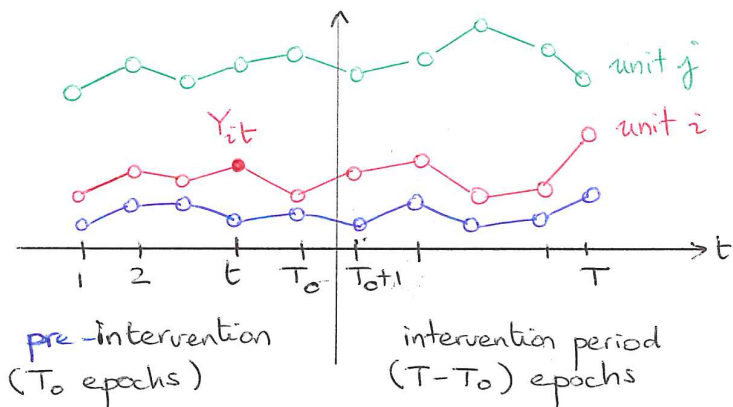


CI = SYNTHETIC CONTROLS

Consider a balanced panel data consisting of n units observed over T time periods. A subset of n_t units receive a treatment at time $T_0+1 \leq T$; while the remaining n_c units are never treated. Let W_{it} denote the treatment status of unit i at time t . Assuming a block treatment assignment, $W_{it} = \mathbb{1}(\{i \geq n_c+1; t \geq T_0+1\})$

Put $Y_{it}(0)$ = control potential outcome of unit i @ time t
 $Y_{it}(1)$ = treatment P.O.

$$Y_{it} = W_{it}Y_{it}(1) + (1-W_{it})Y_{it}(0)$$



We are interested in estimating the average treatment effect on the treated group over the intervention period

$$\begin{cases} \text{ATT}(T_0, T) = \frac{1}{T-T_0} \sum_{t=T_0+1}^T \text{ATT}_t \\ \text{ATT}_t = \mathbb{E}[Y_{it}(1) - Y_{it}(0) \mid W_{it} = 1] \end{cases}$$

In CI = PANEL DATA METHODS, we derived (2)
 a consistent estimator of $\text{ATT}(T_0, T)$ under the parallel trend (A) and no anticipation (B) assumptions (p15/16):

$$\begin{aligned} \hat{\Delta} &= \underset{\Delta, \alpha, \beta}{\text{argmin}} \sum_{i=1}^n \sum_{t=1}^T (Y_{it} - \alpha_i - \beta_t - \Delta W_{it})^2 \\ &= \left\{ \frac{1}{n_t(T-T_0)} \sum_{\substack{i \geq n_c+1 \\ t \geq T_0+1}} Y_{it} - \frac{1}{n_c T_0} \sum_{\substack{i \geq n_c+1 \\ t \leq T_0}} Y_{it} \right\} \\ &\quad - \left\{ \frac{1}{n_c(T-T_0)} \sum_{\substack{i \leq n_c \\ t \geq T_0+1}} Y_{it} - \frac{1}{n_c T_0} \sum_{\substack{i \leq n_c \\ t \leq T_0}} Y_{it} \right\} \\ &= \text{difference-in-differences estimator} \end{aligned}$$

→ $\text{ATT}(T_0, T)$ as $n_t, n_c \rightarrow \infty$ under (A) & (B).

Guarantees (i.e. consistency, bias, asymptotic normality...) of estimators of $\text{ATT}(T_0, T)$ are commonly derived under the assumption of a data generating process. Many panel data methods assume that

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

← additive noise

↑
model component

See Chernozhukov et al (2019) for a discussion. The DID approach assumes that each unit i and each time period have a different offset: $\lambda_{it} = \alpha_i + \beta_t$. When the underlying model departs from this simple additive structure, the parallel trend may break down and

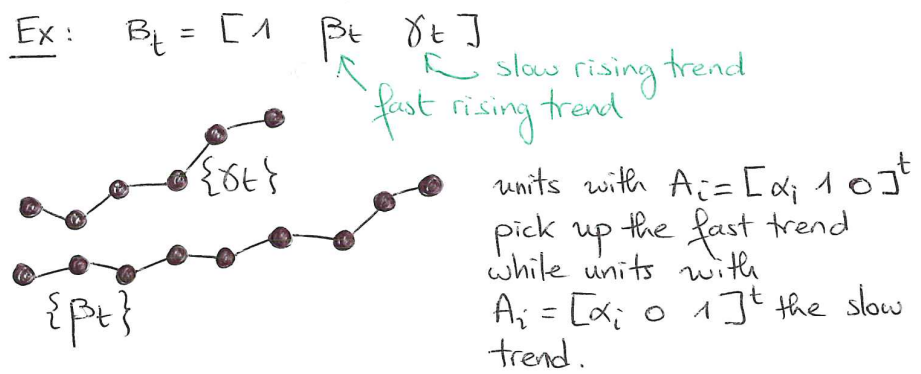
the DID estimator of $ATT(T_0, T)$ may be biased. (3)

