

Causal Machine Learning

Prof. Tzu-Ting Yang
楊子霆

Institute of Economics, Academia Sinica
中央研究院經濟研究所

October 9, 2024

Machine Learning: Main Idea

Machine learning

- Machine learning methods: Use **data-driven algorithms** to predict outcome **Y** given many covariates **X**.
 - There are **many** machine learning methods
 - The best methods vary with the particular data application
 - You can consider **regression** is a type of machine learning methods
- The main goal is **prediction** or **classification**
 - This is useful in some economics applications
 - Forecast economic growth rate using many factors
 - Predict user-rating of products
 - Classify the types of individuals given many socio-economic measures and predict their loan repayment probability

Machine learning and Causal Inference

- For causal inference, machine learning methods help us select the key control variables from high-dimensional data
 - Many covariates \mathbf{X}

Problem of High-Dimensional Data

Problem of High-Dimensional Data

- Consider linear regression model with p potential covariates where p is too large.

$$Y_i = \beta_0 + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i, \quad i = 1, \dots, n$$

- Y_i is observed outcome for individual i
 - X_i^j is observed covariate j for individual i
 - n is sample size
 - p is the number of covariates
- Problem:** p could be much larger than n

Why We Need High-Dimensional Data?

- Why we need so many covariates?
 - Including more covariates can reduce **omitted variable bias** (**selection bias**)
 - Linear regression may predict well if include many relevant covariates
 - The number of covariates increases if we account for non-linearity or interaction effects
- **Example:**
 - Cross-country regressions, where we have only small number of countries, but thousands of macro variables.

Problem of High-Dimensional Data

Example 1

Sala-i-Martin, Xavier (1997), "**I Just Ran Two Million Regressions**", American Economic Review

- The author tries to examine the hypothesis of growth convergence and the determinants of economic growth
 - Whether poorer economies' per capita incomes will tend to grow at faster rates than richer economies
- He finds that a substantial number of variables can be found to be strongly related to growth
- Citations: 3,888 times

Example 1: Test the Convergence Hypothesis

- Examine the relation between GDP growth rate and initial per capita GDP:

$$\underbrace{\text{GrowthRate}}_{Y_i} = \beta_0 + \underbrace{\alpha}_{\text{ATE}} \underbrace{\log(\text{GDP})}_{D_i} + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i$$

- Control a lot of covariates X_i^j
- In their data, they have $p = 60$ covariates, $n = 90$ observations
 - Need to do variable selection
- Test the convergence hypothesis $\alpha < 0$
 - Poor countries catch up with richer countries, conditional on similar characteristics (e.g. institutions, human capital etc.)
 - Prediction from the classical Solow growth model.

Problem of High-Dimensional Data

Example 2

John J. Donohue and Steven D. Levitt (2001), “**The Impact of Legalized Abortion on Crime**”, The Quarterly Journal of Economics

- The authors examine the effect of legalized abortion on crime rate
 - Mechanism:
 - Fewer children at the highest risk of committing crime being born due to the availability of the procedure
- They offer evidence that legalized abortion has contributed significantly to recent crime reductions
- Crime began to fall roughly eighteen years after abortion legalization
- A lot of debates on this issue

Example 2: Effect of Abortion on Crime

- **Goal:** Understand causal effect of d_{it} (abortion) on y_{it} (crime)
- **Problem:** Abortion rates are not randomly assigned
- **Key concern:**
 - States are different for lots of reasons
 - Crime rates in states evolve differently for lots of reasons
 - Factors that are associated to changes in crime rates could also be correlated with the changes in abortion rates
 - For example, share of high school students dropped out

Example 2: Effect of Abortion on Crime

Baseline Model

- Donohue and Levitt (2001) baseline model

$$y_{it} = \alpha_0 d_{it} + \sum_{j=1}^p \beta_j x_{it}^j + \gamma_t + \delta_i + \epsilon_{it}$$

- Sample size: 50 states and 12 years ($n = 600$)
- y_{it} = crime-rate (violent, property, or murder per 1,000 people)
- d_{it} = "effective" abortion rate
- $p = 284$ **possible covariates** in x_{it}^j (see paper)
 - lagged prisoners, lagged police $\times t$, initial income difference, initial income difference $\times t$, initial beer consumption difference $\times t$, average income, average income $\times t$, initial abortion rate
 - γ_t time effects
 - δ_i state effects

Post-Double Selection Method: Main Idea

Using Machine Learning to Improve Causal Inference

Post-Double Selection Method

- Suppose we want to estimate causal effect of treatment D_i on outcome Y_i
- Control p potential covariates where p is too large

$$Y_i = \beta_0 + \alpha D_i + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i, \quad i = 1, \dots, n$$

- **Traditional Variable Selection:**

- 1 Drop all X_i^j that have small coefficients, using model selection devices (classical such as t-tests or modern)
- 2 Run OLS of Y_i on D_i and selected covariates X_i^j

- **Does not work** because fails to eliminate **omitted variable bias** (Leeb and Potscher, 2009).

