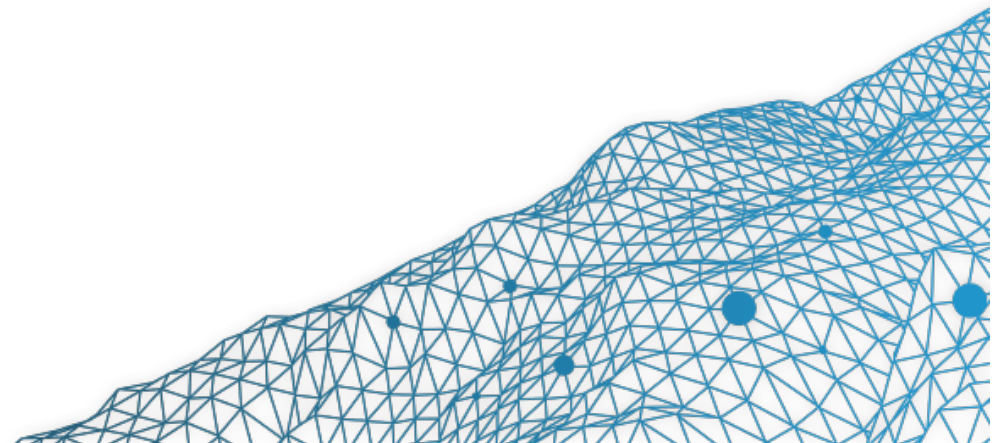


11 Observations, Experiments & Causality

Arthur Charpentier (Université du Québec à Montréal)

Machine Learning & Econometrics

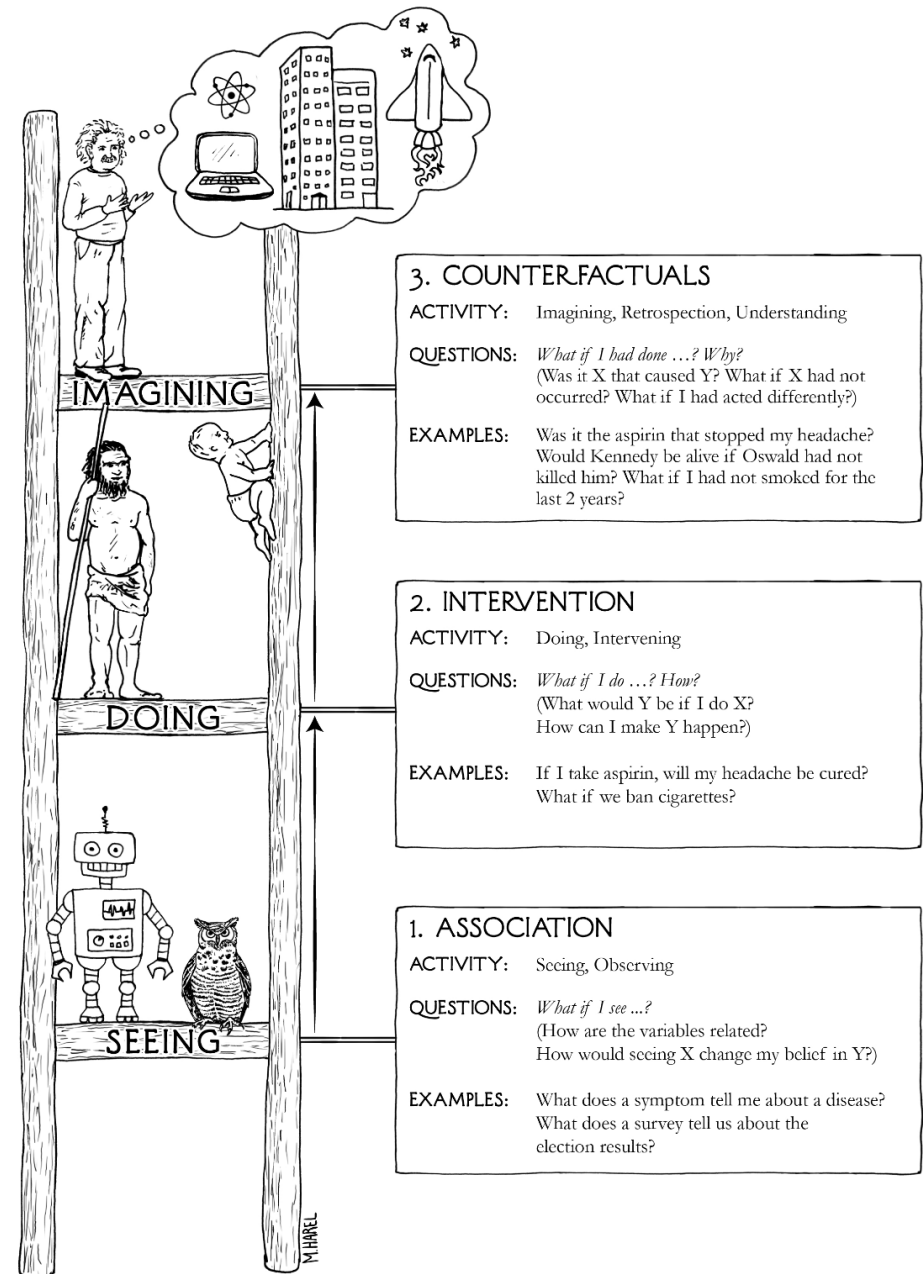
SIDE Summer School - July 2019



Causal Inference

“Social scientists know that large amounts of data will not overcome the selection problems that make causal inference so difficult”, Grimmer (2015, *We Are All Social Scientists Now: How Big Data, Machine Learning, and Causal Inference Work Together*)

(source Pearl & Mackenzie (2018, *The Book of Why*))



Causality and Counterfactuals

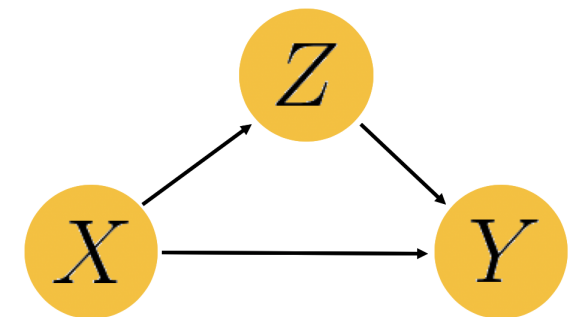
“we may define a cause to be an object, followed by another, and where all the objects similar to the first are followed by objects similar to the second. Or in other words where, if the first object had not been, the second never had existed”,
Hume (1748, [An Enquiry Concerning Human Understanding](#))

Classical conditional : if X occurred, then Y occurred

Counterfactual : if X had not occurred, then Y would not have occurred

“no causation without manipulation”, Holland (1986, [Statistics and Causal Inference](#))

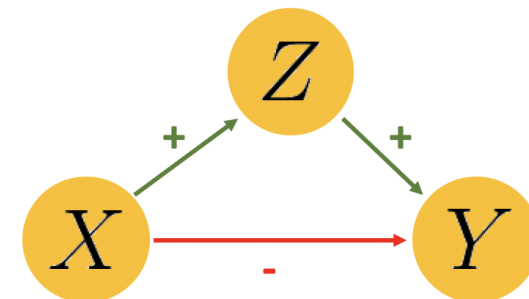
Importance of [Directed Acyclic Graphs](#) (DAG) to describe causal effects between variables



Berkeley Gender Bias

Graduate admissions data from Berkeley, 1973

- men : 8442 applications, 44% admission rate
- women : 4321 applications, 35% admission rate



Discrimination towards women ?

		A	B	C	D	E	F
M	applied	825	560	325	417	191	373
	admitted	62%	63%	37%	33%	28%	6%
F	applied	108	25	593	375	393	341
	admitted	82%	68%	34%	35%	24%	7%

see Bickel *et al.* (1975, [Sex bias in graduate admissions](#))

