

Introduction to Statistical Machine Learning with Applications in Econometrics

Lecture 13: LASSO for Instrumental Variable Models

Instructor: Ma, Jun

Renmin University of China

December 16, 2021

Instrumental variable

- Consider

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + e_i, \\E[e_i] &= 0 \\Cov[X_i, e_i] &\neq 0.\end{aligned}$$

- An instrument is an variable Z_i which satisfies the following conditions:
 1. The IV is exogenous: $Cov[Z_i, e_i] = 0$.
 2. The IV determines the endogenous regressor: $Cov[Z_i, X_i] \neq 0$.
- When an IV variable satisfying those conditions is available, it allows us to estimate the effect of X on Y consistently:

$$\begin{aligned}Cov[Y_i, Z_i] &= \beta_1 Cov[X_i, Z_i] + Cov[e_i, Z_i] \\&= \beta_1 Cov[X_i, Z_i] \implies \beta_1 = \frac{Cov[Y_i, Z_i]}{Cov[X_i, Z_i]}.\end{aligned}$$

Sources of endogeneity

There are several possible sources of endogeneity:

1. Omitted explanatory variables.
2. Simultaneity.
3. Errors in variables.

All result in regressors correlated with the errors.

Omitted explanatory variables

- Suppose that the true model is

$$\log(Wage_i) = \beta_0 + \beta_1 Education_i + \beta_2 Ability_i + V_i,$$

where V_i is uncorrelated with $Education$ and $Ability$.

- Since $Ability$ is unobservable, the econometrician regresses $\log(Wage)$ against $Education$, and $\beta_2 Ability$ goes into the error part:

$$\begin{aligned}\log(Wage_i) &= \beta_0 + \beta_1 Education_i + U_i, \\ U_i &= \beta_2 Ability_i + V_i.\end{aligned}$$

- $Education$ is correlated with $Ability$: we can expect that $\text{Cov}(Education_i, Ability_i) > 0$, $\beta_2 > 0$, and therefore $\text{Cov}(Education_i, U_i) > 0$.

Simultaneity

- Consider the following demand-supply system:

$$\text{Demand: } Q^d = \beta_0^d + \beta_1^d P + U^d,$$

$$\text{Supply: } Q^s = \beta_0^s + \beta_1^s P + U^s,$$

where: Q^d = quantity demanded, Q^s = quantity supplied,
 P = price.

- The quantity and price are determined simultaneously in the equilibrium:

$$Q^d = Q^s = Q.$$

- Note that Q^d and Q^s are not observed separately, we observe only the equilibrium values Q .

$$\begin{aligned}
Q^d &= \beta_0^d + \beta_1^d P + U^d, \\
Q^s &= \beta_0^s + \beta_1^s P + U^s, \\
Q^d &= Q^s = Q.
\end{aligned}$$

- Solving for P , we obtain

$$0 = (\beta_0^d - \beta_0^s) + (\beta_1^d - \beta_1^s) P + (U^d - U^s),$$

or

$$P = -\frac{\beta_0^d - \beta_0^s}{\beta_1^d - \beta_1^s} - \frac{U^d - U^s}{\beta_1^d - \beta_1^s}.$$

- Thus,

$$\text{Cov}(P, U^d) \neq 0 \text{ and } \text{Cov}(P, U^s) \neq 0.$$

The demand-supply equations cannot be estimated by OLS.

- Consider the following labour supply model for married women:

$$Hours_i = \beta_0 + \beta_1 Children_i + \text{Other Factors} + U_i,$$

where *Hours*=hours of work, *Children*=number of children.

- It is reasonable to assume that women decide simultaneously how much time to devote to career and family.
- Thus, while we may be mainly interested in the effect of family size on labour supply, there is another equation:

$$Children_i = \gamma_0 + \gamma_1 Hours_i + \text{Other Factors} + V_i,$$

and *Children* and *Hours* are determined simultaneously in an equilibrium.

- As a result, $\text{Cov}(Children_i, U_i) \neq 0$, and the effect of family size cannot be estimated by OLS.

Errors in variables

- Consider the following model:

$$Y_i = \beta_0 + \beta_1 X_i^* + V_i,$$

where X_i^* is the true regressor.