Review: Omitted Variable Bias

- Suppose the true model is:

$$Y_i = \delta + \alpha D_i + \beta X_i + \epsilon_i$$

- X_i is the observed characteristics (e.g. family income)
- But we estimate this model:

$$Y_i = \delta + \alpha D_i + u_i$$

- where $u_i = \beta X_i + \epsilon_i$

Review: Omitted Variable Bias

- OVB formula:

$$\begin{aligned}\hat{\alpha} &\xrightarrow{p} \alpha + \frac{\text{Cov}(u_i, D_i)}{V(D_i)} \\ &= \alpha + \beta \frac{\text{Cov}(X_i, D_i)}{V(D_i)}\end{aligned}$$

- The difference between estimated treatment effect $\hat{\alpha}$ and true effect α depends on two components:
 - 1 β : The effect of omitted variable X_i on outcome variable Y_i ,
 - 2 $\frac{\text{Cov}(X_i, D_i)}{V(D_i)}$: The relationship between omitted variable X_i and treatment variable D_i

Post-Double Selection Method

Intuition

- Based on OVB formula, Belloni, Chernozhukov, Fernandez-Val and Hansen (2013) propose **Post-Double Selection (PDS)** approach:
 - 1 Use ML methods (LASSO) to select covariates X_i^j that can predict Y_i .
 - 2 Use ML methods (LASSO) to select covariates X_i^j that can predict D_i .
 - 3 Run OLS of Y_i on D_i and the union of covariates selected in steps 1 and 2
- The additional selection step 2 can help eliminate the **omitted variable bias**
- This method addresses selection bias under the assumption of **selection on observables (CIA)**, but may not fully eliminate bias from unobserved factors

Post-Double Selection Method

Formal Illustration

- Suppose the true model is:

$$Y_i = \beta_0 + \alpha D_i + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i$$

Post-Double Selection Method

Formal Illustration

- **Step 1 of PDS:** Use ML methods (LASSO) to select covariates X_i^j that can predict Y_i
 - Denote the set of LASSO-selected covariates by A

$$Y_i = \delta_0 + \sum_{j=1}^p \delta_j X_i^j + \zeta_i$$

- **Step 2 of PDS:** Use ML methods (LASSO) to select covariates X_i^j that can predict D_i
 - Denote the set of LASSO-selected covariates by B

$$D_i = \gamma_0 + \sum_{j=1}^p \gamma_j X_i^j + \varepsilon_i$$

Post-Double Selection Method

Formal Illustration

- **Step 3 of PDS:** Estimate the following OLS regression:

$$Y_i = \pi_0 + \alpha D_i + \sum_{j=1}^g \pi_j U_i^j + v_i$$

- Let U_i^j denote the **union of covariates in A and B**
- The PDS estimator of treatment effect α is the coefficient on D in the above OLS regression

Post-Double Selection Method

Literature

Belloni, Chernozhukov, Fernandez-Val and Hansen (2013) "**Inference on Treatment Effects after Selection amongst High-Dimensional Controls**", Review of Economic Studies

Belloni, Chernozhukov, Fernandez-Val and Hansen (2014)
"High-dimensional Methods and Inference on Structural and Treatment Effects", Journal of Economic Perspectives

- For more details of PDS method, please read above papers

LASSO Estimation: Details

Shrinkage Methods

Main Idea

- Consider linear regression model with p potential covariates where p is too large.

$$Y_i = \beta_0 + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i, \quad i = 1, \dots, n$$

- Shrinkage estimators minimize **sum of squares error (SSE) with a penalty for model size**
 - This shrinks some of unimportant parameter estimates towards zero

Shrinkage Methods

Main Idea

- Depending on algorithm of penalty for model size, there are two popular methods:
 - Ridge regression
 - **Least Absolute Shrinkage and Selection Operator (LASSO) regression**
- The key assumption is **approximate sparsity**:
 - Some of the β_j coefficients are well-approximated by zero, and the approximation error is sufficiently 'small'

LASSO regression

Least Absolute Shrinkage and Selection Operator

- We can estimate the following LASSO regression:

$$y_i = \sum_{j=1}^p \beta_j x_i^j + \epsilon_i, \quad i = 1, \dots, n$$

- Subject to $\sum_{j=1}^p |\beta_j| \leq s$
- s is a small number

LASSO regression

Least Absolute Shrinkage and Selection Operator

- The LASSO estimator a set of $\hat{\beta}_j$ solves the following optimization problem:

$$(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p) = \min_{\beta_j} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_i^j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- where $\lambda \geq 0$ is a **tuning parameter**
- This is equivalent to minimizing SSE subject to $\sum_{j=1}^p |\beta_j| \leq s$ for some s
- The LASSO estimator:
 - Minimizes the **sum of squared errors (SSE)** with a penalty for model complexity
 - Can set some coefficients exactly to zero, effectively selecting a subset of variables

Adjust Scale of Covariates

- However, the magnitude of parameter estimates β_j are related to **scale of covariates**
- Make parameter estimates β_j invariant to the scale of covariates
- We need to **standardize** the covariates X and outcome variable Y by dividing their standard deviations σ_{Xj} and σ_Y
 - The covariates we use in regression will be transformed to:

$$\bullet x_i^j = \frac{X_i^j - \bar{X}^j}{\sigma_{Xi}}$$

- The outcome variable will be transformed to:

$$\bullet y_i = \frac{Y_i - \bar{Y}}{\sigma_Y}$$

LASSO regression

- There is no closed-form solution for LASSO estimator
- **Intuition:**
 - There is a cost to including lots of regressors
 - $\lambda \sum_{j=1}^p |\beta_j|$
 - We can minimize the objective function by throwing out the ones that contribute little to the fit
 - The effect of the penalization is that LASSO sets the $\hat{\beta}_j$ s for some variables to zero
- Key question: **how to choose λ**

How to select λ

How to select λ

Overview

- The **tuning parameter** λ controls the **strength of penalty** and determines a set of covariates
 - Each tuning parameter value λ corresponds to a fitted regression model
- The shrinkage methods allow us to **simplify the model selection problem to a one-dimensional problem**

Three Ways to Choose λ

Overview

1 Cross-validation approach:

- It is a data-driven approach. Choose **tuning parameter** λ that minimize **prediction error**

2 Rigorous approach:

- It is a theory-driven approach. Belloni et al. (2012, *Econometrica*) develop theory and feasible algorithms for the optimal λ .

