

09 - High-Dimensional Confounding Adjustment

ml4econ, HUJI 2023

Itamar Caspi
June 4, 2023 (updated: 2023-06-04)

Replicating this presentation

Use the **pacman** package to install and load packages:

```
if (!require("pacman"))  
  install.packages("pacman")  
  
pacman::p_load(  
  tidyverse,  
  tidymodels,  
  hdm,  
  ggdag,  
  knitr,  
  xaringan,  
)
```

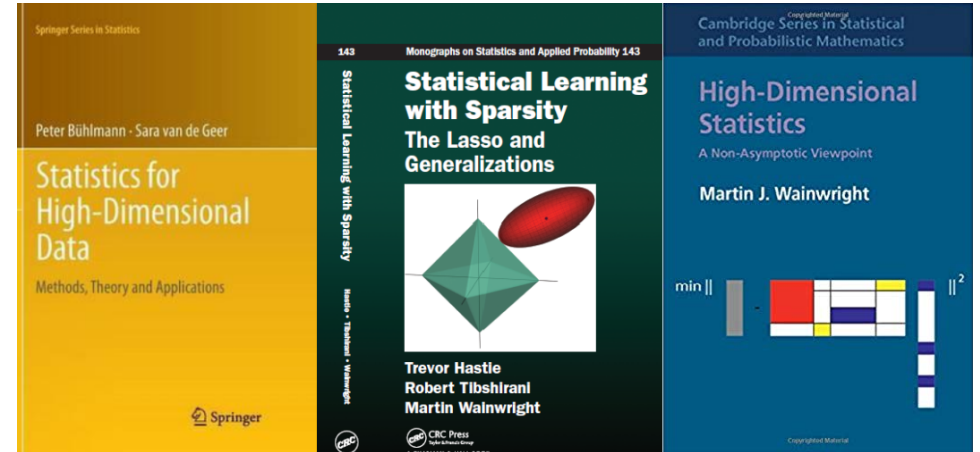
Outline

- Lasso and Variable Selection
- High Dimensional Confoundedness
- Empirical Illustration using `hdm`

Lasso and Variable Selection

Key Lasso Theory Resources

- *Statistical Learning with Sparsity - The Lasso and Generalizations* (Hastie, Tibshirani, and Wainwright), **Chapter 11: Theoretical Results for the Lasso.** (PDF available online)
- *Statistics for High-Dimensional Data - Methods, Theory and Applications* (Bühlmann and van de Geer), **Chapter 7: Variable Selection with the Lasso.**
- *High Dimensional Statistics - A Non-Asymptotic Viewpoint* (Wainwright), **Chapter 7: Sparse Linear Models in High Dimensions**



Guidance vs. Guarantees: Fundamental Differences

- We've primarily relied on *guidance* for our work:
 - Selection of folds in CV
 - Size determination of the holdout set
 - Tuning parameter(s) adjustment
 - Loss function selection
 - Function class selection
- But in causal inference, *guarantees* become vital:
 - Selecting variables
 - Deriving confidence intervals and p -values
- To attain these guarantees, we generally need:
 - Assumptions regarding a "true" model
 - Asymptotic principles, such as $n \rightarrow \infty$, $k \rightarrow ?$

Key Notations in Lasso Literature

Assume $\boldsymbol{\beta}$ is a $k \times 1$ vector with a typical element as β_i .

- ℓ_0 -norm is $\|\boldsymbol{\beta}\|_0 = \sum_{j=1}^k \mathbf{1}_{\{\beta_j \neq 0\}}$, indicating the count of non-zero elements in $\boldsymbol{\beta}$.
- ℓ_1 -norm is $\|\boldsymbol{\beta}\|_1 = \sum_{j=1}^k |\beta_j|$.
- ℓ_2 -norm or Euclidean norm is $\|\boldsymbol{\beta}\|_2 = \left(\sum_{j=1}^k |\beta_j|^2 \right)^{\frac{1}{2}}$.
- ℓ_∞ -norm is $\|\boldsymbol{\beta}\|_\infty = \sup_j |\beta_j|$, signifying the maximum magnitude among $\boldsymbol{\beta}$'s entries.
- Support of $\boldsymbol{\beta}$ is $S \equiv \text{supp}(\boldsymbol{\beta}) = \{j : \beta_j \neq 0, j = 1, \dots, k\}$, the subset of non-zero coefficients.
- Size of the support $s = |S|$ is the count of non-zero elements in $\boldsymbol{\beta}$, namely $s = \|\boldsymbol{\beta}\|_0$

Understanding the Basic Setup of Lasso

The linear regression model is given as:

$$Y_i = \alpha + X_i' \boldsymbol{\beta}^0 + \varepsilon_i, \quad i = 1, \dots, n,$$

$$\mathbb{E} [\varepsilon_i X_i] = 0, \quad \alpha \in \mathbb{R}, \quad \boldsymbol{\beta}^0 \in \mathbb{R}^k.$$

Under the **exact sparsity** assumption, we include only a subset of variables of size $s \ll k$ in the model, where $s \equiv \|\boldsymbol{\beta}\|_0$ represents the sparsity index.

$$\underbrace{\mathbf{X}_S = (X_{(1)}, \dots, X_{(s)})}_{\text{Sparse Variables}}, \quad \underbrace{\mathbf{X}_{S^c} = (X_{(s+1)}, \dots, X_{(k)})}_{\text{Non-Sparse Variables}}$$

Here, S is the subset of active predictors, $\mathbf{X}_S \in \mathbb{R}^{n \times s}$ corresponds to the subset of covariates in the sparse set, and $\mathbf{X}_{S^c} \in \mathbb{R}^{n \times k-s}$ refers to the subset of "irrelevant" non-sparse variables.

