

Introduction to Regression with Measurement Error¹

STA431 Spring 2023

¹See last slide for copyright information.

Overview

- 1 Measurement Error
- 2 Reliability
- 3 Consequences of Ignoring Measurement Error

Measurement Error

- Snack food consumption
- Exercise
- Income
- Cause of death (classification error)

Measurement Error

- Snack food consumption
- Exercise
- Income
- Cause of death (classification error)
- Even amount of drug that reaches animals blood stream in an experimental study.
- Is there anything that is *not* measured with error?

Additive measurement error

A very simple model

$$W = X + e$$

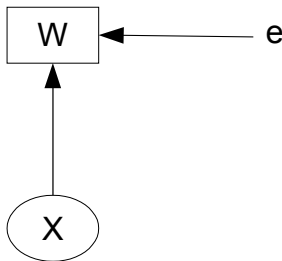
Where $E(X) = \mu_x$, $E(e) = 0$, $Var(X) = \sigma_x^2$, $Var(e) = \sigma_e^2$, and $Cov(X, e) = 0$.

Additive measurement error

A very simple model

$$W = X + e$$

Where $E(X) = \mu_x$, $E(e) = 0$, $Var(X) = \sigma_x^2$, $Var(e) = \sigma_e^2$, and $Cov(X, e) = 0$.



Variance and Covariance

$$W = X + e$$

$$\begin{aligned} \text{Var}(W) &= \text{Var}(X) + \text{Var}(e) \\ &= \sigma_x^2 + \sigma_e^2 \end{aligned}$$

Variance and Covariance

$$W = X + e$$

$$\begin{aligned} \text{Var}(W) &= \text{Var}(X) + \text{Var}(e) \\ &= \sigma_x^2 + \sigma_e^2 \end{aligned}$$

$$\begin{aligned} \text{Cov}(X, W) &= \text{Cov}(X, X + e) \\ &= \text{Cov}(X, X) + \text{Cov}(X, e) \\ &= \sigma_x^2 \end{aligned}$$

Explained Variance

- Variance is an index of unit-to-unit variation in a measurement.
- Explaining unit-to-unit variation is an important goal of Science.

Explained Variance

- Variance is an index of unit-to-unit variation in a measurement.
- Explaining unit-to-unit variation is an important goal of Science.
- How much of the variation in an observed variable comes from variation in the latent quantity of interest, and how much comes from random noise?

Definition of Reliability

Reliability is the squared correlation between the observed variable and the latent variable (true score).

Calculation of Reliability

Squared correlation between observed and true score

$$\rho^2$$

Calculation of Reliability

Squared correlation between observed and true score

$$\rho^2 = \left(\frac{Cov(X, W)}{SD(X)SD(W)} \right)^2$$

Calculation of Reliability

Squared correlation between observed and true score

$$\begin{aligned}\rho^2 &= \left(\frac{Cov(X, W)}{SD(X)SD(W)} \right)^2 \\ &= \left(\frac{\sigma_x^2}{\sqrt{\sigma_x^2} \sqrt{\sigma_x^2 + \sigma_e^2}} \right)^2\end{aligned}$$

Calculation of Reliability

Squared correlation between observed and true score

$$\begin{aligned}\rho^2 &= \left(\frac{Cov(X, W)}{SD(X)SD(W)} \right)^2 \\ &= \left(\frac{\sigma_x^2}{\sqrt{\sigma_x^2} \sqrt{\sigma_x^2 + \sigma_e^2}} \right)^2 \\ &= \frac{\sigma_x^4}{\sigma_x^2(\sigma_x^2 + \sigma_e^2)}\end{aligned}$$

Calculation of Reliability

Squared correlation between observed and true score

$$\begin{aligned}\rho^2 &= \left(\frac{Cov(X, W)}{SD(X)SD(W)} \right)^2 \\ &= \left(\frac{\sigma_x^2}{\sqrt{\sigma_x^2} \sqrt{\sigma_x^2 + \sigma_e^2}} \right)^2 \\ &= \frac{\sigma_x^4}{\sigma_x^2(\sigma_x^2 + \sigma_e^2)} \\ &= \frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2}.\end{aligned}$$

Calculation of Reliability

Squared correlation between observed and true score

$$\begin{aligned}\rho^2 &= \left(\frac{Cov(X, W)}{SD(X)SD(W)} \right)^2 \\ &= \left(\frac{\sigma_x^2}{\sqrt{\sigma_x^2} \sqrt{\sigma_x^2 + \sigma_e^2}} \right)^2 \\ &= \frac{\sigma_x^4}{\sigma_x^2(\sigma_x^2 + \sigma_e^2)} \\ &= \frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2}.\end{aligned}$$

Reliability is the proportion of the variance in the observed variable that comes from the latent variable of interest, and not from random error.

How to estimate reliability from data

- Correlate usual measurement with “Gold Standard?”
- Not very realistic, except maybe for some bio-markers.

How to estimate reliability from data

- Correlate usual measurement with “Gold Standard?”
- Not very realistic, except maybe for some bio-markers.
- One answer: Measure twice.

