

Introduction to Regression with Measurement Error

STA431: Spring 2015

See last slide for copyright information

Measurement Error

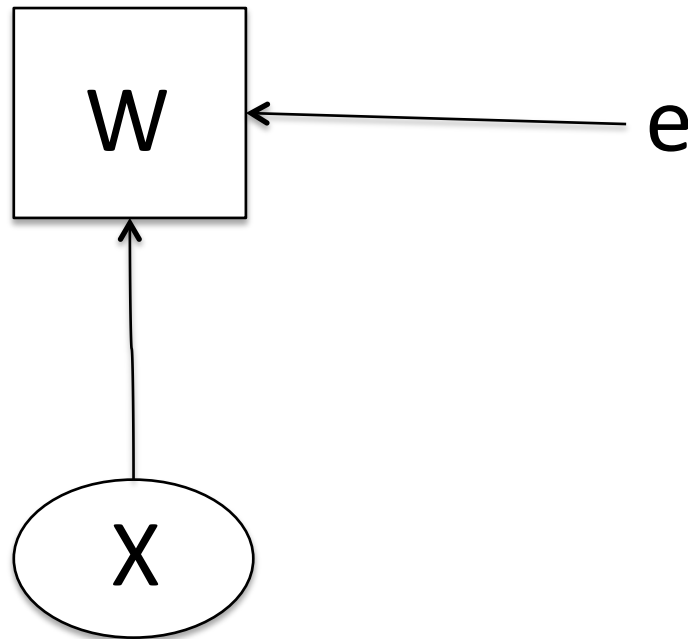
- Snack food consumption
- Exercise
- Income
- Cause of death
- Even amount of drug that reaches animal's blood stream in an experimental study
- Is there anything that is *not* measured with error?

For categorical variables

Classification error is common

Additive measurement error:

$$W = X + e$$



Simple additive model for measurement error: Continuous case

$$W = X + e$$

Where $E(X) = \mu$, $E(e) = 0$, $Var(X) = \sigma_X^2$, $Var(e) = \sigma_e^2$, and $Cov(X, e) = 0$.
Because X and e are uncorrelated,

$$Var(W) = Var(X) + Var(e) = \sigma_X^2 + \sigma_e^2$$

$$\begin{aligned} Cov(X, W) &= E(\overset{c}{\dot{X}} \overset{c}{\dot{W}}) \\ &= E(\overset{c}{\dot{X}} (\overset{c}{\dot{X}} + e)) \\ &= E(\overset{c}{\dot{X}}^2) + E(\overset{c}{\dot{X}})E(e) \\ &= \sigma_X^2 \end{aligned}$$

How much of the variation in the observed variable comes from variation in the quantity of interest, and how much comes from random noise?

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

First, recall

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

$$\text{Var}(X + a) = \text{Var}(X)$$

$$\text{Cov}(X + a, Y + b) = \text{Cov}(X, Y)$$

Reliability

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

$$(\text{Corr}(X, W))^2 = \frac{\sigma_X^2}{\sigma_X^2 + \sigma_e^2}$$

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

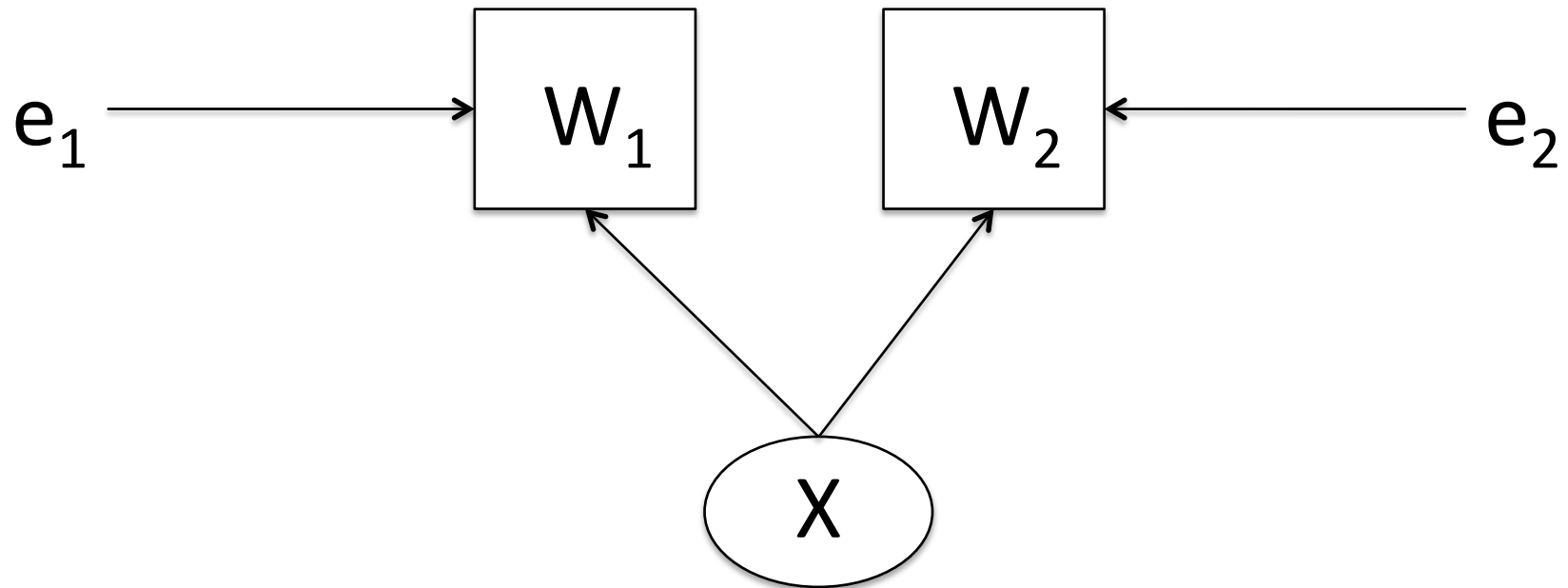
Correlate usual measurement with “Gold Standard?”

