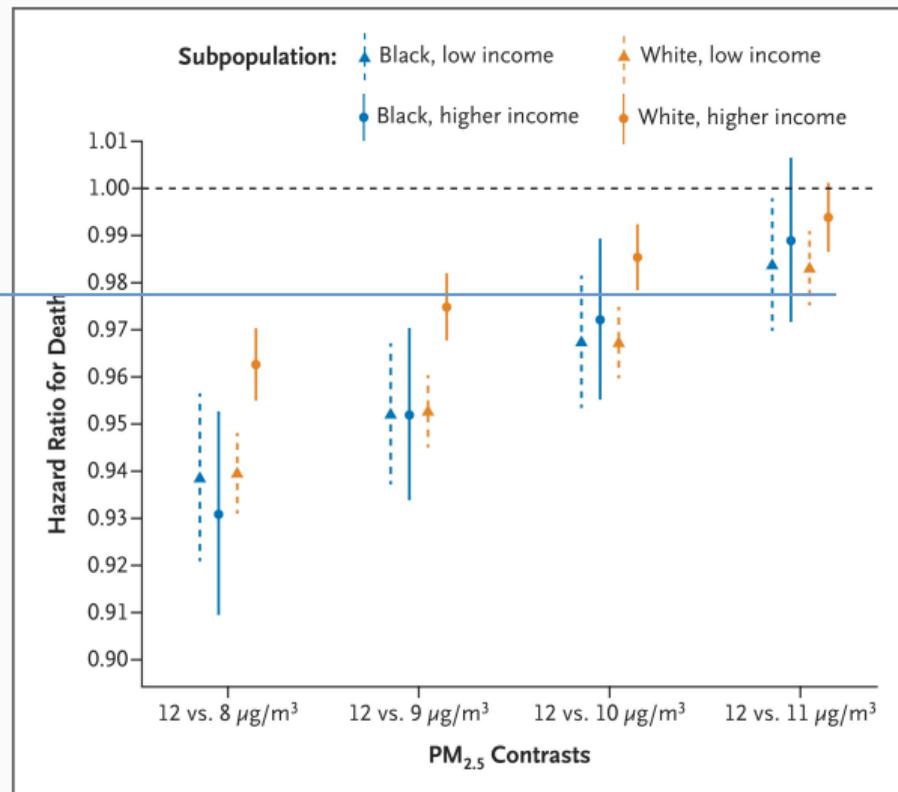


Mathematical Statistics II

STA2212H S LEC9101

Week 10

March 18 2025



SPECIAL ARTICLE

Air Pollution and Mortality at the Intersection of Race and Social Class

Kevin P. Josey, Ph.D., Scott W. Delaney, Sc.D., J.D., Xiao Wu, Ph.D.,
Rachel C. Nethery, Ph.D., Priyanka DeSouza, Ph.D., Danielle Braun, Ph.D.,
and Francesca Dominici, Ph.D.

ABSTRACT

BACKGROUND

From the Departments of Biostatistics (K.P.J., R.C.N., D.B., F.D.) and Environmental Health (S.W.D.), Harvard T.H. Chan School of Public Health, Boston; the Department of Biostatistics, Mailman School of Public Health, Columbia University, New York (X.W.); and the Department of Urban and Regional Planning, University of Colorado Denver, Denver (P.D.). Dr. Dominici can be reached at fdominic@hsph.harvard.edu or at the Department of Biostatistics, Harvard T.H. Chan School of Public Health, 655 Huntington Ave., Bldg. 2, 4th Fl., Boston, MA 02115.

Drs. Josey and Delaney and Drs. Braun and Dominici contributed equally to this article.

This article was published on March 24, 2023, at [NEJM.org](https://www.nejm.org).

N Engl J Med 2023;388:1396–404.
DOI: 10.1056/NEJMsa2300523
Copyright © 2023 Massachusetts Medical Society.

Black Americans are exposed to higher annual levels of air pollution containing fine particulate matter (particles with an aerodynamic diameter of $\leq 2.5 \mu\text{m}$ [$\text{PM}_{2.5}$]) than White Americans and may be more susceptible to its health effects. Low-income Americans may also be more susceptible to $\text{PM}_{2.5}$ pollution than high-income Americans. Because information is lacking on exposure–response curves for $\text{PM}_{2.5}$ exposure and mortality among marginalized subpopulations categorized according to both race and socioeconomic position, the Environmental Protection Agency lacks important evidence to inform its regulatory rulemaking for $\text{PM}_{2.5}$ standards.

