

Instrumental-Variable

Ruben Zamar
Department of Statistics UBC

March 12, 2014

THE LINEAR REGRESSION MODEL

- Ordinary linear regression model

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

THE LINEAR REGRESSION MODEL

- Ordinary linear regression model

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- Design Matrix

$$\mathbf{X}_c = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2p} - \bar{x}_p \\ x_{31} - \bar{x}_1 & x_{32} - \bar{x}_2 & \cdots & x_{3p} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{np} - \bar{x}_p \end{pmatrix}$$

THE LS ESTIMATOR

Ordinary Least Squares Estimate

$$\hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} \quad , \quad \hat{\alpha}_{LS} = \bar{y} - \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}_{LS}$$

THE LS ESTIMATOR

Ordinary Least Squares Estimate

$$\hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} \quad , \quad \hat{\alpha}_{LS} = \bar{y} - \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}_{LS}$$

$$\hat{\boldsymbol{\Sigma}}_{xx} = \frac{1}{n} \mathbf{X}'_c \mathbf{X}_c \quad (\text{sample covariance matrix})$$

$$\hat{\boldsymbol{\Sigma}}_{xy} = \frac{1}{n} \mathbf{X}'_c \mathbf{y}$$

THE LS ESTIMATOR

Ordinary Least Squares Estimate

$$\hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} \quad , \quad \hat{\alpha}_{LS} = \bar{y} - \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}_{LS}$$

$$\hat{\boldsymbol{\Sigma}}_{xx} = \frac{1}{n} \mathbf{X}'_c \mathbf{X}_c \quad (\text{sample covariance matrix})$$

$$\hat{\boldsymbol{\Sigma}}_{xy} = \frac{1}{n} \mathbf{X}'_c \mathbf{y}$$

$$\hat{\boldsymbol{\beta}}_{LS} = \hat{\boldsymbol{\Sigma}}_{xx}^{-1} \hat{\boldsymbol{\Sigma}}_{xy}$$

KEY ASSUMPTIONS

- All covariates are **exogenous**

X and ε are independent

KEY ASSUMPTIONS

- All covariates are **exogenous**

X and ε are independent

- Errors have **zero mean**

$$E(\varepsilon) = \mathbf{0}$$

LS IS UNBIASED

$$E\{\hat{\beta}_{LS}\} = E\{E\{\hat{\beta}_{LS}|X\}\}$$

LS IS UNBIASED

$$\begin{aligned} E\{\widehat{\boldsymbol{\beta}}_{LS}\} &= E\{E\{\widehat{\boldsymbol{\beta}}_{LS}|X\}\} \\ &= E\{(X_c'X_c)^{-1}X_c' E\{\mathbf{y}|X\}\} \end{aligned}$$

LS IS UNBIASED

$$\begin{aligned} E \left\{ \widehat{\boldsymbol{\beta}}_{LS} \right\} &= E \left\{ E \left\{ \widehat{\boldsymbol{\beta}}_{LS} | X \right\} \right\} \\ &= E \left\{ (X_c' X_c)^{-1} X_c' E \{ \mathbf{y} | X \} \right\} \\ &= E \left\{ (X_c' X_c)^{-1} X_c' E \{ \mathbf{1}\alpha + X\boldsymbol{\beta} + \boldsymbol{\varepsilon} | X \} \right\} \end{aligned}$$

LS IS UNBIASED

$$\begin{aligned} E \left\{ \widehat{\boldsymbol{\beta}}_{LS} \right\} &= E \left\{ E \left\{ \widehat{\boldsymbol{\beta}}_{LS} | \mathbf{X} \right\} \right\} \\ &= E \left\{ (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c E \left\{ \mathbf{y} | \mathbf{X} \right\} \right\} \\ &= E \left\{ (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c E \left\{ \mathbf{1}\alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} | \mathbf{X} \right\} \right\} \\ &= E \left\{ \underbrace{(\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{1}\alpha}_{=0} + \underbrace{(\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{X}\boldsymbol{\beta}}_{=I} + (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \underbrace{E \left\{ \boldsymbol{\varepsilon} \right\}}_{=0} \right\} \end{aligned}$$