Simpson's Paradox & Ecological Fallacy

	hosp. A	hosp. B
total	1000	1000
survivors	800	900
deads	200	100
rate (%)	80%	90%

	hosp. A	hosp. B
total	600	900
survivors	590	870
deads	10	30
rate (%)	98%	97%

	hosp. A	hosp. B
total	400	100
survivors	210	30
deads	190	70
rate (%)	53%	30%

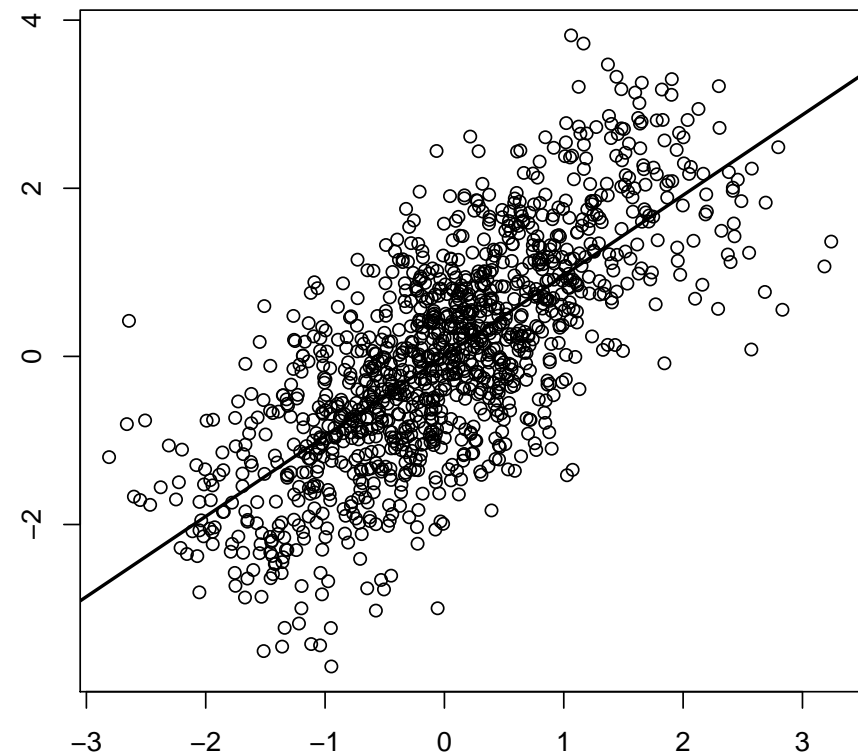
Simpson's Paradox & Ecological Fallacy

Heuristically, it is possible to have

$$\frac{a}{A} \leq \frac{b}{B} \text{ and } \frac{c}{C} \leq \frac{d}{D}$$

and at the same time

$$\frac{a+c}{A+C} \geq \frac{b+d}{B+D}$$



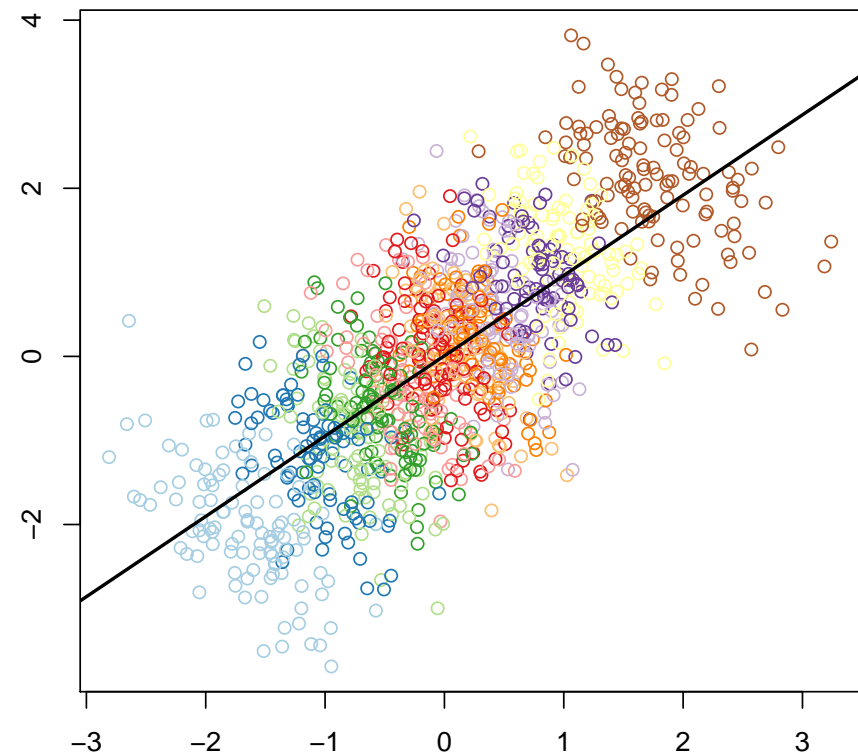
Simpson's Paradox & Ecological Fallacy

Heuristically, it is possible to have

$$\frac{a}{A} \leq \frac{b}{B} \text{ and } \frac{c}{C} \leq \frac{d}{D}$$

and at the same time

$$\frac{a + c}{A + C} \geq \frac{b + d}{B + D}$$



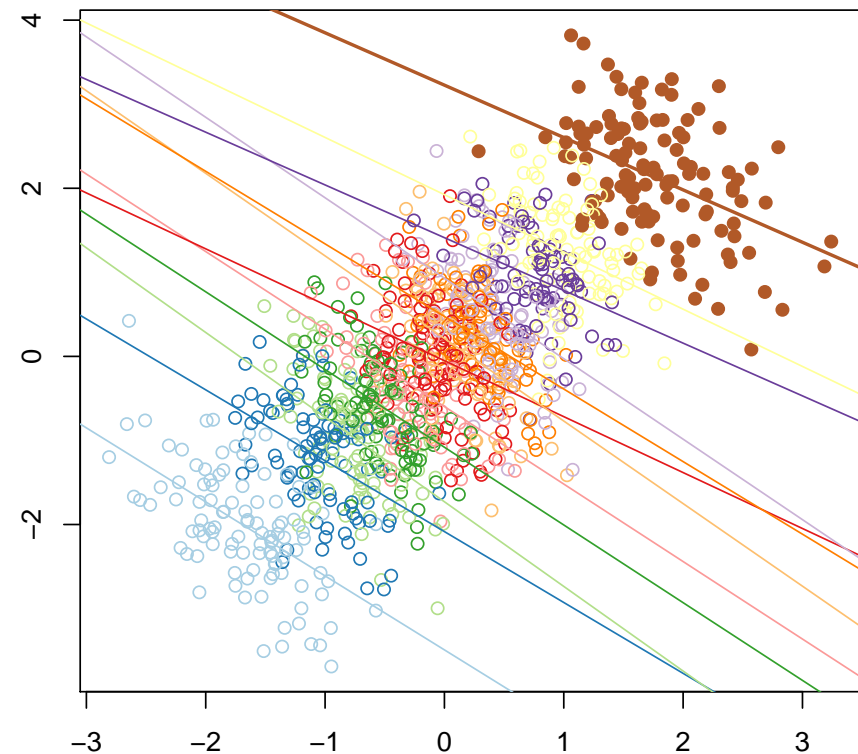
Simpson's Paradox & Ecological Fallacy

Heuristically, it is possible to have

$$\frac{a}{A} \leq \frac{b}{B} \text{ and } \frac{c}{C} \leq \frac{d}{D}$$

and at the same time

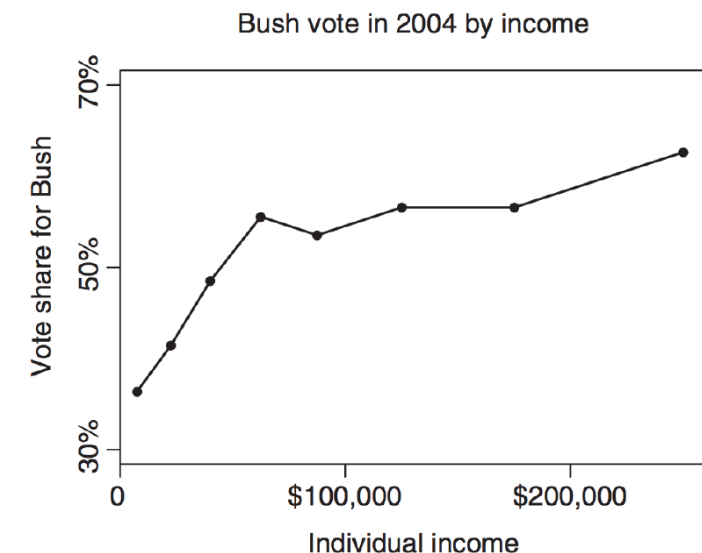
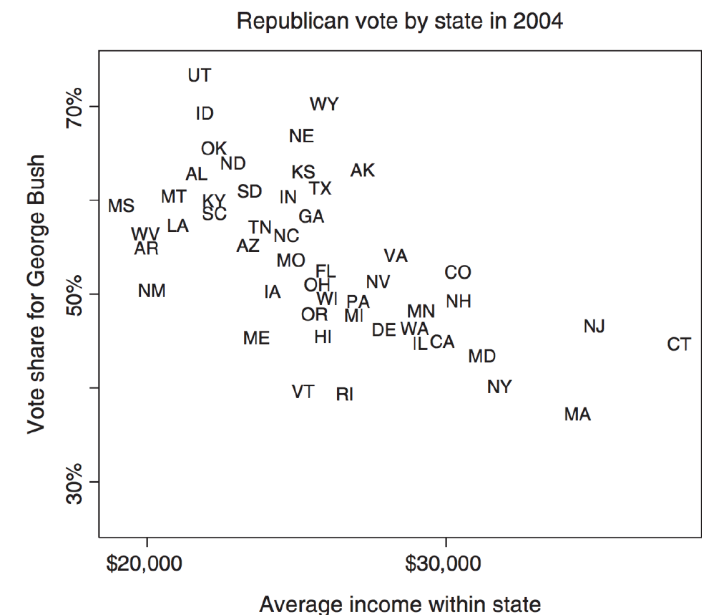
$$\frac{a+c}{A+C} \geq \frac{b+d}{B+D}$$



Simpson's Paradox & Ecological Fallacy

Very important concept in political science

see Gelman (1986, [Red State, Blue State, Rich State, Poor State: Why Americans Vote the Way They Do](#))



Conditional Independence or *could a machine create a DAG ?*

Conditional Independence

$X \perp\!\!\!\perp Y|Z$ if and only if

- $f(x, y|z) = f(x|z) \cdot f(y|z)$
- $f(x, y, z) \cdot f(z) = f(x, z) \cdot f(y, z)$
- $f(y|x, z) = f(y|z)$

i.e. once we know Z , learning the value of x does not provide additional information about X

Partial Correlation

$$\rho_{xy|z} = \frac{\rho_{xy} - \rho_{xz} \cdot \rho_{yz}}{\sqrt{(1 - \rho_{xz}^2) \cdot (1 - \rho_{yz}^2)}}$$

Conditional Independence or *could a machine create a DAG ?*

$$\text{as } X \perp\!\!\!\perp Y \implies \rho_{xy} = 0, \quad X \perp\!\!\!\perp Y|Z \implies \rho_{xy|z} = 0$$

and if (X, Y, Z) is Gaussian, the converse is true

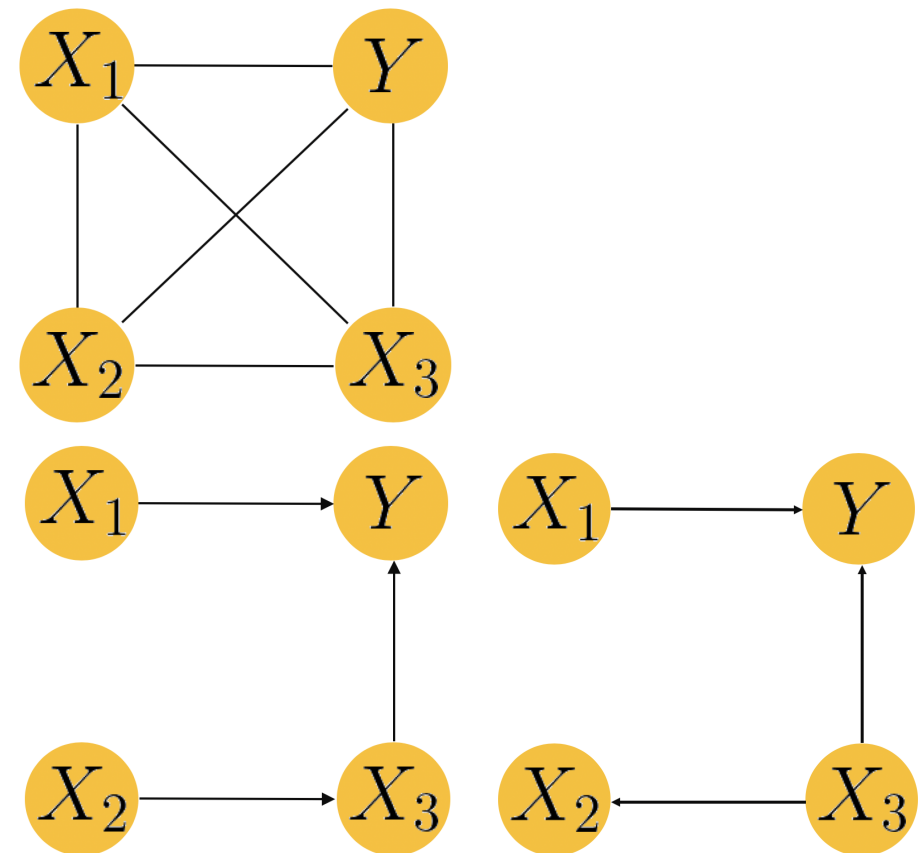
$$\text{as } \rho_{xy} = 0 \implies X \perp\!\!\!\perp Y, \quad \rho_{xy|z} = 0 \implies X \perp\!\!\!\perp Y|Z$$

Conditional Independence or *could a machine create a DAG ?*

Consider variables X_1, X_2, X_3 and Y
We have

- $X_1 \perp\!\!\!\perp X_2$
- $X_1 \perp\!\!\!\perp X_3$
- $X_2 \perp\!\!\!\perp Y | X_3$

Using conditional independence tests,
remove edges, identify colliders and chains
(might not be a unique solution...)



