

Inferencia en Estudios Observacionales

2025-05-07

Intro

Hoy vamos a ver otros dos tópicos importantes relacionados con la inferencia en contextos no experimentales.

También conoceremos un debate que se dio en la literatura académica sobre métodos experimentales y no experimentales, empezando con Lalonde (1986).

Errores Agrupados II

Errores Agrupados II

En la clase pasada, discutimos el papel de errores agrupados en ensayos aleatorios.

En general, si tenemos un experimento a nivel de grupo, los errores tienden a ser mayores. Tenemos menos observaciones independientes de las que imaginamos.

Errores Agrupados II

Ese problema también ocurre en contextos no experimentales.

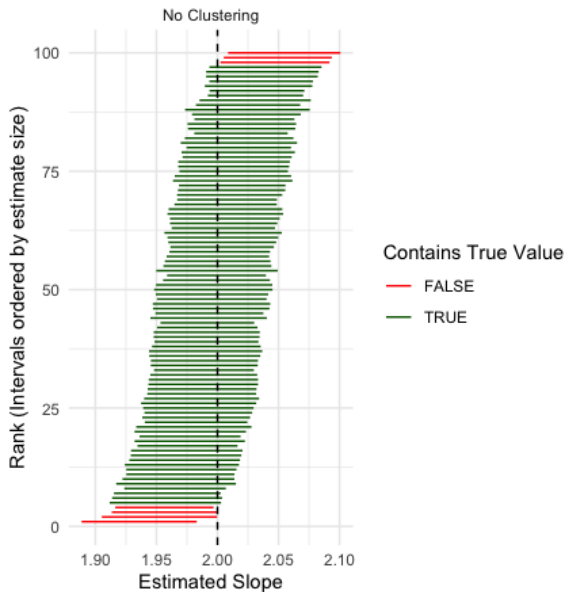
Es muy común que una muestra incluya una estructura jerárquica: diferentes personas en diferentes comunas, o regiones por ejemplo.

Si existen shocks o diferencias no observadas a un nivel más alto, tenemos un problema similar: los errores no son i.i.d., y por tanto el estimador normal de error estándar no es válido.

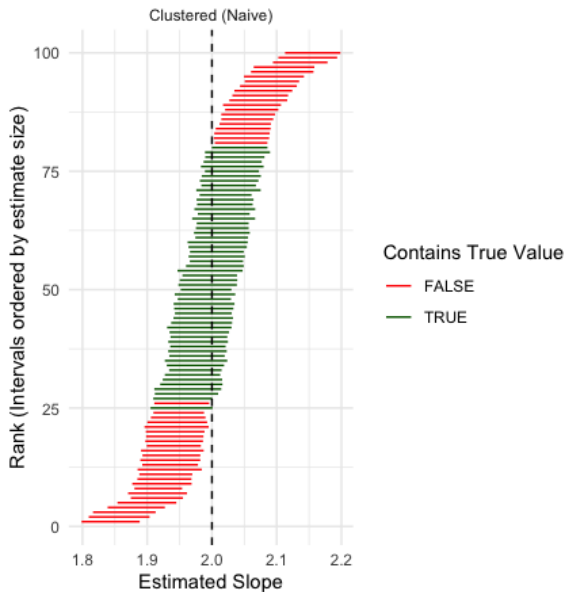
Normalmente eso lleva a subestimar los errores estándar.

Vamos a entender ese problema con una simulación.

95% Confidence Intervals by Method

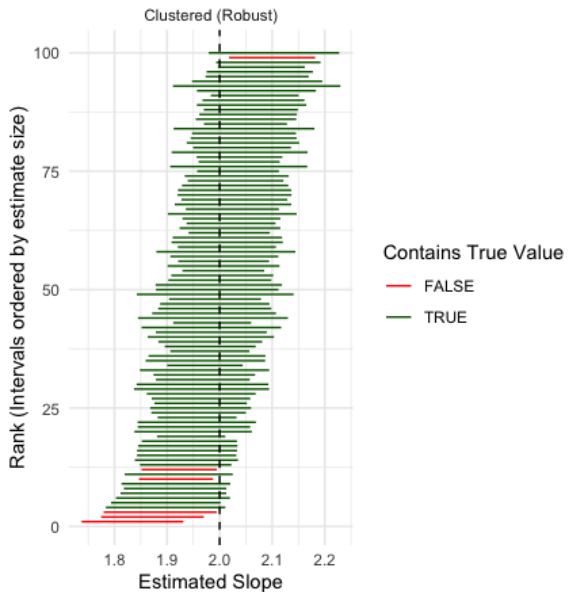


95% Confidence Intervals by Method



Si no llevamos en cuenta la estructura de los errores, la prueba no es válida: rechazamos mucho más que 5% de los casos.

