

Including Measurement Error in the Regression
Model: A First Try¹
STA431 Winter/Spring 2017

¹See last slide for copyright information.

Overview

- 1 Moment Structure Equations
- 2 A first try
- 3 Identifiability
- 4 Parameter Count Rule

Moments and Moment Structure Equations

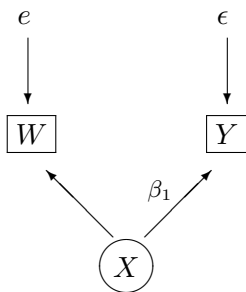
Model $D \sim P_\theta$

- *Moments* of a distribution are quantities such $E(X)$, $E(Y^2)$, $Var(X)$, $E(X^2Y^2)$, $Cov(X, Y)$, and so on.
- *Moment structure equations* are a set of equations expressing moments of the distribution of the observable data in terms of the model parameters. $m = g(\theta)$
- If the moments involved are limited to variances and covariances, the moment structure equations are called *covariance structure equations*.

Important process

- Calculate the moments of the distribution (usually means, variances and covariances) in terms of the model parameters, obtaining a system of moment structure equations. $m = g(\theta)$
- Solve the moment structure equations for the parameters, expressing the parameters in terms of the moments. $\theta = g^{-1}(m)$
- Method of Moments: $\hat{\theta} = g^{-1}(\hat{m})$
- By LLN and Continuous mapping, $\hat{\theta} \xrightarrow{p} \theta$
- So even if we're not going to use the Method of Moments, solving $\theta = g^{-1}(m)$ shows that consistent estimation is possible.

A first try at including measurement error in the explanatory variable



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

Observable data are the pairs (W_i, Y_i) for $i = 1, \dots, n$.
Try to fit the true model.

Details

Make everything normal for simplicity

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

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

where

- X_i is normally distributed with mean μ_x 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 $\omega > 0$
- X_i, e_i, ϵ_i are all independent.

Observable data are the pairs (W_i, Y_i) for $i = 1, \dots, n$.

Model implies that the (W_i, Y_i) are independent bivariate normal

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

$$W_i = \nu + X_i + e_i$$

with

$$E \begin{pmatrix} W_i \\ Y_i \end{pmatrix} = \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} \nu + \mu_x \\ \beta_0 + \beta_1 \mu_x \end{pmatrix},$$

and variance covariance matrix

$$\text{cov} \begin{pmatrix} W_i \\ Y_i \end{pmatrix} = \boldsymbol{\Sigma} = [\sigma_{i,j}] = \begin{pmatrix} \phi + \omega & \beta_1 \phi \\ \beta_1 \phi & \beta_1^2 \phi + \psi \end{pmatrix}.$$

Could we know the parameters if we knew $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$?

Big problem revealed by the moment structure equations

$m = g(\theta)$. Solve to obtain $\theta = g^{-1}(m)$

$$\theta = (\beta_0, \beta_1, \mu_x, \phi, \psi, \nu, \omega)$$

$$\mu_1 = \mu_x + \nu$$

$$\mu_2 = \beta_0 + \beta_1 \mu_x$$

$$\sigma_{1,1} = \phi + \omega$$

$$\sigma_{1,2} = \beta_1 \phi$$

$$\sigma_{2,2} = \beta_1^2 \phi + \psi$$

It is impossible to solve these five equations uniquely for the seven model parameters.

A numerical example

$$\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} \mu_x + \nu \\ \beta_0 + \beta_1 \mu_x \end{pmatrix}$$

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{22} \end{pmatrix} = \begin{pmatrix} \phi + \omega & \beta_1 \phi \\ & \beta_1^2 \phi + \psi \end{pmatrix}$$

	μ_x	β_0	ν	β_1	ϕ	ω	ψ
θ_1	0	0	0	1	2	2	3
θ_2	0	0	0	2	1	3	1

Both θ_1 and θ_2 imply a bivariate normal distribution with mean zero and covariance matrix

$$\Sigma = \begin{bmatrix} 4 & 2 \\ 2 & 5 \end{bmatrix},$$

and thus the same distribution of the sample data.