LS IS UNBIASED

$$\begin{aligned} E \left\{ \widehat{\boldsymbol{\beta}}_{LS} \right\} &= E \left\{ E \left\{ \widehat{\boldsymbol{\beta}}_{LS} | X \right\} \right\} \\ &= E \left\{ (X_c' X_c)^{-1} X_c' E \{ \mathbf{y} | X \} \right\} \\ &= E \left\{ (X_c' X_c)^{-1} X_c' E \{ \mathbf{1}\alpha + X\boldsymbol{\beta} + \boldsymbol{\varepsilon} | X \} \right\} \\ &= E \left\{ \overbrace{(X_c' X_c)^{-1} X_c' \mathbf{1}\alpha}^{=0} + \overbrace{(X_c' X_c)^{-1} X_c' X\boldsymbol{\beta}}^{=I} + (X_c' X_c)^{-1} X_c' \overbrace{E \{ \boldsymbol{\varepsilon} \}}^{=0} \right\} \\ &= \boldsymbol{\beta} \end{aligned}$$

LS IS CONSISTENT

$$\hat{\beta}_{LS} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \longrightarrow \Sigma_{xx}^{-1} \Sigma_{xy}, \quad \text{by LLN}$$

LS IS CONSISTENT

$$\hat{\beta}_{LS} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \longrightarrow \Sigma_{xx}^{-1} \Sigma_{xy}, \quad \text{by LLN}$$

$$\Sigma_{xx}^{-1} \Sigma_{xy} = \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, y)$$

LS IS CONSISTENT

$$\hat{\beta}_{LS} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \longrightarrow \Sigma_{xx}^{-1} \Sigma_{xy}, \quad \text{by LLN}$$

$$\begin{aligned} \Sigma_{xx}^{-1} \Sigma_{xy} &= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, y) \\ &= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \alpha + \mathbf{x}'\boldsymbol{\beta} + \varepsilon) \end{aligned}$$

LS IS CONSISTENT

$$\hat{\beta}_{LS} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \longrightarrow \Sigma_{xx}^{-1} \Sigma_{xy}, \quad \text{by LLN}$$

$$\Sigma_{xx}^{-1} \Sigma_{xy} = \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, y)$$

$$= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \alpha + \mathbf{x}'\beta + \varepsilon)$$

$$= \Sigma_{xx}^{-1} \overbrace{\text{Cov}(\mathbf{x}, \mathbf{x}'\beta)}^{\Sigma_{xx}\beta} + \Sigma_{xx}^{-1} \overbrace{\text{Cov}(\mathbf{x}, \varepsilon)}^0$$

LS IS CONSISTENT

$$\hat{\beta}_{LS} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \longrightarrow \Sigma_{xx}^{-1} \Sigma_{xy}, \quad \text{by LLN}$$

$$\Sigma_{xx}^{-1} \Sigma_{xy} = \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, y)$$

$$= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \alpha + \mathbf{x}'\beta + \varepsilon)$$

$$= \Sigma_{xx}^{-1} \overbrace{\text{Cov}(\mathbf{x}, \mathbf{x}'\beta)}^{\Sigma_{xx}\beta} + \Sigma_{xx}^{-1} \overbrace{\text{Cov}(\mathbf{x}, \varepsilon)}^0$$

$$= \beta$$

ENDOGENEITY

- Some covariates (called **endogenous**) are correlated with the error term

ENDOGENEITY

- Some covariates (called **endogenous**) are correlated with the error term
- This problem arises in several contexts:

ENDOGENEITY

- Some covariates (called **endogenous**) are correlated with the error term
- This problem arises in several contexts:

Errors-in-variables: Some covariates are measured with errors

ENDOGENEITY

- Some covariates (called **endogenous**) are correlated with the error term
- This problem arises in several contexts:

Errors-in-variables: Some covariates are measured with errors

Omitted covariates: Some model covariates are correlated with unobserved predictors

ENDOGENEITY

- Some covariates (called **endogenous**) are correlated with the error term
- This problem arises in several contexts:

Errors-in-variables: Some covariates are measured with errors

Omitted covariates: Some model covariates are correlated with unobserved predictors

Simultaneity: Some covariates simultaneously affect and are affected by the response variable