95% Confidence Intervals by Method



Errores Agrupados II

Si no conocemos el proceso generador de datos, ¿cómo saber cuál es el nivel apropiado para agrupar?

- Es difícil de obtener una regla general.
- Si tenemos un argumento teórico para agrupar en un cierto nivel, lo hagamos.
- Otra opción es agrupar en el mismo nivel que la variación que se utiliza. Por ej:
 - Si estudiamos escuelas con diferentes practicas de ensino, agrupamos en la escuela.
 - Si las practicas son las mismas para escuelas en una misma comuna, o en un mismo barrio, se agrupa en ese nivel.
- Es particularmente importante cuando tenemos datos de un individuo o un grupo en diferentes momentos en el tiempo.
 - Como una unidad tiende a ser correlacionada con si misma en el tiempo, devemos agrupar en el nivel de la unidad.

Una limitación de este ajuste es que funciona bien cuando tenemos un número suficiente de grupos (al menos 50).

Si tenemos menos grupos que eso, necesitamos otro tipo de ajuste.

Bootstrap

Bootstrap es una técnica *extremadamente* útil para hacer inferencia. La ventaja es que el bootstrap funciona *casi siempre*. Y cuando no funciona, normalmente hay una versión un poco diferente que sí funciona.

La idea del bootstrap es que no vamos a intentar modelar la estructura de los datos analíticamente y derivar una fórmula para los errores estándar.

En vez de eso, vamos utilizar los propios datos.

Idealmente, queremos comprender como nuestra estimación se comporta si la muestra aleatoria que extraímos fuera otra muestra de la misma población.

Pero no tenemos otra muestra de la misma población.

Lo que vamos hacer es utilizar la muestra que tenemos para crear una otra muestra *posible* y ver qué pasa con el estimador. Hacemos eso muchas veces y construimos la distribución del estimador.

El truco es cómo construir esas muestras alternativas.

Podemos extrair una submuestra de nuestra muestra. Pero eso no va a funcionar bien porque vamos tener menos observaciones.

Si extraímos una muestra de mismo tamaño, sin remplazo, obtenemos una muestra idéntica. Tanpoco sirve.

Pero si extraímos una muestra de mismo tamaño **con remplazo**, vamos tener una muestra un poco diferente.

- Algunas unidades van a estar multiples veces, y otras no van estar.

El método de bootstrap para estimar el error estándar es:

- 1: Creamos una nueva muestra mediante extracción con remplazo de la muestra que tenemos.
- 2: Estimamos nuestro parámetro de interés.
- 3: Repetimos 1 y 2 varias veces (mil o diez mil).
- 4: El error estándar es simplemente la desviación estándar de los estimadores.

Un ejemplo muy simple:

Imagina que tenemos 4 individuos en la muestra:

i	Y	X
A	4	3
B	2	6
C	10	7
D	6	5

Si hacemos la regresión de Y en X, obtenemos $\beta = 0.97$.

Una muestra alternativa podría ser:

i	Y	X
C	10	7
C	10	7
D	6	5
B	2	6

En ese caso, estimamos $\beta = 2.55$.

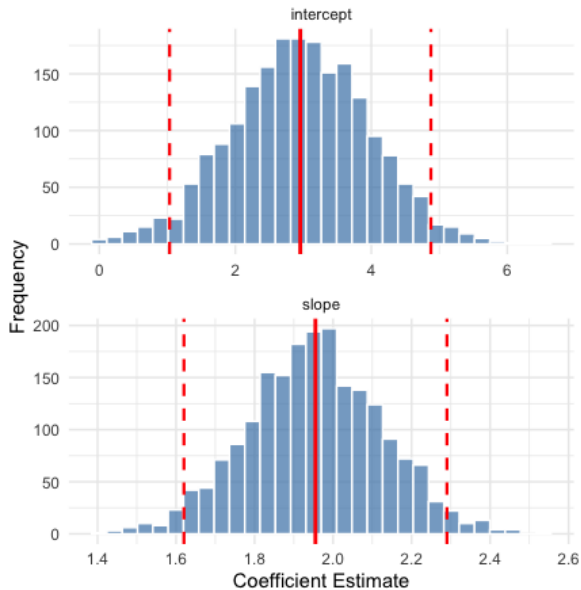
Otra posibilidad:

i	Y	X
A	4	3
A	4	3
B	2	6
B	2	6

Ahora tenemos $\beta = -0.667$.

