

# Lecture 3: Maximum Likelihood Estimator



Mauricio Sarrias

Universidad Católica del Norte

October 2, 2017

## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

## Main idea:

We will show that an OLS estimate of  $\rho$  will be biased and inconsistent.

# Finite and asymptotic properties



Consider the following **pure first order spatial autoregressive model**:

$$\underset{(n \times 1)}{\mathbf{y}} = \rho_0 \underset{(n \times 1)}{\mathbf{W}\mathbf{y}} + \underset{(n \times 1)}{\boldsymbol{\varepsilon}}, \quad (1)$$

where  $\rho_0$  is the true population parameter of the data generating process (DGP). The reduced form for the **pure SLM** in (1) is:

$$\mathbf{y} = (\mathbf{I}_n - \rho_0 \mathbf{W})^{-1} \boldsymbol{\varepsilon}. \quad (2)$$

As a result, the spatial lag term equals:

$$\mathbf{W}\mathbf{y} = \mathbf{W} (\mathbf{I}_n - \rho_0 \mathbf{W})^{-1} \boldsymbol{\varepsilon}. \quad (3)$$

This result will be useful later. Now, recall that if the model is  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , then the OLS estimator is  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ . Then, considering (1) the OLS estimate for  $\rho_0$  is:

$$\hat{\rho}_{OLS} = \left[ \underbrace{(\mathbf{W}\mathbf{y})^\top}_{(1 \times n)} \underbrace{(\mathbf{W}\mathbf{y})}_{(n \times 1)} \right]^{-1} \underbrace{(\mathbf{W}\mathbf{y})^\top}_{(1 \times n)} \underbrace{\mathbf{y}}_{(n \times 1)}. \quad (4)$$

# Finite and asymptotic properties



Substituting the expression for  $\mathbf{y}$  in the population equation (1) into (4) gives us the following **sampling error** equation:

$$\begin{aligned}\hat{\rho}_{OLS} &= \rho_0 + [(\mathbf{W}\mathbf{y})^\top (\mathbf{W}\mathbf{y})]^{-1} (\mathbf{W}\mathbf{y})^\top \boldsymbol{\varepsilon} \\ &= \rho_0 + \left( \sum_{i=1}^n \mathbf{y}_{Li}^2 \right)^{-1} \left( \sum_{i=1}^n \mathbf{y}_{Li} \boldsymbol{\varepsilon}_i \right),\end{aligned}$$

where  $\mathbf{y}_{Li}$  is the  $i$ th element of the spatial lag operator  $\mathbf{W}\mathbf{y} = \mathbf{y}_L$ . Assuming that  $\mathbf{W}$  is nonstochastic, the mathematical expectation of  $\hat{\rho}_{OLS}$  is

$$\begin{aligned}\mathbb{E}(\hat{\rho}_{OLS} | \mathbf{W}) &= \rho_0 + \mathbb{E} \left( [(\mathbf{W}\mathbf{y})^\top (\mathbf{W}\mathbf{y})]^{-1} (\mathbf{W}\mathbf{y})^\top \boldsymbol{\varepsilon} \middle| \mathbf{W} \right) \\ &= \rho_0 + \left( \sum_{i=1}^n \mathbf{y}_{Li}^2 \right)^{-1} \mathbb{E} \left( \sum_{i=1}^n \mathbf{y}_{Li} \boldsymbol{\varepsilon}_i \middle| \mathbf{W} \right).\end{aligned}\tag{5}$$

From (5) it is clear that if the expectation of the last term is zero, then  $\hat{\rho}_{OLS}$  will be unbiased.

# Finite and asymptotic properties



Note that

$$\begin{aligned}\mathbb{E} \left( \sum_{i=1}^n y_{Li} \epsilon_i \middle| \mathbf{W} \right) &= \mathbb{E} \left[ (\mathbf{W}\mathbf{y})^\top \boldsymbol{\varepsilon} \middle| \mathbf{W} \right] \\ &= \mathbb{E} \left[ \boldsymbol{\varepsilon}^\top (\mathbf{I} - \rho \mathbf{W}^\top)^{-1} \mathbf{W}^\top \boldsymbol{\varepsilon} \middle| \mathbf{W} \right] \quad \text{using (3)} \\ &= \mathbb{E} \left[ \boldsymbol{\varepsilon}^\top \mathbf{C}^\top \boldsymbol{\varepsilon} \middle| \mathbf{W} \right] \\ &= \mathbb{E} \left[ \text{tr} (\boldsymbol{\varepsilon}^\top \mathbf{C}^\top \boldsymbol{\varepsilon}) \middle| \mathbf{W} \right] \\ &= \mathbb{E} \left[ \text{tr} (\mathbf{C}^\top \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) \middle| \mathbf{W} \right] \quad \text{since } \text{tr}(ABC) = \text{tr}(BCA) \\ &= \text{tr} (\mathbf{C}) \mathbb{E} (\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \middle| \mathbf{W}) \quad \text{since } \text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A}^\top) \\ &\neq 0,\end{aligned} \tag{6}$$

where  $\mathbf{C} = \mathbf{W} (\mathbf{I} - \rho \mathbf{W})^{-1}$ . Therefore, given the result in (6) we have that  $\mathbb{E} (\hat{\rho}_{OLS} \middle| \mathbf{W}) = \rho_0$  if and only if  $\text{tr}(\mathbf{C}) = 0$ , which occurs if  $\rho_0 = 0$ . If  $\rho = 0$ ,  $\mathbf{C} = \mathbf{W}$ , and  $\text{tr}(\mathbf{C}) = \text{tr}(\mathbf{W}) = 0$  because the diagonal elements of  $\mathbf{W}$  are zeros. In other words, if the true model follows a spatial autoregressive structure, the OLS estimate of  $\rho$  will be biased.

