

# Session IV

## TA Session: DiD and Synth control

Evaluating public policies

**Arthur Heim** (PSE & Cnaf)

March 19, 2023\*

# Outline

## 1 Introduction

What we have seen so far

The kind of question we ask

## 2 Synthetic controls in brief

## 3 First application: Basque terrorism from Abadie and Gardeazabal (2003)

## 4 Cool implementations of Synth controls

## 5 Application: Prison and black male incarceration

# Introduction

## What we have seen so far

- 1 We used potential outcome notations to define causal impacts, and identification strategies
- 2 With randomization, assignment to treatment is independent from potential outcomes hence simple differences identify causal relationships and we may use regressions (not only OLS) to estimate the parameter of interests.
- 3 With panel data or repeated cross sections, we may identify the average treatment effect on the treated if we assume parallel trend i.e. in the absence of treatment, treated and untreated units would have followed the same path.
- 4 When there is only one group that's treated, usual models (regressions with leads and lags in particular) work fine but a strand of recent papers showed that in the multiple group, multiple period setting, TWFE are strongly biased.
- 5 New methods allow either to fix the problem depending on the setting or estimate different parameters

Now, we stay with this parallel trend intuition and move to macro data.

# The kind of question we ask

## What's effect of vaccinal mandate on covid ?

- You remember the "pass vaccinal" ; we would like to know if it works, right ?
- We wouldn't be too happy with using past data for forecasting, right ?
- We wouldn't consider a single country (e.g. Germany) as a counterfactual because they are different and probably on different paths etc.
- How do we find a counterfactual when the whole country is affected ?
- Cross country variations ? Which countr(y)ies should we choose ?
- **What could we do ?**

Idea: use weighted average of several countries to construct a synthetic France : Synthetic controls

# Outline

- ① Introduction
  
- ② Synthetic controls in brief
  - What are synthetic controls
  - Implementation
  
- ③ First application: Basque terrorism from Abadie and Gardeazabal (2003)
  
- ④ Cool implementations of Synth controls
  
- ⑤ Application: Prison and black male incarceration

# Synthetic controls in brief

## What are synthetic controls

- **Intuition :**
  - With more aggregated data, usually no clear comparison “group” to assess the impact of a policy change, even with DiD.
  - Using a weighted mixture of other regions may provide a better counterfactual than a single one, hence the synthetic control
- **Main advantages of the methods compared with regressions :**
  - No need for extrapolation – instead use interpolation
  - Focus on data before the results are known to define counterfactual (researcher less incline to harking)
  - Weights make explicit what units contribute to the counterfactual and by how much
  - Bridge a gap between qualitative and quantitative data (case study)
  - Very practical discussion in Abadie, Diamond, and Hainmueller (2011) about the *Synth* package that implement this method in R.
  - **State of the art nowadays:** Alberto Abadie. 2021. “Using Synthetic Controls: Feasibility, Data Requirements, and Methodological Aspects.” *Journal of Economic Literature* 59, no. 2 (June 1, 2021): 391–425

# Synthetic controls in brief

## What are synthetic controls

- Define an outcome  $Y_{it}$  with  $I + 1$  aggregated units of interest and  $i = 1$  is the treated region.
- Each region is observed over  $T$  periods, treatment occurs at time  $t_0$ .
- Like in Rubin's causal model, treatment effect is defined at time  $t$  by

$$Y_{it}(1) - Y_{it}(0)$$

where, here, the treated unit  $Y_{1t}$  is a realisation of the theoretical value  $Y_{it}(1)$  and  $Y_{it}(0)$  will be evaluated simply as a weighted average of (some) other units:

$$Y_{it}(0) = \sum_{i=2}^{I+1} w_i^* Y_{it}$$

- What's hard ? Choosing which units to keep and how to weight them.

# Synthetic controls in brief

## What are synthetic controls

- Define an outcome  $Y_{it}$  with  $I + 1$  aggregated units of interest and  $i = 1$  is the treated region.
- Each region is observed over  $T$  periods, treatment occurs at time  $t_0$ .
- Like in Rubin's causal model, treatment effect is defined at time  $t$  by

$$Y_{it}(1) - Y_{it}(0)$$

where, here, the treated unit  $Y_{1t}$  is a realisation of the theoretical value  $Y_{it}(1)$  and  $Y_{it}(0)$  will be evaluated simply as a weighted average of (some) other units:

$$Y_{it}(0) = \sum_{i=2}^{I+1} w_i^* Y_{it}$$

- What's hard ? Choosing which units to keep and how to weight them.
- Solution: Data-driven procedure. Optimize an algorithm that choose weights  $w_i^*$  that minimize a distance measure.



# Synthetic controls in brief

## One step back to gain intuition

- There is one obvious synthetic control estimator: put equal weights  $\frac{1}{I}$  to every control units.
- The counterfactual is the simple average of untreated country. But why the proportional weight ?
- Synthetic control is a data driven (i.e. machine learning) procedure to find a weighting scheme that minimize an error term over a **training data set** (pre-treatment period) and use it over a **test set** (post-treatment).
- Same "spirit" as a matching estimator but you do not weight units by the inverse of their treatment probability.