Si repetimos ese proceso, podemos comprender la variación del estimador.
Vamos ver una demostración.

Bootstrap Distributions for Regression Coefficients



Podemos utilizar el bootstrap de una manera parecida para crear intervalos de confianza:

- 1: Creamos una nueva muestra mediante extracción con remplazo de la muestra que tenemos.
- 2: Estimamos nuestro parámetro de interés.
- 3: Repetimos 1 y 2 varias veces (mil o diez mil).
- 4: El intervalo de confianza de 95% es simplemente la parte donde están el 95% de las estimativas.

Para una regresión simple no hay una gran ventaja de usar el bootstrap, por que conocemos bien cómo estimar los errores estándar.

El bootstrap es más recomendado cuando es difícil de obtener una forma analítica para los errores. Por ejemplo:

- Imaginate que estamos haciendo una estimativa con ajustes como en subclasificación. La estimación incluye calcular varias diferencias y ponderarlas por diferentes pesos. No podemos utilizar los errores normales de una regresión.

Pero podemos escribir el estimador como un código, crear nuevas muestras y utilizar el bootstrap.

Podemos cambiar el bootstrap básico para lidiar con diversos problemas.

Para incorporar errores agrupados con un procedimiento de *cluster bootstrap*.

Eso significa que, en vez de extrair nuevas unidades de la muestra con remplazo, una a una, extraímos los grupos enteros.

Por ejemplo: en vez de recrear la muestra por alumno, la recreamos por escuela.

Un de los pocos casos en que el bootstrap **no funciona bien** es con emparejamiento con vecino más cercano.

La razón es que el bootstrap funciona basado en cambios “suaves” en el estimador: si la muestra cambia un poco, queremos que el estimador cambie un poco.

El caso de emparejamiento con vecino más cercano es uno en que el estimador no cambia para cambios pequeños en la muestra, por que el vecino más cercano tiende a ser todavía el más cercano bajo un poco de variación.

Por esa razón, otros estimadores de emparejamiento son más utilizados, como el kernel (que pondera muchos vecinos cercanos) o el de ponderación de probabilidad inversa.

¿Podemos Confiar en Métodos No Experimentales?

En su paper **Evaluating the Econometric Evaluations of Training Programs with Experimental Data**, Lalonde se preguntó si es posible utilizar métodos no experimentales para replicar los resultados de ensayos aleatorizados.

Para ello, utilizó datos de un ensayo aleatorizado de capacitación laboral que ocurrió en los años 70.

- Podemos utilizar el grupo de control experimental para obtener resultados sin riesgo de sesgo.
- Después construimos un grupo de comparación con datos observacionales y comparamos los resultados.

El *National Supported Work Demonstration* (NSW) fue un programa que empleó personas vulnerables por entre 9 y 18 meses.

Incluyó mujeres en un programa llamado AFDC (*Aid for Families with Dependent Children*), y hombres. Ex-criminales, exadictos, personas que no terminaron el *high school*.

TABLE 1 — THE SAMPLE MEANS AND STANDARD DEVIATIONS OF
PRE-TRAINING EARNINGS AND OTHER CHARACTERISTICS FOR
THE NSW AFDC AND MALE PARTICIPANTS