## What about consistency?

Note that we can write:

$$\hat{\rho}_{OLS} = \rho_0 + \left( \frac{1}{n} \sum_{i=1}^N \mathbf{y}_{Li}^2 \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{y}_{Li} \epsilon_i \right). \quad (7)$$

Under ‘some conditions’ we can show that:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{y}_{Li}^2 \rightarrow q, \quad (8)$$

where  $q$  is some finite scalar (We need some assumptions here about  $\rho$  and the structure of the spatial weight matrix ). However, for the second term we obtain

$$\frac{1}{n} \sum_{i=1}^n \mathbf{y}_{Li} \epsilon_i \xrightarrow{p} \mathbb{E}(\mathbf{y}_{Li} \epsilon_i) = \text{tr}(\mathbf{C}) \mathbb{E}(\epsilon \epsilon^\top) \neq 0. \quad (9)$$

- Quadratic form of error terms, which introduces endogeneity.
- Then,  $\hat{\rho}_{OLS}$  is inconsistent, and we need to account for the simultaneity (ML or IV)



## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

- We will perform a simple simulation experiment to assess the properties of the OLS estimator.
- We will assume that the true DGP is:

$$\mathbf{y} = \rho_0 \mathbf{W}\mathbf{y} + \boldsymbol{\varepsilon} \quad (10)$$

where

- the true value  $\rho_0 = 0.7$ ,
- the sample size for each sample is  $n = 225$ ,
- $\boldsymbol{\varepsilon} \sim N(0, 1)$  and  $\mathbf{W}$  is an artificial  $n \times n$  weight matrix,
- The  $\mathbf{W}$  is constructed from a neighbor list for rook contiguity on a  $500 \times 500$  regular lattice.

The syntax for creating the global parameters for the simulation in R is the following:

```
# Global parameters
library("spdep")           # Load package
set.seed(123)              # Set seed
S      <- 100               # Number of simulations
n      <- 225               # spatial units
rho    <- 0.7               # True rho
w      <- cell2nb(sqrt(n), sqrt(n)) # Create artificial W matrix
iw     <- invIrM(w, rho)    # Compute inverse of (I - rho*W)
rho_hat <- vector(mode = "numeric", length = S) # Vector to save results.
```

The loop for the simulation is the following

```
# Loop for simulation
for (s in 1:S) {
  e <- rnorm(n, mean = 0 , sd = 1) # Create error term
  y <- iw %*% e # True DGP
  Wy <- lag.listw(nb2listw(w), y) # Create spatial lag
  out <- lm(y ~ Wy) # Estimate OLS
  rho_hat[s] <- coef(out)["Wy"] # Save results
}
```

Note that since  $\mathbf{W}$  is considered as fixed it is created out of simulation loop.

The results are the following

```
summary(rho_hat)
```

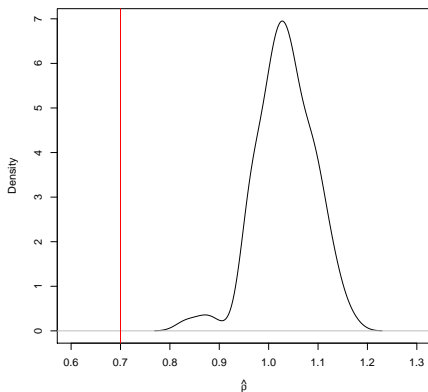
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.8309  0.9981  1.0330  1.0330  1.0750  1.1680
```

- The estimated  $\rho$  ranges from 0.8 to 1.2, that is, the range does not include the true parameter  $\rho_0 = 0.7$ .
- The mean of the estimated parameters is 1, which is very far away from 0.7. We can conclude that the OLS estimator of the pure SLM model is highly biased.

Finally, we can plot the sampling distribution of the estimated parameters in the following way:

