

# Causal inference foundations and applications in environmental health sciences

Jay S. Kaufman, PhD

Dept of Epidemiology, Biostatistics & Occupational Health  
McGill University, Montreal, Quebec, Canada

Sunday August 26<sup>th</sup>  
Shaw Centre, Ottawa  
Room 4  
13:00 - 16:00



# Outline for Sunday August 26th, 2018

Jay Kaufman, McGill University      13:00 - 14:25 (85 min)

Break      14:25 - 14:35 (10 min)

Alexander Keil, UNC-CH      14:35 - 16:00 (85 min)

1<sup>st</sup> Half Topic:      85 minutes => 114 slides

Causation vs. Association, DAGs and Confounding      ( 3 - 21)

Propensity Scores      ( 23 - 48)

IPTW/Marginal Structural Models/G-Formula      ( 49 - 116)

Hernán & Robins. *Causal Inference*, in press.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

Causal inference is necessary for medical and public policy decision-making because we hope to optimize some outcome.

Causal inference is about inherently unobservable things (i.e. the future under different scenarios)

Because we can't directly observe what we want to know, we model it.

**Good models**

**Bad models**

From 1999 to 2009, the number of Americans who fell into a swimming pool and drowned each year is correlated with the number of films in which Nicholas Cage appeared that year.

Shall we reduce the number of pool drownings by keeping Cage off the screen?

Three main inferential targets of these models:

1) Real world in the present

surveillance, descriptive study

2) Real world in the future

clinical prediction model

3) Hypothetical world in the future

causal inference, etiologic study

The inferential target determines the adjustment strategy.

Most people here are interested in 3)

If you are trying to estimate the causal effect of a treatment, your job is to **PREDICT** what would happen in the **FUTURE** if you did **thing A** compared to what would happen if you did **thing B**.

To do this from observational data, you must often adjust statistically for factors that are associated with the treatment and the outcome.

You observe:

$$\Pr(Y|X=x)$$

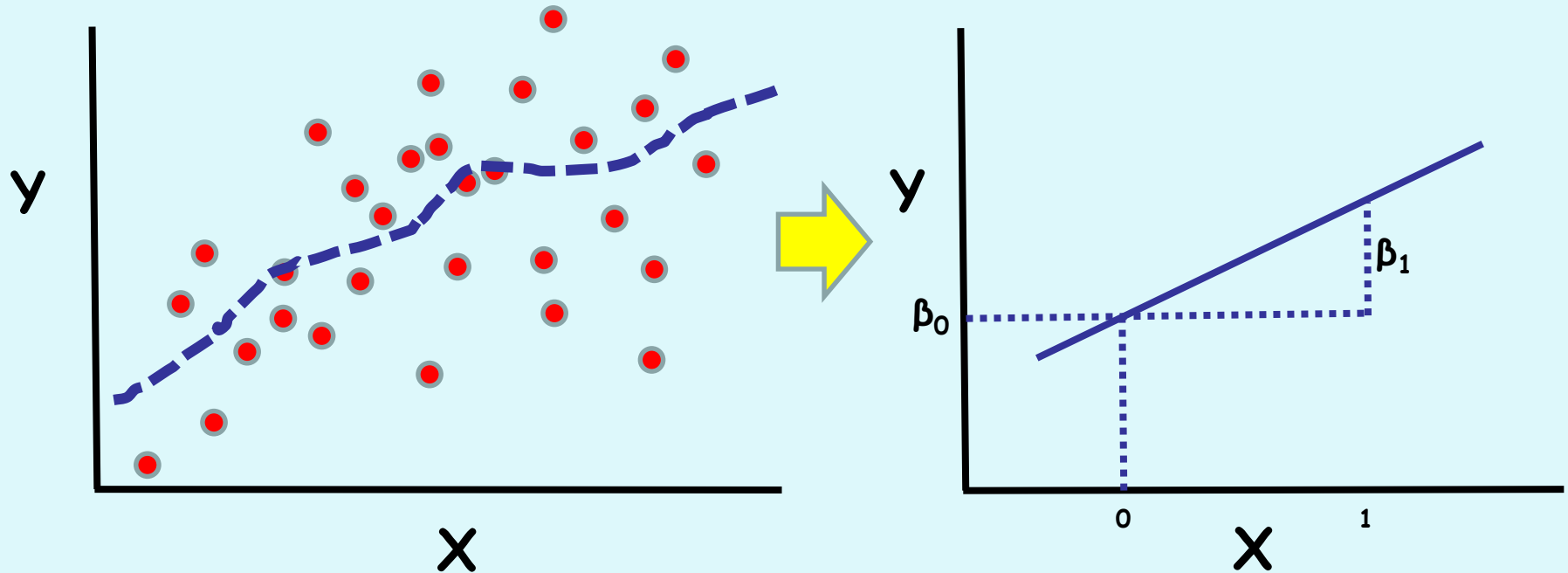
You want to know:

$$\Pr(Y|\text{SET}[X=x])$$



This is the intervention you want to know about, but unfortunately you don't really get to "SET" anything.

Statistical models are used to estimate relationships between variables in observational data sets.



But it is mechanistic knowledge or structural assumptions that allow us to infer causal effects from these relationships (not statistical considerations)

Read:  $\Pr(Y|\text{SET}[X=x])$

as:  $\Pr(Y|\text{SET}[X=x_1])$  versus  $\Pr(Y|\text{SET}[X=x_2])$

$x_1$  and  $x_2$  are the levels at which you intervene to set the treatment; contrast is usually a difference or ratio.

Causal inference from passively observed data requires not only structural identification, but also:

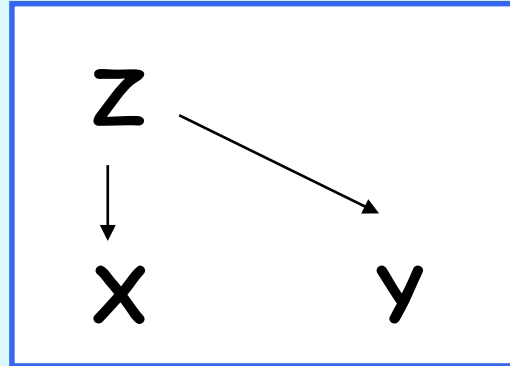
**positivity** (there are sufficient data available on the treatment and outcome in the range of interest)

**consistency** (the way that people came to be treated in the data set is comparable to the way that you plan to treat them in your intervention)

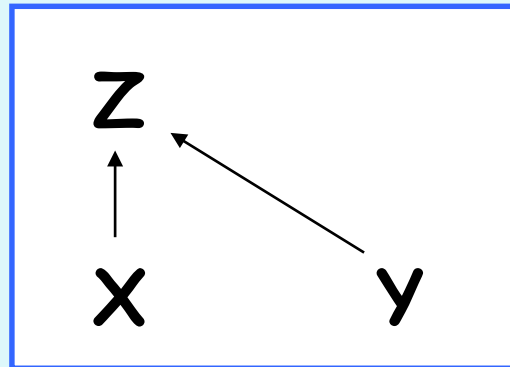
**correct specification of statistical models**

# Three main structural threats to validity:

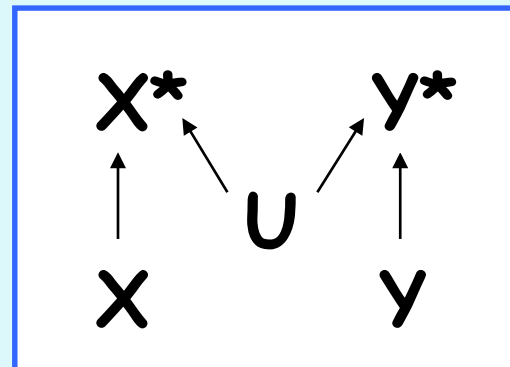
## Confounding Bias



## Selection Bias



## Information Bias





# Identification versus estimation

If you can't get be guaranteed to get the right answer as  $n \rightarrow \infty$ , then you have an identification problem.

DAGs are all about expected values, and therefore are focused entirely on identification.

This has some drawbacks for application to real-world studies with  $n \ll \infty$ .

For example, RCT DAG shows no confounding, but an RCT can, by bad luck, have an imbalanced covariate. This is confounding, even though there is no confounding in the correctly specified DAG.

# Statistical Adjustment:

"The philosophers have only interpreted the world, the point, however, is to change it." --- Karl Marx

There are two kinds of analysis:

descriptive                      and                      etiologic (causal)

If you are doing descriptive analysis, you show a picture of the world as it really is. No "adjustments". Why not?

Because the real world is unadjusted.

If you are doing an etiologic (causal) analysis, your job is to identify what would happen if you intervened on the world in some specific way.

In order to do this from observational data, you must often adjust statistically for factors that are associated with the exposure and outcome under study.

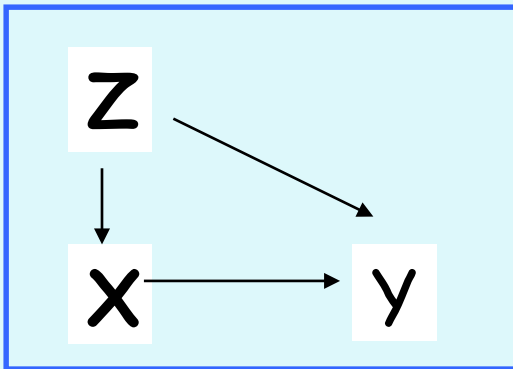
You observe:

$$\Pr(Y|X=x)$$

You want to know:

$$\Pr(Y|\text{SET}[X=x])$$

The adjustment tradition in epidemiology and the social sciences exists to link these two quantities:



$$\Pr(Y|X=x) \neq \Pr(Y|\text{SET}[X=x])$$

BUT!

$$\sum \Pr(Y|X=x, Z=z) \Pr(Z=z) = \Pr(Y|\text{SET}[X=x])$$

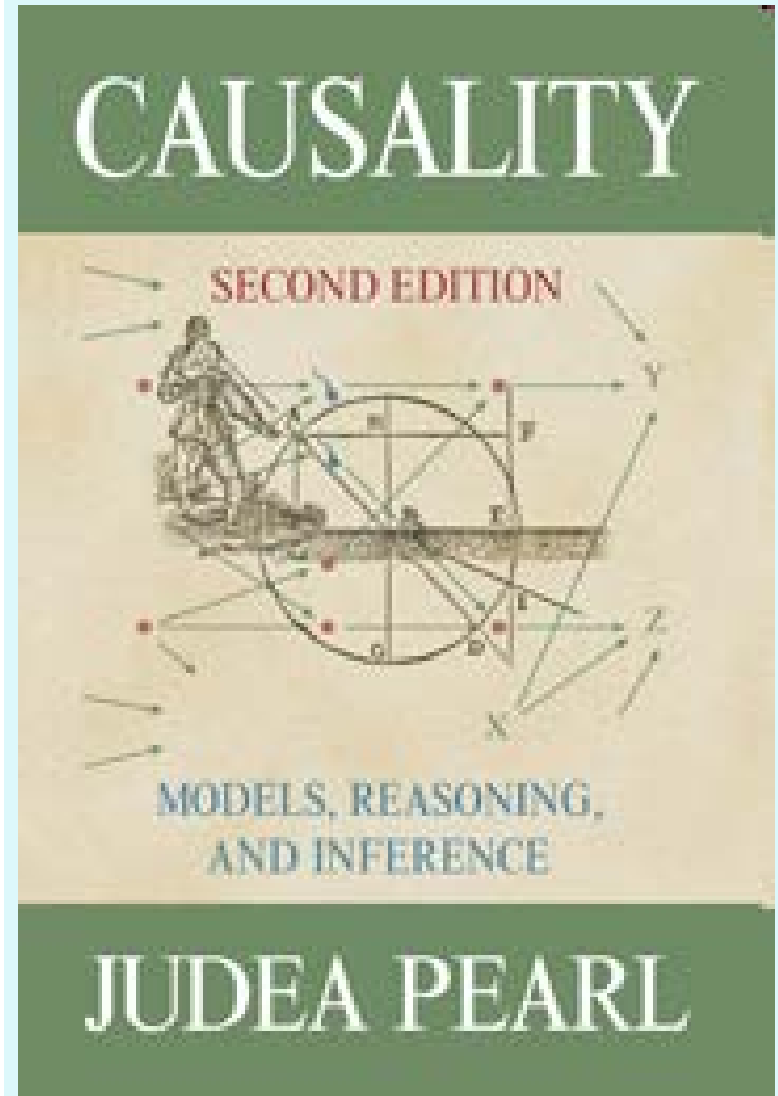
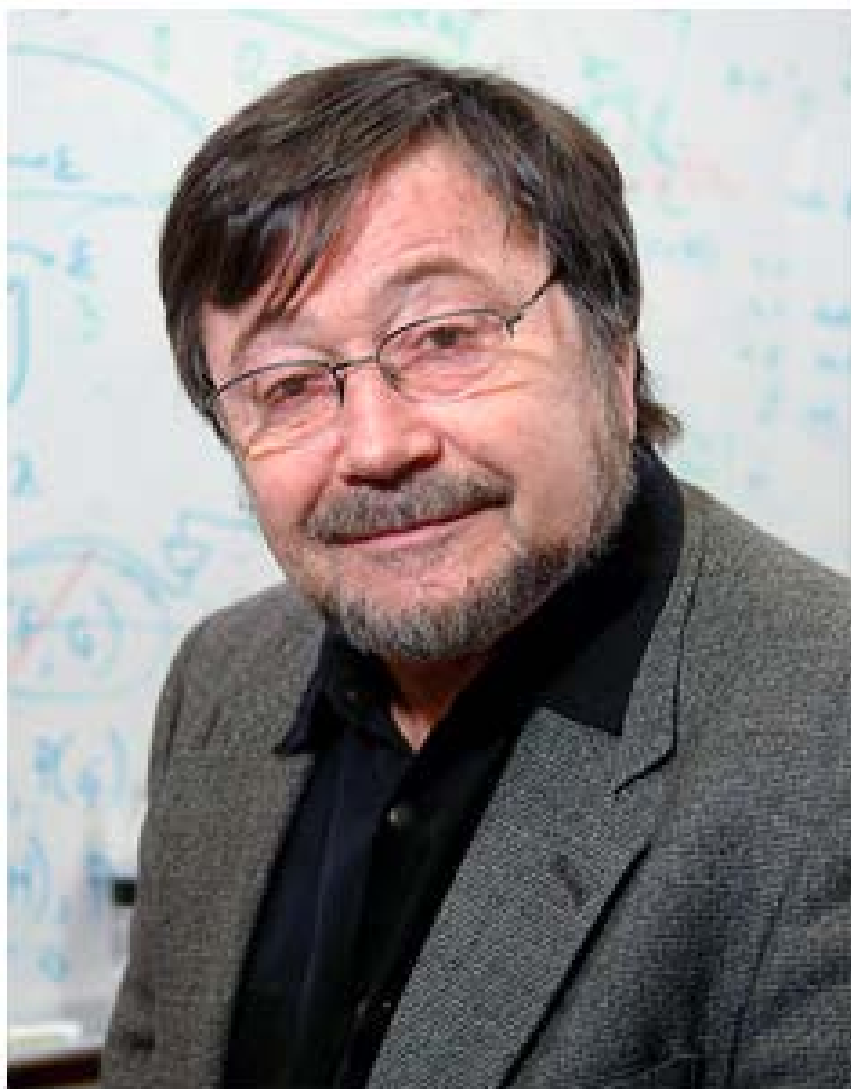
Read:  $\Pr(Y|\text{SET}[X=x])$   
as:  $\Pr(Y|\text{SET}[X=x_1])$  versus  $\Pr(Y|\text{SET}[X=x_2])$

where  $x_1$  and  $x_2$  are two levels at which you can intervene to set the exposure, and the contrast is usually a difference or ratio.

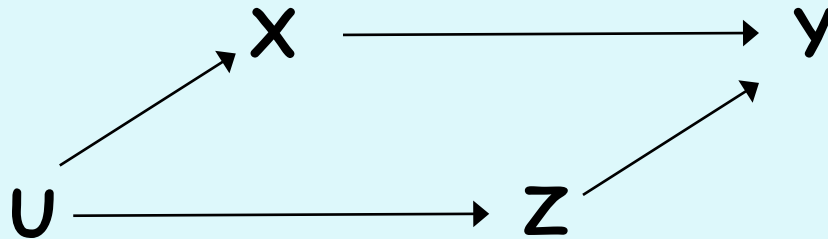
Clearly, quantities intermediate between exposure and outcome are not "confounders", they are just the way that the exposure has the effect that it has.

(See: Kaufman *American J of Law & Medicine* 2017  
Schisterman et al *Epidemiology* 2009, etc)

# Graphical language for encoding subject matter knowledge about causal structure



Compare a graphical model with a typical parametric epidemiologic model, such as logistic regression:



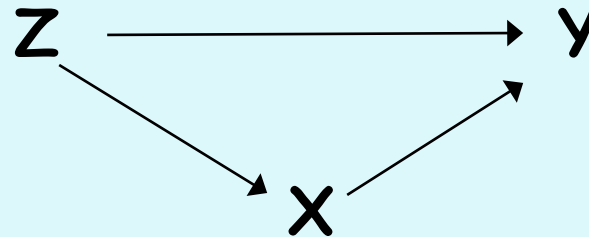
The graphical model asserts only that:  $Y = f(X, Z, \varepsilon_y)$   
and that  $X = f(U, \varepsilon_x)$  and  $Z = f(U, \varepsilon_z)$

The logistic regression model:

$$E(Y|X,Z) = \left[ \frac{e^{(\alpha + \beta_1 X + \beta_2 Z)}}{1 + e^{(\alpha + \beta_1 X + \beta_2 Z)}} \right]$$

makes *MANY* assertions, including the multiplicative interaction of  $X$  and  $Z$ , and the linearity of the  $\ln(\text{odds})$  of  $Y$  across all values of  $X$  and  $Z$ .

Furthermore, a graphical model can represent many assumptions that cannot be encoded in a typical statistical model:



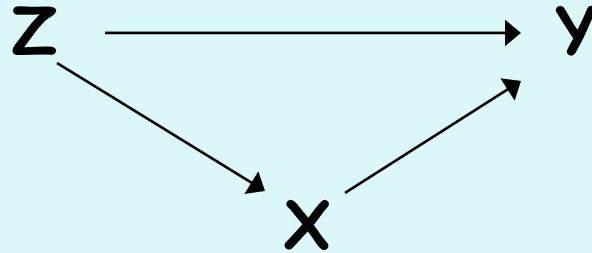
The graphical model asserts that:  $X = f(Z, \varepsilon_X)$

The logistic regression model:

$$E(Y|X,Z) = \left[ \frac{e^{(\alpha + \beta_1 X + \beta_2 Z)}}{1 + e^{(\alpha + \beta_1 X + \beta_2 Z)}} \right]$$

cannot easily represent this constraint, even if it is known by the investigators to be true on subject matter grounds (e.g.,  $Z = \text{SEX}$ ,  $X = \text{SMOKING}$ )

# Confounding



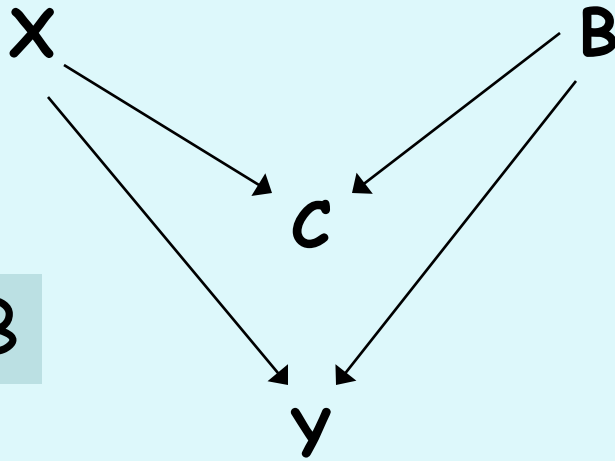
Confounding is a divergence between two kinds of conditional probability distributions of  $Y$ :

the distribution given that we find  $X$  at the value  $x$  (estimable from the data), and the distribution given that we intervene to force  $X$  to take the value  $x$ .

With Pearl's *SET* notation, express confounding as:

$$\Pr(Y = y \mid SET[X = x]) \neq \Pr(Y = y \mid X = x)$$

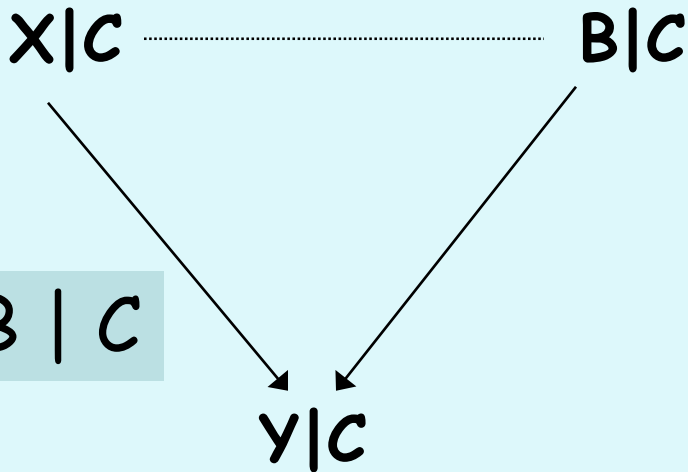




$$X \perp\!\!\!\perp B$$

NO CONFOUNDING  
FOR THE EFFECT  
OF X ON Y

No confounding in the source population when there is no backdoor path, or when a backdoor path is blocked by a collider.



$$X \not\perp\!\!\!\perp B \mid C$$

~~NO~~ CONFOUNDING  
FOR THE EFFECT  
OF X ON Y, GIVEN C

Most causal inference methods assume that you have no unmeasured confounders:

Regression

Propensity Scores

Marginal Structural Models

G-methods (SNMs, G-Formula, etc)

“Quasi-Experimental” Methods use structural assumptions to achieve identification even in the presence of unmeasured confounding:

Instrumental variables

Regression Discontinuity

Fixed Effects Differences in Differences

Some causal inference methods achieve identification based on extrapolation of a parametric model.