A natural generalization is $l_{it} = A_i^t B_t$ for

$A_i, B_t \in \mathbb{R}^k$ [aka FACTOR MODELS]

latent time varying factors
 unknown unit factor loadings

Ex: $A_i = [\alpha_i \ 1]^t$
 $B_t = [1 \ \beta_t]^t$ recovers $l_{it} = \alpha_i + \beta_t$



The $(n \times T)$ matrix of potential outcomes can be modeled as $\underline{Y}(0) = \underline{L} + \underline{\varepsilon}$ $E(\underline{\varepsilon} | \underline{L}) = 0$

$(n \times T)$ $(n \times T)$ $(n \times T)$

for a low rank matrix \underline{L} . Some authors have proposed direct estimation of \underline{L} to recover the ATT, see for example Athey, Bayati, Daudchenko, Imbens, Khosravi (2017) via nuclear norm minimization. Synthetic Control methods address confounding bias without explicitly estimating $\underline{L} \Rightarrow$ indirect approach. SC weight the control units to rebalance the treatment & control groups and create

parallel trends. The motivation being that balancing the treatment & control pre-intervention outcomes should also balance the latent factors A_i . The first section introduces the original SC approach and some of its variants. The next two sections detail two popular generalizations: Synthetic DID and Augmented SC. We conclude with a short discussion on micro vs aggregated data. (4)

I. CANONICAL SC

Recall that $ATT(T_0, T) = \frac{1}{(T - T_0)} \sum_{t=T_0+1}^T ATT_t$ with

$ATT_t = E[Y_{it}(1) | w_{it} = 1] - E[Y_{it}(0) | w_{it} = 1]$

↓ Estimated using ↓ Estimated using SC
 $\bar{Y}_{nt} := \frac{1}{n_t} \sum_{i \in trt} Y_{it}$ Notation: $\hat{Y}_{nt}(0)$

The SC estimator of \hat{ATT}_t is

$\hat{ATT}_t = \bar{Y}_{nt} - \hat{Y}_{nt}(0)$

↑
 The problem reduces to a single treated unit [& $n_c \geq 1$ controls]

Without loss of generality, denote by Y_{nt} the time series observations of the treated unit [clear from context if it corresponds to a sample mean or not].

Observations are $\{Y_{it}\}$ for $i=1, \dots, n$ (5)
 $t=1, \dots, T$

with control units denoted with index $i=1, \dots, n-1$.

Introducing matrix notation:

$$\begin{array}{c} n_c \updownarrow \\ 1 \updownarrow \end{array} \left(\begin{array}{c|cc} & 1 & 1 \\ \hline \underline{Y}_0 & \underline{Y}_{T_0+1} & \underline{Y}_T \\ \hline -\underline{Y}_n^t & ? & ? \end{array} \right)$$

$\xleftarrow{T_0} \quad \xleftarrow{T-T_0}$

SC performs a vertical regression, where each control unit receives a non-negative weight $\hat{w}_i \geq 0$ computed such that the weighted sum of the pre-intervention control observations matches the treatment:

$$Y_{nt} \approx \sum_{i=1}^{n_c} \hat{w}_i Y_{it} \quad \forall t=1, \dots, T_0$$

"the Synthetic Control"

Exact Matching
may not be feasible

There are a multitude of criteria to fit the weights.

Abadie (2003, 2010) [simplex regression]

$$\hat{w} = \underset{\substack{w \geq 0 \\ w^t \mathbf{1} = 1}}{\operatorname{argmin}} \quad \|\underline{Y}_n - \underline{Y}_0^t w\|_2^2 + \lambda \|w\|_2^2$$

$$\hat{Y}_{nt}(0) := \underline{Y}_t^t \hat{w} \quad t \geq T_0+1$$

Remarks (i) Abadie (2003, 2010) do not consider the ridge penalty & a weighted OLS. (6)

(ii) The constraints $w \geq 0$ & $w^t \mathbf{1} = 1$ ensures a sparse solution: many control units receive a 0 weight.