Lasso: The Optimization

The Lasso (Least Absolute Shrinkage and Selection Operator), introduced by Tibshirani (1996), poses the following optimization problem:

$$\min_{\beta_0, \beta} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \|\beta\|_1$$

In this setup, Lasso places a "budget constraint" on the sum of *absolute* values of β 's.

Differing from ridge, the Lasso penalty is linear (shifting from 1 to 2 bears the same weight as moving from 101 to 102).

A major strength of Lasso lies in its ability to perform model selection - it zeroes out most of the β 's in the model, making the solution *sparse*.

Any penalty involving the ℓ_1 norm will achieve this.

Evaluating the Lasso

Suppose β^0 is the true vector of coefficients and $\hat{\beta}$ represents the Lasso estimator. We can evaluate Lasso's effectiveness in several ways:

I. Prediction Quality

$$\text{Loss}_{\text{pred}}(\hat{\beta}; \beta^0) = \frac{1}{N} \|(\hat{\beta} - \beta^0)\mathbf{X}\|_2^2 = \frac{1}{N} \sum_{j=1}^k [(\hat{\beta}_j - \beta_j^0)\mathbf{X}_{(j)}]^2$$

II. Parameter Consistency

$$\text{Loss}_{\text{param}}(\hat{\beta}; \beta^0) = \|\hat{\beta} - \beta^0\|_2^2 = \sum_{j=1}^k (\hat{\beta}_j - \beta_j^0)^2$$

III. Support Recovery (Sparsistency)

For example, score +1 if $\text{sign}(\beta^0) = \text{sign}(\beta_j)$ for all $j = 1, \dots, k$, and 0 otherwise.

Leveraging Lasso for Variable Selection

- Variable selection consistency is crucial for causal inference, considering omitted variable bias.
- Lasso frequently serves as a tool for variable selection.
- The successful selection of the "true" support by Lasso depends heavily on strong assumptions about:
 - Distinguishing between relevant and irrelevant variables.
 - Identifying β .

Critical Assumption #1: Distinguishable Sparse Betas

Lower Eigenvalue: The minimum eigenvalue, λ_{\min} , of the sub-matrix \mathbf{X}_S , should be bounded away from zero.

$$\lambda_{\min} (\mathbf{X}_S' \mathbf{X}_S / N) \geq C_{\min} > 0$$

Linear dependence between the columns of \mathbf{X}_S makes it impossible to identify the true β , even if we *knew* which variables are included in \mathbf{X}_S .

NOTE: The high-dimensional lower eigenvalue condition replaces the low-dimensional rank condition (i.e., that $\mathbf{X}'\mathbf{X}$ is invertible).

Critical assumption #2: Distinguishable active predictors

Irrepresentability condition (Zou ,2006; Zhao and Yu, 2006): There must exist some $\eta > 0$ such that

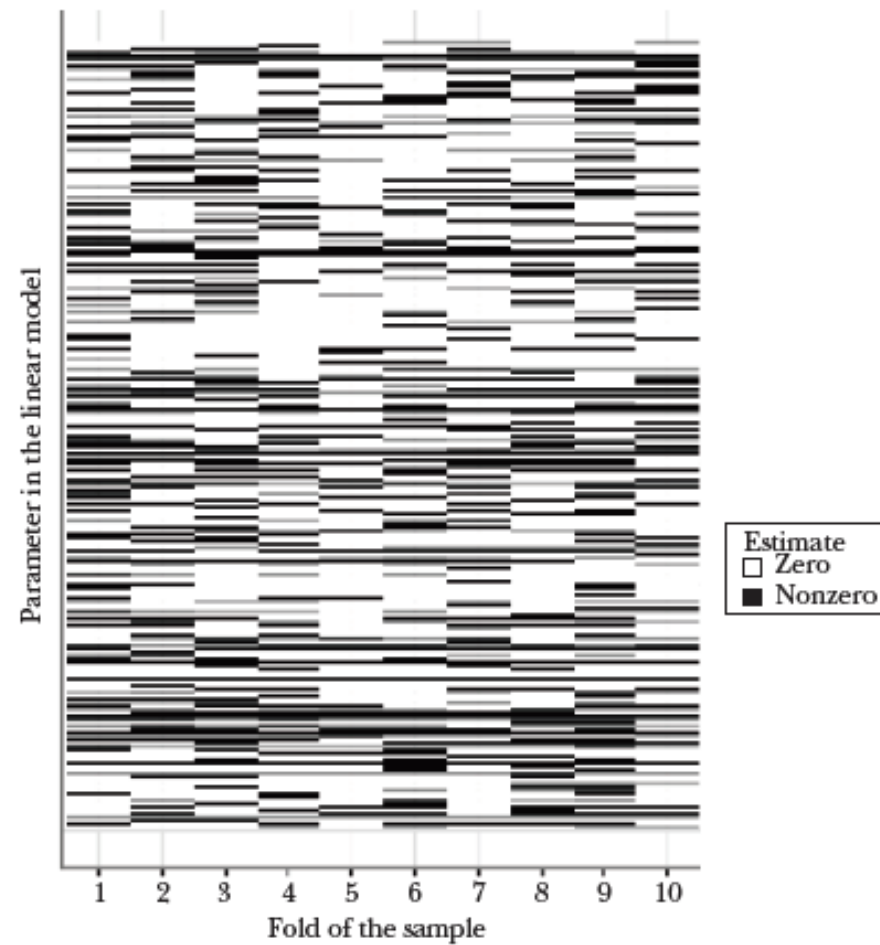
$$\max_{j \in S^c} \left\| (\mathbf{X}'_S \mathbf{X}_S)^{-1} \mathbf{X}'_S \mathbf{x}_j \right\|_1 \leq 1 - \eta$$