Measure twice

Called “equivalent measurements” because **error variance is the same**

$$W_1 = X + e_1$$

$$W_2 = X + e_2,$$

where $E(X) = \mu_x$, $Var(X) = \sigma_x^2$, $E(e_1) = E(e_2) = 0$,
 $Var(e_1) = Var(e_2) = \sigma_e^2$, and X , e_1 and e_2 are all independent.

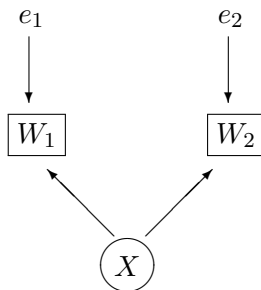
Measure twice

Called “equivalent measurements” because **error variance is the same**

$$W_1 = X + e_1$$

$$W_2 = X + e_2,$$

where $E(X) = \mu_x$, $Var(X) = \sigma_x^2$, $E(e_1) = E(e_2) = 0$,
 $Var(e_1) = Var(e_2) = \sigma_e^2$, and X , e_1 and e_2 are all independent.



Reliability equals the correlation between two equivalent measurements

This is a population correlation

Reliability equals the correlation between two equivalent measurements

This is a population correlation

$$\text{Corr}(W_1, W_2) = \frac{\text{Cov}(W_1, W_2)}{SD(W_1)SD(W_2)}$$

Reliability equals the correlation between two equivalent measurements

This is a population correlation

$$\begin{aligned} \text{Corr}(W_1, W_2) &= \frac{\text{Cov}(W_1, W_2)}{SD(W_1)SD(W_2)} \\ &= \frac{\text{Cov}(X + e_1, X + e_2)}{\sigma_x^2 + \sigma_e^2} \end{aligned}$$

Reliability equals the correlation between two equivalent measurements

This is a population correlation

$$\begin{aligned} \text{Corr}(W_1, W_2) &= \frac{\text{Cov}(W_1, W_2)}{SD(W_1)SD(W_2)} \\ &= \frac{\text{Cov}(X + e_1, X + e_2)}{\sigma_x^2 + \sigma_e^2} \\ &= \frac{\text{Cov}(X, X) + 0 + 0 + 0}{\sigma_x^2 + \sigma_e^2} \end{aligned}$$

Reliability equals the correlation between two equivalent measurements

This is a population correlation

$$\begin{aligned}
 Corr(W_1, W_2) &= \frac{Cov(W_1, W_2)}{SD(W_1)SD(W_2)} \\
 &= \frac{Cov(X + e_1, X + e_2)}{\sigma_x^2 + \sigma_e^2} \\
 &= \frac{Cov(X, X) + 0 + 0 + 0}{\sigma_x^2 + \sigma_e^2} \\
 &= \frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2},
 \end{aligned}$$

which is the reliability.

Estimate the reliability: Measure twice for a sample of size n

With a well-chosen time gap

$$\text{Calculate } r = \frac{\sum_{i=1}^n (W_{i1} - \bar{W}_1)(W_{i2} - \bar{W}_2)}{\sqrt{\sum_{i=1}^n (W_{i1} - \bar{W}_1)^2} \sqrt{\sum_{i=1}^n (W_{i2} - \bar{W}_2)^2}}.$$

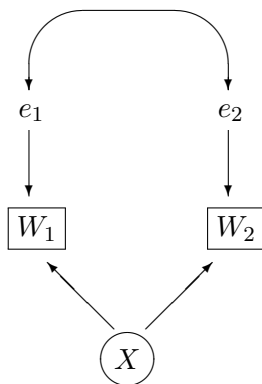
Estimate the reliability: Measure twice for a sample of size n

With a well-chosen time gap

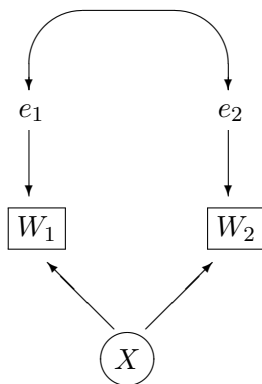
$$\text{Calculate } r = \frac{\sum_{i=1}^n (W_{i1} - \bar{W}_1)(W_{i2} - \bar{W}_2)}{\sqrt{\sum_{i=1}^n (W_{i1} - \bar{W}_1)^2} \sqrt{\sum_{i=1}^n (W_{i2} - \bar{W}_2)^2}}.$$

- Test-retest reliability
- Alternate forms reliability
- Split-half reliability

Omitted variables can cause correlated measurement error



Omitted variables can cause correlated measurement error



Leading to an over-estimate of reliability.

Measurement error in regression analysis

- Mostly we are interested in relationships between latent (true) variables.
- But all we have at best are the true variables measured with error.
- Models like $Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + \epsilon_i$ are mis-specified.

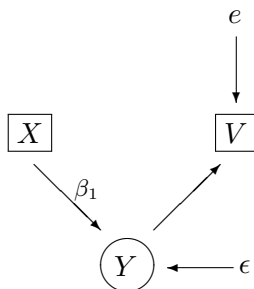
Measurement error in regression analysis

- Mostly we are interested in relationships between latent (true) variables.
- But all we have at best are the true variables measured with error.
- Models like $Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + \epsilon_i$ are mis-specified.
- The most common way of dealing with measurement error in regression is to ignore it.
- What effect does this have on estimation and inference?

