{E}(Y_{i2}(1) - Y_{i2}(0) | W_i = 1)$

Indeed,

$$\text{ATT} = \mathbb{E}(Y_{i2}(1) - Y_{i2}(0) | W_i = 1)$$

$$= \mathbb{E}(Y_{i2}(1) | W_i = 1) - \underbrace{\mathbb{E}(Y_{i2}(0) | W_i = 1)}$$

$$\text{[Under (I)]} \left( \begin{array}{l} \mathbb{E}(Y_{i1}(0) | W_i = 1) \\ + \mathbb{E}(Y_{i2}(0) | W_i = 0) \\ - \mathbb{E}(Y_{i1}(0) | W_i = 0) \end{array} \right)$$

$$= \left\{ \mathbb{E}(Y_{i2}(1) | W_i = 1) - \mathbb{E}(Y_{i1}(0) | W_i = 1) \right\}$$

$$- \left\{ \mathbb{E}(Y_{i2}(0) | W_i = 0) - \mathbb{E}(Y_{i1}(0) | W_i = 0) \right\}$$

$$\text{[Under (II)]} = \mathbb{E}(Y_{i1}(1) | W_i = 1)$$

$$= \left\{ \mathbb{E}(Y_{i2}(1) | W_i = 1) - \mathbb{E}(Y_{i1}(1) | W_i = 1) \right\}$$

$$- \left\{ \mathbb{E}(Y_{i2}(0) | W_i = 0) - \mathbb{E}(Y_{i1}(0) | W_i = 0) \right\}$$

$$= \mathbb{E}(Y_{i2} - Y_{i1} | W_i = 1) - \mathbb{E}(Y_{i2} - Y_{i1} | W_i = 0)$$

$$= \Delta$$

$$= \text{limit of } \hat{\Delta} \text{ (p.13) as } n_t, n_c \rightarrow \infty$$

x Remark: Abuse of notation:  $\begin{pmatrix} Y_{it}(0) = Y_{it}(0,0) \\ Y_{it}(1) = Y_{it}(0,1) \end{pmatrix}$    
  $\downarrow$  pre-period  $\square$    
  $\leftarrow$  test period   
 "Pull trajectories"



x Remark : The ATNT :=  $\mathbb{E}(Y_{i2}(1) - Y_{i2}(0) | W_i = 0)$  (15)  
is identified under

$$(I') \quad \mathbb{E}(Y_{i2}(1) - Y_{i1}(1) | W_i = 1) \\ = \mathbb{E}(Y_{i2}(1) - Y_{i1}(1) | W_i = 0)$$

$$(II') \quad \mathbb{E}(Y_{i1}(0) - Y_{i1}(1) | W_i = 0) = 0$$

i.e. under (I'), (II'),  $\Delta = \text{ATNT}$ .

When  $W_{i2}$  are drawn at random, (I), (I'), (II), (II') are satisfied and we identify the ATE:

$$\begin{aligned} \text{ATE} &= \text{ATT} \times \mathbb{P}(W_{i2} = 1) + \text{ATNT} \times \mathbb{P}(W_{i2} = 0) \\ &= \Delta \times \mathbb{P}(W_{i2} = 1) + \Delta \times \mathbb{P}(W_{i2} = 0) \\ &= \Delta \end{aligned}$$

## III.2. Identification in the general case.

Assumptions (I) and (II) generalize nicely in the multi-period case. Put

$$\text{ATT}(T_0, T) = \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(1) - Y_{it}(0) | W_i = 1)$$

= average ATT over  $t = T_0 + 1, \dots, T$

(A) Parallel trends

$$\begin{aligned} &\frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 1) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 1) \\ &\quad \uparrow \text{test period} = \quad \uparrow \text{pre-test} \\ &\frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 0) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 0) \end{aligned}$$

(B) No Anticipation here again, we abused notation, as these must be understood with a path of treatment (p. 114) (16)  
 $\forall t = 1, \dots, T_0 : Y_{it}(0) = Y_{it}(1) \quad , \quad \forall i \text{ with } W_i = 1$

Result: Under (A) and (B), the OLS estimator  $\hat{\Delta}$  (bottom of page 12) converges to  $\text{ATT}(T_0, T)$  as  $n_t, n_c \rightarrow \infty$ .

Indeed,

$$\text{ATT}(T_0, T) = \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(1) - Y_{it}(0) | W_i = 1)$$

$$= \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(1) | W_i = 1)$$

$$- \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 1)$$

[Under (A)]

$$= \frac{1}{T - T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 1)$$

$$+ \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 0)$$

$$- \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 0)$$

[Under (B)]

$$= \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(1) | W_i = 1)$$

$$= \left\{ \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(1) | W_i = 1) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(1) | W_i = 1) \right\}$$

$$- \left\{ \frac{1}{T - T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 0) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 0) \right\}$$

x Remark = Instead, we may collapse the multi-period data to a single pre-test and test: consider the time aggregated outcomes

(17)

$$\begin{cases} Y_{i,pre} = \sum_{t=1}^{T_0} Y_{it} \\ Y_{i,test} = \sum_{t=T_0+1}^T Y_{it} \end{cases}$$

and compute the diff-in-diff estimator

$$\hat{\Delta} = \frac{1}{n_t} \sum_{i \in tr} (Y_{i,test} - Y_{i,pre}) - \frac{1}{n_c} \sum_{i \in ctr} (Y_{i,test} - Y_{i,pre})$$

The difference with  $\hat{\Delta}$  lies in the extra averaging  $\frac{1}{T_0}$  and  $\frac{1}{T-T_0}$  & therefore in the interpretation of the quantity  $\hat{\Delta}$  and  $\tilde{\Delta}$  identify.

$\hat{\Delta} \rightarrow$  per unit, per epoch  $t$   
 $\tilde{\Delta} \rightarrow$  per unit, over a whole time period.

Because of this,  $\tilde{\Delta}$  tends to have a higher variance than  $\hat{\Delta}$ .

$\rightarrow$  To account for temporal correlation when computing the variance of  $\hat{\Delta}$ , one may cluster errors over time, see p.10. Note that this is implicitly suggested in the two-periods case, where

$$\hat{\Delta} = \frac{1}{n_t} \sum_{i \in tr} (Y_{i2} - Y_{i1}) - \frac{1}{n_c} \sum_{i \in ctr} (Y_{i2} - Y_{i1})$$

$$\downarrow$$

$$\text{var } \hat{\Delta} = \frac{1}{n_t} \text{var} (Y_{i2} - Y_{i1} | W_i = 1) + \frac{1}{n_c} \text{var} (Y_{i2} - Y_{i1} | W_i = 0)$$

Temporal correlation is accounted for here since the variance of the difference  $(Y_{i2} - Y_{i1})$  is computed in each group.

(18)

$$\text{var}_1 (Y_{i2} - Y_{i1}) = \text{var } Y_{i2} + \text{var } Y_{i1} - 2 \text{cov} (Y_{i2}, Y_{i1})$$

↑  
 Shorthand for  
 $\text{var}(\dots | W_i = 1)$

If  $\leq 0$ , the resulting diff-in-diff estimator has smaller variance than the difference estimator. To get this,  $(Y_{i1}, Y_{i2})$  must be sufficiently correlated to counter balance  $\text{var } Y_{i1}$ .