Not very realistic, except maybe for
some bio-markers

Measure twice

$$W_1 = X + e_1$$

$$W_2 = X + e_2,$$



Test-Retest

$$W_1 = X + e_1$$

$$W_2 = X + e_2,$$

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

Equivalent measurements

Test-Retest Reliability

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

$$\begin{aligned}\text{Cov}(W_1, W_2) &= \text{Cov}(\overset{c}{W}_1, \overset{c}{W}_2) \\ &= E(\overset{c}{W}_1 \overset{c}{W}_2) \\ &= E(\overset{c}{X} + e_1)(\overset{c}{X} + e_2) \\ &= E(\overset{c}{X}^2) + 0 + 0 + 0 \\ &= \sigma_X^2, \text{ so}\end{aligned}$$

$$\begin{aligned}\text{Corr}(W_1, W_2) &= \frac{\sigma_X^2}{\sqrt{\sigma_X^2 + \sigma_e^2} \sqrt{\sigma_X^2 + \sigma_e^2}} \\ &= \frac{\sigma_X^2}{\sigma_X^2 + \sigma_e^2}\end{aligned}$$

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

Calculate the sample correlation between

$$W_{1,1}, W_{2,1}, \dots, W_{n,1}$$

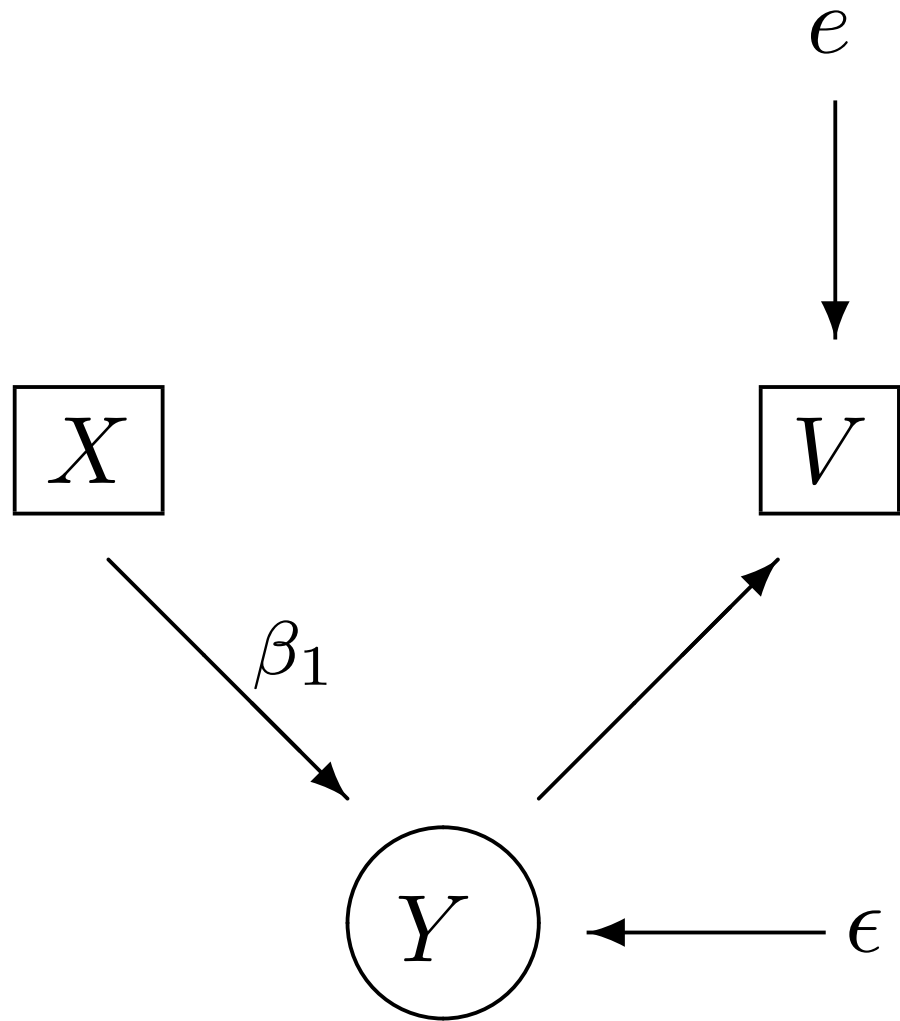
$$W_{1,2}, W_{2,2}, \dots, W_{n,2}$$

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

The consequences of ignoring measurement error in the explanatory (x) variables

First look at measurement error in
the response variable

Measurement error in the response variable



Measurement error in the response variable is a less serious problem:
Re-parameterize

$$Y = \beta_0 + \beta_1 X + \epsilon$$

$$V = \nu + Y + e$$

$$= \nu + (\beta_0 + \beta_1 X + \epsilon) + e$$

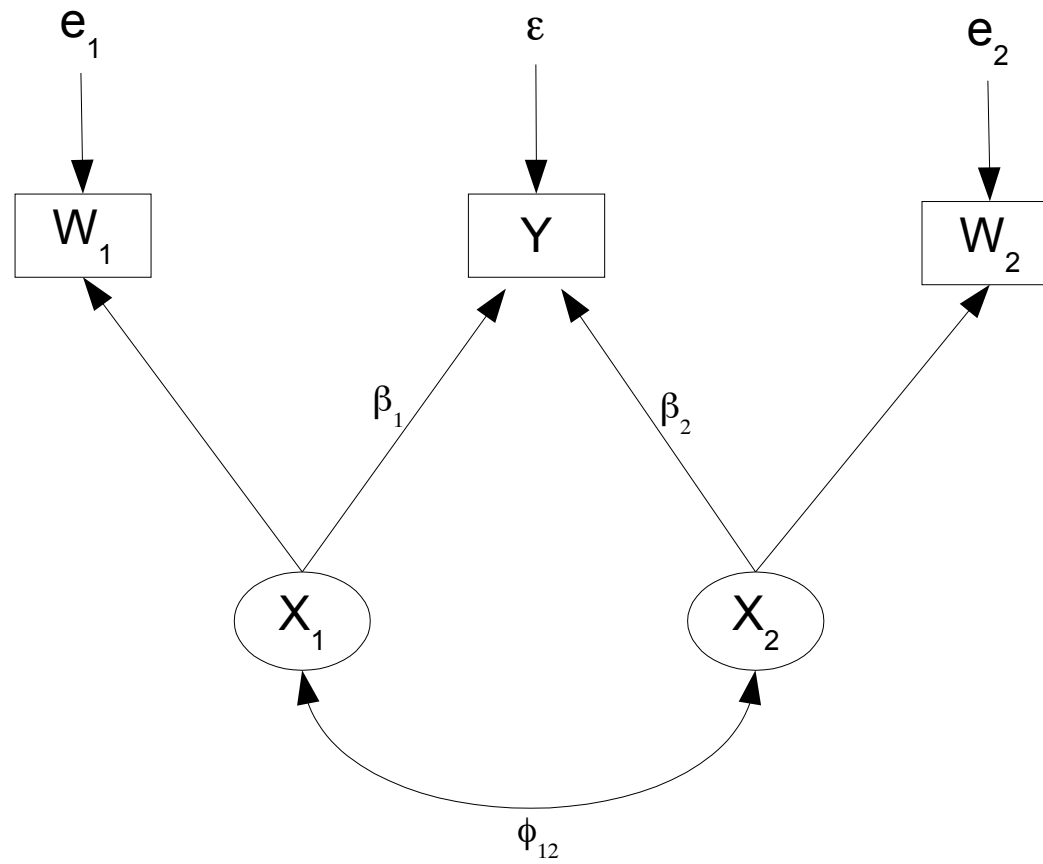
$$= (\nu + \beta_0) + \beta_1 X + (\epsilon + e)$$

$$= \beta'_0 + \beta_1 X + \epsilon'$$

Can't know everything, but all we care about is β_1 anyway.

