

Segment 3: Causal Inference with Regression (in randomized studies)

Section 04: Regression for Causal Inference

Regression for Causal Inference: What's Different?

- ▶ Most often motivated as a *prediction* problem
 - ▶ Model how values of Y differ across values of X
 - ▶ Tailored to comparisons *between* units

$$\Rightarrow Y_i = \beta_0 + \beta_1 Z_i + \beta_2 X_i + \dots + \epsilon_i$$

- ▶ But **causal inference** is all about comparing what would happen if different treatments were applied to the *same units*

$$\Rightarrow \tau_{ATE} = E[Y_i^t - Y_i^c]$$

- ▶ What is the average difference between units with $Z = 0$ and $Z = 1$?
- ▶ What would be the difference if units with $Z = 0$ had actually had $Z = 1$ /?

Causal Inference Doesn't *Need* Unfamiliar Statistics!

Especially with randomization....

- ▶ Can rely on familiar inferential tools, like regression
- ▶ Main difference will be formalizing the underlying structure/assumptions to assess causal validity
 - ▶ What are the units?
 - ▶ What is the treatment?
 - ▶ What are the potential outcomes?
 - ▶ What is the *assignment mechanism*?
 - ▶ What is the role of covariates (blocking)?
- ▶ Much of the above is fairly trivial for randomized studies
- ▶ Eventually, we will consider non-randomized (observational) studies where the above are more critical

Regression Models in Randomized Experiments

- ▶ Common for the analysis of both experimental and observational data
- ▶ Four key features
 1. Models for the *observed* outcomes
 2. Model for the conditional mean only (not the full distribution)
 3. Average Treatment Effect will be a parameter of statistical model
 4. Whether the model accurately describes the conditional mean is immaterial for the large-sample unbiasedness of estimators for the average causal effect
 - ▶ I.e., Randomization guarantees unbiased estimation of the ATE even if the regression model is the “wrong” specification of the conditional mean of Y

Units, Covariates, Treatments, Outcomes

response \rightarrow

$$Y_i = \alpha + \tau Z_i + \mathbf{X}_i \beta + \epsilon_i$$

units, $i = 1, \dots, n$

$\epsilon \in \{0, 1\}$ "treatment"

$\underline{\mathbf{X}}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})$
"pre-treatment"

General guidance: Think of a single "treatment" as distinct from \mathbf{X}

Observed Y vs. Y^z

$$Y_i^{\text{obs}} = \alpha + \tau Z_i + \mathbf{X}_i \beta + \epsilon_i$$

$$E[Y | Z=1] = \alpha + \tau + \mathbf{X} \beta = E[Y^1 | Z=1]$$

$$E[Y | Z=0] = \alpha + \mathbf{X} \beta = E[Y^0 | Z=0]$$

$$E[Y | Z=1] - E[Y | Z=0] = \tau$$

Causal Inference as Prediction

- ▶ Special case of prediction
 - ▶ Predict what would have happened to the same unit(s) under different treatments
 - ▶ Vs. Predict what would happen for a unit with a particular set of features, treatments
- ▶ That is, try to predict the *other* potential outcome
 - ▶ A regression model may help to generate a "close substitute"

Table: Observed Data from the Hypothetical Dietary Experiment

Unit, i	Treatment Z_i	Potential Outcome, Y_i^c	Potential Outcome, Y_i^t	Observed Outcome, Y_i
Audrey	0	140	?	140
Anna	0	140	?	140
Bob	0	150	?	150
Bill	0	150	?	150
Caitlin	1	?	155	155
Cara	1	?	155	155
Dave	1	?	160	160
Doug	1	?	160	160

Handwritten notes: Arrows point from the observed outcomes of the control group (Audrey, Anna, Bob, Bill) to the potential outcomes of the treatment group (Caitlin, Cara, Dave, Doug). A label "predict" is written next to the arrows.

Regression with No Covariates

and randomized treatment assignment

$$Y_i^{obs} = \alpha + \tau Z_i + \epsilon_i$$

- ▶ $\tau = E[Y' - Y^0]$
- ▶ $\hat{\tau}^{ols} = \bar{Y}_t^{obs} - \bar{Y}_c^{obs} = \hat{\tau}_{SATE} = \hat{\tau}_{PATE}$
- ▶ `stan_glm(y ~ z)`
- ▶ Posterior distribution of τ dictates estimation of τ
 - ▶ Posterior median, sd, 95% interval, etc.

$$\begin{aligned}
 E[Y|Z=1] &= \alpha + \tau = E[Y'|Z=1] = E[Y'] \\
 - E[Y|Z=0] &= \alpha = E[Y^0|Z=0] = E[Y^0] \\
 \hline
 \tau &= E[Y'] - E[Y^0]
 \end{aligned}$$

Handwritten notes: An arrow points from the term $E[Y']$ in the first equation to the term $E[Y']$ in the final equation. A label "randomized" is written next to the first equation.