### II.3. Alternative Representation.

Instead of the TWFE model  $Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}$ , the difference-in-differences estimator  $\hat{\Delta}$  (page 12) can be seen to be the OLS estimate of  $\Delta$  in the alternative linear representation

[DID] OLS

$$Y_{it} = a_0 + a_1 \mathbb{1}(i \text{ is in the treatment group}) + a_2 \mathbb{1}(t \geq T_0 + 1) + \Delta \mathbb{1}(t \geq T_0 + 1 \ \& \ i \in tr) + \varepsilon_{it}$$

This representation is particularly suited to compute relative effects and confidence intervals. Note that =

$$\bullet \hat{a}_0 \rightarrow \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it} | W_i = 0) \quad (\text{as } n_{t_1}, n_{t_2} \rightarrow \infty) \quad (19)$$

$$\bullet \hat{a}_1 \rightarrow \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it} | W_i = 0) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it} | W_i = 0)$$

$$\bullet \hat{a}_2 \rightarrow \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it} | W_i = 1) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it} | W_i = 0)$$

$$\Rightarrow \hat{a}_0 + \hat{a}_1 + \hat{a}_2 \rightarrow \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 0)$$

$$- \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 0)$$

$$+ \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(1) | W_i = 1)$$

No anticipation

$$= \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i = 1)$$

// trends

$$= \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i = 1)$$

= reference value for the treatment group in the test period.

$$\Rightarrow \frac{\hat{\Delta}}{\hat{a}_0 + \hat{a}_1 + \hat{a}_2} = \text{relative lift}$$

↑ variance is computed using the delta method.

Remark: Both this linear model and the TWFE allow identification of the ATT under the parallel trend assumption.

## II.4. Checking for pre-trends

(20)

We cannot check for parallel trends in the test period. However, we can check how plausible the assumption holds in pre-test. Consider the following TWFE model

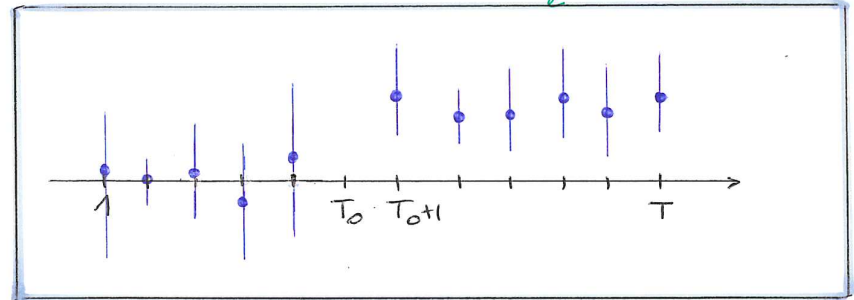
$$Y_{it} = \alpha_i + \beta_t + \sum_{\substack{s=1 \\ s \neq T_0}}^T \gamma_s W_i \mathbb{1}(t=s) + \varepsilon_{it}$$

One coefficient is arbitrarily removed to avoid overspecifying the linear model (here, the last day before intervention)

[see the Appendix for a formal proof] - p.47

& generate an "event study plot"

Coeffs  $\hat{\gamma}_s$  & confidence bounds



### Limitations

- L No guarantees that // trends hold in test period
- L Typically low power (often fail to reject the null)
- L Conditions the analysis on "passing pre-trends" → selection bias



## II.5. Limitations of the TWFE approach.

(21)

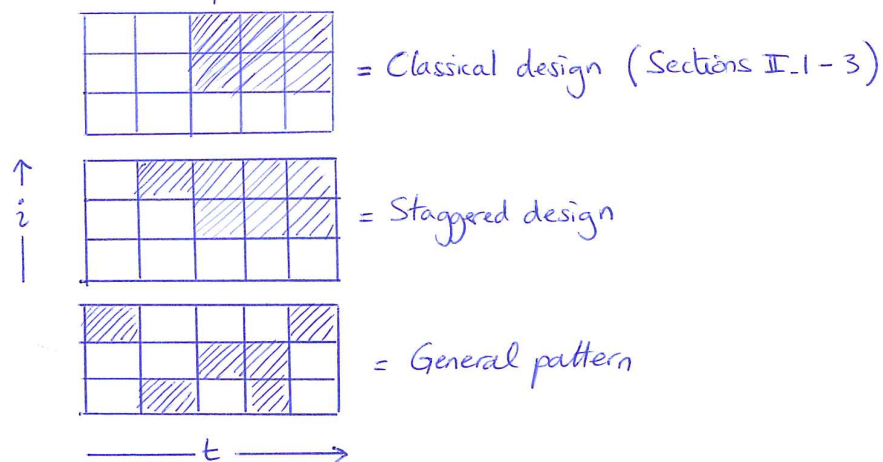
Consider the TWFE model (page 11) with general  $\{W_{it}\}$ :

$$Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}.$$

Let  $\Delta_{it} = Y_{it}(1) - Y_{it}(0) =$  unit  $i$  treatment effect.

When there is no heterogeneity between units & across time,  $E\Delta_{it} = \Delta_0 \forall (i, t)$ , the OLS estimate  $\hat{\Delta}$  of  $\Delta$  in the TWFE model recovers the correct value  $E\hat{\Delta} = \Delta_0$ .

↳ The discussion in this section holds true with general treatment patterns, such as



Issues arise when there is heterogeneity in the treatment effect either across units or time. We state next a simplified version of a result in de Chaisemartin & d'Haultfoeuille (2020), proved under the following assumptions:

(I) Balanced Panel: Observe  $Y_{it} \forall i, \forall t$

(II) Independent individuals

(22)

$(Y_{i1}(0), Y_{i1}(1), W_{i1}, \dots, Y_{iT}(0), Y_{iT}(1), W_{iT})$   
 $i=1, \dots, n$  are independent (time correlation is allowed)

(III) Common Trends

$\forall t \geq 2 \quad E(Y_{it}(0) - Y_{i,t-1}(0))$  independent of  $i$

(IV) Strong Exogeneity

Shocks are independent of the past, present & future treatments

$$E(Y_{it}(0) - Y_{i,t-1}(0) | W_{i1}, \dots, W_{iT}) = E(Y_{it}(0) - Y_{i,t-1}(0)).$$

Under (I), (II), (III), (IV), the OLS estimator  $\hat{\Delta}$  of  $\Delta$  in  $Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}$  satisfies

$$E\hat{\Delta} = E \left[ \sum_{(i,t) | W_{it}=1} r_{it} \Delta_{it} \right]$$

where  $\rightarrow r_{it} = \varepsilon_{it} / \sum_{(i,t) | W_{it}=1} \varepsilon_{it}$

$\rightarrow \varepsilon_{it} =$  residual in  $W_{it} = \gamma + \gamma_i + \delta_t + \varepsilon_{it}$

FE

↳ Weights sum to one; but can be negative

↳ When the treatment effect does not vary across units & time, the OLS estimator  $\hat{\Delta}$  is unbiased for the ATT.

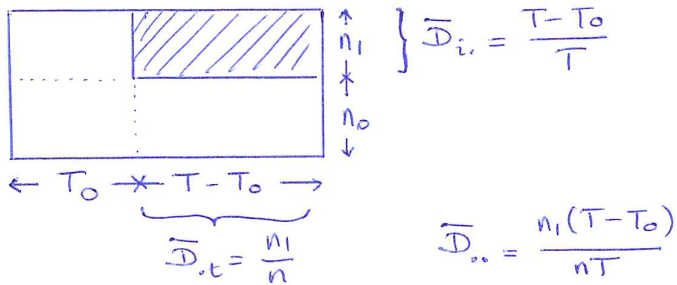
• Example 1: Classical design.

$T_0 =$  # epochs in pre-test ;  $T - T_0 =$  # epochs in test

$n_1 = \#$  treated units,  $n_0 = \#$  control units

(23)

Then one can show that  $\varepsilon_{it} = W_{it} - \bar{W}_{i.} - \bar{W}_{.t} + \bar{W}_{..}$ .



