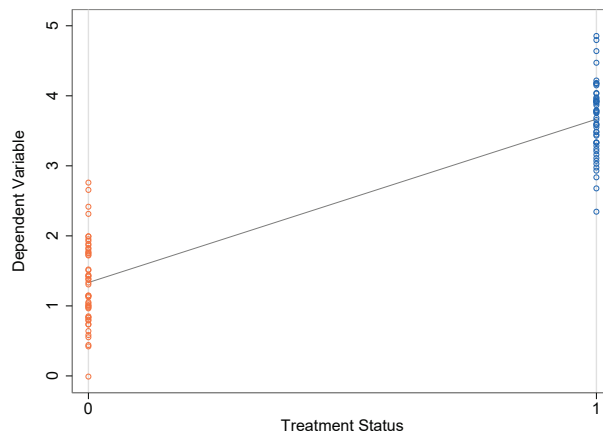




Basic Regression Equation:  $Y_i = \alpha + \beta D_i + \varepsilon_i$



## RCT Regression Specification with Controls

More typical regression specification:

$$Y_{1,i} = \alpha + \beta D_i + \delta X_{0,i} + \gamma Y_{0,i} + \kappa_{strata} + \varepsilon_i$$

We typically include these controls:

- Dummies for randomization strata ( $\kappa_{strata}$ )
- Selected baseline covariates that are not balanced across treatments\*
- Baseline covariates that predict the outcome
  - ▶ Baseline values of outcome variables are (sometimes) most important (ANCOVA)

We do not want to include:

- Controls that could be impacted by treatment (“bad controls” problem)

## “You Don’t Have to Take My Word For It”

The screenshot shows the website for the American Economic Journal: Applied Economics, Volume 11, Number 3, July 2019. The page features a search bar and a list of articles. The articles listed are:

- Front Matter** (pp. 1-4)
- Do 40-Year-Old Facts Still Matter? Long-Run Effects of Federal Oversight under the Voting Rights Act** by Desmond Ang (pp. 1-10)
- Medium- and Long-Term Educational Consequences of Alternative Conditional Cash Transfer Designs: Experimental Evidence from Colombia** by Felipe Barrera-Osorio, Leigh L. Linden and Juan E. Saverda (pp. 14-40)
- Long-Run Effects of Temporary Incentives on Medical Care Productivity** by Pablo A. Celhay, Paul J. Gertler, Paula Giovagnoli and Christel Vermeersch (pp. 51-127)
- Does Teacher Training Actually Work? Evidence from a Large-Scale Randomized Evaluation of a National Teacher Training Program** by Prashant Loyalka, Anna Popova, Guirong Li and Zhaolei Shi (pp. 139-140)
- Reducing Child Mortality in the Last Mile: Experimental Evidence on Community Health Promoters in Uganda** by Martina Björkman Nyqvist, Andrea Guarino, Jakob Svensson and David Yanagizawa-Drott (pp. 161-162)

## “You Don’t Have to Take My Word For It”

American Economic Journal: Applied Economics 2019, 11(3): 128–154  
<https://doi.org/10.1257/app.20170226>

### Does Teacher Training Actually Work? Evidence from a Large-Scale Randomized Evaluation of a National Teacher Training Program<sup>†</sup>

By PRASHANT LOYALKA, ANNA POPOVA, GUIRONG LI, AND ZHAOLEI SHI<sup>✉</sup>

*Despite massive investments in teacher professional development (PD) programs in developing countries, there is little evidence on their effectiveness. We present results of a large-scale, randomized evaluation of a national PD program in China in which teachers were randomized to receive PD; PD plus follow-up; PD plus evaluation of the content of PD content; or no PD. Precise estimates indicate PD and associated interventions failed to improve teacher and student outcomes after one year. A detailed analysis of the causal chain shows teachers find PD content to be overly theoretical, and PD delivery too rote and passive, to be useful. (JEL I21, I28, J24, J45, O15, P36)*