# Synthetic controls in brief

## What are synthetic controls

- Ideally, we would like to construct a synthetic control that resembles the treated unit in all relevant pre-intervention characteristics.
- Formalizing this idea we define  $U_i$  as a  $(r \times 1)$  vector of observed covariates for each unit.
- These variables will commonly consist of a set of predictors of the outcome variable.
- Moreover, we define a  $(T_0 \times 1)$  vector  $K = (k_1, \dots, k_{T_0})'$  that denotes some linear combination of **pre-intervention outcomes**:  
$$\bar{Y}_i^K = \sum_{s=1}^{T_0} k_s Y_{is}.$$
- Linear combinations of pre-intervention outcomes can be used to control for unobserved common factors whose effects vary over time. The user can choose to include as many as  $M$  (linearly independent) combinations of pre-intervention outcomes (with  $M \leq T_0$ ) to control for such unobserved common factors.
- Careful: adding them all increases the risk of overfitting.

# Synthetic controls in brief

## What are synthetic controls

- To implement the synthetic control estimator numerically, we need to define a distance between the synthetic controls unit and the treated unit.
- To do that, we combine the characteristics of the exposed unit in the  $(k \times 1)$  matrix  $X_1 = (U_1', \bar{Y}_1^{K_1}, \dots, \bar{Y}_1^{K_M})'$  and the values of the same characteristics of the control units in the  $(k \times J)$  matrix  $X_0$  with the  $j$ -th row  $(U_j', \bar{Y}_j^{K_1}, \dots, \bar{Y}_j^{K_M})'$ .
- Notice that  $k = r + M$ , controls + pre-outcomes.

# Synthetic controls in brief

## What are synthetic controls

- To create the most similar synthetic control unit, the `synth()` function chooses the vector  $W^*$  to minimize a distance,  $\|X_1 - X_0W\|$ , between  $X_1$  and  $X_0W$ , subject to the weight constraints.
- In particular, following Abadie, Diamond, and Hainmueller (2010), the `synth()` function finds  $W^*$  that minimizes

$$\|X_1 - X_0W\|_V = \sqrt{(X_1 - X_0W)' V (X_1 - X_0W)}$$

- where  $V$  is defined as some  $(k \times k)$  symmetric and positive semidefinite matrix.
- The  $V$  matrix is introduced to allow different weights to the variables in  $X_0$  and  $X_1$  depending on their predictive power on the outcome.
- We still need to choose the weights  $V$ .
- Any different set of values for  $V$  gives another `synth` estimator.

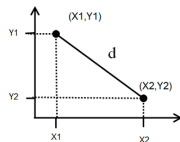
# Synthetic controls in brief

## This distance definition: Intuition ?

$$\|X_1 - X_0W\|_V = \sqrt{(X_1 - X_0W)' V (X_1 - X_0W)}$$

Remember how to calculate the distance Between two point in a N=2 dimension space ?

The "Norm" from above is like a generalization of the Euclidian distance In matrix form and with weights.



$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Figure 1: Distance between 2 coordinates

# Synthetic controls in brief

## Implementation

- **Goal:** we have two sets of weights to determine:
  - $W^*$  : How much importance do we give to each comparison units ?
  - $V$  : How much importance we give characteristics  $X$  to predict  $Y$  ?
- **Solutions in the literature**

- Abadie, Diamond, and Hainmueller (2010) propose to choose  $V$  such that the synthetic control  $W(V)$  minimizes the mean squared prediction error (MSPE) of this synthetic control with respect to  $Y_{1t}^N$

$$\sum_{t \in \mathcal{T}_0} (Y_{1t} - w_2(V)Y_{2t} - \dots - w_{J+1}(V)Y_{J+1t})^2,$$

for some set  $\mathcal{T}_0 \subseteq \{1, 2, \dots, T_0\}$  of **pre-intervention periods**.

- Abadie, Diamond, and Hainmueller (2015) propose to choose the two sets of weights via **out-of-sample validation**.

# Synthetic controls in brief

## Implementation

- **Intuition for both methods**

- The question of choosing  $\mathbf{V} = (v_1, \dots, v_k)$  boils down to assessing the relative importance of each of  $X_{11}, \dots, X_{k1}$  as a predictor of  $Y_{1t}(0)$ . That is, the value  $v_h$  aims to reflect the relative importance of approximating the value of  $X_{h1}$  for predicting  $Y_{1t}(0)$  in the post-intervention period,  $t = T_0 + 1, \dots, T$ .
- $Y_{1t}(0)$  is observed before treatment but not after (where we observe  $Y_{1t}(1)$ ).
- Because  $Y_{1t}(0)$  is not observed for  $t = T_0 + 1, \dots, T$ , we cannot directly evaluate the relative importance of fitting each predictor to approximate  $Y_{1t}(0)$  in the post-intervention period.
- But it is possible to use pre-intervention data to assess the predictive power on  $Y_{1t}(0)$  of the variables  $X_{1j}, \dots, X_{kj}$ .