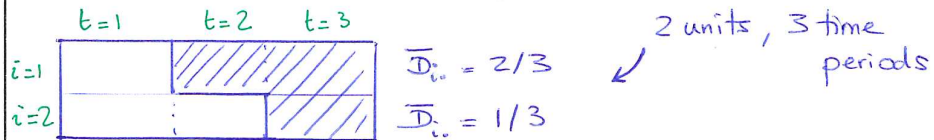
$\varepsilon_{it} = \frac{n_0 T_0}{nT} \quad \forall (i,t) | W_{it} = 1$

$r_{it} = \frac{1}{n_1(T-T_0)}$

$\Rightarrow E\hat{\Delta} = \frac{1}{n_1(T-T_0)} \sum_{(i,t) | W_{it}=1} E\Delta_{it}$ , as required.

Heterogeneity of treatment effects is not an issue here.

• Example 2 = Staggered design



$\varepsilon_{it} = \begin{cases} 1 - 2/3 - 1/2 + 1/2 = 1/3 \\ 1 - 2/3 - 1 + 1/2 = -1/6 \\ 1 - 1/3 - 1 + 1/2 = 1/6 \end{cases}$

$r_{it} = \begin{cases} 1 \\ -1/2 \\ 1/2 \end{cases}$

$\Rightarrow E\hat{\Delta} = 1 \times E\Delta_{1,2} - \frac{1}{2} E\Delta_{1,3} + \frac{1}{2} E\Delta_{2,3}$

$\neq \frac{1}{3} (E\Delta_{1,2} + E\Delta_{1,3} + E\Delta_{2,3})$ .

(24)

Worse,  $E\hat{\Delta}$  may be negative when all  $E\Delta_{it}$  are  $> 0$ .

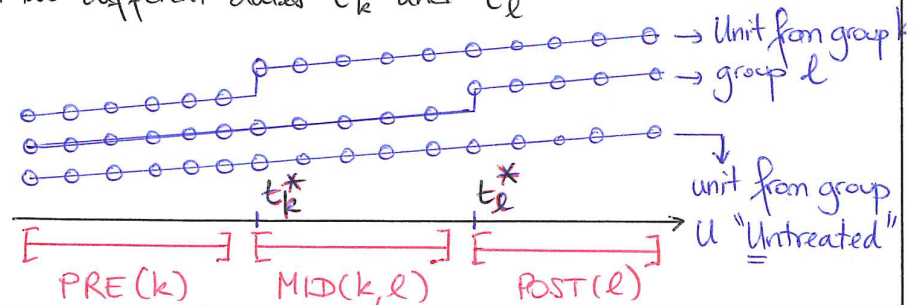
Take e.g.  $E\Delta_{1,3} = 4$  &  $E\Delta_{1,2} = E\Delta_{2,3} = 1$ .

Then  $E\hat{\Delta} = -1/2$ .

↳ Negative weights are of concern only under heterogeneous effects. For example, with  $E\Delta_{1,3} = E\Delta_{1,2} = E\Delta_{2,3} = 1$ ,  $E\hat{\Delta} = 1$ .

x Remark: Since  $\varepsilon_{it} = D_{it} - \bar{D}_{i.} - \bar{D}_{.t} + \bar{D}_{..}$ , negative weights in a staggered design are more likely on early adopters in late time periods.

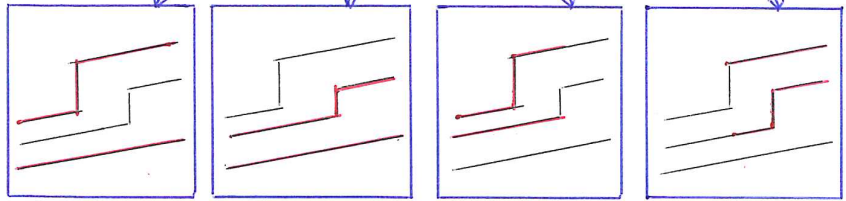
x Remark: Goodman-Bacon (2018) provides helpful intuition why this is happening in a staggered design. He shows that  $\hat{\Delta}$  can be decomposed as a weighted sum of canonical diff-in-diff estimates, where early adopters are mistakenly taken as control units "forbidden comparisons". Specifically, consider  $n$  units whose treatment status turns on at two different dates  $t_k^*$  and  $t_l^*$



Then

$$\hat{\Delta} = s_{ku} \hat{\Delta}_{ku} + s_{lu} \hat{\Delta}_{lu} + s_{kl}^k \hat{\Delta}_{kl}^k + s_{kl}^l \hat{\Delta}_{kl}^l$$

sum to one



Here, the untreated group is correctly used as a control group

Here, the treated group is used as control

$$\hat{\Delta}_{ku} = \left( \bar{Y}_k^{POST(k)} - \bar{Y}_k^{PRE(k)} \right) - \left( \bar{Y}_u^{POST(k)} - \bar{Y}_u^{PRE(k)} \right)$$

$$\hat{\Delta}_{lu} = \left( \bar{Y}_l^{POST(l)} - \bar{Y}_l^{PRE(l)} \right) - \left( \bar{Y}_u^{POST(l)} - \bar{Y}_u^{PRE(l)} \right)$$

↳  $\hat{\Delta}_{ku}$ ;  $\hat{\Delta}_{lu}$  are legitimate did estimators

$$\hat{\Delta}_{kl}^k = \left( \bar{Y}_k^{MID(k,l)} - \bar{Y}_k^{PRE(k)} \right) - \left( \bar{Y}_l^{MID(k,l)} - \bar{Y}_l^{PRE(k)} \right)$$

$$\hat{\Delta}_{kl}^l = \left( \bar{Y}_l^{POST(l)} - \bar{Y}_l^{MID(k,l)} \right) - \left( \bar{Y}_k^{POST(l)} - \bar{Y}_k^{MID(k,l)} \right)$$

↑ acting as treatment group      ↑ acting as control group

⇒ TWFE  $\hat{\Delta}$  is inappropriate in a staggered design.

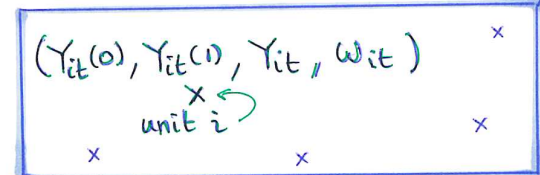
• Weights are non-negative & depend on the subsamples sizes squared & subsample variance. Weights are larger when the two groups are similar in size & when the treatment occurs in the middle of the time window.

III - STRATIFICATION

In Sections I and II, we considered a collection of units (randomized or not), and derived estimators of the ATE/ATT when units are observed over some period of time.

No pre-test + RCT :  $W_{it} = W_i + \{Y_{it}(0), Y_{it}(1)\}$   
 $\forall i, \forall t$

The OLS estimator  $\hat{\Delta}$  of  $\Delta$  in  $Y_{it} = \beta_0 + \Delta W_i + \epsilon_{it}$  is a consistent estimator of  $ATE = \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) - Y_{it}(0))$  (p.10)



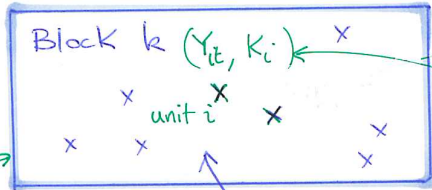
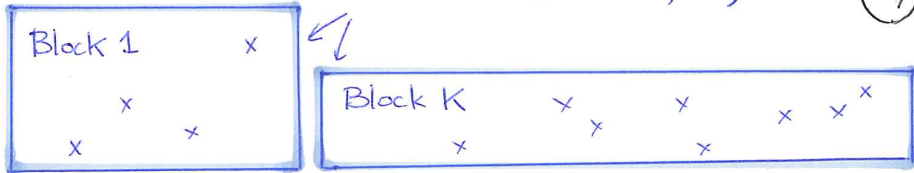
↳ Pre-test + Observational data:  $W_{it}$  are not randomized. The OLS estimator  $\hat{\Delta}$  of  $\Delta$  in  $Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \epsilon_{it}$  is a consistent estimator of  $ATT(T_0, T) = \frac{1}{T-T_0} \sum_{t=T_0+1}^T E(Y_{it}(1) - Y_{it}(0))$  ( $W_i=1$ )

under assumptions (A) Parallel trends & (B) No anticipation (p.18)

Instead of a single collection of units where these results hold, we may consider  $K$  of them (aka strata).