METHODS

We analyzed 623 million person-years of Medicare data from 73 million persons 65 years of age or older from 2000 through 2016 to estimate associations between annual $\text{PM}_{2.5}$ exposure and mortality in subpopulations defined simultaneously by racial identity (Black vs. White) and income level (Medicaid eligible vs. ineligible).

RESULTS

Lower $\text{PM}_{2.5}$ exposure was associated with lower mortality in the full population, but marginalized subpopulations appeared to benefit more as $\text{PM}_{2.5}$ levels decreased. For example, the hazard ratio associated with decreasing $\text{PM}_{2.5}$ from $12 \mu\text{g}$ per cubic meter to $8 \mu\text{g}$ per cubic meter for the White higher-income subpopulation was 0.963 (95% confidence interval [CI], 0.955 to 0.970), whereas equivalent hazard

Today

1. Recap Mar 11 [goodness-of-fit tests](#)
2. Introduction to causal inference
3. Project guidelines:

if you are not sure how to fit your paper into the guidelines contact me

office hour: Tuesday 3-4; Monday 7-8

email: nancym.reid@utoronto.ca

Recap: multinomial goodness of fit statistics

- Pearson's χ^2 test

$$Q = \sum_{j=1}^k \frac{\{Y_j - np_j(\hat{\theta})\}^2}{np_j(\hat{\theta})} \xrightarrow{d} \chi_{k-1-p}^2$$

$\tilde{\theta}$ MLE in multinomial

- Likelihood ratio (deviance) test

$$W = 2 \sum_{j=1}^k Y_j \log \left(\frac{Y_j}{np_j(\tilde{\theta})} \right) \xrightarrow{d} \chi_{k-1-p}^2$$

Recap: Smooth goodness-of-fit statistics

$$K_n = \sup_t |\widehat{F}_n(t) - F_0(t)| \xrightarrow{d} K, \quad \text{pr}(K > x) = 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp(-2j^2 x^2)$$

$$W_n^2 = \int \{\widehat{F}_n(t) - F_0(t)\}^2 dF_0(t) \xrightarrow{d} \sum_{j=1}^{\infty} \frac{Z_j^2}{j^2 \pi^2}$$

$$A_n^2 = \int \frac{\{\widehat{F}_n(t) - F_0(t)\}^2}{F_0(t)\{1 - F_0(t)\}} dF_0(t) \xrightarrow{d} \sum_{j=1}^{\infty} \frac{Z_j^2}{j(j+1)}$$

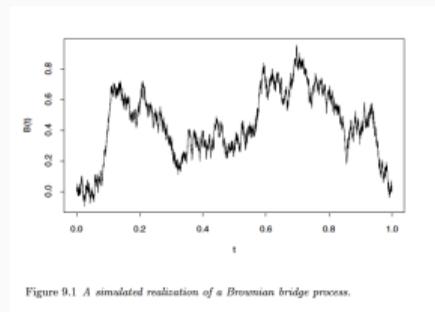


Figure 9.1 A simulated realization of a Brownian bridge process.

Topic Introduction

Goodness-of-Fit tests via Machine Learning

GOF tests:

- Assess how well a model H_0 describes the data
- No alternative model H_1 specified, if it was, likelihood ratio $L(\text{data} | H_1) / L(\text{data} | H_0)$ would provide optimal test (Neyman-Pearson)
- Standard GOF tests in HEP: χ^2 (most frequent), Kolgomorov- Smirnov (seldomly), others..
 - Difficulties arise for multi-dimensional distributions

➔ Machine Learning offers various possibilities 😊

Today's topic!

[link](#)



J. R. Statist. Soc. B (2020)
82, Part 3, pp. 773–795

Goodness-of-fit testing in high dimensional generalized linear models

Jana Janková and Rajen D. Shah,
University of Cambridge, UK

Peter Bühlmann
Eidgenössische Technische Hochschule Zürich, Switzerland

and Richard J. Samworth
University of Cambridge, UK

[Received August 2019, Revised March 2020]

Summary: We propose a family of tests to assess the goodness of fit of a high dimensional

2. Methodology: generalized residual prediction tests

As mentioned in Section 1.1, our generalized residual prediction (GRP) testing methodology relies on an initial fit of the form

$$\hat{\beta} := \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{n} \sum_{i=1}^n \rho(Y_i, x_i^T \beta) + \lambda \|\beta\|_1 \right\}.$$