Different Ways to Balance Covariates

Randomization is motivated in large part for its production of *covariate balance* across treatment groups

- ▶ Some methods can improve balance *by design*
 - ▶ Stratified randomized experiments */blocking*
 - ▶ Paired randomized experiments

(but balance is not *guaranteed* for any individual randomization)

- ▶ Could also attempt to correct covariate imbalance by *analysis*
 - ▶ Separate analysis within subgroups
 - ▶ Regression adjustment

Regression with Covariates

(still with complete randomization)

$$Y_i^{obs} = \alpha + \tau Z_i + X_i \beta + \epsilon_i$$

- ▶ α, β regarded as “nuisance” parameters that do not get a causal interpretation
- ▶ We *do not* assume that the regression function is correctly specified
 - ▶ I.e., that the Y_i^{obs} is actually linear in X_i and Z_i
- ▶ $\hat{\tau}^{ols}$ is unbiased for τ_{SATE} and τ_{PATE}
- ▶ `stan_glm(y ~ z + x1 + x2 + ... + xp)`
- ▶ Posterior distribution of τ dictates estimation of τ
 - ▶ Posterior median, sd, 95% interval, etc.

Benefits of Adjusting for Pre-Treatment Covariates

Under randomization

$$Y_i^{obs} = \alpha + \tau Z_i + X_i\beta + \epsilon_i$$

- ▶ Including \mathbf{X} does not change interpretation of τ
 - ▶ Not generally true under *nonrandomized* treatment
- ▶ Can adjust for *random* imbalances between groups
 - ▶ “Unlucky assignments”
- ▶ Can adjust for *systematic* imbalances between groups
 - ▶ Like blocking factors
- ▶ Ultimately bring the estimate $\hat{\tau}$ closer to the truth in any given sample/randomization
 - ▶ Less variability around truth \Rightarrow smaller standard error of a given estimate

Benefits of Adjustment: Example

Effect of showing children an educational TV show

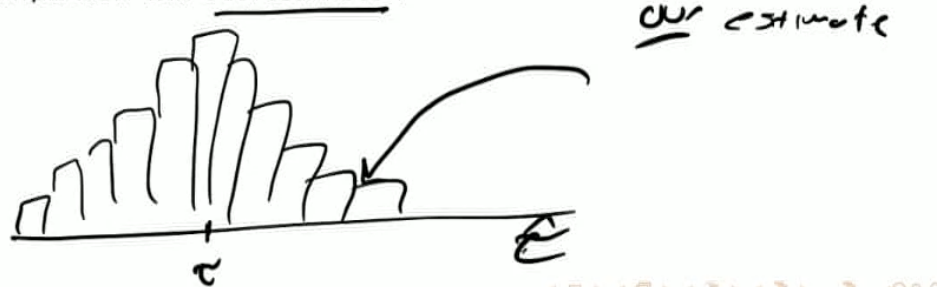
$n = 192$ elementary school classes

Z : New vs. standard educational TV program

Y : Reading ability measured at end of school year

X : Pre-test reading ability

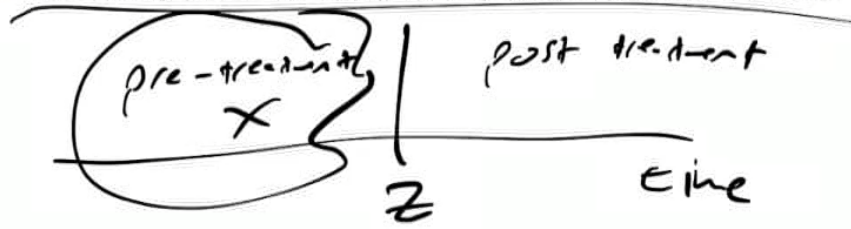
- ▶ Notice that pre-test reading is positively correlated with Y
- ▶ Notice that the average pre-test reading ability is Δ_x higher in the treatment group
 - ▶ “Unlucky” randomization
- ▶ (X positively correlated with Y) + ($\Delta_x \neq 0$) \Rightarrow unadjusted comparison will overestimate τ



Benefits of Adjusting for Pre-Treatment Covariates

Under randomization

- ▶ Same idea as the educational TV example applies to any pre-treatment covariate that may predict Y
- ▶ \Rightarrow Always the chance to improve precision and bring $\hat{\tau}$ "closer to the truth" by adjusting for *predictors* of Y
- ▶ In *practice* we can't measure (or adjust) for *everything*
 - ▶ May need to prioritize which variable to include for adjustment
 - ▶ E.g., based on theory
 - ▶ E.g., based on some measure of imbalance and correlation with Y
- ▶ Reasoning does not apply to any *posttreatment* variable!