Variable	Full National Supported Work Sample			
	AFDC Participants		Male Participants	
	Treatments	Controls	Treatments	Controls
Age	33.37 (7.43)	33.63 (7.18)	24.49 (6.58)	23.99 (6.54)
Years of School	10.30 (1.92)	10.27 (2.00)	10.17 (1.75)	10.17 (1.76)
Proportion High School Dropouts	.70 (.46)	.69 (.46)	.79 (.41)	.80 (.40)
Proportion Married	.02 (.15)	.04 (.20)	.14 (.35)	.13 (.35)
Proportion Black	.84 (.37)	.82 (.39)	.76 (.43)	.75 (.43)
Proportion Hispanic	.12 (.32)	.13 (.33)	.12 (.33)	.14 (.35)
Real Earnings	\$393	\$395	1472	1558
1 year Before	(1,203)	(1,149)	(2656)	(2961)
Training	[43]	[41]	[58]	[63]
Real Earnings	\$854	\$894	2860	3030
2 years Before	(2,087)	(2,240)	(4729)	(5293)
Training	[74]	[79]	[104]	[113]
Hours Worked	90	92	278	274
1 year Before	(251)	(253)	(466)	(458)
Training	[9]	[9]	[10]	[10]
Hours Worked	186	188	458	469
2 years Before	(434)	(450)	(654)	(689)
Training	[15]	[16]	[14]	[15]
Month of Assignment (Jan. 78 = 0)	-12.26 (4.30)	-12.30 (4.23)	-16.08 (5.97)	-15.91 (5.89)
Number of Observations	800	802	2083	2193

Note: The numbers shown in parentheses are the standard deviations and those in the square brackets are the standard errors.

La tabla de equilibrio parece buena.

Las diferencias son pequeñas y ninguna es estadísticamente significativa.

Ahora vamos mirar los grupos de control no experimentales.

Vienen de dos encuestas: el *Panel Study of Income Dynamics* (PSID) y el *Current Population Survey* (CPS).

Con ellas, Lalonde crea 6 grupos de comparación:

- PSID 1: hombres jefes de familia en PSID entre 1975 y 1978, con menos de 55 años y no jubilados.
- PSID 2: ... que no trabajaban en 76.
- PSID 3: ... que no trabajaban en 75 y 76.
- CPS 1: hombres en CPS con menos de 55 años.
- CPS 2: ... que no trabajaban en 76.
- CPS 3: ... y bajo el nivel de pobreza en 75.

TABLE 3—ANNUAL EARNINGS OF NSW MALE TREATMENTS, CONTROLS, AND SIX CANDIDATE COMPARISON GROUPS FROM THE *PSID* AND *CPS-SSA*

Year	Treatments	Controls	Comparison Group ^{a,b}					
			<i>PSID</i> -1	<i>PSID</i> -2	<i>PSID</i> -3	<i>CPS-SSA</i> -1	<i>CPS-SSA</i> -2	<i>CPS-SSA</i> -3
1975	\$3,066 (283)	\$3,027 (252)	19,056 ^a (272)	7,569 (568)	2,611 (492)	13,650 (73)	7,387 (206)	2,729 (197)
1976	\$4,035 (215)	\$2,121 (163)	20,267 (296)	6,152 (601)	3,191 (609)	14,579 (75)	6,390 (187)	3,863 (267)
1977	\$6,335 (376)	\$3,403 (228)	20,898 (296)	7,985 (621)	3,981 (594)	15,046 (76)	9,305 (225)	6,399 (398)
1978	\$5,976 (402)	\$5,090 (227)	21,542 (311)	9,996 (703)	5,279 (686)	14,846 (76)	10,071 (241)	7,277 (431)
Number of Observations	297	425	2,493	253	128	15,992	1,283	305

^aThe Comparison Groups are defined as follows: *PSID*-1: All male household heads continuously from 1975 through 1978, who were less than 55-years-old and did not classify themselves as retired in 1975; *PSID*-2: Selects from the *PSID*-1 group all men who were not working when surveyed in the spring of 1976; *PSID*-3: Selects from the *PSID*-1 group all men who were not working when surveyed in either spring of 1975 or 1976; *CPS-SSA*-1: All males based on Westat's criteria, except those over 55-years-old; *CPS-SSA*-2: Selects from *CPS-SSA*-1 all males who were not working when surveyed in March 1976; *CPS-SSA*-3: Selects from the *CPS-SSA*-1 unemployed males in 1976 whose income in 1975 was below the poverty level.

Con estos grupos, Lalonde testa varios estimadores:

- La diferencia de medias de ingresos
- La diferencia de medias controlando por edad, edad al cuadrado, educación, si uno terminó el high school, y raza.
- La diferencia de *la evolución de los ingresos*.
- La diferencia de la evolución de los ingresos, con controles.
- La diferencia de medias de ingresos controlando por el nivel de ingresos anterior
- La diferencia de medias de ingresos controlando por el nivel de ingresos anterior y otros controles
- La diferencia de medias de ingresos con todos los controles anteriores posibles.