Semi-parametric methods (e.g. propensity scores, IPTW, TMLE, etc) rely less on model form. Letting a computer pick the model reduces “wish bias”.

Non-parametric methods (e.g. matching) require no model at all.

Doubly robust methods require that at least one model be right, but not both.

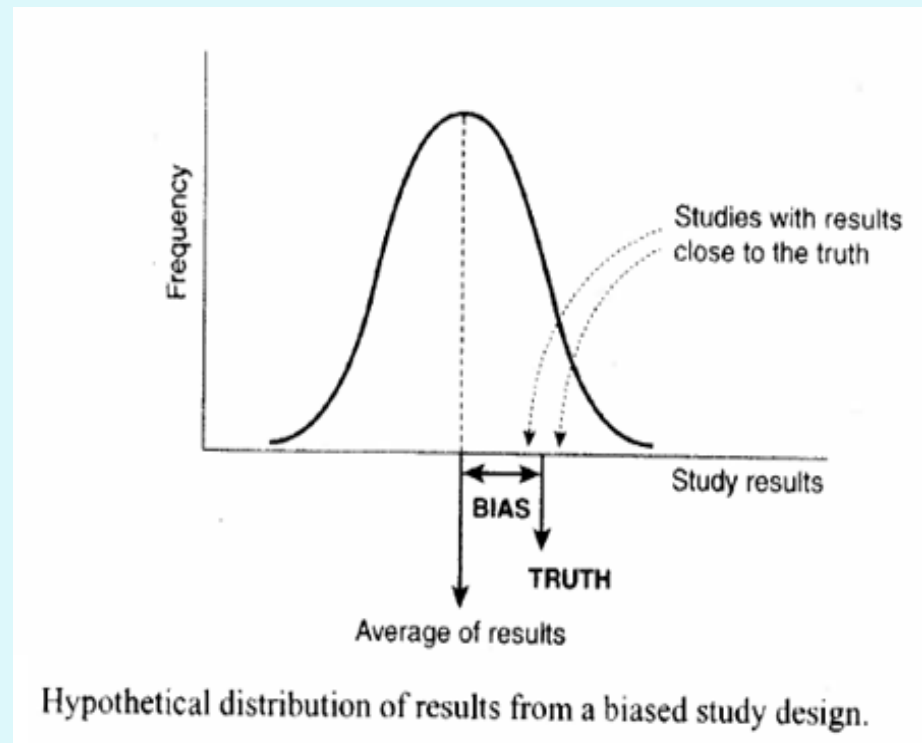
Computer intensive methods (e.g. bootstrapping) reduce reliance on distributional assumptions.

# Confounding is a bias:

Validity and Bias: The epidemiologist's goal is the most **VALID** and **PRECISE** estimate possible of the causal effect of exposure on disease.

Error comes from sampling variability (lack of precision) and bias (lack of validity).

Szklo & Nieto,  
2<sup>nd</sup> edition, 2007



# Big data context:

All usual threats to validity still apply:

confounding bias	}	None are reduced as $n \rightarrow \text{large}$
selection bias		
information bias		

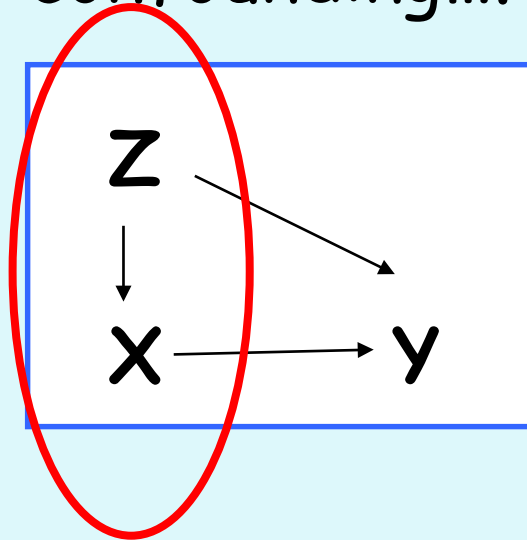
Big data means that random error should ↓

So for precision: BIGGER REALLY IS BETTER

Precision gain may not be “worth the price” if data quantity is negatively correlated with data quality (as often happens).

# Exposure Modeling

OK, now back to confounding....



All confounding arrived via arrows into X.

So modeling receipt of treatment is a way to control for confounding.  $\Pr(X=1|Z)$  is "the propensity score"

Rosenbaum, P.R. & Rubin, D.B. (1983). The Central Role of the Propensity Score in Observational Studies for Causal Effects. *Biometrika*, 70(1), 41-55.

- over 20,000 citations
- causal inference is a missing data problem
- defines a balancing score  $b(z)$ , is a function of the observed covariates  $z$  such that the conditional distribution of  $x$  given  $b(x)$  is the same for treated ( $x = 1$ ) and control ( $x = 0$ ) units, that is:  $z \perp\!\!\!\perp x \mid b(z)$ .



**Identifies the propensity score as the coarsest possible balancing score.**

# "Weakly Ignorable Treatment Assignment" within strata of $Z$

Conditioning on  $Z$  sufficient adjustment to control for all confounding bias if, within each stratum of  $Z$ , observed exposure  $X$  is statistically independent of the potential response ( $Y|SET[X=x]$ ), for each imposed value  $x$ .

Write:

$$Y_x \perp\!\!\!\perp X \mid Z$$

where  $Y_x$  is the potential response of  $Y$  to treatment  $x$   
i.e., ( $Y|SET[X=x]$ )



## Target Populations:

Each subject  $i$  has a pair of potential outcomes:  $Y_i(0)$  and  $Y_i(1)$  for binary treatment that takes values 0 and 1.

Each subject  $i$  receives only one treatment, so 50% of values are missing.

For subject  $i$ , individual effect of treatment is  $Y_i(1) - Y_i(0)$

Average treatment effect (ATE) =  $E[Y_i(1) - Y_i(0)]$

An alternative target population is to consider the effect only in the exposed.

This is called the ATT (average treatment effect for the treated) or ETT (effect of treatment on the treated)

ATT is defined as  $E[Y(1) - Y(0) | X = 1]$ .

In an RCT,  $ACE = ATT$  because, due to randomization, the treated population will not, on average, differ systematically from the overall population.

Applied researchers must decide whether ATE or ATT is of greater utility or interest in their particular context.

e.g. for estimating effectiveness of an intensive, structured smoking cessation program, ATT might make more sense.

In contrast, for effect on smoking cessation of an information brochure given by family physicians to patients who are current smokers, ATE may be of greater interest.

Also possible to define ATU (but not commonly used).

# Propensity Scores

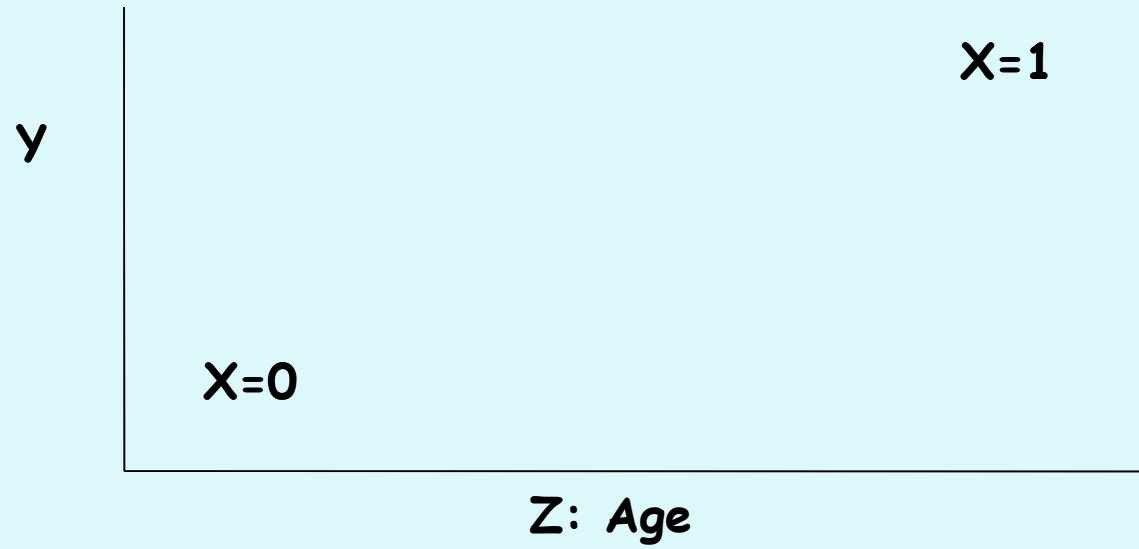
Typically, many background characteristics need to be controlled in a study of an observed exposure. Propensity score technology reduces the collection of background characteristics to a single "composite" characteristic that appropriately summarizes the collection.

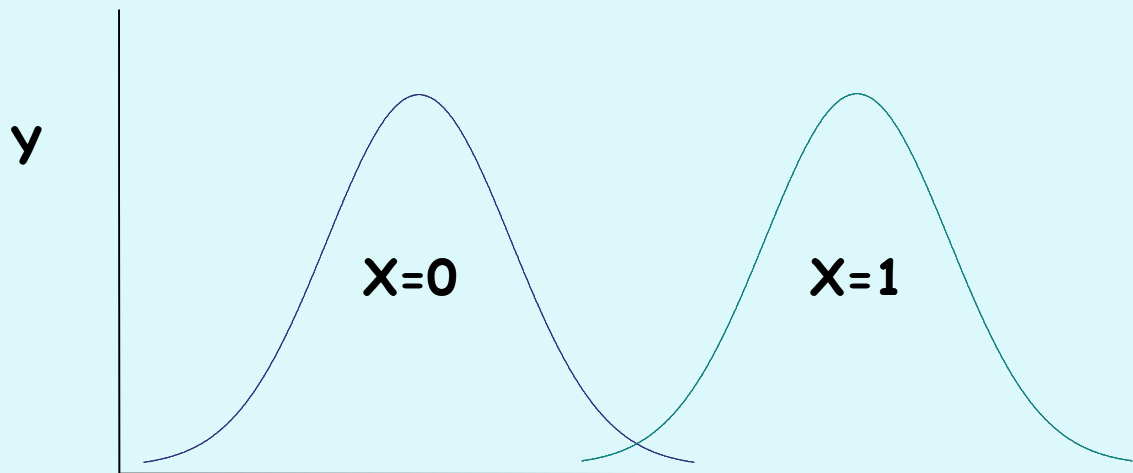
This reduction:

**many characteristics → one composite characteristic**

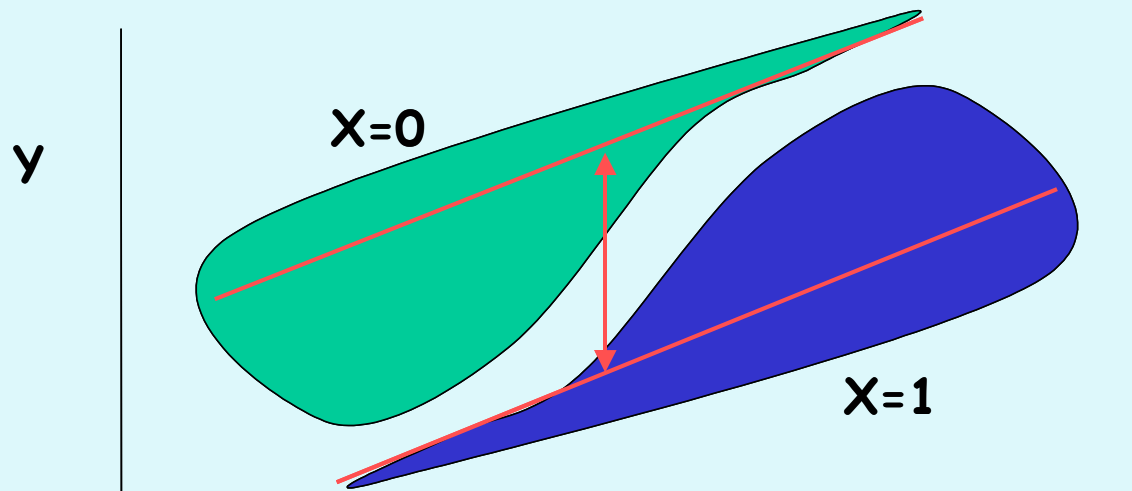
allows the *straightforward* assessment of whether the exposed and unexposed groups overlap enough on background characteristics to allow sensible estimation of causal effects from the data set.

Use to CHECK BALANCE and assess covariate sufficiency.





**Z: Pretreatment covariate (e.g. age)**



**Z: Pretreatment covariate (e.g. age)**

The propensity score is just the estimated  $\Pr(X=1|Z)$ .

How to use this information?

At least 4 common methods have been defined:

- 1) stratification (or subclassification) on the propensity score (R & R's original idea)
- 2) propensity score matching
- 3) covariate adjustment using the propensity score
- 4) inverse probability of treatment weighting (IPTW)

# Which Variable to Include

- **Best to include: ( $\downarrow$  Bias)**  
factors that affect outcome and are correlated with receipt of treatment
- **May be beneficial to include ( $\uparrow$  precision ):**  
factors that affect outcome and are not correlated with receipt of treatment
- **May be harmful to include ( $\downarrow$  precision ):**  
factors that do not affect outcome but are correlated with receipt of treatment
- **Worst to include: ( $\uparrow$  Bias)**  
factors that are affected by the exposure or by the outcome

# Comparing propensity score distributions

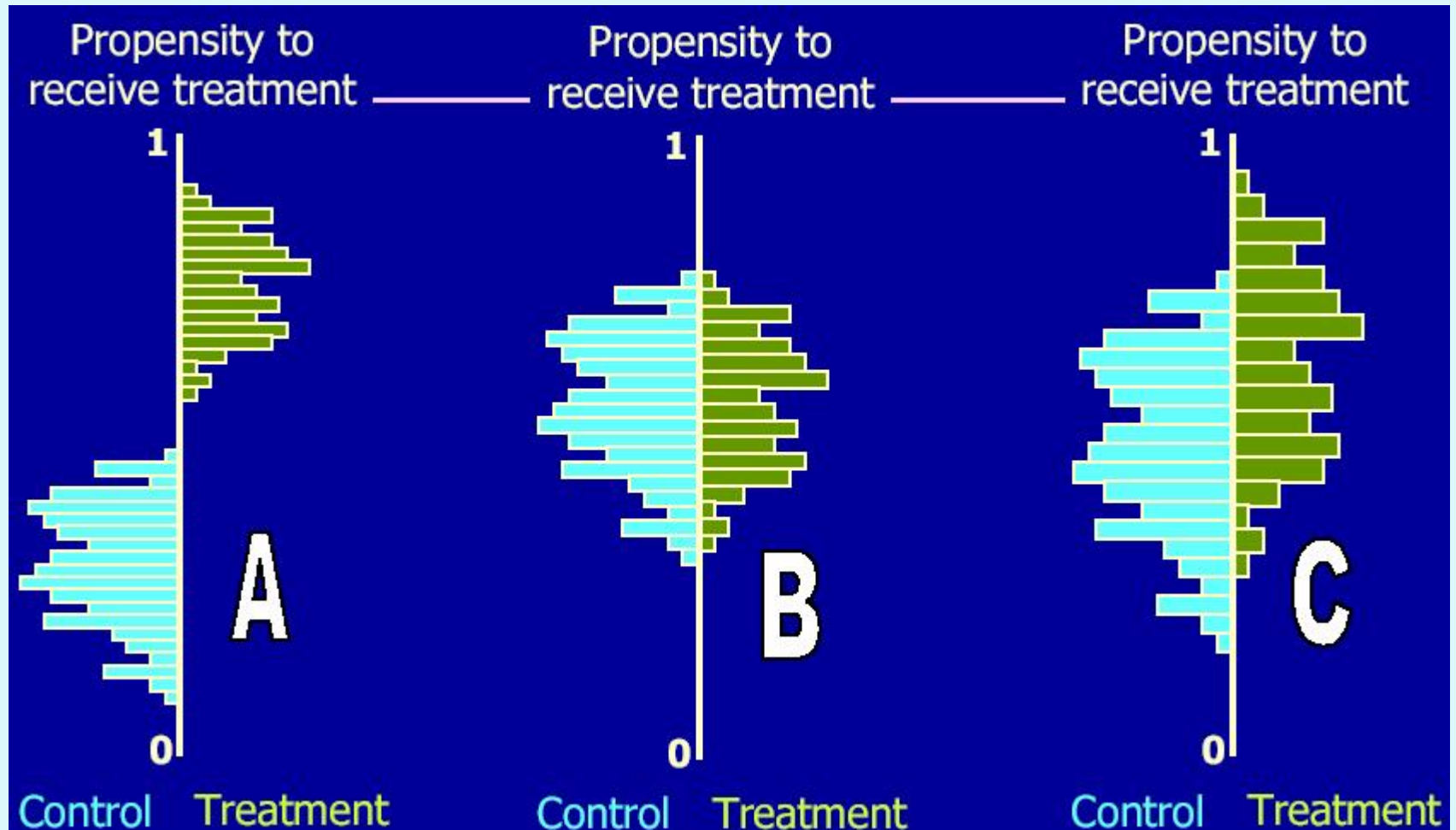
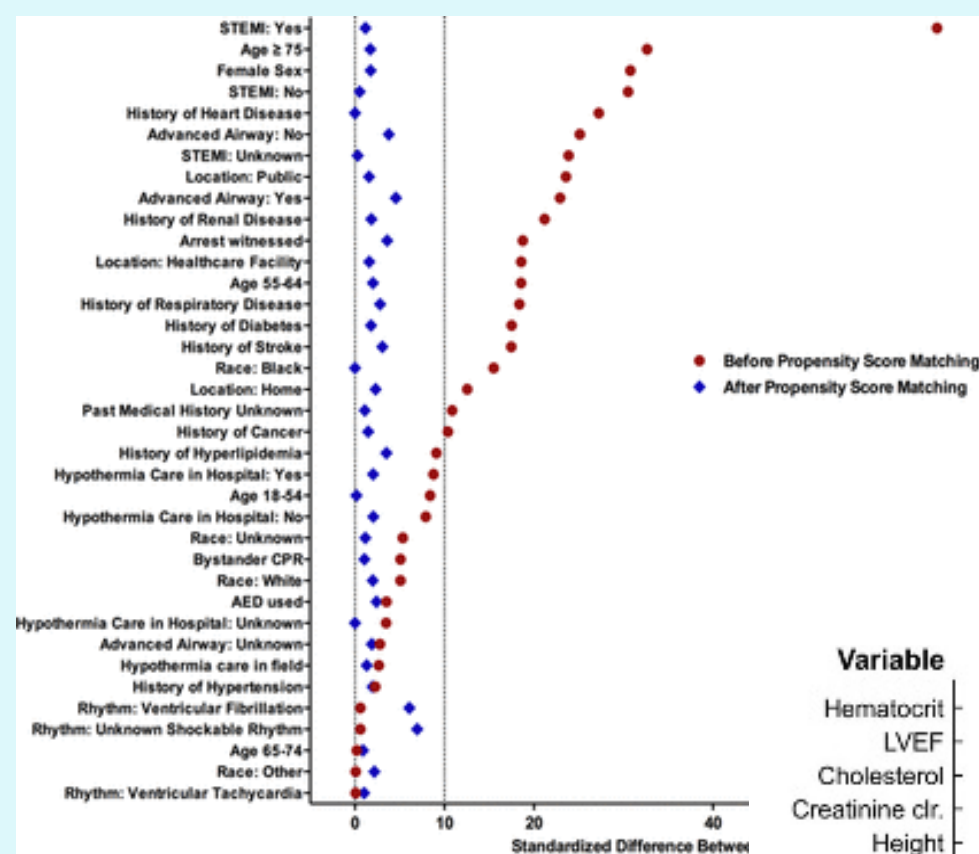


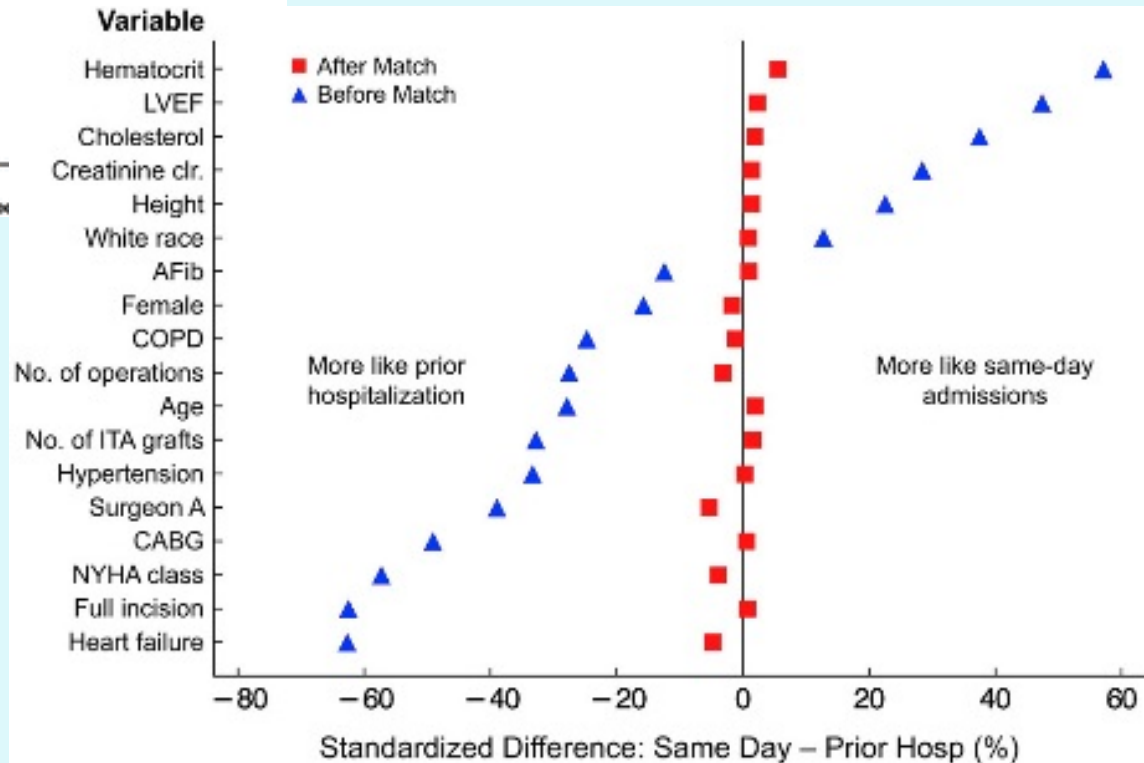
Figure from: Thomas E Love, PhD (<https://cwru.pure.elsevier.com/en/persons/thomas-e-love-2>)



Vyas et al. Early Coronary Angiography and Survival After Out-of-Hospital Cardiac Arrest. *Circulation* 2015.