- The (additive) asymptotic bias due to endogeneity is given by

$$AB = \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \varepsilon)$$

- The (additive) asymptotic bias due to endogeneity is given by

$$\begin{aligned} AB &= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \varepsilon) \\ &= \frac{\text{Cov}(x, \varepsilon)}{\text{Var}(x)} \quad (\text{single linear regression}) \end{aligned}$$

- The (additive) asymptotic bias due to endogeneity is given by

$$\begin{aligned} AB &= \Sigma_{xx}^{-1} \text{Cov}(\mathbf{x}, \varepsilon) \\ &= \frac{\text{Cov}(x, \varepsilon)}{\text{Var}(x)} \quad (\text{single linear regression}) \\ &= \text{Corr}(x, \varepsilon) \frac{sd(\varepsilon)}{sd(x)} \quad (\text{single linear regression}) \end{aligned}$$

ERRORS IN VARIABLES (EIV)

A simple example:

$$y_i = 1 + 2v_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

$$x_i = v_i + e_i \quad (\text{error-in-variable})$$

ERRORS IN VARIABLES (EIV)

A simple example:

$$\begin{aligned}y_i &= 1 + 2v_i + \varepsilon_i, & i = 1, 2, \dots, n \\x_i &= v_i + e_i \quad (\text{error-in-variable})\end{aligned}$$

To fix ideas suppose that

v_i, ε_i and e_i are independent

$$v_i \sim N(0, 1), \quad \varepsilon_i, e_i \sim N(0, 0.5)$$

$$\begin{aligned}y_i &= 1 + 2v_i + \varepsilon_i \\ &= 1 + 2(x_i - e_i) + \varepsilon_i \\ &= 1 + 2x_i + (\varepsilon_i - 2e_i) \\ &= 1 + 2x_i + \Delta_i\end{aligned}$$

$$\Delta_i = \varepsilon_i - 2e_i$$

$$\begin{aligned}y_i &= 1 + 2v_i + \varepsilon_i \\ &= 1 + 2(x_i - e_i) + \varepsilon_i \\ &= 1 + 2x_i + (\varepsilon_i - 2e_i) \\ &= 1 + 2x_i + \Delta_i\end{aligned}$$

$$\Delta_i = \varepsilon_i - 2e_i$$

$$\text{cov}(x_i, \Delta_i) = -2 \times 0.5 = -1 \quad (\text{endogeneity})$$

$$\text{Var}(x_i) = 1.5$$

EIV (continued)

$$\begin{aligned}y_i &= 1 + 2v_i + \varepsilon_i \\ &= 1 + 2(x_i - e_i) + \varepsilon_i \\ &= 1 + 2x_i + (\varepsilon_i - 2e_i) \\ &= 1 + 2x_i + \Delta_i\end{aligned}$$

$$\Delta_i = \varepsilon_i - 2e_i$$

$$\text{cov}(x_i, \Delta_i) = -2 \times 0.5 = -1 \quad (\text{endogeneity})$$

$$\text{Var}(x_i) = 1.5$$

$$\text{AB} = \frac{\text{cov}(x_i, \Delta_i)}{\text{Var}(x_i)} = -\frac{1}{1.5} = -0.667$$

OMITTED VARIABLES

Suppose that the “true” model is

$$y_i = \alpha + \beta_1 \overbrace{x_{1i}}^{\text{included}} + \beta_2 \overbrace{x_{2i}}^{\text{omitted}} + \underbrace{\varepsilon_i}_{\text{independent errors}}$$

OMITTED VARIABLES

Suppose that the “true” model is

$$y_i = \alpha + \beta_1 \overbrace{x_{1i}}^{\text{included}} + \beta_2 \overbrace{x_{2i}}^{\text{omitted}} + \overbrace{\varepsilon_i}^{\text{independent errors}}$$

To fix ideas suppose that

$$\text{Var}(x_{1i}) = \text{Var}(x_{2i}) = 1, \quad \text{var}(\varepsilon_i) = 0.1$$

and

$$\text{Cov}(x_{1i}, x_{2i}) = 0.4,$$

OMITTED VARIABLES (continued)