Conditional Independence or *could a machine create a DAG ?*

Assuming a Gaussian setting, Spirtes *et al.* (2000, [Causation, Prediction, and Search](#)) suggested to test $\rho_{xy|z} = 0$ using

$$z_n = \frac{1}{2} \sqrt{n - \dim[z]} - 3 \log \frac{|1 + \rho_{xy|z}|}{|1 - \rho_{xy|z}|}$$

One can also consider a non-parametric test, based on the distance

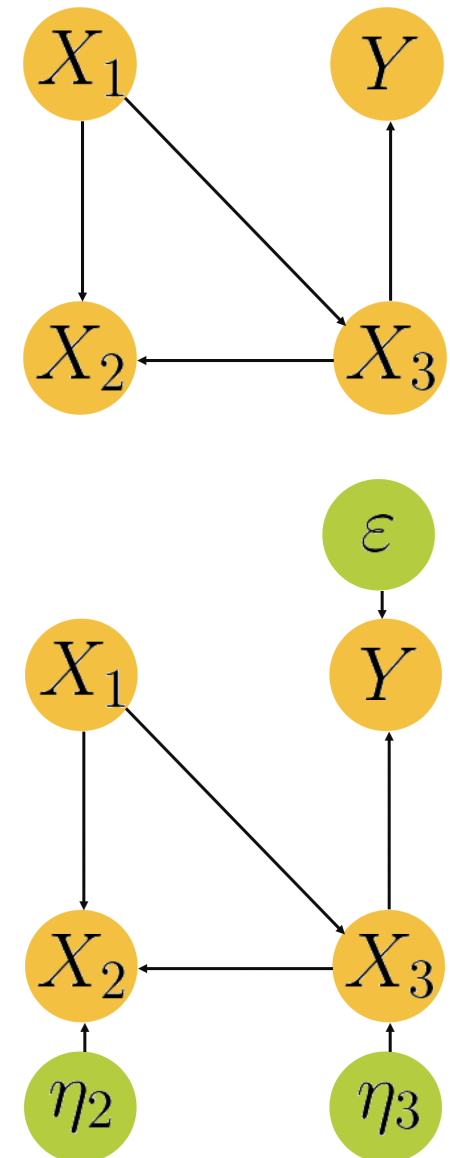
$$d(x, y, z) = \hat{f}(x, y, z) \cdot \hat{f}(z) - \hat{f}(x, z) \cdot \hat{f}(y, z)$$

(main issue here: curse of dimensionality)

From DAG to Structural Econometric Equations

There is a one-to-one mapping between a system of structural equations and a graphical model

- $X_2 = \beta_{23}X_3 + \beta_{21}X_1 + \eta_2$
- $X_3 = \beta_{31}X_1 + \eta_3$
- $Y = \beta_3X_3 + \varepsilon$



Observation and Experiment

Can we use observational data ?

Consider the following dataset (from Imai (2022) - [resume.csv](#))

1	<code>firstname</code>	first name of the fictitious job applicant
2	<code>sex</code>	sex of applicant, $\text{sex} \in \{\text{female}, \text{male}\}$
3	<code>race</code>	race of applicant, $\text{race} \in \{\text{black}, \text{white}\}$
4	<code>callback</code>	whether a call back was made, $\text{call} \in \{0, 1\}$

See

1		<code>firstname</code>	<code>sex</code>	<code>race</code>	<code>callback</code>
2	1	Allison	female	white	0
3	2	Kristen	female	white	0
4	3	Lakisha	female	black	0
5	4	Latonya	female	black	0
6	5	Carrie	female	white	0
7	6	Jay	male	white	0

Observation and Experiment

It is experimental data, collected from an experimental research design, in which a **treatment variable** is manipulated in order to examine its causal effects on an outcome variable.

The treatment refers to the race of a fictitious applicant, implied by the name given on the résumé.

The outcome variable is whether the applicant receives a callback.

We are interested in examining whether or not the résumés with different names yield varying callback rates.

	no call	call	total
black	2278	157	2435
white	2200	235	2435
total	4478	392	4870

Observation and Experiment

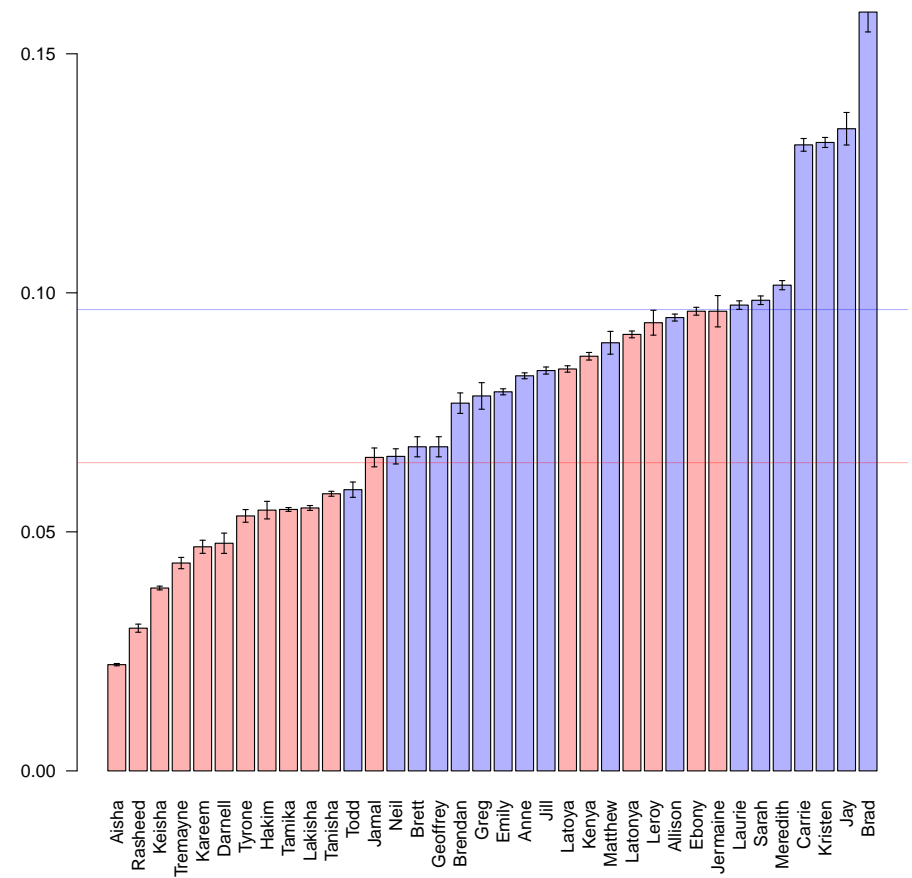
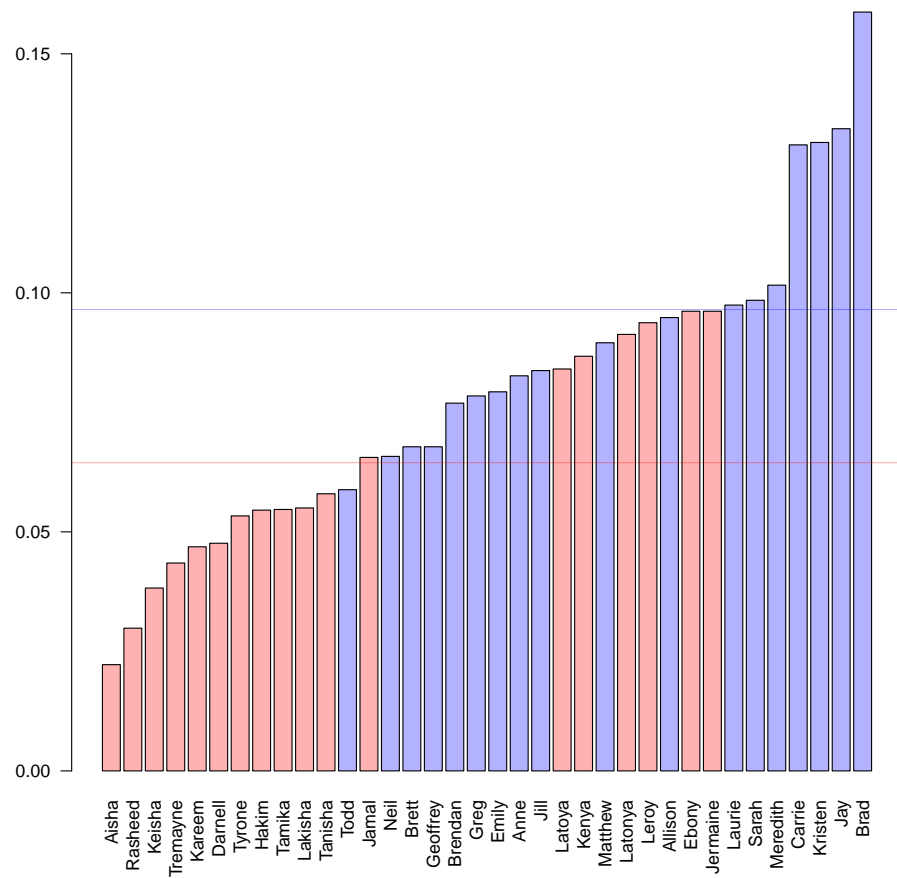
Conditional probabilities

	no call	call	total
black	50.87%	40.05%	50.00%
white	49.13%	59.95%	50.00%
total	100.00%	100.00%	100.00%

	no call	call	total
black	93.55%	6.45%	100.00%
white	90.35%	9.65%	100.00%
total	91.95%	8.05%	100.00%

3.2% points difference...

Observation and Experiment



Looking for a Counterfactual

1	firstname	sex	race	callback	
2	1	Allison	female	white	0

Would the same employer would have called back if the applicant's name were instead a stereotypically African-American.

Unfortunately, we would never observe this counterfactual outcome, because the researchers who conducted this experiment did not send out the same résumé to the same employer using Lakisha as the first name.

Let t denote the treatment (the first name0) which sound African-American ($t = 1$) or not ($t = 0$). Let $y(t)$ denote the response (call back) as a function of the treatment t . We do observe various covariate $\mathbf{x} = (x_1, x_2, \dots, x_k)$ such as the age, the education, etc.

Looking for a Counterfactual

The dataset is here

black-sounding		callback		age	education
i	name t	$y(t = 1)$	$y(t = 0)$	x_1	x_2
1	1	1	?	20	college
2	0	?	0	55	high school
3	0	?	1	40	graduate school

We would like to quantify *ceteris paribus* $y(1) - y(0)$ (called causal effect)

Fundamental problem of causal inference : we cannot observe the counterfactual outcomes

Looking for a Counterfactual

In observational studies, researchers do not conduct an intervention.

But they still want to quantify the impact of a policy. For instance the impact of an increase of minimum wage on unemployment (see Card & Krueger (1994), [minwage.csv](#)))

In 1992, New Jersey (NJ) raised the minimum wage from 4.25 to 5.05 (per hour).