A popular alternative to simplex regression is the elastic net, see e.g. Daudchenko & Imbens (2017)

Daudchenko & Imbens (2017) [elastic net]

$$\hat{w} = \underset{w}{\operatorname{argmin}} \quad \|\underline{Y}_n - \underline{Y}_0^t w\|_2^2 + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$$

$$\hat{Y}_{nt} := \underline{Y}_t^t \hat{w} \quad t \geq T_0+1$$

↑ May add an intercept

* Guarantees under the latent factor model.

$$\forall t \quad Y_{it}(0) = A_i^t B_t + \varepsilon_{it}$$

$$\forall i \quad \Rightarrow Y_{nt} - \hat{Y}_{nt}(0) = Y_{nt} - \sum_{i=1}^{n_c} \hat{w}_i Y_{it}, \quad t \geq T_0+1$$

$$= \left(A_n - \sum_{i=1}^{n_c} \hat{w}_i A_i \right)^t B_t$$

Under some conditions, balancing pre-intervention outcomes will balance the factor loadings

Put $X_i = (Y_{i1}, \dots, Y_{iT_0})^t$. Then:

$$\left(A_n - \sum_{i=1}^{n_c} \hat{w}_i A_i \right) = (B^t B)^{-1} B^t \left(X_n - \sum_{i=1}^{n_c} \hat{w}_i X_i \right) \quad (7)$$

$$- (B^t B)^{-1} B^t \left(\varepsilon_{n,1:T_0} - \sum_{i=1}^{n_c} \hat{w}_i \varepsilon_{i,1:T_0} \right)$$

$$B = \begin{pmatrix} -B_1^t \\ \vdots \\ -B_{T_0}^t \end{pmatrix} \in \mathbb{R}^{T_0 \times k}$$

$$\varepsilon_{i,1:T_0} = (\varepsilon_{i1}, \dots, \varepsilon_{iT_0})^t$$

Under exact matching and some regularity conditions (e.g. $B^t B$ should be invertible), Abadie, Diamond and Hainmuller (2010) provide a bound for $Y_{nt} - \hat{Y}_{nt}(0)$, and show that the bias can be made arbitrarily small as $T_0 \rightarrow \infty$. However, exact matching on pre-intervention outcomes is not always feasible & increasing T_0 does not help in reducing the bias.

The SC estimator of $ATT(T_0, T)$ is

$$\hat{ATT}(T_0, T) = \frac{1}{T-T_0} \sum_{t=T_0+1}^T (Y_{nt} - \hat{Y}_{t(0)})$$

$$= \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{nt} \right\} \leftarrow \text{trt} - \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T \sum_{i=1}^{n_c} \hat{w}_i Y_{it} \right\} \leftarrow \text{cte}$$

= "weighted" difference estimator

The SC estimator is "just" a difference estimator, with proper weighting of the control units. Assuming $w_i = \frac{1}{n_c} \forall i \in \mathcal{C}$

recover the difference estimator [but randomization is required to recover / identify $ATT(T_0, T) = ATE$]. It turns out that the SC point estimate can be computed as the weighted least square estimate in a one-way fixed effect model:

Theorem.

$$\hat{\Delta} = \underset{\beta, \Delta}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t=T_0+1}^T \hat{w}_i (Y_{it} - \beta_t - \Delta W_{it})^2$$

Non-negative, sum to 1 time FE.

$$= \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{nt} - \frac{1}{T-T_0} \sum_{t=T_0+1}^T \sum_{i=1}^{n_c} \hat{w}_i Y_{it}$$

$$= \text{SC estimator of } ATT(T_0, T).$$

[See Appendix A for a proof]

[$w_i = \frac{1}{n_c}$ represents the diff estimator as the OLS solution of a one-way FE model]

The result above shows that SC estimators omit unit FEs. Generalizations with both unit & time FEs yields Synthetic DID (SDID) type of estimators:

Generalization

$$\hat{\Delta} = \underset{\alpha, \beta, \Delta}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t=1}^T \hat{w}_i (Y_{it} - \alpha_i - \beta_t - \Delta W_{it})^2$$

$$= \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{nt} - \frac{1}{T_0} \sum_{t=1}^{T_0} Y_{nt} \right\} \leftarrow \text{trt} \quad (9)$$

$$- \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T \sum_{i=1}^{n_c} \hat{w}_i Y_{it} - \frac{1}{T_0} \sum_{t=1}^{T_0} \sum_{i=1}^{n_c} \hat{w}_i Y_{it} \right\} \leftarrow \text{ctl}$$

= difference-in-differences estimator of $ATT(T_0, T)$ where each control unit receives a weight w_i .