3 Information criteria approach:

- Select the value of λ that minimizes information criterion (AIC, AICc, BIC or EBIC).

Cross-validation approach

Methods of Choosing λ : Cross-validation approach

Overview

- Cross-validation is a simple, intuitive way to evaluate model fit (choose λ) based on **prediction error**
 - Single-split validation, K-fold cross validation, leave-one-out cross validation
 - Divide sample into both a **training data (estimation)** and a **validation data (evaluate prediction)**
 - Computationally more expensive

Methods of Choosing λ : Cross-validation approach

Overview

- Consider two types of data sets

1. Training data set

- Used to estimate a regression model
- Get estimated coefficients $\hat{\beta}_j$

2. Validation data set

- Additional data used to determine how good is the regression model fit
- A test observation (x_i^{j+}, y_i^+) is a previously unseen observation

Methods of Choosing λ : Cross-validation approach

Prediction Error

- We use **training data set** to estimate the following regression using LASSO:

$$y_i = \sum_{j=1}^p \beta_j x_i^j + \epsilon_i, \quad i = 1, \dots, n$$

- Choose a set of **LASSO estimator** $\hat{\beta}_j$ that minimize **SSE with a penalty for model size**
 - $\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_i^j)^2 + \lambda \sum_{j=1}^p |\beta_j|$
 - where $\lambda \geq 0$ is a **tuning parameter**
- We can predict y_i using the estimated regression model $\hat{f}_\lambda()$
 - LASSO: $\hat{y}_i = \hat{f}_\lambda(x_i^j) = \sum_{j=1}^p \hat{\beta}_j x_i^j$

Methods of Choosing λ : Cross-validation approach

Prediction Error

- Then, we use this estimated regression model to predict unseen observations y_i^+ in the **validation data set**
- We can evaluate prediction accuracy based on **prediction errors** using **SSE** or **mean of squares error (MSE)**

$$SSE = \sum_{i=1}^n (y_i^+ - \hat{f}_\lambda(x_i^+))^2$$

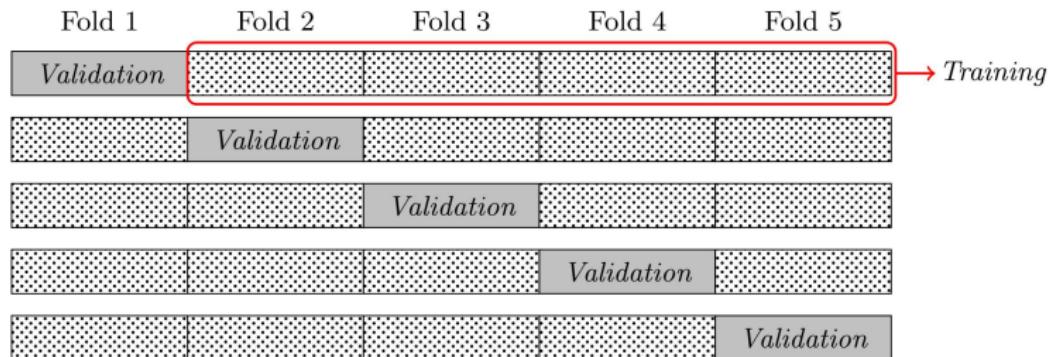
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i^+ - \hat{f}_\lambda(x_i^+))^2$$

- (x_i^+, y_i^+) are previously unseen observations in **validation data set**
- $\hat{f}_\lambda()$ is the estimated regression model based on **training data set**

K-fold Cross Validation

Step 1

- 1 Randomly divide the data set $\{1, \dots, n\}$ into K groups F_1, \dots, F_K of roughly equal size
 - Commonly we split the data into 5 or 10 groups/folds ($K = 5$ or $K = 10$)
 - Consider training on (x_i^j, y_i) , $i \notin F_k$, and validating on (x_i^{j+}, y_i^{+}) , $i \in F_k$



