



# ScPoEconometrics

## Differences-in-Differences

Florian Oswald, Gustave Kenedi and Pierre Villedieu  
SciencesPo Paris  
2024-02-12

# Recap from last week

- Applied inference tools to regression analysis
- *Standard error* of regression coefficients
- *Statistical significance* of regression coefficients

## Today: *Differences-in-differences*

- Exploits changes in policy over time that don't affect everyone
- Need to find (or construct) appropriate control group(s)
- *Key assumption*: parallel trends
- *Empirical application*: impact of *minimum wage* on *employment*



# Evaluation methods

- Multiple regression often does not provide causal estimates because of *selection on unobservables*.
- RCTs are one way to solve this problem but they are often impossible to do.
- Four main causal evaluation methods used in economics:
  - *instrumental variables (IV)*,
  - *propensity-score matching*,
  - *differences-in-differences (DiD)*, and
  - *regression discontinuity designs (RDD)*.
- These methods are used to identify **causal relationships** between treatments and outcomes.
- In this lecture, we will cover a popular and rigorous program evaluation method: **differences-in-differences**.
- Next week we will look at **Instrumental Variables**.



# Differences-in-Differences (DiD)

- Usual starting point: subjects are not randomly allocated to treatment ⚠

## DiD Requirements:

- 2 time periods: before and after treatment.
- 2 groups:
  - *control group*: never receives treatment,
  - *treatment group*: initially untreated and then fully treated.
- Under certain assumptions, control group can be used as the counterfactual for treatment group



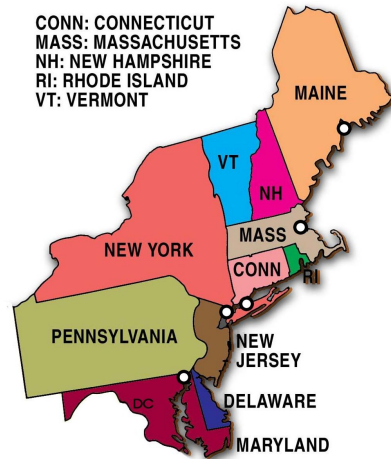
# An Example: Minimum Wage and Employment

- Imagine you are interested in assessing the **causal** impact of increasing the minimum wage on (un)employment.
- Why is this not that straightforward? What should the control group be?
- Seminal 1994 **paper** by prominent labor economists David Card and Alan Krueger entitled "Minimum Wages and Employment: A Case Study of the Fast-Food Industry in New Jersey and Pennsylvania"
- Estimates the effect of an increase in the minimum wage on the employment rate in the fast-food industry. Why this industry?



# Institutional Details

- In the US, there is a national minimum wage, but states can depart from it.
- April 1, 1992: New Jersey minimum wage increases from \$4.25 to \$5.05 per hour.
- Neighboring Pennsylvania did not change its minimum wage level.



Pennsylvania and New Jersey are *very similar*: similar institutions, similar habits, similar consumers, similar incomes, similar weather, etc.



# Card and Krueger (1994): Methodology

- Surveyed 410 fast-food establishments in New Jersey (NJ) and eastern Pennsylvania
- Timing:
  - Survey before NJ MW increase: Feb/March 1992
  - Survey after NJ MW increase: Nov/Dec 1992
- What comparisons do you think they did?

Let's take a closer at their data

```
# install package that contains the cleaned data
remotes::install_github("b-rodrigues/diffindiff")
# load package
library(diffindiff)
# load data
ck1994 <- njmin
```

```
ck1994 %>%
  select(sheet, chain, state, observation, empft, emppt)
head()
```

```
## # A tibble: 6 × 6
##   sheet chain  state      observation  empft emppt
##   <chr> <chr>  <chr>      <chr>      <dbl> <dbl>
## 1 46    bk     Pennsylvania February 1992  30    15
## 2 49    kfc     Pennsylvania February 1992   6.5   6.5
## 3 506   kfc     Pennsylvania February 1992    3     7
## 4 56    wendys Pennsylvania February 1992  20    20
## 5 61    wendys Pennsylvania February 1992    6    26
## 6 62    wendys Pennsylvania February 1992    0    31
```



# Task 1 (10 minutes)

1. Take a look at the dataset and list the variables. Check the variable definitions with ?njmin.
2. Tabulate the number of stores by state and by survey wave (observation). Does it match what's in *Table 1* of the **paper**?
3. Create a full-time equivalent (FTE) employees variable called `empfte` equal to  $\text{empft} + 0.5 * \text{emppt} + \text{nmgrs}$ . `empft` and `emppt` correspond respectively to the number of full-time and part-time employees. `nmgrs` corresponds to the number of managers. This is how Card and Krueger compute their full-time equivalent (FTE) employment variable (p.775 of the paper).
4. Compute the average number of FTE employment, average percentage of FT employees (out of the number of FTE employees), and average starting wage (`wage_st`) by state and by survey wave. Compare your results with *Table 2* of the paper.
5. How different are New Jersey and Pennsylvania's fast-food restaurants before the minimum wage increase?





# Task 1: Solution

## 1. Load the data

```
remotes::install_github("b-rodrigues/diffindiff")  
  
library(diffindiff)  
library(tidyverse)  
  
ck1994 <- njmin
```

## 2. Table of #stores by state and wave

```
library(skimr)  
  
ck1994 %>% group_by(state, observation) %>% summarise(n_stores = n_distinct(sheet))
```

## 3. Create a full-time equivalent (FTE) employees variable called empft

```
library(skimr)  
  
ck1994 = ck1994 %>% ungroup %>% mutate(empfte = empft + 0.5*emppt + nmgrs)
```



# Task 1: Solution

4. Compute average FTE employment, average percentage of FT employees and average starting wage (wage\_st) by state and by survey wave.

```
library(skimr)

did = ck1994 %>% group_by(state, observation) %>% summarise(avg_fte = mean(empfte, na.rm = T),
                                                             pc_ft = mean(empft/empfte, na.rm = T)*100,
                                                             avg_w = mean(wage_st, na.rm = T))
```



# Card and Krueger DiD: Tabular Results

## Average Employment Per Store Before and After the Rise in NJ Minimum Wage

Variables	Pennsylvania	New Jersey
FTE employment before	23.33	20.44
FTE employment after	21.17	21.03
Change in mean FTE employment	-2.17	0.59

## DiD Estimate

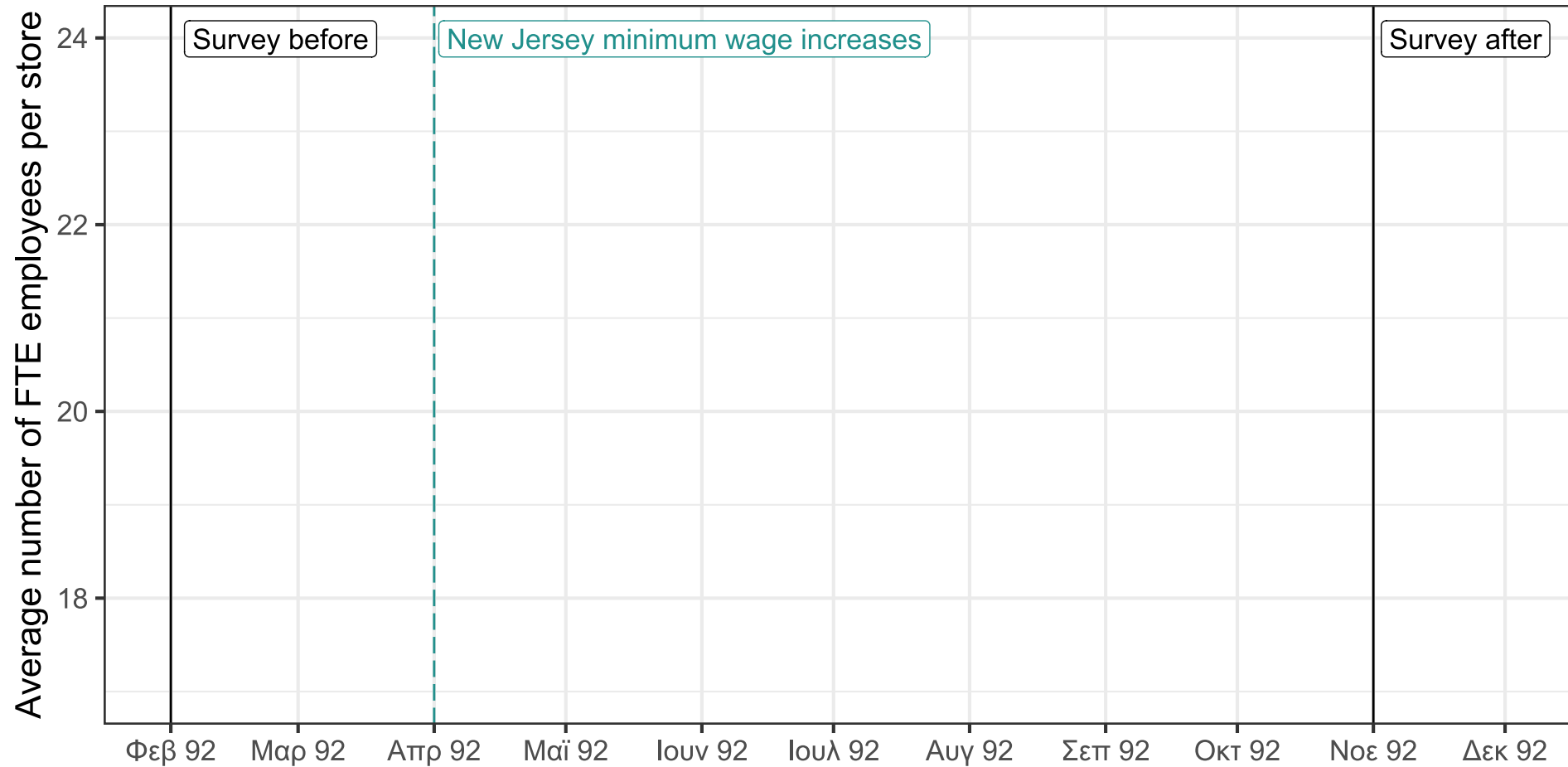
Differences-in-differences causal estimate:  $0.59 - (-2.17) = 2.76$