- Suppose that X_i^* is not directly observable. Instead, we observe X_i that measures X_i^* with an error ε_i :

$$X_i = X_i^* + \varepsilon_i.$$

- Since X_i^* is unobservable, the econometrician has to regress Y_i against X_i .

$$\begin{aligned}X_i &= X_i^* + \varepsilon_i, \\Y_i &= \beta_0 + \beta_1 X_i^* + V_i.\end{aligned}$$

- The model for Y_i as a function of X_i can be written as

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 (X_i - \varepsilon_i) + V_i \\&= \beta_0 + \beta_1 X_i + V_i - \beta_1 \varepsilon_i,\end{aligned}$$

or

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + e_i, \\e_i &= V_i - \beta_1 \varepsilon_i.\end{aligned}$$

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + e_i, \\ e_i &= V_i - \beta_1 \varepsilon_i, \\ X_i &= X_i^* + \varepsilon_i. \end{aligned}$$

- We can assume that

$$\text{Cov} [X_i^*, V_i] = \text{Cov} [X_i^*, \varepsilon_i] = \text{Cov} [\varepsilon_i, V_i] = 0.$$

- However,

$$\begin{aligned} \text{Cov} [X_i, e_i] &= \text{Cov} [X_i^* + \varepsilon_i, V_i - \beta_1 \varepsilon_i] \\ &= \text{Cov} [X_i^*, V_i] - \beta_1 \text{Cov} [X_i^*, \varepsilon_i] \\ &\quad + \text{Cov} [\varepsilon_i, V_i] - \beta_1 \text{Cov} [\varepsilon_i, \varepsilon_i] \end{aligned}$$

- Thus, X_i is endogenous and β_1 cannot be estimated by OLS.

Example: Compulsory schooling laws and return to education

- ▶ Angrist and Krueger, 1991, *QJE*, suggested using school start age policy to estimate β_1 in
$$\log(Wage_i) = \beta_0 + \beta_1 Education_i + \beta_2 Ability_i + V_i.$$
- ▶ We need to find an IV variable Z such that $Cov(Ability_i, Z_i) = 0$ and $Cov(Education_i, Z_i) \neq 0$.
- ▶ They argue that due to compulsory schooling laws, the season of birth variable satisfies the IV conditions:
 - ▶ A child has to attend the school until he reaches a certain drop-out age.
 - ▶ Students born in the first quarter of the year, reach the legal drop-out age before their classmates who were born later in the year.
 - ▶ The quarter of birth dummy variable is correlated with education.
 - ▶ The quarter of birth is uncorrelated with ability.

Example: Sibling-sex composition and labor supply

- ▶ Angrist and Evans, 1998, *AER*, argue that the parents' preferences for a mixed sibling-sex composition can be used to estimate β_1 in $Hours_i = \beta_0 + \beta_1 Children_i + \dots + U_i$.
- ▶ We need to find an IV Z such that $Cov [U_i, Z_i] = 0$ and $Cov (Children_i, Z_i) \neq 0$.
- ▶ Consider a dummy variable that takes on the value one if the sex of the second child matches the sex of the first child.
 - ▶ If the parents prefer a mixed sibling-sex composition, they are more likely to have another child if their first two children are of the same sex.
 - ▶ The same-sex dummy is correlated with the number of children.
 - ▶ Since sex mix is randomly determined, the same sex dummy is exogenous.

Instrumental variable model

- Consider the following model:

$$Y_i = \gamma_0 + \gamma_1 X_{i1} + \dots + \gamma_k X_{ik} + \beta_1 D_{i1} + \dots + \beta_m D_{im} + U_i,$$

where

- Y_i is the dependent variable.
- γ_0 is the coefficient on the constant regressor: $E[U_i] = 0$.
- X_{i1}, \dots, X_{ik} are the k exogenous regressors:

$$\text{Cov}[X_{i1}, U_i] = \dots = \text{Cov}[X_{ik}, U_i] = 0.$$

- D_{i1}, \dots, D_{im} are the m endogenous regressors:

$$\text{Cov}[D_{i1}, U_i] \neq 0, \dots, \text{Cov}[D_{im}, U_i] \neq 0.$$

- ▶ Suppose that the econometrician observes l additional exogenous variables (IVs) Z_{i1}, \dots, Z_{il}
- ▶ We assume that the IVs Z_{i1}, \dots, Z_{il} are excluded from the structural equation:

$$Y_i = \gamma_0 + \gamma_1 X_{i1} + \dots + \gamma_k X_{ik} + \beta_1 D_{i1} + \dots + \beta_m D_{im} + U_i,$$

so we still have $k + 1 + m$ structural coefficients to estimate.

- ▶ The necessary condition for identification is that the number of IVs is at least as large as the number of unknowns or $l \geq m$.

