

STA 2053 Assignment 3 (Random vectors and measurement error)¹

$$\text{cov}(\mathbf{x}) = E \{(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top\} \quad \text{cov}(\mathbf{x}, \mathbf{y}) = E \{(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top\}$$

These questions are not to be handed in. They are practice for the quiz on October 31st.

1. Let \mathbf{x} be a random vector, and let \mathbf{A} and \mathbf{B} be matrices of constants (of the right dimensions).

(a) Show $\text{cov}(\mathbf{Ax}) = \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^\top$

(b) Show $\text{cov}(\mathbf{Ax}, \mathbf{Bx}) = \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{B}^\top$.

2. Let \mathbf{x} be a $p \times 1$ random vector and let \mathbf{y} and \mathbf{z} be $q \times 1$ random vectors. Show that $\text{cov}(\mathbf{x}, \mathbf{y} + \mathbf{z}) = \text{cov}(\mathbf{x}, \mathbf{y}) + \text{cov}(\mathbf{x}, \mathbf{z})$.

3. Let

$$\mathbf{x}_i = \begin{pmatrix} x_{i,1} \\ \vdots \\ x_{i,p} \end{pmatrix} \quad \text{and} \quad \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix}.$$

Let the $p \times p$ matrix $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$. Give a *scalar* formula for element (2, 3) of $\hat{\boldsymbol{\Sigma}}$. If you get stuck, an example with $p = 3$ should help.

4. Let \mathbf{x} and \mathbf{y} be random vectors, and let \mathbf{A} and \mathbf{B} be matrices of constants. Starting from the definition, find $\text{cov}(\mathbf{Ax}, \mathbf{By})$. Show your work. Of course we are assuming that the matrices are the right size.
5. Denote the centered version of a general random vector \mathbf{y} by $\overset{c}{\mathbf{y}} = \mathbf{y} - \boldsymbol{\mu}_y$. Let $\mathbf{L} = \mathbf{A}_1\mathbf{x}_1 + \cdots + \mathbf{A}_m\mathbf{x}_m + \mathbf{b}$, where the \mathbf{A}_j are matrices of constants, and \mathbf{b} is a vector of constants. Show $\overset{c}{\mathbf{L}} = \mathbf{A}_1 \overset{c}{\mathbf{x}}_1 + \cdots + \mathbf{A}_m \overset{c}{\mathbf{x}}_m$.
6. Let \mathbf{x} and \mathbf{y} be $p \times 1$ random vectors. State whether the following is true or false, and show your work. $\text{cov}(\mathbf{x} + \mathbf{y}) = \text{cov}(\mathbf{x}) + \text{cov}(\mathbf{y}) + 2\text{cov}(\mathbf{x}, \mathbf{y})$.

-
7. Suppose we have two equivalent measurements with uncorrelated measurement error:

$$\begin{aligned} W_1 &= X + e_1 \\ W_2 &= X + e_2, \end{aligned}$$

where $E(X) = \mu_x$, $\text{Var}(X) = \sigma_x^2$, $E(e_1) = E(e_2) = 0$, $\text{Var}(e_1) = \text{Var}(e_2) = \sigma_e^2$, and X , e_1 and e_2 are all independent. What if we were to measure the true score X by adding the two imperfect measurements together? Would the result be more reliable?

¹This assignment was prepared by Jerry Brunner, Department of Statistical Sciences, University of Toronto. It is licensed under a [Creative Commons Attribution - ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-sa/3.0/). 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/2053f22>

- (a) Let $S = W_1 + W_2$. Calculate the reliability of S . Recall that reliability is defined as the squared correlation between the true score and the surface measurement.
- (b) Suppose you take n independent measurements (in psychometric theory, these would be called equivalent test items). What is the reliability of $S = \sum_{i=1}^n W_i$? Show your work.
- (c) What is the reliability of $\bar{W}_n = \frac{1}{n} \sum_{i=1}^n W_i$? Show your work.
- (d) What happens to the reliability of S and \bar{W}_n as the number of measurements $n \rightarrow \infty$?
8. 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 $\boldsymbol{\theta}$ for this model?
- (b) Denote the covariance matrix of the observable variables by $\boldsymbol{\Sigma} = [\sigma_{ij}]$. The unique σ_{ij} values are the moments, and there is a covariance structure equation for each one. Calculate the variance-covariance matrix $\boldsymbol{\Sigma}$ 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 ? Why can you answer this question without doing any calculations?

- (f) Go ahead and do the calculation. To what does $\hat{\beta}_n$ converge?
- (g) Are there any points in the parameter space for which $\hat{\beta}_n$ converges to the right answer? Compare your answer to the set of points where β is identifiable.
- (h) Suppose the reliability of W_i were known, 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 $\hat{\beta}_n$? Give the formula for an improved estimator of β .

9. The improved version of $\hat{\beta}_n$ 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 $D_{i,1}$ is score on a self-esteem questionnaire.

- Make a path diagram of this model.
 - Show that $|Corr(D_{i,1}, D_{i,2})| \leq |Corr(F_{i,1}, F_{i,2})|$. That is, measurement error weakens (attenuates) the correlation.
 - Suppose the reliability of $D_{i,1}$ is ρ_1^2 and the reliability of $D_{i,2}$ is ρ_2^2 . If you knew these values, how could you apply ρ_1^2 and ρ_2^2 to $Corr(D_{i,1}, D_{i,2})$, to obtain $Corr(F_{i,1}, F_{i,2})$?
 - 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.
10. 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

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

where $\text{cov} \begin{pmatrix} X_{i,1} \\ X_{i,2} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{12} & \phi_{22} \end{pmatrix}$, $V(\epsilon_i) = \psi$, $V(e_i) = \omega$, 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.”