Yes the essence of differences-in-differences is *that* simple! 😊

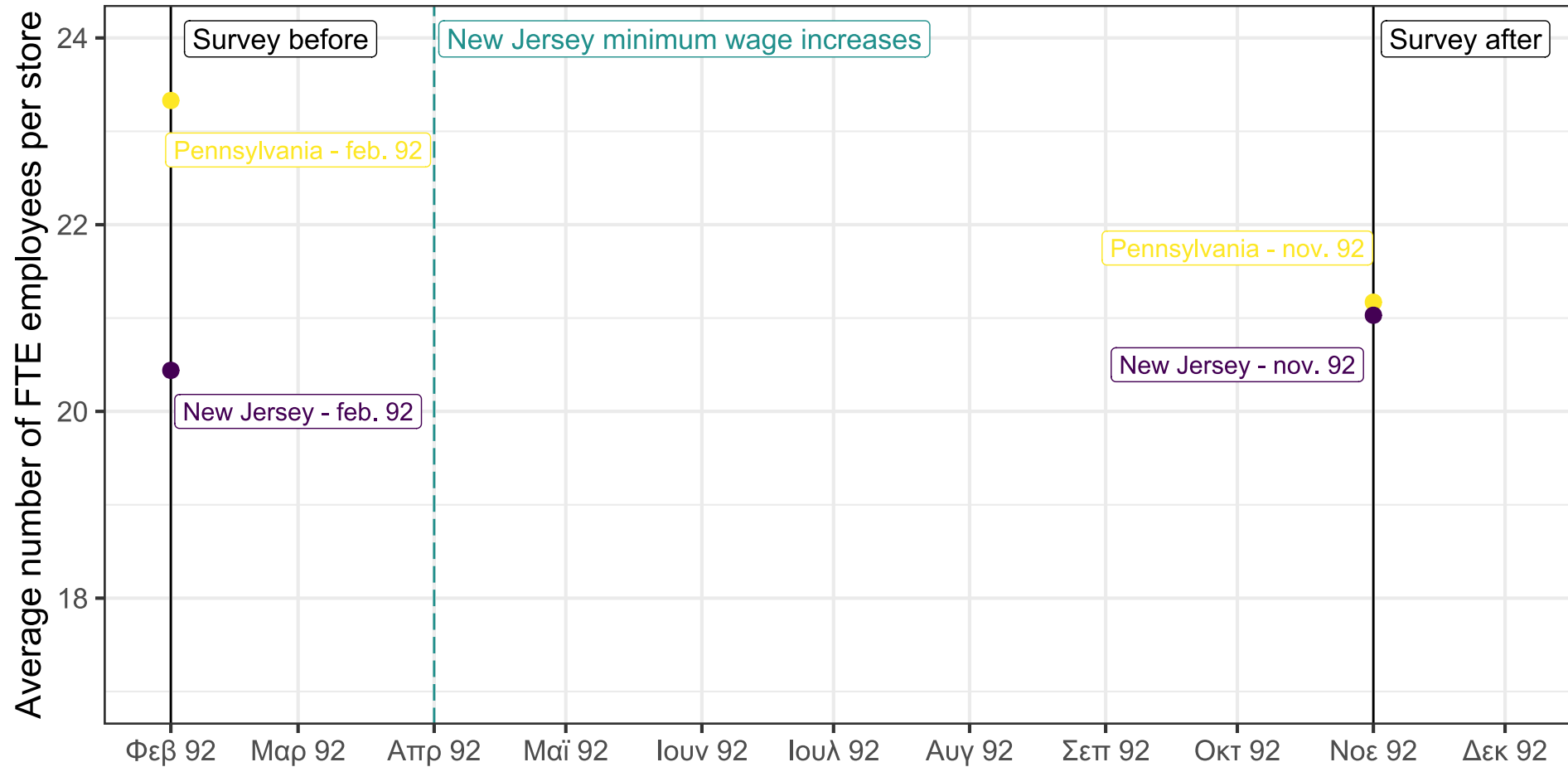
Let's look at these results graphically.