Winger et al. Propensity-score analysis in thoracic surgery: When, why, and an introduction to how. *Journal of Thoracic and Cardiovascular Surgery* 2016



## PS Advantages:

- Non-parametric, good face validity
- Captures interactions among covariates
- Examination of PS distributions exposes non-comparability (Rubin 1997)
- Single propensity score can be used in analysis of effect of exposure on >1 outcome
- Matching & stratification are robust to misspecification of PS
- Multidimensional matching with less loss of potential matches

## PS Disadvantages:

- Balance may not be achieved (missing confounders, deterministic predictors, small sample size)
- Unmeasured confounders - only controlled to the extent that they are correlated with measured confounders
- Residual confounding within PS strata
- Direct adjustment relies on model form (Rubin *PDS* 2004)
- Potential adjusting for variables affected by exposure or other non-confounders
- Modeling  $\Pr(X=1|Z)$  may be hard for rare exposure with many  $Z$
- No effect estimates obtained for covariates
- Generalizability (if overlap is limited)
- Inconvenient for non-binary exposures

# Consider the famous example of the impact of the NSW job-training program, using the "nswre74" dataset

```
. des
```

```
Contains data from C:\Wednesday\nswre74.dta
```

```
obs:      445
vars:      11                      20 May 2013 12:01
size:     10,235
```

storage	display	value		
variable name	type	format	label	variable label
treat	byte	%9.0g	noyes	treatment group?
age	byte	%9.0g		age (years)
ed	byte	%9.0g		years of education
black	byte	%9.0g	noyes	black race?
hisp	byte	%9.0g	noyes	hispanic ethnicity?
married	byte	%9.0g	noyes	married?
nodeg	byte	%9.0g	noyes	no high school degree?
re74	float	%9.0g		1974 earnings
re75	float	%9.0g		1975 earnings
re78	float	%9.0g		1978 earnings
age2	float	%9.0g		age squared

```
Sorted by:
```

# BALANCE DIAGNOSTICS

Propensity score is a balancing score: conditional on true propensity score, distribution of measured baseline covariates is independent of treatment assignment.

But we don't know the true propensity score. Must estimate from data, which requires checking balance.

This is a major advantage of PS over regression models.

**In strata of subjects with similar PS, distribution of measured baseline covariates should be similar between treated and untreated subjects.**

We have a way to find out if our model is adequate.

## BALANCE DIAGNOSTICS

For matched data, compare treated and untreated subjects within the propensity score matched sample.

For IPTW, compare treated and untreated subjects in the inverse weighted sample

For data subclassified on the PS, compare treated and untreated subjects within the same strata of PS.

For 1-to-1 matching, the **standardized difference** is used to compare the means of continuous and binary variables between treatment groups. Multilevel categorical variables are represented with a set of binary indicator variables.

```

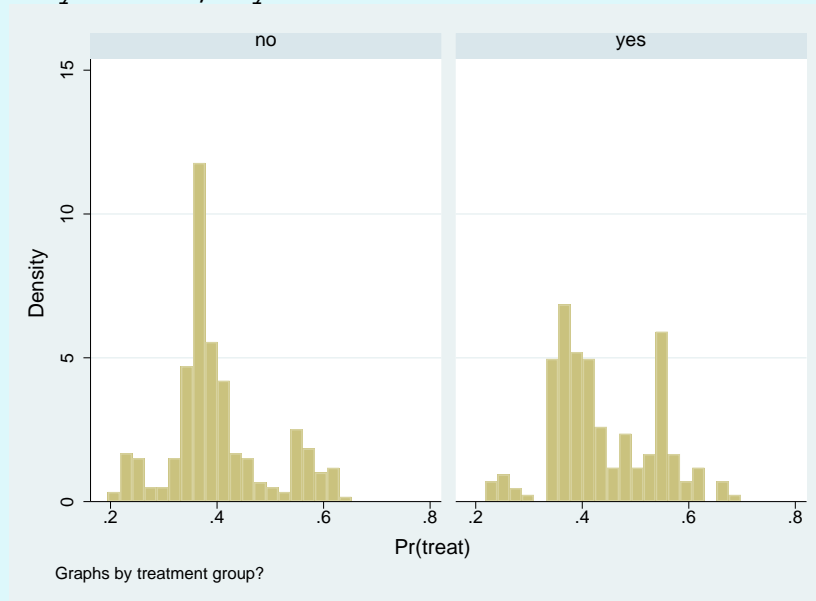
* build propensity score models for nsw ("by hand")

* predict treatment status
logit treat age age2 ed black hisp nodeg married re75 re74

* output predicted probability of treatment
predict ps

* summary statistics for propensity score, by treatment
bysort treat: sum ps, det
histogram ps, by(treat)

```



Good  
overlap

```

* nearest neighbor matching (1:1) without replacement

* randomly order data in case match ties
set seed 123456 // set seed to reproduce results

gen ranorder = runiform()
sort ranorder

* create propensity score
psmatch2 treat age age2 ed black hisp nodeg married re75 re74, logit neighbor(1) noreplacement

```

\* evaluate balance (Rosenbaum and Rubin formula)  
pstest age age2 ed black hisp nodeg married re75 re74, both graph

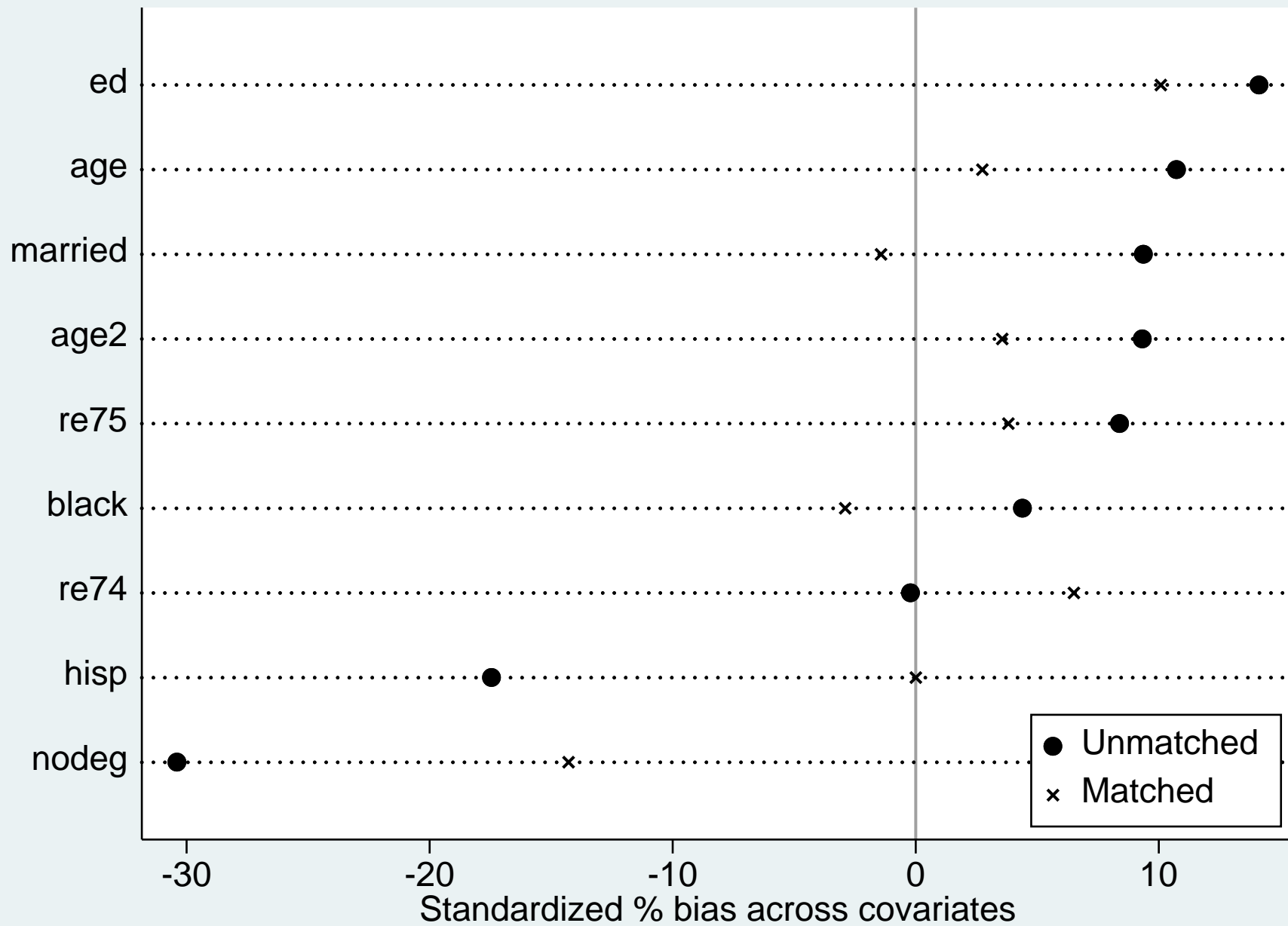
Variable	Unmatched Matched	Mean		%reduct		t-test		V(T)/ V(C)
		Treated	Control	%bias	bias	t	p> t	
age	U	25.816	25.054	10.7		1.12	0.265	1.03
	M	25.816	25.622	2.7	74.5	0.27	0.788	1.12
age2	U	717.39	677.32	9.3		0.97	0.333	1.01
	M	717.39	702.11	3.6	61.9	0.35	0.726	1.13
ed	U	10.346	10.088	14.1		1.50	0.135	1.55*
	M	10.346	10.162	10.1	28.6	0.94	0.346	1.36*
black	U	.84324	.82692	4.4		0.45	0.649	.
	M	.84324	.85405	-2.9	33.8	-0.29	0.772	.
hisp	U	.05946	.10769	-17.5		-1.78	0.076	.
	M	.05946	.05946	0.0	100.0	-0.00	1.000	.
nodeg	U	.70811	.83462	-30.4		-3.22	0.001	.
	M	.70811	.76757	-14.3	53.0	-1.30	0.195	.
married	U	.18919	.15385	9.4		0.98	0.327	.
	M	.18919	.19459	-1.4	84.7	-0.13	0.895	.
re75	U	1532.1	1266.9	8.4		0.87	0.382	1.08
	M	1532.1	1411.7	3.8	54.6	0.35	0.725	0.92
re74	U	2095.6	2107	-0.2		-0.02	0.982	0.74*
	M	2095.6	1750.8	6.5	-2910.6	0.70	0.486	1.12

\* if variance ratio outside [0.75; 1.34] for U and [0.75; 1.34] for M

Sample	Ps R2	LR chi2	p>chi2	MeanBias	MedBias	B	R	%Var
Unmatched	0.028	17.04	0.048	11.6	9.4	40.1*	1.00	40
Matched	0.006	2.92	0.967	5.0	3.6	17.7	1.26	20

\* if B>25%, R outside [0.5; 2]





# Can we improve the model a bit?

```
* generate squared term for 1974 earnings
gen re74_2 = re74*re74
label var re74_2 "1974 earnings squared"
```

```
* randomly order data in case match ties
replace ranorder = runiform()
sort ranorder
```

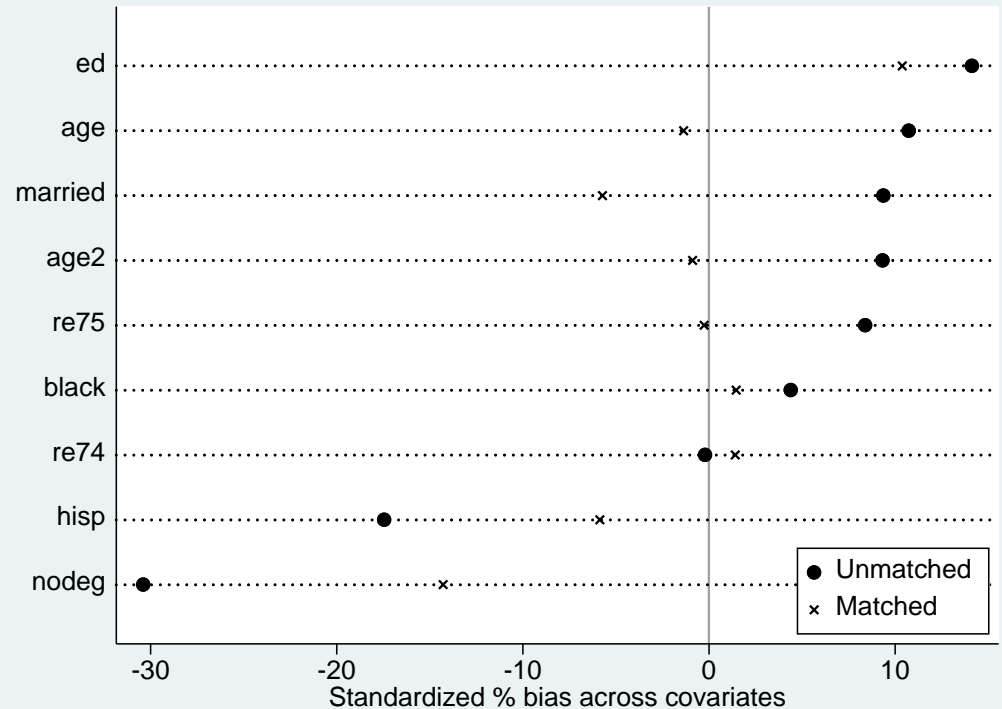
```
* create propensity score
psmatch2 treat age age2 ed black hisp nodeg married re75 re74 re74_2, logit neighbor(1) noreplacement
```

```
* evaluate balance
pstest age age2 ed black hisp nodeg married re75 re74, both graph
```

```
* checking balance by hand...
```

```
Standardized difference for age = 1.4
Standardized difference for age2 = 0.9
Standardized difference for ed = 9.4
Standardized difference for black = 1.5
Standardized difference for hisp = 6.8
Standardized difference for nodeg = 13.0
Standardized difference for married = 5.5
Standardized difference for re75 = 0.3
Standardized difference for re74 = 1.5
```

Better, but "ed" and "nodeg" still don't look so good.



# Treatment effect in final matched sample:

```
regress re78 treat if _weight==1
```

```
Number of obs    =        370
```

re78	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
treat	1762.777	719.3954	2.45	0.015	348.1358	3177.419
_cons	4586.366	508.6894	9.02	0.000	3586.063	5586.669

Tx effect = \$1763 (95% CI: 348, 3177)

## "Doubly robust":

```
regress re78 treat age age2 ed black hisp nodeg married re75 re74 if _weight==1
```

```
Number of obs    =        370
```

re78	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
treat	1712.479	714.1617	2.40	0.017	308.0126	3116.945
age	386.866	328.6266	1.18	0.240	-259.409	1033.141
age2	-5.408275	5.392483	-1.00	0.317	-16.0131	5.196549
ed	383.7983	260.5215	1.47	0.142	-128.5417	896.1384
black	-2242.44	1252.403	-1.79	0.074	-4705.408	220.5279
hisp	190.7841	1829.93	0.10	0.917	-3407.946	3789.514
nodeg	101.479	1113.154	0.09	0.927	-2087.644	2290.601
married	-341.127	960.039	-0.36	0.723	-2229.134	1546.88
re75	.0345891	.154948	0.22	0.823	-.2701306	.3393089
re74	.1187688	.1048129	1.13	0.258	-.0873557	.3248932
_cons	-3871.263	5476.255	-0.71	0.480	-14640.83	6898.306

Tx effect = \$1712 (95% CI: 308, 3117)

# Treatment effects from -psmatch2- command:

```
psmatch2 treat age age2 ed black hisp nodeg married re75 re74 re74_2, logit neighbor(1)
noreplacement outcome(re78)
```

Variable	Sample	Treated	Controls	Difference	S.E.	T-stat
re78	Unmatched	6349.1435	4554.80112	1794.34238	632.853392	2.84
	ATT	6349.1435	4586.36616	1762.77735	719.395391	2.45

Note: S.E. does not take into account that the propensity score is estimated.

```
tab _treat _support
```

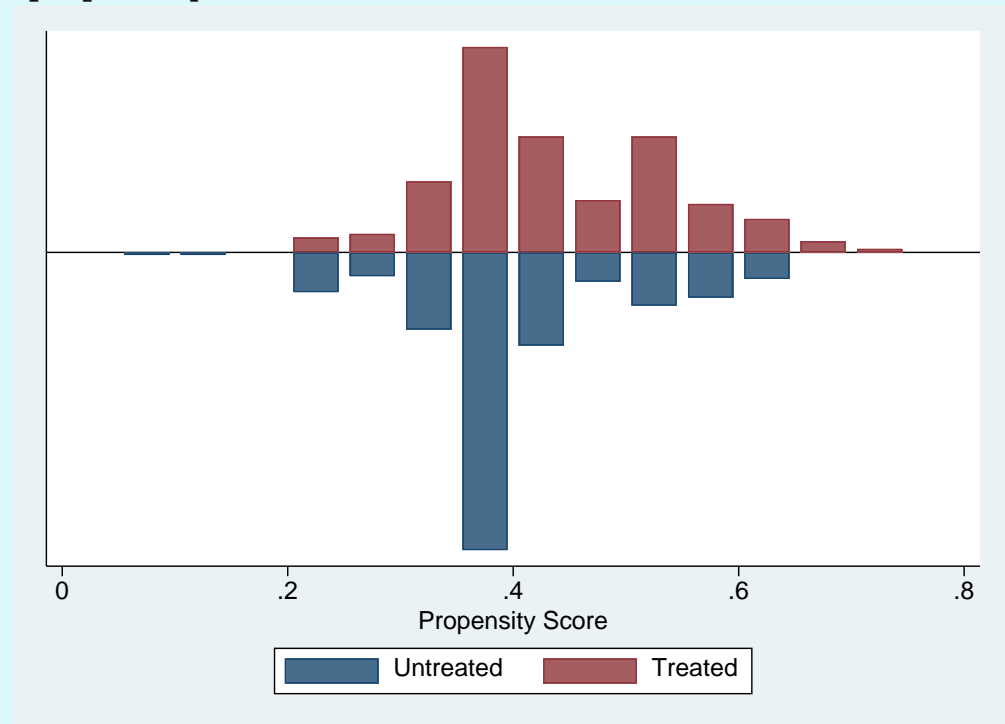
	psmatch2:	
	Common	
Treatment	support	
assignment	On support	Total
Untreated	260	260
Treated	185	185
Total	445	445

```
psgraph, name(pscore1, replace)
```

\* bootstrapped standard error for ATT:

```
. bootstrap r(att), reps(500): psmatch2 treat
logit neighbor(1) noreplacement outcome(re78)
```

	Observed	Bootstrap			Normal-based	
	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
tx effect	1762.777	726.0403	2.43	0.015	339.7644	3185.79



# Using Stata's new "treatment effects" suite of commands:

```
teffects psmatch (re78) (treat age ed black hisp married nodeg re74 age2 re74_2), atet
```

```
Treatment-effects estimation      Number of obs      =          445
Estimator      : propensity-score matching  Matches: requested =          1
Outcome model  : matching                      min =          1
Treatment model: logit                      max =          8
```

re78	AI Robust		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
ATT	1721.515	739.8562	2.33	0.020	271.4232	3171.606

```
teffects psmatch (re78) (treat age ed black hisp married nodeg re74 age2 re74_2)
```

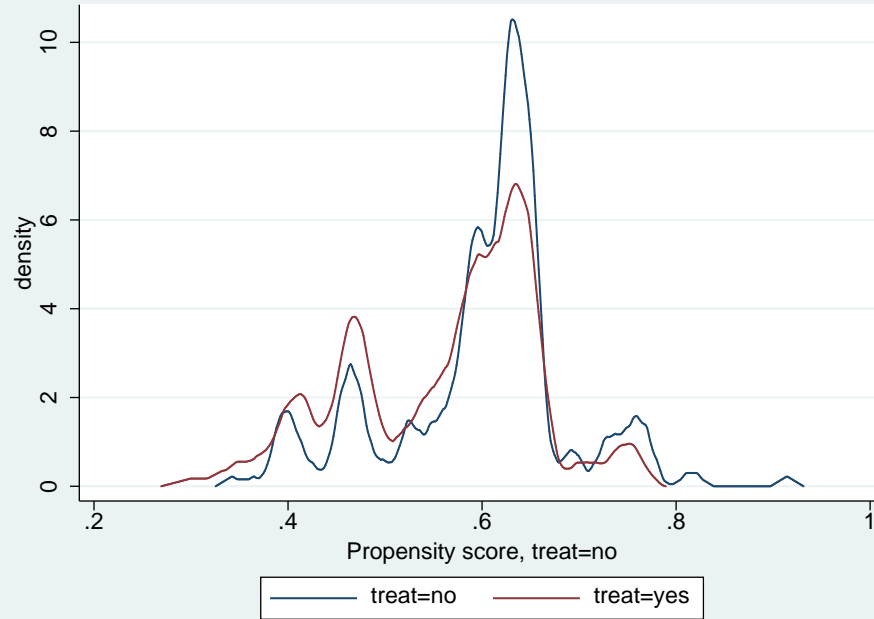
re78	AI Robust		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
ATE	1381.962	608.0156	2.27	0.023	190.2732	2573.651

```
teffects psmatch (re78) (treat age ed black hisp married nodeg re74 age2 re74_2), nneighbor(2) nopvalues
cformat(%9.2f) pformat(%5.2f) sformat(%8.2f) caliper(.2)
```

re78	AI Robust		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
ATE	1539.94	560.53			441.33	2638.55

# Post-estimation commands:

```
teffects overlap
```



```
tebalance summarize, baseline
```

Covariate balance summary

	Raw	Matched
Number of obs =	445	890
Treated obs =	185	445
Control obs =	260	445

	Means		Variances	
	Control	Treated	Control	Treated
age	25.05385	25.81622	49.81176	51.1943
ed	10.08846	10.34595	2.606044	4.042714
black	.8269231	.8432432	.1436739	.1329025

etc. . .