Whenever 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 the explanatory variables



Measurement error in the explanatory variables

- True model

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i$$

$$W_{i,1} = X_{i,1} + e_{i,1}$$

$$W_{i,2} = X_{i,2} + e_{i,2}$$

- Naïve model

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

True Model (More detail)

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i$$

$$W_{i,1} = X_{i,1} + e_{i,1}$$

$$W_{i,2} = X_{i,2} + e_{i,2},$$

where independently for $i = 1, \dots, n$, $E(X_{i,1}) = \mu_1$, $E(X_{i,2}) = \mu_2$,
 $E(\epsilon_i) = E(e_{i,1}) = E(e_{i,2}) = 0$, $Var(\epsilon_i) = \sigma^2$, $Var(e_{i,1}) = \omega_1$,
 $Var(e_{i,2}) = \omega_2$, the errors ϵ_i , $e_{i,1}$ and $e_{i,2}$ are all independent,
 $X_{i,1}$ is independent of ϵ_i , $e_{i,1}$ and $e_{i,2}$,
 $X_{i,2}$ is independent of ϵ_i , $e_{i,1}$ and $e_{i,2}$, and

$$V \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}$$

Reliabilities

- Reliability of W_1 is $\frac{\phi_{11}}{\phi_{11} + \omega_1}$

- Reliability of W_2 is $\frac{\phi_{22}}{\phi_{22} + \omega_2}$

Test X_2 controlling for (holding constant) X_1

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\frac{\partial}{\partial x_2} E(Y) = \beta_2$$

That's the usual conditional model

Unconditional: Test X_2 controlling for X_1

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

$$\begin{aligned} \text{Cov}(X_2, Y) &= \beta_1 \text{Cov}(X_1, X_2) + \beta_2 \text{Var}(X_2) \\ &= \beta_1 \phi_{12} + \beta_2 \phi_{22} \end{aligned}$$

Hold X_1 constant at fixed x_1

$$\text{Cov}(X_2, Y | X_1 = x_1) = \beta_2 \text{Var}(X_2) = \beta_2 \phi_{22}$$

Controlling Type I Error Probability

- Type I error is to reject H_0 when it is true, and there is actually no effect or no relationship
- Type I error is very bad. Maybe that's why it's called an "error of the first kind."
- False knowledge is worse than ignorance.

Simulation study: Use pseudo-random number generation to create data sets

- Simulate data from the true model with $\beta_2=0$
- Fit naïve model
- Test $H_0: \beta_2=0$ at $\alpha = 0.05$ using naïve model
- Is H_0 rejected five percent of the time?

```
rmvn <- function(nn,mu,sigma)
# Returns an nn by kk matrix, rows are independent MVN(mu,sigma)
{
kk <- length(mu)
dsig <- dim(sigma)
if(dsig[1] != dsig[2]) stop("Sigma must be square.")
if(dsig[1] != kk) stop("Sizes of sigma and mu are inconsistent.")
ev <- eigen(sigma,symmetric=T)
sqr1 <- diag(sqrt(ev$values))
PP <- ev$vectors
ZZ <- rnorm(nn*kk) ; dim(ZZ) <- c(kk,nn)
rmvn <- t(PP*%sqr1*%ZZ+mu)
rmvn
}# End of function rmvn
```

```

merereg <- function(beta0=1, beta1=1, beta2=0, sigmasq = 0.5,
                    mu1=0, mu2=0, phi11=1, phi22=1, phi12 = 0.80,
                    rel1=0.80, rel2=0.80, n=200)
#####
# Model is      Y = beta0 + beta1 X1 + beta2 X2 + epsilon
#              W1 = X1 + e1
#              W2 = W2 + e2
# Fit naive model
#              Y = beta0 + beta1 W1 + beta2 W2 + epsilon
# Inputs are
#
#  beta0, beta1 beta2      True regression coefficients
#  sigmasq                 Var(epsilon)
#  mu1                     E(X1)
#  mu2                     E(X2)
#  phi11                   Var(X1)
#  phi22                   Var(X2)
#  phi12                   Cov(X1,X2) = Corr(X1,X1), because
#                          Var(X1) = Var(X2) = 1
#  rel1                    Reliability of W1
#  rel2                    Reliability of W2
#  n                       Sample size
# Note: This function uses rmvn, a multivariate normal random number
#       generator I wrote. The rmultnorm of the package MSBVAR does
#       the same thing but I am having trouble installing it.
#####