# Synthetic controls in brief

## Algorithm in the Synth package

- 1 Divide the pre-intervention periods into a initial training period and a subsequent validation period. For simplicity and concreteness, we will assume that  $T_0$  is even and the training and validation periods span  $t = 1, \dots, t_0$  and  $t = t_0 + 1, \dots, T_0$ , respectively, with  $t_0 = T_0/2$ . In practice, the lengths of the training and validation periods may depend on application-specific factors, such as the extent of data availability on outcomes in the pre-intervention and post-intervention periods, and the specific times when the predictors are measured in the data.
- 2 For every value  $\mathbf{V}$ , let  $\tilde{w}_2(\mathbf{V}), \dots, \tilde{w}_{J+1}(\mathbf{V})$  be the synthetic control weights computed with training period data on the predictors. The mean squared prediction error of this synthetic control with respect to  $Y_{1t}(0)$  in the validation period is:

$$\sum_{t=t_0+1}^{T_0} (Y_{1t} - \tilde{w}_2(\mathbf{V})Y_{2t} - \dots - \tilde{w}_{J+1}(\mathbf{V})Y_{J+1t})^2,$$

- 3 Minimize the mean squared prediction error in the previous equation with respect to  $\mathbf{V}$ .
- 4 Use the resulting  $\mathbf{V}^*$  and data on the predictors for the last  $t_0$  periods before in the intervention,  $t = T_0 - t_0 + 1, \dots, T_0$ , to calculate  $\mathbf{W}^* = \mathbf{W}(\mathbf{V}^*)$ .<sup>7</sup>



# Synthetic controls in brief

## Implementation

- The `synth()` function allows for flexibility in the choice of  $V$ .
- The default behaviour follows Abadie, Diamond, and Hainmueller (2010) and  $V^*$  is chosen among all positive definite and diagonal matrices such that the mean squared prediction error (MSPE) of the outcome variable is minimized over some set of pre-intervention periods.
- The pre-period is a "training set", the post period serves as a "validation set".
- This is where it's a form of *machine learning*: use an algorithm to find the  $V$  matrix that minimize the root-mean square error on the validation set.
- Typical "overfitting/bias" trade-off in this setting :
- You can add a lot of variable in the  $X$  matrix and find weights that predict very well the pre-treatment period but perform very poorly outside of it.

# Outline

- 1 Introduction
- 2 Synthetic controls in brief
- 3 First application: Basque terrorism from Abadie and Gardeazabal (2003)
  - Context
  - First step: preparing data
  - Second step: adjustments and estimation
  - Estimated dopplegänger for Basque region
  - Interpretations
  - Inference
- 4 Cool implementations of Synth controls

# First application: Basque terrorism from Abadie and Gardeazabal (2003)

## Context

- Abadie and Gardeazabal (2003) estimates the impact of terrorism in the Basque country on growth.
- Terrorism started in 1970
- They cannot use a standard DiD method because none of the other Spanish regions followed the same time trend as the Basque Country
- They therefore take a weighted average of other Spanish regions as a synthetic control group

# First step: preparing data

```
# Load the dataset from the Synth package
data(basque)
# Take a look at the DB We'll use gdpicap as dependent and use share of
# each level of education in the population and investments as predictor
# We also use some specific sector shares for special years dataprep:
# prepare data for synth
dataprep.out <- dataprep(foo = basque, predictors = c("school.illit", "school.prim",
"school.med", "school.high", "school.post.high", "invest"), predictors.op = c("mean"),
dependent = c("gdpicap"), unit.variable = c("regionno"), time.variable = c("year"),
special.predictors = list(list("gdpicap", 1960:1969, c("mean")), list("sec.agriculture",
seq(1961, 1969, 2), c("mean")), list("sec.energy", seq(1961, 1969, 2),
c("mean")), list("sec.industry", seq(1961, 1969, 2), c("mean")), list("sec.construction",
seq(1961, 1969, 2), c("mean")), list("sec.services.venta", seq(1961,
1969, 2), c("mean")), list("sec.services.nonventa", seq(1961, 1969,
2), c("mean")), list("popdens", 1969, c("mean"))), treatment.identifier = 17,
controls.identifier = c(2:16, 18), time.predictors.prior = c(1964:1969),
time.optimize.ssr = c(1960:1969), unit.names.variable = c("regionname"),
time.plot = c(1955:1997))
```