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

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

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 $\beta_2 = 0$?

11. 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 $\widehat{\Sigma}$ 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.
12. Here is a one-stage formulation of the double measurement regression model. Independently for $i = 1, \dots, n$, let

$$\begin{aligned}\mathbf{w}_{i,1} &= \mathbf{x}_i + \mathbf{e}_{i,1} \\ \mathbf{v}_{i,1} &= \mathbf{y}_i + \mathbf{e}_{i,2} \\ \mathbf{w}_{i,2} &= \mathbf{x}_i + \mathbf{e}_{i,3}, \\ \mathbf{v}_{i,2} &= \mathbf{y}_i + \mathbf{e}_{i,4}, \\ \mathbf{y}_i &= \boldsymbol{\beta}\mathbf{x}_i + \boldsymbol{\epsilon}_i\end{aligned}$$

where

\mathbf{y}_i is a $q \times 1$ random vector of latent response variables. Because q can be greater than one, the regression is multivariate.

$\boldsymbol{\beta}$ is an $q \times p$ matrix of unknown constants. These are the regression coefficients, with one row for each response variable and one column for each explanatory variable.

\mathbf{x}_i is a $p \times 1$ random vector of latent explanatory variables, with expected value zero and variance-covariance matrix $\boldsymbol{\Phi}_x$, a $p \times p$ symmetric and positive definite matrix of unknown constants.

$\boldsymbol{\epsilon}_i$ is the error term of the latent regression. It is a $q \times 1$ random vector with expected value zero and variance-covariance matrix $\boldsymbol{\Psi}$, a $q \times q$ symmetric and positive definite matrix of unknown constants.

$\mathbf{w}_{i,1}$ and $\mathbf{w}_{i,2}$ are $p \times 1$ observable random vectors, each representing \mathbf{x}_i plus random error.

$\mathbf{v}_{i,1}$ and $\mathbf{v}_{i,2}$ are $q \times 1$ observable random vectors, each representing \mathbf{y}_i plus random error.

$\mathbf{e}_{i,1}, \dots, \mathbf{e}_{i,4}$ are the measurement errors in $\mathbf{W}_{i,1}, \mathbf{V}_{i,1}, \mathbf{W}_{i,2}$ and $\mathbf{V}_{i,2}$ respectively. Joining the vectors of measurement errors into a single long vector \mathbf{e}_i , its covariance matrix may be written as a partitioned matrix

$$\text{cov}(\mathbf{e}_i) = \text{cov} \begin{pmatrix} \mathbf{e}_{i,1} \\ \mathbf{e}_{i,2} \\ \mathbf{e}_{i,3} \\ \mathbf{e}_{i,4} \end{pmatrix} = \begin{pmatrix} \Omega_{11} & \Omega_{12} & \mathbf{0} & \mathbf{0} \\ \Omega_{12}^\top & \Omega_{22} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Omega_{33} & \Omega_{34} \\ \mathbf{0} & \mathbf{0} & \Omega_{34}^\top & \Omega_{44} \end{pmatrix} = \boldsymbol{\Omega}.$$

In addition, the matrices of covariances between $\mathbf{X}_i, \boldsymbol{\epsilon}_i$ and \mathbf{e}_i are all zero.

Collecting $\mathbf{W}_{i,1}, \mathbf{W}_{i,2}, \mathbf{V}_{i,1}$ and $\mathbf{V}_{i,2}$ into a single long data vector \mathbf{D}_i , we write its variance-covariance matrix as a partitioned matrix:

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} & \boldsymbol{\Sigma}_{13} & \boldsymbol{\Sigma}_{14} \\ & \boldsymbol{\Sigma}_{22} & \boldsymbol{\Sigma}_{23} & \boldsymbol{\Sigma}_{24} \\ & & \boldsymbol{\Sigma}_{33} & \boldsymbol{\Sigma}_{34} \\ & & & \boldsymbol{\Sigma}_{44} \end{pmatrix},$$

where the covariance matrix of $\mathbf{W}_{i,1}$ is Σ_{11} , the covariance matrix of $\mathbf{V}_{i,1}$ is Σ_{22} , the matrix of covariances between $\mathbf{W}_{i,1}$ and $\mathbf{V}_{i,1}$ is Σ_{12} , and so on.

- (a) Write the elements of the partitioned matrix Σ in terms of the parameter matrices of the model. Be able to show your work for each one.
 - (b) Prove that all the model parameters are identifiable by solving the covariance structure equations.
 - (c) Give a Method of Moments estimator of Φ_x . There is more than one reasonable answer. Remember, your estimator cannot be a function of any unknown parameters, or you get a zero. For a particular sample, will your estimate be in the parameter space? Mine is.
 - (d) Give a Method of Moments estimator for β . Remember, your estimator cannot be a function of any unknown parameters, or you get a zero. How do you know your estimator is consistent? Use $\widehat{\Sigma} \xrightarrow{p} \Sigma$.
13. For the double measurement regression model of Question 12,
- (a) How many unknown parameters appear in the covariance matrix of the observable variables?
 - (b) How many unique variances and covariances are there in the covariance matrix of the observable variables? This is also the number of covariance structure equations.
 - (c) How many equality constraints does the model impose on the covariance matrix of the observable variables? What are they?
 - (d) Does the number of covariance structure equations minus the number of parameters equal the number of constraints?