$n$  units in total:  $i = 1, \dots, n$



add stratum id  $K_i = k$ .

$n_k$  units in the  $k$ -th block

In each block  $k$ , assume either an RCT, or that parallel trends + no anticipation holds.

To analyze such data and recover ATE,  $ATT(T_0, T)$ , we may

(i) analyze block by block and aggregate the point estimates

(ii) pull all the data together and solve for  $\Delta$  in  $Y_{it} = \beta_0 + \Delta w_i + \varepsilon_{it}$  (RCT) or  $Y_{it} = \alpha_i + \beta_t + \Delta w_{it} + \varepsilon_{it}$  (DiD) using a proper weighting of the units, to account for potential imbalance across blocks  $\rightarrow$  IPW approach.

III.1. Analysis block by block

An analysis block by block requires a substantial amount of data (treated/control units), so that asymptotic results are well approximated in each block.

In the  $k$ -th stratum, let  $\hat{\Delta}_k$  denote the OLS estimator of  $\Delta_k$  in

$$Y_{it} = \beta_{0k} + \Delta_k w_i + \varepsilon_{it}, \quad \{i \mid K_i = k\}$$

(difference est / RCT)  $\xrightarrow{(n_k \text{ units})}$   $t = 1, \dots, T$

or

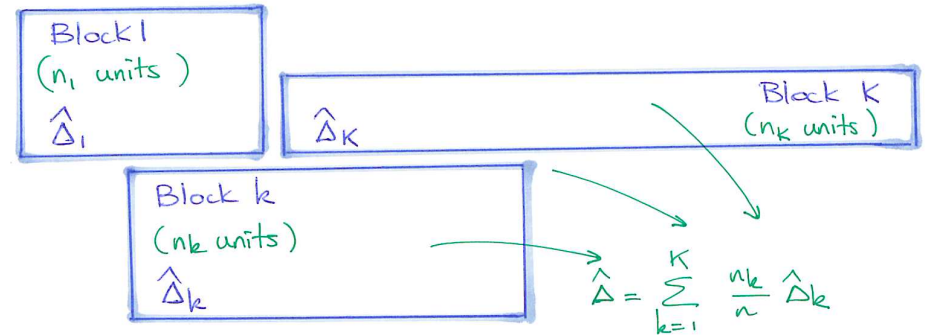
$$Y_{it} = \alpha_{ik} + \beta_{tk} + \Delta_k w_{it} + \varepsilon_{it}, \quad \{i \mid K_i = k\}$$

(diff-in-diff / no anticipation + // trends)  $\xrightarrow{t = 1, \dots, T}$

Let  $n_1 + \dots + n_K = n$ . The aggregate estimator  $\hat{\Delta}_{AGG}$  is defined by

$$\hat{\Delta}_{AGG} = \sum_{k=1}^K \frac{n_k}{n} \hat{\Delta}_k$$

The variance of  $\hat{\Delta}$  can be easily computed using  $\text{var } \hat{\Delta} = \sum_{k=1}^K \left(\frac{n_k}{n}\right)^2 \text{var } \hat{\Delta}_k$ .



Put  $\bullet$   $ATE(k) := \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) - Y_{it}(0) \mid K_i = k)$

$\bullet$   $ATT(k; T_0, T) := \frac{1}{T - T_0} \sum_{t=T_0+1}^T E(Y_{it}(1) - Y_{it}(0) \mid K_i = k, w_i = 1).$

As  $n_{ik}$  ( $=$  # treated units in the  $k$ -th stratum) and  $n_{0k} = n_k - n_{1k} \rightarrow \infty$ , the aggregate estimator

$$\hat{\Delta}_{AGG} \rightarrow \sum_{k=1}^K \pi_k \text{ATE}(k) = \text{ATE} \quad (\text{difference est; RCT})$$

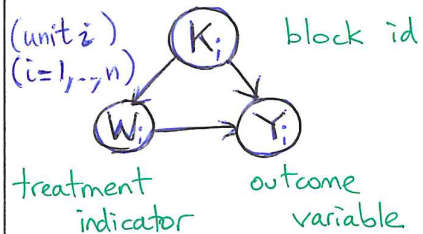
$$\rightarrow \sum_{k=1}^K \pi_k \text{ATT}(k; T_0, T) = \text{ATT}(T_0, T) \quad (\text{did; no anticipation + trends})$$

where  $\frac{n_k}{n} \rightarrow \pi_k$  as  $n_k \rightarrow \infty$ .

### III.2. IPW

We discuss the difference and diff-in-diff estimators separately.

• The difference estimator: in each block  $k=1, \dots, K$ , we perform an RCT. In other words, when pulling the data together, the block / stratum id confounds the treatment effect:



$\Rightarrow$  Consider an IPW approach, weighting unit  $i$  in stratum  $k$  by  $\frac{1}{P(W_i=1 | K_i=k)}$

if  $i$  is in the treatment group, and by  $\frac{1}{1 - P(W_i=1 | K_i=k)}$  if  $i$  is in the control group.

Let  $\hat{\Delta}_{IPW}$  denote the WLS estimate of  $\Delta$  in

$$Y_{it} = \beta_0 + \Delta W_i + \varepsilon_{it}, \quad (i=1, \dots, n; t=1, \dots, T)$$

with weight  $z_i = z(W_i, K_i) = \begin{cases} \frac{1}{P(W_i=1 | K_i=k)} & \text{if } i \in \text{trt} \\ & K_i=k \\ 1 & \text{if } i \in \text{ctrl} \\ & K_i=k \end{cases}$

• Result: The IPW estimator of  $\Delta$  is the weighted difference in means

$$\hat{\Delta}_{IPW} = \frac{\sum_{i|W_i=1,t} z(W_i, K_i) Y_{it}}{\sum_{i|W_i=1,t} z(W_i, K_i)} - \frac{\sum_{i|W_i=0,t} z(W_i, K_i) Y_{it}}{\sum_{i|W_i=0,t} z(W_i, K_i)}$$

[treated units] [control units]

(see proof on the next page)

Then  $\hat{\Delta}_{IPW} \xrightarrow{(n \rightarrow \infty)} \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) - Y_{it}(0)) = \text{ATE}$

Lets consider the treated units in  $\hat{\Delta}_{IPW}$ . The control units are treated similarly. Multiplying the numerator and denominator by  $\frac{1}{nT}$ , we see that

$$\frac{1}{nT} \sum_{i|W_i=1,t} z(W_i, K_i) \xrightarrow{(n \rightarrow \infty)} E\{W z(W, K)\}$$

$$= \sum_k z(1, k) P(W=1, K=k)$$

$$= \sum_k P(K=k) = 1$$

• Similarly,  $\frac{1}{nT} \sum_{i|W_i=0,t} z(W_i, K_i) Y_{it}(1) \xrightarrow{(n \rightarrow \infty)} \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1))$

Proof:  $(\hat{\beta}_0, \hat{\Delta}_{IPW}) = \underset{\beta_0, \Delta}{\operatorname{argmin}} \left\{ \sum_{i=1}^{nT} z(w_i, k_i) (Y_i - \beta_0 - \Delta w_i)^2 \right\}$  (31)

Put  $\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\Delta}_{IPW} \end{pmatrix}$ .