```
plot(density(rho_hat),  
      xlab = expression(hat(rho)),  
      main = "")  
abline(v = rho, col = "red")
```

**Figure:** Distribution of  $\hat{\rho}$



*Notes:* This graph shows the sampling distribution of  $\rho$  estimated by OLS for each sample in the Monte Carlo simulation study. The true DGP follows a pure Spatial Lag Model where the true parameter is  $\rho_0 = 0.7$

1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

2 **Maximum Likelihood Estimator (MLE)**

- **Introduction to MLE**
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance



- The maximum likelihood estimate (MLE) is a way to estimate the value of a **parameter of interest**
- The MLE is the value of  $\theta$  that **maximizes** the likelihood.

# Likelihood Function

- Let  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  the conditional density...
- ... that is, the probability of observing  $y_i|\mathbf{x}_i$ .

The **likelihood function** denoted by capital  $L$  is:

$$L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}) = \prod_{i=1}^n L(\boldsymbol{\theta}; y_i|\mathbf{x}_i) = \prod_{i=1}^n f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

where  $\mathbf{y} = (y_1, \dots, y_n)$ .

- $L(\boldsymbol{\theta}; y_i|\mathbf{x}_i)$  is the likelihood contribution of the  **$i$ -th observation**,
- $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$  is the likelihood function of the **whole sample**.

The likelihood function says that, for any given sample  $\mathbf{y}|\mathbf{X}$ , the likelihood estimation is to find a set of parameters estimates, say  $\hat{\boldsymbol{\theta}}$ , such that this likelihood is maximized.

The **log-likelihood function** is:

$$\ln L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}) = \ln L(\boldsymbol{\theta}) = \ln \underbrace{\left( \prod_{i=1}^n f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \right)}_{f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})} = \sum_{i=1}^n \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

- The log-likelihood function is a monotonically increasing function of  $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$ :
  - Any maximizing value  $\hat{\boldsymbol{\theta}}$  of  $\ln L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$  must also maximize  $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$ .
- Taking logarithms converts products into sums.
  - It allows some simplification in the numerical determination of the MLE.
  - Likelihood values are often extremely small (but can also be extremely large). Numerical optimization of the likelihood highly problematic.
  - Simplification of the study of the properties of the estimator.

# Example:

## Linear Regression



Consider that  $\{y_i, \mathbf{x}_i\}$  is i.i.d, and  $y_i = \mathbf{x}_i' \boldsymbol{\beta}_0 + \epsilon_i$ , where  $\epsilon_i | \mathbf{x}_i \sim N(0, \sigma_0^2)$ . So, with  $\boldsymbol{\theta} = (\boldsymbol{\beta}', \sigma^2)$  and  $\mathbf{w}_i = (y_i, \mathbf{x}_i')'$ , the conditional pdf is

$$\begin{aligned} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) &= \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp \left[ -\frac{(y_i - \mathbf{x}_i' \boldsymbol{\beta}_0)^2}{2\sigma_0^2} \right] \\ &= \phi(y_i - \mathbf{x}_i' \boldsymbol{\beta}_0, \sigma_0^2) \end{aligned}$$

The joint p.d.f of the sample is:

$$\begin{aligned} \prod_{i=1}^n f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) &= [2\pi\sigma_0^2]^{n/2} \exp \left[ -\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)}{2\sigma_0^2} \right] \\ &= \phi(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0, \sigma_0^2 \cdot \mathbf{I}_n) \end{aligned}$$

The parameter space is  $\boldsymbol{\Theta}$  is  $\mathbb{R}^K \times \mathbb{R}_{++}$ , where  $K$  is the dimension of  $\boldsymbol{\beta}$  and  $\mathbb{R}_{++}$  is the set of positive real numbers reflecting the a priori restriction that  $\sigma_0^2 > 0$

1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

2 **Maximum Likelihood Estimator (MLE)**

- Introduction to MLE
- **Maximum Likelihood Estimator**
- Identification
- The Score Function
- The Information Matrix

3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

## Definition (ML Estimator)

The MLE is a value of the parameter vector that maximizes the sample average log-likelihood function:

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

where  $\Theta$  denotes the parameter space in which the parameter vector  $\boldsymbol{\theta}$  lies. Usually  $\Theta = \mathbb{R}^K$ .

By the nature of the objective function, the MLE is the **estimator which makes the observed data most likely to occur**. In other words, the MLE is the best “rationalization” of what we observed.

## Population analogous

$$\mathbb{E} [\ln L(\boldsymbol{\theta}; \mathbf{y} | \mathbf{X})] \equiv \int \ln L(\boldsymbol{\theta}; \mathbf{y} | \mathbf{X}) dF(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta}_0)$$

where  $F(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta}_0)$  is the joint CDF of  $(\mathbf{y}, \mathbf{X})$

- We will assume that the the sample is i.i.d.
- We also assume that we know the **true conditional density** (this is a strong assumption!).

## Assumption: Distribution

The sample  $\{y_i, \mathbf{x}_i\}$  is i.i.d with **true conditional density**  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$ .

# Expected Log-Likelihood Inequality



Is  $\mathbb{E}[\ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X})]$  maximized at  $\boldsymbol{\theta}_0$ ?

**Assumption: Dominance I**

$\mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} |\ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X})|]$  exists.

**Lemma (Expected Log-likelihood Inequality)**

*If Dominance I assumption holds, then*

$$\mathbb{E}[\ln f(y|\mathbf{x}; \boldsymbol{\theta})] \leq \mathbb{E}[\ln f(y|\mathbf{x}; \boldsymbol{\theta}_0)]$$



1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- **Identification**
- The Score Function
- The Information Matrix

3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

- Before employing MLE, it is necessary to check whether the data-generating process is sufficiently informative about the parameters of the model.
- Recall OLS:  $\hat{\beta}$  to be unique  $\mathbf{X}$  must be full-column rank. Otherwise, ...
- The question is: is the population  $\mathbb{E}[\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$  uniquely maximized at  $\boldsymbol{\theta}_0$ ?
  - If there exists another  $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  that maximized  $\mathbb{E}[\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$ , then MLE is not identified.