K-fold Cross Validation

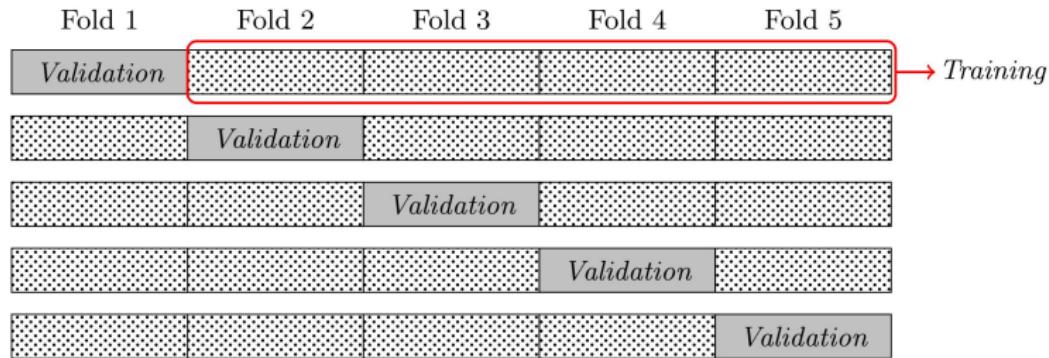
Step 2

- 2 For each value of the tuning parameter $\lambda \in \{\lambda_1, \dots, \lambda_m\}$, we do the followings:
 - Note that each λ has the corresponding regression model so we have m regression models in this case

K-fold Cross Validation

Step 2-1

- 2-1 Compute the regression estimates \widehat{f}_λ^{-k} using the **training data set**

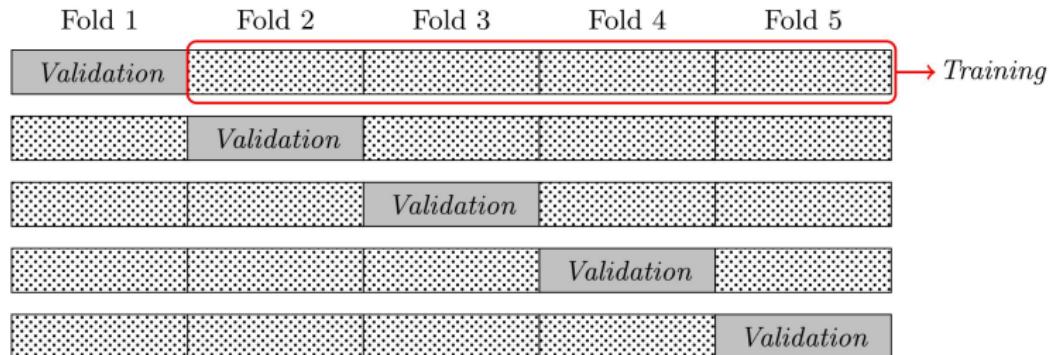


K-fold Cross Validation

Step 2-2

2-2 Compute the prediction error $e_k(\lambda)$ on the each **validation data set**:

$$e_k(\lambda) = \sum_{i \in F_k} (y_i^+ - \hat{f}_\lambda^{-k}(x_i^+))^2$$

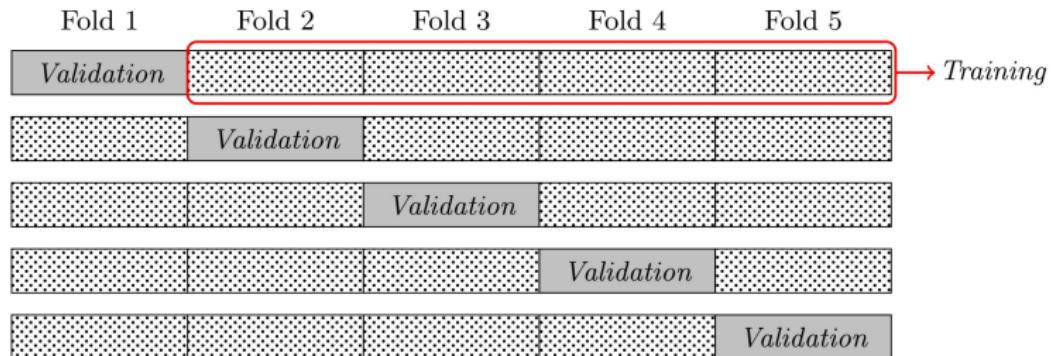


K-fold Cross Validation

Step 2-3

2-3 Then, compute the **average prediction error** over all **validation data sets**

$$CV(\lambda) = \frac{1}{n} \sum_{k=1}^K e_k(\lambda) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in F_k} (y_i^+ - \hat{f}_\lambda^{-k}(x_i^+))^2$$

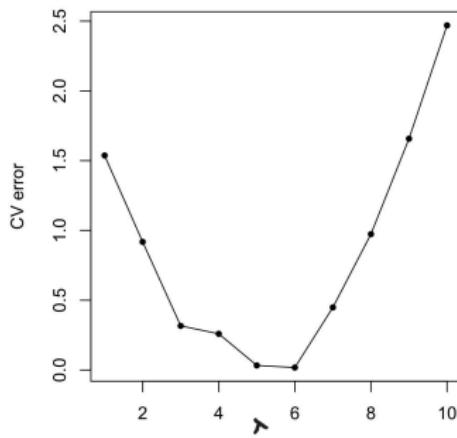


K-fold Cross Validation

Step 2-3

2-3 Having done this, we get a **average prediction error** $CV(\lambda)$

- It is also called **cross-validation error curve**
- This curve is a function of λ



