

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 Rates

	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 rates.
- 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 rate

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 rate.
- 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 rate.
- 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,

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

Single Independent 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 rate
- Loss of power when $\beta_1 \neq 0$
- Measurement error just makes relationship seem weaker than it is. Reassuring, but watch out!

Two Independent 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

$$Var \begin{bmatrix} X_{i,1} \\ X_{i,1} \end{bmatrix} = \begin{bmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{bmatrix}$$

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

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

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 decrease the bias.

Need a statistical model that
includes measurement error

First, random vectors and matrices
(see Appendix A)