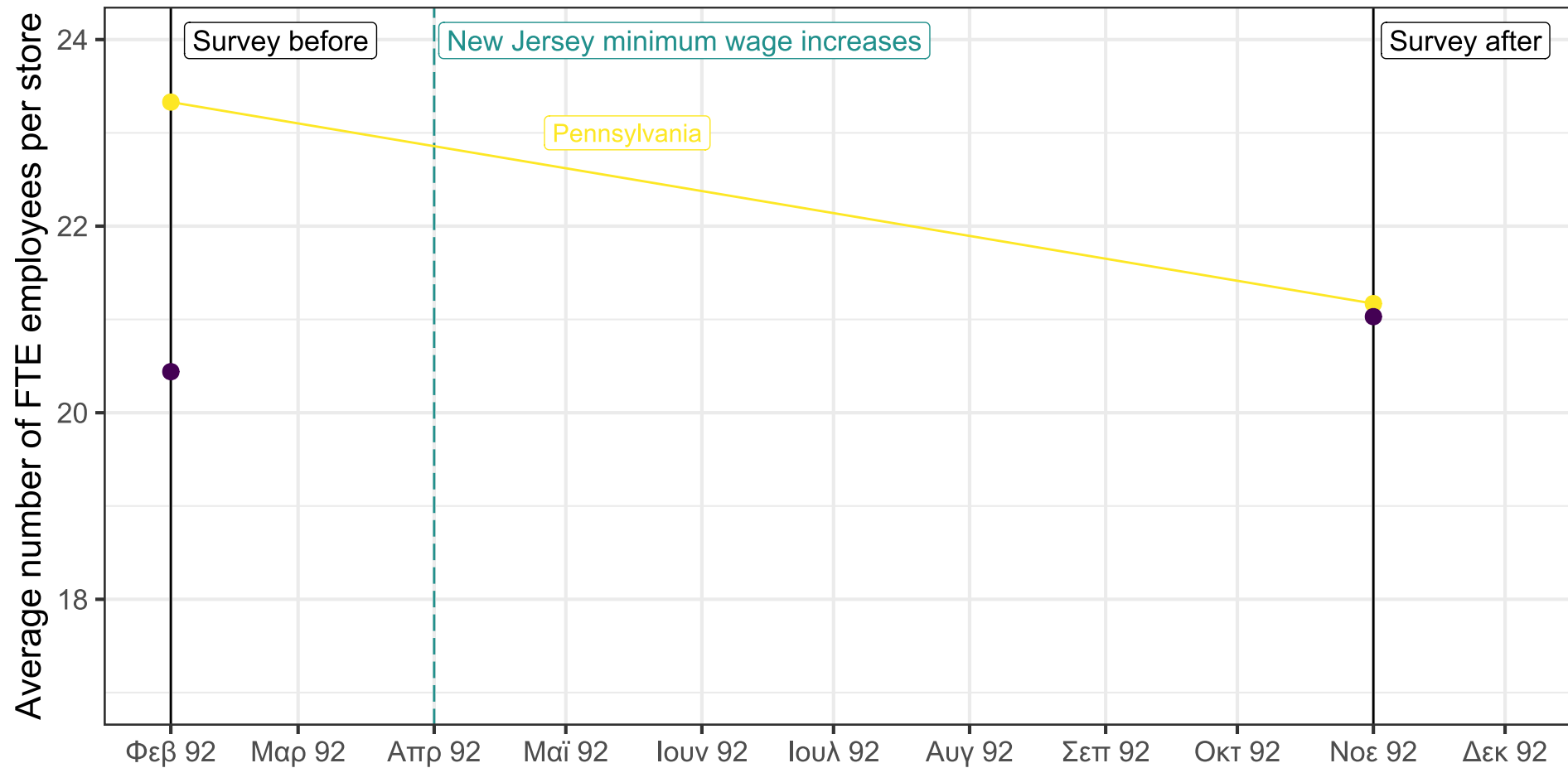
# DiD Graphically



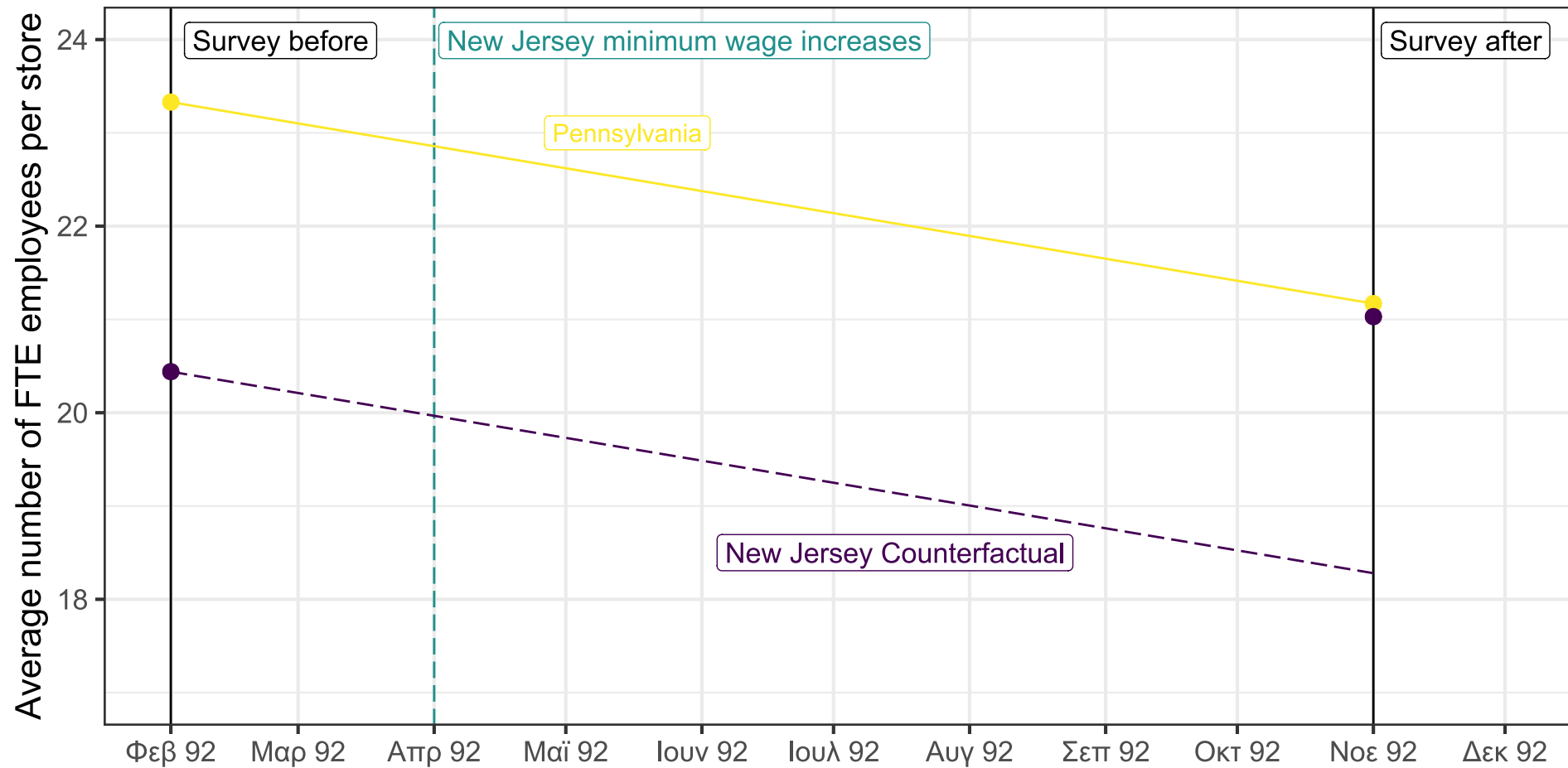
# DiD Graphically



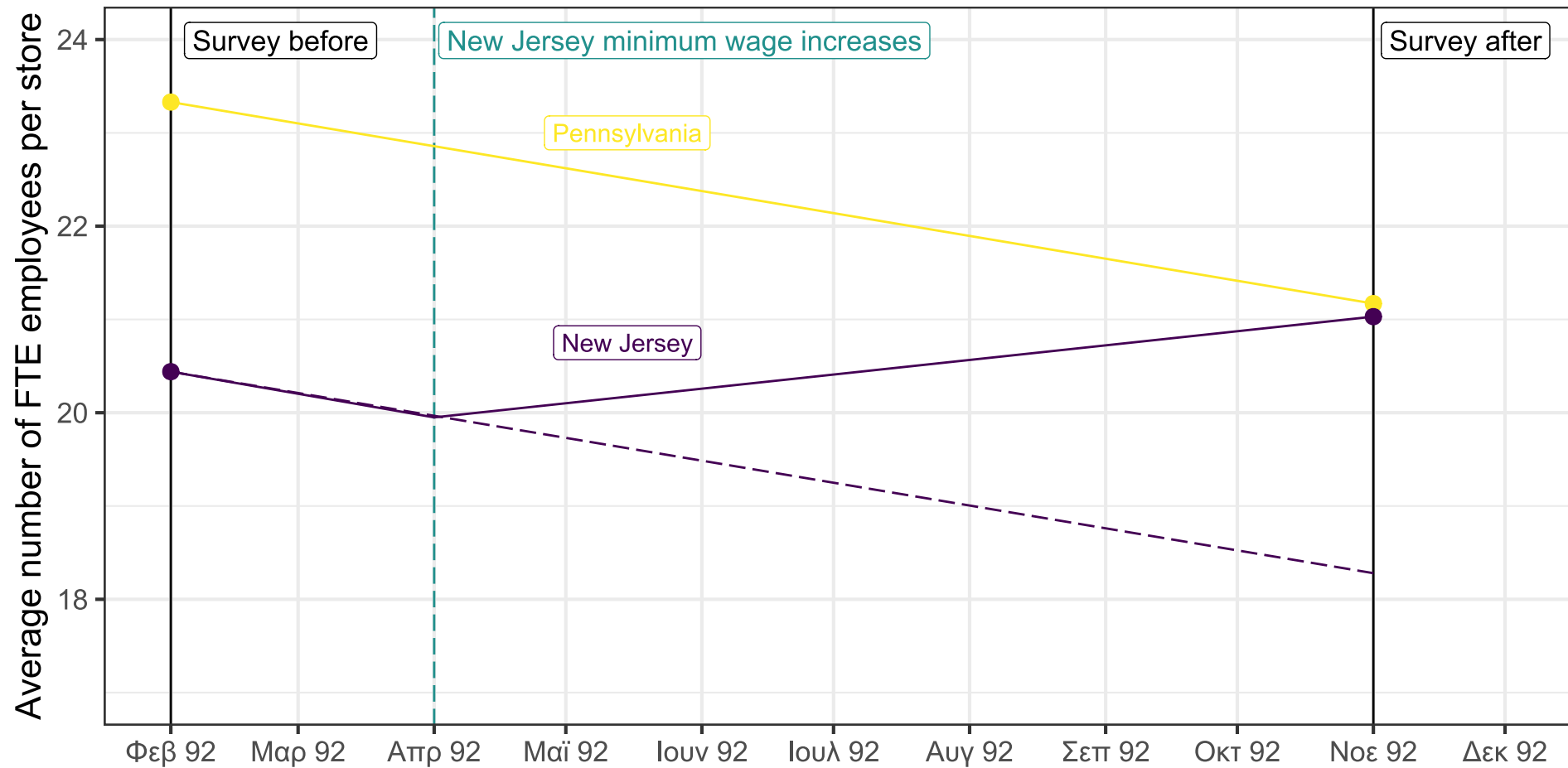
# DiD Graphically



# DiD Graphically

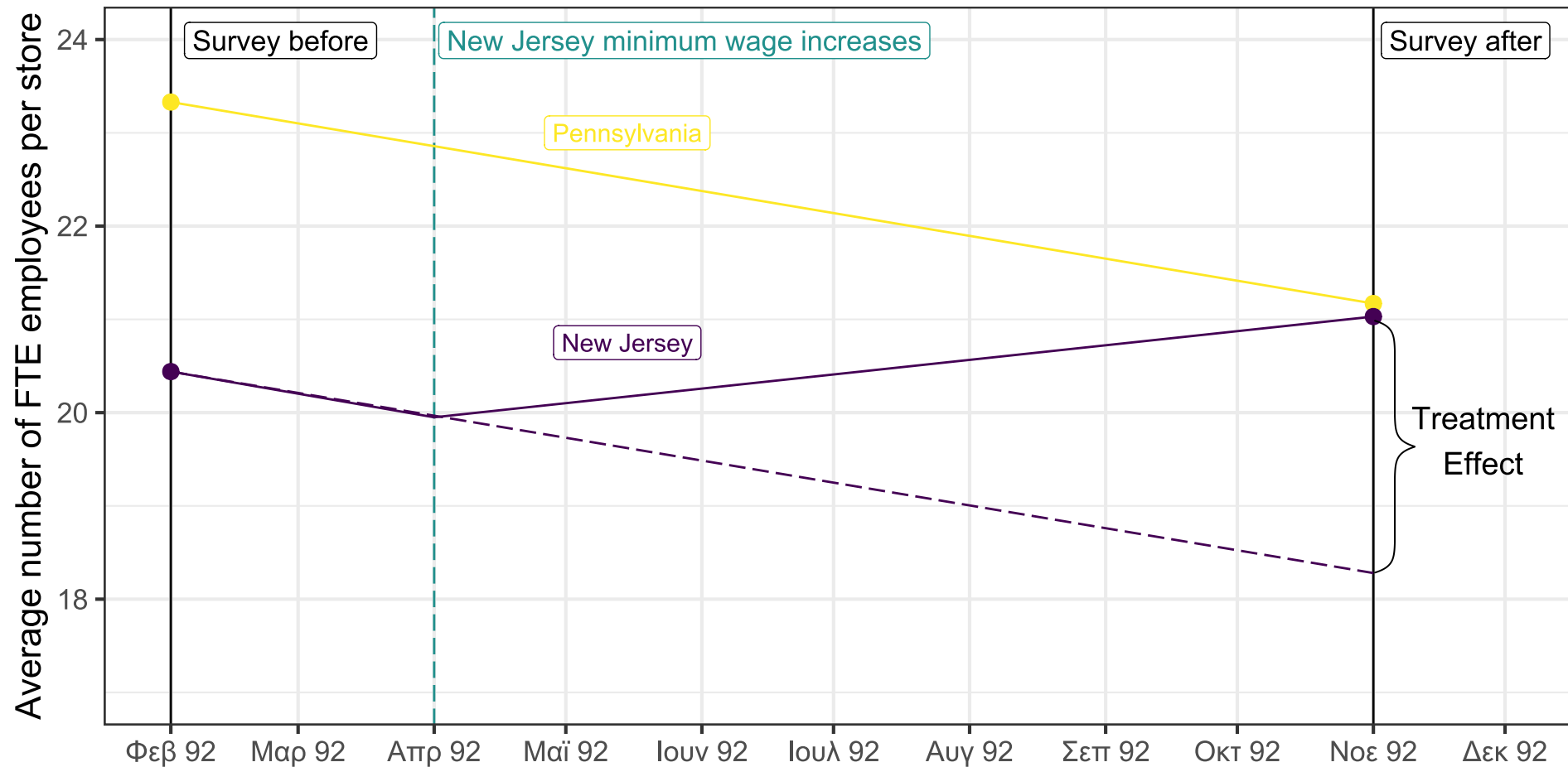