II. SYNTHETIC DID

Arkhangelsky et al (2021) go one step further and introduce weights λ_t to balance pre-intervention time periods with intervention ones:

$$(\hat{\lambda}_0, \hat{\lambda}) = \underset{\substack{\lambda_0, \lambda \\ \mathbb{R} \quad \Lambda}}{\operatorname{argmin}} \sum_{i=1}^{n_c} \left(\lambda_0 + \sum_{t=1}^{T_0} \lambda_t Y_{it} - \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{it} \right)^2$$

\uparrow Horizontal Regression over the control units

where $\Lambda = \left\{ \lambda \in \mathbb{R}_+^T \mid \sum_{t=1}^{T_0} \lambda_t = 1, \lambda_t = \frac{1}{T-T_0}, t \geq T_0+1 \right\}$

\uparrow Ensures a sparse solution.

\Rightarrow Weights λ_t complement weights w_i .

Horizontal (HZ): select in pre-intervention times t similar to the intervention period.

Vertical (VT): select in the control pool units that are similar to the treated unit.

Both unit & time weights are used in a TWFE regression model to estimate the $ATT(T_0, T)$:

Synthetic DID

$$\hat{\Delta} = \underset{\alpha, \beta, \Delta}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t=1}^T \hat{w}_i \hat{\lambda}_t (Y_{it} - \alpha_i - \beta_t - \Delta w_{it})^2$$

$$= \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{nt} - \frac{1}{T_0} \sum_{t=1}^{T_0} \hat{\lambda}_t Y_{nt} \right\} \leftarrow \text{trt}$$

$$- \left\{ \frac{1}{T-T_0} \sum_{t=T_0+1}^T \sum_{i=1}^{n_c} \hat{w}_i Y_{it} - \sum_{t=1}^{T_0} \sum_{i=1}^{n_c} \hat{w}_i \hat{\lambda}_t Y_{it} \right\} \leftarrow \text{ctl}$$

Arkhangelsky et al (2021) fit unit weights using a simplex ridge regression + intercept and time weights using a simplex regression with no penalty + intercept.

The SDID estimator of $ATT(T_0, T)$ takes the form

$$\widehat{ATT}(T_0, T) = \frac{1}{T-T_0} \sum_{t=T_0+1}^T \widehat{ATT}_t$$

$$\widehat{ATT}_t = Y_{nt} - \widehat{Y}_{nt}(0), \quad t \geq T_0+1$$

$$\widehat{Y}_{nt}(0) = \langle \hat{\lambda}, \underline{Y}_n \rangle + \langle \hat{w}, \underline{Y}_t \rangle - \langle \hat{w}, \underline{Y}_0 \hat{\lambda} \rangle$$

The authors provide asymptotic guarantees that $\widehat{ATT}(T_0, T)$ identify $ATT(T_0, T)$ under general regularity conditions & a latent factor model. Assuming a large panel with $n_c \rightarrow \infty, T_0 \rightarrow \infty, n_t(T-T_0) \rightarrow \infty$, they prove

that $\widehat{ATT}(T_0, T)$ is asymptotically normally distributed, $\textcircled{11}$
 centered around $ATT(T_0, T)$, with variance $O\left(\frac{1}{n_c(T-T_0)}\right)$

• Remark: HZ vs VT regression

SDID makes use of information contained in HZ and VT regressions to estimate the ATT. There are cases, however, where both HZ & VT yield the same point estimates, as noted in Shen, Ding, Sekhan & Yu (2023).