```

```

{
# Calculate SD(e1) and SD(e2)
sd1 <- sqrt((phi11-rel1)/rel1)
sd2 <- sqrt((phi22-rel2)/rel2)
# Random number generation
epsilon <- rnorm(n,mean=0,sd=sqrt(sigmasq))
e1 <- rnorm(n,mean=0,sd=sd1)
e2 <- rnorm(n,mean=0,sd=sd2)
# X1 and X2 are bivariate normal. Need rmvn function.
Phi <- rbind(c(phi11,phi12),
             c(phi12,phi22))
X <- rmvn(n, mu=c(mu1,mu2), sigma=Phi) # nx2 matrix
X1 <- X[,1]; X2 <- X[,2]
# Now generate Y, W1 and W2

Y = beta0 + beta1*X1 + beta2*X2 + epsilon
W1 = X1 + e1
W2 = X2 + e2

# Fit the naive model
merereg <- summary(lm(Y~W1+W2))$coefficients
merereg # Returns table of beta-hats, SEs, t-statistics and p-values
} # End function merereg

```

```

> merereg() # All the default values of inputs
              Estimate Std. Error  t value    Pr(>|t|)
(Intercept) 0.9704708 0.05423489 17.893845 3.692801e-43
W1           0.6486972 0.06336434 10.237576 5.385982e-20
W2           0.2079601 0.06201811  3.353216 9.578634e-04
>
> merereg()[3,4] # Just the p-value for H0: beta2=0
[1] 0.0006340172
>
> # H0 rejected twice. Is the function okay?
> merereg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.03946133
> merereg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.2582209
> merereg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.08474088
> merereg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.5182614
> merereg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.2889913

```

```
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.1667587
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.4414364
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.2268087
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.8298779
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.3508289
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.05173589
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.243059
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.8818203
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.3430994
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.4860574
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.9644776
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.09245873
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.04757209
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.7947851
> merreg(rel1=1,rel2=1)[3,4] # No measurement error
[1] 0.8039931
```

Try it with measurement error

```
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.01080889
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.0007349183
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.01884786
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.003615565
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.003421935
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 3.895541e-07
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 3.328842e-07
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.0754436
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 0.0001274642
> merreg()[3,4] # Reliabilities both equal 0.80
[1] 6.900713e-05
```


A **Big** Simulation Study (6 Factors)

- Sample size: $n = 50, 100, 250, 500, 1000$
- $\text{Corr}(X_1, X_2): \phi_{12} = 0.00, 0.25, 0.75, 0.80, 0.90$
- 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

- $5 \times 5 \times 3 \times 5 \times 5 \times 4 = 7,500$ treatment combinations

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 naïve 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

Look at a small part of the results

- Both reliabilities = 0.90
- Everything is normally distributed
- $\beta_0 = 1, \beta_1=1, \beta_2=0$ (H_0 is true)

Weak Relationship between X_1 and Y : 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 independent 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 independent variables. If either is zero, there is no problem.
- Factors affecting severity of the problem are (next slide)

Factors affecting severity of the problem

- 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 independent variables is safe, assuming that the latent independent variables are correlated.

The problem applies to other kinds of regression, and various 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

If X_1 is randomly assigned

- Then it is independent of X_2 : Zero correlation.
- So even if an experimentally manipulated variable is measured (implemented) with error, there will be no inflation of Type I error probability.
- If X_2 is randomly assigned and X_1 is a covariate observed with error (very common), then again there is no correlation between X_1 and X_2 , and so no inflation of Type I error probability.
- Measurement error may decrease the precision of experimental studies, but in terms of Type I error it creates no problems.
- This is good news!

What is going on theoretically?

First, need to look at some large-sample tools

Sample Space Ω , ω an element of Ω

- Observing whether a single individual is male or female:

$$\Omega = \{F, M\}$$

- Pair of individuals and observed their genders in order:

$$\Omega = \{(F, F), (F, M), (M, F), (M, M)\}$$

- Select n people and count the number of females:

$$\Omega = \{0, \dots, n\}$$

- For limits problems, the points in Ω are infinite sequences

Random variables are functions from Ω into the set of real numbers

$$\Pr\{X \in B\} = \Pr(\{\omega \in \Omega : X(\omega) \in B\})$$

Random sample $X_1(\omega), \dots, X_n(\omega)$

$$T = T(X_1, \dots, X_n)$$

$$T = T_n(\omega)$$

Let $n \rightarrow \infty$

To see what happens for large samples

Modes of Convergence

- Almost Sure Convergence
- Convergence in Probability
- Convergence in Distribution

Almost Sure Convergence

We say that T_n converges *almost surely* to T , and write $T_n \xrightarrow{a.s.}$ if

$$\Pr\{\omega : \lim_{n \rightarrow \infty} T_n(\omega) = T(\omega)\} = 1.$$

Acts like an ordinary limit, except possibly on a set of probability zero.