K-fold Cross Validation

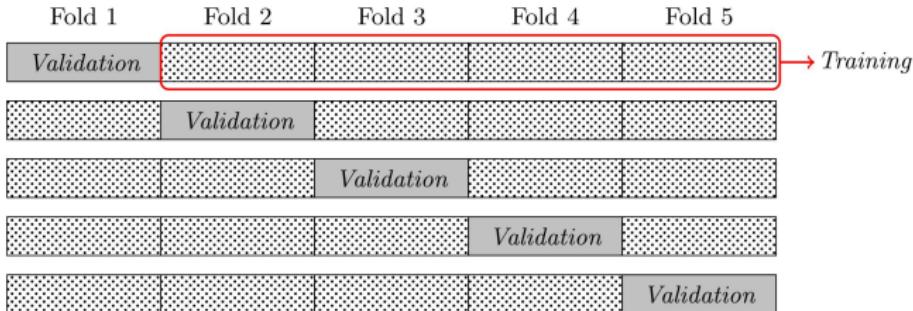
Step 3

- 3 Choose the value of **tuning parameter** λ that minimizes $CV(\lambda)$ curve:

$$\hat{\lambda} = \arg \min_{\lambda \in \{\lambda_1, \dots, \lambda_m\}} CV(\lambda)$$

- When $K = n$, we call this **leave-one-out cross-validation**, because we leave out one data point at a time

Example: 5-fold cross validation



- 1 Pick 1st part of data as the **validation data set** and use other parts of data as the **training data set**
- 2 For each value of the tuning parameter $\lambda \in \{\lambda_1, \dots, \lambda_m\}$, we do the followings:
 - 2-1 Get estimated regression models based on **training data set**
 - 2-2 Use these estimated regression models to predict each **validation data set** and get **prediction errors**
 - 2-3 Do the same thing for $k = 2, \dots, 5$ and get **average prediction errors**
- 3 Choose λ that can minimize **average prediction errors** $CV(\lambda)$

Post-Double Selection Method: STATA Example

STATA Example: Test the Convergence Hypothesis

Overview

Sala-i-Martin, Xavier (1997), "**I Just Ran Two Million Regressions**", American Economic Review

- The author tries to examine the hypothesis of growth convergence and the determinants of economic growth
 - Whether poorer economies' per capita incomes will tend to grow at faster rates than richer economies
- He finds that a substantial number of variables can be found to be strongly related to growth
- Citations: 2,931 times

STATA Example: Test the Convergence Hypothesis

Overview

- Examine the relation between GDP growth rate and initial per capita GDP:

$$\underbrace{\text{GrowthRate}}_{Y_i} = \beta_0 + \underbrace{\alpha}_{\text{ATE}} \underbrace{\log(\text{GDP})}_{D_i} + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i$$

- Control a lot of covariates X_i^j
- In their data, they have $p = 60$ covariates, $n = 90$ observations
 - Need to do variable selection
- Test the convergence hypothesis $\alpha < 0$
 - Poor countries catch up with richer countries, conditional on similar institutions etc.
 - Prediction from the classical Solow growth model.

STATA Example: Test the Convergence Hypothesis

Data and Code

- See **ML.do**
- Use `growth.dta`

STATA Example: Test the Convergence Hypothesis

Examine Data

```
1 codebook  
2 misstable summarize
```

- **Line 1:** codebook is used to describe data contents, which provides information about the variables, such as their labels, value labels, and summary statistics.
- **Line 2:** misstable summarize is used to display a summary of missing values for all variables in the dataset.

STATA Example: Test the Convergence Hypothesis

Examine Data

```
1 duplicates report  
2 duplicates drop
```

- **Line 1:** `duplicates report` is used to report the number of duplicates in the dataset based on all variables.
- **Line 2:** `duplicates drop` is used to remove duplicate observations from the dataset based on all variables.

STATA Example: Test the Convergence Hypothesis

Examine Data

```
1 duplicates report country_id  
2 duplicates report country_id_new  
3 duplicates drop country_id_new,force
```

- **Line 1:** `duplicates report country_id` is used to report the number of duplicates in the dataset based on the variable `country_id`.
- **Line 2:** `duplicates report country_id_new` is used to report the number of duplicates in the dataset based on the variable `country_id_new`.
- **Line 3:** `duplicates drop country_id_new,force` is used to remove duplicate observations from the dataset based on the variable `country_id_new`

STATA Example: Test the Convergence Hypothesis

Control All Covariates

```
1 reg Outcome gdphs465 bmp1l- tot1,r
```

- Control all 60 covariates
- This could result in imprecise estimates of coefficients due to overfitting

STATA Example: Test the Convergence Hypothesis

Control All Covariates

```
. reg Outcome gdphs465 bmp1l- tot1,r
```

Linear regression

	Number of obs	=	90
F(61, 28)	=	8.96	
Prob > F	=	0.0000	
R-squared	=	0.8871	
Root MSE	=	.03074	

Outcome	Robust					
	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
gdphs465	-.009378	.0324211	-0.29	0.775	-.0757896	.0570336
bmp1l	-.0688627	.0408531	-1.69	0.103	-.1525464	.014821
freeop	.080069	.2318888	0.35	0.732	-.3949337	.5550716
freetar	-.4889626	.4124222	-1.19	0.246	-1.333771	.355846
h65	-2.362099	.7377739	-3.20	0.003	-3.87336	-.8508374
hm65	.7071434	.518609	1.36	0.184	-.355179	1.769466
hf65	1.693448	.4580274	3.70	0.001	.7552219	2.631675
p65	.2655267	.1645915	1.61	0.118	-.0716236	.602677
pm65	.1369526	.1310772	1.04	0.305	-.1315469	.4054521
pf65	-.3312669	.1843816	-1.80	0.083	-.7089556	.0464217
s65	.0390793	.1772324	0.22	0.827	-.3239648	.4021234
sm65	-.0306685	.1230053	-0.25	0.805	-.2826334	.2212964
sf65	-.1799173	.1014187	-1.77	0.087	-.387664	.0278294
fert65	.0068808	.0289444	0.24	0.814	-.0524091	.0661708
mort65	-.2334545	.856995	-0.27	0.787	-1.988929	1.52202
life065	-.0149145	.1826972	-0.08	0.936	-.3891527	.3593237
gpop1	.9701846	1.80003	0.54	0.594	-2.71701	4.657379
fert1	.0088382	.0352985	0.25	0.804	-.0634676	.0811439
mort1	.0665629	.7121914	0.09	0.926	-1.392295	1.525421



