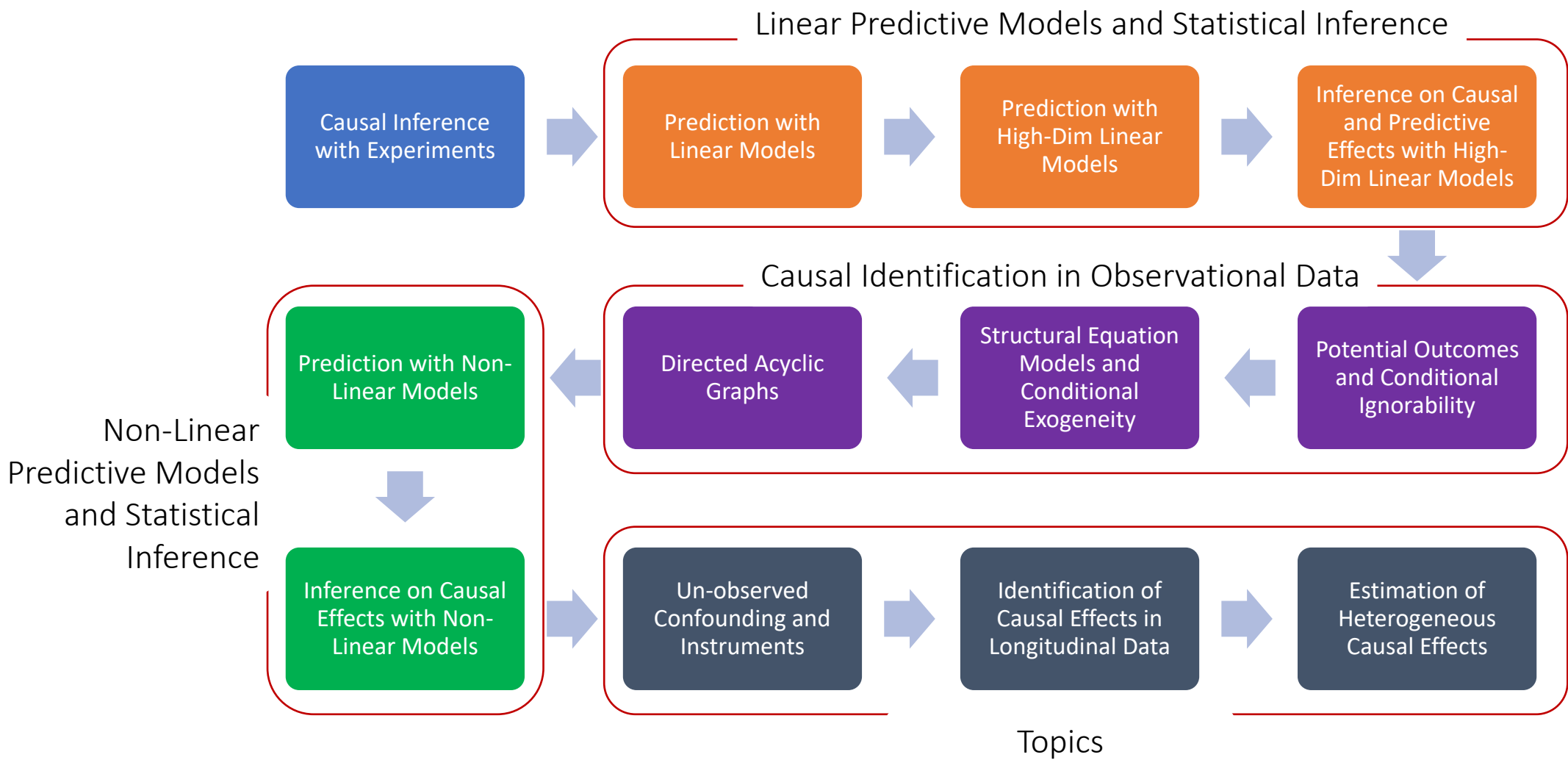


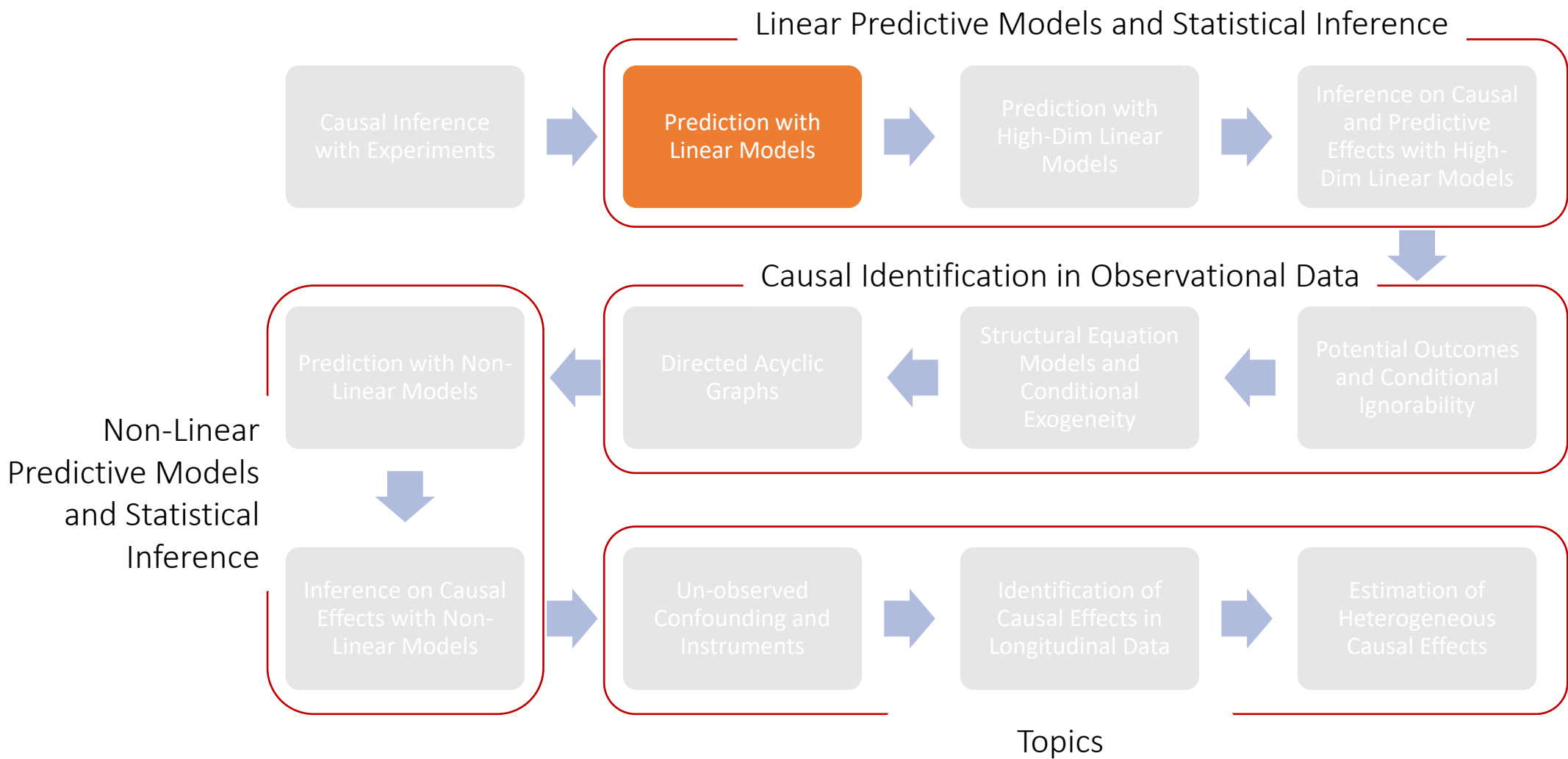
MS&E 228:

Analyzing Experiments with Linear Regression

Vasilis Syrgkanis

MS&E, Stanford





Even if the relationship between outcome Y and covariates X is non-linear, we can always write:

$$Y = \beta'X + \epsilon, \quad E[\epsilon X] = 0$$

The function $\beta'X$ is the Best Linear Predictor (BLP) or equivalently the best linear approximation to the Conditional Expectation Function (CEF) $E[Y|X]$



Even if the relationship between outcome Y and covariates X is non-linear, we can always write:

$$Y_i = \hat{\beta}' X_i + \hat{\epsilon}_i, \quad E_n[\hat{\epsilon} X] = 0$$

The function $\hat{\beta}' X$ is the Best Linear Predictor in sample and $\hat{\beta}$ are the sample regression coefficients



You should expect OLS to produce accurate predictions in the worst-case if the number of variables is small compared to number of samples.