2SLS

- Consider the first-stage projection models:

$$\begin{aligned} D_{i1} &= \pi_{0,1} + \pi_{1,1}Z_{i1} + \dots + \pi_{l,1}Z_{il} \\ &\quad + \pi_{l+1,1}X_{i1} + \dots + \pi_{l+k,1}X_{ik} + V_{i1}, \\ &\quad \vdots \quad \vdots \quad \vdots \\ D_{im} &= \pi_{0,m} + \pi_{1,m}Z_{i1} + \dots + \pi_{l,m}Z_{il} \\ &\quad + \pi_{l+1,m}X_{i1} + \dots + \pi_{l+k,m}X_{ik} + V_{im}, \end{aligned}$$

where $(\pi_{0,1}, \pi_{1,1}, \dots, \pi_{l+k,1})$ are projection coefficients.

- All right-hand side variables are exogenous.
- The first stage coefficients π 's can be estimated consistently by OLS by regressing Y 's against Z 's and X 's.

- After estimating π 's, obtain the fitted values for D 's:

$$\begin{aligned}\widehat{D}_{i1} &= \widehat{\pi}_{0,1} + \widehat{\pi}_{1,1}Z_{i1} + \dots + \widehat{\pi}_{l,1}Z_{il} \\ &\quad + \widehat{\pi}_{l+1,1}X_{i1} + \dots + \widehat{\pi}_{l+k,1}X_{ik}, \\ &\quad \vdots \quad \vdots \quad \vdots \\ \widehat{D}_{im} &= \widehat{\pi}_{0,m} + \widehat{\pi}_{1,m}Z_{i1} + \dots + \widehat{\pi}_{l,m}Z_{il} \\ &\quad + \widehat{\pi}_{l+1,m}X_{i1} + \dots + \widehat{\pi}_{l+k,m}X_{ik}.\end{aligned}$$

- In the second stage, regress (OLS) the dependent variable Y against a constant, X 's, and \widehat{D} 's obtained in the first stage:

$$Y_i = \widehat{\gamma}_0^{2\text{sls}} + \widehat{\gamma}_1^{2\text{sls}}X_{i1} + \dots + \widehat{\gamma}_k^{2\text{sls}}X_{ik} + \widehat{\beta}_1^{2\text{sls}}\widehat{D}_{i1} + \dots + \widehat{\beta}_m^{2\text{sls}}\widehat{D}_{im} + \widehat{U}_i.$$

- One can show that the resulting 2SLS estimators $\widehat{\gamma}_0^{2\text{sls}}, \widehat{\gamma}_1^{2\text{sls}}, \dots, \widehat{\gamma}_k^{2\text{sls}}, \widehat{\beta}_1^{2\text{sls}}, \dots, \widehat{\beta}_m^{2\text{sls}}$ are consistent and asymptotically normal.

2SLS estimation with many IVs

- ▶ We consider the simple model (0 intercept):

$$Y_i = \alpha D_i + U_i$$

$$E[U_i] = 0$$

$$\text{Cov}[D_i, U_i] \neq 0.$$

- ▶ Suppose that we have l IVs $Z_i \in \mathbb{R}^l$ ($Z_i = (Z_{i1}, Z_{i2}, \dots, Z_{il})^\top$) which satisfies $\text{Cov}[U_i, Z_i] = 0$.
- ▶ The first-stage of 2SLS uses the projection model of D_i on Z_i :

$$D_i = Z_i^\top \pi + V_i$$

$$E[Z_i V_i] = 0$$

$$\pi = \underset{a}{\operatorname{argmin}} E \left[(D_i - Z_i^\top a)^2 \right].$$

- ▶ Then,

$$\begin{aligned} Y_i &= \alpha D_i + U_i \\ D_i &= Z_i^\top \pi + V_i \end{aligned} \implies Y_i = \alpha Z_i^\top \pi + \alpha V_i + U_i.$$

Regression of Y_i on $Z_i^\top \pi$ consistently estimates α .

- ▶ \mathbf{Z} : the $n \times l$ matrix of IVs; $\mathbf{D} = (D_1, D_2, \dots, D_n)^\top$;
 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^\top$; $\mathbf{U} = (U_1, U_2, \dots, U_n)^\top$; $\mathbf{V} = (V_1, V_2, \dots, V_n)^\top$.
- ▶ Since π is unknown, we replace it with $\hat{\pi} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{D}$:

$$\hat{\alpha}^{2\text{sls}} = \frac{\mathbf{D}^\top \mathbf{P}_Z \mathbf{Y}}{\mathbf{D}^\top \mathbf{P}_Z \mathbf{D}} = \alpha + \frac{n^{-1} \mathbf{D}^\top \mathbf{P}_Z \mathbf{U}}{n^{-1} \mathbf{D}^\top \mathbf{P}_Z \mathbf{D}},$$

where $\mathbf{P}_Z = \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$.