# DiD Graphically

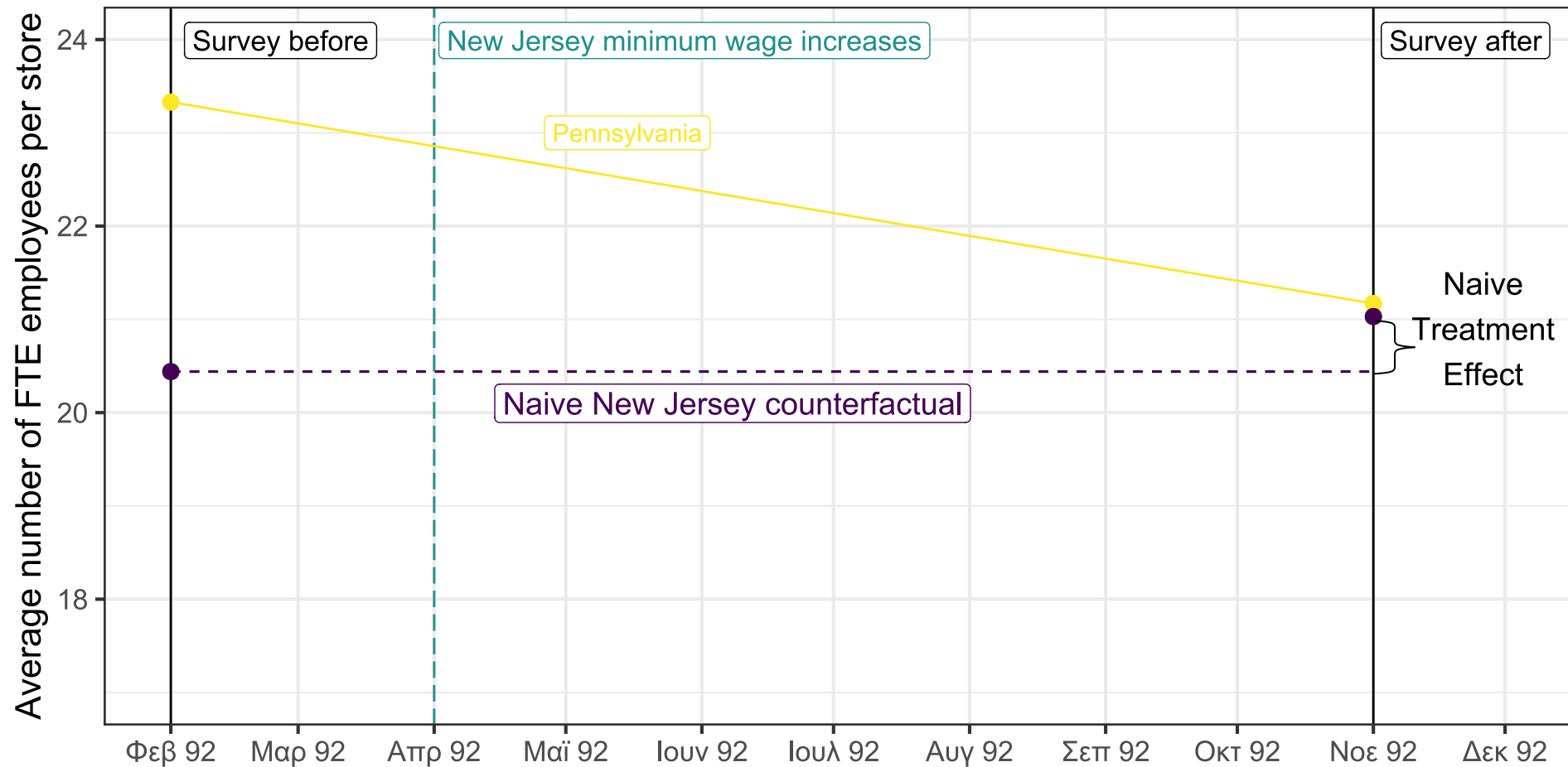




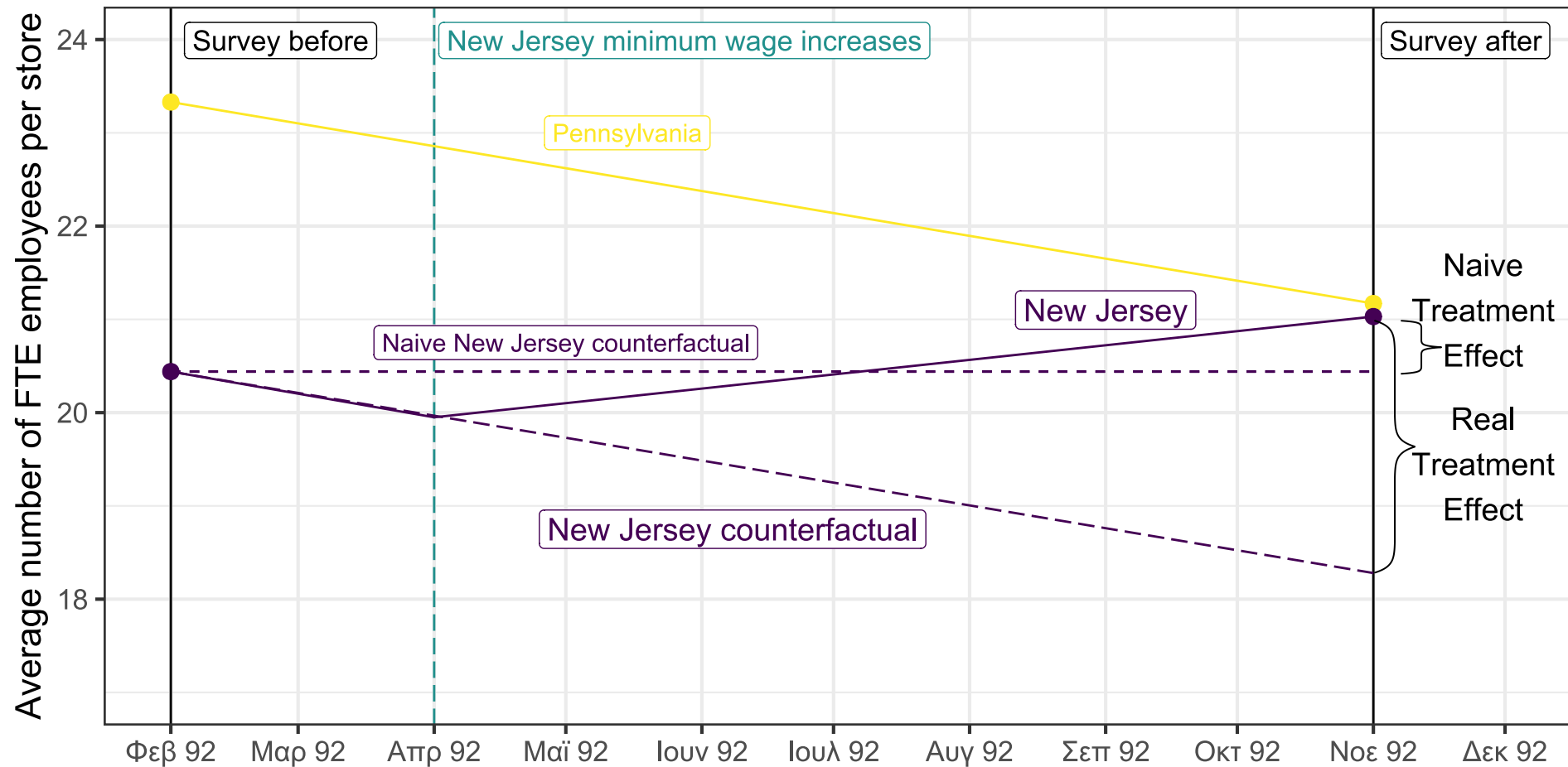
# DiD Graphically



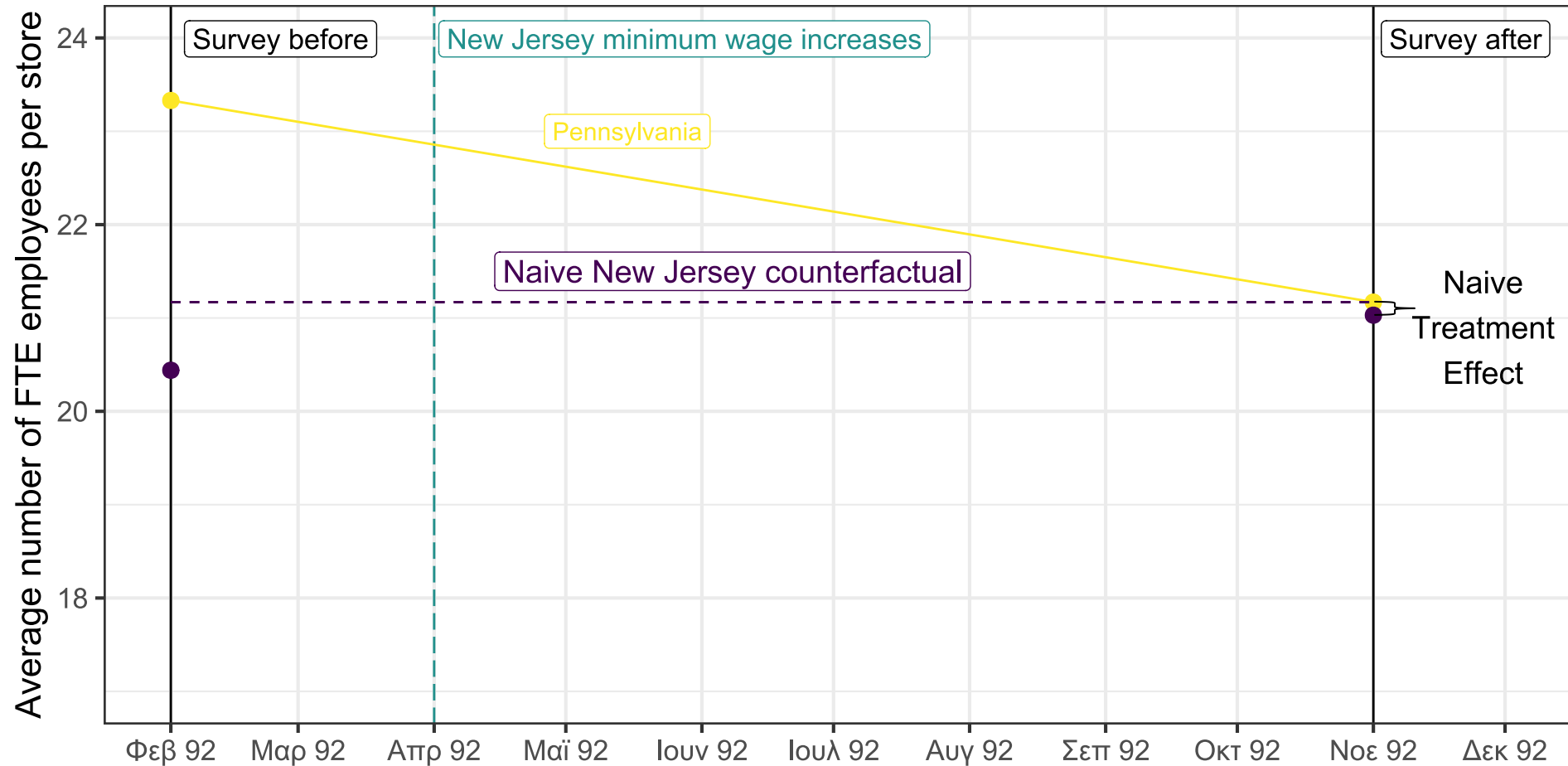
# What if we had done a naive after/before comparison?