# Post-estimation commands:

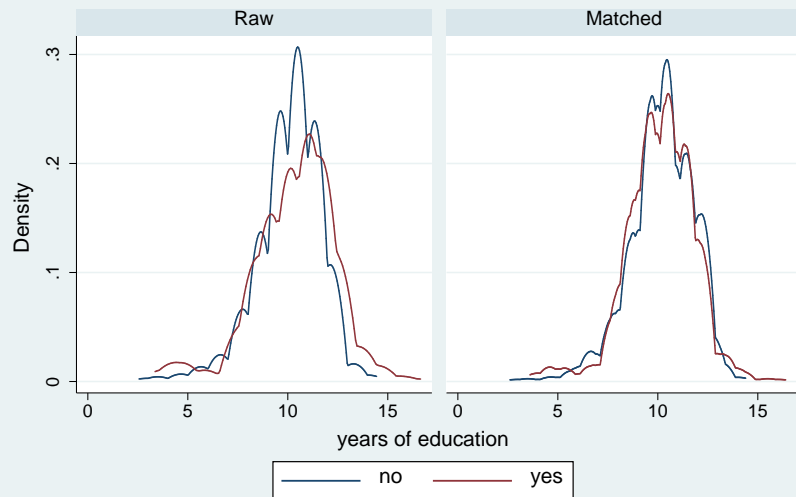
```
tebalance summarize
```

	Standardized differences		Variance ratio	
	Raw	Matched	Raw	Matched
age	.1072771	.0217026	1.027755	.8960536
ed	.1412198	-.0409838	1.551284	1.315371
black	.0438866	.005008	.9250286	.9911161
hisp	-.1745611	-.0240977	.5828804	.930887
married	.0936407	.0181852	1.180212	1.033589
nodeg	-.3039864	.0273399	1.499755	.9626536
re74	-.0021599	-.0554885	.7380953	.8065356
age2	.0932032	.0054639	1.011541	.805725
re74_2	-.0611646	-.0505421	.5038162	.7137261

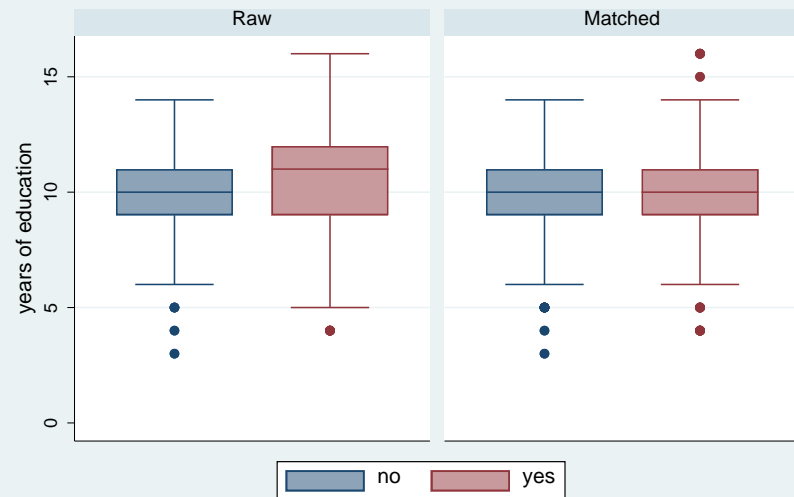
```
tebalance density ed
```

```
tebalance box ed
```

Balance plot



Balance plot



# Easier to specify more flexible models and check balance:

```
teffects psmatch (re78)(treat age ed black hisp married nodeg c.age#(c.age c.ed i.nodeg) c.re74#(c.re74 i.black))
```

Treatment-effects estimation	Number of obs	=	445
Estimator : propensity-score matching	Matches: requested	=	1
Outcome model : matching	min	=	1
Treatment model: logit	max	=	8

re78	Coef.	AI Robust Std. Err.	z	P> z	[95% Conf. Interval]	
ATE	1479.321	711.7633	2.08	0.038	84.29022	2874.351

```
tebalance summarize
```

Covariate balance summary

	Raw	Matched
Number of obs	445	890
Treated obs	185	445
Control obs	260	445

	Standardized differences		Variance ratio	
	Raw	Matched	Raw	Matched
age	.1072771	.0496652	1.027755	.9333797
ed	.1412198	-.0285369	1.551284	1.135913
black	.0438866	-.0416538	.9250286	1.075269
hisp	-.1745611	.0156963	.5828804	1.046103
married	.0936407	-.0234989	1.180212	.9610636
nodeg	-.3039864	.0054267	1.499755	.9926484
age#age	.0932032	.033652	1.011541	.8480046
age#ed	.1554148	.0223087	1.214733	1.016357
nodeg#age	-.2152312	.0118524	1.330088	.951223
re74#re74	-.0611646	-.023542	.5038162	.6547103
black#re74				
no	-.0497845	.1128004	.3545242	1.309552
yes	.0183761	-.0373243	.8288651	.8377786



# Further developments in exposure modeling for confounder control and causal inference:

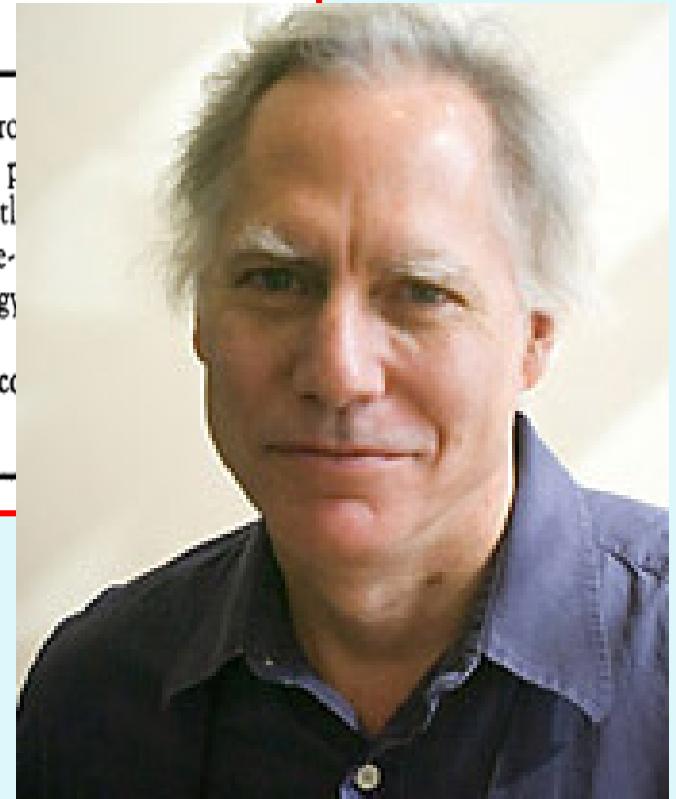
## Marginal Structural Models and Causal Inference in Epidemiology

*James M. Robins,<sup>1,2</sup> Miguel Ángel Hernán,<sup>1</sup> and Babette Brumback<sup>2</sup>*

In observational studies with exposures or treatments that vary over time, standard approaches for adjustment of confounding are biased when there exist time-dependent confounders that are also affected by previous treatment. This paper introduces marginal structural models, a new class of

causal models that allow for improved adjustment in those situations. The proposed structural model can be consistently estimated by a class of estimators, the inverse-proportionally weighted estimators. (Epidemiology

**Keywords:** causality, counterfactuals, epidemiologic methods, longitudinal data, structural models, causal variables



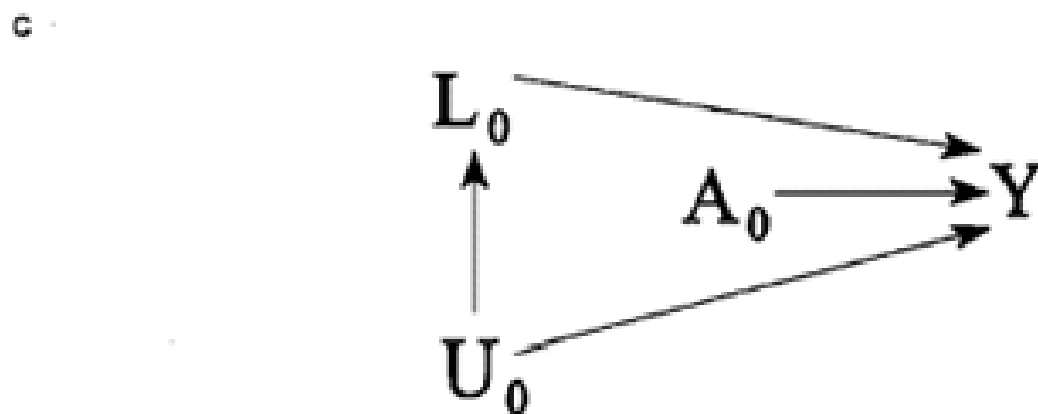
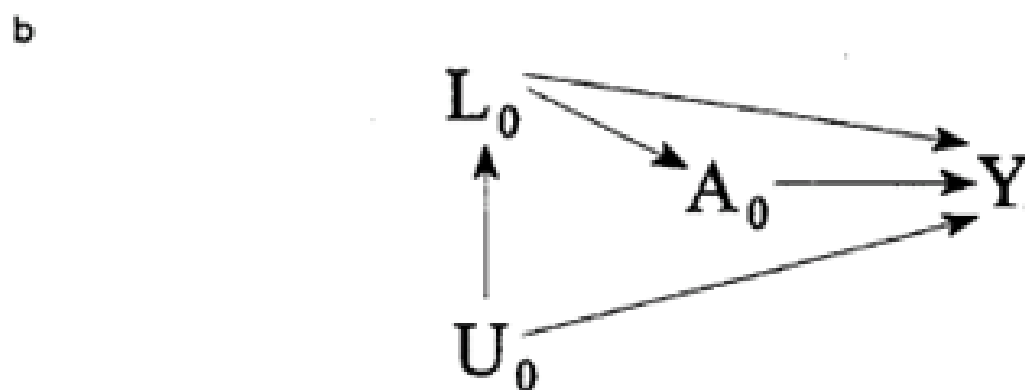
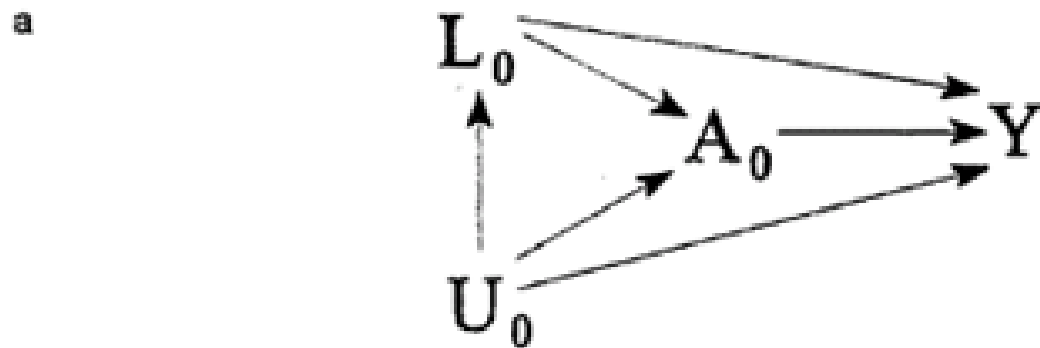


FIGURE 2. Causal graphs for a point exposure  $A_0$ .

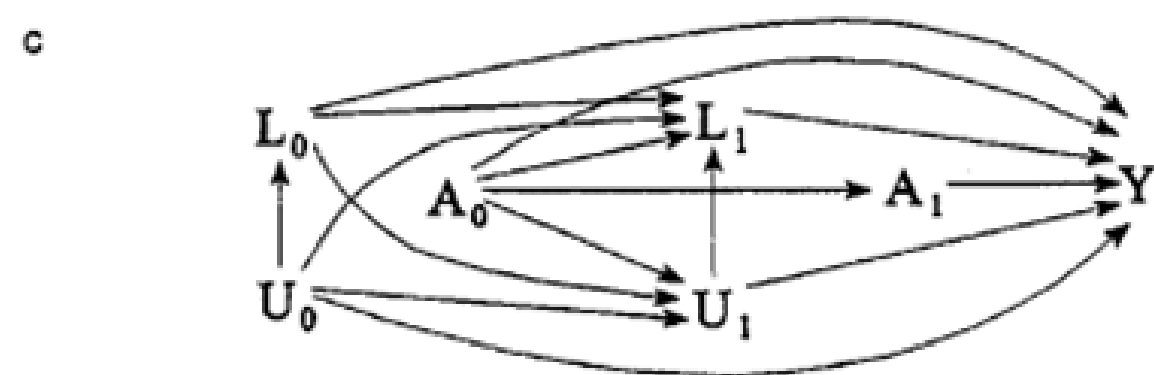
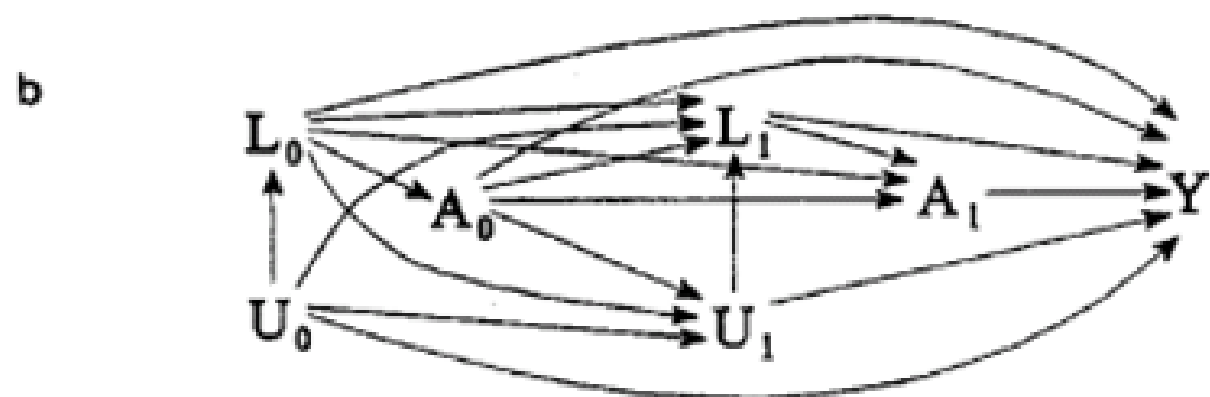
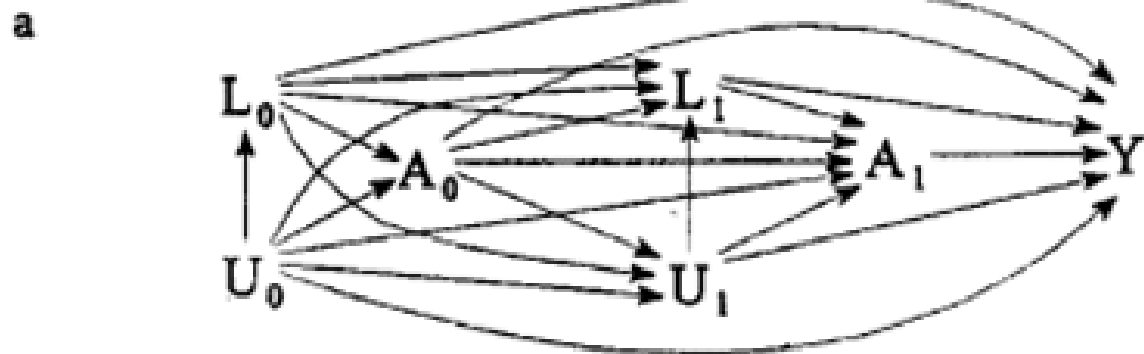


FIGURE 1. Causal graphs for a time-dependent exposure.

The average causal effect is the average over individual causal effects, where the individual effect is the contrast between counterfactual exposure states.

$Y_{a_0=1}$  = the outcome a subject would have if  $a_0 = 1$

$Y_{a_0=0}$  = the outcome a subject would have if  $a_0 = 0$

$Y_{a_0=1} - Y_{a_0=0}$  = individual causal effect on the  
difference scale

If no confounding, then:

$$\Pr[Y_{a_0=1} = 1] - \Pr[Y_{a_0=0} = 1] =$$
$$\text{cRD} = \Pr[Y=1|A_0=1] - \Pr[Y=1|A_0=0]$$

The causal RD, RR, and OR can also be expressed in terms of the parameters of the following linear, log linear, and linear logistic models for the two counterfactual probabilities

$\text{pr}(Y_{a_0=1} = 1)$  and  $\text{pr}(Y_{a_0=0} = 1)$ :

$$\text{pr}[Y_{a_0} = 1] = \psi_0 + \psi_1 a_0 \quad (1)$$

$$\log \text{pr}[Y_{a_0} = 1] = \theta_0 + \theta_1 a_0 \quad (2)$$

$$\text{logit pr}[Y_{a_0} = 1] = \beta_0 + \beta_1 a_0 \quad (3)$$

where  $Y_{a_0}$  is  $Y_{a_0=1}$  if  $a_0 = 1$  and  $Y_{a_0}$  is  $Y_{a_0=0}$  if  $a_0 = 0$ .

causal RD =  $\psi_1$ , causal RR =  $e^{\theta_1}$ , and causal OR =  $e^{\beta_1}$ .

Models 1-3 are saturated MSMs.

*Marginal* because they model the marginal distribution of the counterfactual random variables  $Y_{a0=1}$  and  $Y_{a0=0}$  rather than the joint distribution.

*Saturated*, because each has two unknown parameters so each model places no restriction on the possible values of the two unknown probabilities  $\text{pr}(Y_{a0=1}=1)$  and  $\text{pr}(Y_{a0=0}=1)$ .

These models do not include covariates because they are, by definition, models for causal effects on the entire source population, not models for observed associations.

Crude RD, RR, and OR can also be expressed in terms of the parameters of the following saturated linear, log linear, and linear logistic models for the **observed** outcome  $Y$ .

$$\text{pr}[Y=1 | A_0=a_0] = \psi'_0 + \psi'_1 a_0 \quad (4)$$

$$\log \text{pr}[Y=1 | A_0=a_0] = \theta'_0 + \theta'_1 a_0 \quad (5)$$

$$\text{logit pr}[Y=1 | A_0=a_0] = \beta'_0 + \beta'_1 a_0 \quad (6)$$

These are models for **associations** observed when comparing subpopulations (defined by levels of treatment) of the source population (i.e. not causal effects).

Crude RD =  $\psi'_1$ , crude RR =  $e^{\theta'_1}$ , and crude OR =  $e^{\beta'_1}$ . The parameters of the associational models 4-6 will differ from the parameters of the MSMs 1-3, except when treatment is unconfounded.

Because models 4-6 are models for the observed data, unbiased estimates of the model parameters can be obtained (assuming no selection bias or measurement error).