INTUITION: What's inside $\|\cdot\|_1$ is like regressing \mathbf{x}_j on the variables in \mathbf{X}_S .

- When $\eta = 1$, the sparse and non-sparse variables are orthogonal to each other.
- When $\eta = 0$, we can perfectly reconstruct (some elements of) \mathbf{X}_S using \mathbf{X}_{S^c} .

Thus, the irrepresentability condition roughly states that we can distinguish the sparse variables from the non-sparse ones.

Figure 2
Selected Coefficients (Nonzero Estimates) across Ten LASSO Regressions



Source: [Mullainathan and Spiess \(JEP 2017\)](#).

Setting the Optimal Tuning Parameter

- Throughout this course, we have frequently chosen λ empirically, often by cross-validation, based on its predictive performance.
- In causal analysis, however, the end goal is inference, not prediction. These two objectives often conflict (bias vs. variance).
- Ideally, the choice of λ should provide assurances about the model's performance.
- Generally, for satisfying sparsistency, we set λ such that it selects non-zero β 's with a high probability.

High Dimensional Confoundedness

"Naive" Implementation of the Lasso

Run `glmnet`

```
glmnet(Y ~ DX)
```

where DX is the feature matrix which includes both the treatment D and the features vector X . The estimated coefficients are:

$$\left(\hat{\alpha}, \hat{\tau}, \hat{\beta}'\right)' = \arg \min_{\alpha, \tau \in \mathbb{R}, \beta \in \mathbb{R}^{k+1}} \sum_{i=1}^n (Y_i - \alpha - \tau D_i - \beta' X_i)^2 + \lambda \left(|\tau| + \sum_{j=1}^k |\beta_j| \right)$$

ISSUES:

1. Both $\hat{\tau}$ and $\hat{\beta}$ experience shrinkage, thus biased towards zero.
2. Lasso might eliminate D_i , i.e., shrink $\hat{\tau}$ to zero. The same can happen to relevant confounding factors.
3. The choice of λ is a challenge.

Moving Towards a Solution

To avoid eliminating D_i , we can adjust the Lasso regression:

$$\left(\hat{\alpha}, \hat{\tau}, \hat{\beta}'\right)' = \arg \min_{\alpha, \tau \in \mathbb{R}, \beta \in \mathbb{R}^k} \sum_{i=1}^n \left(Y_i - \alpha - \tau D_i - \beta' X_i\right)^2 + \lambda \left(\sum_{j=1}^k |\beta_j|\right)$$

We can then *debias* the results using the "Post-Lasso" method, i.e., use the Lasso for variable selection, then run OLS with the selected variables.

ISSUES: The Lasso might eliminate features that are correlated with D_i because they are not good predictors of Y_i , leading to *omitted variable bias*.

Problem Solved?

What can go wrong? Distribution of $\sqrt{n}(\hat{\alpha} - \alpha)$ is not what you think

$$y_i = d_i\alpha + x_i\beta + \epsilon_i, \quad d_i = x_i\gamma + v_i$$

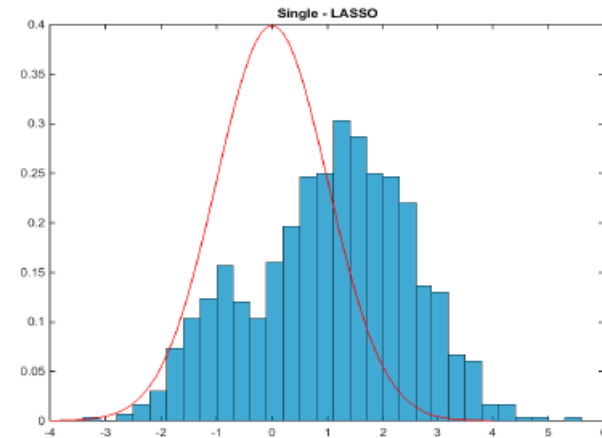
$$\alpha = 0, \quad \beta = .2, \quad \gamma = .8,$$

$$n = 100$$

$$\epsilon_i \sim N(0, 1)$$

$$(d_i, x_i) \sim N\left(0, \begin{bmatrix} 1 & .8 \\ .8 & 1 \end{bmatrix}\right)$$