- ▶ $n^{-1} \mathbf{D}^\top \mathbf{P}_Z \mathbf{D}$ is less variable when n and l are both large. The bias of $\hat{\alpha}^{2\text{sls}}$ mainly depends on the numerator $n^{-1} \mathbf{D}^\top \mathbf{P}_Z \mathbf{U}$.
- ▶ Suppose that $E[\mathbf{UV}^\top \mid \mathbf{Z}] = \sigma_{UV} \mathbf{I}_n$ and $\mathbf{Z}^\top \mathbf{Z} = \mathbf{I}_l$, then

$$E \left[\frac{1}{n} \mathbf{D}^\top \mathbf{P}_Z \mathbf{U} \mid \mathbf{Z} \right] = \sigma_{UV} \frac{l}{n}.$$

- ▶ When the number of IVs is large and comparable to the sample size n , the bias can be substantial.

- ▶ In the context of a small and fixed number of IVs, adding one more IV reduces the variance of the 2SLS estimator.
- ▶ However, if there are too many IVs, the bias becomes non-negligible and we have to selection a small subset of best IVs out of the long list of potential IVs.
- ▶ Under an alternative asymptotic analysis, when the number of IVs l is assumed to be growing $l = l_n \uparrow \infty$ as $n \uparrow \infty$ such that $l_n/n \rightarrow c > 0$, the 2SLS estimator is inconsistent.
- ▶ LASSO is used for data-driven IV selection.

Optimal instrument

- Suppose that $E[U_i | Z_i] = 0$, then for any function f , $\text{Cov}[f(Z_i), U_i] = 0$ and $\zeta_i = f(Z_i)$ can be used as an IV:

$$E[\zeta_i Y_i] = \alpha E[\zeta_i D_i] \implies \hat{\alpha}^{\text{iv}} = \frac{\sum_{i=1}^n \zeta_i Y_i}{\sum_{i=1}^n \zeta_i D_i}.$$

- Denote $\hat{\mathbf{D}} = \mathbf{P}_Z \mathbf{D}$.

$$\hat{\alpha}^{\text{2sls}} = \frac{\mathbf{D}^\top \mathbf{P}_Z \mathbf{Y}}{\mathbf{D}^\top \mathbf{P}_Z \mathbf{D}} = \frac{\hat{\mathbf{D}}^\top \mathbf{Y}}{\hat{\mathbf{D}}^\top \mathbf{D}} = \frac{\sum_{i=1}^n \hat{D}_i Y_i}{\sum_{i=1}^n \hat{D}_i D_i},$$

where $\hat{D}_i = Z_i^\top \hat{\pi}$ and $\hat{\pi}$ are the first-stage OLS coefficients.

- $\hat{\alpha}^{\text{2sls}}$ can be viewed as an IV estimator using estimated projection $Z_i^\top \hat{\pi}$ in lieu of the unknown true projection $Z_i^\top \pi_i$ as the instrument.
- $\hat{\alpha}^{\text{2sls}}$ summarizes the information in all instruments Z_i and uses a single IV $Z_i^\top \pi$.

- Assume that the model is homoskedastic: $E[U_i^2 | D_i] = \sigma^2$. We can show that the optimal IV estimator is the one $\hat{\alpha}^*$ that uses $\zeta_i^* = E[D_i | Z_i]$:

$$\sqrt{n}(\hat{\alpha}^* - \alpha) \rightarrow_d N\left(0, \frac{\sigma^2}{E[(\zeta_i^*)^2]}\right)$$

and

$$\sqrt{n}(\hat{\alpha}^{\text{iv}} - \alpha) \rightarrow_d N\left(0, \frac{\sigma^2 E[\zeta_i^2]}{(E[\zeta_i \zeta_i^*])^2}\right).$$

Approximation to the optimal instrument

- The 2SLS uses a linear projection $Z_i^\top \pi$ to approximate $E[D_i | Z_i]$:

$$\begin{aligned}\pi &= \underset{a}{\operatorname{argmin}} E \left[(D_i - Z_i^\top a)^2 \right] \\ &= \underset{a}{\operatorname{argmin}} E \left[(E[D_i | Z_i] - Z_i^\top a)^2 \right].\end{aligned}$$

- We generate a dictionary $W_i = (W_{i1}, \dots, W_{ip})^\top \in \mathbb{R}^p$:

$$W_i = \left(Z_{i1}, Z_{i2}, \dots, Z_{il}, Z_{i1}^2, Z_{i1}Z_{i2}, \dots, Z_{i1}Z_{il}, Z_{i2}^2, \dots \right),$$

whose dimension p can be larger than n .