In what follows, we refer to $\hat{\beta}$ as the GLM lasso, though it is not essential that the loss function $\rho: \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}$ is the negative log-likelihood that is obtained from a GLM, and indeed this definition incorporates penalized quasi-likelihood estimators, among others. Our general framework for goodness-of-fit testing will also assume that we have available an auxiliary data set $(X_A, Y_A) \in \mathbb{R}^{n_A \times p} \times \mathcal{Y}^{n_A}$ independent of (X, Y) . In the rest of the paper, we take $n_A = n$ for simplicity, although this is not needed for our procedures. Consider the Pearson-type residuals

$$R_i = \frac{Y_i - \mu(x_i^T \tilde{\beta})}{\sqrt{V\{\mu(x_i^T \tilde{\beta})\}}}, \quad i = 1, \dots, n.$$

Here $\tilde{\beta} \in \mathbb{R}^p$ is an additional estimate of β_0 that may be computed by using the auxiliary data set, or in certain circumstances may be taken as $\hat{\beta}$ itself: we discuss these two cases in the following sections. Given the vector R of residuals, the basic form of our test statistic is $w^T R$; here $w \in \mathbb{R}^n$ is a direction that is typically derived by using the auxiliary data set. We describe in detail the construction of such a w in Section 2.1, where the goal is general goodness-of-fit testing.

A further modification of the method can enable us to use multiple directions w to test simultaneously for different departures from the null or to aggregate over different directions derived by using flexible regression methods with different tuning parameters. Given a set $W \subseteq \mathbb{R}^n$ of direction vectors w , our proposed test statistic then takes the form

$$\sup_{w \in W} w^T R.$$

The Annals of Statistics

2022, Vol. 50, No. 5, 2514–2544

<https://doi.org/10.1214/22-AOS2187>

© Institute of Mathematical Statistics, 2022

TESTING GOODNESS-OF-FIT AND CONDITIONAL INDEPENDENCE WITH APPROXIMATE CO-SUFFICIENT SAMPLING

BY RINA FOYGEL BARBER^{1,a} AND LUCAS JANSON^{2,b}

¹*Department of Statistics, University of Chicago, ^arina@uchicago.edu*

²*Department of Statistics, Harvard University, ^bljanson@fas.harvard.edu*

Goodness-of-fit (GoF) testing is ubiquitous in statistics, with direct ties to model selection, confidence interval construction, conditional independence testing, and multiple testing, just to name a few applications. While testing the GoF of a simple (point) null hypothesis provides an analyst great flexibility in the choice of test statistic while still ensuring validity, most GoF tests for composite null hypotheses are far more constrained, as the test statistic must have a tractable distribution over the entire null model space. A notable exception is *co-sufficient sampling* (CSS): resampling the data conditional on a sufficient statistic for the null model guarantees valid GoF testing using any test statistic the analyst chooses. But CSS testing requires the null model to

- randomization; confounding; observational studies; experiments; “correlation is not causation”, Simpson’s ‘paradox’
- counterfactuals; average treatment effect; conditional average treatment effect; ...
- graphical models; directed acyclic graphs; causal graphs; Markov assumptions...