TABLE 5—EARNINGS COMPARISONS AND ESTIMATED TRAINING EFFECTS FOR THE NSW MALE PARTICIPANTS USING COMPARISON GROUPS FROM THE *PSID* AND THE *CPS-SSA*^{a,b}

Name of Comparison Group ^d	Comparison Group Earnings Growth 1975–78 (1)	NSW Treatment Earnings Less Comparison Group Earnings				Difference in Differences: Difference in Earnings Growth 1975–78 Treatments Less Comparisons		Unrestricted Difference in Differences: Quasi Difference in Earnings Growth 1975–78		Controlling for All Observed Variables and Pre-Training Earnings (10)
		Pre-Training Year, 1975		Post-Training Year, 1978		Without Age (6)	With Age (7)	Unad-justed (8)	Ad-justed ^c (9)	
		Unad-justed (2)	Ad-justed ^c (3)	Unad-justed (4)	Ad-justed ^c (5)					
Controls	\$2,063 (325)	\$39 (383)	\$ – 21 (378)	\$886 (476)	\$798 (472)	\$847 (560)	\$856 (558)	\$897 (467)	\$802 (467)	\$662 (506)
<i>PSID</i> -1	\$2,043 (237)	– \$15,997 (795)	– \$7,624 (851)	– \$15,578 (913)	– \$8,067 (990)	\$425 (650)	– \$749 (692)	– \$2,380 (680)	– \$2,119 (746)	– \$1,228 (896)
<i>PSID</i> -2	\$6,071 (637)	– \$4,503 (608)	– \$3,669 (757)	– \$4,020 (781)	– \$3,482 (935)	\$484 (738)	– \$650 (850)	– \$1,364 (729)	– \$1,694 (878)	– \$792 (1024)
<i>PSID</i> -3	(\$3,322 (780)	\$455 (539)	\$455 (704)	\$697 (760)	– \$509 (967)	\$242 (884)	– \$1,325 (1078)	\$629 (757)	– \$552 (967)	\$397 (1103)
<i>CPS-SSA</i> -1	\$1,196 (61)	– \$10,585 (539)	– \$4,654 (509)	– \$8,870 (562)	– \$4,416 (557)	\$1,714 (452)	\$195 (441)	– \$1,543 (426)	– \$1,102 (450)	– \$805 (484)
<i>CPS-SSA</i> -2	\$2,684 (229)	– \$4,321 (450)	– \$1,824 (535)	– \$4,095 (537)	– \$1,675 (672)	\$226 (539)	– \$488 (530)	– \$1,850 (497)	– \$782 (621)	– \$319 (761)
<i>CPS-SSA</i> -3	\$4,548 (409)	\$337 (343)	\$878 (447)	– \$1,300 (590)	\$224 (766)	– \$1,637 (631)	– \$1,388 (655)	– \$1,396 (582)	\$17 (761)	\$1,466 (984)

^aThe columns above present the estimated training effect for each econometric model and comparison group. The dependent variable is earnings in 1978. Based on the experimental data an unbiased estimate of the impact of training presented in col. 4 is \$886. The first three columns present the difference between each comparison group's 1975 and 1978 earnings and the difference between the pre-training earnings of each comparison group and the NSW treatments.

^bEstimates are in 1982 dollars. The numbers in parentheses are the standard errors.

^cThe exogenous variables used in the regression adjusted equations are age, age squared, years of schooling, high school dropout status, and race.

^dSee Table 3 for definitions of the comparison groups.

Los resultados para el grupo experimental son *muy similares* dentre todos los estimadores.

Los resultados no experimentales son altamente variables. Mismo los casos con más restricciones y controles no puede replicar los resultados experimentales.

En general, un resultado *preocupante* para el uso de evaluaciones no experimentales.

Ese problema fue reanalizado por Dehejia y Wahba en su paper *Causal effects in nonexperimental studies: reevaluating the evaluation of training programs*.

Su principal contribución es de utilizar métodos de emparejamiento y tener atención al soporte común (*common support*).

- Más cuidado en elegir controles más comparables (no solo en promedio).

Lo primero que hacen es replicar los resultados de Lalonde, y utilizar una muestra con más información sobre los ingresos pasados.