- We can also use the linear projection $W_i^\top \delta$ to approximate $E[D_i | Z_i]$, where

$$\begin{aligned}\delta &= \underset{b}{\operatorname{argmin}} E \left[(D_i - W_i^\top b)^2 \right] \\ &= \underset{a}{\operatorname{argmin}} E \left[(E[D_i | Z_i] - W_i^\top b)^2 \right].\end{aligned}$$

- It is easy to show that

$$\mathbb{E} \left[(\mathbb{E} [D_i | Z_i] - W_i^\top \delta)^2 \right] < \mathbb{E} \left[(\mathbb{E} [D_i | Z_i] - Z_i^\top \pi)^2 \right].$$

- We assume that if p is very large, then the approximation error is very much close to zero and $W_i^\top \delta$ is the optimal instrument.
- If $p < n$, we can regress D_i on W_i to get the OLS coefficient $\hat{\delta}$ and uses the estimated optimal instrument $W_i^\top \hat{\delta}$.
- This procedure is equivalent to 2SLS using all variables in W_i as instruments. We showed that when p is large, the 2SLS estimator may be substantially biased.
- When $p > n$, the 2SLS estimator is not computable. We are forced to select a subset from W_i .

- We assume that the conditional expectation model

$$D_i = W_i^\top \delta + V_i = \sum_{j=1}^p \delta_j W_{ij} + V_i$$

$$E[V_i | Z_i] = 0$$

is sparse: $l^* = |\mathcal{A}|$ is a small number, where $|\mathcal{A}|$ denotes the number of elements in $\mathcal{A} = \{j : \delta_j \neq 0\}$ ($\delta = (\delta_1, \delta_2, \dots, \delta_p)^\top$), although p is very large.

- The IVs in \mathcal{A} are called the effective IVs. Dropping ineffective IVs would not result in loss of efficiency.
- Clearly, there is no difference between the IV estimator using $W_{i,\mathcal{A}}^\top \delta_{\mathcal{A}}$ ($W_{i,\mathcal{A}} = \{W_{ij} : j \in \mathcal{A}\}$ and $\delta_{\mathcal{A}} = \{\delta_j : j \in \mathcal{A}\}$) and the IV estimator using $W_i^\top \delta$.
- However, we do not know \mathcal{A} (identities of the effective IVs). We use LASSO selection to find them.

Algorithm

1. LASSO regression of D_i against W_i :

$$\left(\widehat{\delta}_{1,\lambda}, \dots, \widehat{\delta}_{p,\lambda}\right) = \underset{b_1, \dots, b_p}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(D_i - \sum_{j=1}^p b_j W_{ij} \right)^2 + \lambda \sum_{j=1}^p |b_j| \right\}.$$

Let $\widehat{\mathcal{A}} = \left\{ j : \widehat{\delta}_{j,\lambda} \neq 0 \right\}$ be the selected controls.

- The dropped IVs are either ineffective or have small coefficients ($\delta_j \propto n^{-1/2}$). In the latter case, it can be shown that such variables do not contribute to the asymptotic variance, so we can drop them without loss of efficiency.

2. Post-LASSO of D_i against $W_{i,\widehat{\mathcal{A}}}$ and get the OLS coefficients

$\left\{ \widehat{\delta}_j^{\text{pl}} : j \in \widehat{\mathcal{A}} \right\}$. Generate the fitted value as the estimated optimal IV: $\widehat{\zeta}_i^* = \sum_{j \in \widehat{\mathcal{A}}} \widehat{\delta}_j^{\text{pl}} W_{ij}$.