When treatment is unconfounded, these same estimates will also be unbiased for the corresponding causal parameters of models 1-3.

Under No Confounding:

causal RD =  $\psi_1$       is equivalent to      crude RD =  $\psi'_1$

causal RR =  $e^{\theta_1}$       is equivalent to      crude RR =  $e^{\theta'_1}$

causal OR =  $e^{\beta_1}$       is equivalent to      crude OR =  $e^{\beta'_1}$



## No Unmeasured Confounders

In the real world of observational epidemiologic data, exposure will be confounded, so crude (unadjusted) model of the data will not have a causal interpretation.

i.e., the crude association parameter will not equal the corresponding causal parameter.

Assume that you have no remaining unmeasured confounders given data on measured confounders  $L_0$ .

Unbiased estimates of the causal parameters  $\psi_1$ ,  $\theta_1$ , and  $\beta_1$  obtained by performing a weighted analysis.

This is the innovation of this article (weighting vs adjustment)

Specifically, using a weighted regression model each subject  $i$  is assigned a weight  $w_i$  equal to the inverse of the conditional probability of receiving the treatment that was actually received.

$$w_i = 1 / \text{pr}[A_0 = a_{0i} | L_0 = l_{0i}]$$

where,  $l_{0i}$  is the observed value of  $L_0$  for subject  $i$ .

The true weights  $w_i$  are unknown but can be estimated from the data in a preliminary logistic regression of  $A_0$  on  $L_0$ . For example, the logistic regression model:

$$\text{logit pr}[A_0 = 1 | L_0 = l_0] = \alpha_0 + \alpha_1 l_0$$

$$\text{logit } \text{pr}[A_0 = 1 | L_0 = l_0] = \alpha_0 + \alpha_1 l_0$$

Then if  $A_0$  is tx,  $L_0$  is the column vector of covariates and  $\alpha_1$  is a row vector of unknown parameters to be estimated, one obtains fitted estimates for  $\alpha_0$  and  $\alpha_1$  via standard logistic regression software. For a subject  $i$  with  $A_0 = 0$  and  $L_0 = l_{0i}$ :

$$w_i = \frac{1}{\text{pr}(A_0 = 0 | L_0 = l_{0i})} = \frac{1}{\left( \frac{e^{\alpha_0 + \alpha_1 l_{0i}}}{1 + e^{\alpha_0 + \alpha_1 l_{0i}}} \right)} = (1 + e^{\alpha_0 + \alpha_1 l_{0i}})$$

$$\text{logit } \text{pr}[A_0 = 1 | L_0 = l_0] = \alpha_0 + \alpha_1 l_0$$

On the other hand, for a subject  $i$  with  $A_0 = 1$  and  $L_0 = l_{0i}$ :

$$w_i = \frac{1}{\text{pr}(A_0 = 1 | L_0 = l_{0i})} = \frac{1 + e^{\alpha_0 + \alpha_1 l_{0i}}}{e^{\alpha_0 + \alpha_1 l_{0i}}} = (1 + e^{-\alpha_0 - \alpha_1 l_{0i}})$$

If there are no unmeasured confounders given data on  $L_0$  ( $L_0$  is a sufficient adjustment set) then one can control confounding by modifying the crude analysis by weighting each subject  $i$  by  $w_i$ .

Denominator of  $w_i$  is the probability that subject  $i$  had his/her own observed treatment. Hence "IPTW".

The effect of weighting is to create a **pseudopopulation** consisting of  $w_i$  copies of each subject  $i$ .

For example, if  $\Pr[A_0=1|L_0=l_0] = 0.25$  for a given treated subject, then  $w_i = 1/0.25 = 4$ , and so the subject contributes four copies of him/herself to the pseudopopulation.

This new pseudopopulation has the following two important properties:

- 1) in the pseudopopulation, unlike the actual population,  $A_0$  is unconfounded by the measured covariates  $L_0$ .
- 2)  $\text{pr}(Y_{a0=1} = 1)$  and  $\text{pr}(Y_{a0=0} = 1)$  in the pseudopopulation are the same as in the true study population so that the causal RD, RR, and OR are the same in both populations.

It follows that one can unbiasedly estimate the causal RD, RR, and OR by a standard **crude** analysis in the pseudopopulation.

## Appendix

Analyze data in Table A1 under the assumption of no unmeasured confounders given  $L_0$ .

TABLE A1. Observed Data from a Point-Treatment Study with Dichotomous Treatment  $A_0$ , Stratified by the Confounder  $L_0$

	$L_0 = 1$		$L_0 = 0$	
	$A_0 = 1$	$A_0 = 0$	$A_0 = 1$	$A_0 = 0$
$Y = 1$	108	24	20	40
$Y = 0$	252	16	30	10
Total	360	40	50	50

Ignore sampling variability and thus the distinction between parameters of the source population and their empirical estimates.

Under the assumption of no unmeasured confounder,  $\text{pr}(Y_{a0=1} = 1)$  is a weighted average of the  $L_0$ -stratum-specific risks among the treated with weights proportional to the distribution of  $L_0$  in the entire study population (i.e. target).

That is,  $\text{pr}(Y_{a0=1} = 1)$  is given by:

$$\sum_{l_0} \text{pr}[Y = 1 | A_0 = 1, L_0 = l_0] \text{pr}[L_0 = l_0]$$

where the sum is over the possible values of  $L_0$ .

The  $L_0$ -standardized risk of outcome in the exposed.



$$\sum_{l_0} pr[Y = 1 | A_0 = 1, L_0 = l_0] pr[L_0 = l_0]$$

Calculating from Table A1, this  $L_0$ -standardized risk in the treated group is estimated as  $pr(Y_{a0=1} = 1) =$

$$(108/360)(0.8) + (20/50)(0.2) = (0.3 * 0.8) + (0.4 * 0.2) = 0.32$$

Similarly,  $pr(Y_{a0=0} = 1)$  is the  $L_0$ -standardized risk in the untreated:

$$\sum_{l_0} pr[Y = 1 | A_0 = 0, L_0 = l_0] pr[L_0 = l_0]$$

which, from Table A1, is

$$(24/40)(0.8) + (40/50)(0.2) = (0.6 * 0.8) + (0.8 * 0.2) = 0.64.$$

If  $\text{pr}(Y_{a0=1} = 1) = 0.32$  and  $\text{pr}(Y_{a0=0} = 1) = 0.64$

It follows that:

$$\text{causal RD} = 0.32 - 0.64 = -0.32$$

$$\text{Causal RR} = 0.32 / 0.64 = 0.50$$

$$\text{Causal OR} = (0.32/0.68) / (0.64/0.36) = 0.26$$

We estimated these causal effects by **standardizing the risks**, not by standardizing the effect parameters.

Note that these differ from the crude parameters computed from Table A2 (ignoring data on the confounder  $L_0$ ):

TABLE A2. Crude Data from the Point-Treatment Study of Table A1

	$A_0 = 1$	$A_0 = 0$
$Y = 1$	128	64
$Y = 0$	282	26
Total	410	90

Thus,  $\psi_1 = -0.32$ ,  $\theta_1 = \log 0.50$ , and  $\beta_1 = \log 0.26$

in models 1-3 differ from the parameters

$$\psi'_1 = -0.40, \theta'_1 = \log 0.044, \text{ and } \beta'_1 = \log 0.18$$

of models 4-6.

It is well known that the causal RD and causal RR (but not the causal OR) are also equal to weighted averages of the stratum-specific RDs and RRs.

The causal RD equals the standardized RD (*sRD*) where:

$$sRD = \sum_{l_0} RD_{l_0} pr[L_0 = l_0]$$

and  $RD_{l_0} = pr[Y = 1 | A_0 = 1, L_0 = l_0] - pr[Y = 1 | A_0 = 0, L_0 = l_0]$  is the risk difference in stratum  $l_0$ .

$$RD_{l_0=0} = (0.4 - 0.8) = -0.4$$

$$RD_{l_0=1} = (0.3 - 0.6) = -0.3$$

$$SRD = -0.4(0.2) + -0.3(0.8) = -0.08 + -0.24 = -0.32$$

$$sRD = \sum_{l_0} RD_{l_0} pr[L_0 = l_0]$$

The traditional approach to estimating the causal RD is to calculate this  $sRD$ .

Likewise this works for the  $sRR$ :

$$RR_{l_0=0} = (0.4 / 0.8) = 0.5$$

$$RR_{l_0=1} = (0.3 - 0.6) = 0.5$$

$$sRR = 0.5(0.2) + 0.5(0.8) = 0.1 + 0.4 = 0.50$$

The IPTW method is an alternative approach to estimation of the causal RD and RR that, in contrast to the approach based on calculating the  $sRD$ , allows generalization to unsaturated MSMs in longitudinal studies with time-varying treatments.

Table A3 displays the data from the study in a different format.

It gives the number of subjects with each of the possible combinations of  $l_0$ ,  $a_0$ , and  $y$ , as well as the weight

$w = 1/\text{pr}[A_0 = a_0|L_0 = l_0]$  associated with each.

Table A3: Inverse Probability of Treatment Weights  $w$  and Composition of the Pseudopopulation in a Point-Treatment Study (Robins et al 2000)

$L_0$	$A_0$	$Y$	N Observed Population	$\text{Pr}(A_0 L_0)$	$w$	N Pseudo Population
1	1	1	108	0.9	1.11	120
1	1	0	252	0.9	1.11	280
1	0	1	24	0.1	10	240
1	0	0	16	0.1	10	160
0	1	1	20	0.5	2	40
0	1	0	30	0.5	2	60
0	0	1	40	0.5	2	80
0	0	0	10	0.5	2	20

The final column of the table represents the number of subjects in the weighted pseudopopulation for each combination of  $(l_0, a_0, y)$ . Note that the  $w_i$  need not be whole numbers or sum to 1, and so the number of subjects in the pseudopopulation can be greater than the number in the actual population.

TABLE A4. Pseudopopulation Created by Inverse Probability of Treatment Weighting from a Point-Treatment Study with Dichotomous Treatment  $A_0$ , Stratified by the Confounder  $L_0$

	$L_0 = 1$		$L_0 = 0$	
	$A_0 = 1$	$A_0 = 0$	$A_0 = 1$	$A_0 = 0$
$Y = 1$	120	240	40	80
$Y = 0$	280	160	60	20
Total	400	400	100	100

Note that  $A_0 \perp\!\!\!\perp L_0$  in this Table!

It can be seen that  $L_0$  and  $A_0$  are unassociated in the pseudopopulation, which implies that  $A_0$  is unconfounded.

The lack of association between  $L_0$  and  $A_0$  implies that in the pseudopopulation, the  $L_0$ -standardized risk in the treated equals the crude risk  $\text{pr}(Y = 1 | A_0 = 1) = 0.32$  and the  $L_0$ -standardized risk in the untreated equals the crude risk  $\text{pr}(Y = 1 | A_0 = 0) = 0.64$ .

This means that we can calculate the crude RD or RR in Table A4 and it will have a causal interpretation:

TABLE A5. Crude Data from the Pseudopopulation of Table A4

	$A_0 = 1$	$A_0 = 0$
$Y = 1$	160	320
$Y = 0$	340	180
Total	500	500



```
. csi 160 320 340 180
```

	Exposed	Unexposed	Total
Cases	160	320	480
Noncases	340	180	520
Total	500	500	1000
Risk	0.32	0.64	0.48
	Point estimate		
Risk difference	-0.32		
Risk ratio	0.50		

It follows under assumption of no unmeasured confounding given  $L_0$ , that crude RD, RR, and OR in the pseudopopulation equal causal RD, RR, and OR in the actual population.

IPTW analysis estimates crude parameters of the pseudo-Population, thus causal parameters of the actual population.

## Stabilized Weights

The probabilities  $\text{pr}[A_0 = a_{0i} | L_0 = l_{0i}]$  may vary greatly between subjects when components of  $L_0$  are strongly associated with  $A_0$ .

This variability can result in extremely large values of the weight  $w_i$  for a few subjects, and these few subjects will contribute a very large number of copies of themselves to the pseudopopulation and thus will dominate the weighted analysis, with the result that the IPTW estimator will have a large variance and a potentially skewed distribution.

For unsaturated MSMs, this variability can be somewhat mitigated by replacing the weight  $w_i$  by the stabilized weight:  $sw_i = \text{pr}[A_0 = a_{0i}] / \text{pr}[A_0 = a_{0i} | L_0 = l_{0i}]$ .

For example, suppose  $A_0$  was unconfounded so that  $A_0$  and  $L_0$  are unassociated and  $\text{pr}[A_0 = a_{0i}] = \text{pr}[A_0 = a_{0i} | L_0 = l_{0i}]$ .

Then  $sw_i = 1$ , and each subject contributes the same weight.

When  $A_0$  is confounded,  $sw_i$  will not be constant but will vary around 1, depending on a subject's value of  $L_0$ .

$sw_i$  will still tend to be much less variable than  $w_i$ .

When using the weight  $sw_i$  rather than the weight  $w_i$  the estimates of the parameters of an MSM remain unbiased and will generally be less variable.

$\text{pr}[A_0 = a_{0i}]$  and  $\text{pr}[A_0 = a_{0i} | L_0 = l_{0i}]$  are unknown and must be estimated.

$\text{pr}[A_0 = a_{0i}]$  can be estimated as the proportion of subjects in the study sample with  $A_0$  equal to  $a_{0i}$ .

Table A3 displayed the appendix data for number of subjects with each of the possible combinations of  $l_0$ ,  $a_0$ , and  $y$ , as well as the weight  $1/\text{pr}[A_0 = a_0 | L_0 = l_0]$  associated with each.

$$\text{pr}[A_0 = 1] = (108+252+20+30)/500 = 410/500 = 0.82$$

$$\text{pr}[A_0 = 0] = (24+16+40+10)/500 = 90/500 = 0.18$$

<u><math>L_0</math></u>	<u><math>A_0</math></u>	<u><math>Y</math></u>	<u><math>N</math></u>	<u><math>\text{pr}(A_0   L_0)</math></u>	<u><math>w</math></u>	<u><math>sw</math></u>	weighted <u>Pseudo N</u>	stabilized <u>Pseudo N</u>
1	1	1	108	0.9	1.11	0.9111	120	98.4
1	1	0	252	0.9	1.11	0.9111	280	229.6
1	0	1	24	0.1	10	1.8	240	43.2
1	0	0	16	0.1	10	1.8	160	28.8
0	1	1	20	0.5	2	1.64	40	32.8
0	1	0	30	0.5	2	1.64	60	49.2
0	0	1	40	0.5	2	0.36	80	14.4
0	0	0	10	0.5	2	0.36	20	3.6

New Stabilized Table A5: Crude Data from the Stabilized Pseudopopulation of Table A4

	<u>A<sub>0</sub>=1</u>	<u>A<sub>0</sub>=0</u>
Y=1	131.2	57.6
Y=0	278.8	32.4
Total	410	90

. csi 1312 576 2788 324

	Exposed	Unexposed	Total
Cases	1312	576	1888
Noncases	2788	324	3112
Total	4100	900	5000
Risk	0.32	0.64	0.3776
	Point estimate		
Risk difference	-0.32		
Risk ratio	0.50		

## Limitations of Marginal Structural Models

IPTW estimators will be biased and thus MSMs should not be used in studies in which at each time  $k$  there is a covariate level  $l_k$  such that all subjects with that level of the covariate are certain to receive the identical treatment  $a_k$ .

For example, suppose that in some stratum, treatment is impossible given the covariates.

In this instance  $\text{pr}[A_0 = a_k | L_0 = l_k] = 0$

Which implies that  $1/\text{pr}[A_0 = a_k | L_0 = l_k] = 1/0 = \text{undefined}$

As probability of tx  $\rightarrow 0$ , weight  $\rightarrow \text{infinity}$ .

## Limitations of Marginal Structural Models

This implies that MSMs should not be used in occupational cohort studies.

For example, take an occupational cohort study in which  $A_k$  is level of exposure to an industrial chemical at time  $k$  and  $L_k = 1$  if subject is off work at time  $k$  and  $L_k = 0$  otherwise. Then all subjects with  $L_k = 1$  have  $A_k = 0$ , because all subjects off work are necessarily unexposed.

So weights go to infinity.

Similarly, in a study of the effect of screening on mortality from cervical cancer, women who have had their cervix operatively removed by time  $k$  so that  $L_k = 0$  cannot receive exposure (screening) at that time, so MSMs should not be used.

# “Doubly Robust” estimator (Robins *JASA* 1994):

Combines IPT weighting with regression adjustment within treatment categories:

Remains a consistent causal estimator if either:

1) the exposure weighting model is correctly specified but the outcome regression models is not,

or...

2) the outcome regression model is correctly specified but the exposure weighting model is not.

Lunceford & Davidian *Stat Med* 2004



# Return to the example NSW data set used earlier:

```
import delimited "C:\Users\jkaufm5\propensity.csv"
quietly logit treat age educ black##married hisp nodegr re74
           re75 u74 u75 age_2 educ_2 re74_2 re75_2
```

```
predict ps
replace ps = (ps*treat) + ((1-ps)*(1-treat))
(260 real changes made)
```

```
gen ipt = 1/ps
sum ipt, detail
```

ipt				
-----				
	Percentiles	Smallest		
1%	1.212611	1.125102		
5%	1.330638	1.127132		
10%	1.480171	1.202106	Obs	445
25%	1.522238	1.210468	Sum of Wgt.	445
50%	1.816627		Mean	1.998664
		Largest	Std. Dev.	.6399586
75%	2.30197	4.985171		
90%	2.857443	5.323654	Variance	.4095471
95%	3.001669	5.432951	Skewness	1.861643
99%	4.28544	5.471707	Kurtosis	8.71183



**\*test standardized difference without weighting\***

```
mean re75_std, over(treat)
```

```
lincom [re75_std]1 - [re75_std]0
```

Mean	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
(1)	.0841478	.0968099	0.87	0.385	-.1061148 .2744104

```
*test standardized difference with weighting*
```

```
mean re75_std [pw=ipt], over(treat)
```

```
lincom [re75_std]1 - [re75_std]0
```

Mean	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
(1)	-.0028258	.0971702	-0.03	0.977	-.1937964 .1881448

```
reg re78 treat
```

Linear regression

Number of obs = 445

re78	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
treat	1794.343	632.8536	2.84	0.005	550.5749	3038.111
_cons	4554.802	408.046	11.16	0.000	3752.856	5356.749

```
reg re78 treat [pw=ipt], cluster(id)
```

Linear regression

Number of obs = 445  
(Std. Err. adjusted for 445 clusters in id)

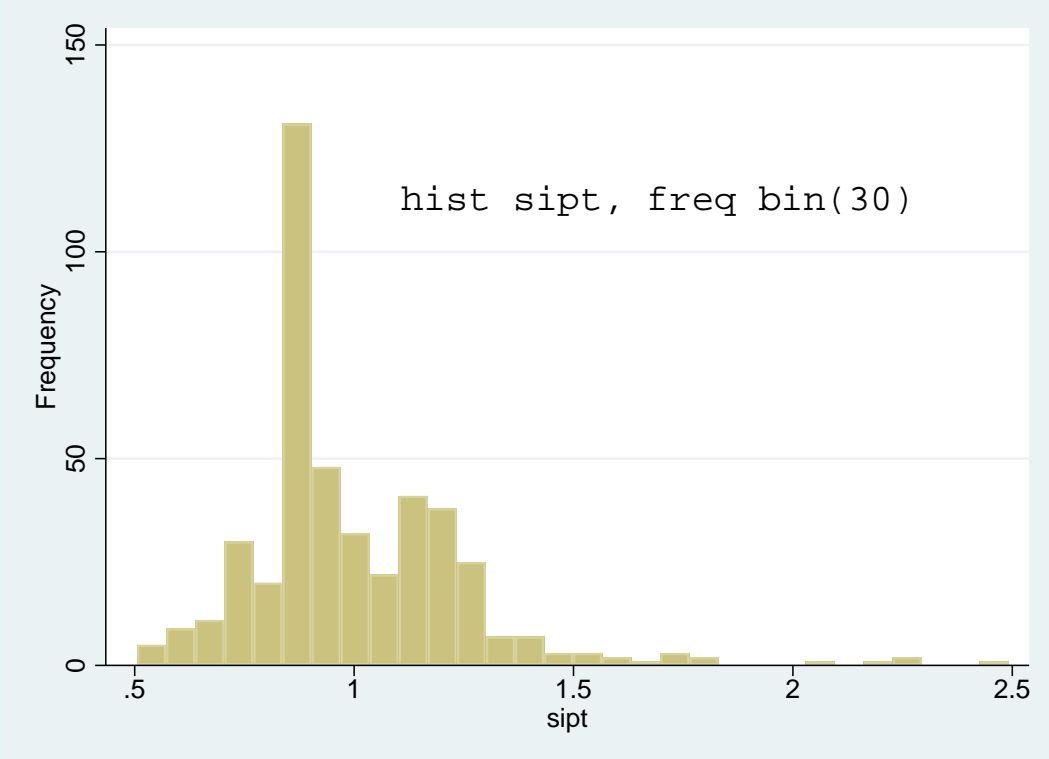
re78	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
treat	1626.067	687.8781	2.36	0.019	274.1651	2977.968
_cons	4540.558	344.7552	13.17	0.000	3863.003	5218.112

sum treat

Variable	Obs	Mean	Std. Dev.	Min	Max
treat	445	.4157303	.4934022	0	1

```
local p = r(mean)
gen sipt = (`p'*ipt)*treat + (((1-`p')*ipt)*(1-treat))
sum sipt
```

Variable	Obs	Mean	Std. Dev.	Min	Max
sipt	445	.9991493	.2533877	.5056734	2.492409



```
reg re78 treat [pw=sipt], cluster(id)
(sum of wgt is 444.62
```

Linear regression Number of obs = 445  
(Std. Err. adjusted for 445 clusters in id)

re78		Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
treat		1626.067	687.8781	2.36	0.019	274.1651	2977.968
_cons		4540.558	344.7552	13.17	0.000	3863.003	5218.112

```
teffects ipw (re78) (treat age educ black##married hisp nodegr re74
re75 u74 u75 age2 educ2 re742 re752), nopvalues
```

Treatment-effects estimation Number of obs = 445  
Estimator : inverse-probability weights  
Outcome model : weighted mean  
Treatment model: logit

re78		Coef.	Robust Std. Err.	[95% Conf. Interval]	
ATE		1626.067	655.6775	340.9623	2911.171
POmean		4540.558	337.0657	3879.921	5201.194

Warning: convergence not achieved

NSW example has point treatment with binary exposure where MSMs have no particular advantage over PS.