- The Book

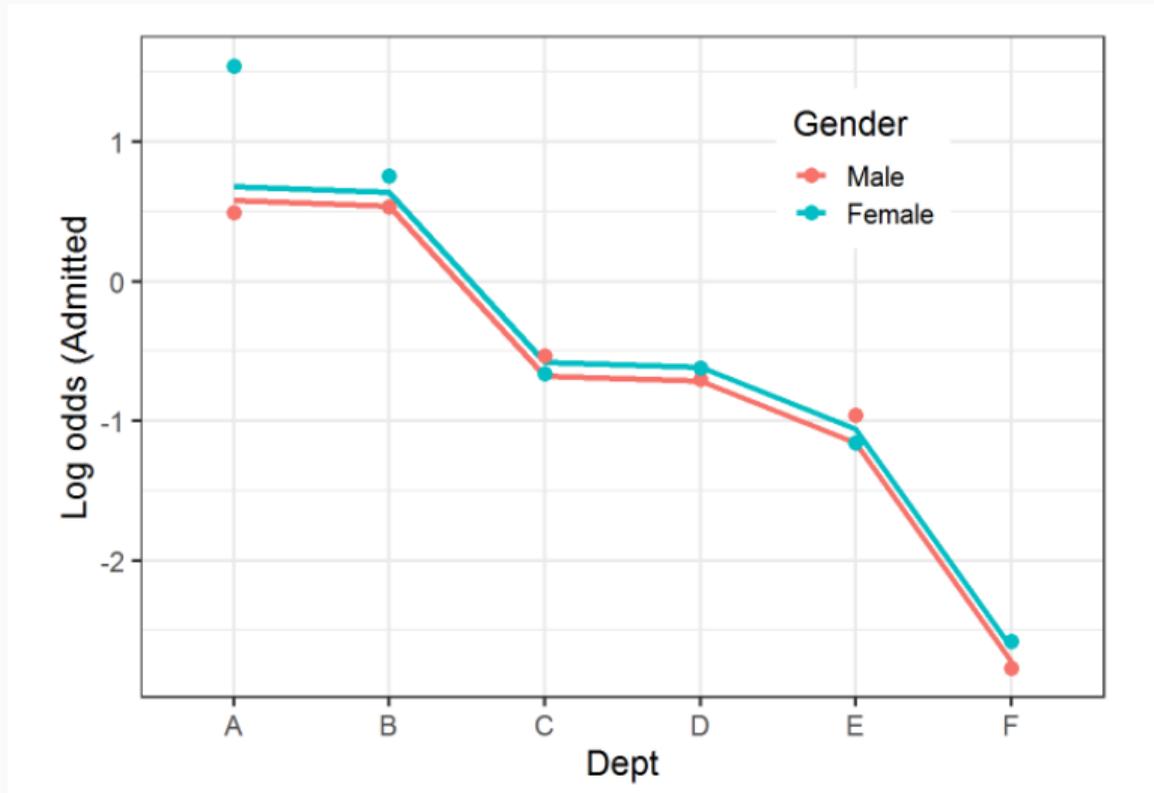


Confounding variables

Major	Number of applicants	Men		Women		
		Number admitted	Percent admitted	Number of applicants	Number admitted	Percent admitted
A	825	512	62	108	89	82
B	560	353	63	25	17	68
C	325	120	37	593	202	34
D	417	138	33	375	131	35
E	191	53	28	393	94	24
F	373	22	6	341	24	7
Total	2691	1198	44	1835	557	30

`data(UCBAdmissions)`

... Confounding variables



race of defendant	death penalty imposed	death penalty not imposed	percentage
white	19	141	11.88%
black	17	149	10.24%

race of defendant	death penalty imposed	death penalty not imposed	percentage
white	19	141	11.88%
black	17	149	10.24%

white victim	race of defendant	death penalty imposed	death penalty not imposed	percentage
	white	19	132	12.58%
	black	11	52	17.46%

black victim	race of defendant	death penalty imposed	death penalty not imposed	percentage
	white	0	9	0%
	black	6	97	5.83%

258

6 · Stochastic Models

Age (years)	Smokers	Non-smokers
Overall	139/582 (24)	230/732 (31)
18–24	2/55 (4)	1/62 (2)
25–34	3/124 (2)	5/157 (3)
35–44	14/109 (13)	7/121 (6)
45–54	27/130 (21)	12/78 (15)
55–64	51/115 (44)	40/121 (33)
65–74	29/36 (81)	101/129 (78)
75+	13/13 (100)	64/64 (100)

Table 6.8 Twenty-year survival and smoking status for 1314 women (Appleton *et al.*, 1996). The smoker and non-smoker columns contain number dead/total (% dead).

- A – binary treatment indicator
- Y – binary outcome
- “ A **causes** Y ” to be distinguished from “ A is associated with Y ”

AoS uses X for tmt
could be continuous

- A – binary treatment indicator AoS uses X for tmt
- Y – binary outcome could be continuous
- “A **causes** Y” to be distinguished from “A is associated with Y”

- introduce **potential outcomes** $Y(0), Y(1)$ AoS C_0, C_1 ; HR Y^a
- **causal treatment effect** $\theta = E(Y(1)) - E(Y(0))$ want to estimate this
- **association** $\alpha = E(Y | A = 1) - E(Y | A = 0)$ have data to estimate α

- **Consistency assumption:** $Y = Y(a)$ we can learn about potential outcome from observed values

Potential outcomes C_0, C_1

X	Y	C_0	C_1
0	4	4	*
0	7	7	*
0	2	2	*
0	8	8	*
1	3	*	3
1	5	*	5
1	8	*	8
1	9	*	9

treatment X , response Y Potential outcomes Y^0, Y^1

Table 2.1

	A	Y	Y^0	Y^1
Rheaia	0	0	0	?
Kronos	0	1	1	?
Demeter	0	0	0	?
Hades	0	0	0	?
Hestia	1	0	?	0
Poseidon	1	0	?	0
Hera	1	0	?	0
Zeus	1	1	?	1
Artemis	0	1	1	?
Apollo	0	1	1	?
Leto	0	0	0	?
Ares	1	1	?	1
Athena	1	1	?	1
Hephaestus	1	1	?	1
Aphrodite	1	1	?	1
Cyclope	1	1	?	1
Persephone	1	1	?	1
Hermes	1	0	?	0
Hebe	1	0	?	0
Dionysus	1	0	?	0