STATA Example: Test the Convergence Hypothesis

Control All Covariates

- The initial per capita GDP has little impact on economic growth rate
- Does NOT support prediction from the classical Solow growth model

Review: Post-Double Selection Method

Formal Illustration

- Suppose the true model is:

$$\underbrace{Y_i}_{\text{GrowthRate}} = \beta_0 + \underbrace{\alpha}_{\text{ATE}} \underbrace{\log(\text{GDP})}_{D_i} + \sum_{j=1}^p \beta_j X_i^j + \epsilon_i$$

Review: Post-Double Selection Method

Formal Illustration

- **Step 1:** Use ML methods (LASSO) to select covariates X_i^j that can predict Y_i
 - Denote the set of LASSO-selected covariates by A

$$\underbrace{GrowthRate}_{Y_i} = \delta_0 + \sum_{j=1}^p \delta_j X_i^j + \zeta_i$$

- **Step 2:** Use ML methods (LASSO) to select covariates X_i^j that can predict D_i
 - Denote the set of LASSO-selected covariates by B

$$\underbrace{\log(GDP)}_{D_i} = \gamma_0 + \sum_{j=1}^p \gamma_j X_i^j + \varepsilon_i$$

Review: Post-Double Selection Method

Formal Illustration

- **Step 3:** Estimate the following OLS regression:

$$\underbrace{Y_i}_{\text{GrowthRate}} = \pi_0 + \underbrace{\alpha}_{\text{ATE}} D_i + \sum_{j=1}^g \pi_j U_i^j + v_i$$

- Let U_i^j denote the **union of covariates in A and B**
- The PDS estimator of treatment effect α is the coefficient on D in the above OLS regression
- The additional selection step 2 can help eliminate the **omitted variable bias**

Review: Three Ways to Choose λ

Overview

- 1 **Cross-validation approach:** It is a data-driven approach. Choose **tuning parameter** λ that minimize **prediction error**
 - STATA command: **cvllasso**
- 2 **Rigorous approach:** It is a theory-driven approach. Belloni et al. (2012, *Econometrica*) develop theory and feasible algorithms for the optimal λ .
 - STATA command: **rlasso**
- 3 **Information criteria approach:** Select the value of λ that minimizes information criterion (AIC, AICc, BIC or EBIC).
 - STATA command: **lasso2**
 - To use the above commands, you need to install **lassopack** package
 - STATA command: **ssc install lassopack**

STATA Example: Test the Convergence Hypothesis

Use Post-Double Selection Method

```
1 pdlasso Outcome gdph465 (bmp1l- tot1),rob
```

- **pdlasso**: Implement double selection method with LASSO to select covariates
- Select covariates from $bmp1l - tot1$ to predict outcome Y_i and treatment variable D_i

STATA Example: Test the Convergence Hypothesis

Use Post-Double Selection Method

```
. pdlasso Outcome gdph465 (bmp1- tot1),rob  
1. (PDS/CHS) Selecting HD controls for dep var Outcome...  
Selected: bmp1  
2. (PDS/CHS) Selecting HD controls for exog regressor gdph465...  
Selected: freetar hm65 sf65 lifee065 humanf65 pop6565
```

Estimation results:

Specification:
Regularization method: lasso
Penalty loadings: heteroskedastic
Number of observations: 90
Exogenous (1): gdph465
High-dim controls (60):
Selected controls (7):
Unpenalized controls (1):
_cons

bmp1 freeop freetar h65 hm65 hf65 p65 pm65 pf65 s65
sm65 sf65 fert65 mort65 lifee065 gpop1 fert1 mort1
invsh41 geetot1 geerec1 gde1 govwb1 govsh41 gvxde41
high65 highhm65 highf65 highhc65 highcm65 highcf65 human65
humanm65 humanf65 hyrf65 hyrm65 no65 nom65 nof65
pinstab1 pop65 worker65 pop1565 pop6565 sec65 secm65
secf65 secc65 seccm65 seccf65 syr65 syrm65 syrf65
teapri65 teasec65 ex1 im1 xr65 tot1

Structural equation:

OLS using CHS lasso-orthogonalized vars

Outcome	Robust					
	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gdph465	-.0278804	.0164206	-1.70	0.090	-.0600642	.0043035