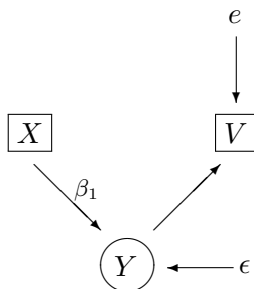
Measurement error in regression analysis

- Mostly we are interested in relationships between latent (true) variables.
- But all we have at best are the true variables measured with error.
- Models like $Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + \epsilon_i$ are mis-specified.
- The most common way of dealing with measurement error in regression is to ignore it.
- What effect does this have on estimation and inference?
- First consider ignoring measurement error just in the response variable.

Measurement error in the response variable



Measurement error in the response variable



True model:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e_i$$

Naive model: $V_i = \beta_0 + \beta_1 X_i + \epsilon_i$

Is $\hat{\beta}_1$ consistent?

Ignoring measurement error in Y

First calculate $Cov(X_i, V_i)$. Under the true model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e_i,$$

Is $\hat{\beta}_1$ consistent?

Ignoring measurement error in Y

First calculate $Cov(X_i, V_i)$. Under the true model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e_i,$$

$$Cov(X_i, V_i) = Cov(X, \beta_1 X_i + \epsilon_i)$$

Is $\hat{\beta}_1$ consistent?

Ignoring measurement error in Y

First calculate $Cov(X_i, V_i)$. Under the true model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e_i,$$

$$\begin{aligned} Cov(X_i, V_i) &= Cov(X, \beta_1 X_i + \epsilon_i) \\ &= \beta_1 \sigma_x^2 \end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(V_i - \bar{V})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(V_i - \bar{V})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\hat{\sigma}_{x,v}}{\hat{\sigma}_x^2}\end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(V_i - \bar{V})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\hat{\sigma}_{x,v}}{\hat{\sigma}_x^2} \\ &\xrightarrow{p} \frac{Cov(X_i, V_i)}{Var(X_i)}\end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(V_i - \bar{V})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
 &= \frac{\hat{\sigma}_{x,v}}{\hat{\sigma}_x^2} \\
 &\xrightarrow{p} \frac{Cov(X_i, V_i)}{Var(X_i)} \\
 &= \frac{\beta_1 \sigma_x^2}{\sigma_x^2}
 \end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

Have $Cov(X_i, V_i) = \beta_1 \sigma_x^2$ and $Var(X_i) = \sigma_x^2$

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(V_i - \bar{V})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
 &= \frac{\hat{\sigma}_{x,v}}{\hat{\sigma}_x^2} \\
 &\xrightarrow{p} \frac{Cov(X_i, V_i)}{Var(X_i)} \\
 &= \frac{\beta_1 \sigma_x^2}{\sigma_x^2} \\
 &= \beta_1
 \end{aligned}$$

Why did it work?

Why did it work?

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e$$

Why did it work?

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e$$

$$= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i$$

Why did it work?

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$V_i = \nu + Y_i + e$$

$$= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i$$

$$= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i)$$

Why did it work?

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\V_i &= \nu + Y_i + e \\&= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i \\&= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i) \\&= \beta'_0 + \beta_1 X_i + \epsilon'_i\end{aligned}$$

Why did it work?

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\V_i &= \nu + Y_i + e \\&= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i \\&= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i) \\&= \beta'_0 + \beta_1 X_i + \epsilon'_i\end{aligned}$$

- This is a re-parameterization.

Why did it work?

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\V_i &= \nu + Y_i + e \\&= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i \\&= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i) \\&= \beta'_0 + \beta_1 X_i + \epsilon'_i\end{aligned}$$

- This is a re-parameterization.
- Most definitely *not* one-to-one.

Why did it work?

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\V_i &= \nu + Y_i + e \\&= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i \\&= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i) \\&= \beta'_0 + \beta_1 X_i + \epsilon'_i\end{aligned}$$

- This is a re-parameterization.
- Most definitely *not* one-to-one.
- (ν, β_0) is absorbed into β'_0 .
- (ϵ_i, e_i) is absorbed into ϵ'_i .

Why did it work?

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\V_i &= \nu + Y_i + e \\&= \nu + (\beta_0 + \beta_1 X_i + \epsilon_i) + e_i \\&= (\nu + \beta_0) + \beta_1 X_i + (\epsilon_i + e_i) \\&= \beta'_0 + \beta_1 X_i + \epsilon'_i\end{aligned}$$

- This is a re-parameterization.
- Most definitely *not* one-to-one.
- (ν, β_0) is absorbed into β'_0 .
- (ϵ_i, e_i) is absorbed into ϵ'_i .
- Can't know everything, but all we care about is β_1 anyway.

Don't Worry

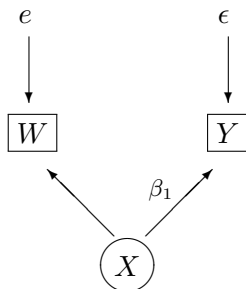
Don't Worry

- If a response variable appears to have no measurement error, assume it does have measurement error but the problem has been re-parameterized.

Don't Worry

- If a response variable appears to have no measurement error, assume it does have measurement error but the problem has been re-parameterized.
- Measurement error in Y is part of ϵ .