Potential outcomes

Table 1.1

	$Y^{a=0}$	$Y^{a=1}$
Rheia	0	1
Kronos	1	0
Demeter	0	0
Hades	0	0
Hestia	0	0
Poseidon	1	0
Hera	0	0
Zeus	0	1
Artemis	1	1
Apollo	1	0
Leto	0	1
Ares	1	1
Athena	1	1
Hephaestus	0	1
Aphrodite	0	1
Cyclope	0	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

Observed outcomes

Table 1.2

	A	Y
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

$$\theta = E(Y(1)) - E(Y(0))$$

risk difference; ratio; odds

also called “ATE” and “ACE”: average treatment/causal effect

$$\alpha = E(Y | A = 1) - E(Y | A = 0) \quad \text{this can be estimated from the data}$$

If A is independent of $(Y(0), Y(1))$, then $\theta = \alpha$

If treatment is **randomly assigned**, then $A \perp (Y(0), Y(1))$

$\perp \equiv$ independent

Example 16.2

X	Y	C_0	C_1
0	0	0	0^*
0	0	0	0^*
0	0	0	0^*
0	0	0	0^*
1	1	1^*	1
1	1	1^*	1
1	1	1^*	1
1	1	1^*	1

$$\theta = \mathbf{0}; \quad \alpha = 1$$

(C_0, C_1) not independent of X

X	Y	C_0	C_1
0	0	0	0^*
1	0	0	0^*
1	0	0	0^*
1	0	0	0^*
1	1	1^*	1
1	1	1^*	1
1	1	1^*	1
1	1	1^*	1

$$\theta = \mathbf{0}, \quad \alpha = 4/7 < 1$$

thought experiment

Potential outcomes

Table 1.1

	$Y^{a=0}$	$Y^{a=1}$
Rheia	0	1
Kronos	1	0
Demeter	0	0
Hades	0	0
Hestia	0	0
Poseidon	1	0
Hera	0	0
Zeus	0	1
Artemis	1	1
Apollo	1	0
Leto	0	1
Ares	1	1
Athena	1	1
Hephaestus	0	1
Aphrodite	0	1
Cyclope	0	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

Observed outcomes

Table 1.2

	A	Y
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

1. A well-understood evidence-based mechanism, or set of mechanisms, that links a cause to its effect
2. two phenomena are linked by a stable association, whose direction is established and which cannot be explained by mutual dependence on some other allowable variable
3. observed association may be linked to causal effect via counterfactuals if
 $(Y(0), Y(1)) \perp A$ not usually testable

- typically have additional explanatory variables (covariates) X

AoS uses Z ; HR use L

- causal effect of treatment when $X = x$

$$\theta(x) = E(Y(1) | X = x) - E(Y(0) | X = x)$$

- marginal causal effect

$$\theta = E_X\{E(Y(1) | X) - E(Y(0) | X)\}$$

- association function

$$r(x) = E(Y | A = 1, X = x) - E(Y | A = 0, X = x)$$

- marginal association

$$E_X\{r(X)\}$$

Table 2.2

	L	A	Y
Rheia	0	0	0
Kronos	0	0	1
Demeter	0	0	0
Hades	0	0	0
Hestia	0	1	0
Poseidon	0	1	0
Hera	0	1	0
Zeus	0	1	1
Artemis	1	0	1
Apollo	1	0	1
Leto	1	0	0
Ares	1	1	1
Athena	1	1	1
Hephaestus	1	1	1
Aphrodite	1	1	1
Cyclope	1	1	1
Persephone	1	1	1
Hermes	1	1	0
Hebe	1	1	0
Dionysus	1	1	0

$$\theta_{L=0}$$

$$\theta_{L=1}$$

$L = 1$ critical condition

$L = 0$ stable condition
conditional randomization

- in observational studies treatment is not randomly assigned $\implies \theta(x) \neq r(x)$
- **No unmeasured confounding:**

$$\{Y(a); a \in \mathcal{A}\} \perp A \mid X$$

can learn about $Y(a)$ even if $A \neq a$ by using observed Y for 'similar' people from $A = a$ group