Recall the notation $\begin{matrix} \uparrow \\ n_c \\ \downarrow \\ 1 \end{matrix} \left(\begin{array}{c|c} Y_0 & Y_T \\ \hline -Y_n^t & ? \end{array} \right) \begin{matrix} \leftarrow T_0 \\ \rightarrow 1 \end{matrix}$

HZ: $\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \| Y_T - Y_0 \alpha \|_2 \quad \hat{Y}_{NT}^{HZ}(0) = Y_n^t \hat{\alpha}$

VT: $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \| Y_n^t - Y_0^t \beta \|_2 \quad \hat{Y}_{NT}^{VT}(0) = Y_T^t \hat{\beta}$

If columns are linearly independent (rank $Y_0 = T_0 < n_c$)
 then $\hat{\alpha} = (Y_0^t Y_0)^{-1} Y_0^t Y_T$

If rows are linearly independent (rank $Y_0 = n_c < T_0$)
 then $\hat{\beta} = (Y_0 Y_0^t)^{-1} Y_0 Y_n^t$

Assuming more generally that rank $Y_0 = R \leq \min(n_c, T_0)$,
 and denoting Y_0^+ the pseudo inverse of Y_0 [some
 properties: $(Y_0^+)^t$ is the pseudo inverse of Y_0^t since
 $Y_0 = U \Sigma V^t = \sum_{l=1}^R \sigma_l u_l v_l^t$ (SVD decomposition) and

$Y_0^+ = \sum_{l=1}^R \frac{1}{\sigma_l} v_l u_l^t = V \Sigma^{-1} U^t$, with $\textcircled{12}$
 $U \in \mathbb{R}^{n_c \times R}$, $V \in \mathbb{R}^{T_0 \times R}$, $\Sigma \in \mathbb{R}^{R \times R} = \operatorname{diag}(\sigma_1, \dots, \sigma_R)$.

Also, $Y_0^+ Y_0 Y_0^+ = Y_0^+$

Then $\hat{\alpha} = Y_0^+ Y_T \Rightarrow \hat{Y}_{NT}^{HZ}(0) = Y_n^t Y_0^+ Y_T$
 $\hat{\beta} = (Y_0^+)^t Y_n^t \Rightarrow \hat{Y}_{NT}^{VT}(0) = Y_T^t (Y_0^+)^t Y_n^t$
 and we see that $\hat{Y}_{NT}^{HZ}(0) = \hat{Y}_{NT}^{VT}(0)$.

Symmetry in HZ and VT regression under l_2 -norm minimization.

In addition, $\hat{\beta}^t Y_0 \hat{\alpha} = Y_n^t Y_0^+ Y_0 Y_0^+ Y_T = Y_n^t Y_0^+ Y_T = \hat{Y}_{NT}^{HZ}(0) = \hat{Y}_{NT}^{VT}(0)$.

So that the LS estimate $\langle \hat{\alpha}, Y_0^t \hat{\beta} \rangle$ is expressed in terms of both $\hat{\alpha}$ and $\hat{\beta}$.

III - AUGMENTED SC (ASC)

ASC starts with the canonical SC estimator of the weights & makes use of an outcome model to estimate & correct for the bias due to imperfect pre-intervention fit:

$\hat{Y}_{nt}(0) = \hat{M}_{nt}(0) + \sum_{i=1}^{n_c} \hat{w}_i (Y_{it} - \hat{M}_{it}(0))$
 ↑ outcome model estimate ↑ SC weight bias from the outcome model

- Similar to an AIPW estimator
- ASC belongs to the class of Doubly Robust estimators.

(13)

Ben-Michael et al (2021) discuss in depth the properties of this estimator for a class of linear outcome models fitted using Ridge Regression:

$$\hat{M}_{it}(0) = \sum_{t=1}^{T_0} \hat{\rho}_t Y_{it}$$

$i=1, \dots, n$ \leftarrow pre-intervention outcomes only

where

$$(\hat{\rho}_0, \hat{\rho}) = \underset{\substack{\hat{\rho}_0 \in \mathbb{R}, \hat{\rho} \in \mathbb{R}^{T_0}}} {\operatorname{argmin}} \sum_{i=1}^{n_c} \left(\sum_{t=1}^{T_0} \rho_t Y_{it} - \frac{1}{T-T_0} \sum_{t=T_0+1}^T Y_{it} \right)^2 + \lambda \|\rho\|_2^2$$

\leftarrow (+ ρ_0) \leftarrow Average outcome of all unit i in the intervention period

\leftarrow With or without intercept.