## Second step: adjustments and estimation

```
### In the paper, they make a few adjustments to the source data: I
### replicate them here. 1. combine highest and second highest schooling
### category and eliminate highest category
dataprep.out$X1["school.high", ] <- dataprep.out$X1["school.high", ] + dataprep.out$X1["school.post.high",
]
dataprep.out$X1 <- as.matrix(dataprep.out$X1[-which(rownames(dataprep.out$X1) ==
"school.post.high"), ])
dataprep.out$X0["school.high", ] <- dataprep.out$X0["school.high", ] + dataprep.out$X0["school.post.high",
]
dataprep.out$X0 <- dataprep.out$X0[-which(rownames(dataprep.out$X0) == "school.post.high"),
]

# 2. make total and compute shares for the schooling categories
lowest <- which(rownames(dataprep.out$X0) == "school.illit")
highest <- which(rownames(dataprep.out$X0) == "school.high")

dataprep.out$X1[lowest:highest, ] <- (100 * dataprep.out$X1[lowest:highest,
])/sum(dataprep.out$X1[lowest:highest, ])

dataprep.out$X0[lowest:highest, ] <- 100 * scale(dataprep.out$X0[lowest:highest,
], center = FALSE, scale = colSums(dataprep.out$X0[lowest:highest, ]))

# run synth
synth.out <- synth(data.prep.obj = dataprep.out)
```

## Second step: adjustments and estimation

### What happened inside the function

Let's get the  $W$  matrix

```
# Get result tables  
synth.tables <- synth.tab(dataprep.res = dataprep.out, synth.res = synth.out)  
# Look at the W matrix  
synth.tables$tab.w %>%  
  kbl()
```

	w.weights	unit.names	unit.numbers
2	0.000	Andalucia	2
3	0.000	Aragon	3
4	0.000	Principado De Asturias	4
5	0.000	Baleares (Islas)	5
6	0.000	Canarias	6
7	0.000	Cantabria	7
8	0.000	Castilla Y Leon	8
9	0.000	Castilla-La Mancha	9
10	0.851	Cataluna	10
11	0.000	Comunidad Valenciana	11
12	0.000	Extremadura	12
13	0.000	Galicia	13
14	0.149	Madrid (Comunidad De)	14
15	0.000	Murcia (Region de)	15
16	0.000	Navarra (Comunidad Foral De)	16
18	0.000	Rioja (La)	18

## Second step: adjustments and estimation

### What happened inside the function

- Only 2 Spanish regions were picked by the algorithm to be used in the synthetic Basque country: Catalonia, Madrid (Comunidad De).
- Let's plot the weights given to the predictor and see what was picked and how much importance was given:

```
Vs <- cbind(v.weights = unlist(synth.tables$tab.v), keyName = unlist(labels(synth.tables[["tab.v"]])
  as.data.frame(.) %>%
  mutate(v.weights = as.numeric(v.weights), keyName = as.factor(keyName))
# plot weights
plotweights <- ggplot(Vs) + geom_col(aes(x = keyName, y = v.weights, fill = v.weights)) +
  coord_flip() + theme(legend.position = "none") + xlab("Predictors")
```

## Second step: adjustments and estimation

### What happened inside the function

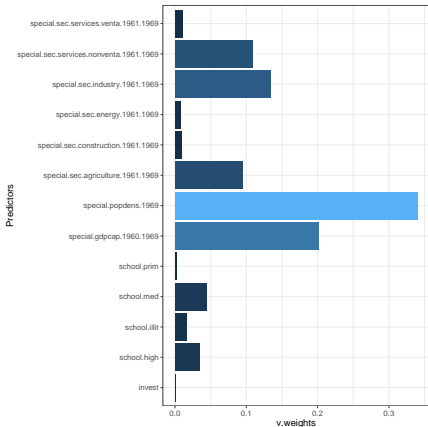
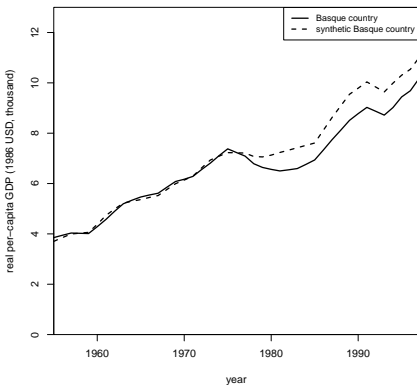


Figure 2: Weights of the predictors

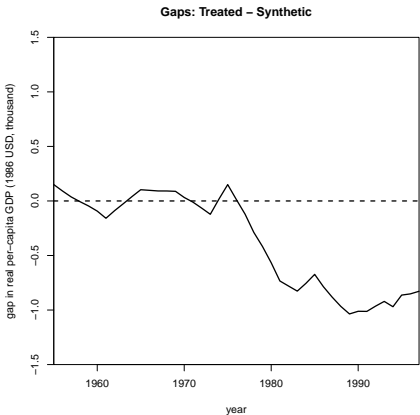


# Estimated doppelgänger for Basque region



**Figure 3:** Observed and doppelgänger Basque region using Synthetic control replicating Abadie and Gardeazabal (2003)

# Estimated dopplegänger for Basque region



**Figure 4:** Effect of terrorism on Basque GDP estimated using Synthetic control replicating Abadie and Gardeazabal (2003)