Measurement error in a single explanatory variable



True model:

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\ W_i &= X_i + e_i,\end{aligned}$$

Naive model: $Y_i = \beta_0 + \beta_1 W_i + \epsilon_i$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \text{ and } W_i = X_i + e_i$$

Have $Cov(W_i, Y_i) = \beta_1 \sigma_x^2$ and $Var(W_i) = \sigma_x^2 + \sigma_e^2$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \text{ and } W_i = X_i + e_i$$

Have $Cov(W_i, Y_i) = \beta_1 \sigma_x^2$ and $Var(W_i) = \sigma_x^2 + \sigma_e^2$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (W_i - \bar{W})(Y_i - \bar{Y})}{\sum_{i=1}^n (W_i - \bar{W})^2}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ and $W_i = X_i + e_i$

Have $Cov(W_i, Y_i) = \beta_1 \sigma_x^2$ and $Var(W_i) = \sigma_x^2 + \sigma_e^2$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (W_i - \bar{W})(Y_i - \bar{Y})}{\sum_{i=1}^n (W_i - \bar{W})^2} \\ &= \frac{\hat{\sigma}_{w,y}}{\hat{\sigma}_w^2}\end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ and $W_i = X_i + e_i$

Have $Cov(W_i, Y_i) = \beta_1 \sigma_x^2$ and $Var(W_i) = \sigma_x^2 + \sigma_e^2$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (W_i - \bar{W})(Y_i - \bar{Y})}{\sum_{i=1}^n (W_i - \bar{W})^2} \\ &= \frac{\hat{\sigma}_{w,y}}{\hat{\sigma}_w^2} \\ &\xrightarrow{p} \frac{Cov(W_i, Y_i)}{Var(W_i)}\end{aligned}$$

Target of $\hat{\beta}_1$ as $n \rightarrow \infty$

$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ and $W_i = X_i + e_i$

Have $Cov(W_i, Y_i) = \beta_1 \sigma_x^2$ and $Var(W_i) = \sigma_x^2 + \sigma_e^2$

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{\sum_{i=1}^n (W_i - \bar{W})(Y_i - \bar{Y})}{\sum_{i=1}^n (W_i - \bar{W})^2} \\
 &= \frac{\hat{\sigma}_{w,y}}{\hat{\sigma}_w^2} \\
 &\xrightarrow{p} \frac{Cov(W_i, Y_i)}{Var(W_i)} \\
 &= \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)
 \end{aligned}$$

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)$$

$$W_i = X_i + e_i$$

- $\hat{\beta}_1$ converges to β times the reliability of W_i .
- It's inconsistent.

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)$$

$$W_i = X_i + e_i$$

- $\hat{\beta}_1$ converges to β times the reliability of W_i .
- It's inconsistent.
- Because the reliability is less than one, it's asymptotically biased toward zero.
- The worse the measurement of X_i , the more the asymptotic bias.

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)$$

$$W_i = X_i + e_i$$

- $\hat{\beta}_1$ converges to β times the reliability of W_i .
- It's inconsistent.
- Because the reliability is less than one, it's asymptotically biased toward zero.
- The worse the measurement of X_i , the more the asymptotic bias.
- Sometimes called “attenuation” (weakening).
- If a good estimate of reliability is available from another source, one can “correct for attenuation.”

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)$$

$$W_i = X_i + e_i$$

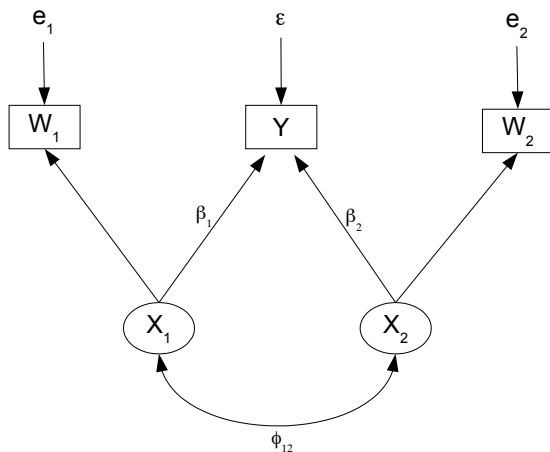
- $\hat{\beta}_1$ converges to β times the reliability of W_i .
- It's inconsistent.
- Because the reliability is less than one, it's asymptotically biased toward zero.
- The worse the measurement of X_i , the more the asymptotic bias.
- Sometimes called “attenuation” (weakening).
- If a good estimate of reliability is available from another source, one can “correct for attenuation.”
- When $H_0 : \beta_1 = 0$ is true, it's not a serious problem.