- Excluyen los que no tienen información de ingresos en 74.
- También excluyen los que entraron en el programa más tarde y tenían ingresos cuando el programa empezó.

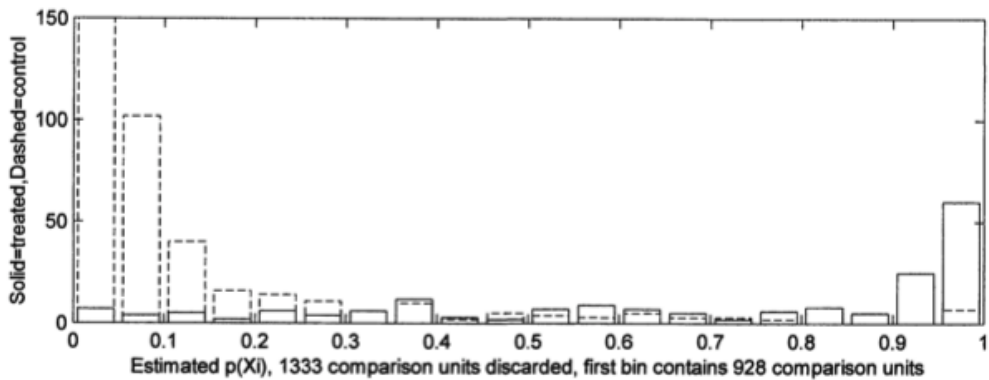
Ellos muestran que los resultados son un poco diferentes, pero la conclusión no cambia solo por la diferencia de muestra.

Comparison group	A. Lalonde's original sample					C. RE74 subsample (results use RE74)				
	NSW treatment earnings less comparison group earnings 1978		Unrestricted differences in differences: Quasi-difference in earnings growth 1975-1978		Controlling for all variables ^f	NSW treatment earnings less comparison group earnings 1978		Unrestricted differences in differences: Quasi-difference in earnings growth 1975-1978		Controlling for all variables ^f
	Unadjusted ^b	Adjusted ^c	Unadjusted ^d	Adjusted ^e		Unadjusted ^b	Adjusted ^c	Unadjusted ^d	Adjusted ^e	
	(1)	(2)	(3)	(4)		(1)	(2)	(3)	(4)	
NSW	886 (472)	798 (472)	879 (467)	802 (468)	820 (468)	1,794 (633)	1,688 (636)	1,750 (632)	1,672 (638)	1,655 (640)
PSID-1	-15,578 (913)	-8,067 (990)	-2,380 (680)	-2,119 (746)	-1,844 (762)	-15,205 (1155)	-879 (931)	-582 (841)	218 (866)	731 (886)
PSID-2	-4,020 (781)	-3,482 (935)	-1,364 (729)	-1,694 (878)	-1,876 (885)	-3,647 (960)	94 (1042)	721 (886)	907 (1004)	683 (1028)
PSID-3	697 (760)	-509 (967)	629 (757)	-552 (967)	-576 (968)	1,070 (900)	821 (1100)	1,370 (897)	822 (1101)	825 (1104)
CPS-1	-8,870 (562)	-4,416 (577)	-1,543 (426)	-1,102 (450)	-987 (452)	-8,498 (712)	-8 (572)	-78 (537)	739 (547)	972 (550)
CPS-2	-4,195 (533)	-2,341 (620)	-1,649 (459)	-1,129 (551)	-1,149 (551)	-3,822 (671)	615 (672)	-263 (574)	879 (654)	790 (658)
CPS-3	-1,008 (539)	-1 (681)	-1,204 (532)	-263 (677)	-234 (675)	-635 (657)	1,270 (798)	-91 (641)	1,326 (796)	1,326 (798)

En seguida, estiman el propensity score para cada grupo de comparación y eliminan controles fuera del soporte común.

Los modelos para estimar el PS son modelos logit que incluyen edad, edad al cuadrado, educación, educación al cuadrado, casado, no tiene diploma, raza, ingresos en 74 y 75 y sus cuadrados, desempleo en 74 y 75. Además, también incluyen:

- PSID: desempleo en 74 * negro.
- CPS: educación * ingresos en 74, y edad al cubo.