# First application: Basque terrorism from Abadie and Gardeazabal (2003)

## Interpretations

- The algorithm uses data before 1970 to choose two sets of weights for regions and predictors from a pool of control regions to compute an average of GDP per capita that is "as close as" possible as the observed Basque region before 1970.
- This model is then used to predict counterfactual for Basque country in the absence of terrorism.
- This prediction is **out of sample**, we use parameters obtained from past data to predict post 1970  $Y_0$ .
- As with Dif-in-Dif, the causal interpretation relies on the parallel trend assumption.
- Here, the argument is that it is not plausible that factors that produce a tight fit before would diverge afterwards
- Now, how do we know if it's **this particular** curve (sets of weights) ?  
How do we test significance ?

# Inference

Conventional statistical inference is difficult because we typically have two time series

- 2T observations
- strong serial correlation and too few clusters

Alternative: permutation tests

- run placebo SC on all units in the donor pool
- compute the treatment effect for each placebo
- compare placebos to the estimated treatment effect
- compute empirical p-value

# Inference

- You need to load

```
library("SCTools")
```

to use the function *generate.placebos*

```
##### Inference on synthetic control relies on the estimation of the same  
##### model on placebo states and compare the ratio between pre/post MSPE  
##### (for instance)
```

```
placebos <- generate.placebos(dataprep.out, synth.out, Sigf.ipop = 5, strategy = "multicore")  
placeboplots <- plot_placebos(placebos)
```

# Inference

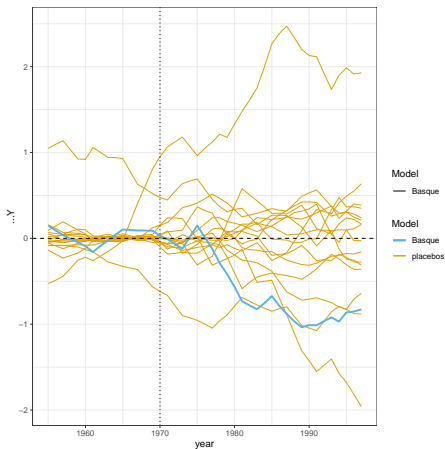


Figure 5: Placebo estimates of the synth model on other control units

# Inference

## Permutation test: how to:

- Iteratively apply the synthetic control method to each country/state in the donor pool and obtain a distribution of placebo effects.
- Calculate the RMSPE for each placebo for the pre-treatment period:

$$RMSPE = \left( \frac{1}{T - T_0} \sum_{t=T_0+t}^T \left( Y_{1t} - \sum_{j=2}^{J+1} w_j^* Y_{jt} \right)^2 \right)^{\frac{1}{2}}$$

- Calculate the RMSPE for each placebo for the post-treatment period (similar equation but for the post-treatment period).
- Compute the ratio of the post- to pre-treatment RMSPE.
- Sort this ratio in descending order from greatest to highest.
- Calculate the treatment unit's ratio in the distribution as  $p = \text{RANK} / \text{TOTAL}$

# Inference

## Exact p-value

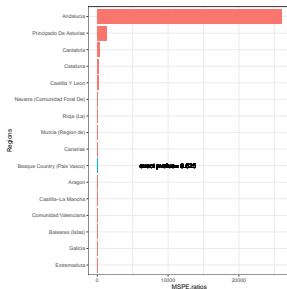


Figure 6: Bar chart of the Post/Pre MSE ratio

The placebo tests tell us the synthetic estimation is not clearly an outlier compared with placebo estimates so actually, we can't say that the effect is different from alternative permutations.

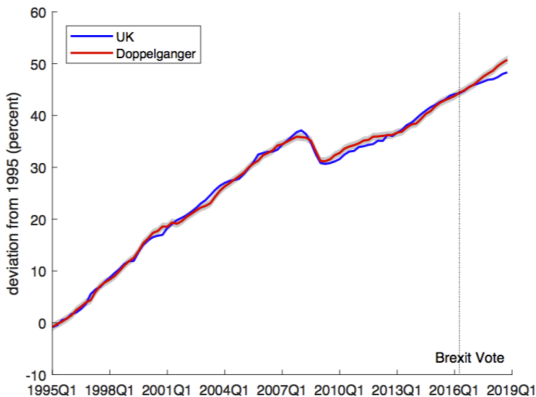


# Outline

- 1 Introduction
- 2 Synthetic controls in brief
- 3 First application: Basque terrorism from Abadie and Gardeazabal (2003)
- 4 Cool implementations of Synth controls**
  - Impact of Brexit on growth (Born et al. 2019)
  - mandatory COVID-19 certificates on vaccine uptake (Mills and Rüttenauer 2022)
- 5 Application: Prison and black male incarceration