Its predictions converge to the predictions of the BLP in the population



Almost always measure predictive performance of your estimated model on a held-out sample

Predictive effect β_1 of *target variable* is the coefficient in a *simple one variable regression*



$$\left(\begin{array}{c} \text{part of outcome} \\ \text{(un-explained by other)} \end{array} \right) \sim \left(\begin{array}{c} \text{part of target} \\ \text{(un-explained by other)} \end{array} \right)$$

Coefficient of D in $\text{OLS}(y \sim D, W)$ is mathematically equivalent in samples to

$$y_{\text{res}} = y - \text{OLS}(y \sim W).\text{predict}(W)$$

$$D_{\text{res}} = D - \text{OLS}(D \sim W).\text{predict}(W)$$



Coefficient of D_{res} in $\text{OLS}(y_{\text{res}} \sim D_{\text{res}})$

If we want an interval that roughly contains the predictive effect with probability α , we can use

$$CI(\alpha) := \left[\hat{\beta}_1 - z_{1-\frac{\alpha}{2}} \hat{\sigma}_n, \hat{\beta}_1 + z_{1-\frac{\alpha}{2}} \hat{\sigma}_n \right]$$

$$\hat{\sigma}_n := \frac{1}{\sqrt{n}} \sqrt{\frac{E_n[\hat{\epsilon}^2 \check{D}^2]}{E_n[\check{D}^2]^2}}$$



e.g. for 95% confidence interval, $z_{1-\frac{\alpha}{2}} \approx 1.96$

Revisit Covariate Adjustment for Effect Inference in Experiments

Co-variates for Precision

- Even if we are only interested on ATE covariates can be valuable for precision
- Suppose variance of y is large but can be explained largely by W
- Then we can use W to remove all the explained variation from y
- Then perform our ATE analysis on the remnant variation
- This is oftentimes performed in practice via ordinary linear regression of y on the vector $(1, D, W)$ (after centering W , i.e. $E[W] = 0$)

Is this consistent?

- Suppose that the conditional expectation function (CEF) of the outcome is indeed linear, with $(D, 1, W)$

$$E[Y | D, W] = D\alpha + \alpha_0 + W'\beta$$

- Then note that

$$\begin{aligned} E[Y(0)] &= E[E[Y|D = 0, W]] = \alpha_0 \\ E[Y(1)] &= E[E[Y|D = 1, W]] = \alpha + \alpha_0 \end{aligned}$$

- Baseline outcome is coefficient associated with the intercept 1
- Average effect is coefficient associated with treatment D

Is it consistent if $E[Y | D, W]$ is not linear?

Simple Case

- Suppose we run linear regression of $Y \sim (D, 1)$
- Is the coefficient associated with D , converging to the ATE?
- Let's look at the population BLP:

$$Y = a_0 + aD + \epsilon, \quad E[\epsilon(D, 1)] = 0$$

- Let's examine the normal equations:

$$E[\epsilon D] = 0 \Rightarrow E[\epsilon | D = 1] = 0 \Rightarrow E[Y | D = 1] = a_0 + a$$

$$E[\epsilon] = 0 \Rightarrow E[\epsilon | D = 0] = 0 \Rightarrow E[Y | D = 0] = a_0$$

- From these two we uncover that:

$$a = E[Y | D = 1] - E[Y | D = 0] = E[Y(1) - Y(0)] = \delta$$

What if we add covariates?