1	<code>chain</code>	name of fast-food restaurant chain
2	<code>location</code>	location of restaurants (centralNJ, northNJ, PA, shoreNJ, southNJ)
3	<code>wage_before</code>	wage before minimum-wage increase
4	<code>wage_after</code>	wage after minimum-wage increase
5	<code>full_before</code>	number of full-time employees
6	<code>full_after</code>	(before and after)
7	<code>part_before</code>	number of part-time employees
8	<code>part_after</code>	(before and after)

Looking for a Counterfactual

One can look at a counterfactual, with a neighboring state - Pennsylvania (PA) - that will be our control group

This would be a cross-section comparison design

	NJ		PA	
	mean	below \$5.5	mean	below \$5.5
before	\$4.61	91.06%	\$4.65	94.03%
after	\$5.08	0.34%	\$4.61	95.52%

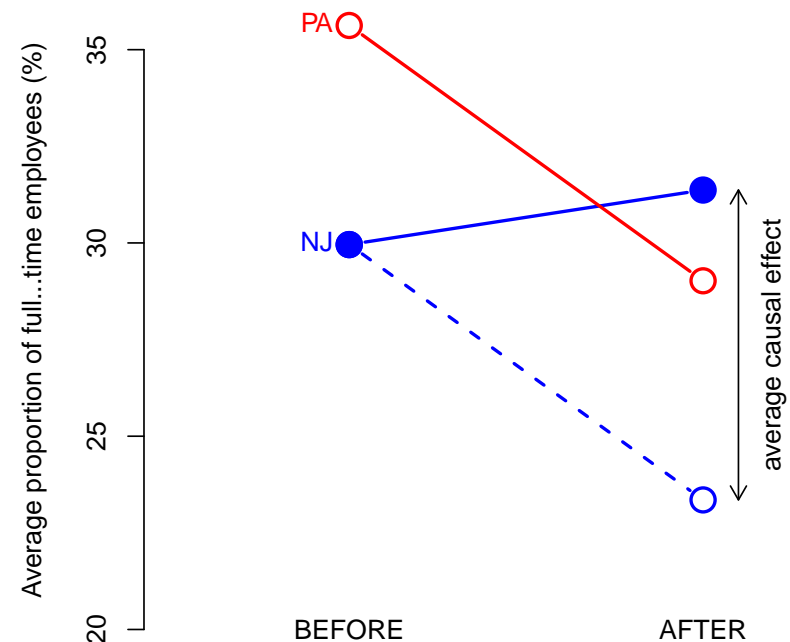
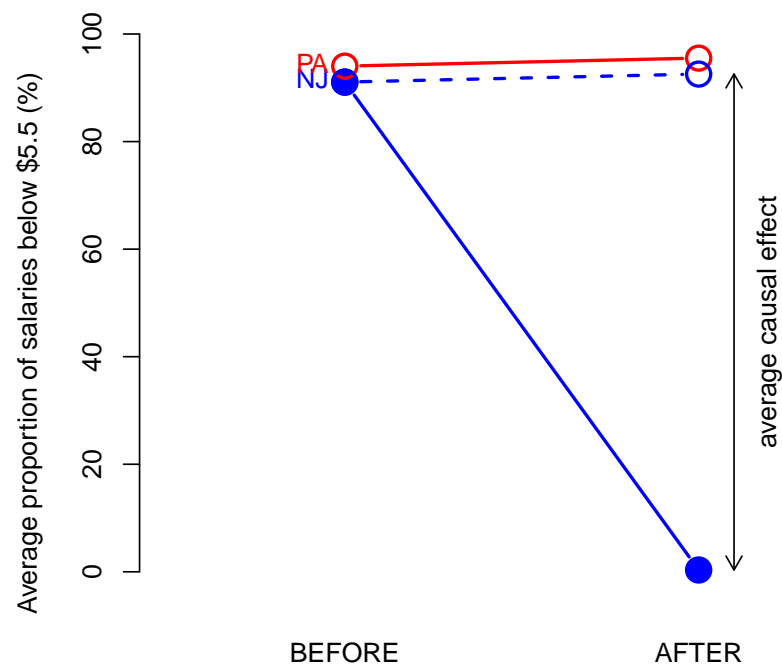
Looking for a Counterfactual

	NJ			PA		
	partial	full	proportion	partial	full	proportion
before	5423	2319	29.9%	1291	714	35.6%
after	5351	2446	31.4%	1339	547	29.0%

Difference-in-differences estimate,

$$DID = \underbrace{\bar{y}_{t=1}^{\text{after}} - \bar{y}_{t=1}^{\text{before}}}_{\text{difference in the treatment group}} - \underbrace{\bar{y}_{t=0}^{\text{after}} - \bar{y}_{t=0}^{\text{before}}}_{\text{difference in the control group}}$$

Observation and Experiment



(the counterfactual outcome for the treatment group has a time trend parallel to that of the control group)

Randomized Experiments

n individuals, that are either treated ($t_i = 1$) or not ($t_i = 0$). We observe outcome y_i for covariates \mathbf{x}_i .

We want to study **potential outcomes** $y_i(1)$ and $y_i(0)$

	turnout				
	$y_i(1)$	$y_i(0)$	t_i	$x_{1,i}$	$x_{2,i}$
1	y_1	?	1	$x_{1,1}$	$x_{2,1}$
2	?	y_2	0	$x_{1,2}$	$x_{2,2}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
n	y_n	?	1	$x_{1,n}$	$x_{2,n}$

The **causal effect** is $y_i(1) - y_i(0)$

Assume no simultaneity, no interference between individuals and same treatment

Randomized Experiments

ATE - Average Treatment Effect

Average Treatment Effect $\mathbb{E}[Y(1) - Y(0)]$

Sample Average Treatment Effect $\frac{1}{n} \sum_{i=1}^n [y_i(1) - y_i(0)]$

Assume randomization of the treatment assignment, n_1 treated, n_0 non treated

Assumption : $(Y(1), Y(0)) \perp\!\!\!\perp T$

Crude estimator (difference in means), $\hat{\tau} = \frac{1}{n_1} \sum_{i:t_i=1} y_i - \frac{1}{n_0} \sum_{i:t_i=0} y_i$, or

$$\hat{\tau} = \sum_{i=1}^n \frac{t_i y_i}{n_1} - \frac{(1 - t_i) y_i}{n_0}$$

Then $\mathbb{E}[\hat{\tau}] = \mathbb{E}[Y(1) - Y(0)]$

Randomized Experiments

Intuition with a simple regression, $y_i(t_i) = \alpha + \beta t_i + \varepsilon_i$, where $\mathbb{E}[\varepsilon] = 0$.

Causal effect is measured by $y_i(1) - y_i(0) = \beta$

More generally, assume that ε is $\varepsilon(t)$ - with $\mathbb{E}[\varepsilon(T)] = 0$, then

$y_i(1) - y_i(0) = \beta + \varepsilon_i(1) - \varepsilon_i(0)$. Thus $\beta = \mathbb{E}[Y_i(1) - Y_i(0)]$ (and $\beta = \mathbb{E}[Y_i(1)]$)

and $\alpha = \mathbb{E}[Y_i(0)]$

One can consider a more general model, with $y_i(t_i) = \alpha + \beta t_i + \varepsilon(t_i)$ where $\mathbb{E}[\varepsilon(t)] = 0$.

Then $\beta = \mathbb{E}[Y_i(1) - Y_i(0)]$ and $\alpha = \mathbb{E}[Y_i(0)]$ (as previously)

Randomized Experiments

In this model - introduced in Neyman (1923, [Sur les applications de la théorie des probabilités aux expériences agricoles](#)) - (called Rubin-Neyman) let $\sigma_t^2 = \text{Var}[\varepsilon(t)]$

$$\text{Since } \hat{\tau} = \frac{1}{n_1} \sum_{i:t_i=1} y_i - \frac{1}{n_0} \sum_{i:t_i=0} y_i = \frac{1}{n_1} \sum_{i=1}^n t_i \cdot y_i - \frac{1}{n_0} \sum_{i=1}^n (1 - t_i) \cdot y_i$$

The variance estimate of $\hat{\tau}$ is usually biased

$$\mathbb{E} \left(\frac{\hat{\sigma}^2}{\sum (t_i - \bar{t})^2} \right) - \left(\frac{\sigma_0^2}{n_0} + \frac{\sigma_1^2}{n_1} \right) = \frac{(n_1 - n_0)(n - 1)}{n_1 n_0 (n - 2)} (\sigma_1^2 - \sigma_0^2)$$

so bias is null when either $n_1 = n_0$ or $\sigma_1^2 = \sigma_0^2$.

But generally biased (even asymptotically).

(Cluster) Randomized Experiments

One can also consider some cluster randomized experiments : treatment is at cluster level, see Bland (2004, [Cluster randomised trials in the medical literature](#)) or Boruch *et al.* (2004, [Estimating the effects of interventions that are deployed in many places: place-randomized trials](#))

Consider clusters $j = 1, \dots, m$, outcome is $y_{i,j}$, with treatment $t_j \in \{0, 1\}$.

Assume random assignment, i.e. $(Y_{i,j}(1), Y_{i,j}(0)) \perp\!\!\!\perp T_j$.

And for convenience, assume that n_j 's are equal...

$$\hat{\tau} = \frac{1}{m_1} \sum_{j:t_j=1} \bar{y}_j - \frac{1}{m_0} \sum_{j:t_j=0} \bar{y}_j \text{ where } \bar{y}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} y_{i,j}$$

(Cluster) Randomized Experiments

As before, $\mathbb{E}[\hat{\tau}] = \mathbb{E}[Y(1) - Y(0)]$

Exact variance is $\text{Var}[\hat{\tau}] = \frac{\text{Var}[\bar{Y}(1)]}{m_1} + \frac{\text{Var}[\bar{Y}(0)]}{m_0}$,

where $\text{Var}[\bar{Y}(t)] = \frac{\sigma_t^2}{n} [1 + (n-1)\rho_t] \leq \sigma_t^2$

(due to (possible) intraclass correlation ρ_t)

Observational Studies & Propensity Score

In the regression discontinuity design, we want to find an (arbitrary) cutoff point c that determines the treatment assignment, i.e. $t_i = \mathbf{1}_{x_i \geq c}$.

We want to estimate $\mathbb{E}[Y_i(1) - Y_i(0) | X_i = c]$

Identification of the Average Treatment Effect :

Assumption 1: Overlap, $\forall x, \mathbb{P}[T = 1 | X = x] \in (0, 1)$

Assumption 2: Ignorability $(Y(1), Y(0)) \perp\!\!\!\perp T | X = x, \forall x$

Set $\mu(t, x) = \mathbb{E}[Y_i(t) | T = t, X = x]$

Crude estimator (difference in means),

$$\hat{\tau} = \frac{1}{n} \sum_{i=1}^n \hat{\mu}(1, x) - \hat{\mu}(0, x)$$

Observational Studies & Regression Discontinuity

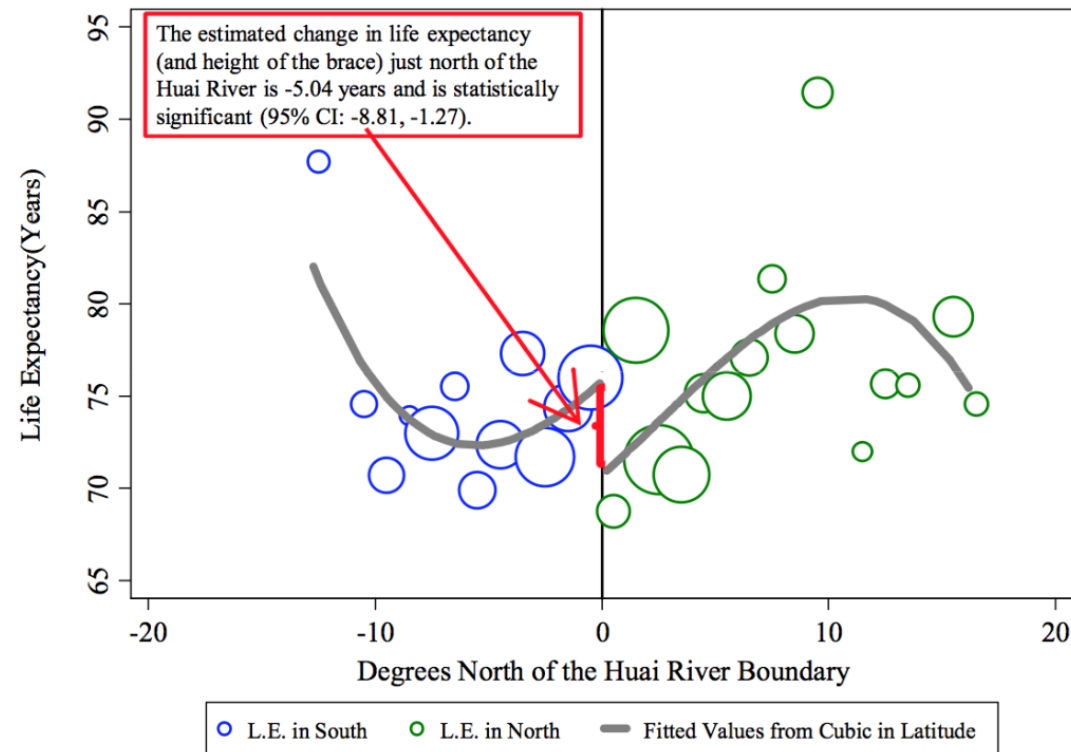


Fig. 3. The plotted line reports the fitted values from a regression of life expectancy on a cubic in latitude using the sample of DSP locations, weighted by the population at each location.

