

STA 2101 Assignment 6¹

The questions on this assignment are not to be handed in. They are practice for Quiz Six on Friday November 1st.

1. In the following regression model, the explanatory variables X_1 and X_2 are random variables. The true model is

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

independently for $i = 1, \dots, n$, where $\epsilon_i \sim N(0, \sigma^2)$.

The mean and covariance matrix of the explanatory variables are given by

$$E \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \text{and} \quad \text{Var} \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}$$

The explanatory variables $X_{i,1}$ and $X_{i,2}$ are independent of ϵ_i .

Unfortunately $X_{i,2}$, which has an impact on Y_i and is correlated with $X_{i,1}$, is not part of the data set. Since $X_{i,2}$ is not observed, it is absorbed by the intercept and error term, as follows.

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\ &= (\beta_0 + \beta_2 \mu_2) + \beta_1 X_{i,1} + (\beta_2 X_{i,2} - \beta_2 \mu_2 + \epsilon_i) \\ &= \beta'_0 + \beta_1 X_{i,1} + \epsilon'_i. \end{aligned}$$

The primes just denote a new β_0 and a new ϵ_i . It was necessary to add and subtract $\beta_2 \mu_2$ in order to obtain $E(\epsilon'_i) = 0$. And of course there could be more than one omitted variable. They would all get swallowed by the intercept and error term, the garbage bins of regression analysis.

- (a) What is $Cov(X_{i,1}, \epsilon'_i)$?
- (b) Calculate the variance-covariance matrix of $(X_{i,1}, Y_i)$ under the true model. Is it possible to have non-zero covariance between $X_{i,1}$ and Y_i when $\beta_1 = 0$?
- (c) Suppose we want to estimate β_1 . The usual least squares estimator is

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_{i,1} - \bar{X}_1)(Y_i - \bar{Y})}{\sum_{i=1}^n (X_{i,1} - \bar{X}_1)^2}.$$

You may just use this formula; you don't have to derive it. Is $\hat{\beta}_1$ a consistent estimator of β_1 if the true model holds? Answer Yes or no and show your work. You may use the consistency of the sample variance and covariance without proof.

¹This assignment 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 L^AT_EX source code is available from the course website: <http://www.utstat.toronto.edu/~brunner/oldclass/2101f19>

- (d) Are there *any* points in the parameter space for which $\widehat{\beta}_1 \xrightarrow{p} \beta_1$ when the true model holds?
2. Independently for $i = 1, \dots, n$, let $Y_i = \beta X_i + \epsilon_i$, where $X_i \sim N(\mu, \sigma_x^2)$ and $\epsilon_i \sim N(0, \sigma_\epsilon^2)$. Because of omitted variables that influence both X_i and Y_i , we have $Cov(X_i, \epsilon_i) = c \neq 0$.
- (a) The least squares estimator of β is $\frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}$. Is this estimator consistent? Answer Yes or No and prove your answer.
- (b) Give the parameter space for this model. There are some constraints on c .
- (c) First consider points in the parameter space where $\mu \neq 0$. Give an estimator of β that converges almost surely to the right answer for that part of the parameter space. If you are not sure how to proceed, try calculating the expected value and covariance matrix of (X_i, Y_i) .
- (d) What happens in the rest of the parameter space — that is, where $\mu = 0$? Is a consistent estimator possible there? So we see that parameters may be identifiable in some parts of the parameter space but not all.
3. We know that omitted explanatory variables are a big problem, because they induce non-zero covariance between the explanatory variables and the error terms ϵ_i . The residuals have a lot in common with the ϵ_i terms in a regression model, though they are not the same thing. A reasonable idea is to check for correlation between explanatory variables and the ϵ_i values by looking at the correlation between the residuals and explanatory variables.

Accordingly, for a multiple regression model with an intercept so that $\sum_{i=1}^n e_i = 0$, calculate the sample correlation r between explanatory variable j and the residuals e_1, \dots, e_n . Use this formula for the correlation: $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$. Simplify. What can the sample correlations between residuals and x variables tell you about the correlation between ϵ and the x variables?

4. This question explores the consequences of ignoring measurement error in the response variable. Independently for $i = 1, \dots, n$, let

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

where $Var(X_i) = \phi$, $E(X_i) = \mu_x$, $Var(e_i) = \omega$, $Var(\epsilon_i) = \psi$, and X_i, e_i, ϵ_i are all independent. The explanatory variable X_i is observable, but the response variable Y_i is latent. Instead of Y_i , we can see V_i , which is Y_i plus a piece of random noise. Call this the *true model*.