► selection done by
Lasso



Reject $H_0 : \alpha = 0$ (the truth) of no effect about 50% of the time

Source: <https://stuff.mit.edu/~vchern/papers/Chernozhukov-Saloniki.pdf>

Solution: Double-selection Lasso (Belloni, et al., REStud 2013)

Step 1: Perform two Lasso regressions: Y_i on X_i and D_i on X_i :

$$\hat{\gamma} = \arg \min_{\gamma \in \mathbb{R}^{p+1}} \sum_{i=1}^n (Y_i - \gamma' X_i)^2 + \lambda_{\gamma} \left(\sum_{j=2}^p |\gamma_j| \right)$$
$$\hat{\delta} = \arg \min_{\delta \in \mathbb{R}^{q+1}} \sum_{i=1}^n (D_i - \delta' X_i)^2 + \lambda_{\delta} \left(\sum_{j=2}^q |\delta_j| \right)$$

Step 2: Refit the model using OLS, but only include the \mathbf{X} 's that were significant predictors of both Y_i and D_i .

Step 3: Proceed with the inference using standard confidence intervals.

The tuning parameter λ is set in a way that ensures non-sparse coefficients are correctly selected with high probability.

Does it Work?

Double Selection Works

$$y_i = d_i\alpha + x_i\beta + \epsilon_i, \quad d_i = x_i\gamma + v_i$$

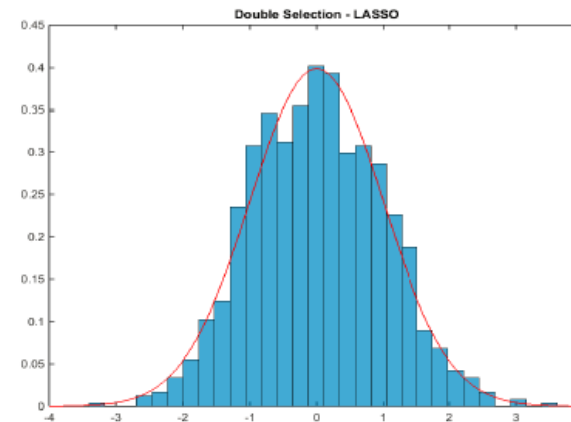
$$\alpha = 0, \quad \beta = .2, \quad \gamma = .8,$$

$$n = 100$$

$$\epsilon_i \sim N(0, 1)$$

$$(d_i, x_i) \sim N\left(0, \begin{bmatrix} 1 & .8 \\ .8 & 1 \end{bmatrix}\right)$$

► double selection
done by Lasso



Reject $H_0 : \alpha = 0$ (the truth) about 5% of the time (nominal size = 5%)

Source: <https://stuff.mit.edu/~vchern/papers/Chernozhukov-Saloniki.pdf>

Statistical Inference

Uniform Validity of the Double Selection

Theorem (Belloni, Chernozhukov, Hansen: WC 2010, ReStud 2013)

Uniformly within a class of approximately sparse models with restricted isometry conditions

$$\sigma_n^{-1} \sqrt{n}(\hat{\alpha} - \alpha_0) \rightarrow_d N(0, 1),$$

where σ_n^2 is conventional variance formula for least squares. Under homoscedasticity, semi-parametrically efficient.

- ▶ Model selection mistakes are asymptotically negligible due to double selection.
- ▶ Analogous result also holds for *endogenous* models, see Chernozhukov, Hansen, Spindler, *Annual Review of Economics*, 2015.

Source: <https://stuff.mit.edu/~vchern/papers/Chernozhukov-Saloniki.pdf>

Intuition: Partialling-out regression

Consider two methods for estimating the effect of X_{1i} (a scalar) on Y_i , while adjusting for X_{2i} :

Alternative 1: Run

$$Y_i = \alpha + \beta X_{1i} + \gamma X_{2i} + \varepsilon_i$$

Alternative 2: First, run Y_i on X_{2i} and X_{1i} on X_{2i} and keep the residuals, i.e., run

$$Y_i = \gamma_0 + \gamma_1 X_{2i} + u_i^Y, \quad \text{and} \quad X_{1i} = \delta_0 + \delta_1 X_{2i} + u_i^{X_1},$$

and keep \hat{u}_i^Y and $\hat{u}_i^{X_1}$. Next, run

$$\hat{u}_i^Y = \beta^* \hat{u}_i^{X_1} + v_i.$$

According to the **Frisch-Waugh-Lovell (FWL) Theorem**, $\hat{\beta} = \hat{\beta}^*$.