- By the BLP decomposition of Y using $(D, 1, W)$

$$Y = D\alpha + \alpha_0 + \beta'W + \epsilon, \quad E[\epsilon(D; 1; W)] = 0$$

- Note that the quantity:

$$U = \beta'W + \epsilon$$

- Also satisfies

$$\begin{aligned} E[U(D; 1)] &= \beta' E[W(D; 1)] + E[\epsilon(D; 1)] \\ &= \beta' E[W] E[(D; 1)] = 0 \end{aligned}$$

= 0 by orthogonality of ϵ

By independence of W
and D

Since $E[W] = 0$ (de-meanned W)

Is this consistent? Beyond Linear CEF

- We can write

$$Y = D\alpha + \alpha_0 + U, \quad E[U(D; 1)] = 0$$

- where $U = \beta'W + \epsilon$
- Thus α, α_0 solve the Normal Equations of Y on $D, 1$
- Thus α, α_0 are the BLP of Y using $(D, 1)$
- The parameters a, a_0 are the same as the parameters in the population BLP if we did not use covariates!
- We already saw that in that case a is the ATE!

Is this consistent? Beyond Linear CEF

- The crucial property we used was that $U = \beta'W + \epsilon$ is uncorrelated with $(D, 1)$
- This holds whenever W is un-correlated with $(D, 1)$
- So, if we run OLS with any set of controls W that are un-correlated with $(D, 1)$, then the population parameters a, a_0 associated with $(D, 1)$, remain un-changed!



The coefficient associated with treatment D in OLS with co-variate adjustment is always consistent for the treatment effect, when run on data from a randomized experiment, as-long-as covariates are de-meanned. The true relationship of outcome with covariates does not need to be linear.

Analysis of Variance (ANOVA)

Quality of In-Sample Estimate

- What if we run a Linear Regression with n samples
- How accurate is the parameter \hat{a} in the in-sample BLP?
- Does the accuracy change or improve if we use W or not?

Variance of Estimate

- The parameter α is the “predictive effect” of D
- Using the characterization of the estimate of the predictive effect:

$$\sqrt{n}(\hat{\alpha} - \alpha) \overset{a}{\sim} N(0, V_{\alpha})$$

$$V_{\alpha} = \frac{E[\epsilon^2 \tilde{D}^2]}{E[\tilde{D}^2]^2}$$

- ϵ is population residual outcome:

$$\epsilon = y - D\alpha - \alpha_0 - W'\beta \quad E[\epsilon(D, 1, W)] = 0$$

- \tilde{D} is residual treatment (removing whatever is linearly predictable from $(1, W)$)

$$\tilde{D} = D - E[D]$$

Heteroskedasticity Robust Variance

- Variance formula

$$V_{\alpha} = \frac{E[\epsilon^2 \tilde{D}^2]}{E[\tilde{D}^2]^2}$$

- is valid even when the linear CEF assumption is violated
- Important to note that this formula is known as the “heteroskedasticity robust variance formula” (HCO)
- Many software packages make the simplification that the residual ϵ is independent of D, W , leading to $V_{\alpha} = E[\epsilon^2]/E[\tilde{D}^2]$. This is incorrect in most cases!

Variance without adjustment

- If we don't use W , then the formula remains the same, but the residual changes (note: population parameters α, α_0 are the same)
$$\bar{\epsilon} = Y - D\alpha - \alpha_0 = W'\beta + \epsilon$$
- The residual treatment \tilde{D} is the same, since W is independent of D
- Comparing the two variances, only the numerators change:

$$\begin{aligned} \underbrace{E[\bar{\epsilon}^2 \tilde{D}^2]}_{\text{Numerator when not using } W} &= E[(W'\beta + \epsilon)^2 \tilde{D}^2] \\ &= E[(W'\beta)^2 \tilde{D}^2] + E[\epsilon^2 \tilde{D}^2] + 2E[\beta' W \epsilon \tilde{D}^2] \\ &= E[(W'\beta)^2 \tilde{D}^2] + \underbrace{E[\epsilon^2 \tilde{D}^2]}_{\text{Numerator when using } W} \end{aligned}$$

If this term is zero, then we know that using W can only decrease variance of estimate

When can we claim that
interaction term is zero?

$$E[\epsilon W \tilde{D}^2] = 0$$

One Case: Strong Assumption

- Suppose that CEF is linear, i.e.

$$E[Y \mid D, W] = a_0 + a D + \beta' W$$

- Then we have that:

$$Y = a_0 + a D + \beta' W + \epsilon, \quad E[\epsilon \mid D, W] = 0$$

- Then we can verify that the cross-term is zero:

$$E[\epsilon W \tilde{D}^2] = E[E[\epsilon \mid D, W] W \tilde{D}^2] = 0$$

- Hence, when CEF is linear then variance of OLS estimate with extra co-variates (adjusted) is weakly smaller than two-means estimate (unadjusted)

What if CEF is Non-Linear?