Parameter Identifiability

- No matter how large the sample size, it will be impossible to decide between θ_1 and θ_2 , because they imply exactly the same probability distribution of the observable data.
- The problem here is that the parameters of the regression are not identifiable.

Definitions

Think of $D_i \sim \text{Poisson}(\lambda)$ and $D_i \sim \text{Poisson}(\lambda_1 + \lambda_2)$.

If the probability distribution of the observable data is a one-to-one function of the parameter (vector), the parameter (vector) is said to be identifiable.

- The probability distribution of the data is always a function of the parameter.
- If the parameter is also a function of the probability distribution, the function is one-to-one and the parameter is identifiable.
- That is, if the parameter can somehow be recovered from the distribution of the data, it is identifiable.

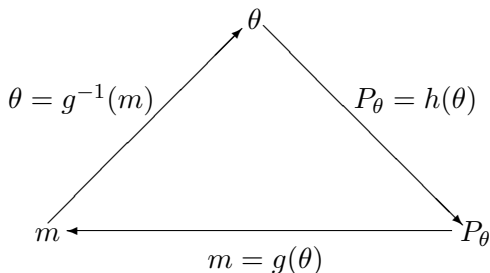
Regression model with no measurement error

Example of proving identifiability

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{X}_i + \epsilon_i$$

- The mean and covariance matrix of $\mathbf{D}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix}$ are functions of the probability distribution (calculate expected values).
- To get Method of Moments estimates, we solved for the parameters from the mean and covariance matrix of \mathbf{D}_i .
- Therefore the parameters are a function of the probability distribution.
- So they are identifiable.
- This is the way it goes in general.

Identification of parameters from the moments



- $m = g(\theta)$ are the moment structure equations.
- $\theta = g^{-1}(m)$ is the solution of the moment structure equations.
- In this course, parameters will be identified from $\mathbf{m} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ (usually just $\boldsymbol{\Sigma}$), or not at all.

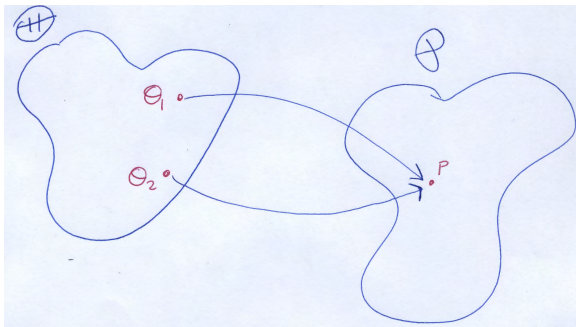
Identification from the moments $\boldsymbol{\mu}(\boldsymbol{\theta})$ and $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ or not at all

- If the distributions are normal, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are all there is.
- If the distributions are unknown, we still have $(\overline{\mathbf{D}}_n, \widehat{\boldsymbol{\Sigma}}_n) \xrightarrow{p} (\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
- If the parameters can be recovered from $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, they can be estimated based on $\overline{\mathbf{D}}_n$ and $\widehat{\boldsymbol{\Sigma}}_n$.
- If the parameters cannot be recovered from $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, we are out of luck.
- So in practice, identifiability means identifiability from the moments.
- Usually just $\boldsymbol{\Sigma}$.

Non-identifiability

Parameter is identifiable if the probability distribution of the observable data is a one-to-one function of the parameter.

If two different parameter values yield the same distribution of the data, the parameter is not identifiable.



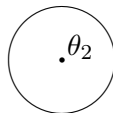
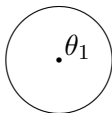
Identifiability is a big concept

- It means *knowability* of the parameters from the distribution of the data.
- We will do simple proofs that show whether certain information can be known.
- Call it the **algebra of the knowable**.

Theorem

If the parameter vector is not identifiable, consistent estimation is impossible.