STATA Example: Test the Convergence Hypothesis

Use Post-Double Selection Method

OLS with PDS-selected variables and full regressor set

Outcome	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]
gdph465	-.0500059	.0150732	-3.32	0.001	-.0795488 -.0204629
bmp11	-.0782423	.0157799	-4.96	0.000	-.1091703 -.0473144
freetar	-.5746764	.2538159	-2.26	0.024	-1.072146 -.0772064
hm65	.0511529	.0538366	0.95	0.342	-.0543649 .1566707
sf65	-.0470218	.0487002	-0.97	0.334	-.1424725 .0484288
lifee065	.2122794	.054255	3.91	0.000	.1059415 .3186173
humanf65	-.000376	.0035354	-0.11	0.915	-.0073052 .0065531
pop6565	.1343893	.2301626	0.58	0.559	-.3167211 .5854996
_cons	-.4064513	.1830995	-2.22	0.026	-.7653198 -.0475828

Standard errors and test statistics valid for the following variables only:
gdph465

STATA Example: Test the Convergence Hypothesis

Use Post-Double Selection Method

- Instead of 60 variables, this method only select 7 variables that are correlated with outcome Y_t or treatment variable D_t
 - Higher initial per capita GDP could lead to lower GDP growth rate in the later years
 - Poor countries do catch up with richer countries
- Support prediction from the classical Solow growth model

STATA Example: Test the Convergence Hypothesis

Use Post-Double Selection Method

```
1 ** step 1:  
2 rlasso Outcome bmp11- tot1,rob  
3  
4 ** step 2:  
5 rlasso gdpsh465 bmp11- tot1,rob  
6  
7 ** step 3:  
8 reg Outcome gdpsh465 bmp11 freetar hm65 sf65 lifee065  
    humanf65 pop6565,r
```

- You can implement PDS method by yourself
- PDS uses rigorous approach to select optimal λ .
 - STATA command: **rlasso**

STATA Example: Visualize Data

Scatter Plot

```
1 graph twoway scatter Outcome gdpsh465, title("GDP  
    Growth Rate (2000) vs. log(GDP) (1965)") ///  
2 xtitle("log(GDP) in 1965") ///  
3 ytitle("GDP Growth Rate in 2000") ///  
4 graphregion(color(white)) xlabel(5(1)10) ylabel  
    (-0.1(0.1)0.2)  
5  
6 graph export "$pic\gdp_scatter.png", replace width  
    (3000)
```

- **pdslasso**: Implement double selection method with LASSO to select covariates
- Select covariates from $bmp1 - tot1$ to predict outcome Y_i and treatment variable D_i

STATA Example: Visualize Data

Scatter Plot

```
1 graph twoway scatter Outcome gdpsh465, title("GDP  
2 Growth Rate (2000) vs. log(GDP) (1965)") ///  
3 xtitle("log(GDP) in 1965") ///  
4 ytitle("GDP Growth Rate in 2000") ///  
5 graphregion(color(white)) xlabel(5(1)10) ylabel  
6 (-0.1(0.1)0.2)  
graph export "$pic\gdp_scatter.png", replace width  
(3000)
```

- **graph twoway scatter**

- Generates a scatter plot with 'Outcome' as the y-axis and 'gdpsh465' as the x-axis.

- **graphregion(color(white))**

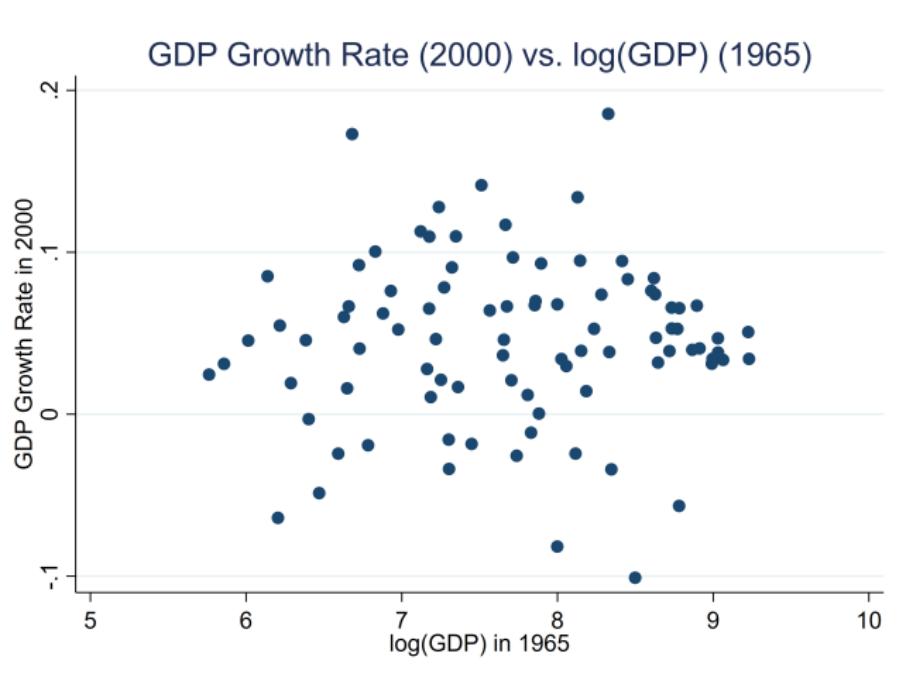
- Sets the background color of the graph to white for a clean look.