Guarantees of Double-selection Lasso (VERY Wonkish)

Approximate Sparsity Consider the following regression model:

$$Y_i = f(W_i) + \varepsilon_i = X_i' \beta^0 + r_i + \varepsilon_i, \quad 1, \dots, n$$

where r_i is the approximation error.

Under *approximate sparsity*, it is assumed that $f(W_i)$ can be approximated sufficiently well (up to r_i) by $X_i' \beta^0$, while using only a small number of non-zero coefficients.

Restricted Sparse Eigenvalue Condition (RSEC) This condition puts bounds on the number of variables outside the support the Lasso can select. Relevant for the post-lasso stage.

Regularization Event The tuning parameter λ is to a value that it selects to correct model with probability of at least p , where p is set by the user. Further assumptions regarding the quantile function of the maximal value of the gradient of the objective function at β^0 , and the error term (homoskedasticity vs. heteroskedasticity). See Belloni et al. (2012) for further details.

Additional Extensions of Double-selection

1. **Other Function Classes (Double-ML):** Chernozhukov et al. (AER 2017) proposed using other function classes, such as applying random forest for $Y \sim X$ and regularized logit for $D \sim X$.
2. **Instrumental Variables:** Techniques involving instrumental variables have been developed by Belloni et al. (Ecta 2012) and Chernozhukov et al. (AER 2015). For further understanding, please refer to the problem set.
3. **Heterogeneous Treatment Effects:** Heterogeneous treatment effects have been studied by Belloni et al. (Ecta 2017). We'll explore this topic more thoroughly next week.
4. **Panel Data:** Consideration for panel data was made by Belloni et al. (JBES 2016).

Evidence on the Applicability of Double-Lasso

"Machine Labor" (Angrist and Frandsen, 2022 JLE):

- ML can be useful for regression-based causal inference using Lasso.
 - Post-double-selection (PDS) Lasso offers data-driven sensitivity analysis.
 - ML-based instrument selection can improve on 2SLS, but split-sample IV and limited information maximum likelihood (LIML) perform better.
 - ML might not be optimal for Instrumental Variables (IV) applications in labor economics. This is due to the creation of artificial exclusion restrictions potentially resulting in inaccurate findings.
-

More from "Labor Machine"

quality of PDS bias mitigation depends on design features like regressor variance and the extent of OVB. Moreover, we have examined a scenario in which OLS with full-dictionary control is feasible and effectively removes OVB. Even so, **PDS seems a useful tool for sensitivity analysis in a regression context, where analysts may choose from an abundance of possible control variables.** Findings where the target causal estimate remains reasonably stable while the list of selected controls varies from one routine to another reinforce claims of robustness.

It is worth emphasizing that a causal interpretation of the ML estimates in table 2 turns on a maintained conditional independence assumption. **ML methods do not create quasi-experimental variation. Rather, ML uses data to pick from among a large set of modeling options founded on a common identifying assumption.** This facilitates estimation in high-dimensional control scenarios and may increase precision (although that is not the finding here). We have also noted considerable sensitivity to implementation details, specifically to software choice and lasso penalty determination.

Source: Angrist and Frandsen (2022).

Table 2
Post-lasso Estimates of Elite College Effects

	Double Selection (PDS)			Outcome Selection			All Controls
	Plug-In (1)	CV λ (2)	cvarlasso (3)	Plug-In (7)	CV λ (8)	cvarlasso (9)	OLS (7)
A. Private School Effects							
Estimate	.038 (.040)	.020 (.039)	.040 (.041)	.046 (.041)	.043 (.043)	.042 (.043)	.017 .039
Number of controls	18	100	112	10	35	50	303
B. Effects of School-Average SAT Score/100							
Estimate	-.009 (.020)	-.013 (.018)	-.009 (.019)	-.008 (.020)	-.009 (.019)	-.008 (.019)	-.012 (.018)
Number of controls	24	151	58	10	34	43	303
C. Effects of Attending Schools Rated Highly Competitive or Better							
Estimate	.068 (.033)	.051 (.033)	.073 (.033)	.076 (.031)	.080 (.032)	.082 (.032)	.053 .033
Number of controls	17	185	106	10	34	43	303

DML: A Cautionary Tale

Hünormund, Louw, and Caspi (2023 JCI):

- DML is highly sensitive to a few "bad controls" in the covariate space, leading to potential bias.
- This bias varies depending on the theoretical causal model, raising questions about the practicality of data-driven control variable selection.

