

# Selección de Modelos y Regularización

## Ciencia de Datos y Econometría Aplicada

Ignacio Sarmiento-Barbieri

Universidad de los Andes

September 8, 2025

# Potential outcomes and treatment effects



# Potential outcomes and treatment effects



\$\$\$

\$43M

# Potential outcomes and treatment effects



\$\$\$

\$43M

# Potential outcomes and treatment effects



\$\$\$

\$43M

$Y(1)$



\$

\$700K = \$42.3M

$Y(0)$  = Treatment effect

counterfactual

# Potential outcomes and treatment effects



\$\$\$

\$43M

$Y(1)$

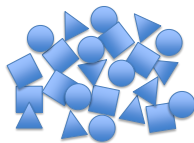


\$

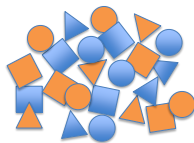
\$700K = \$42.3M

$Y(0)$  = Treatment effect

# Potential outcomes and treatment effects

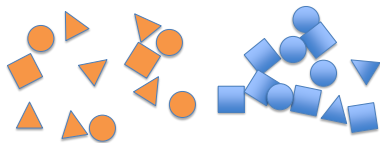


# Potential outcomes and treatment effects

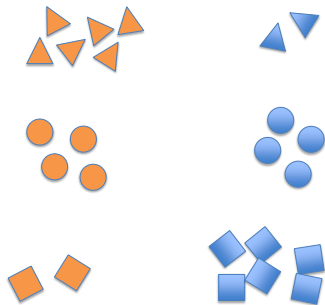




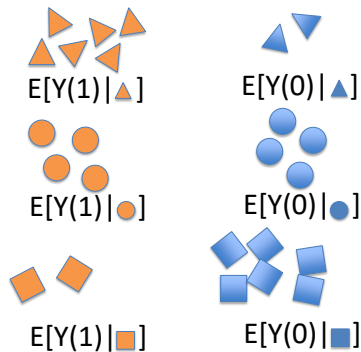
# Potential outcomes and treatment effects



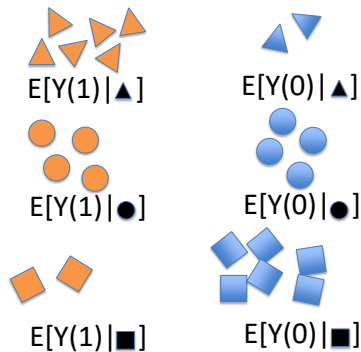
# Potential outcomes and treatment effects



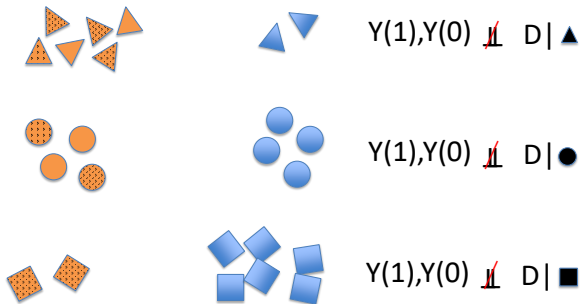
# Potential outcomes and treatment effects



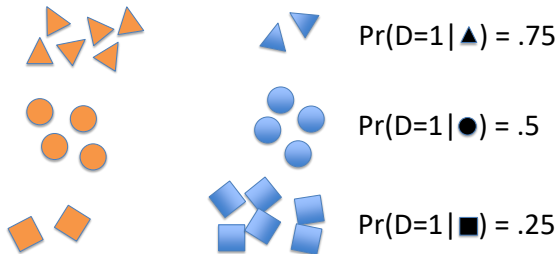
# Potential outcomes and treatment effects



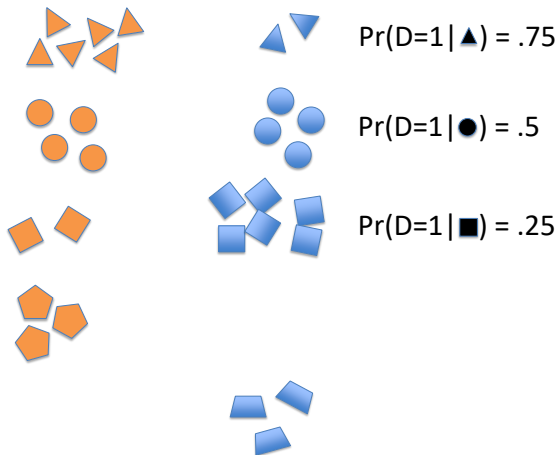
# Potential outcomes and treatment effects