- Let $\theta_1 \neq \theta_2$ but $P_{\theta_1}(d_n) = P_{\theta_2}(d_n)$ for all n .
- So the distribution of $T_n = T_n(D_1, \dots, D_n)$ is identical for θ_1 and θ_2 .
- Suppose T_n is a consistent estimator of θ .
- Then $T_n \xrightarrow{p} \theta_1$ and $T_n \xrightarrow{p} \theta_2$.



- Impossible.

Identifiability of *functions* of the parameter vector

- If a function $g(\boldsymbol{\theta})$ can be recovered from the distribution of the observable data, that function of the parameter vector is said to be identifiable.
- This applies to individual parameters and subsets of the parameters.
- Frequently, not everything can be known, but informative *functions* of the parameter are knowable.

Some sample questions will be based on this model:

Independently for $i = 1, \dots, n$, let $W_i = X_i + e_i$, where

- $X_i \sim N(\mu_x, \phi)$
- $e_i \sim N(0, \omega)$
- X_i and e_i are independent.
- Only W_i is observable (X_i is a latent variable).

How does this fit the definition of a *model*?

Sample questions

Let $W_i = X_i + e_i$, where

- $X_i \sim N(\mu_x, \phi)$
- $e_i \sim N(0, \omega)$
- X_i and e_i are independent.
- Only W_i is observable (X_i is a latent variable).

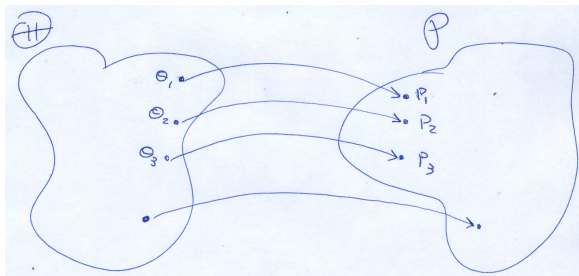
In the following questions, you may use the fact that the normal distribution corresponds uniquely to the pair (μ, σ^2) .

- 1 What is the parameter vector θ ?
- 2 What is the parameter space Θ ?
- 3 What is the probability distribution of the observable data?
- 4 Give the moment structure equations.
- 5 Either prove that the parameter is identifiable, or show by an example that it is not. A simple numerical example is best.
- 6 Give two *functions* of the parameter vector that are identifiable.

A Useful Equivalent Definition of Identifiability

Equivalent to P_θ is a one-to-one function of θ

- Suppose a statistical model implies $\mathbf{D} \sim P_\theta, \theta \in \Theta$. If no two points in Θ yield the same probability distribution, then the parameter θ is said to be identifiable.
- That is, identifiability means that $\theta_1 \neq \theta_2$ implies $P_{\theta_1} \neq P_{\theta_2}$.



Pointwise identifiability

As opposed to global identifiability

- Frequently, parameters will be identifiable in some parts of the parameter space but not others.
- The parameter is said to be identifiable at a point θ_0 if no other point in Θ yields the same probability distribution as θ_0 .
- That is, $\theta \neq \theta_0$ implies $P_\theta \neq P_{\theta_0}$ for all $\theta \in \Theta$.

If the parameter is identifiable at every point in Θ , it is identifiable according to the earlier definitions.

Determining identifiability in practice

- In practice, identifiability means that the moment structure equations can be solved uniquely for the parameters. Uniquely means there is only one solution.
- This is a strictly mathematical issue, though it has huge implications for statistical estimation and inference.

Proving identifiability

- You can explicitly solve the moment structure equations.
- You can use theorems.
- We will develop a collection of identifiability rules.
- These are really simple theorems about the existence of unique real solutions to equations.
- They are not well-known to mathematicians because they are too specific to be interesting.
- We will be able to look at a path diagram and verify that the parameters are identifiable. Usually.

Proving that a parameter is *not* identifiable

- You can carefully describe the set of points in the parameter space that yield the same mean and covariance matrix. It's a lot of work, even for small models.
- You can produce a numerical example of two different points that yield the same mean and covariance matrix. That settles it, but it can still be a lot of work for big models.
- You can use a big theorem.