Suppose variable x_2 is omitted and we fit the model

$$y_i = \alpha + \beta x_{1i} + \Delta_i$$

OMITTED VARIABLES (continued)

Suppose variable x_2 is omitted and we fit the model

$$y_i = \alpha + \beta x_{1i} + \Delta_i$$

$$\Delta_i = \beta_2 x_{2i} + \varepsilon_i$$

$$\text{cov}(\Delta_i, x_{1i}) = 0.4\beta_2 \quad (\text{endogeneity})$$

OMITTED VARIABLES (continued)

Suppose variable x_2 is omitted and we fit the model

$$y_i = \alpha + \beta x_{1i} + \Delta_i$$

$$\Delta_i = \beta_2 x_{2i} + \varepsilon_i$$

$$\text{cov}(\Delta_i, x_{1i}) = 0.4\beta_2 \quad (\text{endogeneity})$$

$$AB = \frac{\text{cov}(x_{1i}, \Delta_i)}{\text{Var}(x_{1i})} = 0.4\beta_2$$

SOME EXAMPLES OF ENDOGENEITY

- **Risk factors of coronary heart disease**

Carrol, Enciclopedia of Biostatistics, 2005: patients's long-term blood pressure may be measured with error (**errors-in-variables**)

SOME EXAMPLES OF ENDOGENEITY

- **Risk factors of coronary heart disease**

Carrol, Enciclopedia of Biostatistics, 2005: patients's long-term blood pressure may be measured with error (**errors-in-variables**)

- **Effect of smoking on birth weight**

Permutt and Hebel, Biometrics, 1989: smoking by pregnant mothers may be correlated with unobserved socioeconomic factors which affect their infant birth weight (**omitted covariates**)

SOME EXAMPLES OF ENDOGENEITY

- **Risk factors of coronary heart disease**

Carrol, Enciclopedia of Biostatistics, 2005: patients's long-term blood pressure may be measured with error (**errors-in-variables**)

- **Effect of smoking on birth weight**

Permutt and Hebel, Biometrics, 1989: smoking by pregnant mothers may be correlated with unobserved socioeconomic factors which affect their infant birth weight (**omitted covariates**)

- **Atherosclerotic cardiovascular disease**

Holmes et al., PloS ONE, 2010: C-reactive protein (CPR) may **cause** atherogenesis or may **be produced** by inflammation within atherosclerotic plaque. Thus CPR may be a cause and a result of atherosclerosis (**simultaneity**)

THE METHOD OF INSTRUMENTAL VARIABLES

- First introduced in **Appendix B** of a book authored by **Philip G. Wright**, “The Tariff on Animal and Vegetable Oils”, published in **1928**.

- First introduced in **Appendix B** of a book authored by **Philip G. Wright**, “The Tariff on Animal and Vegetable Oils”, published in **1928**.
 - Some people attribute the authorship of Appendix B to Philip’s eldest son, **Sewal Wright**

- First introduced in **Appendix B** of a book authored by **Philip G. Wright**, “The Tariff on Animal and Vegetable Oils”, published in **1928**.
 - Some people attribute the authorship of Appendix B to Philip’s eldest son, **Sewal Wright**
- The method consists of using **“instruments”** (variables not included in the model) to consistently estimate the **regression coefficients of endogenous covariates**

- First introduced in **Appendix B** of a book authored by **Philip G. Wright**, “The Tariff on Animal and Vegetable Oils”, published in **1928**.
 - Some people attribute the authorship of Appendix B to Philip’s eldest son, **Sewal Wright**
- The method consists of using **“instruments”** (variables not included in the model) to consistently estimate the **regression coefficients of endogenous covariates**
- Very useful to investigate possible **causal relations** in observational studies

INSTRUMENTAL VARIABLES

- “**Instruments**” are variables

$$\mathbf{z}_i = \begin{pmatrix} z_{1i} \\ z_{2i} \\ \vdots \\ z_{qi} \end{pmatrix} \quad i = 1, 2, \dots, n, \quad q \geq p$$

INSTRUMENTAL VARIABLES

- “**Instruments**” are variables

$$\mathbf{z}_i = \begin{pmatrix} z_{1i} \\ z_{2i} \\ \vdots \\ z_{qi} \end{pmatrix} \quad i = 1, 2, \dots, n, \quad q \geq p$$