- This is satisfied if (**conditional density identification**):

$$f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \neq f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0) \quad \forall \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$$

## Definition (Global Identification)

The parameter vector  $\theta_0$  is globally identified in  $\Theta$  if, for every  $\theta_1 \in \Theta$ ,  $\theta \neq \theta_1$  implies that:

$$\Pr [f(y_i|\mathbf{x}_i; \theta) \neq f(y_i|\mathbf{x}_i; \theta)] > 0$$

## Assumption: Global Identification

Every parameter vector  $\theta_0 \in \Theta$  is globally identified.

## Lemma (Strict Expected Log-Likelihood Inequality)

*Under the Assumptions of **Distribution**, **Dominance I** and **Global Identification**, then*

$$\boldsymbol{\theta} \neq \boldsymbol{\theta}_0 \implies \mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta})] < \mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta}_0)]$$

## Proof.

Let  $\mathbf{w} = (y, \mathbf{x}')'$  and define

$$a(\mathbf{w}) \equiv f(y|\mathbf{x}; \boldsymbol{\theta})/f(y|\mathbf{x}; \boldsymbol{\theta}_0)$$

First, WTS that  $a(\mathbf{w}) \neq 1$  with positive probability, so that  $a(\mathbf{w})$  is nonconstant random variable (so, we can apply Jensen's Inequality).

$$\begin{aligned} a(\mathbf{w}) \neq 1 &\iff f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0) \\ \Pr[a(\mathbf{w}) \neq 1] &\iff \Pr[f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0)] \end{aligned}$$

But, by Global Identification:

$$\Pr[f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0)] > 0 \implies \Pr[a(\mathbf{w}) \neq 1] > 0$$

Now, WTS  $\mathbb{E}[\log a(\mathbf{w})] < \log \{\mathbb{E}[a(\mathbf{w})]\}$ . We use the strict version of **Jensen's inequality** which states that if  $c(x)$  is a strictly concave function and  $x$  is nonconstant random variable, then  $\mathbb{E}[c(x)] < c[\mathbb{E}(x)]$  □

## Proof.

Set  $c(x) = \log(x)$ , since  $\log(x)$  is strictly concave and  $a(\mathbf{w})$  is non-constant. Therefore

$$\mathbb{E}[\log a(\mathbf{w})] < \log \{\mathbb{E}[a(\mathbf{w})]\}$$

Now, WTS that  $\mathbb{E}(a(\mathbf{w})) = 1$ . Note that the conditional mean of  $a(\mathbf{w})$  equals 1 because:

$$\begin{aligned}\mathbb{E}[a(\mathbf{w})|\mathbf{x}] &= \int a(\mathbf{w})f(y|\mathbf{x};\boldsymbol{\theta}_0)dy \\ &= \int \frac{f(y|\mathbf{x};\boldsymbol{\theta})}{f(y|\mathbf{x};\boldsymbol{\theta}_0)}f(y|\mathbf{x};\boldsymbol{\theta}_0)dy \\ &= \int f(y|\mathbf{x};\boldsymbol{\theta})dy \\ &= 1\end{aligned}$$

By the Law of Total Expectations  $\mathbb{E}[a(\mathbf{w})] = 1$ . Combining the results:

$$\mathbb{E}[\log(a(\mathbf{w}))] < \log(1) = 0$$

But  $\log(a(\mathbf{w})) = \log f(y|\mathbf{x};\boldsymbol{\theta}) - \log f(y|\mathbf{x};\boldsymbol{\theta}_0)$ . □

## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- **The Score Function**
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

- Note that the MLE is the solution to a maximization problem.
- Therefore as any optimization problem, we need the first and second order conditions.
- The problem is that sometimes the FOC do not have a closed form solution.



## Assumption: Integrability

The pdf  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  is twice continuously differentiable in  $\boldsymbol{\theta}$  for all  $\boldsymbol{\theta} \in \Theta$ . Furthermore, the support  $\mathcal{S}(\boldsymbol{\theta})$  of  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  does not depend on  $\boldsymbol{\theta}$ , and differentiation and integration are interchangeable in the sense that

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) &= \int_{\mathcal{S}} \frac{\partial}{\partial \boldsymbol{\theta}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \\ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) &= \int_{\mathcal{S}} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta})\end{aligned}$$

and

$$\begin{aligned}\frac{\partial \mathbb{E} [\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i = x_i]}{\partial \boldsymbol{\theta}} &= \mathbb{E} \left[ \left. \frac{\partial \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right| \mathbf{x}_i = x_i \right] \\ \frac{\partial^2 \mathbb{E} [\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i = x_i]}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} &= \mathbb{E} \left[ \left. \frac{\partial^2 \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right| \mathbf{x}_i = x_i \right]\end{aligned}$$

where all terms exists. In this case, we denote the support of  $F(y)$  simply by  $\mathcal{S}$ .

## Definition (Score Function)

The score function is defined as the vector of first partial derivatives of the log-likelihood function with respect to the parameter vector  $\theta$ :

$$\mathbf{s}(\mathbf{w}, \theta) = \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \theta)}{\partial \theta} = \begin{pmatrix} \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \theta)}{\partial \theta_1} \\ \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \theta)}{\partial \theta_K} \end{pmatrix}$$

The score vector for observation  $i$  is:

$$\mathbf{s}(\mathbf{w}_i; \theta) = \frac{\partial \ln f(y_i|\mathbf{x}_i; \theta)}{\partial \theta}$$

Because of the additivity of terms in the log-likelihood function, we can write:

$$\mathbf{s}(\mathbf{w}, \theta) = \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \theta)$$

## Lemma (Score Identity)

*Under Integrability and Distribution Assumption:*

$$\mathbb{E}[\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})] = \mathbf{0}$$

- We have to be clear whether we are speaking about the score of a single observation  $\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})$  or the score of the sample  $\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})$ .
- Since under random sampling,  $\mathbf{s}(\mathbf{w}, \boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})$ , it is sufficient to establish that  $\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})] = \mathbf{0}$

## Proof.

First, we derive an integral property of pdf. Because we are **assuming**  $F(y|x; \boldsymbol{\theta})$  is a proper cdf.,

$$\int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) = \int_{\mathcal{S}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i = 1 \quad (11)$$

$\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}$ . Given **differentiability**, we can differentiate both sides of this equality with respect to  $\boldsymbol{\theta}$ :