Then  $\hat{\beta} = (X^t Z X)^{-1} X^t Z Y$ ,

where  $Z = \operatorname{diag}\{z(w_i, k_i)\}$   
( $nT \times nT$ )

$X = \begin{pmatrix} 1 & w_1 \\ \vdots & \vdots \\ 1 & w_{nT} \end{pmatrix}$ ,  $Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_{nT} \end{pmatrix}$   
( $nT \times 2$ )      ( $nT \times 1$ )

After calculations,  $X^t X = \begin{pmatrix} \sum z_i & \sum z_i w_i \\ \sum z_i w_i & \sum z_i w_i^2 \end{pmatrix}$

$z_i = z(w_i, k_i) \rightarrow i=1, \dots, nT$

$X^t Z Y = \begin{pmatrix} \sum z_i Y_i \\ \sum z_i w_i Y_i \end{pmatrix}$

$\hat{\beta} = \begin{pmatrix} \frac{\sum z_i Y_i (1 - w_i)}{\sum z_i (1 - w_i)} \\ \frac{(\sum z_i)(\sum z_i w_i Y_i) - (\sum z_i w_i)(\sum z_i Y_i)}{(\sum z_i w_i)(\sum z_i (1 - w_i))} \end{pmatrix}$

$\Rightarrow \hat{\beta}_0 = \frac{\sum_{i|w_i=0} z_i Y_i}{\sum_{i|w_i=0} z_i}$  ;  $\hat{\beta}_0 + \hat{\Delta}_{IPW} = \frac{\sum_{i|w_i=1} z_i Y_i}{\sum_{i|w_i=1} z_i}$  ■

For convenience we denote the  $nT$  observations using a single subscript

x Remark =

$\hat{\Delta}_{AGG} = \sum_{k=1}^K \left( \frac{n_k}{n} \right) \hat{\Delta}_k$  (page 22)

$= \sum_{k=1}^K \frac{n_k}{n} \left( \frac{1}{n_k T} \sum_{t=1}^T \sum_{i|w_i=1} Y_{it} - \frac{1}{n_k T} \sum_{t=1}^T \sum_{i|w_i=0} Y_{it} \right)$

focusing on the first term,

$= \frac{1}{nT} \sum_{t=1}^T \sum_{k=1}^K \sum_{i|w_i=1} \sum_{K_i=k} \left( \frac{n_k}{n_k} \right) Y_{it}$

$= \frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} \sum_{K_i=k} P(w_i=1 | K_i=k) Y_{it}$   
 $= z(w_i, k_i)$

$= \frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} z(w_i, k_i) Y_{it}$

In addition, note that

$\frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} z(w_i, k_i) = \frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} \frac{n_k}{n_k} = 1$   
 $= \sum_{k=1}^K \sum_{i|w_i=1} \sum_{K_i=k} \frac{n_k}{n_k} = \sum_{k=1}^K n_k = n$

Thus  $\hat{\Delta}_{AGG} = \frac{\frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} z(w_i, k_i) Y_{it}}{\frac{1}{nT} \sum_{t=1}^T \sum_{i|w_i=1} z(w_i, k_i)}$  - similar term for control / similar term for control

$= \hat{\Delta}_{IPW}$  .  $\Rightarrow$  Point estimates  $\hat{\Delta}_{AGG}$  and  $\hat{\Delta}_{IPW}$  coincide. ■



x Summary for the difference estimator

- Set-up = •  $k = 1, \dots, K$  blocks
  - In each block, randomize units completely at random:  $\{Y_{it}(0), Y_{it}(1)\} \perp W_i \mid K_i = k$

• Analysis block per block (page 22)  
 $\hat{\Delta}_k$  = OLS estimator of  $\Delta_k$  in  $Y_{it} = \beta_{0k} + \Delta_k W_i + \varepsilon_{it}$  for units  $i$  in the  $k$ -th stratum

$\hat{\Delta}_{AGG} = \sum_{k=1}^K \frac{n_k}{n} \hat{\Delta}_k$  ← Need large  $n_k$  to estimate the variance of  $\hat{\Delta}_k$  accurately

• IPW estimator (page 24) All data pulled together  
 $\hat{\Delta}_{IPW}$  = WLS estimator of  $\Delta$  in  $Y_{it} = \beta_0 + \Delta W_i + \varepsilon_{it}$ , where observation  $Y_{it}$  has weight  $z_i(W_i, K_i)$

$$= \begin{cases} 1/P(W_i=1 \mid K_i=k) & \text{if } i \text{ is treated} \\ 1/P(W_i=0 \mid K_i=k) & \text{if } i \text{ is in control.} \end{cases}$$

Then •  $\hat{\Delta}_{IPW} = \frac{1}{nT} \sum_{t=1}^T \sum_{i:W_i=1} z_i(W_i, K_i) Y_{it} - \frac{1}{nT} \sum_{t=1}^T \sum_{i:W_i=0} z_i(W_i, K_i) Y_{it}$   
 $= \hat{\Delta}_{AGG}$  (page 26)

• Also,  
 $\hat{\Delta}_{AGG} = \hat{\Delta}_{IPW} \xrightarrow{n \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E(Y_{it}(1) - Y_{it}(0)) = ATE$

\*  $n_k$  = # units in  $k$ -th block  
 $w_k$  = # treated units in the  $k$ -th block

In the summary page 27, we see that the analysis block per block requires the number  $n_k$  of units to be large in each block to estimate the variance of  $\hat{\Delta}_k$  (and thus of  $\hat{\Delta}_{AGG}$ ) accurately.

This condition can be relaxed with the IPW estimator: units  $(Y_{it}, K_i)$  are assumed to be drawn at random, and we only require to observe many of them to get an accurate estimate of the variance of  $\hat{\Delta}_{IPW}$  (even when clustering errors over  $t = 1, \dots, T$ ).

⇒ With small blocks, the IPW estimator is preferred over the aggregated estimator.

• The difference-in-differences estimator

We begin by stating a general result on the WLS solution of the diff-in-diff linear model, with general weights:

\* Result: The WLS estimate  $\hat{\beta}_3$  of  $\beta_3$  in  $Y_{it} = \beta_0 + \beta_1 \mathbb{1}(i \in \text{trt}) + \beta_2 \mathbb{1}(t > T_0) + \beta_3 \mathbb{1}(i \in \text{trt}, t > T_0) + \varepsilon_{it}$  with weight  $z_i$  (constant weight  $\forall t$ ) is

$$\hat{\beta}_3 = \left\{ \frac{\sum_{i \in \text{trt}} \sum_{t > T_0+1} z_i Y_{it}}{\sum_{i \in \text{trt}} \sum_{t > T_0+1} z_i} - \frac{\sum_{i \in \text{trt}} \sum_{t \leq T_0} z_i Y_{it}}{\sum_{i \in \text{trt}} \sum_{t \leq T_0} z_i} \right\} \left\{ \frac{\sum_{i \in \text{ctr}} \sum_{t > T_0+1} z_i Y_{it}}{\sum_{i \in \text{ctr}} \sum_{t > T_0+1} z_i} - \frac{\sum_{i \in \text{ctr}} \sum_{t \leq T_0} z_i Y_{it}}{\sum_{i \in \text{ctr}} \sum_{t \leq T_0} z_i} \right\}$$

← treated units (red arrow)  
 ← control units (green arrow)  
 test (blue arrow pointing to  $t > T_0+1$ )  
 pre-test (blue arrow pointing to  $t \leq T_0$ )

x Case I = Randomized treatment

Units are randomized (completely at random) in each block. We use IPW weights to recover the ATE.

$$z_i = \begin{cases} 1 / P(W_i=1 | K_i=k) & \text{if } i \text{ is treated} \\ 1 / P(W_i=0 | K_i=k) & \text{if } i \text{ is in control.} \end{cases}$$