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} \right)$$

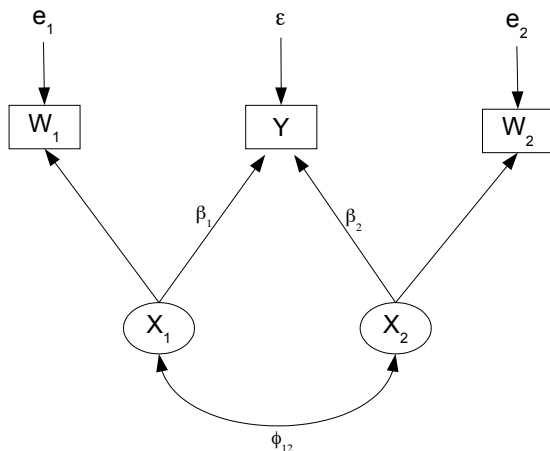
$$W_i = X_i + e_i$$

- $\hat{\beta}_1$ converges to β times the reliability of W_i .
- It's inconsistent.
- Because the reliability is less than one, it's asymptotically biased toward zero.
- The worse the measurement of X_i , the more the asymptotic bias.
- Sometimes called “attenuation” (weakening).
- If a good estimate of reliability is available from another source, one can “correct for attenuation.”
- When $H_0 : \beta_1 = 0$ is true, it's not a serious problem.
- False sense of security?

Measurement error in two explanatory variables



Measurement error in two explanatory variables



Want to assess the relationship of X_2 to Y , *controlling* for X_1 by testing $H_0 : \beta_2 = 0$.

Statement of the model

Independently for $i = 1, \dots, n$

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\W_{i,1} &= X_{i,1} + e_{i,1} \\W_{i,2} &= X_{i,2} + e_{i,2},\end{aligned}$$

Statement of the model

Independently for $i = 1, \dots, n$

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\ W_{i,1} &= X_{i,1} + e_{i,1} \\ W_{i,2} &= X_{i,2} + e_{i,2}, \end{aligned}$$

where

$$E(X_{i,1}) = \mu_1, E(X_{i,2}) = \mu_2, E(\epsilon_i) = E(e_{i,1}) = E(e_{i,2}) = 0,$$

$$Var(\epsilon_i) = \psi, Var(e_{i,1}) = \omega_1, Var(e_{i,2}) = \omega_2,$$

The errors $\epsilon_i, e_{i,1}$ and $e_{i,2}$ are all independent,

$X_{i,1}$ and $X_{i,2}$ are independent of $\epsilon_i, e_{i,1}$ and $e_{i,2}$, and

$$cov \begin{pmatrix} X_{i,1} \\ X_{i,1} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}.$$

Statement of the model

Independently for $i = 1, \dots, n$

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\ W_{i,1} &= X_{i,1} + e_{i,1} \\ W_{i,2} &= X_{i,2} + e_{i,2}, \end{aligned}$$

where

$$E(X_{i,1}) = \mu_1, E(X_{i,2}) = \mu_2, E(\epsilon_i) = E(e_{i,1}) = E(e_{i,2}) = 0,$$

$$Var(\epsilon_i) = \psi, Var(e_{i,1}) = \omega_1, Var(e_{i,2}) = \omega_2,$$

The errors $\epsilon_i, e_{i,1}$ and $e_{i,2}$ are all independent,

$X_{i,1}$ and $X_{i,2}$ are independent of $\epsilon_i, e_{i,1}$ and $e_{i,2}$, and

$$cov \begin{pmatrix} X_{i,1} \\ X_{i,1} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}.$$

Note

- Reliability of W_1 is $\frac{\phi_{11}}{\phi_{11} + \omega_1}$.
- Reliability of W_2 is $\frac{\phi_{22}}{\phi_{22} + \omega_2}$.

True Model versus Naive Model

True model:

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\W_{i,1} &= X_{i,1} + e_{i,1} \\W_{i,2} &= X_{i,2} + e_{i,2},\end{aligned}$$

Naive model: $Y_i = \beta_0 + \beta_1 W_{i,1} + \beta_2 W_{i,2} + \epsilon_i$

True Model versus Naive Model

True model:

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\W_{i,1} &= X_{i,1} + e_{i,1} \\W_{i,2} &= X_{i,2} + e_{i,2},\end{aligned}$$

Naive model: $Y_i = \beta_0 + \beta_1 W_{i,1} + \beta_2 W_{i,2} + \epsilon_i$

- Fit the naive model.
- See what happens to $\hat{\beta}_2$ as $n \rightarrow \infty$ when the true model holds.

True Model versus Naive Model

True model:

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\ W_{i,1} &= X_{i,1} + e_{i,1} \\ W_{i,2} &= X_{i,2} + e_{i,2}, \end{aligned}$$

Naive model: $Y_i = \beta_0 + \beta_1 W_{i,1} + \beta_2 W_{i,2} + \epsilon_i$

- Fit the naive model.
- See what happens to $\hat{\beta}_2$ as $n \rightarrow \infty$ when the true model holds.
- Start by calculating $cov(\mathbf{d}_i) = cov \begin{pmatrix} W_{i,1} \\ W_{i,2} \\ Y_i \end{pmatrix}$.

Covariance matrix of the observable data

Covariance matrix of the observable data

$$\Sigma = cov \begin{pmatrix} W_{i,1} \\ W_{i,2} \\ Y_i \end{pmatrix}$$

Covariance matrix of the observable data

$$\begin{aligned}
\Sigma &= cov \begin{pmatrix} W_{i,1} \\ W_{i,2} \\ Y_i \end{pmatrix} \\
&= \begin{pmatrix} \omega_1 + \phi_{11} & \phi_{12} & \beta_1 \phi_{11} + \beta_2 \phi_{12} \\ \phi_{12} & \omega_2 + \phi_{22} & \beta_1 \phi_{12} + \beta_2 \phi_{22} \\ \beta_1 \phi_{11} + \beta_2 \phi_{12} & \beta_1 \phi_{12} + \beta_2 \phi_{22} & \beta_1^2 \phi_{11} + 2 \beta_1 \beta_2 \phi_{12} + \beta_2^2 \phi_{22} + \psi \end{pmatrix}
\end{aligned}$$

What happens to $\hat{\beta}_2$ as $n \rightarrow \infty$?