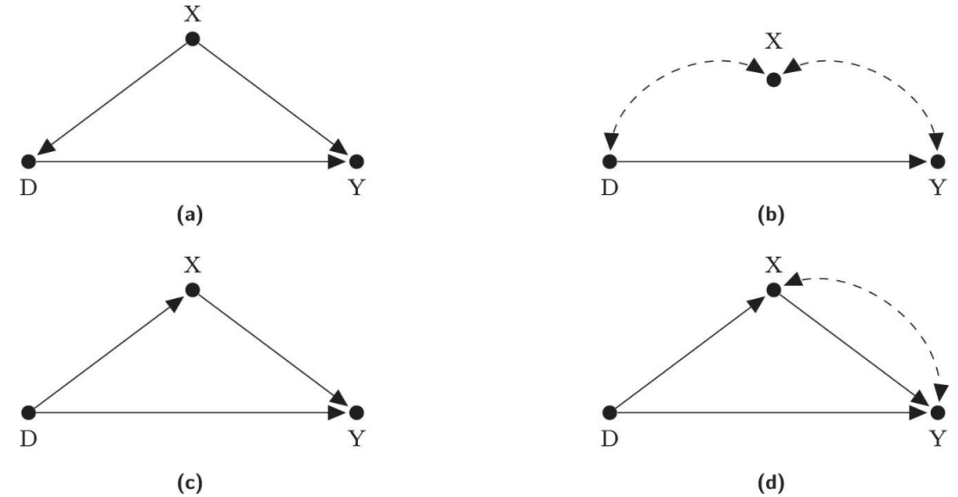


Figure 1

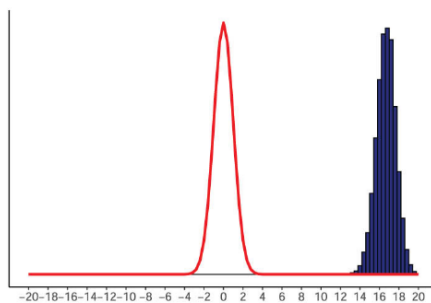
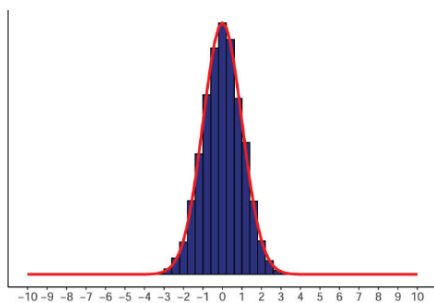
Directed acyclic graphs representing different structural causal models. (a) Good control, (b) M-graph, (c) mediator, and (d) confounded mediator.

Simulation Results

Double Machine Learning

Naïve Lasso

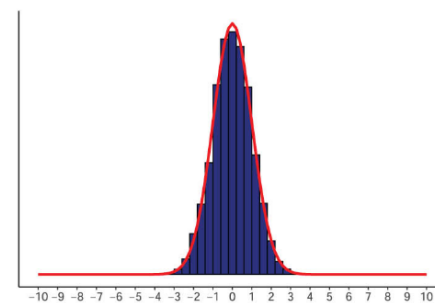
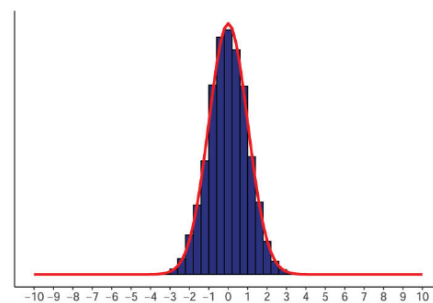
(a)



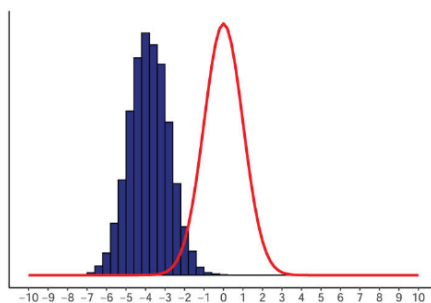
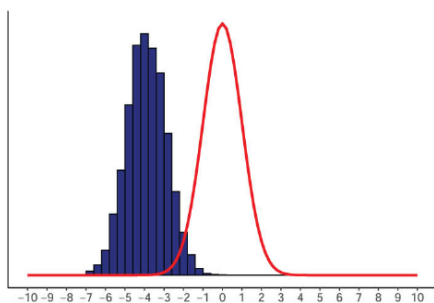
Double Machine Learning

Naïve Lasso

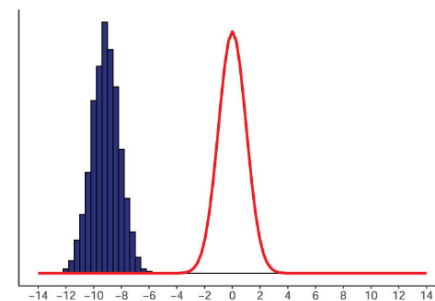
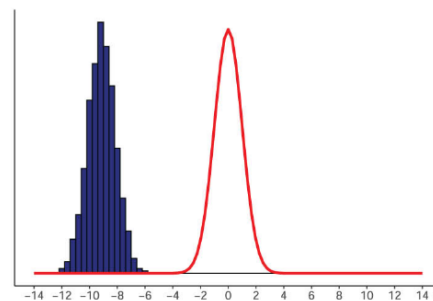
(c)



(b)



(d)