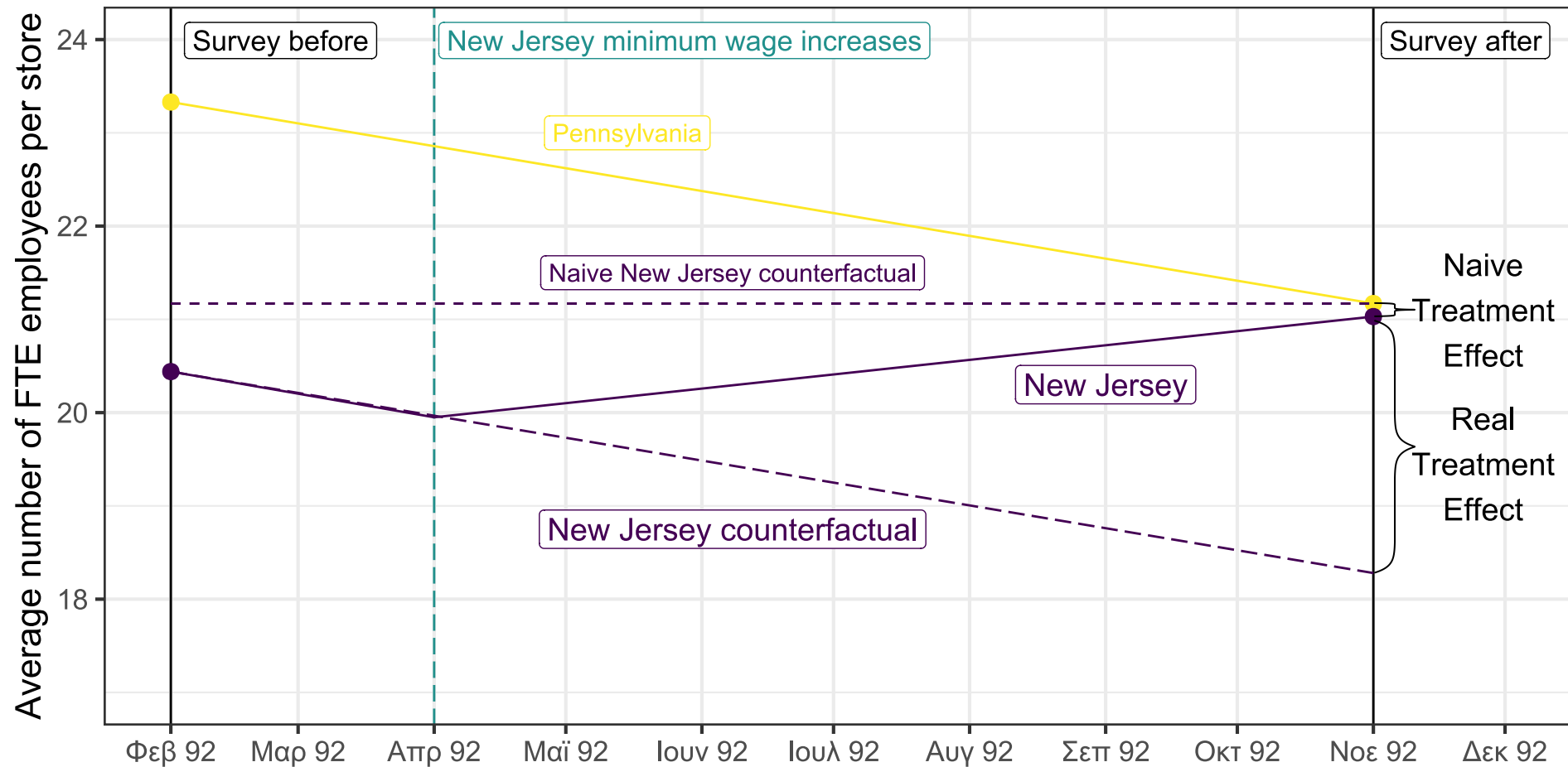
# What if we had done a naive after/before comparison?



# What if we had done a naive after NJ/PA comparison?



# What if we had done a naive after NJ/PA comparison?



# Estimation

# DiD in Regression Form

- In practice, DiD is usually estimated on more than 2 periods (4 observations)
- There are more data points before and after the policy change

3 ingredients:

1. **Treatment dummy variable:**  $TREAT_s$  where the  $s$  subscript reminds us that the treatment is at the state level
2. **Post-treatment periods dummy variables:**  $POST_t$  where the  $t$  subscript reminds us that this variable varies over time
3. **Interaction term between the two:**  $TREAT_s \times POST_t$  🖐 the *coefficient on this term is the DiD causal effect!*



# DiD in Regression Form

## Treatment dummy variable

$$TREAT_s = \begin{cases} 0 & \text{if } s = \text{Pennsylvania} \\ 1 & \text{if } s = \text{New Jersey} \end{cases}$$

## Post-treatment periods dummy variable

$$POST_t = \begin{cases} 0 & \text{if } t < \text{April 1, 1992} \\ 1 & \text{if } t \geq \text{April 1, 1992} \end{cases}$$

## Which observations correspond to $TREAT_s \times POST_t = 1$ ?

- Let's put all these ingredients together:

$$EMP_{st} = \alpha + \beta TREAT_s + \gamma POST_t + \delta(TREAT_s \times POST_t) + \varepsilon_{st}$$

- $\delta$ : causal effect of the minimum wage increase on employment





# Understanding the Regression

$$EMP_{st} = \alpha + \beta TREAT_s + \gamma POST_t + \delta(TREAT_s \times POST_t) + \varepsilon_{st}$$

We have the following:

$$\mathbb{E}(EMP_{st} \mid TREAT_s = 0, POST_t = 0) = \alpha$$

$$\mathbb{E}(EMP_{st} \mid TREAT_s = 0, POST_t = 1) = \alpha + \gamma$$

$$\mathbb{E}(EMP_{st} \mid TREAT_s = 1, POST_t = 0) = \alpha + \beta$$

$$\mathbb{E}(EMP_{st} \mid TREAT_s = 1, POST_t = 1) = \alpha + \beta + \gamma + \delta$$

$$\begin{aligned} & [\mathbb{E}(EMP_{st} \mid TREAT_s = 1, POST_t = 1) - \mathbb{E}(EMP_{st} \mid TREAT_s = 1, POST_t = 0)] - \\ & [\mathbb{E}(EMP_{st} \mid TREAT_s = 0, POST_t = 1) - \mathbb{E}(EMP_{st} \mid TREAT_s = 0, POST_t = 0)] = \delta \end{aligned}$$



# Understanding the Regression

$$EMP_{st} = \alpha + \beta TREAT_s + \gamma POST_t + \delta(TREAT_s \times POST_t) + \varepsilon_{st}$$

In table form:

	Pre mean	Post mean	$\Delta(\text{post} - \text{pre})$
Pennsylvania (PA)	$\alpha$	$\alpha + \gamma$	$\gamma$
New Jersey (NJ)	$\alpha + \beta$	$\alpha + \beta + \gamma + \delta$	$\gamma + \delta$
$\Delta(\text{NJ} - \text{PA})$	$\beta$	$\beta + \delta$	$\delta$

This table generalizes to other settings by substituting *Pennsylvania* with *Control* and *New Jersey* with *Treatment*



## Task 2 (10 minutes)

1. Create a dummy variable, *treat*, equal to FALSE if state is Pennsylvania and TRUE if New Jersey.
2. Create a dummy variable, *post*, equal to FALSE if observation is February 1992 and TRUE otherwise.
3. Estimate the following regression model. Do you obtain the same results as in slide 9?

$$empfte_{st} = \alpha + \beta treat_s + \gamma post_t + \delta(treat_s \times post_t) + \varepsilon_{st}$$



# Task 2: Solution

1. Create a dummy variable, `treat`, equal to FALSE if state is Pennsylvania and TRUE if New Jersey.

```
library(skimr)

ck1994 = ck1994 %>% mutate(treat = case_when(state == "Pennsylvania" ~ F,
                                             state == "New Jersey" ~ T))
```

2. Create a dummy variable, `post`, equal to FALSE if observation is February 1992 and TRUE otherwise.

```
library(skimr)

ck1994 = ck1994 %>% mutate(post = case_when(observation == "February 1992" ~ F,
                                             observation == "November 1992" ~ T))
```

3. Estimate the following regression model. Do you obtain the same results as in slide 9?

```
library(skimr)

did_reg = lm(empfte ~ treat*post, data = ck1994)

summary(did_reg)
```



# Identifying Assumptions

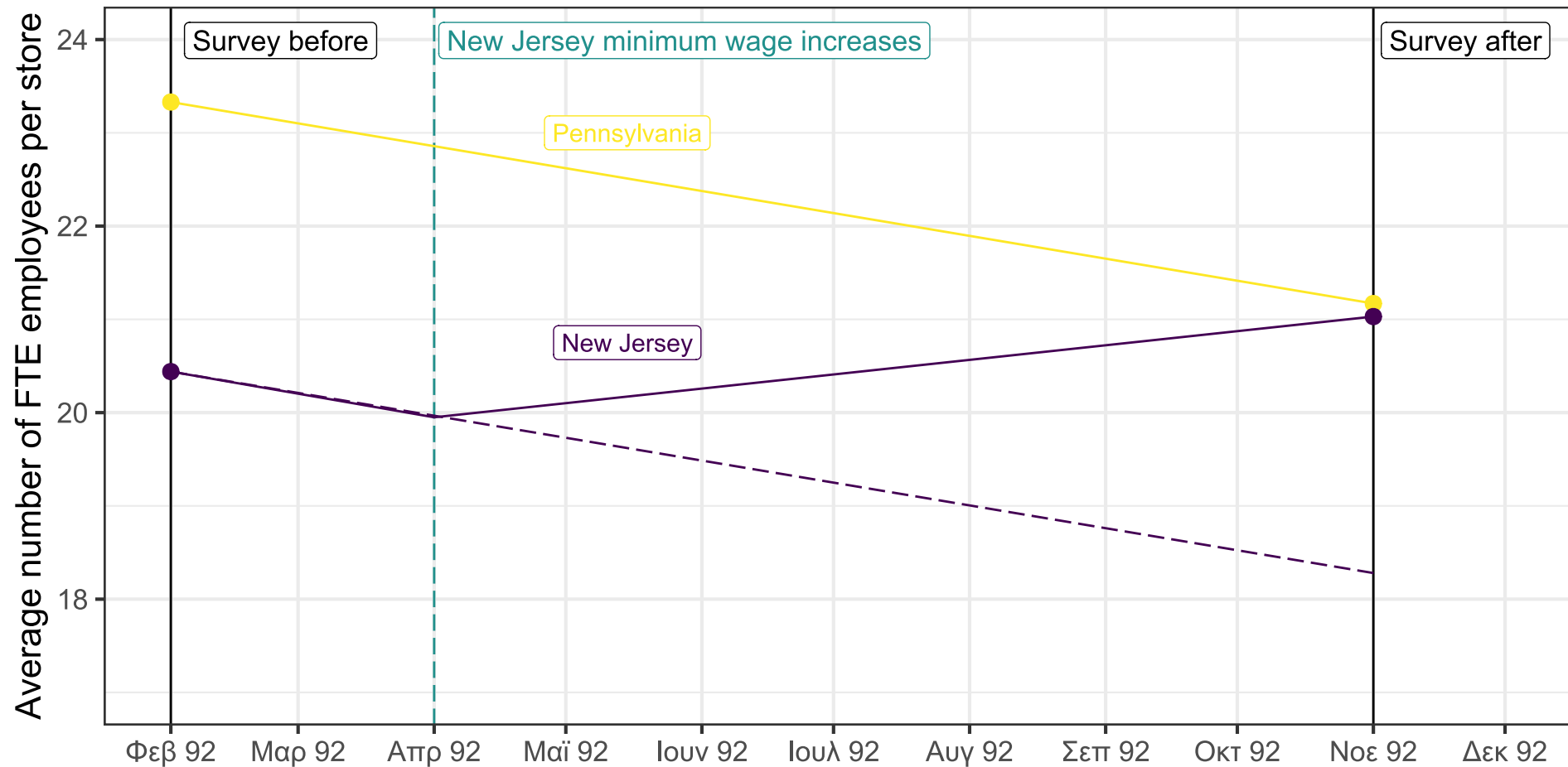
# DiD Crucial Assumption: Parallel Trends

**Common or parallel trends assumption:** absent any minimum wage increase, Pennsylvania's fast-food employment trend would have been what we should have expected to see in New Jersey.