The Parameter Count Rule

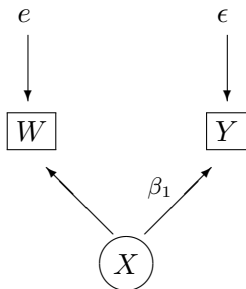
For establishing non-identifiability

Suppose identifiability is to be decided based on a set of moment structure equations. If there are more parameters than equations, the set of points where the parameter vector is identifiable occupies a set of volume zero in the parameter space.

- Note that the empty set has volume zero.
- The parameter count rule is really a theorem about the existence of unique real solutions to systems of equations.
- The moment structure equations need to have derivatives and mixed partial derivatives of all orders, but they usually do.

Back to the example

Trying to include measurement error in the model



- Recall the first attempt to include measurement error in the model.
- There were five moment structure equations in seven unknown parameters.
- The model failed the parameter count rule.
- Game over.

Again: The Parameter Count Rule

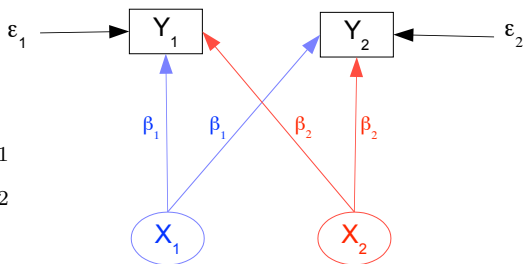
Suppose identifiability is to be decided based on a set of moment structure equations. If there are more parameters than equations, the set of points where the parameter vector is identifiable occupies a set of volume zero in the parameter space.

- So a necessary condition for parameter identifiability is that there be at least as many moment structure equations as parameters.
- There may be points in the parameter space where the parameter is identifiable, but if so, that set of points has volume zero.
- It's not a sufficient condition. There can be more equations than unknown parameters, and still no unique solution.
- Failure of the parameter count rule means that it's impossible to identify the whole parameter vector.
- Useful functions of the parameters may be identifiable, maybe including what you really want to know.
- Maximum likelihood estimation depends on identifiability of the entire parameter vector (usually).

Example

To illustrate the parameter count rule.

There are two latent explanatory variables and two observable response variables.



$$Y_1 = \beta_1 X_1 + \beta_2 X_2 + \epsilon_1$$

$$Y_2 = \beta_1 X_1 + \beta_2 X_2 + \epsilon_2$$

where

- X_1 , X_2 , ϵ_1 and ϵ_2 are independent normal random variables with expected value zero, and
- $Var(X_1) = Var(X_2) = 1$, $Var(\epsilon_1) = \psi_1$ and $Var(\epsilon_2) = \psi_2$.
- Only Y_1 and Y_2 are observable.

The parameter vector is $\theta = (\beta_1, \beta_2, \psi_1, \psi_2)$.

Calculate the covariance matrix of $(Y_1, Y_2)^\top$

Expected value is (zero, zero)

$$Y_1 = \beta_1 X_1 + \beta_2 X_2 + \epsilon_1$$

$$Y_2 = \beta_1 X_1 + \beta_2 X_2 + \epsilon_2,$$

$$\begin{aligned}\Sigma &= \begin{pmatrix} \sigma_{1,1} & \sigma_{1,2} \\ \sigma_{1,2} & \sigma_{2,2} \end{pmatrix} \\ &= \begin{pmatrix} \beta_1^2 + \beta_2^2 + \psi_1 & \beta_1^2 + \beta_2^2 \\ \beta_1^2 + \beta_2^2 & \beta_1^2 + \beta_2^2 + \psi_2 \end{pmatrix}\end{aligned}$$

Covariance structure equations

$$\theta = (\beta_1, \beta_2, \psi_1, \psi_2)$$

$$\sigma_{1,1} = \beta_1^2 + \beta_2^2 + \psi_1$$

$$\sigma_{1,2} = \beta_1^2 + \beta_2^2$$

$$\sigma_{2,2} = \beta_1^2 + \beta_2^2 + \psi_2$$

- Three equations in 4 unknowns, so the model fails.
- Parameter count rule does *not* say that a solution is impossible.
- It says that *the set of points in the parameter space where there is a unique solution (so the parameters are all identifiable) occupies a set of volume zero.*
- Are there any such points at all?