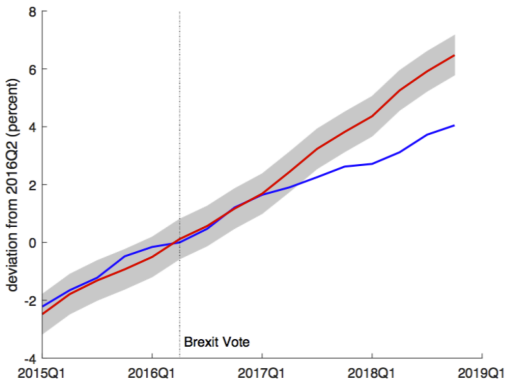
# Cool implementations of Synth controls

## Impact of Brexit on growth (Born et al. 2019)



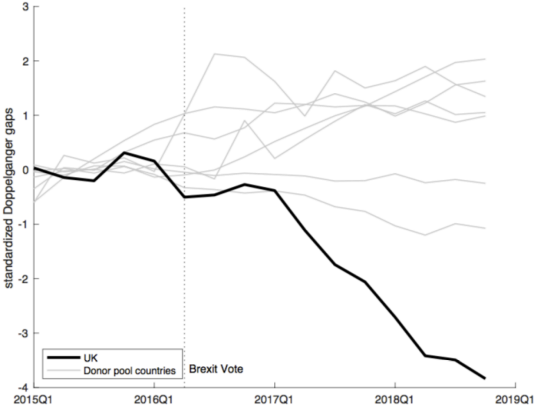
# Cool implementations of Synth controls

## Impact of Brexit on growth (Born et al. 2019)



# Cool implementations of Synth controls

## Impact of Brexit on growth (Born et al. 2019)





# Outline

- ① Introduction
- ② Synthetic controls in brief
- ③ First application: Basque terrorism from Abadie and Gardeazabal (2003)
- ④ Cool implementations of Synth controls
- ⑤ Application: Prison and black male incarceration  
Context and motivations

## Application: Prison and black male incarceration

All data and context come from (Cunningham 2018, chapter 10.) and researches he made a while ago.

### Context and motivations

In 1980, the Texas Department of Corrections (TDC) lost a major civil action lawsuit, *Ruiz v. Estelle*; Ruiz was the prisoner who brought the case, and Estelle was the warden. The case argued that TDC was engaging in unconstitutional practices related to overcrowding and other prison conditions. Texas lost the case, and as a result, was forced to enter into a series of settlements. To amend the issue of overcrowding, the courts placed constraints on the number of inmates who could be placed in cells. To ensure compliance, TDC was put under court supervision until 2003.

## Application: Prison and black male incarceration

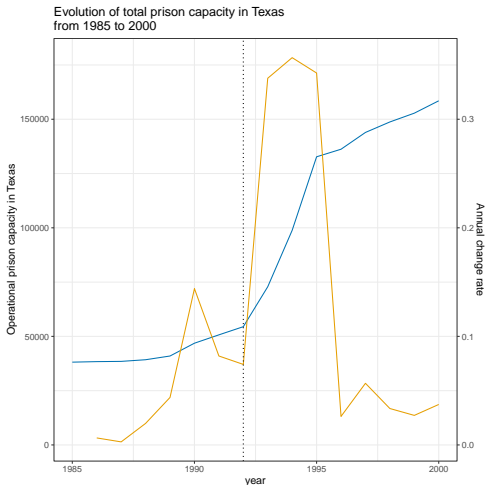
### Context and motivations

Given these constraints, the construction of new prisons was the only way that Texas could keep arresting as many people as its police departments wanted to without having to release those whom the TDC had already imprisoned. If it didn't build more prisons, the state would be forced to increase the number of people to whom it granted parole. That is precisely what happened; following *Ruiz v. Estelle*, Texas used parole more intensively. But then, in the late 1980s, Texas Governor Bill Clements began building prisons. Later, in 1993, Texas Governor Ann Richards began building even more prisons. Under Richards, state legislators approved \$1 billion for prison construction, which would double the state's ability to imprison people within three years.



# Context and motivations

## How much more inmates ?



## Application: Prison and black male incarceration

### What should happen ?

- Just because you have more prisons does not mean the incarceration rate should increase, right ?
- But because the state was using parole to comply with the regulation following the Ruiz vs Estelle case, this is what happened.
- **Research question** : How did the construction of new prisons affected the black male rate of incarceration ?