## “You Don’t Have to Take My Word For It”

We estimate the ATEs using the following ordinary least squares regression model:<sup>22</sup>

$$(1) \quad Y_{ij} = \alpha_0 + \alpha_1 D_j + X_{ij} \alpha + \tau_k + \varepsilon_{ij},$$

where  $Y_{ij}$  is the outcome of interest measured at endline for student  $i$  in school  $j$ ;  $D_j$  is one or more dummies indicating the treatment assignment of school  $j$ ;  $X_{ij}$  is a vector of baseline control variables; and  $\tau_k$  is a set of block fixed effects. In all specifications,  $X_{ij}$  includes the baseline value of the dependent variable whenever this is available. We also estimate treatment effects with an expanded set of baseline controls (we call these our “covariate-adjusted” regressions). For student-level outcomes, this expanded set of controls includes student age, student gender, parent educational attainment, a household asset index, class size, teacher gender, teacher age, teacher experience, teacher education level, a teacher certification dummy, a teacher major in math dummy, and teacher rank. For outcomes measured at the teacher level, student controls are omitted.

## We Want to Include Controls that Predict $T_i$ or $Y_i$

Controls should be orthogonal to treatment status (because we randomized):

- In practice, a control might be correlated with treatment
  - ▶ In small(ish) samples, we may see some differences in observables/covariates between treatment and comparison groups (especially as the number of covariates increases)
- Can be important to show estimated impacts are “robust” to inclusion of covariates
- Adding in one imbalanced covariate can undo randomization (remember, we are now regressing residual of  $Y$  on residual of  $D$ )
  - ▶ This is (one of the reasons) why **stratification** is desirable

Baseline covariates are orthogonal to treatment, so they should not impact coefficient

- Controls help if they explain residual variation in  $Y$

## What Is Machine Learning?

A set of extensions to the standard econometric toolkit (read: “OLS”) aimed at improving predictive accuracy, particularly when data sets are sparse (many variables, most are garbage)

- Subset selection, shrinkage methods (ridge regression, lasso) for covariate selection
- Regression trees, random forests, causal forests to identify treatment effect heterogeneity

Machine learning introduces new tools, relabels existing tools

- **training data/sample/examples**: your (baseline) data
- **features**: covariates

Main focus is on predicting  $Y$ , not testing hypotheses about causal impact of  $\beta$

⇒ ML “results” about  $\beta$  may not be robust or causally identified

## Can We Improve on OLS?

A standard linear model is not (always) the best way to predict  $Y$ :

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$

Can we improve on OLS?

- When  $p > N$ , OLS is not feasible
- When  $p$  is large relative to  $N$ , model may be prone to over-fitting
- OLS explains both structural and spurious relationships in data

Extensions to OLS identify “strongest” predictors of  $Y$

- Strength of correlation vs. (out-of-sample) robustness

Assumption: exact or approximate **sparsity**

## Best Subset Selection

A **best subset selection** algorithm:

- For each  $k = 1, 2, \dots, p$ 
  - ▶ Fit all models containing exactly  $k$  covariates
  - ▶ Identify the “best” in terms of  $R^2$
- Choose the **best subset** based on cross-validation, adjusted  $R^2$ , etc.
  - ▶ Need to address the fact that  $R^2$  always increases with  $k$

When  $p$  is large, best subset selection is not feasible

- Forward and backward subset selection may work poorly when covariates are correlated

## Best Subset Selection

In OLS, we seek to minimize:

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

Best subset selection can be expressed as: choose  $\beta$  to minimize

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \text{ subject to } \sum_{j=1}^p I(\beta_j \neq 0) \leq s$$

where  $s$  is the number of regressors/predictors/features/covariates

- ⇒ But we solve it algorithmically, not analytically
- ⇒ When  $p$  is large, finding the best subset is hard

## Shrinkage Operators: Ridge Regression

Ridge regression solves a closely related minimization problem:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq s$$

or, equivalently,

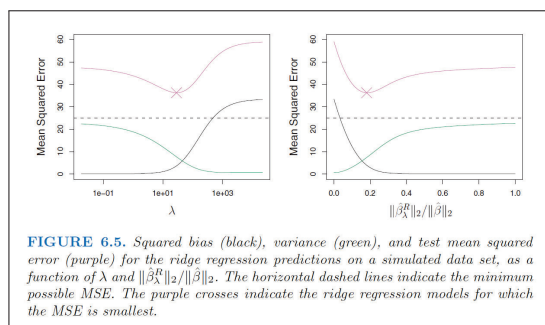
$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

for some **tuning parameter**  $\lambda \geq 0$

Ridge regression shrinks OLS coefficients toward zero

- Shrinkage is more or less proportional, so ridge regression does not identify a subset of regressors to include in the regression model (it just down-weights some relative to others)