- under the assumption of no unmeasured confounding, marginal causal effect

$$E(Y(a)) = \int E(Y \mid A = a, X = x) dF_X(x)$$

can be estimated by the association function

$$\hat{E}(Y(a)) = \frac{1}{n} \sum_{i=1}^n \hat{r}(a, X_i) = \hat{\beta}_0 + \hat{\beta}_1 a + \hat{\beta}_2 \bar{X}_n$$

causal reg function \equiv adjusted treatment effect

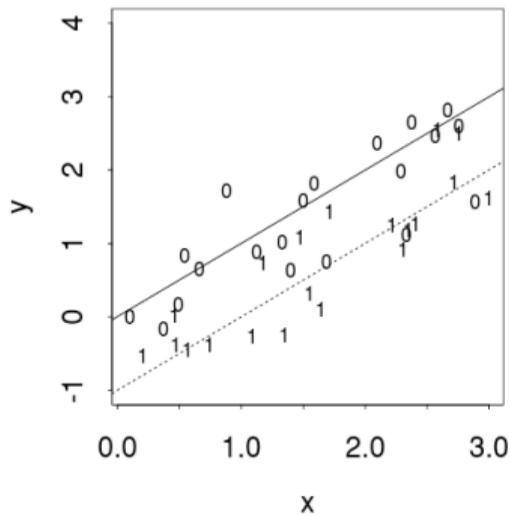
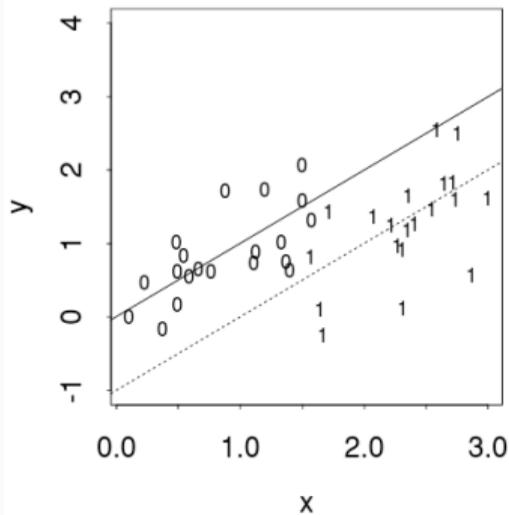


Figure 9.2 Simulated results from experiments to compare the effect of a treatment T on a response Y that varies with a covariate X . The lines show the mean response for $T = 0$ (solid) and $T = 1$ (dotted). Left: the effect of T is confounded with dependence on X . Right: the experiment is balanced, with random allocation of T dependent on X .

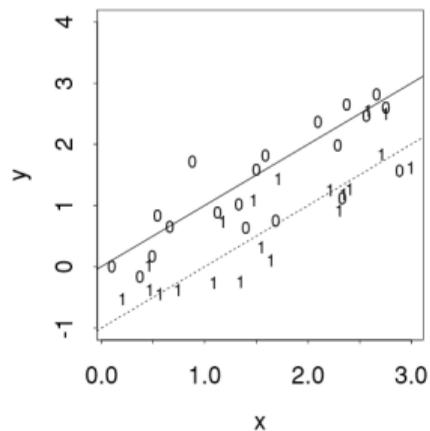
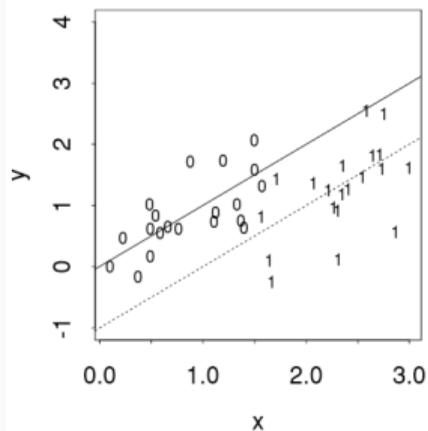


Figure 9.2 Simulated results from experiments to compare the effect of a treatment T on a response Y that varies with a covariate X . The lines show the mean response for $T = 0$ (solid) and $T = 1$ (dots). Left: the effect of T is confounded with dependence on X . Right: the experiment is balanced, with random allocation of T dependent on X .