3. Estimate α using $\widehat{\zeta}_i^*$ as the IV:

$$\widehat{\alpha}^* \left(\widehat{\mathcal{A}} \right) = \frac{\sum_{i=1}^n \widehat{\zeta}_i^* Y_i}{\sum_{i=1}^n \widehat{\zeta}_i^* D_i}.$$

Model with controls

- The structural model with controls $X_i = (X_{i1}, X_{i2}, \dots, X_{ik})^\top$:

$$\begin{aligned} Y_i &= \alpha D_i + X_i^\top \beta + U_i \\ \text{E}[U_i \mid X_i, Z_i] &= 0. \end{aligned}$$

- The intercept is typically one of the elements in X_i .
- Controls X_i have to be included in the first stage. Consider 2SLS and the following projection models:

$$\begin{aligned} D_i &= Z_i^\top \pi + X_i^\top \gamma + V_i \\ \text{E}\left[V_i \begin{pmatrix} Z_i \\ X_i \end{pmatrix}\right] &= 0 \end{aligned}$$

and

$$\begin{aligned} D_i &= Z_i^\top \tilde{\pi} + \tilde{V}_i \\ \text{E}[\tilde{V}_i Z_i] &= 0. \end{aligned}$$

- It is easy to show that $\tilde{\pi} = \Theta\gamma$, where

$$\Theta = (E[Z_i Z_i^\top])^{-1} E[Z_i X_i^\top]$$

$$\tilde{V}_i = D_i - Z_i^\top \tilde{\pi} = V_i + (X_i^\top - Z_i^\top \Theta) \gamma.$$

\tilde{V}_i is not correlated with Z_i but it is correlated with X_i .

- If we drop X_i from the first stage,

$$Y_i = \alpha D_i + X_i^\top \beta + U_i$$

$$D_i = Z_i^\top \tilde{\pi} + \tilde{V}_i$$

$$\implies Y_i = \alpha (Z_i^\top \tilde{\pi} + \tilde{V}_i) + X_i^\top \beta + U_i = \alpha (Z_i^\top \tilde{\pi}) + X_i^\top \beta + \alpha \tilde{V}_i + U_i.$$

The residual $\alpha \tilde{V}_i + U_i$ is correlated with X_i .

- Regression of Y_i against $Z_i^\top \tilde{\pi}$ and X_i does not give consistent estimator for α .

- The 2SLS can be written as an IV estimator:

$$\begin{aligned} \begin{pmatrix} \widehat{\alpha}^{2\text{sls}} \\ \widehat{\beta}^{2\text{sls}} \end{pmatrix} &= \left(\sum_{i=1}^n \begin{pmatrix} \widehat{D}_i \\ X_i \end{pmatrix} \begin{pmatrix} D_i \\ X_i \end{pmatrix}^\top \right)^{-1} \left(\sum_{i=1}^n \begin{pmatrix} \widehat{D}_i \\ X_i \end{pmatrix} Y_i \right), \\ &= \left(\sum_{i=1}^n \begin{pmatrix} \widehat{D}_i \\ X_i \end{pmatrix} \begin{pmatrix} D_i \\ X_i \end{pmatrix}^\top \right)^{-1} \left(\sum_{i=1}^n \begin{pmatrix} \widehat{D}_i \\ X_i \end{pmatrix} Y_i \right), \end{aligned}$$

where $\widehat{D}_i = Z_i^\top \widehat{\pi} + X_i^\top \widehat{\gamma}$ denotes the first-stage fitted value.

- The optimal IV: $\zeta_i^* = E[D_i | X_i, Z_i]$ and the optimal IV estimator:

$$\begin{pmatrix} \widehat{\alpha}^* \\ \widehat{\beta}^* \end{pmatrix} = \left(\sum_{i=1}^n \begin{pmatrix} \zeta_i^* \\ X_i \end{pmatrix} \begin{pmatrix} D_i \\ X_i \end{pmatrix}^\top \right)^{-1} \left(\sum_{i=1}^n \begin{pmatrix} \zeta_i^* \\ X_i \end{pmatrix} Y_i \right).$$

- We need to approximate ζ_i^* .

Many IVs and few controls

- ▶ The conditional expectation model for D_i :

$$\begin{aligned} D_i &= W_i^\top \delta + X_i^\top \gamma + V_i \\ E[V_i \mid X_i, Z_i] &= 0, \end{aligned}$$

where the dictionary W_i contains many polynomials of Z_i and interactions between Z_i and X_i . In this case, we need selection over W_i but need no selection over the controls X_i .

- ▶ In the first-stage regression, we force inclusion of X_i by assigning no penalty weights to their coefficients.
- ▶ In the second stage, we run IV regression by using the post-LASSO fitted value as the IV.

Algorithm

1. LASSO regression of D_i against W_i and X_i :

$$\begin{aligned} & \left(\widehat{\delta}_{1,\lambda}, \dots, \widehat{\delta}_{p,\lambda}, \widehat{\gamma}_{1,\lambda}, \dots, \widehat{\gamma}_{k,\lambda} \right) \\ &= \underset{b_1, \dots, b_p, d_1, \dots, d_k}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(D_i - \sum_{j=1}^p b_j W_{ij} - \sum_{j=1}^k d_j X_{ij} \right)^2 + \lambda \sum_{j=1}^p |b_j| \right\}. \end{aligned}$$

Let $\widehat{\mathcal{A}} = \{j : \widehat{\delta}_{j,\lambda} \neq 0\}$ be the selected controls.