Proceeding as before,

$$\hat{\beta}_3 \xrightarrow{(n \rightarrow \infty)} \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T E Y_{it}(1) - \frac{1}{T_0} \sum_{t=1}^{T_0} E Y_{it}(1) \right\} - \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T E Y_{it}(0) - \frac{1}{T_0} \sum_{t=1}^{T_0} E Y_{it}(0) \right\} = \frac{1}{T-T_0} \left\{ \sum_{t=T_0+1}^T E (Y_{it}(1) - Y_{it}(0)) \right\}$$

since  $\forall t=1, \dots, T_0$ ,  $E Y_{it}(1) = E Y_{it}(0)$  due to randomization (assuming no anticipation)

→ Note that when checking balance of the KPI on the treatment & control groups, one must use weighted averages (using IPW weights) since

$$\begin{aligned} E Y_{it}(1) &= E E Y_{it}(1) | K_i \\ &\approx \sum_{k=1}^K \frac{n_k}{n} \frac{1}{n_k} \sum_{\substack{i | W_i=1 \\ K_i=k}} Y_{it} \\ &= \frac{1}{n} \sum_{i | W_i=1} \left( \frac{n_k}{n_k} \right)^{-1} Y_{it} \quad \blacksquare \end{aligned}$$

x Case II = Observation Study

Assume // trend & no anticipation in each block:

[cond // trend] [two time periods]

$$E(Y_{i1}(0) - Y_{i0}(0) | W_i=0, K_i=k) = E(Y_{i1}(0) - Y_{i0}(0) | W_i=1, K_i=k)$$

[no anticipation]  $Y_{i0}(0) = Y_{i0}(1)$  [in pre-test]

↓ Again, abusing notation ( $Y_{it}(0) = Y_{it}(0,0)$  & test subscript 1)  $Y_{it}(1) = Y_{it}(0,1)$  - full trajectory -

Theorem Abadie (2005) [two time periods]

$$\begin{aligned} ATT &:= E(Y_{i1}(1) - Y_{i1}(0) | W_i=1) \\ &= E(Y_{i1} - Y_{i0} | W_i=1) \\ &\quad - E\{ \alpha(K_i)(Y_{i1} - Y_{i0}) | W_i=0 \} \end{aligned}$$

observed quantities → the ATT is identified.

with  $\alpha(k) = \frac{P(W_i=1 | K_i=k)}{P(W_i=1)} \times \frac{P(W_i=0)}{P(W_i=0 | K_i=k)}$

→ Holds with general vector / real-valued covariates  $K_i$ .

proof =

• First term =  $E(Y_{i1} - Y_{i0} | W_i=1) = E(Y_{i1}(1) - Y_{i0}(1) | W_i=1)$

$$= E(Y_{i1}(1) - Y_{i0}(0) | W_i=1)$$

(no anticipation)

(37)

$$\begin{aligned} \bullet \text{ Second term} &= E\{(Y_{i1} - Y_{i0})\alpha(K_i) | W_i=0\} \\ &= E\{(Y_{i1}(0) - Y_{i0}(0))\alpha(K_i) | W_i=0\} \\ &= E\left[E\left\{\frac{\quad}{\quad} | W_i=0, K_i\right\} | W_i=0\right] \\ &= E\left[\alpha(K_i) \underbrace{E\{(Y_{i1}(0) - Y_{i0}(0)) | W_i=0, K_i\}}_{\psi(K_i)}\right] \\ &= E[\alpha(K_i)\psi(K_i) | W_i=0] \\ &= \sum_k \alpha(k)\psi(k) \underbrace{P(K_i=k | W_i=0)} \end{aligned}$$

definition of  $\alpha(k)$

$$\begin{aligned} &\downarrow \frac{P(W_i=0 | K_i=k)}{P(W_i=0)} P(K_i=k) \\ &= \sum_k \psi(k) \frac{P(W_i=1 | K_i=k)}{P(W_i=1)} P(K_i=k) \\ &= \sum_k \psi(k) P(K_i=k | W_i=1) \end{aligned}$$

Conditional // trend assumption

$$\begin{aligned} &\downarrow = \sum_k E\{(Y_{i1}(0) - Y_{i0}(0)) | W_i=1, K_i=k\} \times P(K_i=k | W_i=1) \\ &= E\left[E\{(Y_{i1}(0) - Y_{i0}(0)) | W_i=1, K_i\} | W_i=1\right] \\ &= E\{Y_{i1}(0) - Y_{i0}(0) | W_i=1\} \quad \blacksquare \end{aligned}$$

In addition, note that

(38)

$$\begin{aligned} E(\alpha(K_i) | W_i=0) &= \sum_k \alpha(k) E(K_i=k | W_i=0) \\ &= \sum_k \frac{P(W_i=0)}{P(W_i=1)} \frac{P(W_i=1 | K_i=k)}{P(W_i=0 | K_i=k)} \\ &\quad \times P(K_i=k | W_i=0) \\ &= \frac{P(K_i=k)}{P(W_i=0)} \\ &= \sum_k \frac{P(W_i=0)}{P(W_i=1)} P(W_i=1 | K_i=k) \frac{P(K_i=k)}{P(W_i=0)} \\ &= \sum_k P(K_i=k | W_i=1) \\ &= 1 \end{aligned}$$

and likewise in the treatment group.

Consequence: The ATT can be estimated using WLS:

$$Y_{it} = \beta_0 + \beta_1 \mathbb{1}(i \text{ is trt}) + \beta_2 \mathbb{1}(t=1) + \beta_3 \mathbb{1}(i \text{ is trt}, t=1) + \varepsilon_{it}$$

(t=0,1  
i=1,...,n)

$$z_i = \begin{cases} \alpha(k) & \text{if } i \text{ in block } k \text{ is in control} \\ 1 & \text{if } i \text{ is treated.} \end{cases}$$

The WLS estimate of  $\hat{\beta}_3$  identifies the ATT.



The result generalizes to multi-periods assuming

(39)

(A) Conditional Parallel Trends

$$\frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i=1, K_i=k) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i=1, K_i=k)$$

$$=$$

$$\frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(0) | W_i=0, K_i=k) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it}(0) | W_i=0, K_i=k)$$

(B) No anticipation

$$Y_{it}(\underbrace{0, \dots, 0}_{T_0}, \underbrace{0, \dots, 0}_{T-T_0}) = Y_{it}(\underbrace{0, \dots, 0}_{T_0}, \underbrace{1, \dots, 1}_{T-T_0}) \quad \forall t=1, \dots, T_0$$

( $\Leftrightarrow Y_{it}(0) = Y_{it}(1)$ ,  $t \leq T_0$ )  $\forall i$  with  $W_i=1$

Then

$$ATT(T_0, T) = \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it}(1) - Y_{it}(0) | W_i=1)$$

P.O.  $\nearrow$

$$= \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}(Y_{it} | W_i=1) - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}(Y_{it} | W_i=1) \right\}$$

observed quantities  $\nearrow$

$$- \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T \mathbb{E}[\alpha(K_i) Y_{it} | W_i=0] - \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{E}[\alpha(K_i) Y_{it} | W_i=0] \right\}$$

& proceed as before to estimate the ATT using weighted least squares:  $Y_{it} = \beta_0 + \beta_1 \mathbb{1}(trt) + \beta_2 \mathbb{1}(test) + \beta_3 \mathbb{1}(trt, test) + \varepsilon_{it}$

weight  $z_i \nearrow$

\* Summary for the did estimator

(40)

- Setup:  $k=1, \dots, K$  blocks
    - In each block, assume parallel trends & no anticipation
- [p. 30  $\rightarrow$  two time periods  
p. 33  $\rightarrow$  multi-periods]

• Analysis block per block (p. 22)

$\hat{\Delta}_k = \text{OLS estimator of } \Delta_k \text{ in } Y_{it} = \alpha_k + \beta t_k + \Delta W_{it} + \varepsilon_{it}$   
( $\forall i | K_i=k$ );  $t=1, \dots, T$

$$\hat{\Delta}_{AGG} = \sum_{k=1}^K \frac{n_k}{n} \hat{\Delta}_k$$

• IPW estimator

$\hat{\Delta}_{IPW} = \text{WLS estimator of } \Delta \text{ in (p. 28)}$

$$Y_{it} = \beta_0 + \beta_1 \mathbb{1}(W_i=1) + \beta_2 \mathbb{1}(t \geq T_0+1) + \Delta \mathbb{1}(W_i=1, t \geq T_0+1) + \varepsilon_{it}$$

with weight  $z_i$ .