Causal effect $\equiv 1$

Left: $\bar{y}_1 - \bar{y}_0 = 0.2 \pm 0.3$

Right: $\bar{y}_1 - \bar{y}_0 = -1.2 \pm 0.3$

adjust for covariate: $y = \beta_0 + \beta_1 x + \delta t + \epsilon$

Left: $\hat{\delta} = -0.7 \pm 0.3$ Right: $\hat{\delta} = -1.25 \pm 0.16$

right randomized within pairs; matched on x

“Bradford-Hill guidelines” Evidence that an observed association is causal is strengthened if:

- the association is strong
- the association is found consistently over a number of independent studies
- the association is specific to the outcome studied
- the observation of a potential cause occurs earlier in time than the outcome
- there is a dose-response relationship
- there is subject-matter theory that makes a causal effect plausible
- the association is based on a suitable natural experiment

see also AoS §16.3

260 16. Causal Inference

	Y = 1	Y = 0	Y = 1	Y = 0
X = 1	.1500	.2250	.1000	.0250
X = 0	.0375	.0875	.2625	.1125
	Z = 1 (men)		Z = 0 (women)	

The marginal distribution for (X, Y) is

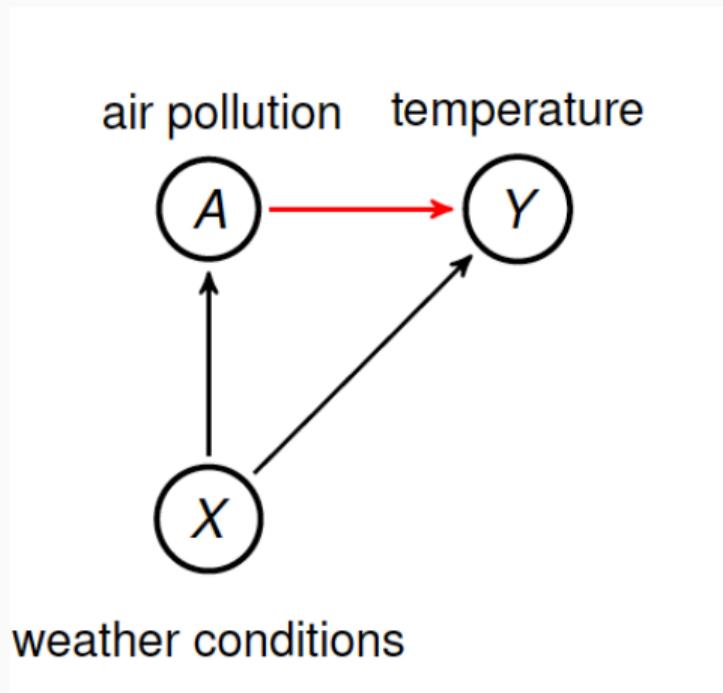
	Y = 1	Y = 0	
X = 1	.25	.25	.50
X = 0	.30	.20	.50
	.55	.45	1

From these tables we find that,

$$\begin{aligned} \mathbb{P}(Y = 1|X = 1) - \mathbb{P}(Y = 1|X = 0) &= -0.1 \\ \mathbb{P}(Y = 1|X = 1, Z = 1) - \mathbb{P}(Y = 1|X = 0, Z = 1) &= 0.1 \\ \mathbb{P}(Y = 1|X = 1, Z = 0) - \mathbb{P}(Y = 1|X = 0, Z = 0) &= 0.1. \end{aligned}$$

To summarize, we *seem* to have the following information:

confusion of causal effect
with association



- assume no unmeasured confounding

- want to estimate

$$E(Y(1) | X) - E(Y(0) | X)$$

causal regression function

- or possibly $E_X\{E(Y(1) | X) - E(Y(0) | X)\}$

marginal effect of A

- regression model

$$E(Y | X, A) = \beta_0 + \beta_1 A + \beta_2 X$$

- or something more complex

$$E(Y | X, A) = f(X, A)$$

- estimand **average causal effect** or **average treatment effect (ATE)**

$$E\{Y(1)\} - E\{Y(0)\}$$

estimand: something we estimate

- under the linear model $E(Y | X, A) = \beta_0 + \beta_1 A + \beta_2 X$, the ATE is β_1
if the linear model is correct