$$\mathbf{0} = \int_{\mathcal{S}} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i \quad (12)$$

This equation states how changes in  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  resulting from changes in  $\boldsymbol{\theta}$  are restricted by (11). We can rewrite (12) as

$$\begin{aligned} \mathbf{0} &= \int_{\mathcal{S}} \frac{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i \\ \mathbf{0} &= \int_{\mathcal{S}} \frac{1}{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \underbrace{dF(y_i|\mathbf{x}_i; \boldsymbol{\theta})}_{f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i} \end{aligned} \quad (13)$$



## Proof.

Now we interpret this integral equation as an expectation. Consider:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &\equiv \frac{1}{f(y_i | \mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) &\equiv \frac{1}{f(y_i | \mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &\equiv \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta})\end{aligned}\tag{14}$$

Then, substituting into (13)

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \mathbf{0}$$

This holds for any  $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ , in particular, for  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ . Setting  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ , we obtain:

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \mathbf{0}$$

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) | \mathbf{x}] = \mathbf{0}$$

Then, by Law of Total Expectations, we obtain the desired result.

# What if the support depend on $\theta$ ?

In this case the support is  $\mathcal{S}(\theta) = A(\theta) \leq y \leq B(\theta)$ . By definition:

$$\int_{A(\theta)}^{B(\theta)} f(y|x; \theta) dy = 1$$

Now, using the Leibnitz's theorem gives:

$$\frac{\partial \int_{A(\theta)}^{B(\theta)} f(y|x; \theta) dy}{\partial \theta} = 0$$

$$\int_{A(\theta)}^{B(\theta)} \frac{\partial f(y|x; \theta)}{\partial \theta} dy + f(B(\theta)|\theta) \frac{\partial B(\theta)}{\partial \theta} - f(A(\theta)|\theta) \frac{\partial A(\theta)}{\partial \theta} = 0$$

To interchange the operations of differentiation and integration we need the second and third terms go to zero. The **necessary condition** is that

$$\lim_{y \rightarrow A(\theta)} f(y|x; \theta) = 0$$

$$\lim_{y \rightarrow B(\theta)} f(y|x; \theta) = 0$$

**Sufficient conditions** are that the support does not depend on the parameter, which means that  $\partial A(\theta)/\partial \theta = \partial B(\theta)/\partial \theta = 0$  or that the density is zero at the terminal points.

1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

2 **Maximum Likelihood Estimator (MLE)**

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- **The Information Matrix**

3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

- Since we are doing an optimization analysis, we need the Hessian Matrix.

$$\mathbf{H}(\mathbf{w}; \boldsymbol{\theta}) = \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} = \begin{pmatrix} \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1^2} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_K} \\ \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_1} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_K} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K^2} \end{pmatrix}$$

If the log-likelihood function is concave in  $\boldsymbol{\theta}$ ,  $\mathbf{H}(\mathbf{w}; \boldsymbol{\theta})$  is said to be negative definite. In the scalar case, for  $K = 1$ , this simply means that the second derivative of the log-likelihood function is negative.



Because of the additivity of terms in the log-likelihood function:

$$\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) \quad \text{where} \quad \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) = \frac{\partial^2 \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}$$

## Remark

It is important to keep in mind that both the score and Hessian depend on the sample and are therefore random variables (they differ in repeated samples).

- To analyze the variance and the limiting distribution of the ML estimator, we require some results on the **Fisher information matrix**.
- It is very related to the Hessian matrix.
- The information matrix of a sample is simply defined as the negative expectation of the Hessian Matrix:

$$\mathbf{I}(\boldsymbol{\theta}) = -\mathbb{E} [\mathbf{H}(\mathbf{w}, \boldsymbol{\theta})]$$

- Why is it useful?
  - It can be used to assess whether the likelihood function is “well behaved” (Identification)
  - Important result: the information matrix is the inverse of the variance of the maximum likelihood estimator.
  - Cramér Rao lower bound.

## Information matrix equality

The information matrix can be derived in two ways, either as minus the expected Hessian, or alternative as the variance of the score function, both evaluated at the true parameter  $\theta_0$

## Assumption: Finite Information

$\text{Var} \left[ \frac{\partial}{\partial \boldsymbol{\theta}} \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \right] \equiv \text{Var} [\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})]$  exists.

## Lemma (Information Identity)

*Under **Distribution, Differentiability and Finite Information Assumption**:*

$$\mathbb{E} \left[ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \right] = - \text{Var} [\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})]$$

Proof: (Homework)

Note the following:

$$\begin{aligned}\text{Var} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] &= \mathbb{E} \left[ \underbrace{\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)}_{(K \times 1)} \underbrace{\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)'}_{(1 \times K)} \right] + \underbrace{\mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)]}_{=0} \mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)]' \\ &= \mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)']\end{aligned}$$

Therefore we can write:

$$-\mathbf{I}(\boldsymbol{\theta}_0) = \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = -\text{Var} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = -\mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)']$$

# Example



Recall that:

$$\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = -0.5 \log(2\pi\sigma^2) - \frac{(y_i - \mathbf{x}_i' \boldsymbol{\beta})^2}{2\sigma^2}$$

We have:

$$\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) = \begin{pmatrix} \frac{1}{\sigma^2} \mathbf{x}_i \cdot \hat{\epsilon}_i \\ -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} \hat{\epsilon}_i^2 \end{pmatrix}$$