The real benefit of MSMs is for generalizations such as longitudinal data or continuous exposure.

Turn next to example in Hernán & Robins Chapter 12

Causal question “what is the average causal effect of smoking cessation on body weight gain?”

Use data from the NHEFS (National Health and Nutrition Examination Survey Data I Epidemiologic Follow-up Study)

Dataset:

[https://cdn1.sph.harvard.edu/wp-content/uploads/sites/1268/2012/10/nhefs\\_stata.zip](https://cdn1.sph.harvard.edu/wp-content/uploads/sites/1268/2012/10/nhefs_stata.zip)

Goal is to estimate the average causal effect of smoking cessation (treatment)  $A$  on weight gain (outcome)  $Y$ .

Use 1566 cigarette smokers aged 25-74 years who had NEHFS baseline visit and follow-up about 10 years later.

Treated  $A = 1$  if they reported having quit smoking before follow-up visit, and untreated  $A = 0$  otherwise.

Weight gain  $Y$  measured (in kg) as the body weight at follow-up visit minus body weight at baseline.

Most people gained weight, but quitters gained more weight on average.

Average weight gain was 4.5 kg in quitters, and 2.0 kg in non-quitters.

$E[Y^{a=1}]$  is mean weight gain that would have been observed if all individuals had quit smoking before follow-up visit.

$E[Y^{a=0}]$  is mean weight gain that would have been observed if all individuals in the population had not quit smoking.

ACE on the difference scale is  $E[Y^{a=1}] - E[Y^{a=0}]$

The associational difference  $E[Y|A=1] - E[Y|A=0] = 4.5 - 2.0 = 2.3$  kg is generally different from the causal difference  $E[Y^{a=1}] - E[Y^{a=0}]$

Hernán & Robins select 9 potential baseline confounders: sex (0: M, 1: F), age (years), race (0: white, 1: other), education (5 categories), intensity and duration of smoking (cigarettes per day and years of smoking), physical activity in daily life (3 categories), recreational exercise (3 categories), and weight (in kg).



IP weighting creates a pseudo-population in which the arrow from the confounders  $L$  to the treatment  $A$  is removed.

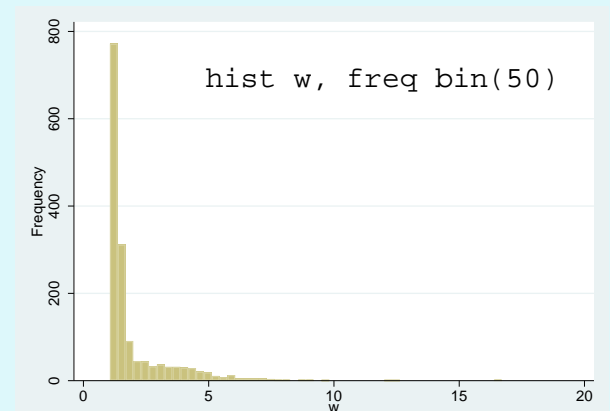
If the confounders  $L$  are sufficient to block all backdoor paths from  $A$  to  $Y$ , then all confounding is eliminated in the pseudo-population.

Then association between  $A$  and  $Y$  in the pseudo-population consistently estimates the causal effect of  $A$  on  $Y$ .

Use logistic regression:

```
logit qsmk sex race c.age##c.age ib(last).education c.wt71##c.wt71
c.smokeintensity##c.smokeintensity ib(last).exercise
ib(last).active c.smokeyrs##c.smokeyrs
predict p_qsmk, pr
```

```
gen w=.
replace w=1/p_qsmk if qsmk==1
replace w=1/(1-p_qsmk) if qsmk==0
summarize w, det
```



The estimated IP weights ranged from 1.05 to 16.7, with mean = 2.00.

Since average weight was 2, the effective sample size doubled.

Association is causation in the pseudopopulation, so fit an associational model in the inverse weighted data:

```
regress wt82_71 qsmk [pweight=w], cluster(seqn)
```

Number of obs = 1,566

wt82_71	Coef.	Robust SE	t	P> t	[95% Conf. Interval]	
qsmk	3.440535	.5258294	6.54	0.000	2.409131	4.47194
_cons	1.779978	.2248742	7.92	0.000	1.338892	2.221065

```
tab sex qsmk [aw=w]
```

Causal estimate = 3.4 kg (95CI 2.4, 4.5)

sex	quit smoking between baseline and 1982			Total
	No smokin	Smoking c		
0	382.51458	382.52246		765.03704
1	401.62064	399.34231		800.96296
Total	784.13522	781.86478		1,566

A  $\perp\!\!\!\perp$  L in pseudopopulation

Why not create a pseudo-population in which  $A$  and  $Y$  are independent but in which the average weight is 1 instead of 2.

## Stabilized weight!

This pseudo-population would be the same size as the study population.

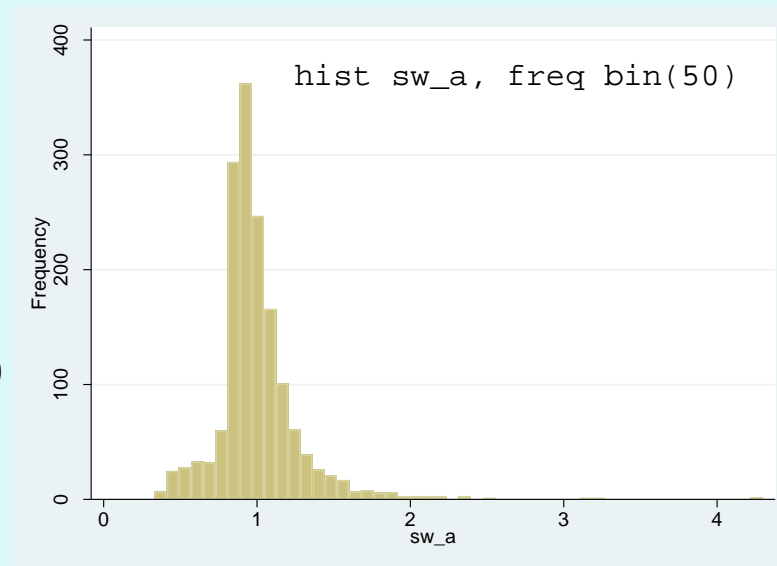
Assign to the treated subjects a probability  $\Pr[A=1]$  of receiving treatment.

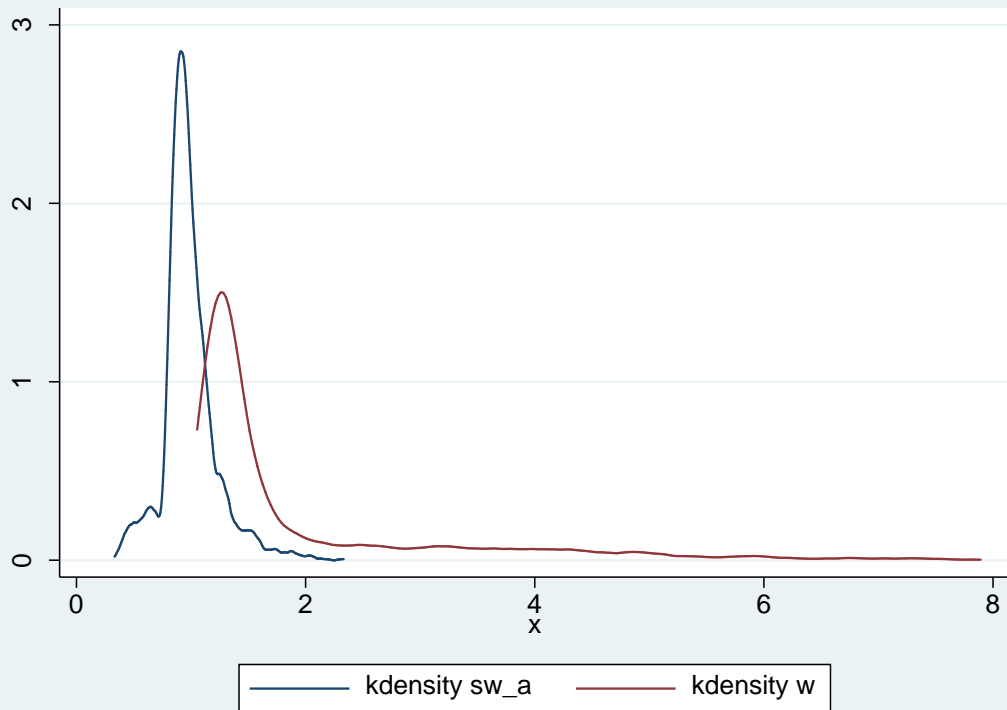
Assign to untreated subjects a probability  $\Pr[A=0]$  of not receiving treatment.

```
/* estimation of denominator of ip weights*/  
logit qsmk sex race c.age##c.age ib(last).education ib(last).active  
c.wt71##c.wt71 c.smokeintensity##c.smokeintensity c.smokeyrs##c.smokeyrs  
ib(last).exercise  
predict pd_qsmk, pr
```

```
/* estimation of numerator of ip weights*/  
logit qsmk  
predict pn_qsmk, pr  
  
gen sw_a=.  
replace sw_a=pn_qsmk/pd_qsmk if qsmk==1  
replace sw_a=(1-pn_qsmk)/(1-pd_qsmk) if qsmk==0  
summarize sw_a, det
```

Mean weight = 1.00





```
twoway (kdensity sw_a if w < 8) (kdensity w if w < 8)
```

```
regress wt82_71 qsmk [pweight=sw_a], cluster(seqn)
(sum of wgt is 1,564.19025221467)
```

Number of obs = 1,566

-----							
wt82_71	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]		
qsmk	3.440535	.5258294	6.54	0.000	2.409131	4.47194	
_cons	1.779978	.2248742	7.92	0.000	1.338892	2.221065	
-----							

In fact, not much change in CI width in this example.

## The advantages of working with an unconditional estimate:

Estimate the causal effect of quitting smoking  $A$  on the risk of death  $D$  by 1992.

$A=1$  defined as a smoker who quit

$D=1$  defined as a participant who died by 1992

Model is:  $\ln \left( \frac{\Pr(D^a=1)}{\Pr(D^a=0)} \right) = \beta_0 + \beta_1 A$

Therefore, the effect parameter is  $e^{\beta_1}$  = the causal odds ratio of death contrasting quitting versus not quitting.

Estimate this causal model by fitting an associational model in the pseudo-population:

$$\ln \left( \frac{\Pr(D = 1|A)}{\Pr(D = 0|A)} \right) = b_0 + b_1 A$$

First, estimate the stabilized weights `sw_a`

```
logit qsmk sex race c.age##c.age ib(last).education  
      c.smokeintensity##c.smokeintensity c.smokeyrs##c.smokeyrs ib(last).exercise  
      ib(last).active c.wt71##c.wt71
```

Then obtain propensity scores:

```
predict pd_qsmk, pr
```

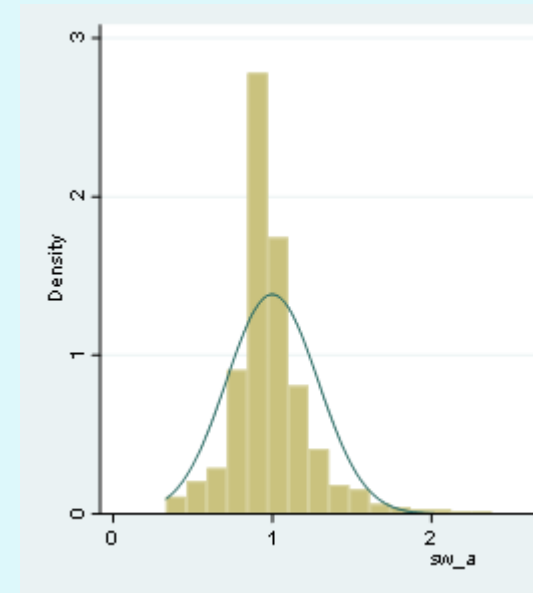
Then get the numerator of the stabilized weight:

```
logit qsmk  
predict pn_qsmk, pr
```

Then construct the stabilized weight as  $\Pr(A) / \Pr(A|L)$ :

```
gen sw_a=.  
replace sw_a=pn_qsmk/pd_qsmk if qsmk==1  
replace sw_a=(1-pn_qsmk)/(1-pd_qsmk) if qsmk==0  
sum sw_a, det  
hist sw_a, norm
```

We confirm that weights are symmetrically distributed around a mean of 1.00



Now fit the weighted logistic regression model, which estimates the associational parameters in the pseudopopulation:

```
logistic death qsmk [pweight=sw_a], cluster(seqn) cformat(%6.2f)
```

Logistic regression				Number of obs		=	1,566	
-----								
death	Odds Ratio	Robust SE	z	P> z	[95% Conf. Interval]			
-----+-----								
qsmk	1.03	0.16	0.19	0.848	0.76	1.40		
_cons	0.23	0.02	-18.88	0.000	0.19	0.26		
-----								

Note: \_cons estimates baseline odds.

Therefore we conclude that on the odds ratio scale, the effect of quitting smoking on short-term mortality is null: OR = 1.0 (95% CI: 0.8, 1.4).

This is very different from the crude association:

```
logistic death qsmk, cformat(%6.2f)
```

-----						
death	Odds Ratio	Stand Err.	z	P> z	[95% Conf. Interval]	
-----+-----						
qsmk	1.40	0.20	2.39	0.017	1.06	1.86
-----						

Observed OR = 1.4 (95% CI: 1.1, 1.9).

This causal estimate is not very exciting (because it's null), but it is still good to note that any GLM will work here, not just logistic.

This allows us to avoid the OR (which has no causal interpretation due to non-collapsibility concerns explained earlier).

For example, use binomial regression to get causal RR:

```
binreg death qsmk [pweight=sw_a], rr cluster(seqn) cformat(%6.2f)
```

```
Generalized linear models                No. of obs      =          1,566
Variance function: V(u) = u*(1-u)        [Bernoulli]
Link function      : g(u) = ln(u)         [Log]
```

-----						
	Semirobust					
death		Risk Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
-----+-----						
qsmk		1.02	0.13	0.19	0.848	0.80 1.32
_cons		0.18	0.01	-26.28	0.000	0.16 0.21
-----						

Note: \_cons estimates baseline risk.

For a large effect magnitude, the OR and RR would be very different.



Or risk difference:

```
binreg death qsmk [pweight=sw_a], rd cluster(segn) cformat(%6.2f)
```

Generalized linear models	No. of obs	=	1,566
Variance function: $V(u) = u*(1-u)$	[Bernoulli]		
Link function : $g(u) = u$	[Identity]		

	Semirobust					
death	Risk Diff.	Std. Err.	z	P> z	[95% Conf. Interval]	
qsmk	0.00	0.02	0.19	0.849	-0.04	0.05
_cons	0.18	0.01	15.52	0.000	0.16	0.21

Or the standardized absolute risks in the exposed and unexposed:

```
binreg death i.qsmk [pweight=sw_a], rd cluster(seqn) cformat(%6.4f)
margins i.qsmk, cformat(%6.2f)
```

```
Adjusted predictions      Number of obs      =      1,566
Expression   : Predicted mean death, predict()
```

	Delta-method					
	Margin	Std. Err.	z	P> z	[95% CI]	
qsmk=0	0.1839	0.0118	15.52	0.000	0.1606	0.2071
qsmk=1	0.1884	0.0208	9.05	0.000	0.1476	0.2292

---

## Education Corner

# An introduction to g methods

**Ashley I Naimi<sup>1\*</sup>, Stephen R Cole<sup>2</sup> and Edward H Kennedy<sup>3</sup>**

<sup>1</sup>Department of Epidemiology, University of Pittsburgh, <sup>2</sup>Department of Epidemiology, University of North Carolina at Chapel Hill and <sup>3</sup>Department of Statistics, Carnegie Mellon University

\*Corresponding author. Department of Epidemiology University of Pittsburgh 130 DeSoto Street Pittsburgh, PA 15261 [ashley.naimi@pitt.edu](mailto:ashley.naimi@pitt.edu)

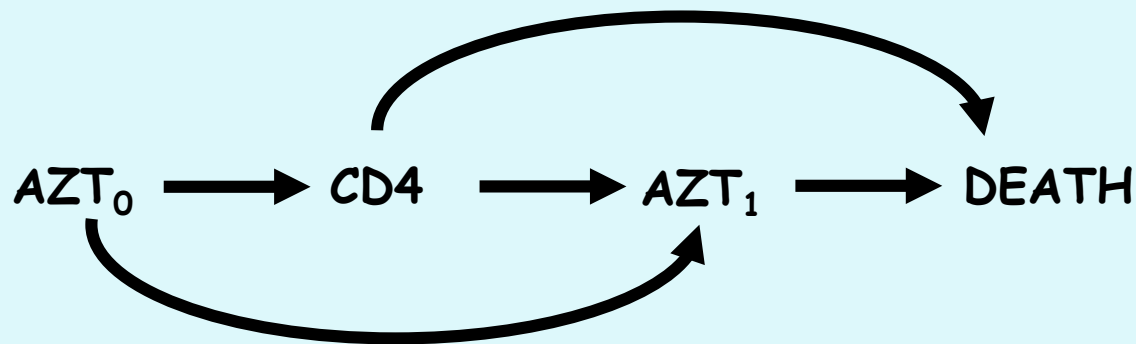


Robins' named his models "g methods" for "general", to enable the identification and estimation of the effects of all kinds of exposures or intervention plans.

A family of methods that include the g formula, marginal structural models, and structural nested models.

Traditional regression models require no feedback between time-varying treatments and time-varying confounders.

Before 1986, no solution to this general problem.



Naimi example concerns the effect of treatment for HIV on (continuous) CD4 count.

Table 1 presents data from a hypothetical cohort study ( $A=1$  for treated,  $A=0$  otherwise).

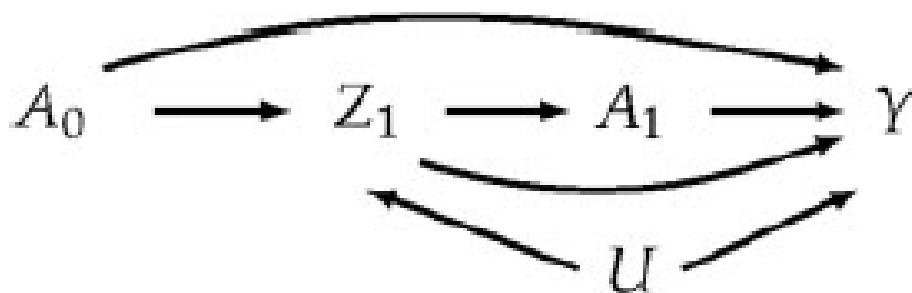
$A_0$	$Z_1$	$A_1$	$Y$	$N$
0	0	0	87.29	209,271
0	0	1	112.11	93,779
0	1	0	119.65	60,654
0	1	1	144.84	136,293
1	0	0	105.28	134,781
1	0	1	130.18	60,789
1	1	0	137.72	93,903
1	1	1	162.83	210,527

Treatment measured at baseline ( $A_0$ ) and at follow up ( $A_1$ ).

One covariate is elevated HIV viral load ( $Z=1$  means  $>200$  copies/ml,  $Z=0$  otherwise)

Outcome  $Y$  is CD4 count at the end of follow up in cells/mm<sup>3</sup>.

$N$  is large, so no sampling variability issues or measures.



**Figure 1.** Causal diagram representing the relation between anti-retroviral treatment at time 0 ( $A_0$ ), HIV viral load just prior to the second round of treatment ( $Z_1$ ), anti-retroviral treatment status at time 1 ( $A_1$ ), the CD4 count measured at the end of follow-up ( $Y$ ), and an unmeasured common cause ( $U$ ) of HIV viral load and CD4.