Interested in $H_0 : \beta_2 = 0$

$$\hat{\beta}_2 = \frac{\hat{\sigma}_{11}\hat{\sigma}_{23} - \hat{\sigma}_{12}\hat{\sigma}_{13}}{\hat{\sigma}_{11}\hat{\sigma}_{22} - \hat{\sigma}_{12}^2}$$

What happens to $\hat{\beta}_2$ as $n \rightarrow \infty$?

Interested in $H_0 : \beta_2 = 0$

$$\begin{aligned}\hat{\beta}_2 &= \frac{\hat{\sigma}_{11}\hat{\sigma}_{23} - \hat{\sigma}_{12}\hat{\sigma}_{13}}{\hat{\sigma}_{11}\hat{\sigma}_{22} - \hat{\sigma}_{12}^2} \\ &\xrightarrow{p} \frac{\sigma_{11}\sigma_{23} - \sigma_{12}\sigma_{13}}{\sigma_{11}\sigma_{22} - \sigma_{12}^2}\end{aligned}$$

What happens to $\hat{\beta}_2$ as $n \rightarrow \infty$?

Interested in $H_0 : \beta_2 = 0$

$$\begin{aligned}
 \hat{\beta}_2 &= \frac{\hat{\sigma}_{11}\hat{\sigma}_{23} - \hat{\sigma}_{12}\hat{\sigma}_{13}}{\hat{\sigma}_{11}\hat{\sigma}_{22} - \hat{\sigma}_{12}^2} \\
 &\xrightarrow{p} \frac{\sigma_{11}\sigma_{23} - \sigma_{12}\sigma_{13}}{\sigma_{11}\sigma_{22} - \sigma_{12}^2} \\
 &= \frac{\beta_1\omega_1\phi_{12} + \beta_2(\omega_1\phi_{22} + \phi_{11}\phi_{22} - \phi_{12}^2)}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}
 \end{aligned}$$

What happens to $\hat{\beta}_2$ as $n \rightarrow \infty$?

Interested in $H_0 : \beta_2 = 0$

$$\begin{aligned}
 \hat{\beta}_2 &= \frac{\hat{\sigma}_{11}\hat{\sigma}_{23} - \hat{\sigma}_{12}\hat{\sigma}_{13}}{\hat{\sigma}_{11}\hat{\sigma}_{22} - \hat{\sigma}_{12}^2} \\
 &\xrightarrow{p} \frac{\sigma_{11}\sigma_{23} - \sigma_{12}\sigma_{13}}{\sigma_{11}\sigma_{22} - \sigma_{12}^2} \\
 &= \frac{\beta_1\omega_1\phi_{12} + \beta_2(\omega_1\phi_{22} + \phi_{11}\phi_{22} - \phi_{12}^2)}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2} \\
 &\neq \beta_2
 \end{aligned}$$

Inconsistent.

When $H_0 : \beta_2 = 0$ is true

When $H_0 : \beta_2 = 0$ is true

$$\hat{\beta}_2 \xrightarrow{p} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

When $H_0 : \beta_2 = 0$ is true

$$\hat{\beta}_2 \xrightarrow{p} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

So $\hat{\beta}_2$ goes to the wrong target unless

- There is no relationship between X_1 and Y , or
- There is no measurement error in W_1 , or
- There is no correlation between X_1 and X_2 .

When $H_0 : \beta_2 = 0$ is true

$$\hat{\beta}_2 \xrightarrow{p} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

So $\hat{\beta}_2$ goes to the wrong target unless

- There is no relationship between X_1 and Y , or
- There is no measurement error in W_1 , or
- There is no correlation between X_1 and X_2 .

Also, the t statistic for $H_0 : \beta_2 = 0$ goes to plus or minus ∞ and the p -value $\xrightarrow{p} 0$. Remember, H_0 is true.

How bad is it for finite sample sizes?

$$\hat{\beta}_2 \xrightarrow{P} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

How bad is it for finite sample sizes?

$$\hat{\beta}_2 \xrightarrow{P} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

A big simulation study (Brunner and Austin, 2009) with six factors

How bad is it for finite sample sizes?

$$\hat{\beta}_2 \xrightarrow{P} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

A big simulation study (Brunner and Austin, 2009) with six factors

- Sample size: $n = 50, 100, 250, 500, 1000$
- $Corr(X_1, X_2)$: $\phi_{12} = 0.00, 0.25, 0.75, 0.80, 0.90$
- Proportion of variance in Y explained by X_1 : $0.25, 0.50, 0.75$
- Reliability of W_1 : $0.50, 0.75, 0.80, 0.90, 0.95$
- Reliability of W_2 : $0.50, 0.75, 0.80, 0.90, 0.95$
- Distribution of latent variables and error terms: Normal, Uniform, t , Pareto.