See Imbens & Lemieux (2008, [Regression Discontinuity Designs](#)).

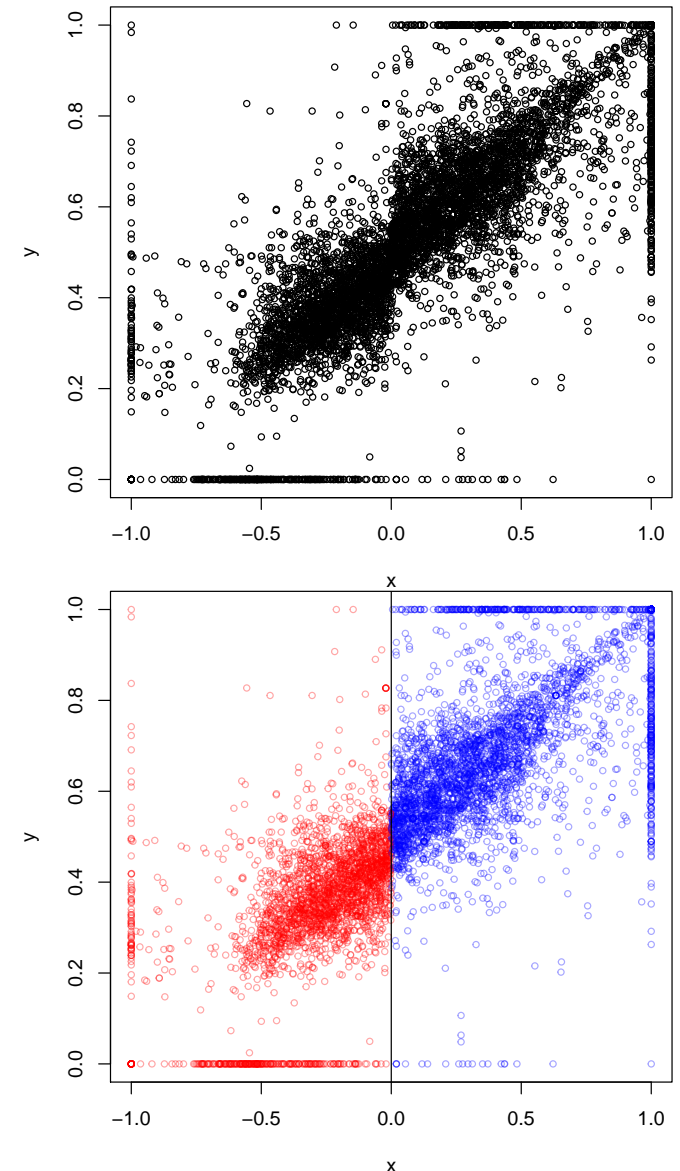
Observational Studies & Regression Discontinuity

Consider the dataset from Lee (2008, [Randomized experiments from non-random selection in U.S. House elections](#)) - “*How would the Democratic party have performed in period 2, had they not held the seat (i.e. had the Democrats lost the election in period 1)*”

```
1 > library(RDDtools)
2 > data(Lee2008)
```

We want to test if there is a discontinuity
(in 0, x was rescaled here)

- with parametric tools
- with nonparametric tools



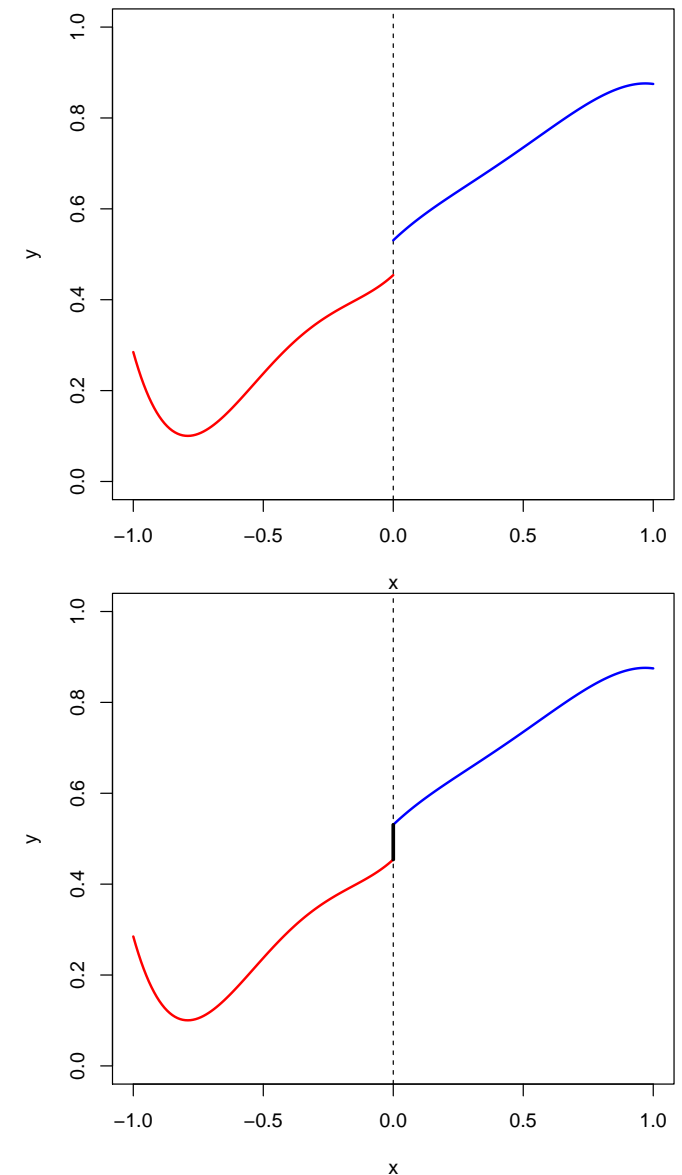
Observational Studies & Regression Discontinuity

Use some 4th order polynomial, on each part

```

1 > idx1 = (Lee2008$x>0)
2 > reg1 = lm(y~poly(x,4),data=Lee2008[idx1,])
3 > idx2 = (Lee2008$x<0)
4 > reg2 = lm(y~poly(x,4),data=Lee2008[idx2,])
5 > s1=predict(reg1,newdata=data.frame(x=0))
6 > s2=predict(reg2,newdata=data.frame(x=0))
7 > abs(s1-s2)
8           1
9 0.07659014

```

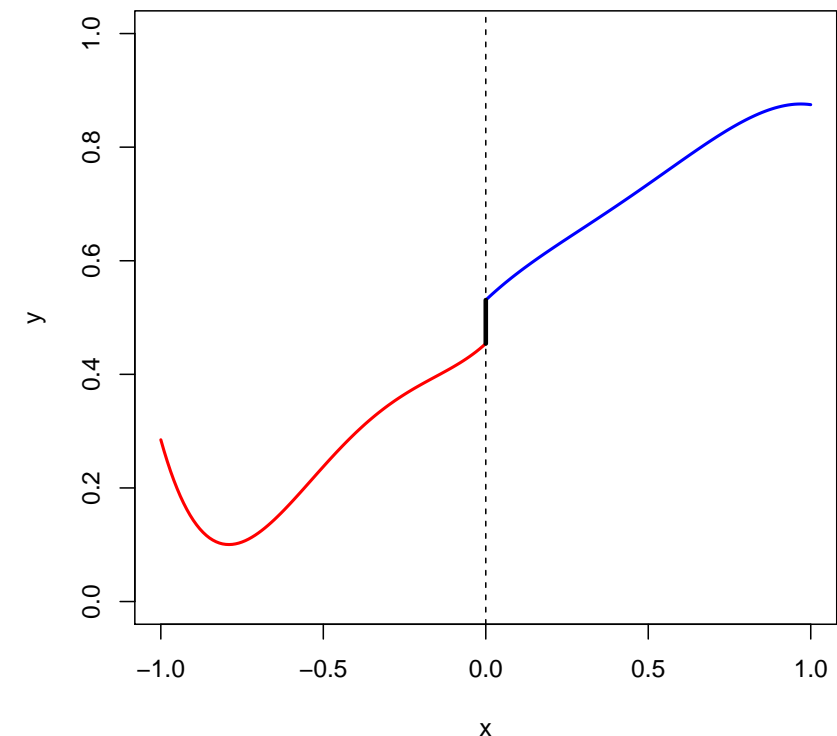


Observational Studies & Regression Discontinuity

```

1 > reg_para <- RDDreg_lm(RDDdata(y =
    Lee2008$y, x = Lee2008$x, cutpoint
    = 0), order = 4)
2 > reg_para
3 ### RDD regression: parametric ###
4 Polynomial order: 4
5 Slopes: separate
6 Number of obs: 6558 (left: 2740,
    right: 3818)
7
8 Coefficient:
9 Estimate Std. Error t value Pr(>|t|)
10 D 0.076590 0.013239 5.7851 7.582e-09
    ***