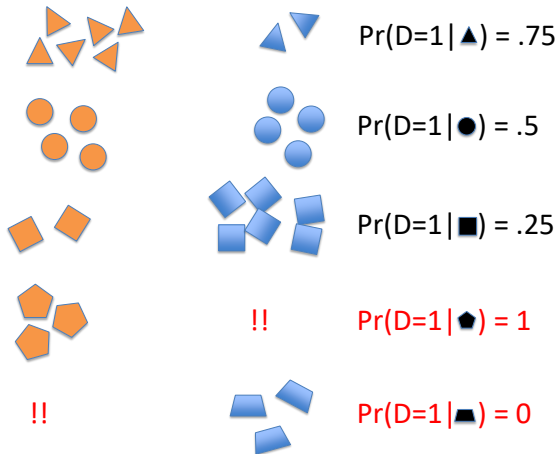
# Potential outcomes and treatment effects



# Potential outcomes and treatment effects



# Potential outcomes and treatment effects





# Basic causal inference summary

- ▶ Target :

$$ATE = E[Y_i(1) - Y_i(0)] = E[\tau_i]$$

- ▶ Key identifying assumption:

$$(Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i | X_i$$

- ▶ Estimation:

- ▶ Multiple linear regression (OLS)

$$Y_i = \beta_0 + \tau D_i + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \varepsilon$$

- ▶ Matching
- ▶ Propensity score methods
- ▶ Machine-assisted:
  - ▶ Post-Double Selection Lasso

# Model Selection When the Goal is Causal Inference

Let's start with the following model

$$y_i = \alpha + \beta D_i + g(X_i) + \zeta_i \quad (1)$$

were

- ▶  $D_i$  is the treatment/policy variable of interest,
- ▶  $X_i$  is a set controls
- ▶  $E[\zeta_i | D_i, X_i] = 0$

# Model Selection When the Goal is Causal Inference

- ▶ Traditional approach: researcher selects  $X_i$
- ▶ Problem: mistakes can occur.
- ▶ Same if they use an “automatic” model selection approach.
- ▶ It can leave out potentially important variables with small coefficients but non zero coefficients out

# Model Selection When the Goal is Causal Inference

- ▶ The omission of such variables then generally contaminates estimation and inference results based on the selected set of variables. (e.g. OVB)
- ▶ The validity of this approach is delicate because it relies on perfect model selection.
- ▶ Because model selection mistakes seem inevitable in realistic settings, it is important to develop inference procedures that are robust to such mistakes.
- ▶ Solution here: Lasso

# Model Selection When the Goal is Causal Inference

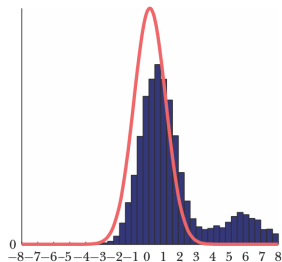
- ▶ Using Lasso is useful for prediction
- ▶ However, naively using Lasso to draw inferences about model parameters can be problematic.
- ▶ Part of the difficulty is that these procedures are designed for prediction, not for inference
- ▶ Leeb and Pötscher 2008 show that methods that tend to do a good job at prediction can lead to incorrect conclusions when inference is the main objective
- ▶ This observation suggests that more desirable inference properties may be obtained if one focuses on model selection over the predictive parts of the economic problem

# Inference with Selection among Many Controls

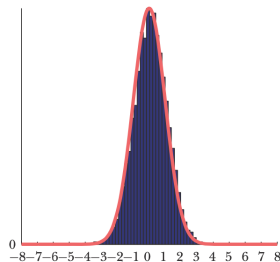
Figure 1

**The “Double Selection” Approach to Estimation and Inference versus a Naive Approach: A Simulation from Belloni, Chernozhukov, and Hansen (forthcoming)**  
(distributions of estimators from each approach)

A: A Naive Post-Model Selection Estimator



B: A Post-Double-Selection Estimator



Source: Belloni, Chernozhukov, and Hansen (forthcoming).

Notes: The left panel shows the sampling distribution of the estimator of  $\alpha$  based on the first naive procedure described in this section: applying LASSO to the equation  $y_i = d_i + x_i' \theta_j + r_{ji} + \zeta_i$  while forcing the treatment variable to remain in the model by excluding  $\alpha$  from the LASSO penalty. The right panel shows the sampling distribution of the “double selection” estimator (see text for details) as in Belloni, Chernozhukov, and Hansen (forthcoming). The distributions are given for centered and studentized quantities.