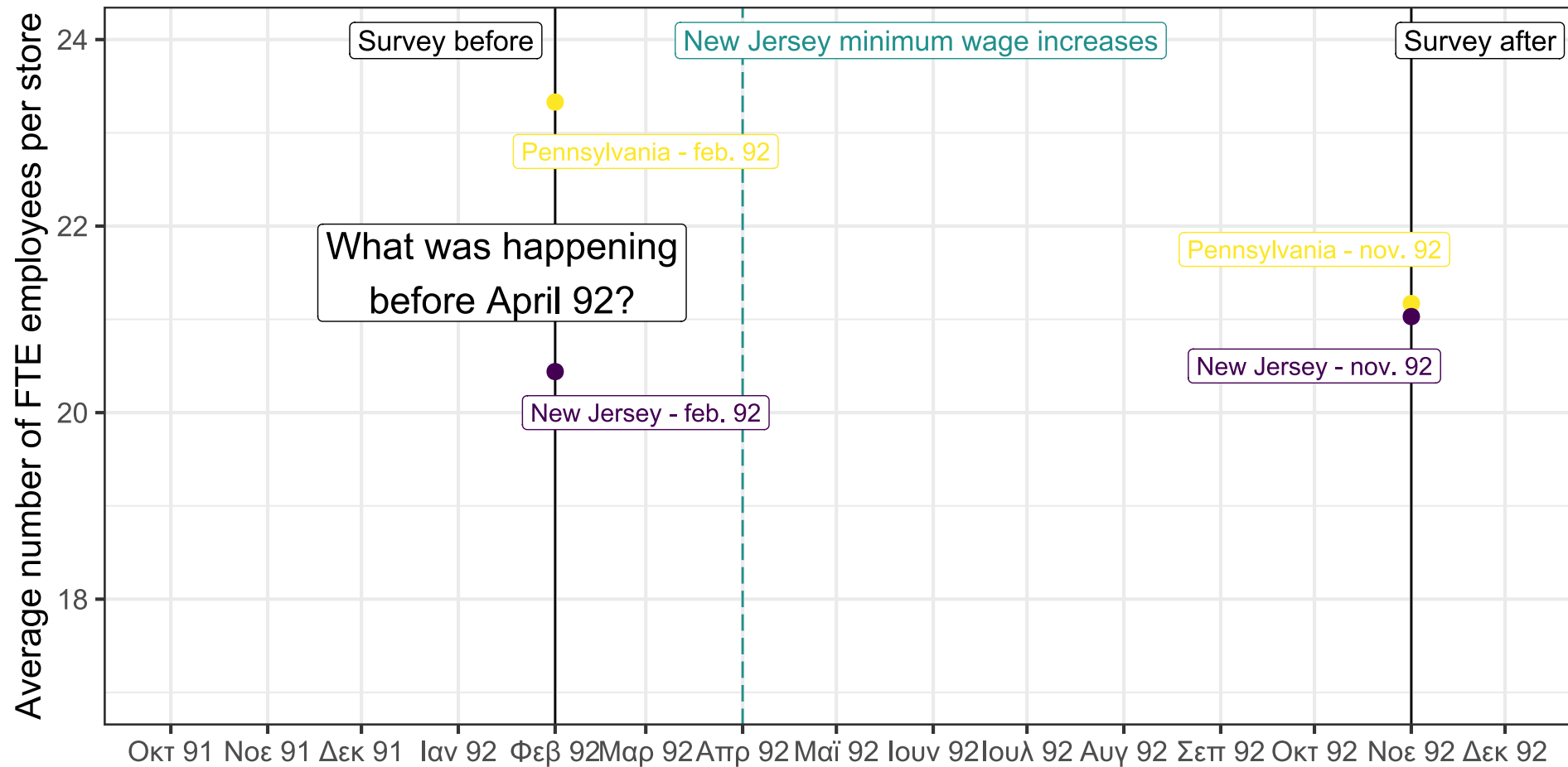
- This assumption states that Pennsylvania's fast-food employment trend between February and November 1992 provides a reliable counterfactual employment trend New Jersey's fast-food industry *would have experienced* had New Jersey not increased its minimum wage.
- Impossible to completely validate or invalidate this assumption.
- *Intuitive check:* compare trends before policy change (and after policy change if no expected medium-term effects)