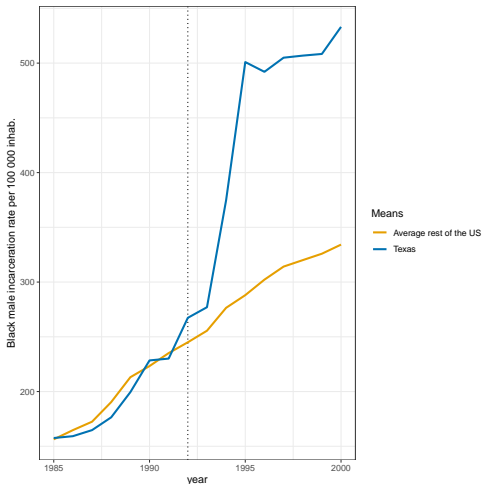
# Application: Prison and black male incarceration

## Context and motivations

- 1 Try and get a sense of what is happening by plotting black male population incarcerated in texas vs the rest of the USA (variable bmp)
- 2 Use synthetic control to estimate counterfactual texas bmprate and see whether building more prison increases black male incarceration rate.
- 3 Use the *Main\_code\_TA3.R* file to guide your work

# What I expect in the end

## Evolution of the dependent variable



# What I expect in the end

## Prepare synth data

```
texas <- read_data("texas.dta") %>%
  as.data.frame(.)

dataprep_out <- dataprep(
  foo = texas, #database texas
  predictors = c("poverty", "income"), #main predictor are poverty and income
  predictors.op = "mean", # operator we want to use is the mean
  time.predictors.prior = 1985:1993, #We predict from 1985 to 1993, and
  #Following Cuningham and other, we add special predictor in a list : bmprison from 1988, 1990:1992)
  special.predictors = list(
    list("bmprison", c(1988, 1990:1992), "mean"),
    list("alcohol", 1990, "mean"),
    list("aidscapita", 1990:1991, "mean"),
    list("black", 1990:1992, "mean"),
    list("perc1519", 1990, "mean")),
  dependent = "bmprison",
  unit.variable = "statefip",
  unit.names.variable = "state",
  time.variable = "year",
  treatment.identifier = 48, #Texas is the 48 state in the list "statefip"
  controls.identifier = c(1,2,4:6,8:13,15:42,44:47,49:51,53:56),
  time.optimize.ssr = 1985:1993,
  time.plot = 1985:2000
)
# Now we can run the synth command
synth_out <- synth(data_prep_obi = dataprep_out)
```

# What I expect in the end

## Results of the Synth estimation

```
# Get result tables  
synth_tables <- synth.tab(dataprep.res = dataprep_out, synth.res = synth_out)
```

- Only 3 US states were picked by the algorithm to be used in the synthetic Texas: California, Florida, Louisiana.
- Let's plot the weights given to the predictor and see what was picked and how much importance was given:

```
Vs2 <- cbind(v.weights = unlist(synth_tables$tab.v), keyName = unlist(labels(synth_tables[["tab.v"]  
  as.data.frame(.) %>%  
  mutate(v.weights = as.numeric(v.weights), keyName = as.factor(keyName))  
# plot weights  
plotweights2 <- ggplot(Vs2) + geom_col(aes(x = keyName, y = v.weights, fill = v.weights)) +  
  coord_flip() + theme(legend.position = "none") + xlab("Predictors")
```

# What I expect in the end

## What happened inside the function

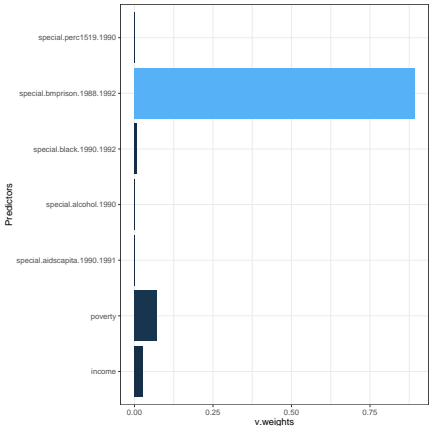
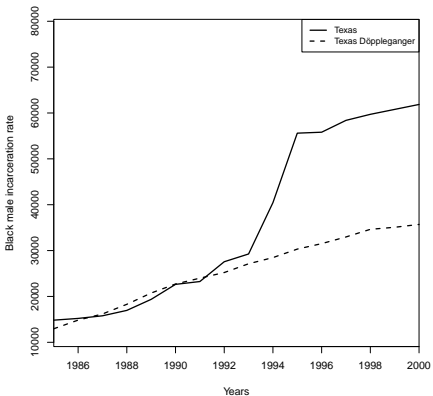


Figure 9: Weights of the predictors

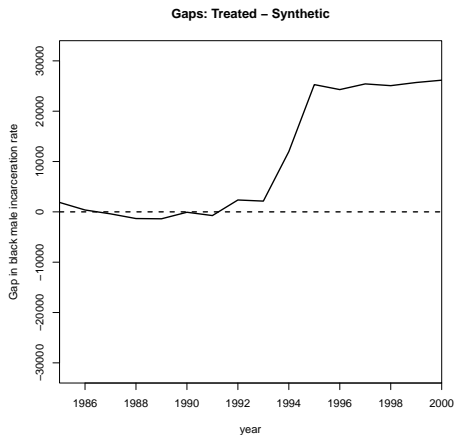
# What I expect in the end



**Figure 10:** Observed and doppelgänger Texas black male incarceration rate per 10 000 inhabitants.



# What I expect in the end



**Figure 11:** Effect of building prisons on black male incarceration rate

# What I expect in the end

## Interpretations

- The algorithm uses data before 1992 to choose two sets of weights for states and predictors from a pool of control states to compute an average black male incarceration rate "as close as" possible as the observed one for Texas region before 1992.
- This model is then used to predict counterfactual for Texas in the absence of massive prison construction.
- This prediction is **out of sample**, we use parameters obtained from past data to predict post 1992  $Y_0$ .
- As with Dif-in-Dif, the causal interpretation relies on the parallel trend assumption.
- Here, the argument is that it is not plausible that factors that produce a tight fit before would diverge afterwards
- Now, let's estimate placebo models on other states and make the plot.