Example: Heterogeneous Effects

- Assume $E[W] = 0$ and suppose that:

$$E[Y | D, X] = D \underbrace{\alpha}_{\text{ATE}} + \alpha_0 + D \underbrace{W' \gamma}_{\text{effect modifier}} + W' \beta$$

- Equivalently:

$$Y = Da + a_0 + DW' \gamma + W' \beta + v, \quad E[v | D, W] = 0$$

- What OLS $Y \sim D, 1, W$ is estimating is coefficients $\tilde{\alpha}, \tilde{\alpha}_0, \tilde{\beta}$, the solution to:

$$E \left[(Y - D\tilde{\alpha} - \tilde{\alpha}_0 - W'\tilde{\beta}) \begin{pmatrix} D \\ 1 \\ W \end{pmatrix} \right] = 0$$

- Since W is un-correlated with $D, 1$ by our previous analysis, $\tilde{\alpha}, \tilde{\alpha}_0$ are the same as a, a_0 , i.e. we are still recovering the correct ATE.

What if CEF is Non-Linear?

Example: Heterogeneous Effects

- Assume $E[W] = 0$ and suppose that:

$$Y = Da + a_0 + DW'\gamma + W'\beta + v, \quad E[v | D, W] = 0$$

- What OLS $Y \sim D, 1, W$ is estimating is coefficients $\tilde{a}, \tilde{a}_0, \tilde{\beta}$, the solution to:

$$E \left[(Y - D\tilde{a} - \tilde{a}_0 - W'\tilde{\beta}) \begin{pmatrix} D \\ 1 \\ W \end{pmatrix} \right] = 0$$

- But $\tilde{\beta} = \beta + E[D]\gamma$ by examining the third normal equation:

$$E[(DW'\gamma + W'\beta - W'\tilde{\beta})W] = 0 \Rightarrow E[W'](\gamma E[D] + \beta - \tilde{\beta}) = 0$$

- So, the residual of OLS of $Y \sim D, 1, W$ is:

$$\epsilon = Y - Da - a_0 - W'\tilde{\beta} = (D - E[D])W'\gamma + v, \quad E[v | D, W] = 0$$

- And the cross-term in the variance analysis becomes:

$$E[W\epsilon\tilde{D}^2] = E[\tilde{D}^3]E[WW']\gamma \neq 0$$

OLS with Interactive Terms

- What if instead we run OLS including interaction terms [Lin'13]

$$Y \sim D, 1, W, DW$$

- Since $E[W] = 0$ and $W \perp\!\!\!\perp D$, we have (DW, W) are un-correlated with D

$$E[DDW] = E[D^2]E[W] = 0$$

- Thus, by the analysis we did before, in the absence of any model assumptions, the coefficient of D and of the intercept, recover the true ATE a and the true mean baseline outcome a_0

Let's try it out!

OLS with Interactive Terms

- These interactive terms enforce the residual ϵ of OLS to satisfy a stronger orthogonality property

$$E \left[\epsilon \begin{pmatrix} D \\ 1 \\ W \\ DW \end{pmatrix} \right] = 0$$

- We can conclude the stronger un-correlatedness properties:

$$E[\epsilon DW] = 0 \Rightarrow E[\epsilon W \mid D = 1] = 0$$

$$E[\epsilon W] = 0 \Rightarrow E[\epsilon W \mid D = 0] = 0$$

- From which we can argue that the cross-term in the variance is zero

$$E[W \epsilon \tilde{D}^2] = E[E[W \epsilon \mid D] \tilde{D}^2] = 0$$

Even if you only care about ATE, if you have p covariates and $p \ll n$ run OLS with interactive terms (after de-meaning covariates)!

Guaranteed improved precision