How bad is it for finite sample sizes?

$$\hat{\beta}_2 \xrightarrow{P} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

A big simulation study (Brunner and Austin, 2009) with six factors

- Sample size: $n = 50, 100, 250, 500, 1000$
- $Corr(X_1, X_2)$: $\phi_{12} = 0.00, 0.25, 0.75, 0.80, 0.90$
- Proportion of variance in Y explained by X_1 : $0.25, 0.50, 0.75$
- Reliability of W_1 : $0.50, 0.75, 0.80, 0.90, 0.95$
- Reliability of W_2 : $0.50, 0.75, 0.80, 0.90, 0.95$
- Distribution of latent variables and error terms: Normal, Uniform, t , Pareto.

There were $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations.

Simulation study procedure

Within each of the $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations,

- 10,000 random data sets were generated

Simulation study procedure

Within each of the $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations,

- 10,000 random data sets were generated
- For a total of 75 million data sets

Simulation study procedure

Within each of the $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations,

- 10,000 random data sets were generated
- For a total of 75 million data sets
- All generated according to the true model, with $\beta_2 = 0$.
- Fit naive model, test $H_0 : \beta_2 = 0$ at $\alpha = 0.05$.

Simulation study procedure

Within each of the $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations,

- 10,000 random data sets were generated
- For a total of 75 million data sets
- All generated according to the true model, with $\beta_2 = 0$.
- Fit naive model, test $H_0 : \beta_2 = 0$ at $\alpha = 0.05$.
- Proportion of times H_0 is rejected is a Monte Carlo estimate of the Type I Error Probability.

Simulation study procedure

Within each of the $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations,

- 10,000 random data sets were generated
- For a total of 75 million data sets
- All generated according to the true model, with $\beta_2 = 0$.
- Fit naive model, test $H_0 : \beta_2 = 0$ at $\alpha = 0.05$.
- Proportion of times H_0 is rejected is a Monte Carlo estimate of the Type I Error Probability.
- It should be around 0.05.

Look at a small part of the results

Look at a small part of the results

- Both reliabilities = 0.90
- Everything is normally distributed
- $\beta_0 = 1$, $\beta_1 = 1$ and of course $\beta_2 = 0$.

Table 1 of Brunner and Austin (2009, p.39)

Canadian Journal of Statistics, Vol. 37, Pages 33-46, Used without permission

TABLE 1: Estimated Type I error rates when independent variables and measurement errors are all normal, and reliability of W_1 and W_2 both equal 0.90.

| N | Correlation between X_1 and X_2 | | | | |
|--|-------------------------------------|---------------------|--------|--------|--------|
| | 0.0 | 0.2 | 0.4 | 0.6 | 0.8 |
| 25% of variance in Y is explained by X_1 | | | | | |
| 50 | 0.0476 [†] | 0.0505 [†] | 0.0636 | 0.0715 | 0.0913 |
| 100 | 0.0504 [†] | 0.0521 [†] | 0.0834 | 0.0940 | 0.1294 |
| 250 | 0.0467 [†] | 0.0533 [†] | 0.1402 | 0.1624 | 0.2544 |
| 500 | 0.0468 [†] | 0.0595 [†] | 0.2300 | 0.2892 | 0.4649 |
| 1,000 | 0.0505 [†] | 0.0734 | 0.4094 | 0.5057 | 0.7431 |
| 50% of variance in Y is explained by X_1 | | | | | |
| 50 | 0.0460 [†] | 0.0520 [†] | 0.0963 | 0.1106 | 0.1633 |
| 100 | 0.0535 [†] | 0.0569 [†] | 0.1461 | 0.1857 | 0.2837 |
| 250 | 0.0483 [†] | 0.0625 | 0.3068 | 0.3731 | 0.5864 |
| 500 | 0.0515 [†] | 0.0780 | 0.5323 | 0.6488 | 0.8837 |
| 1,000 | 0.0481 [†] | 0.1185 | 0.8273 | 0.9088 | 0.9907 |
| 75% of variance in Y is explained by X_1 | | | | | |
| 50 | 0.0485 [†] | 0.0579 [†] | 0.1727 | 0.2089 | 0.3442 |
| 100 | 0.0541 [†] | 0.0679 | 0.3101 | 0.3785 | 0.6031 |
| 250 | 0.0479 [†] | 0.0856 | 0.6450 | 0.7523 | 0.9434 |
| 500 | 0.0445 [†] | 0.1323 | 0.9109 | 0.9635 | 0.9992 |
| 1,000 | 0.0522 [†] | 0.2179 | 0.9959 | 0.9998 | 1.0000 |

[†]Not significantly different from 0.05, Bonferroni corrected for 7,500 tests.