- (a) Make a path diagram of the true model.

- (b) Strictly speaking, the distributions of X_i, e_i and ϵ_i are unknown parameters because they are unspecified. But suppose we are interested in identifying just the Greek-letter parameters. Does the true model pass the test of the Parameter Count Rule? Answer Yes or No and give the numbers.
- (c) Calculate the variance-covariance matrix of the observable variables as a function of the model parameters. Show your work.
- (d) Suppose that the analyst assumes that V_i is that same thing as Y_i , and fits the naive model $V_i = \beta_0 + \beta_1 X_i + \epsilon_i$, in which

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

Assuming the *true* model (not the naive model), is $\hat{\beta}_1$ a consistent estimator of β_1 ? Answer Yes or No and show your work.

- (e) Why does this prove that β_1 is identifiable?
5. This question explores the consequences of ignoring measurement error in the explanatory variable when there is only one explanatory variable. Independently for $i = 1, \dots, n$, let

$$\begin{aligned} Y_i &= \beta X_i + \epsilon_i \\ W_i &= X_i + e_i \end{aligned}$$

where all random variables are normal with expected value zero, $Var(X_i) = \phi > 0$, $Var(\epsilon_i) = \psi > 0$, $Var(e_i) = \omega > 0$ and ϵ_i, e_i and X_i are all independent. The variables W_i and Y_i are observable, while X_i is latent. Error terms are never observable.

- (a) What is the parameter vector θ for this model?
- (b) Denote the covariance matrix of the observable variables by $\Sigma = [\sigma_{ij}]$. The unique σ_{ij} values are the moments, and there is a covariance structure equation for each one. Calculate the variance-covariance matrix Σ of the observable variables, expressed as a function of the model parameters. You now have the covariance structure equations.
- (c) Does this model pass the test of the parameter count rule? Answer Yes or No and give the numbers.
- (d) Are there any points in the parameter space where the parameter β is identifiable? Are there infinitely many, or just one point?
- (e) The naive estimator of β is

$$\hat{\beta}_n = \frac{\sum_{i=1}^n W_i Y_i}{\sum_{i=1}^n W_i^2}.$$

Is $\hat{\beta}_n$ a consistent estimator of β ? Answer Yes or No. To what does $\hat{\beta}_n$ converge?

- (f) Are there any points in the parameter space for which $\widehat{\beta}_n$ converges to the right answer? Compare your answer to the set of points where β is identifiable.
- (g) Suppose the reliability of W_i were known², or to be more realistic, suppose that a good estimate of the reliability were available; call it r_{wx}^2 . How could you use r_{wx}^2 to improve $\widehat{\beta}_n$? Give the formula for an improved estimator of β .
6. The improved version of $\widehat{\beta}_n$ in the last question is an example of *correction for attenuation* (weakening) caused by measurement error. Here is the version that applies to correlation. Independently for $i = 1, \dots, n$, let

$$\begin{aligned} D_{i,1} &= F_{i,1} + e_{i,1} \\ D_{i,2} &= F_{i,2} + e_{i,2} \end{aligned} \quad \text{cov} \begin{pmatrix} F_{i,1} \\ F_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix} \quad \text{cov} \begin{pmatrix} e_{i,1} \\ e_{i,2} \end{pmatrix} = \begin{pmatrix} \omega_1 & 0 \\ 0 & \omega_2 \end{pmatrix}$$

To make this concrete, it would be natural for psychologists to be interested in the correlation between intelligence and self-esteem, but what they want to know is the correlation between *true* intelligence and *true* self-esteem, not just the between score on an IQ test and score on a self-esteem questionnaire. So for subject i , let $F_{i,1}$ represent true intelligence and $F_{i,2}$ represent true self-esteem, while $D_{i,1}$ is the subject's score on an intelligence test and $D_{i,2}$ is score on a self-esteem questionnaire.