Under a linear outcome model, the ASC estimator of $Y_{nt}(0)$, $t \geq T_0+1$, reduces to:

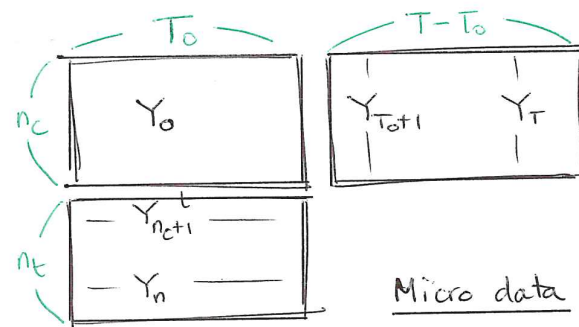
$$\hat{Y}_{nt}(0) = \langle \hat{\rho}, Y_n \rangle + \langle \hat{\omega}, Y_t \rangle - \langle \hat{\omega}, Y_0 \hat{\rho} \rangle$$

\leftarrow Same expression as the SDID estimator. The weights ($\hat{\rho}$ or $\hat{\lambda}$) are not fitted the same way (RR vs simplex regression).

*Remark: Setting $\hat{\rho} \leftarrow \frac{1}{T_0} \mathbb{1}$ (resp. $\hat{\lambda}$ for SDID) and $\hat{\omega} \leftarrow \frac{1}{n_c} \mathbb{1}$ yields the DID estimator. (see also page 10)

Take Away: with randomization.
 SC is to diff
 what ASC/SDID is to diff-in-diff.

IV - MICRO VS AGGREGATED DATA.



($n_c + n_t = n$)

Suppose that the data is clustered

- individuals grouped by city / country
- articles grouped by type (jeans / shoes / ...)

We may consider aggregating units for each cluster before fitting a SC or DID.

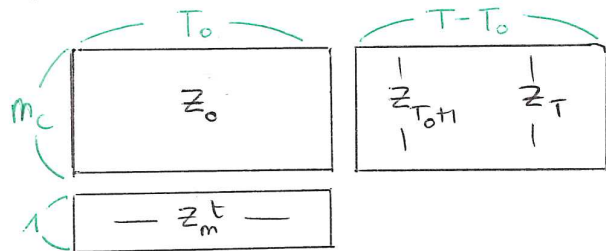
Consider a single treatment cluster & m_c control clusters.

Put
$$z_{jt} = \frac{1}{n_j} \sum_{i \in [n_j]} Y_{it} \quad t=1, \dots, T \quad (15)$$

$$j=1, \dots, m_c+1$$

where $[n_j] = \{ \text{units } i=1, \dots, n \text{ belonging to cluster } j \}$
 $n_j = \# \text{ units in } [n_j]$.

• Aggregated data.



$(m = m_c + 1)$

We are still interested in estimating a unit-level ATT:

$$ATT_t = \mathbb{E} [Y_{it}(1) - Y_{it}(0) | \omega_{it} = 1]$$

$$\downarrow$$

$$\widehat{ATT}_t = z_{mt} - \widehat{z}_{mt}(0)$$

where

$$[\text{ASC/SDID}] \widehat{z}_{mt}(0) = \langle \widehat{\lambda}, z_m \rangle + \langle \widehat{\omega}, z_t \rangle - \langle \widehat{\omega}, z_0 \widehat{\lambda} \rangle$$

$$\widehat{ATT}_t = \left\{ z_{mt} - \sum_{t'=1}^{T_0} \widehat{\lambda}_{t'} z_{mt'} \right\} \leftarrow \text{trt group}$$

$$- \left\{ \sum_{j=1}^{m_c} \widehat{\omega}_j z_{jt} - \sum_{j=1}^{m_c} \sum_{t'=1}^{T_0} \widehat{\omega}_j \widehat{\lambda}_{t'} z_{jt'} \right\}$$

$\leftarrow \text{control group.}$

We recover DID with $\widehat{\lambda} \leftarrow \frac{1}{T_0} \mathbb{1} \quad (16)$

$$\widehat{\omega} \leftarrow \left(\frac{n_1}{n_c}, \dots, \frac{n_{m_c}}{n_c} \right)^t$$

since

$$\widehat{ATT}_t = \left\{ \frac{1}{n_t} \sum_{i=n_c+1}^n Y_{it} - \frac{1}{T_0} \sum_{t'=1}^{T_0} \left(\frac{1}{n_t} \sum_{i=n_c+1}^n Y_{it'} \right) \right\}$$

$$- \left\{ \sum_{j=1}^{m_c} \frac{n_j}{n_c} \left(\frac{1}{n_j} \sum_{i \in [n_j]} Y_{it} \right) - \frac{1}{T_0} \sum_{j=1}^{m_c} \frac{n_j}{n_c} \sum_{t'=1}^{T_0} \left(\frac{1}{n_j} \sum_{i \in [n_j]} Y_{it'} \right) \right\}$$