Special Case: Gain Scores

Setting: One covariate is a "pre-treatment" or "pre-test" or "baseline" version of Y

Can create a "gain score", $g_i = Y_i - X_i$

$$g_i = \alpha + \tau Z_i + \varepsilon_i \rightarrow \hat{\tau} = \bar{g}^E - \bar{g}^C$$

$$Y_i - X_i = \alpha + \tau Z_i + \varepsilon_i \rightarrow Y_i = \alpha + \tau Z_i + X_i + \varepsilon_i$$

\uparrow
 $\beta \equiv 1$
 unnecessary
 assumption

Varying Treatment Effects and Interactions

General Consideration: The treatment effect may vary across units with different X

$$\tau_{CATE|X} = E[Y' - Y^0 | X=x]$$

$$Y_i^{obs} = \alpha + \theta Z_i + X_i \beta + Z_i X_i \gamma + \epsilon_i$$

Treatment effect:

$$E[Y | Z=1] = \alpha + \theta + X\beta + X\gamma$$

$$E[Y | Z=0] = \alpha + X\beta$$

$$E[Y | Z=1] - E[Y | Z=0] = \theta + X\gamma$$

$$\hat{\tau}_{CATE} = \sum_{i=1}^n (\theta + \gamma X_i)$$

"treatment effect"
= function of X

Tau hat should be multiplied by a factor of $1/n$

Unbiasedness Under Randomization

Even when the model is "wrong"

$$Y_i = \alpha + \tau Z_i + X_i \beta + Z_i X_i \gamma + \epsilon_i$$

$$E[Y | Z=1] = \underbrace{\alpha + \tau + (\beta + \gamma)X}_{X\beta_c} = E[Y' | Z=1]$$

$$E[Y | Z=0] = \underbrace{\alpha + \beta X}_{X\beta_c} = E[Y^0 | Z=0]$$

$$E[Y^{mis} | Z=1] = E[Y^c | Z=1] \stackrel{\text{randomization}}{=} E[Y^c | Z=0] = X\hat{\beta}_c$$

$$E[Y^{mis} | Z=0] = E[Y^c | Z=0] = E[Y^c | Z=1] = X\hat{\beta}_c$$

$$\hat{\tau} = \frac{1}{n} \sum_{i=1}^n z_i (\underbrace{y_i^{\text{obs}}}_{y^e} - \underbrace{x_i \hat{\beta}_c}_{y^c_{\text{miss}}}) + (1 - z_i) (\underbrace{x_i \hat{\beta}_c}_{y^e_{\text{miss}}} - \underbrace{y_i^{\text{obs}}}_{y^c})$$

$$\hat{\tau}_c = \frac{1}{n_c} \sum_{i: z_i=0} (x_i \hat{\beta}_c - y_i^{\text{obs}})$$

$$\begin{aligned} \hat{\tau}_c &= \bar{x}_c \hat{\beta}_c - \bar{y}_c^{\text{obs}} (-\bar{x}_c \hat{\beta}_c + \bar{y}_c^{\text{obs}}) \\ &= \bar{y}_c^{\text{obs}} - \bar{y}_c^{\text{obs}} + \hat{\beta}_c (\underbrace{\bar{x}_c - \bar{x}_c}_{\text{"balance"}}) \\ &= 0 \text{ on average} \end{aligned}$$

$$\hat{\tau}_t = \frac{1}{n_t} \sum_{i: z_i=1} (y_i^{\text{obs}} - x_i \hat{\beta}_c)$$

$$= \bar{y}_t^{\text{obs}} - \bar{x}_t \hat{\beta}_c + \bar{x}_c \hat{\beta}_c - \bar{y}_c^{\text{obs}}$$

$$= \bar{y}_t^{\text{obs}} - \bar{y}_c^{\text{obs}} + \hat{\beta}_c (\underbrace{\bar{x}_c - \bar{x}_t}_{=0 \text{ on average}})$$

$$\hat{\tau} = \frac{n_c}{n} \hat{\tau}_c + \frac{n_t}{n} \hat{\tau}_t$$

Summary: Regression for Completely Randomized Experiments

- ▶ Easy to incorporate covariates
- ▶ Unbiased point estimates without covariates
- ▶ Consistent point estimates for the ATE with covariates, *even when model is not true*
 - ▶ Point estimates biased in finite samples
 - ▶ Works in large samples
- ▶ Model for *observed outcomes*
- ▶ Causal estimand (ATE) is a parameter (or parameters) in the statistical model