- (a) Make a path diagram of this model.
- (b) Show that $|\text{Corr}(D_{i,1}, D_{i,2})| \leq |\text{Corr}(F_{i,1}, F_{i,2})|$. That is, measurement error weakens (attenuates) the correlation.
- (c) Suppose the reliability of $D_{i,1}$ is ρ_1^2 and the reliability of $D_{i,2}$ is ρ_2^2 . How could you apply ρ_1^2 and ρ_2^2 to $\text{Corr}(D_{i,1}, D_{i,2})$, to obtain $\text{Corr}(F_{i,1}, F_{i,2})$?
- (d) You obtain a sample correlation between IQ score and self-esteem score of $r = 0.25$, which is disappointingly low. From other data, the estimated reliability of the IQ test is $r_1^2 = 0.90$, and the estimated reliability of the self-esteem scale is $r_2^2 = 0.75$. Give an estimate of the correlation between true intelligence and true self-esteem. The answer is a number.
7. This is a simplified version of the situation where one is attempting to “control” for explanatory variables that are measured with error. People do this all the time, and it doesn't work. Independently for $i = 1, \dots, n$, let

$$\begin{aligned} Y_i &= \beta_1 X_{i,1} + \beta_2 X_{i,2} + \epsilon_i \\ W_i &= X_{i,1} + e_i, \end{aligned}$$

²As a reminder, the reliability of an observed measurement is the proportion of its variance that comes from the “true” latent variable it is measuring. Here, the reliability of W_i is $\frac{\phi}{\phi + \omega}$.

where $V \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}$, $V(\epsilon_i) = \psi$, $V(e_1) = \omega$, all the expected values are zero, and the error terms ϵ_i and e_i are independent of one another, and also independent of $X_{i,1}$ and $X_{i,2}$. The variable $X_{i,1}$ is latent, while the variables W_i , Y_i and $X_{i,2}$ are observable. What people usually do in situations like this is fit a model like $Y_i = \beta_1 W_i + \beta_2 X_{i,2} + \epsilon_i$, and test $H_0 : \beta_2 = 0$. That is, they ignore the measurement error in variables for which they are “controlling.”

(a) Suppose $H_0 : \beta_2 = 0$ is true. Does the ordinary least squares estimator

$$\widehat{\beta}_2 = \frac{\sum_{i=1}^n W_i^2 \sum_{i=1}^n X_{i,2} Y_i - \sum_{i=1}^n W_i X_{i,2} \sum_{i=1}^n W_i Y_i}{\sum_{i=1}^n W_i^2 \sum_{i=1}^n X_{i,2}^2 - (\sum_{i=1}^n W_i X_{i,2})^2}$$

converge to the true value of $\beta_2 = 0$ as $n \rightarrow \infty$ everywhere in the parameter space? Answer Yes or No and show your work.

(b) Under what conditions (that is, for what values of other parameters) does $\widehat{\beta}_2 \xrightarrow{p} 0$ when $\beta_2 = 0$?

8. Finally we have a solution, though as usual there is a little twist. Independently for $i = 1, \dots, n$, let

$$\begin{aligned} Y_i &= \beta X_i + \epsilon_i \\ V_i &= Y_i + e_i \\ W_{i,1} &= X_i + e_{i,1} \\ W_{i,2} &= X_i + e_{i,2} \end{aligned}$$

where

- Y_i is a latent variable.
- V_i , $W_{i,1}$ and $W_{i,2}$ are all observable variables.
- X_i is a normally distributed *latent* variable with mean zero and variance $\phi > 0$.
- ϵ_i is normally distributed with mean zero and variance $\psi > 0$.
- e_i is normally distributed with mean zero and variance $\omega > 0$.
- $e_{i,1}$ is normally distributed with mean zero and variance $\omega_1 > 0$.
- $e_{i,2}$ is normally distributed with mean zero and variance $\omega_2 > 0$.
- X_i , ϵ_i , e_i , $e_{i,1}$ and $e_{i,2}$ are all independent of one another.

(a) Make a path diagram of this model.

(b) What is the parameter vector θ for this model?

(c) Does the model pass the test of the Parameter Count Rule? Answer Yes or No and give the numbers.

- (d) Calculate the variance-covariance matrix of the observable variables as a function of the model parameters. Show your work.
- (e) Is the parameter vector identifiable at every point in the parameter space? Answer Yes or No and prove your answer.
- (f) Some parameters are identifiable, while others are not. Which ones are identifiable?
- (g) If β (the parameter of main interest) is identifiable, propose a Method of Moments estimator for it and prove that your proposed estimator is consistent.
- (h) Suppose the sample variance-covariance matrix $\hat{\Sigma}$ is

	W1	W2	V
W1	38.53	21.39	19.85
W2	21.39	35.50	19.00
V	19.85	19.00	28.81

Give a reasonable estimate of β . There is more than one right answer. The answer is a number. (Is this the Method of Moments estimate you proposed? It does not have to be.) **Circle your answer.**

- (i) Describe how you could re-parameterize this model to make the parameters all identifiable, allowing you do maximum likelihood.