2. Run post LASSO of D_i against the instruments in $W_{i,\widehat{\mathcal{A}}} = \{W_{ij} : j \in \widehat{\mathcal{A}}\}$ and X_i to get OLS coefficients $\{\widehat{\delta}_j^{\text{pl}} : j \in \widehat{\mathcal{A}}\} \cup \{\widehat{\gamma}_1^{\text{pl}}, \dots, \widehat{\gamma}_k^{\text{pl}}\}$. Construct

$$\widehat{\zeta}_i^* = \sum_{j=1}^p \widehat{\delta}_j^{\text{pl}} W_{ij} + \sum_{j=1}^k \widehat{\gamma}_j^{\text{pl}} X_{ij}.$$

3. Estimate (α, β) by using $\widehat{\zeta}_i^*$ as the IV:

$$\begin{pmatrix} \widehat{\alpha}^* \\ \widehat{\beta}^* \end{pmatrix} = \left(\sum_{i=1}^n \begin{pmatrix} \widehat{\zeta}_i^* \\ X_i \end{pmatrix} \begin{pmatrix} D_i \\ X_i \end{pmatrix}^\top \right)^{-1} \left(\sum_{i=1}^n \begin{pmatrix} \widehat{\zeta}_i^* \\ X_i \end{pmatrix} Y_i \right).$$

Partialling out

- ▶ Let $\mathbf{M}_X = \mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ be the projection matrix on the space that is orthogonal to the column space of \mathbf{X} : $\mathbf{M}_X \mathbf{X} = \mathbf{0}$.
- ▶ Write

$$\begin{aligned}Y_i &= \alpha D_i + X_i^\top \beta + U_i \\D_i &= \mathbf{W}_i^\top \delta + X_i^\top \gamma + V_i\end{aligned}$$

in the matrix form

$$\begin{aligned}\mathbf{Y} &= \alpha \mathbf{D} + \mathbf{X}\beta + \mathbf{U} \\ \mathbf{D} &= \mathbf{W}\delta + \mathbf{X}\gamma + \mathbf{V}.\end{aligned}$$

- ▶ Multiply both sides by \mathbf{M}_X to get

$$\begin{aligned}\widetilde{\mathbf{Y}} &= \alpha \widetilde{\mathbf{D}} + \widetilde{\mathbf{U}} \\ \widetilde{\mathbf{D}} &= \widetilde{\mathbf{W}}\delta + \widetilde{\mathbf{V}},\end{aligned}$$

where $\widetilde{\mathbf{Y}} = \mathbf{M}_X \mathbf{Y}$, $\widetilde{\mathbf{D}} = \mathbf{M}_X \mathbf{D}$, $\widetilde{\mathbf{W}} = \mathbf{M}_X \mathbf{W}$, $\widetilde{\mathbf{U}} = \mathbf{M}_X \mathbf{U}$ and $\widetilde{\mathbf{V}} = \mathbf{M}_X \mathbf{V}$.

- ▶ By transforming $(\mathbf{Y}, \mathbf{D}, \mathbf{W})$ into the residuals against \mathbf{X} , we have another numerically equivalent way to compute the IV estimator.

The partialling out algorithm

1. Run LASSO regression of $\tilde{\mathbf{D}} = (\tilde{D}_1, \dots, \tilde{D}_n)^\top$ against $\tilde{\mathbf{W}}$ (\tilde{W}_{ij} denotes its ij -th element of the $n \times p$ matrix $\tilde{\mathbf{W}}$) to get $(\hat{\delta}_{1,\lambda}, \dots, \hat{\delta}_{p,\lambda})^\top$. Let $\hat{\mathcal{A}} = \{j : \hat{\delta}_{j,\lambda} \neq 0\}$ be the selected controls.
2. Run post LASSO regression of $\tilde{\mathbf{D}}$ against the IVs in $\hat{\mathcal{A}}$ to get OLS coefficients $\{\hat{\delta}_j^{\text{pl}} : j \in \hat{\mathcal{A}}\}$. Construct the estimated optimal IV $\hat{\zeta}_i^* = \sum_{j \in \hat{\mathcal{A}}} \hat{\delta}_j^{\text{pl}} \tilde{W}_{ij}$.
3. Estimate α by using $\hat{\zeta}_i^*$ as the IV:

$$\hat{\alpha}^* = \frac{\sum_{i=1}^n \hat{\zeta}_i^* \tilde{Y}_i}{\sum_{i=1}^n \hat{\zeta}_i^* \tilde{D}_i},$$

where $\tilde{\mathbf{Y}} = (\tilde{Y}_1, \dots, \tilde{Y}_n)^\top$.