- \mathbf{z}_i is **correlated** with the endogenous covariates

INSTRUMENTAL VARIABLES

- “**Instruments**” are variables

$$\mathbf{z}_i = \begin{pmatrix} z_{1i} \\ z_{2i} \\ \vdots \\ z_{qi} \end{pmatrix} \quad i = 1, 2, \dots, n, \quad q \geq p$$

- \mathbf{z}_i is **correlated** with the endogenous covariates
- \mathbf{z}_i is **uncorrelated** with the error term

INSTRUMENTAL VARIABLES

- “**Instruments**” are variables

$$\mathbf{z}_i = \begin{pmatrix} z_{1i} \\ z_{2i} \\ \vdots \\ z_{qi} \end{pmatrix} \quad i = 1, 2, \dots, n, \quad q \geq p$$

- \mathbf{z}_i is **correlated** with the endogenous covariates
- \mathbf{z}_i is **uncorrelated** with the error term
- rank of $\text{Cov}(\mathbf{z}_i) = q$

INSTRUMENTAL VARIABLES

- “**Instruments**” are variables

$$\mathbf{z}_i = \begin{pmatrix} z_{1i} \\ z_{2i} \\ \vdots \\ z_{qi} \end{pmatrix} \quad i = 1, 2, \dots, n, \quad q \geq p$$

- \mathbf{z}_i is **correlated** with the endogenous covariates
- \mathbf{z}_i is **uncorrelated** with the error term
- rank of $\text{Cov}(\mathbf{z}_i) = q$
- rank of $\text{Cov}(\mathbf{z}_i, \mathbf{x}_i) = p$

THE CENTERED INSTRUMENT MATRIX

$$Z_c = \begin{pmatrix} z_{11} - \bar{z}_1 & z_{12} - \bar{z}_2 & \cdots & z_{1p} - \bar{z}_p \\ z_{21} - \bar{z}_1 & z_{22} - \bar{z}_2 & \cdots & z_{2p} - \bar{z}_p \\ z_{31} - \bar{z}_1 & z_{32} - \bar{z}_2 & \cdots & z_{3p} - \bar{z}_p \\ \vdots & \vdots & & \vdots \\ z_{n1} - \bar{z}_1 & z_{n2} - \bar{z}_2 & \cdots & z_{np} - \bar{z}_p \end{pmatrix}$$

- First, project the model covariates on the space generated by the instruments

$$\hat{X}_c = \overbrace{\left[Z_c (Z_c' Z_c)^{-1} Z_c' \right]}^{\text{projection matrix}} X_c = P_z X_c$$

- First, project the model covariates on the space generated by the instruments

$$\hat{X}_c = \overbrace{\left[Z_c (Z_c' Z_c)^{-1} Z_c' \right]}^{\text{projection matrix}} X_c = P_z X_c$$

- \hat{X}_c represents the part of X_c uncorrelated with ε

THE METHOD (CONTINUED)

- Second, apply ordinary LS but using the design matrix \widehat{X}_c instead of X_c :

THE METHOD (CONTINUED)

- Second, apply ordinary LS but using the design matrix \widehat{X}_c instead of X_c :

$$\begin{aligned}\widehat{\beta}_{IV} &= \left(\widehat{X}_c' \widehat{X}_c\right)^{-1} \widehat{X}_c' \mathbf{y} \\ &= \left(X_c' P_Z X_c\right)^{-1} X_c' P_Z \mathbf{y} \\ &= \left(X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' X_c\right)^{-1} X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' \mathbf{y}\end{aligned}$$

THE METHOD (CONTINUED)

- Second, apply ordinary LS but using the design matrix \widehat{X}_c instead of X_c :

$$\begin{aligned}\widehat{\beta}_{IV} &= \left(\widehat{X}_c' \widehat{X}_c\right)^{-1} \widehat{X}_c \mathbf{y} \\ &= \left(X_c' P_Z X_c\right)^{-1} X_c' P_Z \mathbf{y} \\ &= \left(X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' X_c\right)^{-1} X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' \mathbf{y} \\ &= \left(\widehat{\Sigma}_{xz} \widehat{\Sigma}_{zz}^{-1} \widehat{\Sigma}_{zx}\right)^{-1} \widehat{\Sigma}_{xz} \widehat{\Sigma}_{zz}^{-1} \widehat{\Sigma}_{zy}\end{aligned}$$