Based on the DAG,  $E(Y|A_0, A_1, Z)$  may be composed of several parts:

effects of  $A_0$ ,  $Z$ , and  $A_1$ ;

3 two-way interactions between  $A_0$ ,  $Z$ , and  $A_1$ ;

1 three-way interactions between  $A_0$ ,  $Z$ , and  $A_1$ ;

Naimi focuses on the ACE of always taking treatment ( $a_0=a_1=1$ ) compared to never taking treatment ( $a_0=a_1=0$ ).

$$\begin{aligned}\varphi &= E(Y^{a_0=1, a_1=1}) - E(Y^{a_0=0, a_1=0}) \\ &= E(Y^{a_0=1, a_1=1} - Y^{a_0=0, a_1=0})\end{aligned}$$

This ACE is a marginal effect because it does not compute contrasts within covariate strata and then combine these (like regression or Mantel-Haenszel)

Write the effect as

$$E(Y^{a_0, a_1}) - E(Y^{0,0}) = \varphi_0 a_0 + \varphi_1 a_1 + \varphi_2 a_0 a_1$$

i.e. that the ACE  $\varphi$  may be composed of two exposure main effects  $\varphi_0$  and  $\varphi_1$  their two-way interaction  $\varphi_2$ .

This marginal effect differs from a conditional effect that is calculated within covariate strata.

For example, a conditional effect would make sense if you were interested in effect measure modification by  $Z$ .

Conditioning on  $Z=1$  would answer the question: what is the effect of  $A_0$  and  $A_1$  in those with high viral load?

But for now, effect of interest is just  $\varphi = \varphi_0 + \varphi_1 + \varphi_2$ .

# Standard Methods

```
input a0 z a1 y n
0 0 0 87.29 209271
0 0 1 112.11 93779
0 1 0 119.65 60654
0 1 1 144.84 136293
1 0 0 105.28 134781
1 0 1 130.18 60789
1 1 0 137.72 93903
1 1 1 162.83 210527
end
```

```
expand n
gen avga = (a0+a1)/2
reg y avga, nohead cformat(%6.2f)
reg y avga z, nohead cformat(%6.2f)
reg y a0, nohead cformat(%6.2f)
reg y a0 z, nohead cformat(%6.2f)
reg y a1, nohead cformat(%6.2f)
reg y a1 z, nohead cformat(%6.2f)
```



## Standard Methods

Model	Estimate of $\beta_1$
$\beta_0 + \beta_1(A_0 + A_1)/2$	60.9
$\beta_0 + \beta_1(A_0 + A_1)/2 + \beta_2Z$	42.6
$\beta_0 + \beta_1A_0$	27.1
$\beta_0 + \beta_1A_0 + \beta_2Z$	18.0
$\beta_0 + \beta_1A_1$	38.9
$\beta_0 + \beta_1A_1 + \beta_2Z$	25.0

In the first model,  $\beta_1 = 60.9$  cells/mm<sup>3</sup> is the crude difference in mean CD4 count for the always treated compared to the never treated.

Second model gives 42.6 cells/mm<sup>3</sup> for the Z-adjusted difference in mean CD4 count for the same contrast.

## Estimate the MSM using IPTW:

MSM is:  $E(Y^{a_0, a_1}) = \beta_0 + \varphi_0 a_0 + \varphi_1 a_1 + \varphi_2 a_0 a_1$

where  $\beta_0 = E(Y^{0,0})$  is an intercept parameter and the  
ACE  $\varphi = E(Y^{1,1} - Y^{0,0}) = \varphi_0 + \varphi_1 + \varphi_2$

First obtain predicted probabilities of the observed tx.

There are 2  $A_1$  values (1 vs 0) for each of 4 levels of  $Z \times A_0$ .

Additionally, there are 2 possible  $A_0$  values (1 vs 0).

Four possible exposure regimes ( $A_0, A_1$ ): 00, 10, 01, and 11

For each  $Z$  value, obtain from the table the predicted probability of the exposure that was actually received.

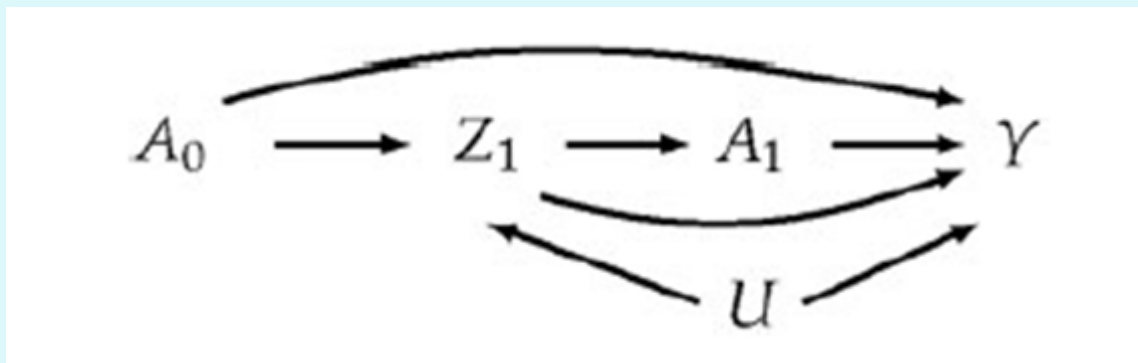
## Estimate the MSM using IPTW (cont):

No variables affect  $A_0$ , so this probability is 0.5 for all individuals in the sample.

According to the DAG,  $A_1$  is not directly affected by  $A_0$

Therefore the  $Z$ -specific probabilities of  $A_1$  are constant across levels of  $A_0$  (because  $A_1 \perp\!\!\!\perp A_0 \mid Z$ )

In a DAG where  $A_0$  affects  $A_1$ , the  $Z$ -specific probabilities of  $A_1$  would vary across levels of  $A_0$ .



# Estimate the MSM using IPTW (cont):

\* IPTW

```
quietly sum a1
```

```
scalar pa1 = r(mean)
```

```
quietly sum a1 if z==0
```

```
scalar pa1z0 = r(mean)
```

```
quietly sum a1 if z==1
```

```
scalar pa1z1 = r(mean)
```

```
scalar w000 = (0.5 * (1-pa1z0))
```

```
scalar w001 = (0.5 * (pa1z0))
```

```
scalar w010 = (0.5 * (1-pa1z1))
```

```
scalar w011 = (0.5 * (pa1z1))
```

```
scalar w100 = (0.5 * (1-pa1z0))
```

```
scalar w101 = (0.5 * (pa1z0))
```

```
scalar w110 = (0.5 * (1-pa1z1))
```

```
scalar w111 = (0.5 * (pa1z1))
```

```
scalar sw000 = (0.5*(1-pa1))/(0.5 * (1-pa1z0))
```

```
scalar sw001 = (0.5*pa1)/(0.5 * (pa1z0))
```

```
scalar sw010 = (0.5*(1-pa1))/(0.5 * (1-pa1z1))
```

```
scalar sw011 = (0.5*pa1)/(0.5 * (pa1z1))
```

```
scalar sw100 = (0.5*(1-pa1))/(0.5 * (1-pa1z0))
```

```
scalar sw101 = (0.5*pa1)/(0.5 * (pa1z0))
```

```
scalar sw110 = (0.5*(1-pa1))/(0.5 * (1-pa1z1))
```

```
scalar sw111 = (0.5*pa1)/(0.5 * (pa1z1))
```

**Table 4.** Stabilized inverse probability weights and Pseudo-population obtained by using inverse probability weights

$A_0$	$Z_1$	$A_1$	$Y$	$sw$	Pseudo $N$
0	0	0	87.23	0.72	151222.84
0	0	1	112.23	1.62	151680.46
0	1	0	119.79	1.62	98110.06
0	1	1	144.78	0.72	98789.40
1	0	0	105.25	0.72	97395.08
1	0	1	130.25	1.62	98321.62
1	1	0	137.80	1.62	151884.02
1	1	1	162.80	0.72	152596.51

```

gen sw = .
replace sw = (0.5*(1-pa1))/(0.5 * (1-pa1z0)) if a0==0 & z==0 & a1==0
replace sw = (0.5*pa1)/(0.5 * (pa1z0)) if a0==0 & z==0 & a1==1
replace sw = (0.5*(1-pa1))/(0.5 * (1-pa1z1)) if a0==0 & z==1 & a1==0
replace sw = (0.5*pa1)/(0.5 * (pa1z1)) if a0==0 & z==1 & a1==1
replace sw = (0.5*(1-pa1))/(0.5 * (1-pa1z0)) if a0==1 & z==0 & a1==0
replace sw = (0.5*pa1)/(0.5 * (pa1z0)) if a0==1 & z==0 & a1==1
replace sw = (0.5*(1-pa1))/(0.5 * (1-pa1z1)) if a0==1 & z==1 & a1==0
replace sw = (0.5*pa1)/(0.5 * (pa1z1)) if a0==1 & z==1 & a1==1

reg y i.a0##i.a1 [pw=sw], nohead cformat(%6.2f)

```

y	Coef.
1.a0	25.02
1.a1	25.00
a0#a1	
1 1	-0.01
_cons	100.02

thus  $ACE \varphi = E(Y^{1,1} - Y^{0,0}) = \varphi_0 + \varphi_1 + \varphi_2 = 25 + 25 + 0 = 50$

Weighting the observed data by the inverse of the probability of observed exposure yields a “pseudo-population” in which treatment at the second time point ( $A_1$ ) is no longer related to and thus no longer confounded by previous viral load ( $Z$ ).

Weighting a regression model for the outcome by the inverse probability of treatment enables us to account for the fact that  $Z$  both confounds  $A_1$  and is affected by  $A_0$ .

```
table a1 a0 z [pweight = sw], contents(freq )
```

-----					
a1		z and a0			
		----- 0 -----		----- 1 -----	
		0	1	0	1
-----+-----					
0		151,222	97,394.8	98,106.2	151,886
1		151,681	98,321.9	98,789.1	152,596
-----					

## Now Replicate this Estimate Using the G-Formula

Start with a mathematical representation of the data generating mechanism (the joint density of the observed data).

Factor the joint density in a way that respects the temporal ordering of the data by conditioning each variable on its history.

$$f(y, a_1, z, a_0) = f(y|a_1, z, a_0)\Pr(a_1|z, a_0)\Pr(z|a_0)\Pr(a_0)$$



Marginalize over the distribution of  $A_1$ ,  $Z$  and  $A_0$  to get the marginal mean of  $Y$ :

$$E(Y) = \sum_{a_1, z, a_0} E(Y|a_1, z, a_0) \Pr(a_1|z, a_0) \Pr(z|a_0) \Pr(a_0)$$

Under intervention on  $A_1$  and  $A_0$ , however, these are certain to take values of 1 or 0, and so these probabilities  $\Pr(a_1|z, a_0)$  and  $\Pr(a_0)$  are struck from the factorization:

This leaves only:

$$E(Y^{a_1, a_0}) = \sum_z E(Y|a_1, z, a_0) \Pr(z|a_0)$$

This equation is the so-called “g-formula”

In the Naimi example therefore, the expected CD4 count under exposure vector 0,0 can be calculated as:

$$E(Y^{0,0}) = E(Y|A_1 = 0, Z = 1, A_0 = 0) \Pr(Z = 1|A_0 = 0) + \\ E(Y|A_1 = 0, Z = 0, A_0 = 0) \Pr(Z = 0|A_0 = 0)$$

Weighting the observed outcome's conditional expectation by the conditional probability that  $Z=z$  accounts for the fact that  $Z$  is itself affected by  $A_0$ , but also confounds the effect of  $A_1$  on  $Y$ .

Likewise for exposure vector 1,1:

$$E(Y^{1,1}) = E(Y|A_1 = 1, Z = 1, A_0 = 1) \Pr(Z = 1|A_0 = 1) + \\ E(Y|A_1 = 1, Z = 0, A_0 = 1) \Pr(Z = 0|A_0 = 1)$$

And similarly for any other combination of  $A_0$  and  $A_1$

```

* g-formula
quietly sum y if a1==0 & z==0 & a0==0
scalar y000 = r(mean)
quietly sum y if a1==0 & z==0 & a0==1
scalar y001 = r(mean)
quietly sum y if a1==0 & z==1 & a0==0
scalar y010 = r(mean)
quietly sum y if a1==1 & z==0 & a0==0
scalar y100 = r(mean)
quietly sum y if a1==0 & z==1 & a0==1
scalar y011 = r(mean)
quietly sum y if a1==1 & z==1 & a0==0
scalar y110 = r(mean)
quietly sum y if a1==1 & z==0 & a0==1
scalar y101 = r(mean)
quietly sum y if a1==1 & z==1 & a0==1
scalar y111 = r(mean)
quietly sum z if a0==1
scalar z1 = r(mean)
quietly sum z if a0==0
scalar z0 = r(mean)

scalar Y00 = (y010*z0)+(y000*(1-z0))
scalar Y11 = (y111*z1)+(y101*(1-z1))
scalar Y01 = (y011*z1)+(y001*(1-z1))
scalar Y10 = (y110*z0)+(y100*(1-z0))

disp "Y00 = ", Y00
disp "Y01 = ", Y01
disp "Y10 = ", Y10
disp "Y11 = ", Y11

```

Using the data in the table yields:

$$E(Y^{0,0}) = 100.0$$

$$E(Y^{1,0}) = 125.0$$

$$E(Y^{0,1}) = 125.0$$

$$E(Y^{1,1}) = 150.0$$

One can model each probability from the data in this way, potentially with many covariate adjustments.

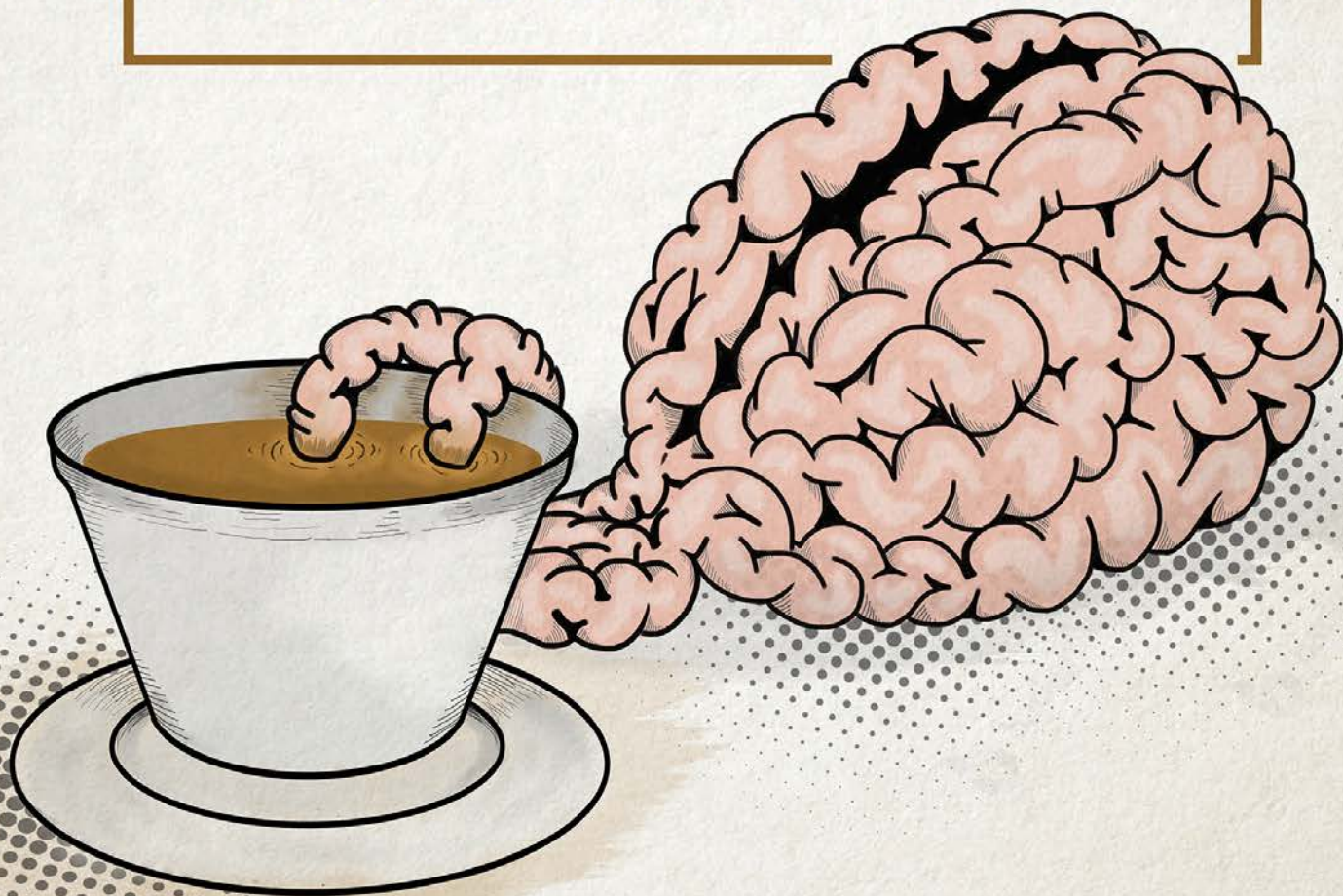
$$\text{Therefore } \varphi = 150.0 - 100.0 = 50.0$$

This is the average causal effect of treatment (always on versus always off ) on CD4 count.

Same as the estimate obtained via IPTW.

The Naimi paper also uses a 3<sup>rd</sup> method (g-estimation).

# COFFEEBREAK



ISEE WORKSHOP- R CODE

# PROPENSITY SCORE MATCHING

p.39-48

#### loading needed packages ####

```
library(data.table)
library(psych)
library(dplyr)
library(MatchIt)
library(Matching)
library(cobalt)
library(Hmisc)
library(WeightIt)
library(Zelig)
```

**Note:** If you do not have these packages installed please install them by using the `install.packages("packagename")` command



```
## importing data, save dta file as a csv or use package  
foreign for STATA version 5-12  
## or package readstata13 for version 13 ##
```

```
data<-read.csv("C:/...../nswre74.csv", header=T, sep=",")  
str(data)
```

```
#### Lable variables ####
```

```
label(data$treat) <- "treatment group"  
label(data$age) <- "age(years)"  
label(data$age2) <- "age squared"  
label(data$ed)<-"years of education"  
label(data$black)<-"black race"  
label(data$hisp)<-"hispanic ethnicity"  
label(data$married)<-"married"  
label(data$nodeg)<-"no high school degree"  
label(data$re74)<-"1974 earnings"  
label(data$re75)<-"1975 earnings"  
label(data$re78)<-"1978 earnings"
```

```
##### means and standard deviations of pre-treatment covariates  
#####
```

```
#Changing the 1 and 2 to 0 and 1 for the factor variables
```

```
data1 <-data.table(data)  
data1$aeg<-as.numeric(data1$age)  
data1$treat<-ifelse(data1$treat=="yes",1,0)  
data1$black<-ifelse(data1$black=="yes",1,0)  
data1$hisp<-ifelse(data1$hisp=="yes",1,0)  
data1$married<-ifelse(data1$married=="yes",1,0)  
data1$nodeg<-ifelse(data1$nodeg=="yes",1,0)
```

```
#### table of means of covariates across treated/control groups
```

```
data2 <- setnames(data1[, sapply(.SD, function(x)  
list(mean=round(mean(x), 3), sd=round(sd(x), 3))),  
by=data1$treat], c("treatment group", sapply(names(data1)[-1],  
paste0, c(".mean", ".SD"))))
```

```
data2
```



```

treatment group age.mean age.SD ed.mean ed.SD black.mean black.SD hisp.mean hisp.SD married.mean
1:           1  25.816  7.155  10.346 2.011      0.843    0.365      0.059   0.237      0.189
2:           0  25.054  7.058  10.088 1.614      0.827    0.379      0.108   0.311      0.154
married.SD nodeg.mean nodeg.SD re74.mean re74.SD re75.mean re75.SD re78.mean re78.SD age2.mean
1:    0.393     0.708    0.456 2095.574 4886.620 1532.055 3219.251 6349.144 7867.402 717.395
2:    0.361     0.835    0.372 2107.027 5687.906 1266.909 3102.982 4554.801 5483.836 677.315
age2.SD aeg.mean aeg.SD
1: 431.252   25.816   7.155
2: 428.784   25.054   7.058
> |

```

```
##### Balance by hand #####
```

```
# age (years)
```

```
a<-subset(data2` treatment group`==1)
```

```
m_var_t<-a$age.mean
```

```
sd_var_t<-a$age.SD
```

```
b<-subset(data2, ` treatment group`==0)
```

```
m_var_c<-b$age.mean
```

```
sd_var_c<-b$age.SD
```

```
d_var<-(100*abs(m_var_t - m_var_c)) / sqrt((sd_var_t^2))
```

```
list<-list("Mean for treated"=m_var_t, "Standard deviation for treated"=sd_var_t,
           "Mean for control"=m_var_c, "Standard deviation for control"=sd_var_c,
           "Standardize difference"=d_var)
```

```
list
```

```
#### repeat this code for each variable
```

#### Age (years) ####

\$`Mean for treated`  
[1] 25.816  
\$`Standard deviation for treated`  
[1] 7.155  
\$`Mean for control`  
[1] 25.054  
\$`Standard deviation for control`  
[1] 7.058  
\$`Standardize difference`  
[1] 10.6499

## Age squared ####

\$`Mean for treated`  
[1] 717.395  
\$`Standard deviation for treated`  
[1] 431.252  
\$`Mean for control`  
[1] 677.315  
\$`Standard deviation for control`  
[1] 428.784  
\$`Standardize difference`  
[1] 9.29387

## Years of education####

\$`Mean for treated`  
[1] 10.346  
\$`Standard deviation for treated`  
[1] 2.011  
\$`Mean for control`  
[1] 10.088  
\$`Standard deviation for control`  
[1] 1.614  
\$`Standardize difference`  
[1] 12.82944