All the usual rules apply.

Strong Law of Large Numbers

$$\overline{X}_n \xrightarrow{a.s.} \mu$$

The only condition required for this to hold is the existence of the expected value.

Let X_1, \dots, X_n be independent and identically distributed random variables; let X be a general random variable from this same distribution, and $Y=g(X)$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n g(X_i) &= \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{a.s.} E(Y) \\ &= E(g(X)) \end{aligned}$$

So for example

$$\frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{a.s.} E(X^k)$$

$$\frac{1}{n} \sum_{i=1}^n U_i^2 V_i W_i^3 \xrightarrow{a.s.} E(U^2 V W^3)$$

That is, sample moments converge almost surely to population moments.

Convergence in Probability

We say that T_n converges *in probability* to T , and write $T_n \xrightarrow{P} T$ if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P\{|T_n - T| < \epsilon\} = 1$$

Almost Sure Convergence \Rightarrow Convergence in Probability

Strong Law of Large Numbers \Rightarrow Weak Law of Large Numbers

Convergence in Distribution

Denote the cumulative distribution functions of T_1, T_2, \dots by $F_1(t), F_2(t), \dots$ respectively, and denote the cumulative distribution function of T by $F(t)$.

We say that T_n converges *in distribution* to T , and write $T_n \xrightarrow{d} T$ if for every point t at which F is continuous,

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

Central Limit Theorem says

$$Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z \sim N(0, 1)$$

Connections among the Modes of Convergence

- $T_n \xrightarrow{a.s.} T \Rightarrow T_n \xrightarrow{P} T \Rightarrow T_n \xrightarrow{d} T.$
- If a is a constant, $T_n \xrightarrow{d} a \Rightarrow T_n \xrightarrow{P} a.$

Consistency

$T_n = T_n(X_1, \dots, X_n)$ is a statistic estimating a parameter θ

The statistic T_n is said to be *consistent* for θ if $T_n \xrightarrow{P} \theta$.

$$\lim_{n \rightarrow \infty} P\{|T_n - \theta| < \epsilon\} = 1$$

The statistic T_n is said to be *strongly consistent* for θ if $T_n \xrightarrow{a.s.} \theta$.

Strong consistency implies ordinary consistency.

Consistency is great but it's not enough

- It means that as the sample size becomes indefinitely large, you (probably) get as close as you like to the truth.
- It's the least we can ask. Estimators that are not consistent are completely unacceptable for most purposes.

$$T_n \xrightarrow{a.s.} \theta \Rightarrow U_n = T_n + \frac{100,000,000}{n} \xrightarrow{a.s.} \theta$$

Consistency of the Sample Variance

$$\begin{aligned}\hat{\sigma}_n^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2\end{aligned}$$

By SLLN, $\bar{X}_n \xrightarrow{a.s.} \mu$ and $\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{a.s.} E(X^2) = \sigma^2 + \mu^2$

Because the function $g(x, y) = x - y^2$ is continuous,

$$\hat{\sigma}_n^2 = g\left(\frac{1}{n} \sum_{i=1}^n X_i^2, \bar{X}_n\right) \xrightarrow{a.s.} g(\sigma^2 + \mu^2, \mu) = \sigma^2 + \mu^2 - \mu^2 = \sigma^2$$

Consistency of the Sample Covariance

$$\hat{\sigma}_{1,2} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{X}_n \bar{Y}_n$$

By SLLN, $\bar{X}_n \xrightarrow{a.s.} E(X)$, $\bar{Y}_n \xrightarrow{a.s.} E(Y)$, and $\frac{1}{n} \sum_{i=1}^n X_i Y_i \xrightarrow{a.s.} E(XY)$

Because the function $g(x, y, z) = x - yz$ is continuous,

$$\begin{aligned} \hat{\sigma}_{1,2} &= g\left(\frac{1}{n} \sum_{i=1}^n X_i Y_i, \bar{X}_n, \bar{Y}_n\right) \xrightarrow{a.s.} g(E(XY), E(X), E(Y)) \\ &= E(XY) - E(X)E(Y) = Cov(X, Y) \\ &= \sigma_{1,2} \end{aligned}$$