Try to solve for the parameters

$$\theta = (\beta_1, \beta_2, \psi_1, \psi_2)$$

Covariance structure equations:

$$\sigma_{1,1} = \beta_1^2 + \beta_2^2 + \psi_1$$

$$\sigma_{1,2} = \beta_1^2 + \beta_2^2$$

$$\sigma_{2,2} = \beta_1^2 + \beta_2^2 + \psi_2$$

- $\psi_1 = \sigma_{1,1} - \sigma_{1,2}$
- $\psi_2 = \sigma_{2,2} - \sigma_{1,2}$
- So those *functions* of the parameter vector are identifiable.
- What about β_1 and β_2 ?

Can we solve for β_1 and β_2 ?

$$\theta = (\beta_1, \beta_2, \psi_1, \psi_2)$$

$$\sigma_{1,1} = \beta_1^2 + \beta_2^2 + \psi_1$$

$$\sigma_{1,2} = \beta_1^2 + \beta_2^2$$

$$\sigma_{2,2} = \beta_1^2 + \beta_2^2 + \psi_2$$

- $\sigma_{1,2} = 0$ if and only if Both $\beta_1 = 0$ and $\beta_2 = 0$.
- The set of points where all four parameters can be recovered from the covariance matrix is *exactly* the set of points where the parameter vector is identifiable.
- It is

$$\{(\beta_1, \beta_2, \psi_1, \psi_2) : \beta_1 = 0, \beta_2 = 0, \psi_1 > 0, \psi_2 > 0\}$$

- A set of infinitely many points in \mathbb{R}^4
- A set of volume zero, as the theorem says.

Suppose $\beta_1^2 + \beta_2^2 \neq 0$

This is the case “almost everywhere” in the parameter space.

The set of infinitely many points $\{(\beta_1, \beta_2, \psi_1, \psi_2)\}$ such that

- $\psi_1 = \sigma_{1,1} - \sigma_{1,2}$
- $\psi_2 = \sigma_{2,2} - \sigma_{1,2}$
- $\beta_1^2 + \beta_2^2 = \sigma_{1,2}$

Substitute back into

$$\text{cov} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} \beta_1^2 + \beta_2^2 + \psi_1 & \beta_1^2 + \beta_2^2 \\ \beta_1^2 + \beta_2^2 & \beta_1^2 + \beta_2^2 + \psi_2 \end{pmatrix}$$

And see they all produce the covariance matrix

$$\Sigma = \begin{pmatrix} \sigma_{1,1} & \sigma_{1,2} \\ \sigma_{1,2} & \sigma_{2,2} \end{pmatrix}$$

And hence the same bivariate normal distribution of $(Y_1, Y_2)^\top$.

Why are there infinitely many points in this set?

$\{(\beta_1, \beta_2, \psi_1, \psi_2)\}$ such that

- $\psi_1 = \sigma_{1,1} - \sigma_{1,2}$
- $\psi_2 = \sigma_{2,2} - \sigma_{1,2}$
- $\beta_1^2 + \beta_2^2 = \sigma_{1,2} \neq 0$

Because $\beta_1^2 + \beta_2^2 = \sigma_{1,2}$ is the equation of a circle with radius $\sqrt{\sigma_{1,2}}$.

Maximum likelihood estimation

$$\boldsymbol{\theta} = (\beta_1, \beta_2, \psi_1, \psi_2)$$

$$\begin{aligned}L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= |\boldsymbol{\Sigma}|^{-n/2} (2\pi)^{-np/2} \exp -\frac{n}{2} \left\{ \text{tr}(\widehat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1}) + (\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \right\} \\L(\boldsymbol{\Sigma}) &= |\boldsymbol{\Sigma}|^{-n/2} (2\pi)^{-n} \exp -\frac{n}{2} \left\{ \text{tr}(\widehat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1}) + \bar{\mathbf{x}}^\top \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}} \right\}\end{aligned}$$