```



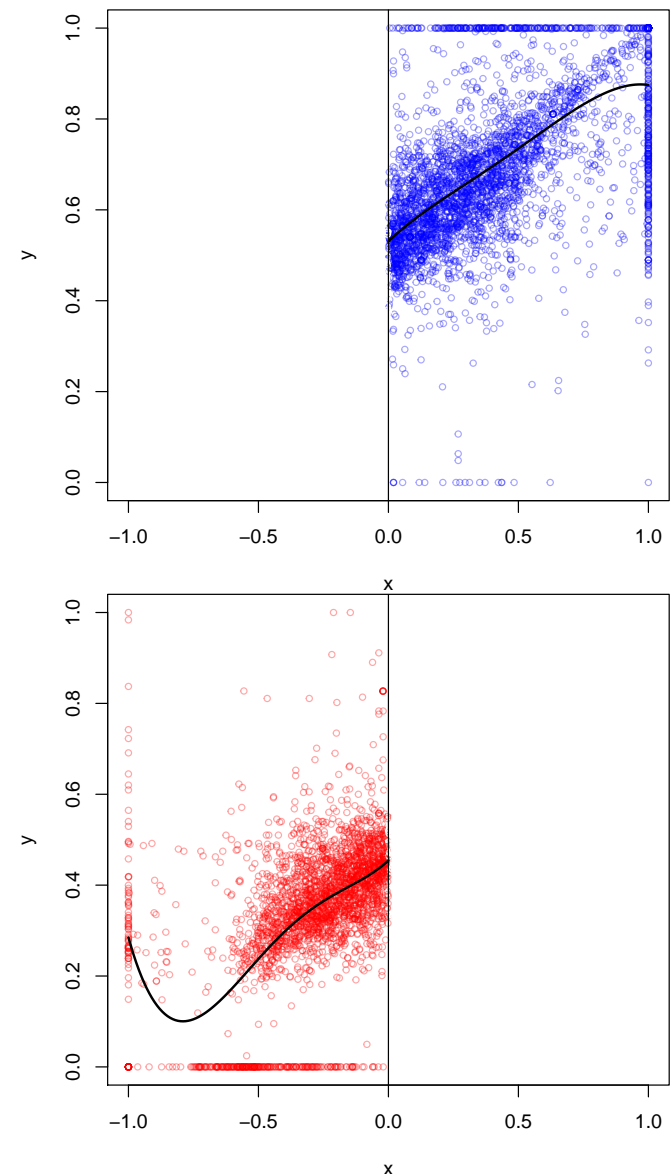
Observational Studies & Regression Discontinuity

or use a simple local regression, see [Imbens & Kalyanaraman \(2012\)](#).

```

1 > reg1 = ksmooth(Lee2008$x[idx1],
  Lee2008$y[idx1], kernel = "normal",
  bandwidth = 0.1)
2 > reg2 = ksmooth(Lee2008$x[idx2],
  Lee2008$y[idx2], kernel = "normal",
  bandwidth = 0.1)
3 > s1 = reg1$y[1]
4 > s2 = reg2$y[length(reg2$y)]
5 > abs(s1-s2)
6 [1] 0.09883813

```

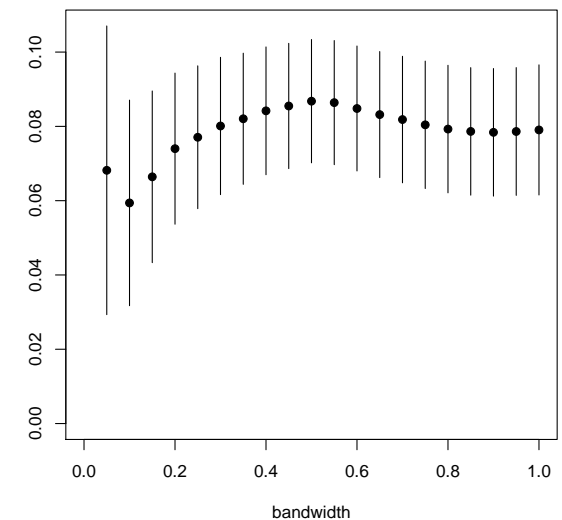
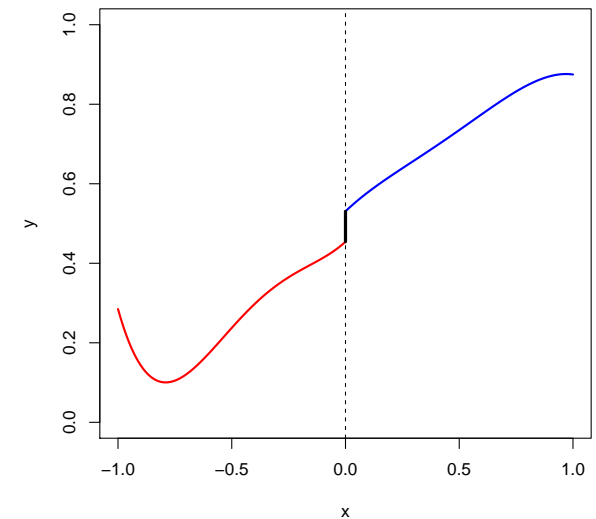


Observational Studies & Regression Discontinuity

```

1 > reg_nonpara <- RDDreg_np(RDDobject = Lee2008_
   rdd, bw = .1)
2 > print(reg_nonpara)
3 ### RDD regression: nonparametric local linear
4   Bandwidth: 0.1
5   Number of obs: 1209 (left: 577, right: 632)
6
7   Coefficient:
8   Estimate Std. Error z value Pr(>|z|)
9 D 0.059397 0.014119 4.207 2.588e-05 ***

```



Observational Studies & Propensity Score

Propensity Score

The Propensity Score is the probability to receive the treatment $\pi(x) = \mathbb{P}[T = 1|X = x]$

Assumption: Balancing property $T \perp\!\!\!\perp X|\pi(X)$

Assumption: Exogeneity given the propensity score $(Y(1), Y(0)) \perp\!\!\!\perp T|\pi(x), \forall x$

$$\hat{\tau} = \sum_{i=1}^n \frac{t_i y_i}{n_1} - \frac{(1 - t_i) y_i}{n_0}$$

consider

$$\hat{\tau} = \sum_{i=1}^n \frac{t_i y_i}{n \hat{\pi}(x_i)} - \frac{(1 - t_i) y_i}{n(1 - \hat{\pi}(x_i))}$$

Observational Studies & Propensity Score

Balancing condition: $\mathbb{E} \left[\frac{TY}{\pi(X)} - \frac{(1-T)Y}{1-\pi(X)} \right] = 0$

See Lunceford & Davidian (2004, [Stratification and Weighting Via the Propensity Score](#))

Consider also Robins's estimator,

$$\hat{\tau} = \left(\sum_{i=1}^n \frac{\hat{\mu}(1, x_i)}{n} + \frac{t_i(y_i - \hat{\mu}(1, x_i))}{n\hat{\pi}(x_i)} \right) - \left(\sum_{i=1}^n \frac{\hat{\mu}(0, x_i)}{n} + \frac{(1-t_i)(y_i - \hat{\mu}(0, x_i))}{n(1-\hat{\pi}(x_i))} \right)$$

Observational Studies & Matching

Consider a simple linear regression model, $y_i = \alpha + \beta t_i + \varepsilon_i$, then

$$\hat{\beta} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i(1) - \hat{y}_i(0))$$

where

$$\hat{y}_i(1) = \begin{cases} y_i & \text{if } t_i = 1 \\ \frac{1}{n_1} \sum_{j=1}^n t_j \cdot y_j & \text{if } t_i = 0 \end{cases}$$

$$\hat{y}_i(0) = \begin{cases} \frac{1}{n_0} \sum_{j=1}^n (1 - t_j) \cdot y_j & \text{if } t_i = 1 \\ y_i & \text{if } t_i = 0 \end{cases}$$

Observational Studies & Matching

Consider a simple fixed effect regression model, $y_{i,t} = \alpha_i + \beta t_{i,t} + \varepsilon_{i,t}$

$$\hat{\beta} = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T (\hat{y}_{i,t}(1) - \hat{y}_{i,t}(0))$$

where, if $n_{1,i} = \sum_{\tau=1}^T t_{i,\tau}$ and $n_{0,i} = T - n_{1,i} = \sum_{\tau=1}^T (1 - t_{i,\tau})$

$$\hat{y}_{i,t}(1) = \begin{cases} y_{i,t} & \text{if } t_{i,t} = 1 \\ \frac{1}{n_{1,i}} \sum_{\tau=1}^T t_{i,\tau} \cdot y_{i,\tau} & \text{if } t_{i,t} = 0 \end{cases}$$

$$\hat{y}_{i,t}(0) = \begin{cases} \frac{1}{n_{0,i}} \sum_{\tau=1}^T (1 - t_{i,\tau}) \cdot y_{i,\tau} & \text{if } t_{i,t} = 1 \\ y_{i,t} & \text{if } t_{i,t} = 0 \end{cases}$$

Observational Studies & Matching

$\hat{\beta}$ is actually the weighted least square estimate,

$$\hat{\beta} = \underset{(\alpha, \beta)}{\operatorname{argmin}} \left\{ \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \omega_{i,t} (y_{i,t} - \alpha_i - \beta t_{i,t}) \right\}$$

with $\omega_{i,t} = T/n_{t,i}$

More generally, we can have heterogeneity, i.e. $y_{i,t} = \alpha_i + \beta t_{i,t} + \gamma^\top \mathbf{x}_{i,t} + \varepsilon_{i,t}$

Let $\pi(\mathbf{x})$ denote the propensity score $\mathbb{P}[T = 1 | \mathbf{X} = \mathbf{x}]$.

Observational Studies & Matching

We have the following (ex-post) interpretation

$$y_{i,t} - \hat{\gamma}^\top \mathbf{x}_{i,t} = \alpha_i + \beta t_{i,t} + \varepsilon_{i,t}$$

so use $\hat{\beta}$ the weighted least square estimate,

$$\hat{\beta} = \underset{(\alpha, \beta)}{\operatorname{argmin}} \left\{ \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \omega_{i,t} (y_{i,t}^* - \alpha_i - \beta t_{i,t}) \right\}$$

with $\omega_{i,t} = T/n_{t,i}$ and

$$y_{i,t}^* = \begin{cases} \frac{1}{\sum_{\tau} t_{i,\tau} \hat{\pi}(\mathbf{x}_{i,\tau})^{-1}} \sum_{\tau=1}^T t_{i,\tau} \hat{\pi}(\mathbf{x}_{i,\tau})^{-1} \cdot y_{i,\tau} & \text{if } t_{i,t} = 1 \\ \frac{1}{\sum_{\tau} (1 - t_{i,\tau}) (1 - \hat{\pi}(\mathbf{x}_{i,\tau}))^{-1}} \sum_{\tau=1}^T (1 - t_{i,\tau}) (1 - \hat{\pi}(\mathbf{x}_{i,\tau}))^{-1} \cdot y_{i,\tau} & \text{if } t_{i,t} = 0 \end{cases}$$

Observational Studies & Matching

This is just the general case of the simple 2 time-period difference in differences

$$y_{i,\tau}(t) = \alpha_i + \beta t + \gamma \tau + \varepsilon_{i,\tau}$$

where τ is the time, $\tau \in \{0, 1\}$. Hence

$$y_{i,0}(0) = \alpha_i + \varepsilon_{i,0}$$

$$y_{i,1}(0) = \alpha_i \gamma + \varepsilon_{i,1}$$

$$y_{i,1}(1) = \alpha_i \beta + \gamma + \varepsilon_{i,1}$$

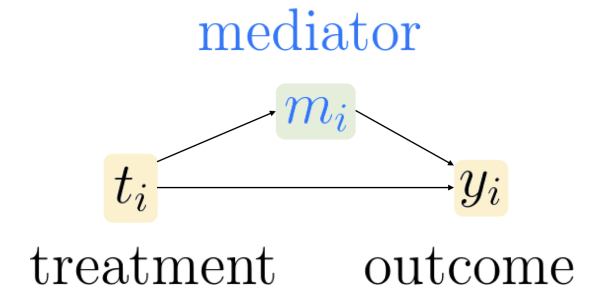
$$\text{Assumption: } \mathbb{E}[y_{i,1}(0) - y_{i,0}(0) | T_{i,1} = t] = \gamma$$

$$\text{i.e.: } \mathbb{E}[\varepsilon_{i,1} - \varepsilon_{i,0} | T_{i,1} = t] = 0$$

Imai *et al.* (2011, [Unpacking the Black Box of Causality](#)).