THE METHOD (CONTINUED)

- Second, apply ordinary LS but using the design matrix \widehat{X}_c instead of X_c :

$$\begin{aligned}\widehat{\beta}_{IV} &= \left(\widehat{X}_c' \widehat{X}_c\right)^{-1} \widehat{X}_c' \mathbf{y} \\ &= \left(X_c' P_Z X_c\right)^{-1} X_c' P_Z \mathbf{y} \\ &= \left(X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' X_c\right)^{-1} X_c' Z_c \left(Z_c' Z_c\right)^{-1} Z_c' \mathbf{y} \\ &= \left(\widehat{\Sigma}_{xz} \widehat{\Sigma}_{zz}^{-1} \widehat{\Sigma}_{zx}\right)^{-1} \widehat{\Sigma}_{xz} \widehat{\Sigma}_{zz}^{-1} \widehat{\Sigma}_{zy} \quad \Leftarrow \text{KEY OBSERVATION}\end{aligned}$$

- $\hat{\beta}_{IV}$ is **consistent** and **asymptotically normal** under mild regularity assumptions

STATISTICAL PROPERTIES

- $\hat{\beta}_{IV}$ is **consistent** and **asymptotically normal** under mild regularity assumptions
- But data may contain outliers in the
 - **response variable** y_i
 - **covariates** x_i
 - **instruments** z_i

- $\hat{\beta}_{IV}$ is **not robust** (sensitive to outliers in the response, covariates and/or instruments)

STATISTICAL PROPERTIES

- $\hat{\beta}_{IV}$ is **not robust** (sensitive to outliers in the response, covariates and/or instruments)
 - $\hat{\beta}_{IV}$ has **unbounded influence** function

STATISTICAL PROPERTIES

- $\hat{\beta}_{IV}$ is **not robust** (sensitive to outliers in the response, covariates and/or instruments)
 - $\hat{\beta}_{IV}$ has **unbounded influence** function
 - $\hat{\beta}_{IV}$ has **zero breakdown point**

- $\hat{\beta}_{IV}$ is **not robust** (sensitive to outliers in the response, covariates and/or instruments)
 - $\hat{\beta}_{IV}$ has **unbounded influence** function
 - $\hat{\beta}_{IV}$ has **zero breakdown point**
 - $\hat{\beta}_{IV}$ **fails robustness-tests** in simulations

- $\hat{\beta}_{IV}$ is **not robust** (sensitive to outliers in the response, covariates and/or instruments)
 - $\hat{\beta}_{IV}$ has **unbounded influence** function
 - $\hat{\beta}_{IV}$ has **zero breakdown point**
 - $\hat{\beta}_{IV}$ **fails robustness-tests** in simulations
 - $\hat{\beta}_{IV}$ **has poor performance in some real data examples** due to outliers.

APPLICATION OF INSTRUMENTAL VARIABLES IN MEDICAL RESEARCH

OVERVIEW

THE FRAMINGHAM HEART STUDY

- Prior work suggests that **high blood pressure** can increase the size of the **left atrial dimension** (Vaziri et al., 1995; Milan et al., 2009; Benjamin et al., 1995).

OVERVIEW

THE FRAMINGHAM HEART STUDY

- Prior work suggests that **high blood pressure** can increase the size of the **left atrial dimension** (Vaziri et al., 1995; Milan et al., 2009; Benjamin et al., 1995).
- We use IV to measure the effect of *long-term systolic blood pressure* (SBP) on left atrial size (LAD).