(A) If treatment is randomized, (p. 29)

$$z_i = \begin{cases} 1 / \mathbb{P}(W_i=1 | K_i=k) & \text{if } i \text{ is treated (in } k\text{-th stratum)} \\ 1 / \mathbb{P}(W_i=0 | K_i=k) & \text{if } i \text{ is in control} \end{cases}$$

(B) If treatment is not randomized (under // trends & no anticipation)

$$z_i = z(W_i, K_i) = \begin{cases} 1 & \text{if } W_i=1, K_i=k \\ \alpha(k) & \text{if } W_i=0, K_i=k \end{cases} \quad (\text{p. 30})$$

$$\hat{\Delta}_{AGG}, \hat{\Delta}_{IPW} \xrightarrow{n \rightarrow \infty} ATT(T_0, T) [= ATE \text{ if } trt \text{ is randomized}]$$

x Appendix A.1. We show that the OLS estimator of  $\Delta$  in  $Y_{it} = \alpha_i + \beta_t + \Delta W_{it} + \varepsilon_{it}$  is the difference-in-differences  $\hat{\Delta}$  on page 12. We proceed in three steps

- ↳ Step I = [toolbox] the Frisch-Waugh Theorem
- ↳ Step II = One-way Fixed Effects
- ↳ Step III = Two-way Fixed Effects.

• Step I: The FW theorem

[Reg 1]  $Y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$   
 $(n \times 1) \quad (n \times k_1) \quad (n \times k_2) \quad (n \times 1)$

[Reg 2]  $M_1 Y = M_1 X_2 \beta_2 + u$   
 $(n \times n) \quad (n \times 1) \quad (n \times k_2) \quad (n \times 1)$

where  $M_1 = I - X_1 (X_1^t X_1)^{-1} X_1^t$   
 = projection matrix onto the  $\perp$  column space of  $X_1$

Then

- (i) OLS estimates of  $\beta_2$  in [Reg 1] and in [Reg 2] are identical
- (ii) OLS residuals from [Reg 1] and [Reg 2] are identical.

↳ Why is this interesting? When  $k_1$  is large, we can solve numerically a simpler problem if one is interested in estimating  $\beta_2$ .

proof =

[Reg 1] Let  $\hat{\beta}_1, \hat{\beta}_2 = \text{OLS estimates of } \beta_1, \beta_2$

[Reg 2] Put  $\tilde{\beta}_2 = (X_2^t M_1 X_2)^{-1} X_2^t M_1 Y = \text{OLS estimate of } \beta_2$ .  
 (Note:  $M_1$  is idempotent)

Note that

$$Y = P_X Y + (I - P_X) Y = X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 + M_X Y \quad (*)$$

↳ where  $M_X = I - P_X$ ;  $X = [X_1 \ X_2]$   
 $n \times (k_1 + k_2)$

$$M_X = I - X(X^t X)^{-1} X^t$$

↳ pre-multiplying by  $X_2^t M_1$ :

$$X_2^t M_1 Y = X_2^t M_1 X_1 \hat{\beta}_1 + X_2^t M_1 X_2 \hat{\beta}_2 + X_2^t M_1 M_X Y$$

= 0 since  $M_1$  wipes off  $X_1$

taking the transpose:

$$M_X M_1^t X_2 = M_X^t X_2 = 0$$

since

$$P_1 P_X = P_X P_1 = P_1$$

$$\Rightarrow M_X M_1 = (I - P_X)(I - P_1)$$

$$= I - P_X - P_1 + P_X P_1$$

$$= I - P_X = M_X$$

since  $M_X$  wipes off all columns in  $X$

⇒ It follows that

$\hat{\beta}_2 = (X_2^t M_1 X_2)^{-1} X_2^t M_1 Y = \tilde{\beta}_2$ ; which concludes the first part of the theorem.

Pre-multiplying (\*) by  $M_1$  yields

$$M_1 Y = M_1 X_2 \hat{\beta}_2 + M_1 X Y$$

↑  
Regressand in [Reg 2]

↑  
 $= M_1 X_2 \tilde{\beta}_2$  since we just proved that  $\hat{\beta}_2 = \tilde{\beta}_2$

It follows that  $M_1 X Y$  is the vector of residuals in [Reg 2]. But it is immediate to see that  $M_1 X Y$  is also the vector of residuals in [Reg 1]. This concludes the second part of the thm.

• Step II = One-Way Fixed Effects. Baltagi

$$Y_{it} = \alpha_i + \Delta W_{it} + \varepsilon_{it}$$

vector notation

$$Y = Z_\alpha \alpha + W \Delta + \Sigma$$

(nTx1) (nTxn) (nx1) (nTx1) (1x1) (nTx1)

$$Y = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{iT} \\ \vdots \\ Y_{n1} \\ \vdots \end{pmatrix} \begin{matrix} \uparrow \text{unit 1} \\ \downarrow \text{unit n} \end{matrix}$$

$$\alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \text{FE}$$

$$W = \begin{pmatrix} W_{11} \\ \vdots \\ W_{iT} \\ \vdots \end{pmatrix}$$

Put  $\mathbb{1}_T = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^T$ . Then  $Z_\alpha = I_n \otimes \mathbb{1}_T$

Kronecker product

$$Z_\alpha = \begin{bmatrix} 1 & & & 0 \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix}$$

↑ n    ↑ T

We are interested in estimating  $\Delta$  using FW.

$$\Rightarrow P_\alpha := Z_\alpha (Z_\alpha^t Z_\alpha)^{-1} Z_\alpha^t = I_n \otimes \bar{J}_T$$

where  $\bar{J}_T = \frac{1}{T} \mathbb{1}_T \mathbb{1}_T^t$

$$= \begin{bmatrix} \bar{J}_T & & \\ & \ddots & \\ & & \bar{J}_T \end{bmatrix} = \text{time-averages for each unit}$$

$P_\alpha Y$  has elements  $\frac{1}{T} \sum_{t=1}^T Y_{it} = \bar{Y}_i$ .

$\Rightarrow Q_\alpha = I - P_\alpha$  plays the same role as  $M_1$  in FW (Step I) = deviations from the mean

$Q_\alpha Y$  has elements  $Y_{it} - \bar{Y}_i$ .

FW theorem with  $Q_\alpha$  implies that  $\hat{\Delta} = \underbrace{(W^t Q_\alpha W)^{-1}}_{\substack{\text{= a scalar; much smaller} \\ \text{dimension than inverting an } (n+1) \times (n+1) \\ \text{matrix}}} W^t Y$

In addition, we can recover the FE using:

$$Y_{it} = \alpha_i + \Delta W_{it} + \varepsilon_{it}$$

$$\bar{Y}_i = \alpha_i + \Delta \bar{W}_i + \bar{\varepsilon}_i$$

Once  $\hat{\Delta}$  is computed, we get  $\hat{\alpha}_i = \bar{Y}_i - \hat{\Delta} \bar{W}_i$ .