## OLS is BLUE, But Ridge Regression Has Lower MSE



**FIGURE 6.5.** Squared bias (black), variance (green), and test mean squared error (purple) for the ridge regression predictions on a simulated data set, as a function of  $\lambda$  and  $\|\hat{\beta}_\lambda^R\|_2 / \|\hat{\beta}\|_2$ . The horizontal dashed lines indicate the minimum possible MSE. The purple crosses indicate the ridge regression models for which the MSE is smallest.

Source: Witten et al. (2023)

**Gauss-Markov Theorem:** OLS is best linear unbiased estimator (BLUE) of  $Y$

- Estimators that are (a little) biased can generate better predictions

## Shrinkage Operators: Lasso

Lasso (Least Absolute Shrinkage and Selection Operator) seeks to minimize:

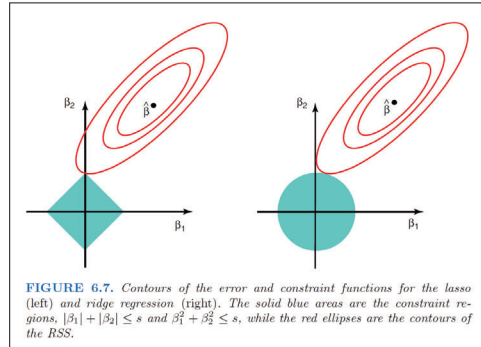
$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

for some **tuning parameter**  $\lambda \geq 0$

**Lasso combines benefits of subset selection, ridge regression**

- Less computationally intensive than subset selection
- Sets some coefficients to 0  $\rightarrow$  identifies parsimonious model
- Better than ridge regression when most covariates are garbage

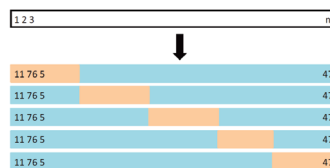
## Lasso Sets Some Coefficients to Zero



Source: Witten et al. (2023)

The lasso constraint region has sharp corners  $\Rightarrow$  some coefficients set to 0

## Three Approaches to Choosing $\lambda$ (1/3)



$k$ -fold cross-validation (default with Stata's `lasso` command):

- Randomly sort observations in  $k$  groups
- For each group  $k$ , estimate lasso on on rest of sample and predict MSE using observations in  $k$  (the hold-out sample); average to get  $MSE(\lambda)$
- Iterate over  $\lambda$  values to choose  $\lambda$  that minimizes MSE



## Three Approaches to Choosing $\lambda$ (2/3)

Bayesian Information Criterion (BIC): function of  $n$ , RSS, number of parameters in the model

$$BIC = \frac{1}{n} (RSS + \log(n)d\hat{\sigma}^2)$$

where  $d$  is the number of covariates and  $\hat{\sigma}^2$  is an unbiased estimate of the variance

- One of several examples of approaches based on (essentially) adjusted  $R^2$
- `bic` is the only one of these options available from Stata's `lasso` command

## Three Approaches to Choosing $\lambda$ (3/3)

Belloni and Chernozhukov (2011), Belloni et al. (2012): alternative approach to choosing  $\lambda$

- Chooses  $\lambda$  iteratively based on data
- Errs on the side of choosing fewer controls to avoid over-fitting
- Allows for heteroskedasticity

Three approaches may generate very different sets of controls

- Costs of too many/too few may vary across empirical contexts

## Using Lasso to Choose Covariates in an RCT

Conceptually, we want to control for covariates that

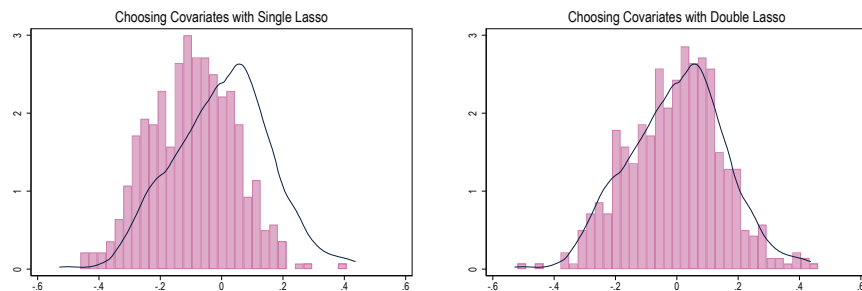
- Predict treatment (i.e. are slightly imbalanced by chance)
- Predict the outcome conditional on treatment (to increase power/precision)
- Over-fitting, noisy controls can make our impact estimates less precise