-

$$\hat{E}(Y(a)) = \frac{1}{n} \sum_{i=1}^n \hat{E}(Y | A = a, X_i)$$

- recovers $\hat{\beta}_1$ in a linear model

- treat $Y_i(1)$ as missing data, if $A_i = 0$ (and v.v.)
- write

$$E(Y(a)) = E \left\{ \frac{1\{A = a\}Y}{\text{pr}(A = a | x)} \right\}$$

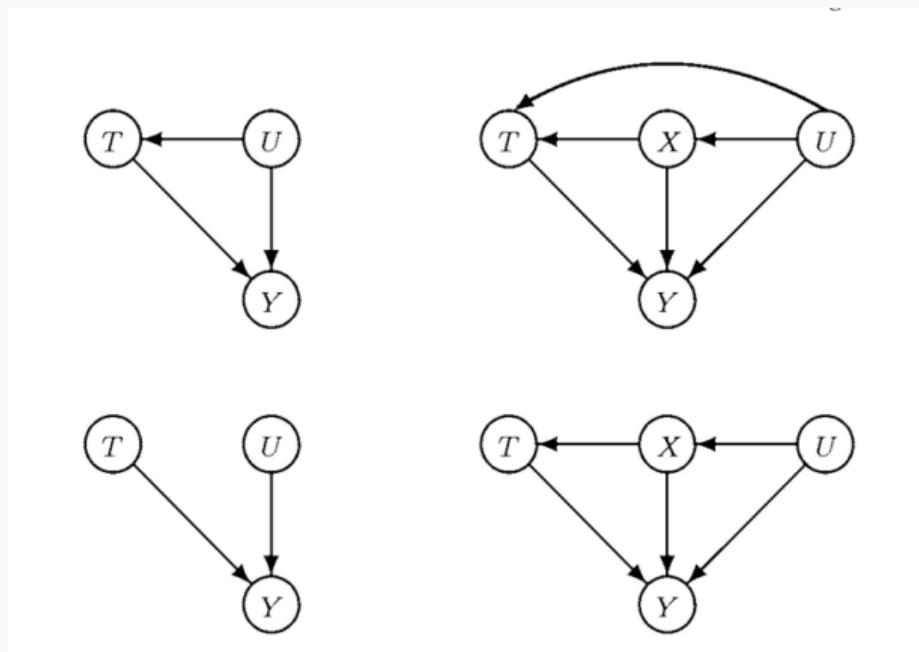
- model $\text{pr}(A = a | X)$, e.g. by logistic regression
- doubly robust estimator

of $E(Y(1))$

$$\hat{\mu}^{AIPW} = \frac{1}{n} \sum_{i=1}^n \frac{A_i Y_i}{\hat{\text{pr}}(A = 1 | X_i)} + \left\{ 1 - \frac{A_i}{\hat{\text{pr}}(A = 1 | X_i)} \right\} \hat{E}(Y(1))$$

graphs can be useful for clarifying dependence relations among random variables

Fig 9.1 SM



276 17. Directed Graphs and Conditional Independence

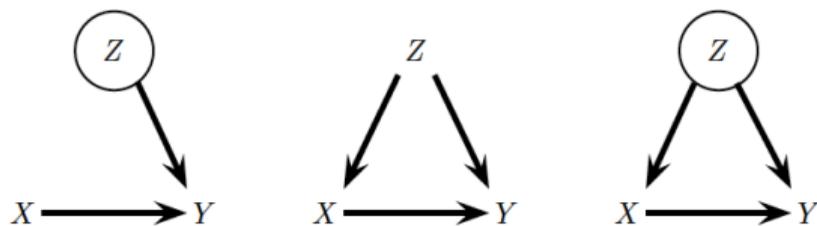


FIGURE 17.11. Randomized study; Observational study with measured confounders; Observational study with unmeasured confounders. The circled variables are unobserved.

randomized study

observational study $E(Y | x) = \int E(Y | X, Z = z) dF_Z(z)$ unobserved confounder: $\theta \neq \alpha$

SPECIAL ARTICLE

Air Pollution and Mortality at the Intersection of Race and Social Class

Kevin P. Josey, Ph.D., Scott W. Delaney, Sc.D., J.D., Xiao Wu, Ph.D.,
Rachel C. Nethery, Ph.D., Priyanka DeSouza, Ph.D., Danielle Braun, Ph.D.,
and Francesca Dominici, Ph.D.

ABSTRACT

BACKGROUND

From the Departments of Biostatistics (K.P.J., R.C.N., D.B., F.D.) and Environmental Health (S.W.D.), Harvard T.H. Chan School of Public Health, Boston; the Department of Biostatistics, Mailman School of Public Health, Columbia University, New York (X.W.); and the Department of Urban and Regional Planning,

Black Americans are exposed to higher annual levels of air pollution containing fine particulate matter (particles with an aerodynamic diameter of $\leq 2.5 \mu\text{m}$ [$\text{PM}_{2.5}$]) than White Americans and may be more susceptible to its health effects. Low-income Americans may also be more susceptible to $\text{PM}_{2.5}$ pollution than high-income Americans. Because information is lacking on exposure–response curves for $\text{PM}_{2.5}$ exposure and mortality among marginalized subpopulations categorized