Experimental / Observational Studies & Mediation

Assume that there is a possible mediator m (possible indirect effects), see Baron & Kenny (1986, [The Moderator-Mediator Variable Distinction](#))



Brader, Valentino & Suhay (2008, [What Triggers Public Opposition to Immigration?](#))

Treatment $t \in \{\text{caucasian, latino}\}$ - randomized experiment

Outcome y is preference over immigration policy - measured

Mediation m is anxiety - measured

Bertrand & Mullainathan (2008, [Are Emily and Greg More Employable than Lakisha and Jamal?](#))

Treatment $t \in \{\text{white name, black name}\}$ - randomized experiment

Outcome $y \{\text{callback, no callback}\}$ - measured

Experimental / Observational Studies & Mediation

Mediation m is perceived qualifications of applicants

(1) Regress y on t : $y_i = \alpha_1 + \beta_1 t_i + \boldsymbol{\xi}_1^\top \mathbf{x}_i + \varepsilon_{1,i}$

(2) Regress m on t : $m_i = \alpha_2 + \beta_2 t_i + \boldsymbol{\xi}_2^\top \mathbf{x}_i + \varepsilon_{2,i}$

(3) Regress y on m and t : $y_i = \alpha_3 + \beta_3 t_i + \gamma_3 m_i + \boldsymbol{\xi}_3^\top \mathbf{x}_i + \varepsilon_{3,i}$

for some control variables \mathbf{x}

β_1 is ATE, and can be decomposed in $\beta_1 = \underbrace{\beta_3}_{\text{direct}} + \underbrace{\gamma\beta_2}_{\text{indirect}}$

Experimental / Observational Studies & Mediation

Here m might depend on t , $m_i(t_i)$ and y depends on m and t , $y_i(t_i, m_i(t_i))$.

- Indirect (mediation) causal effect is $\delta_i(t) = y_i(t, m_i(1)) - y_i(t, m_i(0))$
(identification pb)

- Direct causal effect is $\zeta_i(t) = y_i(1, m_i(t)) - y_i(0, m_i(t))$

Total causal effect is $\tau_i = y_i(1, m_i(1)) - y_i(0, m_i(0))$ and

$$\tau = \frac{[\delta_i(0) + \zeta_i(0)] + [\delta_i(1) + \zeta_i(1)]}{2}$$

$\delta_i(t)$ measures counterfactuals about treatment-induced mediator values

- Controlled direct effects: $\xi_i(t, m, w) = y_i(t, m) - y_i(t, w)$
- Interaction effects: $\xi_i(1, m, w) - \xi_i(0, m, w)$

See `mediation` library

Experimental / Observational Studies & Mediation

In Brader, Valentino & Suhay (2008, [What Triggers Public Opposition to Immigration?](#))

Randomized treatment: t is randomized given \mathbf{x} , $(Y, M) \perp\!\!\!\perp T|X$

Sequential ignorability: m is randomized given \mathbf{x} and t , $Y \perp\!\!\!\perp M|T, X$

Then

$$\bar{\delta}(t) = \int \mathbb{E}[Y_i|m, t, \mathbf{x}_i] (d\mathbb{P}[m|1, \mathbf{x}_i] - d\mathbb{P}[m|0, \mathbf{x}_i]) d\mathbb{P}(\mathbf{x}_i)$$

$$\bar{\zeta}(t) = \int (\mathbb{E}[Y_i|m, 1, \mathbf{x}_i] - \mathbb{E}[Y_i|m, 0, \mathbf{x}_i]) d\mathbb{P}[m|t, \mathbf{x}_i] d\mathbb{P}(\mathbf{x}_i)$$

Use linear structural equations,

$$\begin{cases} m_i = \alpha_2 + \beta_2 t_i + \boldsymbol{\xi}_2^\top \mathbf{x}_i + \varepsilon_{2,i} \\ y_i = \alpha_3 + \beta_3 t_i + \gamma_3 m_i + \boldsymbol{\xi}_3^\top \mathbf{x}_i + \varepsilon_{3,i} \end{cases}$$

Separate least square to estimate $\hat{\beta}_2$ and $\hat{\gamma}$ and then [Sobel test](#) for significance

Experimental / Observational Studies & Mediation

In Bertrand & Mullainathan (2008, [Are Emily and Greg More Employable than Lakisha and Jamal?](#)), crossover design

- (1) standard randomized experiment: send Jamal's CV, record y
- (2) change treatment to the opposite status, keep mediator fixed: send CV (same qualification) as Greg, record y

similar to Hainmueller & Hiscox (2010, [Attitudes toward Highly Skilled and Low-skilled Immigration](#))

Baron-Kenny Procedure

- (1) regress y on t (significant relationship)
- (2) regress m on t (significant relationship)
- (3) regress y on m and t (significant relationship between y and m)

Experimental / Observational Studies & Mediation

E.g. binary mediator $m \in \{0, 1\}$,

$$\text{logit}[\mathbb{E}[M_i|t_i, \mathbf{x}_i]] = \alpha_2 + \beta_2 t_i + \boldsymbol{\xi}_2^\top \mathbf{x}_i$$

$$\text{logit}[\mathbb{E}[Y_i|m_i, t_i, \mathbf{x}_i]] = \alpha_3 + \beta_3 t_i + \gamma_3 m_i + \boldsymbol{\xi}_3^\top \mathbf{x}_i +$$

Cannot use $\beta_2\gamma$, or $\beta_1 - \beta_3$ if

$$\text{logit}[\mathbb{E}[Y_i|t_i, \mathbf{x}_i]] = \alpha_1 + \beta_1 t_i + \boldsymbol{\xi}_1^\top \mathbf{x}_i$$

Assume that $\rho_{2,3} = \text{Corr}[\varepsilon_2, \varepsilon_3]$. Sequential ignorability means $\rho_{2,3} = 0$.

$$\bar{\delta}(0) = \bar{\delta}(1) = \frac{\beta_2 \sigma_1}{\sigma_2} \left(\rho_{1,2} - \rho_{2,3} \sqrt{\frac{1 - \rho_{1,2}^2}{1 - \rho_{2,3}^2}} \right)$$

Thus, $\bar{\delta}(t) = 0$ means $\rho_{2,3} = \rho_{1,2}$.

Experimental / Observational Studies & Mediation

```
1 > fit_m <- lm(mediator ~ treat + x)
2 > fit_y <- lm(y ~ treat + mediator + x)
```

For the mediation analysis use

```
1 > med <- mediation::mediate(fit_m, fit_y, treat="treat", mediator="
    mediator")
```

and for the sensitivity analysis

```
1 > mediation::medsens(med)
2 > plot(medsens(med), "rho")
```

Problem, sometimes m is difficult to manipulate.

Use instrumental variables

Causal Trees

Prediction with Machine Learning techniques ?

Get the smallest mean-squared error in a test set...

Given a candidate $\hat{\mu}(\mathbf{x})$, use $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}(\mathbf{x}_i))^2$

With a tree, $\hat{\mu}(\mathbf{x})$ is the sample mean of y_i 's within leaf $\ell(\mathbf{x})$

The in-sample goodness of fit measure is the **mse**, $\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}(\mathbf{x}_i))^2$

Try to minimize $\text{mse} + \lambda \cdot \text{number of leaves}$

Select λ with lowest out-of-sample goodness of fit measure, or use cross-validation.

```
1 > tree <- rpart(y ~ x, method="anova")
2 > i <- which.min(tree$cptable[, "xerror"])
3 > tree_2 <- prune(tree, cp=tree$cptable[i, "CP"])
```

Causal Trees

Consider our causal model, and let τ_i denote the treatment effect $\tau_i = y_i(1) - y_i(0)$

Consider possible heterogeneity, set $\mu(t, \mathbf{x}) = \mathbb{E}[Y(t)|T = t, \mathbf{X} = \mathbf{x}]$ and $\tau(\mathbf{x}) = \mu(1, \mathbf{x}) - \mu(0, \mathbf{x})$

To estimate $\tau(\mathbf{x})$ a partition tree can be more interesting than a linear predictor.

Approach 1 : analyze the two groups separately

- estimate $\hat{\mu}(1, \mathbf{x})$ on the sub-dataset where $t_i = 1$
- estimate $\hat{\mu}(0, \mathbf{x})$ on the sub-dataset where $t_i = 0$
- use propensity score weighting
- (use within group cross-validation to tune parameters)
- prediction is $\hat{\tau}(\mathbf{x}) = \hat{\mu}(1, \mathbf{x}) - \hat{\mu}(0, \mathbf{x})$

Approach 2 : estimate $\mu(t, \mathbf{x})$ on both covariates

Instruments

Recall that in a regression model, $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$, with endogeneous variables $\mathbb{E}(\mathbf{x}_i^\top \varepsilon_i) \neq \mathbf{0}$ and least squares estimators are not convergent ($\text{plim} \hat{\boldsymbol{\beta}} \neq \boldsymbol{\beta}$ as $n \rightarrow \infty$)

Classical motivations

- variable omission : the true model is $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \mathbf{z}_i^\top \boldsymbol{\gamma} + \varepsilon_i$, then

$$\text{plim} \hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + \mathbb{E}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbb{E}(\mathbf{X}^\top \mathbf{Z}) \boldsymbol{\gamma} \neq \boldsymbol{\beta}$$

See Mincer equation, where y is the wage, x the number of years of study, but there might be some ability bias

- measurement error : the true model is $y_i = \mathbf{x}_i^{*\top} \boldsymbol{\beta} + \varepsilon_i$ where variables are measured with error, $\mathbf{x}_i = \mathbf{x}_i^* + \boldsymbol{\eta}_i$ where $\mathbb{E}[\boldsymbol{\eta}_i | \mathbf{x}_i^*] = \mathbf{0}$ and $\boldsymbol{\eta}_i \perp \varepsilon_i$. Thus, $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \nu_i$ where $\nu_i = \varepsilon_i - \boldsymbol{\beta}^\top \boldsymbol{\eta}_i$. Here, $\mathbb{E}[\nu_i \mathbf{x}_i] = -\boldsymbol{\beta} \text{Var}[\boldsymbol{\eta}] \neq \mathbf{0}$, and

$$\text{plim} \hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + \mathbb{E}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbb{E}(\mathbf{X}^\top \boldsymbol{\eta}) \boldsymbol{\gamma} \neq \boldsymbol{\beta}$$

- simultaneity bias

Causal Trees

Assume here that $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + b\tilde{x}_i + \varepsilon_i$, where

- $\mathbb{E}(\mathbf{x}_i^\top \varepsilon_i) = \mathbf{0}$
- $\mathbb{E}(\tilde{x}_i^\top \varepsilon_i) \neq 0$

We will use instruments \mathbf{z} that

- should be exogeneous, $\mathbb{E}(\mathbf{z}_i^\top \varepsilon_i) = \mathbf{0}$
- should add information to \mathbf{x} , in the sense that in the regression $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \mathbf{z}_i^\top \boldsymbol{\gamma} + \varepsilon'_i$, $\boldsymbol{\gamma} \neq \mathbf{0}$

Set

$$\hat{x} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \tilde{x}$$

and observe that

$$\hat{x} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{x} = \mathbf{x}$$

Causal Trees

Thus, define the **filtrated model**

$$y_i = \hat{\mathbf{x}}_i^\top \boldsymbol{\beta} + b\hat{x}_i + \hat{\varepsilon}_i = \mathbf{x}_i^\top \boldsymbol{\beta} + b\hat{x}_i + \hat{\varepsilon}_i$$

where $\hat{\varepsilon} = \varepsilon + (\tilde{x} - \hat{x})$.