$$\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) = \begin{pmatrix} -\frac{1}{\sigma^2} \mathbf{x}_i \mathbf{x}_i' & -\frac{1}{\sigma^4} \mathbf{x}_i \cdot \hat{\epsilon}_i \\ -\frac{1}{\sigma^4} \mathbf{x}_i' \cdot \hat{\epsilon}_i & \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} \hat{\epsilon}_i^2 \end{pmatrix}$$

$$\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})' = \begin{pmatrix} \frac{1}{\sigma^4} \mathbf{x}_i \mathbf{x}_i' \hat{\epsilon}_i^2 & -\frac{1}{2\sigma^4} \mathbf{x}_i \cdot \hat{\epsilon}_i + \frac{1}{2\sigma^6} \mathbf{x}_i \cdot \hat{\epsilon}_i^3 \\ -\frac{1}{2\sigma^4} \mathbf{x}_i' \cdot \hat{\epsilon}_i + \frac{1}{2\sigma^6} \mathbf{x}_i' \cdot \hat{\epsilon}_i^3 & \frac{1}{4\sigma^4} - \frac{1}{2\sigma^6} \hat{\epsilon}_i^2 + \frac{1}{4\sigma^8} \hat{\epsilon}_i^4 \end{pmatrix}$$

where  $\mathbf{w}_i = (y_i, \mathbf{x}_i')'$ ,  $\boldsymbol{\theta} = (\boldsymbol{\theta}', \sigma^2)'$  and  $\hat{\epsilon}_i \equiv y_i - \mathbf{x}_i' \boldsymbol{\beta}$

So for  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$  the  $\widehat{\epsilon}_i$  in these expressions can be replaced by  $\epsilon_i$ . In the linear regression model,  $\mathbb{E}(\epsilon_i|\mathbf{x}_i) = 0$ . Also, since  $\epsilon_i \sim N(0, \sigma_0^2)$ , we have  $\mathbb{E}(\epsilon_i^3) = 0$  and  $\mathbb{E}(\epsilon_i^4) = 3\sigma_0^4$ . Using these relations, we have:

$$-\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})'] = \begin{pmatrix} \frac{1}{2\sigma_0^2} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') & \mathbf{0} \\ \mathbf{0}' & \frac{1}{2\sigma_0^4} \end{pmatrix}$$

If  $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$  is nonsingular, then  $\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$  is nonsingular.

## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance



- For OLS estimator consistency can be shown by finding the sampling error function and applying LLN.
- This cannot be done for nonlinear estimator such as MLE since closed form solution for finite sample estimators do not exist.

## Question

How can we proceed?

Using some LLN we know that:

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})] \quad (15)$$

That is, the sample average log-likelihood function converges to the expected log-likelihood for any value of  $\boldsymbol{\theta}$ . Recall that:

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

$$\boldsymbol{\theta}_0 \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})]$$

We would like to say that, given that

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})], \text{ then } \hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$$

We might be able to do this using the **continuous mapping theorem**.

- Let  $X_n = \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$ ,
- and  $g(\cdot) = \arg \max_{\boldsymbol{\theta} \in \Theta} (\cdot)$

Then we would like to say that if  $X_n \xrightarrow{P} X$  then  $g(X_n) \xrightarrow{P} g_0(X)$ . In words:

If the sample average of the log likelihood function is close to the true expected value of the log likelihood function, then we would expect that  $\hat{\boldsymbol{\theta}}_n$  will be close to the maximum of the expected likelihood (as  $n$  increases without bound)

However, we cannot do that!

# What is the problem?



- The problem is that the argument of the  $\arg \max(\cdot)$  is a function of  $\theta$ , not a real vector:  
 $\theta \in \Theta$ 
  - The concept of convergence in probability was defined for **sequence of random variables**
- Therefore, we need to define what we mean by the probability limit of **sequence of random functions**, as opposed to a sequence of random variables:

Convergence for sequence of random variables  $\implies X_n = X_n(\omega), \omega \in \Omega$

Convergence for sequence of random function  $\implies Q_n = Q_n(\omega, \theta), \omega \in \Omega$

## Example

In ML estimation, the log-likelihood is a function of the sample data (a random vector that depends on  $\omega$ ) and of a parameter  $\theta$ . By increasing the sample size, we obtain a sequence of log-likelihoods that depend on  $\omega$  and  $\theta$ .

How is the distance between two functions over a set containing an infinite number of possible comparisons at different values of  $\theta$  measured?

- IOW, since we are dealing with convergence on a **function space** we need to define when two functions are close to one another.
- To reduce the infinite dimensional character of a function to a one-dimensional concept of convergence, we take the supremum of the absolute difference of the function values over all  $\theta$  in  $\Theta$

## Definition (Uniform Convergence in Probability)

The sequence of real-valued functions  $\{Q_n(\boldsymbol{\theta})\}$  converges uniformly in probability to the limit function  $Q_0(\boldsymbol{\theta})$  if  $\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| \xrightarrow{P} 0$ . We will say that  $Q_n(\boldsymbol{\theta}) \xrightarrow{P} Q_0(\boldsymbol{\theta})$  **uniformly**.

Another way to express uniform convergence in probability is:

$$\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| = o_p(1)$$

IOW, instead of requiring that the distance  $|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})|$  converge in probability to 0 for each  $\boldsymbol{\theta}$ , we require convergence of  $\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})|$ , which is the maximum distance that can be found by ranging over the space parameters.

Extending the concept to random vectors is straightforward. Now suppose that  $\{Q_n(\boldsymbol{\theta})\}$  is a sequence of  $K \times 1$  random vectors that depend both on the data and on the parameter  $\boldsymbol{\theta} \in \Theta$ . This sequence of **random vectors** is uniformly convergent in probability to  $Q_0(\boldsymbol{\theta})$  if and only if

$$\sup_{\boldsymbol{\theta} \in \Theta} \|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})\| = o_p(1)$$

where  $\|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})\|$  denotes the Euclidean norm of the vector  $Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})$ . By taking the supremum over  $\boldsymbol{\theta}$  we obtain another random quantity that does not depend on  $\boldsymbol{\theta}$ .



## Definition (Pointwise Convergence in probability)

The sequence of real-valued functions  $\{Q_n(\boldsymbol{\theta})\}$  converges pointwise in probability if and only if  $|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| \xrightarrow{P} 0$  **for each**  $\boldsymbol{\theta} \in \Theta$