- **graph export**

- Exports the graph as a PNG file named

STATA Example: Visualize Data

Scatter Plot



STATA Example: Visualize Data

Scatter Plot: residuals

```
1 * Regress Outcome on control variables and get  
  residuals  
2 regress Outcome bmp11 freetar hm65 sf65 lifee065  
    humanf65 pop6565  
3 predict res_Outcome, residuals  
4  
5 * Regress gdpsht465 on control variables and get  
  residuals  
6 regress gdpsht465 bmp11 freetar hm65 sf65 lifee065  
    humanf65 pop6565  
7 predict res_gdpsht465, residuals
```

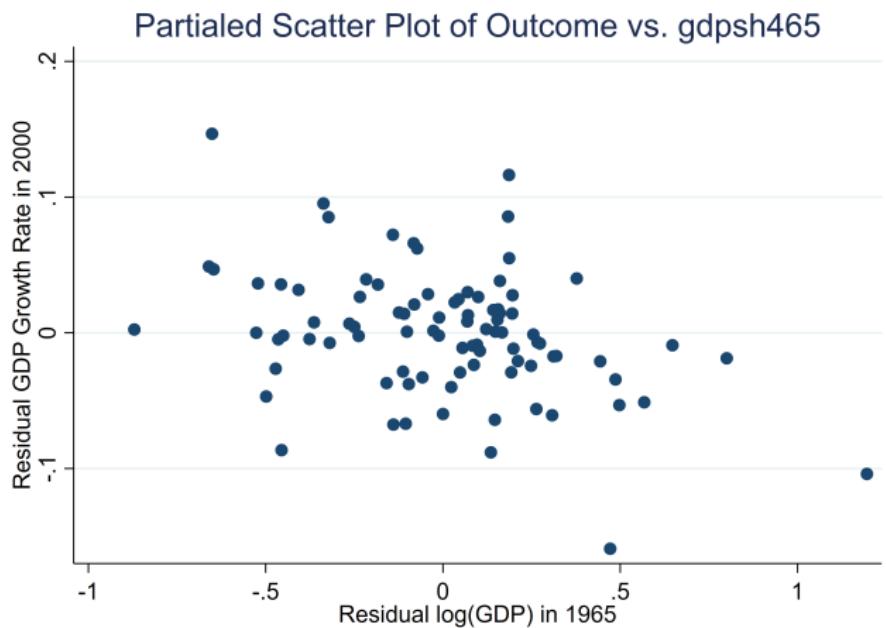
STATA Example: Visualize Data

Scatter Plot: residuals

```
1 * Scatter plot of the residuals
2 scatter res_Outcome res_gdpsh465, ///
3 title("Partialled Scatter Plot of Outcome vs. gdpsh465
") ///
4 xtitle("Residual log(GDP) in 1965") ///
5 ytitle("Residual GDP Growth Rate in 2000") ///
6 graphregion(color(white)) xlabel(-1(0.5)1) ylabel
(-0.1(0.1)0.2)
7
8 graph export "$pic\gdp_scatter_r.png", replace width
(3000)
```

STATA Example: Visualize Data

Scatter Plot: residuals



Post-Double Selection Method: R Example

R Example: Test the Convergence Hypothesis

Data and Code

- See **ML.R**
- Use `growth.dta`
- Install the following packages:
 - `haven`
 - `hdm`

R Example: Test the Convergence Hypothesis

Install Packages

```
1 install.packages("hdm")
2
3 library(hdm)
```

- Install and load package **hdm**
 - This package can help you implement post-double selection method

R Example: Test the Convergence Hypothesis

Read Data

```
1 GrowthData <- read_dta(paste0(rawdata, "/growth.dta"))
   )
2
3 # Checking dimensions of the dataset
4 dim(GrowthData)
5
6 varnames = colnames(GrowthData)
```

- **dim(GrowthData)**: Displays the dimensions of the loaded data, ensuring it is loaded correctly.
- **varnames**: Stores the names of all variables in the dataset for reference.

R Example: Test the Convergence Hypothesis

Create Sample for Analysis

```
1 y = GrowthData[, 1, drop = F]
2 d = GrowthData[, 2, drop = F]
3 X = as.matrix(GrowthData)[, -c(1, 2)]
4 varnames = colnames(GrowthData)
```

- Load dataset and define outcome *y*, treatment variable *d*, and covariates *x*
 - **y**: Dependent variable, representing the first column of the dataset.
 - **d**: Treatment or key independent variable, taken from the second column.
 - **X**: Matrix of covariates, includes all columns except the first two.

R Example: Test the Convergence Hypothesis

PDS Estimation

```
1 doublesel.effect = rlassoEffect(x = X, y = y, d =
2   d, method = "double selection")
3 summary(doublesel.effect)
```

- Implement double selection estimation and report the treatment effect

Recommended Resources for Self-Learning

- Machine Learning & Causal Inference: A Short Course
 - <https://www.gsb.stanford.edu/faculty-research/centers-initiatives/sil/research/methods/ai-machine-learning/short-course>

Recommended Resources for Self-Learning

- NBER Summer Institute 2013: Econometric Methods for High-Dimensional Data
 - https://www.nber.org/econometrics_minicourse_2013/
- One free textbook: An Introduction to Statistical Learning
 - Website:
 - <http://faculty.marshall.usc.edu/gareth-james/>
 - Video lectures:
https://www.youtube.com/playlist?list=PL0g0ngHtcqbPTlZzRHA2ocQZqB1D_qZ5V

Recommended Resources for Self-Learning

- the Stata Lasso Page
 - <https://statalasso.github.io/>