##Black race ####

\$`Mean for treated`  
[1] 0.843  
\$`Standard deviation for treated`  
[1] 0.365  
\$`Mean for control`  
[1] 0.827  
\$`Standard deviation for control`  
[1] 0.379  
\$`Standardize difference`  
[1] 4.383562

# Hispanic ethnicity

\$`Mean for treated`  
[1] 0.059  
\$`Standard deviation for treated`  
[1] 0.237  
\$`Mean for control`  
[1] 0.108  
\$`Standard deviation for control`  
[1] 0.311  
\$`Standardize difference`  
[1] 20.67511

## No high school education

\$`Mean for treated`  
[1] 0.708  
\$`Standard deviation for treated`  
[1] 0.456  
\$`Mean for control`  
[1] 0.835  
\$`Standard deviation for control`  
[1] 0.372  
\$`Standardize difference`  
[1] 27.85088

## Married ##

\$`Mean for treated`  
[1] 0.189  
\$`Standard deviation for treated`  
[1] 0.393  
\$`Mean for control`  
[1] 0.154  
\$`Standard deviation for control`  
[1] 0.361  
\$`Standardize difference`  
[1] 8.905852

## 1975 Earnings##

\$`Mean for treated`  
[1] 1532.055  
\$`Standard deviation for treated`  
[1] 3219.251  
\$`Mean for control`  
[1] 1266.909  
\$`Standard deviation for control`  
[1] 3102.982  
\$`Standardize difference`  
[1] 8.236264

## 1974 Earnings ####

\$`Mean for treated`  
[1] 2095.574  
\$`Standard deviation for treated`  
[1] 4886.62  
\$`Mean for control`  
[1] 2107.027  
\$`Standard deviation for control`  
[1] 5687.906  
\$`Standardize difference`  
[1] 0.2343747

```
##### Build propensity score models for nsw ("by hand") #####
```

```
### Predict treatment status####
```

```
fit1<-glm(treat~age+age2+ed+black+hisp+
          nodeg+married+re74+re75,
          data=data1, family=binomial(link="logit") )
summary(fit1)
```

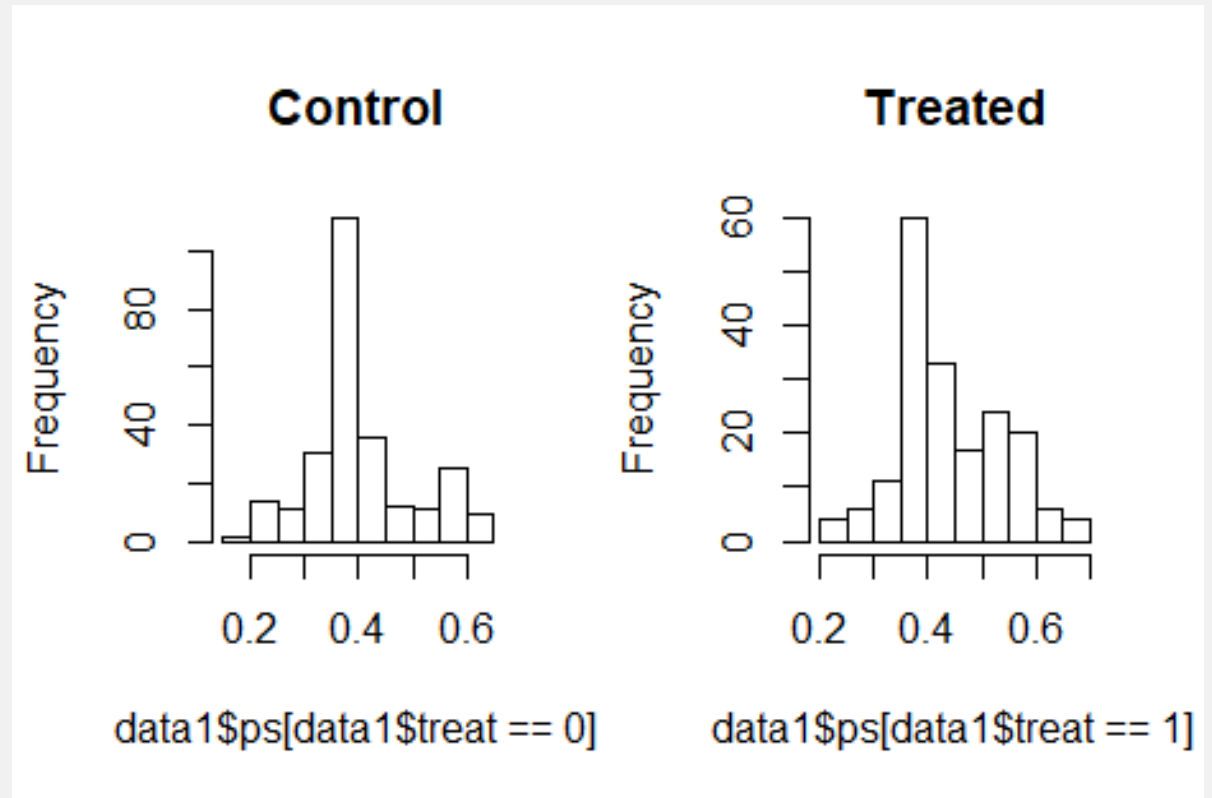
```
## output predicted probability of treatment
data1$ps <-predict(fit1,data1, type = "response")
```

```
### Summary statistics for propensity score, by treatment
```

```
describeBy(data1$ps, group =data1$treat)
par(mfrow = c(1, 2))
hist(data1$ps[data1$treat==0])
hist(data1$ps[data1$treat==1])
```

```
## nearest neighbor matching (1:1) without replacement ##
```

```
##randomly order data in case match ties
set.seed(123456)
data1$ranorder<-runif(445)
data3<-data1[order(ranorder),]
```



```
## propensity score model
fit2<-
glm(treat~age+age2+ed+black+hisp+
nodeg+married+re74+re75,
      data=data3,
family=binomial(link="logit") )
```

```
#### create propensity scores ##
match1<-matchit(fit2, data3, method =
"nearest",replace=F)
```

```
## Evaluate balance
summary(match1)
love.plot(bal.tab(match1), threshold = .1)
MatchBalance(fit2, data3, nboots=500)
```

```
Call:
matchit(formula = fit2, data = data3, method = "nearest", replace = F)
```

Summary of balance for all data:

	Means Treated	Means Control	SD Control	Mean Diff	eQQ Med	eQQ Mean	eQQ Max
distance	0.4377	0.4001	0.0934	0.0376	0.029	0.0383	0.106
age	25.8162	25.0538	7.0577	0.7624	1.000	0.9405	7.000
age2	717.3946	677.3154	428.7844	40.0792	43.000	56.0757	721.000
ed	10.3459	10.0885	1.6143	0.2575	0.000	0.4054	2.000
black	0.8432	0.8269	0.3790	0.0163	0.000	0.0162	1.000
hisp	0.0595	0.1077	0.3106	-0.0482	0.000	0.0486	1.000
nodeg	0.7081	0.8346	0.3722	-0.1265	0.000	0.1243	1.000
married	0.1892	0.1538	0.3615	0.0353	0.000	0.0378	1.000
re74	2095.5737	2107.0266	5687.9056	-11.4530	0.000	487.9811	8412.999
re75	1532.0553	1266.9090	3102.9821	265.1463	0.000	367.6137	2110.260

Summary of balance for matched data:

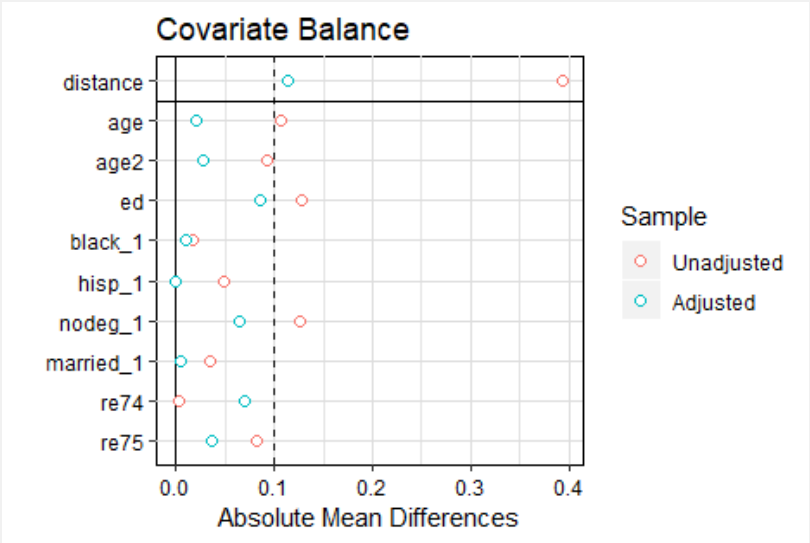
	Means Treated	Means Control	SD Control	Mean Diff	eQQ Med	eQQ Mean	eQQ Max
distance	0.4377	0.4268	0.0911	0.0109	0.0052	0.0120	0.0591
age	25.8162	25.6703	6.8099	0.1459	0.0000	0.4703	7.0000
age2	717.3946	705.0865	407.8264	12.3081	0.0000	32.5459	721.0000
ed	10.3459	10.1730	1.7139	0.1730	0.0000	0.2919	2.0000
black	0.8432	0.8541	0.3540	-0.0108	0.0000	0.0108	1.0000
hisp	0.0595	0.0595	0.2371	0.0000	0.0000	0.0000	0.0000
nodeg	0.7081	0.7730	0.4200	-0.0649	0.0000	0.0649	1.0000
married	0.1892	0.1946	0.3970	-0.0054	0.0000	0.0054	1.0000
re74	2095.5737	1756.0642	4597.6400	339.5095	0.0000	486.6042	9319.1500
re75	1532.0553	1415.2466	3348.4559	116.8087	0.0000	323.9474	2510.5200

Percent Balance Improvement:

	Mean Diff.	eQQ Med	eQQ Mean	eQQ Max
distance	71.0402	82.2363	68.5551	44.2504
age	80.8563	100.0000	50.0000	0.0000
age2	69.2905	100.0000	41.9607	0.0000
ed	32.8220	0.0000	28.0000	0.0000
black	33.7580	0.0000	33.3333	0.0000
hisp	100.0000	0.0000	100.0000	100.0000
nodeg	48.7264	0.0000	47.8261	0.0000
married	84.7059	0.0000	85.7143	0.0000
re74	-2864.3820	0.0000	0.2822	-10.7708
re75	55.9456	0.0000	11.8783	-18.9673

Sample sizes:

	Control	Treated
All	260	185
Matched	185	185
Unmatched	75	0
Discarded	0	0



```
## modified nearest neighbor matching (1:1) without replacement (optional)##
```

```
##generate squared term for 1974 earnings
```

```
data3$re74_2<-data3$re74*data3$re74
```

```
####randomly order data in case match ties
```

```
data3$ranorder<-runif(445)
```

```
data4<-data3[order(ranorder),]
```

```
####create propensity score
```

```
fit3<-glm(treat~age+age2+ed+black+hisp+ nodeg +married+re74+re75+re74_2,  
          data=data4, family=binomial(link="logit") )
```

```
match2<-matchit(fit3, data4, method = "nearest", distance = "logit",replace=F)
```

```
####evaluate balance
```

```
summary(match2)
```

```
love.plot(bal.tab(match2), threshold = .1)
```

```
MatchBalance(fit3, data4, nboots=500)
```

```
##generating dataset with only matched data
```

```
m.data2 <- match.data(match2,distance ="pscore")
```

```
##### treatment effect in final matched sample####
```

```
fit5<-lm(re78~treat, data=m.data2)
```

```
summary(fit5)
```

```
fit6<-lm(re78~treat+age+age2+ed+black+hisp+nodeg  
+married+re75+re74,data=m.data2)
```

```
summary(fit6)
```

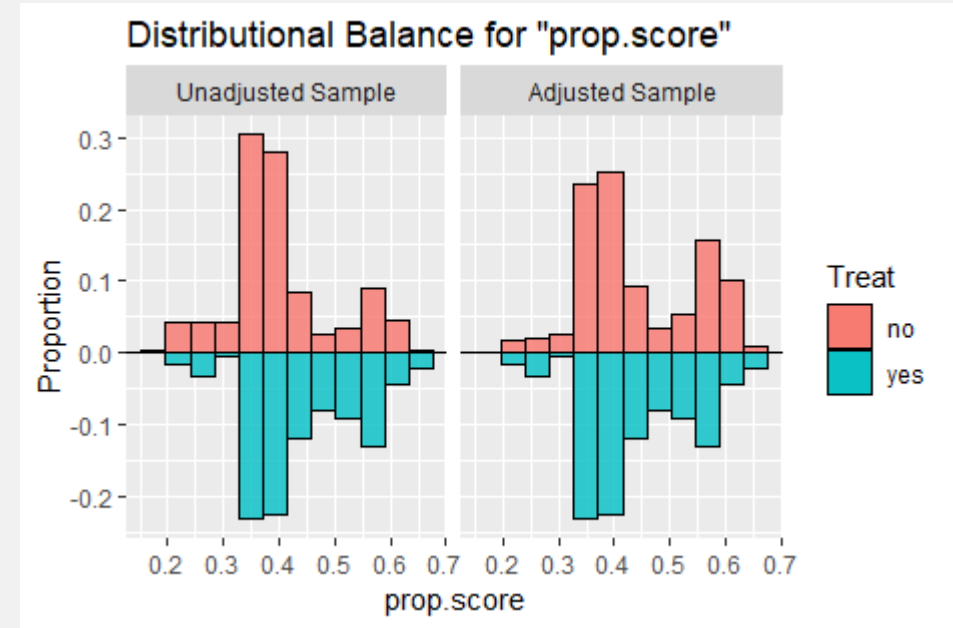
```
## examine support
```

```
covs0 <- subset(data, select = -c(treat, re78))
```

```
W.out <- weightit(treat ~ covs0, data = data,  
                 method = "ps", estimand = "ATT")
```

```
bal.tab(W.out)
```

```
bal.plot(W.out, var.name = "prop.score", which = "both",  
         type = "histogram", mirror = TRUE)
```



#### bootstrapped standard error for ATT##

```
bootstrap.se(match2,Y=data4$re78, max.iter=500)
```

#### To get ATT and ATE you can use a combination of the zelig package and the matchit package

See example: <https://r.iq.harvard.edu/docs/matchit/2.4-15/Examples2.html>



# PROPENSITY SCORE WEIGHTING

p.81-97

# ##### Smoking IPTW Example #####

```
library(estimatr)
library(questionr)
library(glmML)
```

```
aa <- read.csv("C:/.../nhefs.csv",sep="," , header=T)
```

```
#### setting up data ####
```

```
aa$exer<-as.factor(aa$exercise)
```

```
aa$educ<-as.factor(aa$education)
```

```
aa$act<-as.factor(aa$active)
```

```
aa$age2<-aa$age*aa$age
```

```
aa$wt7l_2<-aa$wt7l*aa$wt7l
```

```
aa$smokyr2<-aa$smokeyr*aa$smokeyr
```

```
aa$smokeint2<-aa$smokeintensity*aa$smokeintensity
```

```
#### generating propensity score ####
```

```
propensity <- glm(qsmk ~ sex + age+age2 + race + relevel(educ, ref = 5)+wt7l+wt7l_2+smokeintensity+smokeint2+
  relevel(exer, ref = 2) +relevel(act, ref = 2)+smokeyr+smokyr2, data = aa, family = binomial("logit"))
```

```
summary(propensity)
```

```
aa$p_qsmk<-predict(propensity, data=aa, type="response")
```

```
## Creating weights ####
```

```
aa$w<-ifelse(aa$qsmk==1, 1/aa$p_qsmk, 1/(1-aa$p_qsmk))
```

```
summary(aa$w)
```

```
hist(aa$w)
```

#### #### Regressions with weights ###

```
fit1<-lm_robust(wt82_71~qsmk, data=aa, weight=w, cluster=seqn)  
summary(fit1)
```

#### #### stabilized weight ####

##### ## Estimate the denominator

```
prop2<- glm(qsmk ~ sex + race+age+age2 + relevel(educ, ref = 5)+wt71+wt71_2+smokeintensity+smokeint2+  
            relevel(exer, ref = 2) +relevel(act, ref = 2)+smokeyrs+smokyrs2, data = aa, family = binomial("logit"))
```

```
aa$pd_qsmk<-predict(prop2,data=aa, type="response" )
```

#### #### Estimation of the numerator of ip weights ##

```
fit2<-glm(qsmk~1, data=aa, family = binomial("logit"))  
summary(fit2)  
aa$pn_qsmk<-predict(fit2, data=aa, type="response")
```

#### #### Creating stabilized weights ####

```
aa$sw_a<-ifelse(aa$qsmk==1, aa$pn_qsmk/aa$pd_qsmk, (1-aa$pn_qsmk)/(1-aa$pd_qsmk))  
summary(aa$sw_a)
```

#### #### Running regression again with sw weights##

```
fit3<-lm_robust(wt82_71~qsmk, data=aa, weight=sw_a, cluster=seqn)  
summary(fit3)
```

G-METHODS  
CODE FOR NAIMI'S PAPER

p.98-116

```
# Arrange data into wide format
```

```
a0<-c(0,0,0,0,1,1,1,1)
z<-c(0,0,1,1,0,0,1,1)
a1<-c(0,1,0,1,0,1,0,1)
y<-c(87.29,112.11,119.65,144.84,105.28,130.18,137.72,162.83)
N<-c(209271,93779,60654,136293,134781,60789,93903,210527)
D<-NULL
for(i in 1:8){
  d<-
  data.frame(cbind(rep(a0[i],N[i]),rep(z[i],N[i]),rep(a1[i],N[i]),rep(y[i],N[i])))
  D<-rbind(D,d)
}
names(D)<-c("a0","z","a1","y")
dim(D)
```

```
##### Standard models #####
```

```
D$avga<-(D$a0+D$a1)/2
fit1<-glm(y~avga, data=D)
summary(fit1)
fit2<-glm(y~avga+z, data=D)
summary(fit2)
fit3<-glm(y~a0, data=D)
summary(fit3)
fit4<-glm(y~a0+z, data=D)
summary(fit4)
fit5<-glm(y~a1, data=D)
summary(fit5)
fit6<-glm(y~a1+z, data=D)
summary(fit6)
```

##### IPTW #####

```
pal<-mean(D$a1)
pal
k<-subset(D, z==0)
palz0<-mean(k$a1)
palz0
l<-subset(D, z==1)
palzl<-mean(l$a1)
palzl
```

```
w000 <- (0.5 * (1-palz0))
w00l <- (0.5 * (palz0))
w0l0 <- (0.5 * (1-palzl))
wl00 <- (0.5 * (1-palz0))
wl0l <- (0.5 * (palz0))
wl10 <- (0.5 * (1-palzl))
wl1l <- (0.5 * (palzl))
```

```
sw000 <- (0.5*(1-pal))/(0.5 * (1-palz0))
sw00l <- (0.5*pal)/(0.5 * (palz0))
sw0l0 <- (0.5*(1-pal))/(0.5 * (1-palzl))
sw0ll <- (0.5*pal)/(0.5 * (palzl))
swl00 <- (0.5*(1-pal))/(0.5 * (1-palz0))
swl0l <- (0.5*pal)/(0.5 * (palz0))
swll0 <- (0.5*(1-pal))/(0.5 * (1-palzl))
swlll <- (0.5*pal)/(0.5 * (palzl))
```

```
D$sw<-ifelse(D$a0==0 & D$z==0 & D$a1==0,sw000,
             ifelse(D$a0==0 & D$z==0 & D$a1==1,sw00l,
             ifelse(D$a0==0 & D$z==1 & D$a1==0,sw0l0,
             ifelse(D$a0==0 & D$z==1 & D$a1==1,sw0ll,
             ifelse(D$a0==1 & D$z==0 & D$a1==0,sw000,
             ifelse(D$a0==1 & D$z==0 & D$a1==1,sw00l,
             ifelse(D$a0==1 & D$z==1 & D$a1==0,sw0l0,sw0ll))))))
```

```
fit7<-glm(y~a0*a1, data=D, weights = sw)
summary(fit7)
```

##### g-formula #####

```
a<-subset(D, a1==0 & z==0 & a0==0)
y000<-mean(a$y)
b<-subset(D, a1==0 & z==0 & a0==1)
y001<-mean(b$y)
c<-subset(D, a1==0 & z==1 & a0==0)
y010<-mean(c$y)
d<-subset(D, a1==1 & z==0 & a0==0)
y100<-mean(d$y)
e<-subset(D, a1==0 & z==1 & a0==1)
y011<-mean(e$y)
f<-subset(D, a1==1 & z==1 & a0==0)
y110<-mean(f$y)
g<-subset(D, a1==1 & z==0 & a0==1)
y101<-mean(g$y)
h<-subset(D, a1==1 & z==1 & a0==1)
y111<-mean(h$y)
i<-subset(D, a0==1)
z1<-mean(i$z)
j<-subset(D, a0==0)
z0<-mean(j$z)
```

$$Y_{00} <- (y_{010} * z_0) + (y_{000} * (1 - z_0))$$
$$Y_{11} <- (y_{111} * z_1) + (y_{101} * (1 - z_1))$$
$$Y_{01} <- (y_{011} * z_1) + (y_{001} * (1 - z_1))$$
$$Y_{10} <- (y_{110} * z_0) + (y_{100} * (1 - z_0))$$

Y00  
Y01  
Y10  
Y11