- Our sample includes **164** male subjects from the original cohort who underwent their 16th biennial examination and satisfy the inclusion criteria employed by Vaziri et al. (1995):
 - patient doesn't have a history of heart disease,
 - patient is not taking cardiovascular medication
 - patient has complete clinical data

STATISTICAL MODEL

Following Vaziri et al. (1995), we consider the following model:

Variable	Description	Variable Status
LAD	Left atrial dimension	Response
SBP₁₆	long-term systolic blood pressure (Examination # 16)	Main covariate Endogenous (measured with error)
BMI	Body mass index	Control covariate Exogenous
AGE	Subject's age	Control covariate Exogenous
SBP₁₅	long-term systolic blood pressure (Examination # 15)	Instrument

RESULTS

	SBP	BMI	AGE
COEFF.	0.011	3.195	0.023
P-VALUE	0.684	<0.001	0.685

APPLICATION OF INSTRUMENTAL VARIABLES IN LABOR ECONOMICS

APPLICATION OF RIV TO LABOR UNIONS DATA

- Our dataset contains 181 industries in the United States for the year 2003.

APPLICATION OF RIV TO LABOR UNIONS DATA

- Our dataset contains 181 industries in the United States for the year 2003.
- Unionization rates are obtained from the Union Membership and Coverage Database

APPLICATION OF RIV TO LABOR UNIONS DATA

- Our dataset contains 181 industries in the United States for the year 2003.
- Unionization rates are obtained from the Union Membership and Coverage Database
- The fraction of male workers comes from the Census Population Survey

APPLICATION OF RIV TO LABOR UNIONS DATA

- Our dataset contains 181 industries in the United States for the year 2003.
- Unionization rates are obtained from the Union Membership and Coverage Database
- The fraction of male workers comes from the Census Population Survey
- The remaining variables are constructed using data from Compustat

APPLICATION OF RIV TO LABOR UNIONS DATA

- Our dataset contains 181 industries in the United States for the year 2003.
- Unionization rates are obtained from the Union Membership and Coverage Database
- The fraction of male workers comes from the Census Population Survey
- The remaining variables are constructed using data from Compustat

EFFECT OF LABOR UNIONS ON FIRMS' PROFITS

Variable Name	Variable Description	Variable Status
ROA	Accounting profits scaled by assets	Response variable
UNION	Fraction of workers that are unionized	Endogenous , main covariate
GROWTH	Sales growth rate	Exogenous , control covariate
CAPEX	Capital expenditures scaled by assets	Exogenous , control covariate

EFFECT OF LABOR UNIONS ON FIRMS' PROFITS

Variable Name	Variable Description	Variable Status
SIZE	Fixed assets (log scale)	Exogenous , control covariate
LEV	Total debt scaled by assets	Exogenous , control covariate
KL	Capital-labor ratio	Exogenous , control covariate
MALE	Fraction of male workers	Instrument

- Why is UNION endogenous?

- Why is UNION endogenous?
- There are two possible reasons

- Why is UNION endogenous?
- There are two possible reasons
 1. UNION and ROA are likely to be simultaneously determined

- Why is UNION endogenous?
- There are two possible reasons
 1. UNION and ROA are likely to be simultaneously determined
 - Industry higher performance may stimulate unionization to get a better profits share

- Why is UNION endogenous?
- There are two possible reasons
 1. UNION and ROA are likely to be simultaneously determined
 - Industry higher performance may stimulate unionization to get a better profits share
 - It is conjectured that higher unionization cause industry performance to decline

- Why is UNION endogenous?
- There are two possible reasons
 1. UNION and ROA are likely to be simultaneously determined
 - Industry higher performance may stimulate unionization to get a better profits share
 - It is conjectured that higher unionization cause industry performance to decline
 2. There may be unobserved industry features that determine performance and are correlated with unionization

- Why is MALE a reasonable instrument for UNION?

- Why is MALE a reasonable instrument for UNION?
 - males have more permanent attachment to the labor market and to internal job ladders

- Why is MALE a reasonable instrument for UNION?
 - males have more permanent attachment to the labor market and to internal job ladders
 - ⇒ males derive higher wage/non-wage benefits from unionizing

- Why is MALE a reasonable instrument for UNION?
 - males have more permanent attachment to the labor market and to internal job ladders
 - ⇒ males derive higher wage/non-wage benefits from unionizing
 - ⇒ males are more likely to become unionized.

RESULTS

	UNION	GROWTH	CAPEX	SIZE	LEV	KL
coeff.	-0.271	0.126	0.197	0.003	0.094	0.004
p-value	0.128	0.013	0.000	0.312	0.005	0.313