MOM is consistent, usually

$$m = f(\theta)$$

$$\theta = g^{-1}(m)$$

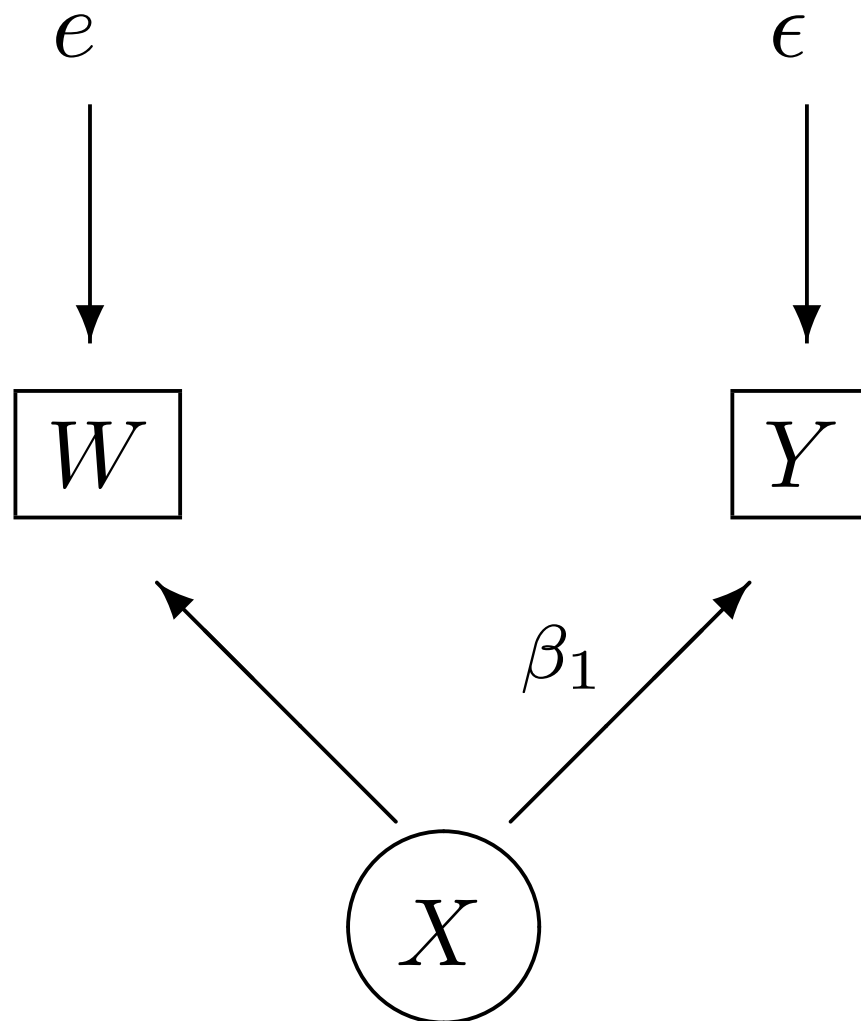
$$\hat{\theta}_n = g^{-1}(\hat{m}_n)$$

By SLLN, $\hat{m}_n \xrightarrow{a.s.} m$

By continuous mapping, $\hat{\theta}_n = g^{-1}(\hat{m}_n) \xrightarrow{a.s.} g^{-1}(m) = \theta$

Provided g^{-1} is continuous at the true parameter value.

True Regression model: Single explanatory variable measured with error



Single Explanatory Variable

- True model

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

- Naive model

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

where independently for $i = 1, \dots, n$, $Var(X_i) = \sigma_X^2$, $Var(e_i) = \sigma_e^2$, and X_i, e_i, ϵ_i are all independent.