### Post-double-selection lasso:

- Use lasso to identify baseline covariates that predict treatment
- Use lasso to identify baseline covariates that predict  $Y$  (controlling for  $T$ )
- Include all of the above in regression estimates of treatment effects

Including predictors of  $T$  matters a lot when treatment is not randomly assigned, less in RCTs

## Post-Double-Selection Approach Is Important in Observational Settings



Using lasso to address selection bias through **post-double-selection**:

- Using lasso to predict  $Y$  after controlling for treatment can lead to omitted variable bias since covariates that are correlated with  $T$  will be dropped (since  $T$  is always included)
- When treatment is randomized, should not be a major issue

## Using Lasso to Choose Covariates in an RCT: Example

Variable	Description
act_any	RCT arm: assigned ACT any subsidized price (less than 500 KSh)
c_act	Episode care: Treated with ACT
b_h_edu	Baseline: years of education of household head
b_knowledge	Baseline: Knows only mosquitoes transmit malaria
b_hh_size	Baseline: Number of household members
b_acres	Baseline: Acres of land owned by household head
b_dist_km	Baseline: Distance (km) from household to study chemist
b_h_age_imputed	Baseline: Age of head (missing replaced by sample mean)
b_h_age_missing	Baseline: Age of head missing

## Preparing the Data

Variable	Obs.	Mean	SD	Min.	Max.
act_any	575	0.7391304	0.4394912	0	1
c_act	575	0.3408696	0.4744143	0	1
b_h_edu	572	5.300699	3.940068	0	16
b_knowledge	575	0.5391304	0.4989005	0	1
b_hh_size	575	5.466087	2.48658	1	14
b_acres	462	2.270779	2.697221	0	34
b_dist_km	574	1.667154	.9307161	.0332327	3.982576
b_h_age_imputed	575	39.2537	15.3422	17	88
b_h_age_missing	575	0.0504348	0.2190309	0	1

## Using Stata's lasso Command

```
lasso linear c_act (act_any) b_*, sel(cv)
...
display "e(othervars_sel)"
b_h_edu b_knowledge b_acres b_dist_km b_h_age_imputed b_h_age_missing

lasso linear c_act (act_any) b_*, sel(cv)
display "e(othervars_sel)"
... no covariates selected ...
```

## Using Stata's lasso Command: Results

	<b>OLS</b>	<b>PDSL</b>
	(1)	(2)
Treatment	0.196	0.185
	(0.045)	(0.044)
	[4.40]	[4.24]
Covariates that predict $T$	No	Yes
Covariates that predict $Y$	No	Yes
Covariates included	0	6

Standard errors in parentheses; t-statistics in brackets.

## Using Stata's lasso Command: Results

	OLS	PDSL
	(1)	(2)
Treatment	0.196	0.185
	(0.045)	(0.044)
	[4.40]	[4.24]
Covariates that predict $T$	No	Yes
Covariates that predict $Y$	No	Yes
Covariates included	0	6
<b>Residual variance</b>	<b>0.2178</b>	<b>0.2056</b>

Standard errors in parentheses; t-statistics in brackets.

## When Are Covariates Important?

$$MDE = (t_{1-\kappa} + t_{\alpha/2}) \sqrt{\frac{1}{P(1-P)}} \sqrt{\frac{\hat{\sigma}^2}{N}} \sqrt{1 + (n_g - 1) \tilde{\rho}}$$

Covariates can increase statistical power substantially when outcomes are serially correlated

- Covariates have limited impact when outcomes vary a lot over time
- Malaria example: residual variance 94 percent of original
- Children's storybooks example: residual variance 46 percent of original