# ECON 251 Discussion Linear Probability Model (LPM) + Diff-in-Diff (DD)

Elird Haxhiu

Fall 2022

## Last discussion of the semester!

- Congratulations!
- Technically almost done...
- Still just the beginning
- Course evaluations





#### We are living in the future. Holy shit. #rstats #rvest #ChatGPT

\$			?	) Help	≡
Playgro	und				
Load a p	oreset		$\sim$	/	
Save	View code	Share		٥	

Ŷ

Scrape the table from this website in R https://www.lineups.com/nfl/nfl-targets

#### library(<mark>rvest</mark>)

url <- "https://www.lineups.com/<u>nfl/nfl-</u> targets"

webpage <- read\_html(url)

targets\_tbl <- html\_table(html\_nodes(webpage, "table")[[1]])

targets\_tbl





We are living in the future. #ChatGPT



v buckenham @v21

#### as I was saying: Oh no. Undergrads

meta.stackoverflow.com/questions/4218... because of this thin

14:18 \$ \$ Playground Load a preset... Playground View code Save Load a preset... Scrape the table from th Save https://www.lineups.com library(rvest) is it a problem fo normally distribu url <- "https://www.lineu targets" Yes, it is a proble not normally dist webpage <- read\_html(u the error term ha targets tbl <- html tabl if this assumptic "table")[[1]])

View c ALT estimates may b 35 Retweets precision. 0

The primary problem is that while the answers which ChatGPT produces have a high rate of being incorrect, they typically look like they might be good and the answers are very easy to produce. There are also many people trying out ChatGPT to create answers, without the expertise or willingness to verify that the answer is correct prior to posting. Because such answers are so easy to produce, a large number of people are posting a lot of answers. The volume of these answers (thousands) and the fact that the answers often require a detailed read by someone with at least some subject matter expertise in order to determine that the answer is actually bad has effectively swamped our volunteer-based quality curation infrastructure.

0

£

6:23 AM · Dec 5, 2022

10 Quote Tweets 377 Likes

1J



#ChatGPT

We are living in the future.



v buckenham @v21

#### as I was saying: Oh no. Undergrads

meta.stackoverflow.com/questions/4218... because of this thin

\$ Playground Load a preset... View code Save Scrape the table from t https://www.lineups.co library(rvest)

url <- "https://www.line targets"

webpage <- read html(

targets tbl <- html tab "table")[[1]])

targets\_tbl

14:18	
\$	
Playgro	ound
Load a p	preset
Save	View c
is it a pro normally	oblem fo / distribu
Yes, it is	a proble
not norm	hally dist
the error	term ha
if this as	sumptic 6:23
estimate	es may b
precisio	n

The primary problem is that while the answers which ChatGPT produces have a high rate of being incorrect, they typically look like they might



**Oliver Emberton** @oliveremberton

Oh and it's (currently) completely free, and requires nothing more than a web browser:

...

7

#### chat.openai.com/chat

0

so ex	8:46 PM · Dec	c 4, 2022			
is Int	55 Retweets	<b>3</b> Quote Tweets	<b>710</b> Likes		
)22					

£

5 Retweets 377 Likes

1J

0

# Outline

- 1. Linear probability model (LPM) for discrete outcomes
- 2. Review HW3 solutions
- 3. Difference-in-differences in practice

## Continuous Outcomes $Y \in \mathbb{R}$

• Independence + continuous Y gives usual "slope interpretation"

 $Y = \beta_0 + \beta_1 X + U$ 

$$E[Y|X] = E[\beta_0 + \beta_1 X + U|X]$$
  
=  $\beta_0 + \beta_1 X + E[U|X]$   
=  $\beta_0 + \beta_1 X$ 

$$\Rightarrow \beta_1 = \frac{\partial}{\partial X} E[Y|X]$$

# Binary Outcomes $Y \in \{0,1\}$

• Independence + binary Y gives "change in prob(Y=1)" interpretation

 $Y = \beta_0 + \beta_1 X + U$ 

$$E[Y|X] = P[Y = 1|X] \cdot 1 + P[Y = 1|X] \cdot 0$$
  
= P[Y = 1|X]

$$\Rightarrow \beta_1 = \frac{\partial}{\partial X} P[Y = 1|X]$$

# Linear Probability Model (LPM)

- OLS estimates of linear model with binary outcome
- LPM is nice because...
  - 1. Easy to estimate
  - 2. Easy to interpret
- LPM is problematic since
  - 1. Predicted values of outcome can be outside of [0,1] interval
  - 2. Does not make sense for X to change P[Y = 1|X] linearly