Empirical Relevance

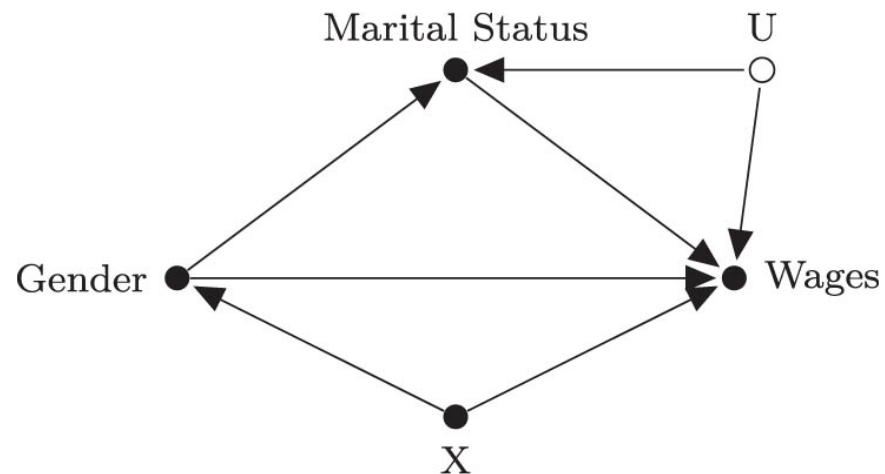


Table 2

Effect of gender on log wages using PSID data from [25] (standard errors in parentheses)

Wave	1981	1990	1999	2007	2009	2011
OLS	−0.249 (0.016)	−0.137 (0.014)	−0.158 (0.016)	−0.168 (0.015)	−0.157 (0.015)	−0.145 (0.016)
DML	−0.268 (0.017)	−0.139 (0.015)	−0.158 (0.016)	−0.164 (0.016)	−0.157 (0.016)	−0.136 (0.017)
DML incl. <i>Marital status</i>	−0.270 (0.022)	−0.154 (0.019)	−0.173 (0.020)	−0.190 (0.019)	−0.179 (0.020)	−0.163 (0.021)

Empirical Illustration using hdm

Introducing the hdm R Package

"**High-Dimensional Metrics**" (hdm) by Victor Chernozhukov, Chris Hansen, and Martin Spindler is an R package for estimation and quantification of uncertainty in high-dimensional approximately sparse models.

[*] A Stata module named **Lassopack** offers a comprehensive set of programs for regularized regression in high-dimensional contexts..]

Illustration: Testing for Growth Convergence

The standard empirical model for growth convergence is represented by the equation:

$$Y_{i,T} = \alpha_0 + \alpha_1 Y_{i,0} + \sum_{j=1}^k \beta_j X_{ij} + \varepsilon_i, \quad i = 1, \dots, n,$$

where

- $Y_{i,T}$ national growth rates in GDP per capita for the periods 1965-1975 and 1975-1985.
- $Y_{i,0}$ is the log of the initial level of GDP at the beginning of the specified decade.
- X_{ij} covariates which might influence growth.

The growth convergence hypothesis implies that $\alpha_1 < 0$.

Growth Data

To test the **growth convergence hypothesis**, we will employ the Barro and Lee (1994) dataset.

```
data("GrowthData")
```

The data features macroeconomic information for a substantial group of countries over various decades. Specifically,

- n equals 90 countries
- k equals 60 country features

While these numbers may not seem large, the quantity of covariates is substantial compared to the sample size. Hence, **variable selection** is crucial!

Let's Have a Look

```
GrowthData %>%  
  as_tibble %>%  
  head(2)
```

```
## # A tibble: 2 x 63  
##   Outcome interc~1 gdpsh~2 bmp1l freeop freetar h65 hm65 hf65 p65 pm65 pf65 s65  
##   <dbl> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
## 1 -0.0243      1  6.59 0.284 0.153 0.0439 0.007 0.013 0.001 0.29 0.37 0.21 0.04  
## 2  0.100      1  6.83 0.614 0.314 0.0618 0.019 0.032 0.007 0.91 1 0.65 0.16  
## # ... with 48 more variables: fert65 <dbl>, mort65 <dbl>, lifee065 <dbl>, gpop1 <dbl>, fert  
## # mort1 <dbl>, invsh41 <dbl>, geetot1 <dbl>, geerec1 <dbl>, gde1 <dbl>, govwb1 <dbl>,  
## # govsh41 <dbl>, gvxdxe41 <dbl>, high65 <dbl>, highm65 <dbl>, highf65 <dbl>, highc65 <dbl>  
## # highcm65 <dbl>, highcf65 <dbl>, human65 <dbl>, humanm65 <dbl>, humanf65 <dbl>, hyr65 <dbl>  
## # hyrm65 <dbl>, hyrf65 <dbl>, no65 <dbl>, nom65 <dbl>, nof65 <dbl>, pinstab1 <dbl>, pop65  
## # worker65 <dbl>, pop1565 <dbl>, pop6565 <dbl>, sec65 <dbl>, secm65 <dbl>, secf65 <dbl>,  
## # secc65 <dbl>, seccm65 <dbl>, seccf65 <dbl>, syr65 <dbl>, syrm65 <dbl>, syrf65 <dbl>, ..
```