From the rank condition, $(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}$ exists

The filtrated model is exogeneous, in the sens that $\text{plim} \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{x}}_i \hat{\varepsilon}_i = 0$

The IV estimate is here $\hat{\boldsymbol{\beta}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y}$

More generally,

$$y_i = \underbrace{\mathbf{x}_i^\top \boldsymbol{\beta}}_{\text{exogeneous}} + \underbrace{\tilde{\mathbf{x}}_i^\top \boldsymbol{\gamma}}_{\text{endogeneous}} + \varepsilon_i$$

Causal Trees

Consider instruments \mathbf{z}_i and set $\mathbf{z}_i^* = (\mathbf{z}_i, \mathbf{x}_i)$, such that they

- should be exogeneous, $\mathbb{E}[\mathbf{z}_i^{*\top} \varepsilon_i] = \mathbf{0}$
- should add information to \mathbf{x} , in the sense that $\text{rank}(\mathbb{E}[\mathbf{z}_i^{*\top} \mathbf{x}_i]) = \dim(\mathbf{x}) + \dim(\tilde{\mathbf{x}})$ (rank condition)

Note that if $\text{rank} < \dim(\mathbf{x}) + \dim(\tilde{\mathbf{x}})$ the model is under-identified

Note that if $\text{rank} > \dim(\mathbf{x}) + \dim(\tilde{\mathbf{x}})$ the model is over-identified

indirect least square

The rank condition means that $\mathbb{E}[\mathbf{Z}^{*\top} \mathbf{X}]$ can be inverted, and

$$\beta = \mathbb{E}[\mathbf{Z}^{*\top} \mathbf{X}]^{-1} \mathbb{E}[\mathbf{Z}^{*\top} \mathbf{Y}]$$

and the empirical version is

$$\hat{\beta} = [\mathbf{Z}^{*\top} \mathbf{X}]^{-1} \mathbf{Z}^{*\top} \mathbf{y}$$

Causal Trees

More generally, we can have some over-identified model

We need to find a matrix A such that $A\mathbb{E}[\mathbf{Z}^{\star\top}\mathbf{X}]$ can be inverted, and then

$$\beta = \mathbb{E}[A\mathbf{Z}^{\star\top}\mathbf{X}]^{-1}\mathbb{E}[A\mathbf{Z}^{\star\top}Y]$$

and the empirical version is

$$\hat{\beta}_A = [A\mathbf{Z}^{\star\top}\mathbf{X}]^{-1}A\mathbf{Z}^{\star\top}\mathbf{y}$$

Then, for any A such that this estimator exists $\hat{\beta}_A$ is convergent and asymptotically Gaussian.

Causal Trees

There is an optimal matrix A^* such that the asymptotic variable of $\hat{\beta}_{A^*}$ is minimal, and it is

$$A^* = \mathbb{E}[\mathbf{X}^\top \mathbf{Z}^*]^{-1} \mathbb{E}[\mathbf{Z}^{*\top} \mathbf{Z}^*]$$

Set $A_n = \mathbf{X}^\top \mathbf{Z}^* (\mathbf{Z}^{*\top} \mathbf{Z}^*)^{-1}$, then $A_n \rightarrow A^*$, and then, the empirical estimator is the so-called double-least-square estimate

$$\hat{\beta} = [\mathbf{X}^\top \mathbf{Z}^* (\mathbf{Z}^{*\top} \mathbf{Z}^*)^{-1} \mathbf{Z}^{*\top} \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{Z}^* (\mathbf{Z}^{*\top} \mathbf{Z}^*)^{-1} \mathbf{Z}^{*\top} \mathbf{y}$$

or

$$\hat{\beta} = [\Pi_{\mathbf{Z}^*} \mathbf{X}^\top \Pi_{\mathbf{Z}^*} \mathbf{X}]^{-1} \Pi_{\mathbf{Z}^*} \mathbf{X}^\top \mathbf{y} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}} \mathbf{y}$$

which is called double-least-square estimate since

- we consider the regression of \mathbf{x} on \mathbf{z}^*
- we consider the regression of y on $\tilde{\mathbf{x}}$

Sometimes, instruments are poor proxys for endogeneous variables, and are called weak instruments.

Appendix: Regression Trees

One can use the **within variance** (within a leaf)

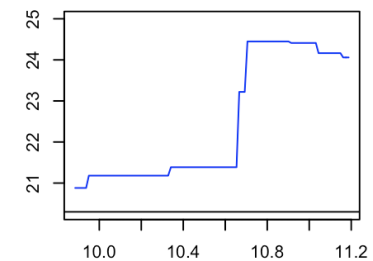
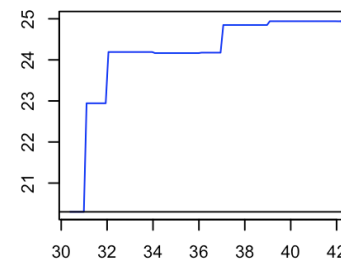
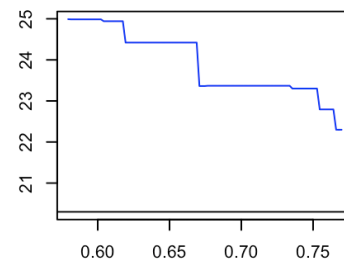
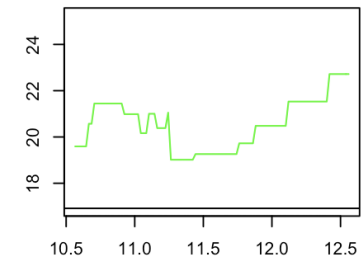
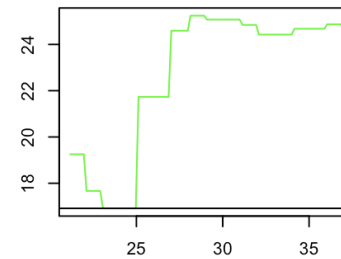
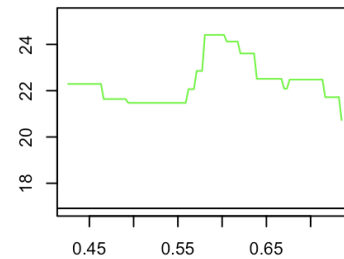
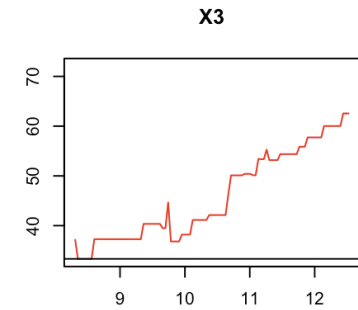
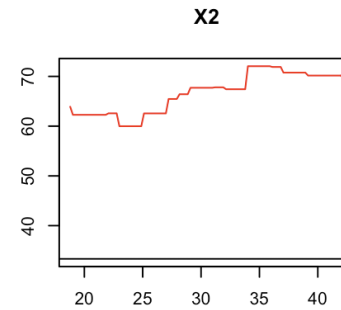
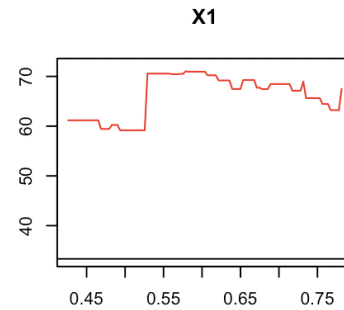
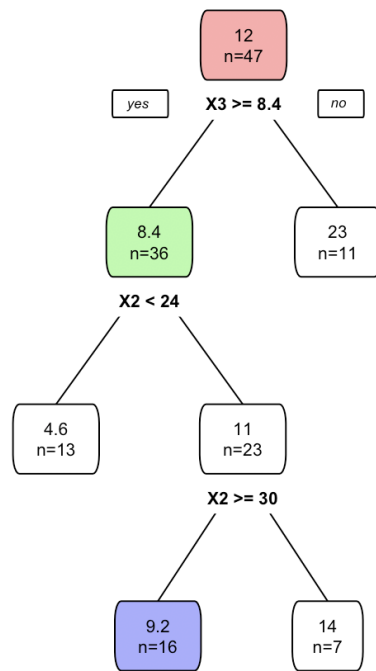
- within a leaf the total sum of squares is $\sum_i (y_i - \bar{y})^2$
- if we split in two parts, $\sum_{i:L} (y_i - \bar{y}_L)^2 + \sum_{i:R} (y_i - \bar{y}_R)^2$

The split is chosen to maximize $\sum_i (y_i - \bar{y})^2 - \sum_{i:L} (y_i - \bar{y}_L)^2 - \sum_{i:R} (y_i - \bar{y}_R)^2$

(one can use Student t -test for pruning)

See the application on the **Chicago** dataset, with 3 explanatory variables

Appendix: Regression Trees



Appendix: Regression Trees

One can use the **entropy** (that can be related to **Kullback-Leibler** distance)

Entropy

For a random variable X the entropy is $H(X) = -\mathbb{E}_X [\log f(X)]$

Natural extensions are the joint entropy $H(X, Y) = -\mathbb{E}_{X, Y} [\log f(X, Y)]$

and the conditional entropy $H(Y|X) = -\mathbb{E}_{X, Y} [\log f(Y|X)]$

Since $H(X) \neq H(Y)$, $H(Y|X) \neq H(X|Y)$.

Mutual Information

Mutual information of (X, Y) is $I(X, Y) = \mathbb{E}_{X, Y} \left[\log \frac{f(X, Y)}{f(X)f(Y)} \right]$

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$$

Appendix: Regression Trees

$I(X, Y) \geq 0$, and $I(X, Y) = 0$ if and only if $X \perp\!\!\!\perp Y$

thus $H(X) \leq H(X|Y)$ and $H(X) = H(X|Y)$ if and only if $X \perp\!\!\!\perp Y$

Kullback-Leibler

For two distributions f, g $KL(f\|g) = \mathbb{E}_f \left[\log \frac{f(X)}{g(X)} \right] = \int \log \frac{f(x)}{g(x)} f(x) dx$

Hence, $I(X, Y) = KL(f\|f^\perp)$

Thus, Kullback-Leibler divergence can be called **relative entropy**.

Appendix: Regression Trees

For a Gaussian vector $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, the joint entropy is

$$h(\mathbf{Z}) = \frac{1}{2} \log [(2\pi)^d |\boldsymbol{\Sigma}|]$$

If $\mathbf{Z}^* \in \operatorname{argmax}\{H(\mathbf{Z})\}$ s.t. $\operatorname{Var}[\mathbf{Z}] = \boldsymbol{\Sigma}$, then $\mathbf{Z}^* \sim \mathcal{N}(\mathbb{E}[\mathbf{Z}^*], \boldsymbol{\Sigma})$

Cross entropy

For distributions f, g , $CE(f|g) = -\mathbb{E}_f [\log g(X)] = -\int \log[g(x)]f(x)dx$

$$CE(f|g) = H(f) + KL(f||g)$$