× Remark: Equivalently, we may use the representation  $Y_{it} = \mu + \alpha_i + \Delta W_{it} + \varepsilon_{it}$ ; imposing  $\sum \alpha_i = 0$  for identifiability.

$$\Rightarrow \bar{Y}_i = \mu + \alpha_i + \Delta \bar{W}_i + \bar{\varepsilon}_i$$

$$\bar{Y}_{..} = \mu + \Delta \bar{W}_{..} + \bar{\varepsilon}_{..}$$

$$\Rightarrow \hat{\beta} = \bar{Y}_{..} - \hat{\Delta} \bar{W}_{..}$$

$$\Rightarrow \hat{\alpha}_i = \bar{Y}_i - \hat{\beta} - \hat{\Delta} \bar{W}_i$$



Step III = Two-Way FE

(45)

$$Y_{it} = \alpha_i + \beta_t + \Delta w_{it} + \varepsilon_{it}$$

$$Y = Z_\alpha \alpha + Z_\beta \beta + W \Delta + \varepsilon$$

$nT \times 1$

$$Z_\alpha = I_n \otimes \mathbb{1}_T = \begin{matrix} \begin{matrix} \begin{matrix} \uparrow & \downarrow \\ 1 & 0 \\ 0 & 1 \\ \dots & \dots \\ 0 & 1 \end{matrix} & \begin{matrix} \uparrow & \downarrow \\ T \\ T \\ T \\ T \end{matrix} \\ \left. \begin{matrix} \leftarrow n \rightarrow \end{matrix} \right\} n \text{ times} \end{matrix}$$

$$Z_\beta = \mathbb{1}_n \otimes I_T = \begin{matrix} \begin{matrix} \begin{matrix} \uparrow & \downarrow \\ 1 & \dots & 1 \\ 1 & \dots & 1 \\ \dots & \dots & \dots \\ 1 & \dots & 1 \end{matrix} & \begin{matrix} \uparrow & \downarrow \\ T \\ T \\ T \\ T \end{matrix} \\ \left. \begin{matrix} \leftarrow T \rightarrow \end{matrix} \right\} n \text{ times} \end{matrix}$$

Apply FW with

$$Q := I_n \otimes I_T - I_n \otimes \bar{J}_T - \bar{J}_n \otimes I_T + \bar{J}_n \otimes \bar{J}_T$$

$$nT \times nT = E_n \otimes E_T \text{ where } E_n := I_n - \bar{J}_n$$

$$E_T := I_T - \bar{J}_T$$

One can check that Q wipes out the unit & time FE (and the intercept if there is one). Moreover,  $QY = \tilde{Y}$  where  $\tilde{Y}_{it} := Y_{it} - \bar{Y}_{i.} - \bar{Y}_{.t} + \bar{Y}_{..}$   
= double de-meaned

$$\bar{Y}_{..} = \frac{1}{nT} \sum_{i,t} Y_{it}$$

We get  $\hat{\Delta} = (W^T Q W)^{-1} W^T Q Y$

(46)

A scalar to invert instead of an  $(nT+1) \times (nT+1)$  matrix.

We recover the FE using

$$Y_{it} = \mu + \alpha_i + \beta_t + \Delta w_{it} + \varepsilon_{it}$$

$$\bar{Y}_{i.} = \mu + \alpha_i + \Delta \bar{w}_{i.} + \bar{\varepsilon}_{i.}$$

$$\bar{Y}_{.t} = \mu + \beta_t + \Delta \bar{w}_{.t} + \bar{\varepsilon}_{.t}$$

$$\bar{Y}_{..} = \mu + \Delta \bar{w}_{..} + \bar{\varepsilon}_{..}$$

$$\Rightarrow \begin{cases} \hat{\mu} = \bar{Y}_{..} - \hat{\Delta} \bar{w}_{..} \\ \hat{\alpha}_i = (\bar{Y}_{i.} - \bar{Y}_{..}) - (\bar{w}_{i.} - \bar{w}_{..}) \hat{\Delta} \\ \hat{\beta}_t = (\bar{Y}_{.t} - \bar{Y}_{..}) - (\bar{w}_{.t} - \bar{w}_{..}) \hat{\Delta} \end{cases}$$

x Remark = We can re-express  $\hat{\Delta}$  in a simplified way:

$$\tilde{Y} := QY$$

$$\tilde{W} := QW$$

$$\text{Then } \hat{\Delta} = (W^T Q W)^{-1} W^T Q Y \rightarrow Q^T = Q$$

$$= (W^T Q^T W)^{-1} W^T Q^T Y$$

$$= (\tilde{W}^T \tilde{W})^{-1} \tilde{W}^T Y$$

$$= \frac{\sum_{i,t} \tilde{w}_{it} Y_{it}}{\sum_{i,t} \tilde{w}_{it} w_{it}}$$

$n_1 = \#$  treated units

$$\tilde{w}_{it} = w_{it} - \bar{w}_{i.} - \bar{w}_{.t} + \bar{w}_{..} = \frac{n_1(T-T_0)}{nT}$$

$$= \frac{T-T_0}{T} \mathbb{1}(i \in \text{trt}) \quad L = \frac{n_1}{n} \mathbb{1}(t \geq T_0+1)$$

$$\sum_{i,t} \tilde{w}_{it} Y_{it}$$

$$= \sum_{i,t} w_{i,t} Y_{it} = \sum_{t \geq T_0+1} \sum_{i \in \text{trt}} Y_{it} =: S_{\text{trt, test}}$$

$$- \sum_{i,t} \bar{w}_i \cdot Y_{it} = \frac{T-T_0}{T} \sum_{t \geq T_0+1} \sum_{i \in \text{trt}} Y_{it} = \frac{T-T_0}{T} (S_{\text{trt, pre}} + S_{\text{trt, test}})$$

$$- \sum_{i,t} \bar{w}_{i,t} Y_{it} = \frac{n_1}{n} \sum_{t \geq T_0+1} \sum_{i=1}^n Y_{it} = \frac{n_1}{n} (S_{\text{trt, test}} + S_{\text{cke, test}})$$

$$+ \sum_{i,t} \bar{w}_{i,t} Y_{it} = \frac{n_1(T-T_0)}{nT} (S_{\text{trt, pre}} + S_{\text{trt, test}} + S_{\text{cke, pre}} + S_{\text{cke, test}})$$

likewise, one can show that  $\sum_{i,t} w_{it} \tilde{w}_{it} = \frac{T_0(T-T_0)n_1(n-n_1)}{nT}$

It follows that  $\hat{\Delta} = \text{did}$ , as required.  $\square$

\* Appendix A.2 To see why the model with all coeffs  $\gamma_s$  is overspecified, note that  $Y = Z_\alpha \alpha + Z_\beta \beta + W \gamma + \varepsilon$ , where

$$W = \begin{matrix} \begin{matrix} 1 & \backslash & 1 \\ 1 & \backslash & 1 \\ 0 & & \\ 0 & & \\ 0 & & \end{matrix} & \left. \begin{matrix} \text{trt unit} \\ \text{trt unit} \\ \text{cke unit} \\ \dots \\ \text{cke unit} \end{matrix} \right\} \begin{matrix} n_1 \text{ of them} \\ n - n_1 \text{ of them} \end{matrix} & = \begin{pmatrix} | & & | \\ \omega_1 & \dots & \omega_T \\ | & & | \end{pmatrix} \\ & & \text{where } \omega_1 + \dots + \omega_T = \begin{pmatrix} \vdots \\ 0 \\ \vdots \\ 0 \end{pmatrix} \end{matrix}$$

$\updownarrow T$

$$Z_\alpha = \begin{pmatrix} | & & | \\ z_{\alpha_1} & \dots & z_{\alpha_{n_1}} \\ | & & | \end{pmatrix}$$

$\updownarrow T$

$$= z_{\alpha_1} + \dots + z_{\alpha_{n_1}}$$

$\updownarrow T$

= columns of the  $n_1$  treated units in  $Z_\alpha$

$Q(\omega_1 + \dots + \omega_T) = Q(z_{\alpha_1} + \dots + z_{\alpha_{n_1}}) = 0 \Rightarrow W^T Q W$  is not invertible.