Data Processing

Rename the response and "treatment" variables:

```
df <-  
  GrowthData %>%  
  rename(YT = Outcome, Y0 = gdpsh465)
```

Transform the data to vectors and matrices (to be used in the `rlassoEffect()` function)

```
YT <- df %>% select(YT) %>% pull()  
  
Y0 <- df %>% select(Y0) %>% pull()  
  
X <- df %>%  
  select(-c("Y0", "YT")) %>%  
  as.matrix()  
  
Y0_X <- df %>%  
  select(-YT) %>%  
  as.matrix()
```

Estimation of the Convergence Parameter α_1

Method 1: OLS

```
ols <- lm(YT ~ ., data = df)
```

Method 2: Naive (rigorous) Lasso

```
naive_Lasso <- rlasso(x = Y0_X, y = YT)
```

Does the Lasso drop Y_0 ?

```
naive_Lasso$beta[2]
```

```
## Y0  
## 0
```

Unfortunately, yes...

Estimation of the Convergence Parameter α_1

Method 3: Partialling-out Lasso

```
part_Lasso <-  
  rlassoEffect(  
    x = X, y = YT, d = Y0,  
    method = "partialling out"  
  )
```

Method 4: Double-selection Lasso

```
double_Lasso <-  
  rlassoEffect(  
    x = X, y = YT, d = Y0,  
    method = "double selection"  
  )
```

Tidying the Results

OLS

```
ols_tbl <- tidy(ols) %>%
  filter(term == "Y0") %>%
  mutate(method = "OLS") %>%
  select(method, estimate, std.error)
```

Naive Lasso

```
naive_Lasso_tbl <- tibble(method = "Naive Lasso",
                           estimate = NA,
                           std.error = NA)
```

Partialling-out Lasso

[illegible]

Double-selection Lasso

[illegible]

Summary of the Convergence Test

```
bind_rows(ols_tbl, naive_Lasso_tbl, part_Lasso_tbl, double_Lasso_tbl) %>%  
  kable(digits = 3, format = "html")
```

method	estimate	std.error
OLS	-0.009	0.030
Naive Lasso	NA	NA
Partialling-out Lasso	-0.050	0.014
Double-selection Lasso	-0.050	0.016

The use of double-selection and partialling-out methods lead to significantly **more precise estimates** and lend support to the **conditional convergence hypothesis**.

An Advanced R Package: DoubleML

- The Python and R packages `{DoubleML}` offer a modern implementation of the double / debiased machine learning framework.
- For more details, visit the [Getting Started](#) and [Examples](#) sections.
- The package is constructed on the `{mlr3}` ecosystem.



```
slides %>% end()
```

 [Source code](#)

Selected References

Ahrens, A., Hansen, C. B., & Schaffer, M. E. (2019). lassopack: Model selection and prediction with regularized regression in Stata.

Angrist, Joshua D, and Alan B Krueger. 1991. "Does Compulsory School Attendance Affect Schooling and Earnings?" *The Quarterly Journal of Economics*, 106(4): 979–1014.

Angrist, J. D., & Frandsen, B. (2022). Machine labor. *Journal of Labor Economics*, 40(S1), S97-S140.

Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen. 2012. Sparse Models and Methods for Optimal Instruments With an Application to Eminent Domain. *Econometrica* 80(6): 2369–2429.

Belloni, A., & Chernozhukov, V. (2013). Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2), 521–547.

Belloni, A., Chernozhukov, V., & Hansen, C. (2013). Inference on treatment effects after selection among high-dimensional controls. *Review of Economic Studies*, 81(2), 608–650.

Selected References

Belloni, A., Chernozhukov, V., & Hansen, C. (2014). High-Dimensional Methods and Inference on Structural and Treatment Effects. *Journal of Economic Perspectives*, 28(2), 29–50.

Chernozhukov, V., Hansen, C., & Spindler, M. (2015). Post-selection and post-regularization inference in linear models with many controls and instruments. *American Economic Review*, 105(5), 486–490.

Chernozhukov, V., Hansen, C., & Spindler, M. (2016). hdm: High-Dimensional Metrics. *The R Journal*, 8(2), 185–199.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., & Newey, W. (2017). Double/debiased/Neyman machine learning of treatment effects. *American Economic Review*, 107(5), 261–265.

Selected References

Dale, Stacy Berg, and Alan B Krueger. 2002. "Estimating the Payoff to Attending a More Selective College: An Application of Selection on Observables and Unobservables." *The Quarterly Journal of Economics*, 117(4): 1491–1527.

Hünermund, P., Louw, B., & Caspi, I. (2023). Double machine learning and automated confounder selection: A cautionary tale. *Journal of Causal Inference*, 11(1), 20220078.

Mullainathan, S. & Spiess, J., 2017. Machine Learning: An Applied Econometric Approach. *Journal of Economic Perspectives*, 31(2), pp.87–106.

Van de Geer, S. A., & Bühlmann, P. (2009). On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3, 1360–1392.

Zhao, P., & Yu, B. (2006). On Model Selection Consistency of Lasso. *Journal of Machine Learning Research*, 7, 2541–2563.