Uniform convergence is stronger than pointwise convergence.

Now we present the **uniform LLN** to study sequences of random functions which is analogous to the Chebychev's LLN for averages of random variables.

## Theorem (Uniform LLN)

Suppose that  $Q(\boldsymbol{\theta}, U)$  is continuous function over  $\boldsymbol{\theta} \in \Theta$ , a closed and bounded subset of  $\mathbb{R}^p$ , and that  $\{U_n\}$  is a sequence of i.i.d. random variables with cdf  $F_U(u)$ . If  $\mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} \|Q(\boldsymbol{\theta}; U)\|]$  exists, then

- 1  $\mathbb{E}[Q(\boldsymbol{\theta}; U)]$  is continuous over  $\boldsymbol{\theta} \in \Theta$  and,
- 2  $\frac{1}{n} \sum_{i=1}^n Q(\boldsymbol{\theta}; u_i) \xrightarrow{P} \mathbb{E}[Q(\boldsymbol{\theta}; U)]$  uniformly.

The following Theorem makes the connection between the uniform convergence of  $\frac{1}{n} \sum_{i=1}^n Q(\boldsymbol{\theta}; u_i)$  to  $\mathbb{E}[Q(\boldsymbol{\theta}; U)]$  and the convergence of  $\hat{\boldsymbol{\theta}}_n$  to  $\boldsymbol{\theta}_0$  using the assumption of **compact parameter space**.

## Theorem (Consistency of Maxima with Compact Parameter Space)

Suppose that:

- 1 (compact parameter space)  $\Theta \subset \mathbb{R}^p$  is a closed and bounded parameter space,
- 2 (uniform convergence)  $Q_n(\theta)$  is a sequence of function that convergence in probability uniformly to a function  $Q_0(\theta)$  on  $\Theta$ ,
- 3 (continuity)  $Q_n(\theta)$  is continuous in  $\theta$  for any data  $(\mathbf{w}_1, \dots, \mathbf{w}_n)$ ,
- 4 (identification)  $Q_0(\theta)$  is uniquely maximized at  $\theta_0 \in \Theta$

then  $\hat{\theta}_n \equiv \arg \max_{\theta \in \Theta} Q_n(\theta)$  converges in probability to  $\theta_0$ .

## Theorem (Consistency of Maxima without Compactness)

Suppose that:

- 1 (interior)  $\theta_0$  is an element of the interior of a convex parameter space  $\Theta$ ,
- 2 (pointwise convergence)  $Q_n(\theta)$  converges in probability to  $Q_0(\theta)$  for all  $\theta \in \Theta$ ,
- 3 (concavity)  $Q_n(\theta)$  is concave over the parameter space for any data  $(\mathbf{w}_1, \dots, \mathbf{w}_n)$ ,
- 4 (identification)  $Q_0(\theta)$  is uniquely maximized at  $\theta_0 \in \Theta$

then, as  $n \rightarrow \infty$ ,  $\hat{\theta}_n$  exists with probability approaching 1 and  $\hat{\theta}_n \xrightarrow{P} \theta_0$ .

## Theorem (Consistency of conditional ML with compact parameter)

Let  $\{y_i, \mathbf{x}_i\}$  be i.i.d with conditional density  $f(y_i|\mathbf{x}_i; \theta_0)$  and let  $\hat{\theta}$  be the conditional ML estimator, which maximizes the average log conditional likelihood:

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \theta)$$

Suppose the model is correctly specified so that  $\theta_0$  is in  $\Theta$ . Suppose that

- 1 (Compactness) the parameter space  $\Theta$  is compact subset of  $\mathbb{R}^K$ ,
- 2  $f(y_i|\mathbf{x}_i; \theta)$  is continuous in  $\theta$  for all  $(y_i, \mathbf{x}_i)$ ,
- 3  $f(y_i|\mathbf{x}_i; \theta)$  is measurable in  $(y_i, \mathbf{x}_i)$  for all  $\theta \in \Theta$  (so  $\hat{\theta}$  is well-defined random variable),
- 4 (identification)  $\Pr [f(y_i|\mathbf{x}_i; \theta) \neq f(y_i|\mathbf{x}_i; \theta_0)] > 0$  for all  $\theta \neq \theta_0$  in  $\Theta$ ,
- 5 (dominance)  $\mathbb{E} [\sup_{\theta \in \Theta} |\log f(y_i|\mathbf{x}_i; \theta)|] < \infty$  (note: the expectation is over  $y_i$  and  $\mathbf{x}_i$ )

Then  $\hat{\theta} \xrightarrow{p} \theta_0$

## Sketch of Proof.

We would like to apply **Consistency of Maxima with Compact Parameter Space Theorem**. In this case, let  $Q(\boldsymbol{\theta}; U) = \log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$ . Now we verify that the condition of the theorem are met:

- $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  is continuous,
- Compactness states that  $\Theta$  is a closed and bounded subset of  $E^K$ ,
- $(y_i, \mathbf{x}_i)$  are i.i.d with conditional density  $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$ ,
- Dominance I states that  $\mathbb{E} \left[ \sup_{\boldsymbol{\theta} \in \Theta} |\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})| \right]$  exists.

Therefore,  $\mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$  is continuous and

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})] \quad (16)$$

uniformly. Let  $Q_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$  and  $Q_0(\boldsymbol{\theta}) = \mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$ . Under the additional assumption of Likelihood Identification, we can invoke the **strict expected log-likelihood inequality**:  $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  that  $\mathbb{E} [\log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] < \mathbb{E} [\log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}_0)]$ . This implies that  $Q_0(\boldsymbol{\theta})$  is uniquely maximized at  $\boldsymbol{\theta}_0$ . Therefore

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \boldsymbol{\theta}_0$$



## Theorem (Consistency of conditional ML without Compactness)

Let  $\{y_i, \mathbf{x}_i\}$  be i.i.d with conditional density  $f(y_i|\mathbf{x}_i; \theta_0)$  and let  $\hat{\theta}$  be the conditional ML estimator, which maximizes the average log conditional likelihood:

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \theta)$$

Suppose the model is correctly specified so that  $\theta_0$  is in  $\Theta$ . Suppose that

- 1 the true parameter vector  $\theta_0$  is an element of the interior of convex parameter space  $\Theta \subset \mathbb{R}^K$ ,
- 2  $\log f(y_i|\mathbf{x}_i; \theta)$  is concave in  $\theta$  for all  $(y_i, \mathbf{x}_i)$ ,
- 3  $\log f(y_i|\mathbf{x}_i; \theta)$  is measurable in  $(y_i, \mathbf{x}_i)$ ,
- 4 (identification)  $\Pr [f(y_i|\mathbf{x}_i; \theta) \neq f(y_i|\mathbf{x}_i; \theta_0)] > 0$  for all  $\theta \neq \theta_0$  in  $\Theta$ ,
- 5  $\mathbb{E} [|\log f(y_i|\mathbf{x}_i; \theta)|] < \infty$  (i.e.,  $\mathbb{E} [\log f(y_i|\mathbf{x}_i; \theta)]$  exists and is finite) for all  $\theta \in \Theta$  (note: the expectation is over  $y_i$  and  $\mathbf{x}_i$ )

Then,  $n \rightarrow \infty$ ,  $\hat{\theta}$  exists with probability approaching 1 and  $\hat{\theta}_n \xrightarrow{p} \theta_0$



## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- **Asymptotic Normality**
- Estimation of Variance