$\begin{matrix} z_{mt} & & z_{mt'} \\ \leftarrow & & \leftarrow \\ z_{jt} & & z_{jt'} \end{matrix}$

$$= \left\{ \frac{1}{n_t} \sum_{i \in \text{trt}} Y_{it} - \frac{1}{n_t T_0} \sum_{t'=1}^{T_0} \sum_{i \in \text{trt}} Y_{it'} \right\}$$

$$- \left\{ \frac{1}{n_c} \sum_{i \in \text{ctl}} Y_{it} - \frac{1}{n_c T_0} \sum_{t'=1}^{T_0} \sum_{i \in \text{ctl}} Y_{it'} \right\}.$$

$$\& \frac{1}{T-T_0} \sum_{t=T_0+1}^T \widehat{ATT}_t = \text{DID}.$$

\Rightarrow No loss of generality in aggregating data per cluster.

Appendix A : Proof of the Theorem page 8 .

Model $Y_{it} = \beta_t + \Delta w_{it} + \varepsilon_{it}$ in matrix notation :

$$\underline{X} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 1 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \\ 0 & 0 & \dots & 1 \end{bmatrix} \begin{matrix} \left. \begin{matrix} 0 \\ \vdots \\ 0 \end{matrix} \right) n \\ \left. \begin{matrix} 1 \\ 0 \\ \vdots \\ 0 \end{matrix} \right) n \\ \left. \begin{matrix} \vdots \\ 0 \\ \vdots \\ 0 \end{matrix} \right) n \end{matrix} \equiv \begin{bmatrix} X_1 & x \\ X_2 & x \\ \dots & \dots \\ X_{T-T_0} & x \end{bmatrix}$$

$\underline{X} \in \mathbb{R}^{n(T-T_0) \times (T-T_0+1)}$

where $\begin{matrix} t \\ \downarrow \\ 1 \\ \downarrow \\ 1 \end{matrix}$

$X_t \equiv \begin{bmatrix} 0 & 1 & 0 \\ & & 1 \end{bmatrix} \in \mathbb{R}^{n \times (T-T_0)}$

$x = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n$

$\underline{\beta} = (\beta_{T_0+1}, \dots, \beta_T, \Delta)^t \in \mathbb{R}^{T-T_0+1}$

$\underline{Y} = (\underbrace{Y_{1,T_0+1}, \dots, Y_{n,T_0+1}}_n \mid \dots \mid \underbrace{Y_{1,T}, \dots, Y_{n,T}}_n)^t \in \mathbb{R}^{n(T-T_0)}$

$\underline{Y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$

\uparrow defined similarly

Weight Matrix

$\underline{W} = \begin{bmatrix} \underline{w} & 0 \\ 0 & \underline{w} \\ & & & 0 \\ & & & & \underline{w} \end{bmatrix} \in \mathbb{R}^{n(T-T_0) \times n(T-T_0)}$

$\underline{w} = \begin{pmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & & w_n \end{pmatrix} \in \mathbb{R}^{n \times n}$

with $\sum_{i=1}^n w_i = 1 ; w_n = 1$.

Then :

$\underline{X}^t \underline{W} \underline{X} = \begin{pmatrix} 2 \mathbb{I}_{(T-T_0) \times (T-T_0)} & \mathbb{1}_{T-T_0} \\ \mathbb{1}_{T-T_0}^t & T-T_0 \end{pmatrix}$

$(\underline{X}^t \underline{W} \underline{X})^{-1} = \frac{1}{T-T_0} \begin{pmatrix} \frac{1}{2} [(T-T_0) \mathbb{I}_{T-T_0} + \mathbb{J}_{T-T_0}] & -\mathbb{1}_{T-T_0} \\ -\mathbb{1}_{T-T_0}^t & 2 \end{pmatrix}$

$\underline{X}^t \underline{W} \underline{Y} = \left(\sum_{i=1}^n w_i Y_{i,T_0+1}, \dots, \sum_{i=1}^n w_i Y_{i,T}, \sum_{t=T_0+1}^T w_n Y_{nt} \right)^t$

The result follows noticing that $\hat{\Delta}$ is the last element in the vector $(\underline{X}^t \underline{W} \underline{X})^{-1} \underline{X}^t \underline{W} \underline{Y}$.