Few IVs and many controls

- In the model

$$Y_i = \alpha D_i + X_i^\top \beta + U_i$$

$$E[U_i \mid X_i, Z_i] = 0$$

$$D_i = Z_i^\top \pi + X_i^\top \gamma + V_i$$

$$E[V_i \mid X_i, Z_i] = 0,$$

the dimension of X_i is large but the number of IVs is small and we do not use its polynomials and interactions to approximate the optimal IV $E[D_i \mid X_i, Z_i]$.

- This is the case, for example, when there is only one binary instrument (a dummy variable, all polynomials are equal) and we do not use its interactions with X_i .
- We assume the outcome equation is sparse: the set of relevant controls $\mathcal{A} = \{j : \beta_j \neq 0\}$ is small.
- In this case, we need to perform LASSO selection over X_i only.
- We apply the partialling out approach by using LASSO and post LASSO over over X_i .

The partialling out algorithm

1. Perform LASSO and post LASSO of D_i against X_i to generate the residual $\widetilde{D}_i^{\text{pl}}$.
2. Perform LASSO and post LASSO of Y_i against X_i to generate the residual $\widetilde{Y}_i^{\text{pl}}$.
3. Perform LASSO and post LASSO of Z_{ij} against X_i to generate the residual $\widetilde{Z}_{ij}^{\text{pl}}$, for $j = 1, 2, \dots, l$.
4. Run OLS of $\widetilde{D}_i^{\text{pl}}$ against $\widetilde{Z}_{i1}^{\text{pl}}, \dots, \widetilde{Z}_{il}^{\text{pl}}$ to get the OLS coefficients $(\widehat{\pi}_1, \dots, \widehat{\pi}_l)$ and the estimated optimal IV $\widehat{\zeta}_i^* = \sum_{j=1}^l \widehat{\pi}_j \widetilde{Z}_{ij}$.
5. Estimate α by

$$\widehat{\alpha}^* = \frac{\sum_{i=1}^n \widehat{\zeta}_i^* \widetilde{Y}_i^{\text{pl}}}{\sum_{i=1}^n \widehat{\zeta}_i^* \widetilde{D}_i^{\text{pl}}}.$$

Many IVs and many controls

- In the model

$$Y_i = \alpha D_i + X_i^\top \beta + U_i$$

$$E[U_i \mid X_i, Z_i] = 0$$

$$D_i = W_i^\top \delta + X_i^\top \gamma + V_i$$

$$E[V_i \mid X_i, Z_i] = 0,$$

where the dictionary W_i contains high-dimensional transformations (polynomials) of the primitive instruments Z_i and interactions of Z_i and X_i .

- Both W_i and X_i are high-dimensional. We need to perform LASSO selections over both.
- The previous partialling out procedure is not practically implementable, since the LASSO and post-LASSO partialling out of many controls from many IVs is computationally hard.
- We do LASSO selection on W_i first to estimate the equation $D_i = W_i^\top \delta + X_i^\top \gamma + V_i$ and find the effective IVs. We then partial out the effects from X_i .

The partialling out algorithm

1. Perform LASSO and post LASSO of D_i against X_i to generate the residual $\widetilde{D}_i^{\text{pl}}$.
2. Perform LASSO and post LASSO of Y_i against X_i to generate the residual $\widetilde{Y}_i^{\text{pl}}$.
3. Perform LASSO and post LASSO of D_i against W_i and X_i to generate the fitted value $\widehat{\zeta}_i^*$ (estimated optimal IV). This step selects a subset from W_i .
4. Perform LASSO and post LASSO of $\widehat{\zeta}_i^*$ against X_i to partial out the effect from X_i and get the residual $\widetilde{\zeta}_i^{\text{pl}}$.
5. Estimate α by

$$\widehat{\alpha}^* = \frac{\sum_{i=1}^n \widetilde{\zeta}_i^{\text{pl}} \widetilde{Y}_i^{\text{pl}}}{\sum_{i=1}^n \widetilde{\zeta}_i^{\text{pl}} \widetilde{D}_i^{\text{pl}}}.$$