# Linear Probability Model (LPM)

- LPM is problematic because...
  - 1. Predicted values of outcome can be outside of [0,1] interval
  - 2. Does not make sense for X to change P[Y = 1|X] linearly
  - 3. Homoskedasticity is always violated

$$Var(Y|X) = E[Y^{2}|X] - E(Y|X)^{2}$$
  
=  $[P(Y = 1|X) \cdot 1^{2} + P(Y = 0|X) \cdot 0^{2}] - [P(Y = 1|X)]^{2}$   
=  $P(Y = 1|X) - P(Y = 1|X)^{2}$   
=  $P(Y = 1|X)[1 - P(Y = 1|X)]$ 

### Alternative Estimators = assume U distribution

• Probit = assume that  $U \sim N(0,1)$ 

$$P(Y = 1|X) = \Phi(\beta_0 + \beta_1 X)$$

where  $\Phi(u) \coloneqq P[U \le u] = F_U(u)$  denotes the standard normal CDF

• Logit = assume that *U* follows logistic distribution with PDF

$$f_U(u) = \frac{1}{1 + e^{-u}}$$

## Review HW3 solutions

# Outline

- 1. Linear probability model (LPM) for discrete outcomes
- 2. Review HW3 solutions
- 3. Difference-in-differences in practice

Difference-in-differences = compare *Y* change of units exposed to some policy *T* with *Y* change of unexposed

2 periods (before/after) and 2 groups (treated/control)

- $Y_{it} \coloneqq$  outcome of interest
- $P_t \coloneqq 1\{t \text{ is after treatment occurs}\}$
- $T_i \coloneqq 1\{i \text{ is treated/exposed}\}$

$$Y_{it} = \beta_0 + \beta_1 P_t + \beta_2 T_i + \beta_3 [P_t \cdot T_i] + U_i$$

	Before	After	After – Before
Control	$\beta_0$	$\beta_0 + \beta_1$	$eta_1$
Treated	$\beta_0 + \beta_2$	$\beta_0 + \beta_1 + \beta_2 + \beta_3$	$\beta_1 + \beta_3$
Treat – Control	$\beta_2$	$\beta_2 + \beta_3$	$\beta_3$





# Parallel Trends Assumption = exposed units Y without policy T would have changed like unexposed units Y

- PTA is an untestable assumption, just like OLS exogeneity or IV exogeneity
- However, if we have access to more data before policy, we can assess how likely it is to hold in practice... commonly known as "checking for pre-trends"
- One reason why people seem to like DD... visual check of identifying assumption!



#### Quasi-Market Competition in Public Service Provision: User Sorting and Cream-Skimming

#### Thorbjørn Sejr Guul<sup>\*,†</sup>, Ulrik Hvidman<sup>†</sup>, Hans Henrik Sievertsen<sup>‡</sup>

\*TrygFonden's Centre for Child Research; <sup>†</sup>Aarhus University; <sup>‡</sup>VIVE – The Danish Center for Social Science Research, IZA, University of Bristol

Address correspondence to the author at tsg@ps.au.dk.

#### Abstract

Quasi-markets that introduce choice and competition between public service providers are intended to improve quality and efficiency. This article demonstrates that quasi-market competition may also affect the distribution of users. First, we develop a simple theoretical framework that distinguishes between user sorting and cream-skimming as mechanisms through which quasi-markets may lead to high-ability users becoming more concentrated among one group of providers and low-ability users among a different group. Second, we empirically examine the impact of a nationwide quasimarket policy that introduced choice and activity-based budgeting into Danish public high schools. We exploit variation in the degree of competition that schools were exposed to, based on the concentration of providers within a geographical area. Using a differences-in-differences design-and register data containing the full population of students over a 9-year period (N = 207,394)—we show that the composition of students became more concentrated in terms of intake grade point average after the reform in high-competition areas relative to low-competition areas. These responses in high-competition regions appear to be driven both by changes in user sorting on the demand side and by cream-skimming behavior among public providers on the supply side.

#### In Stata! Code from Hans Henrik Sivertsen



#### . reg segregation Post Treated PostXTreated

Sourc	e	SS	df	MS	Number of obs	=	18
					· F(3, 14)	=	129.38
Mode	.006	163964	3	.002054655	Prob > F	=	0.0000
Residua	1 .000	222337	14	.000015881	. R-squared	=	0.9652
					<ul> <li>Adj R-squared</li> </ul>	=	0.9577
Tota	1 .006	386301	17	.000375665	Root MSE	=	.00399

segregation	Coefficient	Std. err.	t	P> t
Post	.002788	.0026733	1.04	0.315
Treated	.010465	.0028179	3.71	0.002
PostXTreated	.031551	.0037806	8.35	0.000
_cons	.01043	.0019926	5.23	0.000

#### In Stata! Code from Hans Henrik Sivertsen