Weak Relationship between X_1 and Y : $\text{Var} = 25\%$

| N | Correlation Between X_1 and X_2 | | | | |
|------|-------------------------------------|---------|---------|---------|---------|
| | 0.00 | 0.25 | 0.75 | 0.80 | 0.90 |
| 50 | 0.04760 | 0.05050 | 0.06360 | 0.07150 | 0.09130 |
| 100 | 0.05040 | 0.05210 | 0.08340 | 0.09400 | 0.12940 |
| 250 | 0.04670 | 0.05330 | 0.14020 | 0.16240 | 0.25440 |
| 500 | 0.04680 | 0.05950 | 0.23000 | 0.28920 | 0.46490 |
| 1000 | 0.05050 | 0.07340 | 0.40940 | 0.50570 | 0.74310 |

Moderate Relationship between X_1 and Y : Var = 50%

| N | Correlation Between X_1 and X_2 | | | | |
|------|-------------------------------------|---------|---------|---------|---------|
| | 0.00 | 0.25 | 0.75 | 0.80 | 0.90 |
| 50 | 0.04600 | 0.05200 | 0.09630 | 0.11060 | 0.16330 |
| 100 | 0.05350 | 0.05690 | 0.14610 | 0.18570 | 0.28370 |
| 250 | 0.04830 | 0.06250 | 0.30680 | 0.37310 | 0.58640 |
| 500 | 0.05150 | 0.07800 | 0.53230 | 0.64880 | 0.88370 |
| 1000 | 0.04810 | 0.11850 | 0.82730 | 0.90880 | 0.99070 |

Strong Relationship between X_1 and Y : Var = 75%

| N | Correlation Between X_1 and X_2 | | | | |
|------|-------------------------------------|---------|---------|---------|---------|
| | 0.00 | 0.25 | 0.75 | 0.80 | 0.90 |
| 50 | 0.04850 | 0.05790 | 0.17270 | 0.20890 | 0.34420 |
| 100 | 0.05410 | 0.06790 | 0.31010 | 0.37850 | 0.60310 |
| 250 | 0.04790 | 0.08560 | 0.64500 | 0.75230 | 0.94340 |
| 500 | 0.04450 | 0.13230 | 0.91090 | 0.96350 | 0.99920 |
| 1000 | 0.05220 | 0.21790 | 0.99590 | 0.99980 | 1.00000 |

Marginal Mean Type I Error Probabilities

| Base Distribution | | | | |
|--|------------|------------|------------|------------|
| normal | Pareto | t Distr | uniform | |
| 0.38692448 | 0.36903077 | 0.38312245 | 0.38752571 | |
| | | | | |
| Explained Variance | | | | |
| 0.25 | 0.50 | 0.75 | | |
| 0.27330660 | 0.38473364 | 0.48691232 | | |
| | | | | |
| Correlation between Latent Independent Variables | | | | |
| 0.00 | 0.25 | 0.75 | 0.80 | 0.90 |
| 0.05004853 | 0.16604247 | 0.51544093 | 0.55050700 | 0.62621533 |
| | | | | |
| Sample Size n | | | | |
| 50 | 100 | 250 | 500 | 1000 |
| 0.19081740 | 0.27437227 | 0.39457933 | 0.48335707 | 0.56512820 |
| | | | | |
| Reliability of W_1 | | | | |
| 0.50 | 0.75 | 0.80 | 0.90 | 0.95 |
| 0.60637233 | 0.46983147 | 0.42065313 | 0.26685820 | 0.14453913 |
| | | | | |
| Reliability of W_2 | | | | |
| 0.50 | 0.75 | 0.80 | 0.90 | 0.95 |
| 0.30807933 | 0.37506733 | 0.38752793 | 0.41254800 | 0.42503167 |

Summary

- Ignoring measurement error in the explanatory variables can seriously inflate Type I error probabilities.

Summary

- Ignoring measurement error in the explanatory variables can seriously inflate Type I error probabilities.
- The poison combination is measurement error in the variable for which you are “controlling,” and correlation between latent explanatory variables.
- If either is zero, there is no problem.

$$\hat{\beta}_2 \xrightarrow{p} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

Summary

- Ignoring measurement error in the explanatory variables can seriously inflate Type I error probabilities.
- The poison combination is measurement error in the variable for which you are “controlling,” and correlation between latent explanatory variables.
- If either is zero, there is no problem.

$$\hat{\beta}_2 \xrightarrow{p} \frac{\beta_1 \omega_1 \phi_{12}}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2) - \phi_{12}^2}$$

- Factors affecting severity of the problem are (next slide)

Factors affecting severity of the problem

Problem of inflated Type I error probability

- As the correlation between X_1 and X_2 increases, the problem gets worse.
- As the correlation between X_1 and Y increases, the problem gets worse.
- As the amount of measurement error in X_1 increases, the problem gets worse.
- As the amount of measurement error in X_2 increases, the problem gets *less* severe.
- As the sample size increases, the problem gets worse.
- Distribution of the variables does not matter much.

As the sample size increases, the problem gets worse

For a large enough sample size, no amount of measurement error in the explanatory variables is safe, assuming that the latent explanatory variables are correlated.

Other kinds of regression, other kinds of measurement error

Other kinds of regression, other kinds of measurement error

- Logistic regression
- Proportional hazards regression in survival analysis
- Log-linear models: Test of conditional independence in the presence of classification error
- Median splits
- Even converting X_1 to ranks inflates Type I Error probability.

Moral of the story

Use models that allow for measurement error in the explanatory variables.

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The \LaTeX source code is available from the course website:

<http://www.utstat.toronto.edu/brunner/oldclass/431s23>