Con esta función de PS, ellos utilizan 7 estimadores:

- La diferencia simple
- La diferencia con controles
- La diferencia controlando por el propensity score (y su cuadrado)
- Con estratificación por valores del PS (subclasificación)
- Con estratificación por valores del PS (subclasificación) más controles
- Con matching en el vecino más próximo por el PS
- Con matching en el vecino más próximo por el PS más controles

Table 3. Estimated Training Effects for the NSW Male Participants Using Comparison Groups From PSID and CPS

	NSW earnings less comparison group earnings		NSW treatment earnings less comparison group earnings, conditional on the estimated propensity score					
	(1) Unadjusted	(2) Adjusted ^a	Quadratic in score ^b (3)	Stratifying on the score			Matching on the score	
				(4) Unadjusted	(5) Adjusted	(6) Observations ^c	(7) Unadjusted	(8) Adjusted ^d
NSW	1,794 (633)	1,672 (638)						
PSID-1 ^e	−15,205 (1,154)	731 (886)	294 (1,389)	1,608 (1,571)	1,494 (1,581)	1,255	1,691 (2,209)	1,473 (809)
PSID-2 ^f	−3,647 (959)	683 (1,028)	496 (1,193)	2,220 (1,768)	2,235 (1,793)	389	1,455 (2,303)	1,480 (808)
PSID-3 ^f	1,069 (899)	825 (1,104)	647 (1,383)	2,321 (1,994)	1,870 (2,002)	247	2,120 (2,335)	1,549 (826)
CPS-1 ^g	−8,498 (712)	972 (550)	1,117 (747)	1,713 (1,115)	1,774 (1,152)	4,117	1,582 (1,069)	1,616 (751)
CPS-2 ^g	−3,822 (670)	790 (658)	505 (847)	1,543 (1,461)	1,622 (1,346)	1,493	1,788 (1,205)	1,563 (753)
CPS-3 ^g	−635 (657)	1,326 (798)	556 (951)	1,252 (1,617)	2,219 (2,082)	514	587 (1,496)	662 (776)

Los estimadores con propensity score tienen resultados mucho mejores.

- Un detalle es que sus errores son muy grandes, entonces así mismo no estaríamos seguros del efecto.

En general, un resultado muy positivo para los métodos no experimentales.

- Motivó bastante interés en esos métodos.

Pero, hubo una re-re-evaluación.

Smith and Todd revisit ese debate en su paper: **Does matching overcome LaLonde's critique of non-experimental estimators?**

Esencialmente, su contribución es testar si los métodos de propensity score son *siempre* mejores, o si ciertas elecciones de WB influenciaran sus resultados.

Su primero resultado importante es aplicar los métodos de PS de WB para (1) la muestra original de Lalonde, y (2) una muestra sin una restricción que WB utilizaran (eliminando individuos que empezaran más tarde con ingresos positivos antes del programa).

Eses resultados no tienen el pequeño sesgo que WB encontraron. Pelo contrario, tienen un sesgo muy grande.

Conclusión: métodos de PS pueden ser afetados por pequeños cambios en la muestra, y no siempre funcionan tan bien.

Segundo, Smith and Todd utilizan la misma muestra que WB, pero construyen el PS con solamente las variables que tenía Lalonde.

Otra vez, los resultados son malos. El sesgo es grande, y los resultados varían muchísimo con cambios en el estimador.

Conclusión: métodos de PS también pueden ser afectados por exactamente cuales variables están disponibles.

Tercero, Smith and Todd implementan otros estimadores de PS con aún más estructura. En particular, un estimador que mezcla PS con DID (que vamos estudiar más adelante). Ese estimador resulta un poco mejor, pero todavía tiene problemas.

Finalmente, Smith and Todd muestran que, cuando se utiliza la misma muestra que WB, y se utiliza métodos más simples, como regresiones con los mismos controles, los resultados son bastante buenos.

Conclusión: la mejora en los resultados de WB fue más por causa de sus restricciones en la muestra que por utilizar métodos de matching.

Basicamente, cuando WB excluyeron personas sin ingresos pasados, seleccionaran una muestra con un problema de sesgo *más fácil* de solucionar.

La conclusión general es que propensity score no es una solución general.

Cuando tenemos mucha información y de buena calidad, métodos simples resultan bien. Cuándo tenemos pocas variables, y los resultados son medidos de manera inconsistente, incluso métodos avanzados encuentran problemas.

Importante comprender los datos y sus problemas, y elegir métodos apropiados.