# What I expect in the end

- You need to load

```
library("SCTools")
```

to use the function *generate.placebos*

```
##### Inference on synthetic control relies on the estimation of the same  
##### model on placebo states and compare the ratio between pre/post MSPE  
##### (for instance)
```

```
placebosBM <- generate.placebos(dataprep_out, synth_out, Sigf.ipop = 3, strategy = "multicore")  
placeboplotsBM <- plot_placebos(placebosBM)
```

# Context and motivations

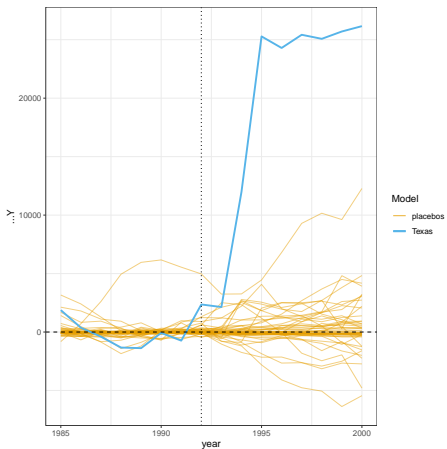
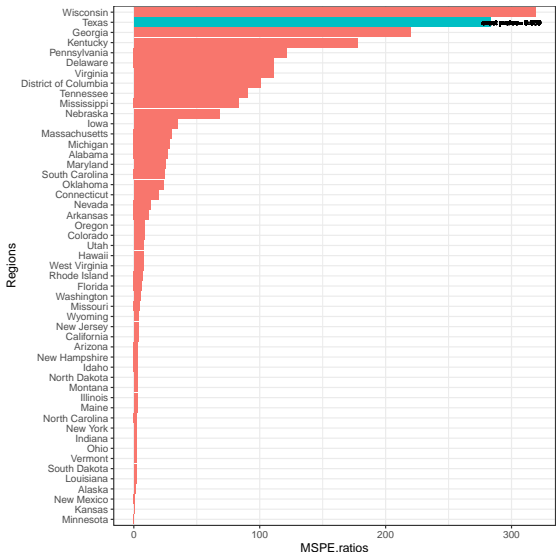


Figure 12: Placebo estimates of the synth model on other control states

# Context and motivations



# What I expect in the end

## Interpretations

- The placebo estimations show that Texas is a clear outlier. Out of all the permutation there's only one placebo estimation whose root mean square prediction error is smaller.
- From our analysis, it seems plausible that the building of prisons from 1992 in Texas dramatically and causally increased the incarceration of black males.
- That's it. That's what we just showed.

# Bibliography I

- ▶ Abadie, Alberto. 2021. "Using Synthetic Controls: Feasibility, Data Requirements, and Methodological Aspects." *Journal of Economic Literature* 59, no. 2 (June 1, 2021): 391–425.
- ▶ Abadie, Alberto, Alexis Diamond, and Jens Hainmueller. 2010. "Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California's Tobacco Control Program." *Journal of the American Statistical Association* 105, no. 490 (June 1, 2010): 493–505.
- ▶ ———. 2011. "**Synth** : An R Package for Synthetic Control Methods in Comparative Case Studies." *Journal of Statistical Software* 42 (13).
- ▶ ———. 2015. "Comparative Politics and the Synthetic Control Method: COMPARATIVE POLITICS AND THE SYNTHETIC CONTROL METHOD." *American Journal of Political Science* 59, no. 2 (February): 495–510.
- ▶ Abadie, Alberto, and Javier Gardeazabal. 2003. "The Economic Costs of Conflict: A Case Study of the Basque Country." *American Economic Review* 93, no. 1 (February 1, 2003): 113–132.
- ▶ Born, Benjamin, Gernot J Müller, Moritz Schularick, and Petr Sedláček. 2019. "The Costs of Economic Nationalism: Evidence from the Brexit Experiment\*." *The Economic Journal* 129, no. 623 (October 1, 2019): 2722–2744.
- ▶ Cunningham, Scott. 2018. *Causal Inference: The Mixtape*.
- ▶ Mills, Melinda C, and Tobias Rüttenauer. 2022. "The Effect of Mandatory COVID-19 Certificates on Vaccine Uptake: Synthetic-Control Modelling of Six Countries." *The Lancet Public Health* 7, no. 1 (January): e15–e22.