# Parallel Trends: Graphically

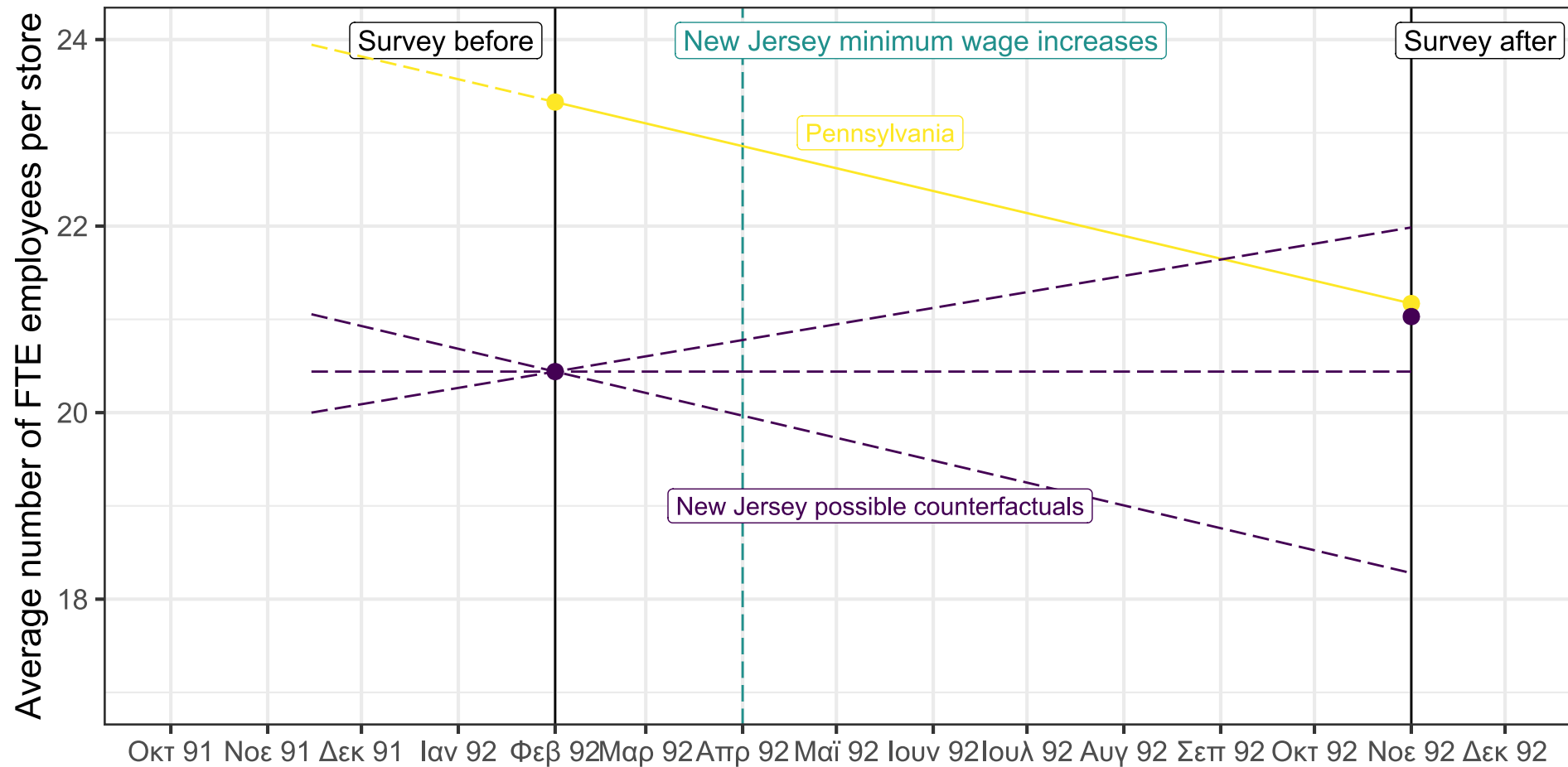


# Checking the parallel trends assumption

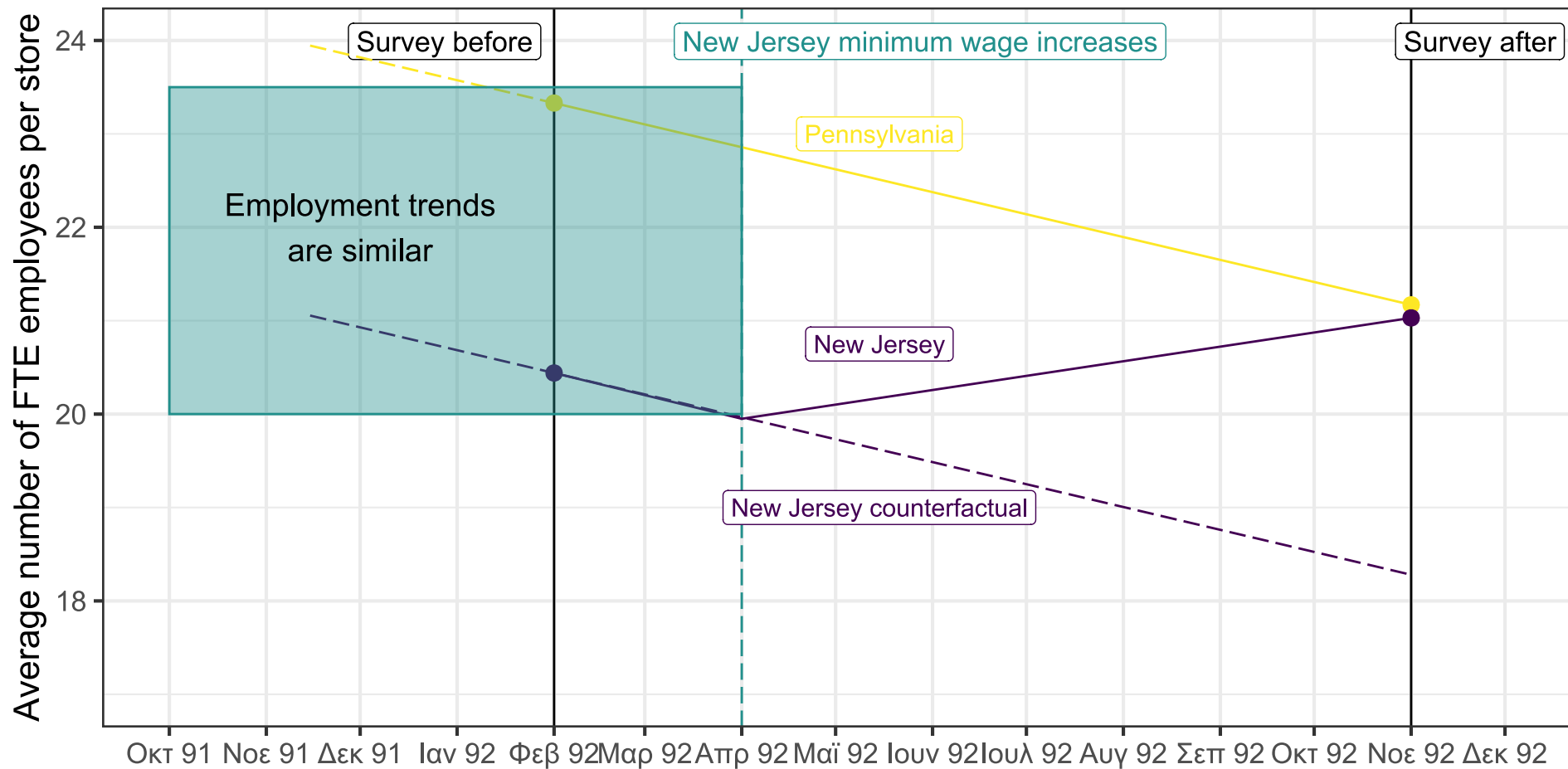




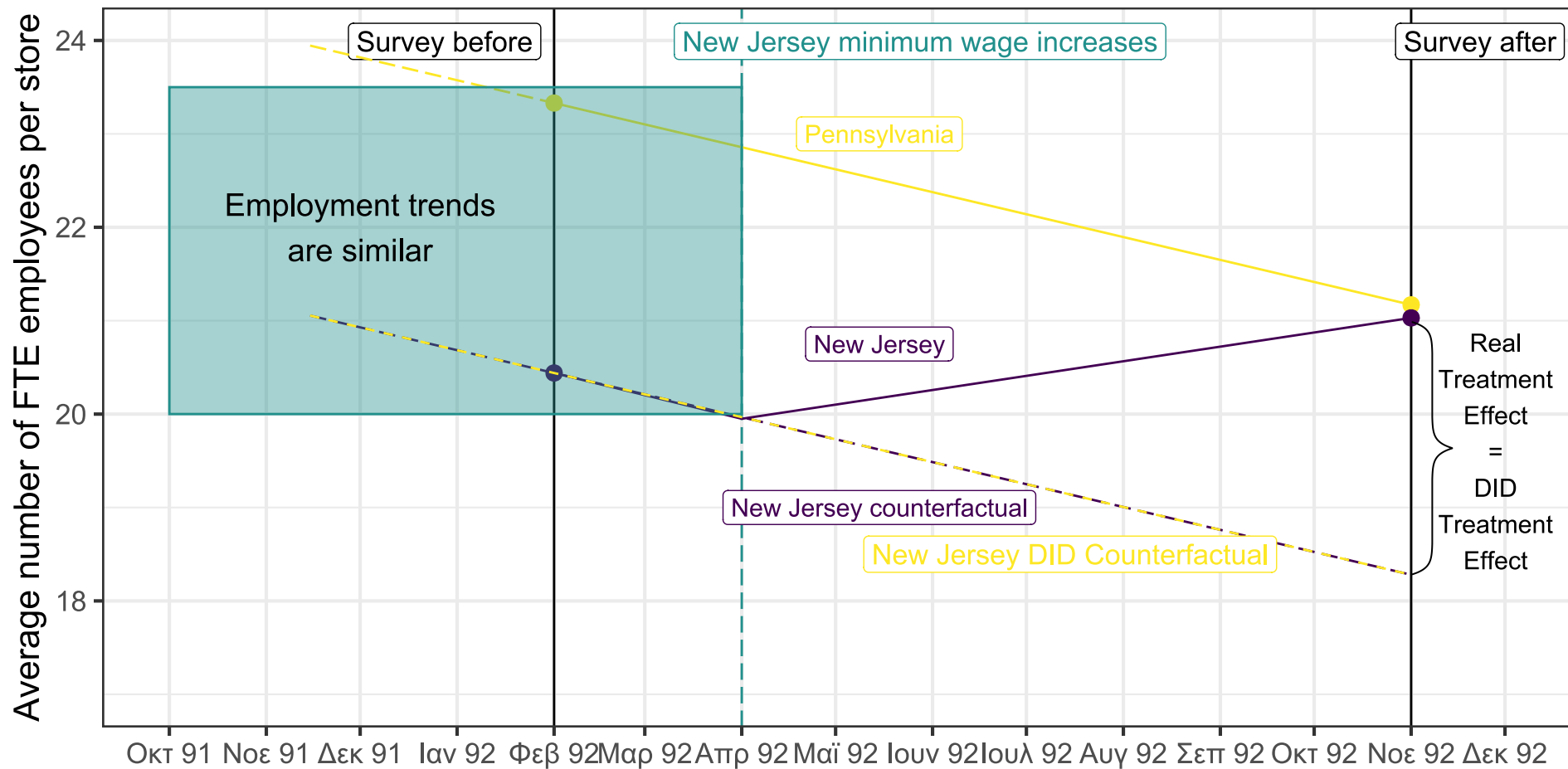
# Checking the parallel trends assumption