Least squares estimate of β_1 for the Naïve Model

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

$$= \frac{\hat{\sigma}_{w,y}}{\hat{\sigma}_w^2}$$

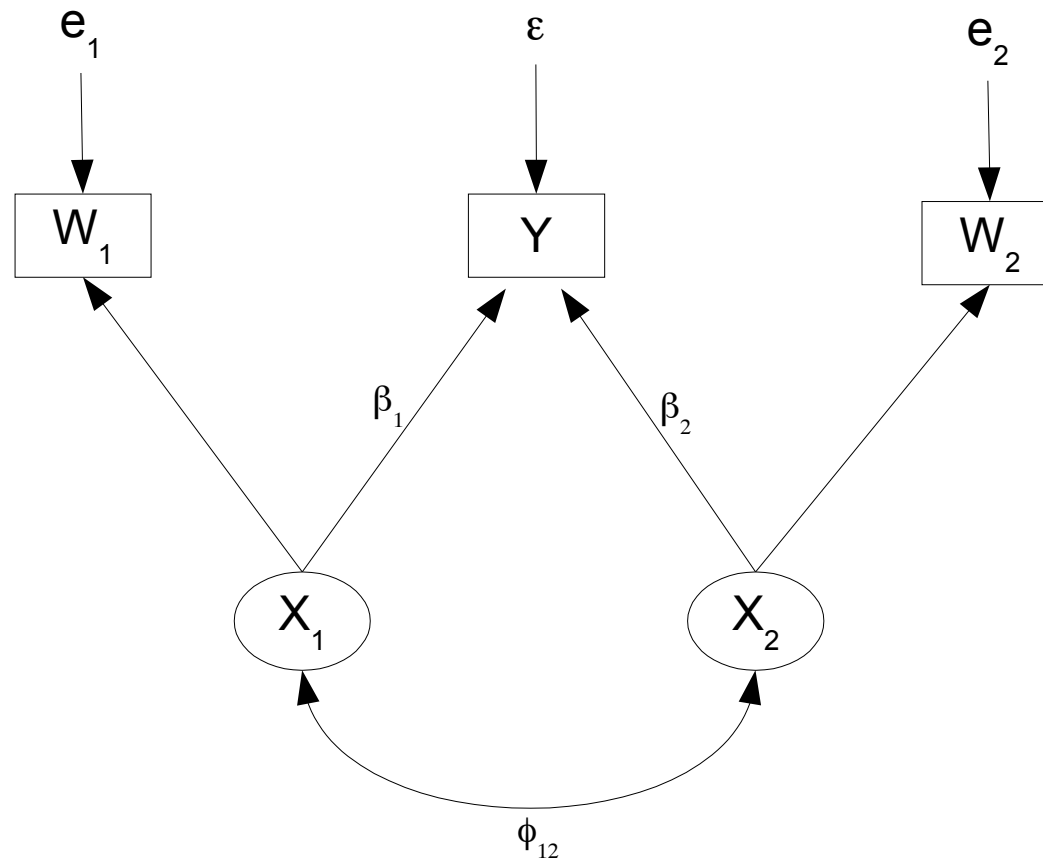
$$\xrightarrow{a.s.} \frac{Cov(W, Y)}{Var(W)}$$

$$= \beta_1 \left(\frac{\sigma_X^2}{\sigma_X^2 + \sigma_e^2} \right)$$

$$\hat{\beta}_1 \xrightarrow{a.s.} \beta_1 \left(\frac{\sigma_X^2}{\sigma_X^2 + \sigma_e^2} \right)$$

- Goes to the true parameter times reliability of W .
- Asymptotically biased toward zero, because reliability is between zero and one.
- No asymptotic bias when $\beta_1=0$.
- No inflation of Type I error probability
- Loss of power when $\beta_1 \neq 0$
- Measurement error just makes relationship seem weaker than it is. Reassuring, but watch out!

Two explanatory variables with error



Two explanatory variables, $\beta_2=0$

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i$$

$$W_{i,1} = X_{i,1} + e_{i,1}$$

$$W_{i,2} = X_{i,2} + e_{i,2},$$

where independently for $i = 1, \dots, n$, $E(X_{i,1}) = \mu_1$, $E(X_{i,2}) = \mu_2$,
 $E(\epsilon_i) = E(e_{i,1}) = E(e_{i,2}) = 0$, $Var(\epsilon_i) = \sigma^2$, $Var(e_{i,1}) = \omega_1$,
 $Var(e_{i,2}) = \omega_2$, the errors ϵ_i , $e_{i,1}$ and $e_{i,2}$ are all independent,
 $X_{i,1}$ is independent of ϵ_i , $e_{i,1}$ and $e_{i,2}$,
 $X_{i,2}$ is independent of ϵ_i , $e_{i,1}$ and $e_{i,2}$, and

$$V \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}$$

Least squares estimate of β_2 for the Naïve Model when true $\beta_2 = 0$

$$\begin{aligned}\hat{\beta}_2 &\xrightarrow{a.s.} \frac{\beta_1 \phi_{1,2} \omega_1}{(\phi_{1,1} + \omega_1)(\phi_{2,2} + \omega_2)} \\ &= \left(\frac{\omega_1}{\phi_{1,1} + \omega_1} \right) \left(\frac{\beta_1 \phi_{1,2}}{\phi_{2,2} + \omega_2} \right)\end{aligned}$$

Combined with estimated standard error going almost surely to zero,
Get t statistic for $H_0: \beta_2 = 0$ going to $\pm\infty$, and p-value going almost
Surely to zero, unless

Combined with estimated standard error going almost surely to zero, get t statistic for $H_0: \beta_2 = 0$ going to $\pm\infty$, and p-value going almost surely to zero, unless

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

$$\widehat{\beta}_2 \xrightarrow{a.s.} \left(\frac{\omega_1}{\phi_{1,1} + \omega_1} \right) \left(\frac{\beta_1 \phi_{1,2}}{\phi_{2,2} + \omega_2} \right)$$

And, anything that increases $Var(W_2)$ will make the problem less severe.

Need a statistical model that
includes measurement error

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. These Powerpoint slides are available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/431s15>