## Theorem (Asymptotic Normality of Conditional ML)

Let  $\mathbf{w} \equiv (y_i, \mathbf{x}_i')'$  be i.i.d. Suppose the conditions of either Theorem 14 or 15 are satisfied, so that  $\hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$ . Suppose, in addition, that:

- 1  $\boldsymbol{\theta}_0$  is in interior of  $\Theta$ ,
- 2  $f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0)$  is twice continuously differentiable in  $\boldsymbol{\theta}$  for all  $(y_i, \mathbf{x}_i)$ ,
- 3  $\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbf{0}$  and  $-\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)']$ ,
- 4 (local dominance condition on the Hessian) for some neighborhood  $\mathcal{N}$  of  $\boldsymbol{\theta}_0$ ,

$$\mathbb{E} \left[ \sup_{\boldsymbol{\theta} \in \mathcal{N}} \|\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta})\| \right] < \infty$$

so that for any consistent estimator  $\tilde{\boldsymbol{\theta}}$ ,  $\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \tilde{\boldsymbol{\theta}}) \xrightarrow{P} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$

- 5  $\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$  is nonsingular.

Then:

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}), \mathbf{V} = -\{\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]\}^{-1} = \{\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)']\}^{-1}$$

## Proof.

The objective function is:

$$Q_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

Given that  $f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0)$  is twice continuously differentiable in  $\boldsymbol{\theta}$ , and given that  $\boldsymbol{\theta}_0$  is in interior of  $\Theta$ , then the maximum likelihood estimator satisfies

$$\frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} = \mathbf{s}(\mathbf{w}; \hat{\boldsymbol{\theta}}) = \mathbf{0}$$

We need to know about the behavior of the gradient around the true parameter. Expand this set of equations in a Taylor series around the true parameters  $\boldsymbol{\theta}_0$ . We will use the mean value theorem to truncate the Taylor series at the second term,

$$\begin{aligned} \underbrace{\frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}}}_{(K \times 1)} &= \underbrace{\frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}}_{(K \times 1)} + \underbrace{\frac{\partial \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}}_{(K \times K)} \underbrace{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)}_{(K \times 1)} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) + \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \end{aligned}$$

where  $\bar{\boldsymbol{\theta}} = \alpha \hat{\boldsymbol{\theta}} + (1 - \alpha) \boldsymbol{\theta}_0$  for some  $\alpha \in (0, 1)$

## Proof.

So,

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right]^{-1} \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right)$$

We know that

$$\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0 \implies \bar{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$$

By uniform LLN, we know that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) \xrightarrow{p} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta})]$$

Then, applying our Lemma:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}_n) \xrightarrow{p} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$$

since  $\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$  exists. Finally, using **probability limit continuity** and **nonsingular information**, then:

$$\left[ \frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}_n) \right]^{-1} \xrightarrow{p} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1}$$

## Proof.

Since  $(y_i, \mathbf{x}_i)$  are i.i.d.  $\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\right)$  is the sum of variables  $\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)$ . The score identity lemma implies that  $\mathbb{E}(\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)) = \mathbf{0}$ , and the Information Identity implies that

$$\text{Var}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)'] = -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$$

The Lindberg-Levy CLT therefore implies:

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right) \xrightarrow{d} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)])$$



## Proof.

Then:

$$\begin{aligned}\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) &= \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right]^{-1} \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right) \\ &\xrightarrow{d} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]) \\ &= \mathbf{N} \left[ \mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right] \\ &= \mathbf{N} \left[ \mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right]\end{aligned}$$



## 1 What are the consequences of applying OLS to SLM?

- Finite and Asymptotic Properties
- Illustration of bias

## 2 Maximum Likelihood Estimator (MLE)

- Introduction to MLE
- Maximum Likelihood Estimator
- Identification
- The Score Function
- The Information Matrix

## 3 Asymptotic Properties

- Consistency
- Asymptotic Normality
- Estimation of Variance

# Variance Estimation

For large but finite samples, we can therefore write the approximate distribution of  $\hat{\boldsymbol{\theta}}_n$  as

$$\hat{\boldsymbol{\theta}} \stackrel{a}{\sim} \text{N} \left[ \boldsymbol{\theta}_0, [\mathbf{I}(\boldsymbol{\theta}_0)]^{-1} \right]$$

we have three potential estimators of  $\mathbf{I}(\boldsymbol{\theta}_0)$ :

- The empirical mean of minus the Hessian,

$$\hat{\mathbf{V}}^1 = \left( \frac{1}{n} \sum_{i=1}^n -\mathbf{H}(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \right)^{-1}$$

- The empirical variance of the score:

$$\hat{\mathbf{V}}^2 = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{w}_i, \hat{\boldsymbol{\theta}})' \right)^{-1}$$

- Minus the expected Hessian evaluated at  $\hat{\boldsymbol{\theta}}$ :

$$\hat{\mathbf{V}}^3 = \left( -\mathbb{E} \left[ \mathbf{H}(\mathbf{w}, \hat{\boldsymbol{\theta}}) \right] \right)^{-1}$$



- Evaluated at a  $\theta \in \Theta$ , each estimator converges in probability uniformly to its expectation.
- Because  $\hat{\theta}_n \xrightarrow{p} \theta_0$ , evaluated at  $\hat{\theta}_n$  each estimator converges in probability to  $\mathbf{I}(\theta_0)$ .
- Because matrix inversion is a continuous transformation, the inverse of each matrix is also a consistent estimator for the variance matrix of the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_n - \theta_0)$