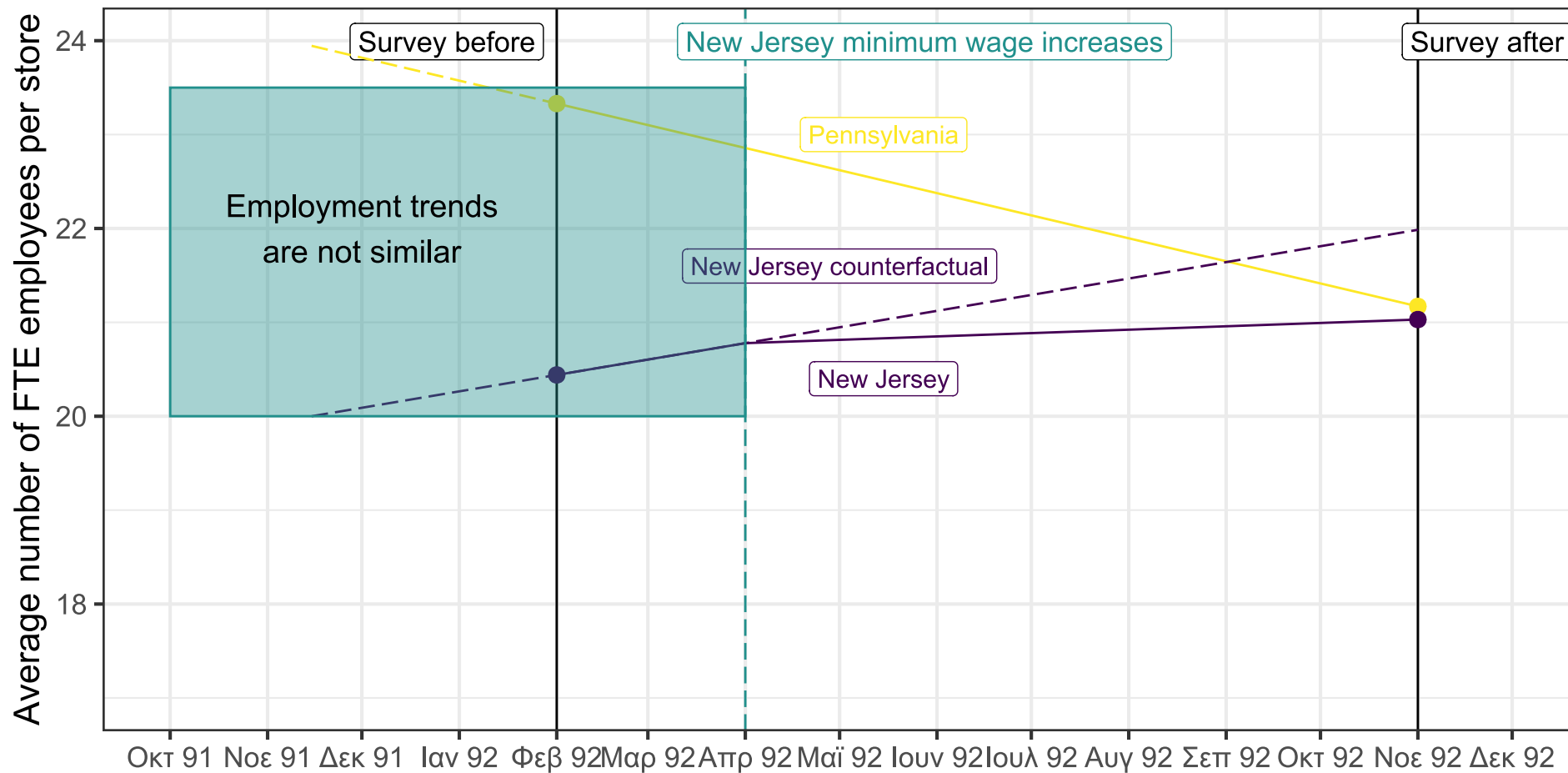
# Parallel trends assumption → Verified



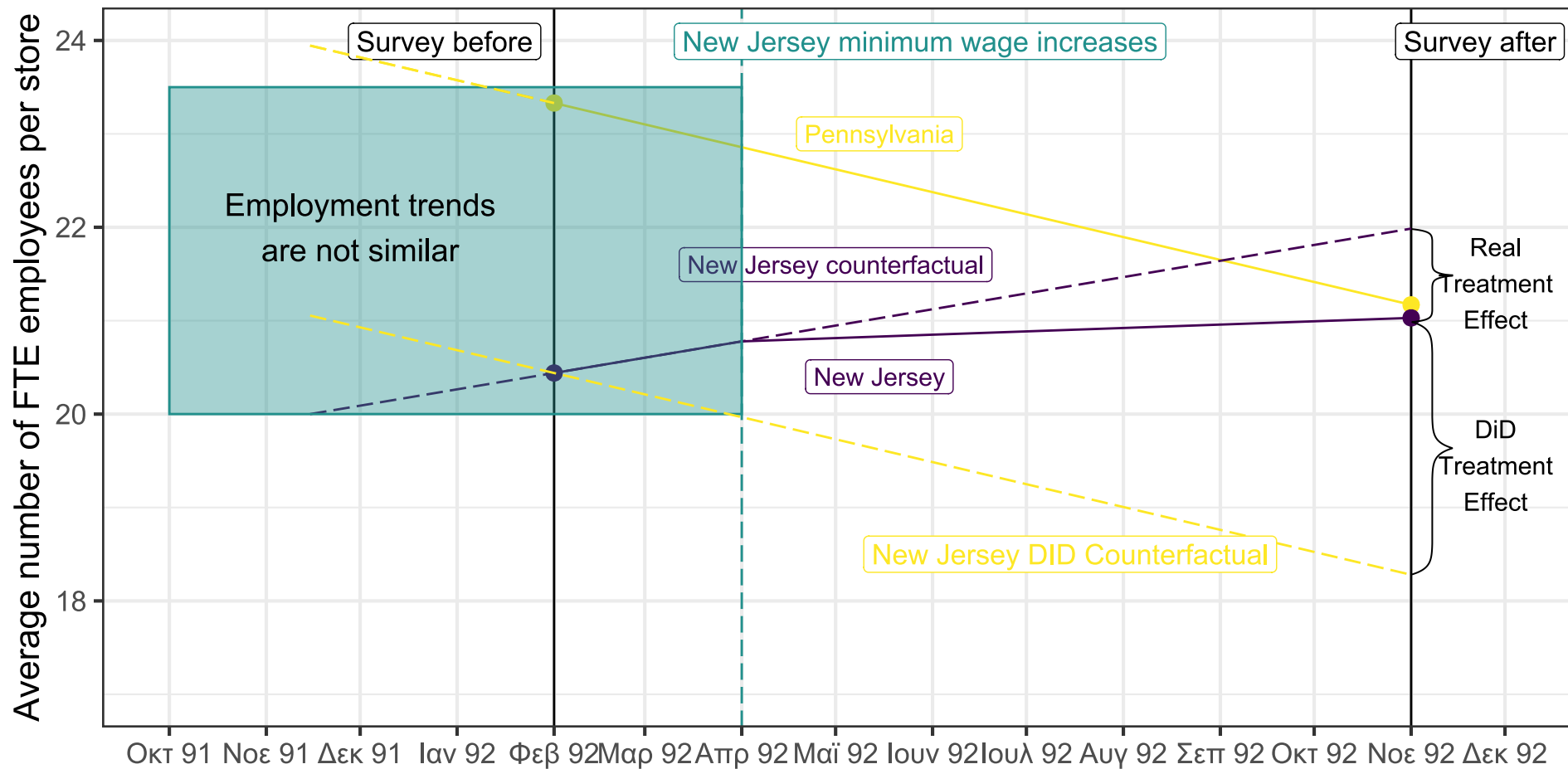
# Parallel trends assumption → Verified



# Parallel trends assumption → Not verified ❌

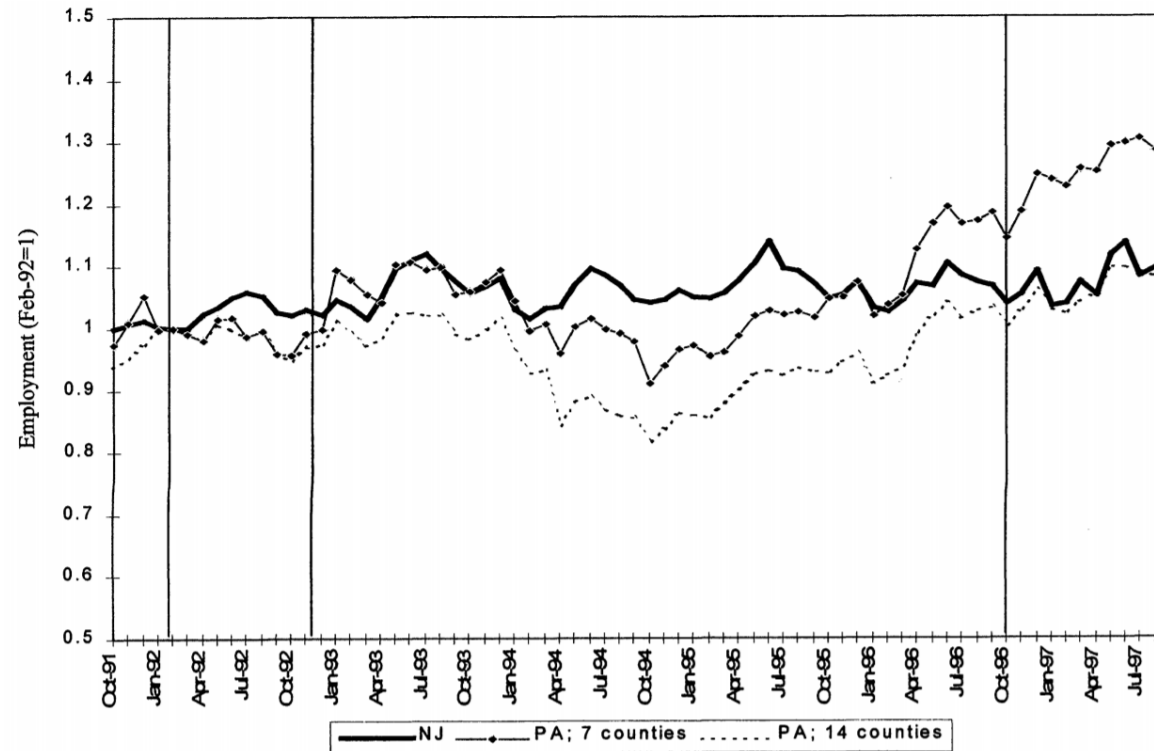


# Parallel trends assumption → Not verified ✖



# Parallel Trends Assumption: Card and Krueger (2000)

Here is the actual trends for Pennsylvania and New Jersey



- Is the common trend assumption likely to be verified?



# Parallel Trends Assumption: Formally

Let:

- $Y_{ist}^1$ : fast food employment at restaurant  $i$  in state  $s$  at time  $t$  if there is a high state MW;
- $Y_{ist}^0$ : fast food employment at restaurant  $i$  in state  $s$  at time  $t$  if there is a low state MW;

These are potential outcomes, you can only observe one of the two.

The key assumption underlying DiD estimation is that, in the no-treatment state, restaurant  $i$ 's outcome in state  $s$  at time  $t$  is given by:

$$\mathbb{E}[Y_{ist}^0 | s, t] = \gamma_s + \lambda_t$$

2 implicit assumptions:

1. **Selection bias**: relates to fixed state characteristics ( $\gamma$ )
2. **Time trend**: same time trend for treatment and control group ( $\lambda$ )



# Parallel Trends Assumption: Formally

Outcomes in the comparison group:

$$\mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Feb}] = \gamma_{PA} + \lambda_{Feb}$$

$$\mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Nov}] = \gamma_{PA} + \lambda_{Nov}$$

$$\begin{aligned} & \mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Nov}] - \mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Feb}] \\ &= \gamma_{PA} + \lambda_{Nov} - (\gamma_{PA} + \lambda_{Feb}) \\ &= \lambda_{Nov} - \lambda_{Feb} \end{aligned}$$





# Parallel Trends Assumption: Formally

Outcomes in the comparison group:

$$\mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Feb}] = \gamma_{PA} + \lambda_{Feb}$$

$$\mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Nov}] = \gamma_{PA} + \lambda_{Nov}$$

$$\begin{aligned} & \mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Nov}] - \mathbb{E}[Y_{ist} | s = \text{Pennsylvania}, t = \text{Feb}] \\ &= \gamma_{PA} + \lambda_{Nov} - (\gamma_{PA} + \lambda_{Feb}) \\ &= \underbrace{\lambda_{Nov} - \lambda_{Feb}}_{\text{time trend}} \end{aligned}$$

→ the comparison group allows to estimate the *time trend*.



# Parallel Trends Assumption: Formally

Let  $\delta$  denote the true impact of the minimum wage increase:

$$\mathbb{E}[Y_{ist}^1 - Y_{ist}^0 | s, t] = \delta$$

Outcomes in the treatment group:

$$\mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Feb}] = \gamma_{NJ} + \lambda_{Feb}$$

$$\mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Nov}] = \gamma_{NJ} + \delta + \lambda_{Nov}$$

$$\begin{aligned} \mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Nov}] - \mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Feb}] \\ &= \gamma_{NJ} + \delta + \lambda_{Nov} - (\gamma_{NJ} + \lambda_{Feb}) \\ &= \delta + \lambda_{Nov} - \lambda_{Feb} \end{aligned}$$



# Parallel Trends Assumption: Formally

Let  $\delta$  denote the true impact of the minimum wage increase:

$$\mathbb{E}[Y_{ist}^1 - Y_{ist}^0 | s, t] = \delta$$

Outcomes in the treatment group:

$$\mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Feb}] = \gamma_{NJ} + \lambda_{Feb}$$

$$\mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Nov}] = \gamma_{NJ} + \delta + \lambda_{Nov}$$

$$\begin{aligned} & \mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Nov}] - \mathbb{E}[Y_{ist} | s = \text{New Jersey}, t = \text{Feb}] \\ &= \gamma_{NJ} + \delta + \lambda_{Nov} - (\gamma_{NJ} + \lambda_{Feb}) \\ &= \delta + \underbrace{\lambda_{Nov} - \lambda_{Feb}}_{\text{time trend}} \end{aligned}$$



# Parallel Trends Assumption: Formally

Therefore we have:

$$\mathbb{E}[Y_{ist}|s = \text{PA}, t = \text{Nov}] - \mathbb{E}[Y_{ist}|s = \text{PA}, t = \text{Feb}] = \underbrace{\lambda_{\text{Nov}} - \lambda_{\text{Feb}}}_{\text{time trend}}$$

$$\mathbb{E}[Y_{ist}|s = \text{NJ}, t = \text{Nov}] - \mathbb{E}[Y_{ist}|s = \text{NJ}, t = \text{Feb}] = \delta + \underbrace{\lambda_{\text{Nov}} - \lambda_{\text{Feb}}}_{\text{time trend}}$$


$$\begin{aligned} DD &= \mathbb{E}[Y_{ist}|s = \text{NJ}, t = \text{Nov}] - \mathbb{E}[Y_{ist}|s = \text{NJ}, t = \text{Feb}] \\ &\quad - \left( \mathbb{E}[Y_{ist}|s = \text{PA}, t = \text{Nov}] - \mathbb{E}[Y_{ist}|s = \text{PA}, t = \text{Feb}] \right) \\ &= \delta + \lambda_{\text{Nov}} - \lambda_{\text{Feb}} - (\lambda_{\text{Nov}} - \lambda_{\text{Feb}}) \\ &= \delta \end{aligned}$$



END

---

 [florian.oswald@sciencespo.fr](mailto:florian.oswald@sciencespo.fr)

 Slides

 Book

 @ScPoEcon

 @ScPoEcon

---