Can write likelihood as either $L(\boldsymbol{\Sigma})$ or $L(\boldsymbol{\Sigma}(\boldsymbol{\theta})) = L_2(\boldsymbol{\theta})$.

$$\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \begin{pmatrix} \beta_1^2 + \beta_2^2 + \psi_1 & \beta_1^2 + \beta_2^2 \\ \beta_1^2 + \beta_2^2 & \beta_1^2 + \beta_2^2 + \psi_2 \end{pmatrix}$$

Likelihood $L_2(\boldsymbol{\theta})$ has non-unique maximum

- $L(\boldsymbol{\Sigma})$ has a unique maximum at $\boldsymbol{\Sigma} = \widehat{\boldsymbol{\Sigma}}$.
- For every positive definite $\boldsymbol{\Sigma}$ with $\sigma_{1,2} \neq 0$, there are infinitely many $\boldsymbol{\theta} \in \Theta$ which produce that $\boldsymbol{\Sigma}$, and have the same height of the likelihood.
- This includes $\widehat{\boldsymbol{\Sigma}}$.
- So there are infinitely many points $\boldsymbol{\theta}$ in Θ with $L_2(\boldsymbol{\theta}) = L(\widehat{\boldsymbol{\Sigma}})$.
- A circle in \mathbb{R}^4 .

A circle in \mathbb{R}^4 where the likelihood is maximal

$\{(\beta_1, \beta_2, \psi_1, \psi_2)\} \subset \mathbb{R}^4$ such that

- $\psi_1 = \hat{\sigma}_{1,1} - \hat{\sigma}_{1,2}$
- $\psi_2 = \hat{\sigma}_{2,2} - \hat{\sigma}_{1,2}$
- $\beta_1^2 + \beta_2^2 = \hat{\sigma}_{1,2}$

Some Questions

Remembering that if the model is true,

- $\psi_1 = \sigma_{1,1} - \sigma_{1,2}$
- $\psi_2 = \sigma_{2,2} - \sigma_{1,2}$
- $\beta_1^2 + \beta_2^2 = \sigma_{1,2}$

What would happen in the numerical search for $\hat{\theta}$ if ...

- $\hat{\sigma}_{1,2} > \hat{\sigma}_{1,1}$?
- $\hat{\sigma}_{1,2} > \hat{\sigma}_{2,2}$?
- $\hat{\sigma}_{1,2} < 0$?

These could not *all* happen, but one of them could. When numerical maximum likelihood search leaves the parameter space, it may indicate that the model is incorrect. Or it might be just a bad starting value.

Testing hypotheses about θ

Some hypotheses are testable if the model is true, but direct likelihood ratio tests are out. All the theory depends on a unique maximum.

Remember,

$$\text{cov} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} \beta_1^2 + \beta_2^2 + \psi_1 & \beta_1^2 + \beta_2^2 \\ \beta_1^2 + \beta_2^2 & \beta_1^2 + \beta_2^2 + \psi_2 \end{pmatrix}$$

- How would you test $H_0 : \beta_1 = \beta_2 = 0$?
- If you did a large-sample likelihood ratio test, what would the degrees of freedom be?

Lessons from this example

- A parameter may be identifiable at some points but not others.
- Identifiability at infinitely many points is possible even if there are more unknowns than equations. But this can only happen on a set of volume zero.
- Some parameters and functions of the parameters may be identifiable even when the whole parameter vector is not.
- Lack of identifiability can produce multiple maxima of the likelihood function – even infinitely many.
- A model whose parameter vector is not identifiable may still be falsified by empirical data.
- Numerical maximum likelihood search may leave the parameter space. This may be a sign that the model is false. It can happen when the parameter is identifiable, too.
- Some hypotheses may be testable when the parameter is not identifiable, but these will be hypotheses about functions of the parameter that are identifiable in the part of the parameter space where the null hypothesis is true. $H_0 : \beta_1 = \beta_2 = 0$

Copyright Information

This slide show 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**. 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/431s17>