

# Maximum Estimation

Paolo Zacchia

Econometric Theory

Lecture 11

# Overview of Maximum Estimation

- This Lecture introduces Maximum Estimation: a family of estimation frameworks for linear and non-linear models.
- All these frameworks are based upon the optimization (that is, maximization or minimization) of an objective function.
- Hence, a more general name for the family of frameworks is that of **Extremum Estimation**.
- However, **Maximum Estimation** is more commonly used.
- **Maximum Likelihood Estimation** (MLE) is one special case of Maximum Estimation, which is emphasized here.
- Lecture 12 focuses on another special case: the Generalized Method of Moments.

# Criteria and M-Estimators

- Suppose that the probabilistic assumptions of a **structural** econometric model with variables  $(\mathbf{y}_i, \mathbf{z}_i, \varepsilon_i)$  make  $\boldsymbol{\theta}_0$  *in the population* the *only* solution to a maximization problem as:

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{Q}_0(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [q(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta})]$$

where  $q(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta})$  is a **criterion** for  $i = 1, \dots, N$ .

- Write the latter hereinafter as  $q(\mathbf{x}_i; \boldsymbol{\theta})$ , where  $\mathbf{x}_i = (\mathbf{y}_i, \mathbf{z}_i)$ .
- With i.i.d. data, it is simply  $\mathcal{Q}_0(\boldsymbol{\theta}) = \mathbb{E} [q(\mathbf{x}_i; \boldsymbol{\theta})]$ .
- An **M-Estimator**  $\hat{\boldsymbol{\theta}}_M$  for  $\boldsymbol{\theta}_0$  maximizes the sample analog of  $\mathcal{Q}_0(\boldsymbol{\theta})$ , written as  $\hat{\mathcal{Q}}_N(\boldsymbol{\theta})$ .

$$\hat{\boldsymbol{\theta}}_M = \arg \max_{\boldsymbol{\theta} \in \Theta} \hat{\mathcal{Q}}_N(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N q(\mathbf{x}_i; \boldsymbol{\theta})$$

## Example: OLS as an M-Estimator

- Let the CEF of some endogenous variable  $Y_i$ , given some  $K$  exogenous variables  $\mathbf{x}_i$ , be linear.

$$\mathbb{E}[Y_i | \mathbf{x}_i] = \mathbf{x}_i^T \boldsymbol{\beta}_0$$

- In such a case, the “true” parameter vector  $\boldsymbol{\beta}_0$  is shown to minimize the “limiting” mean squared error (MSE).

$$\boldsymbol{\beta}_0 = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^K} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N - \mathbb{E} \left[ \left( Y_i - \mathbf{x}_i^T \boldsymbol{\beta} \right)^2 \right]$$

- The OLS estimator is the minimizer of the sample analog.

$$\hat{\boldsymbol{\beta}}_{OLS} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^K} - \frac{1}{N} \sum_{i=1}^N \left( y_i - \mathbf{x}_i^T \boldsymbol{\beta} \right)^2$$

- This is an M-Estimator for  $q(Y_i, \mathbf{x}_i; \boldsymbol{\beta}) = - \left( Y_i - \mathbf{x}_i^T \boldsymbol{\beta} \right)^2$ .

## Example: NLLS as an M-Estimator

- Let the CEF be now non-linear, governed by some function denoted as  $h(\mathbf{x}_i; \boldsymbol{\theta})$ .

$$\mathbb{E}[Y_i | \mathbf{x}_i] = h(\mathbf{x}_i; \boldsymbol{\theta}_0)$$

- The result linking here  $\boldsymbol{\beta}_0$  with the MSE is here as follows.

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta} \in \Theta} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N - \mathbb{E} \left\{ [Y_i - h(\mathbf{x}_i; \boldsymbol{\theta})]^2 \right\}$$

- This lets express the **Non-Linear Least Squares** (NLLS) estimator as the solution to the sample analog of the above problem, so long as  $h(\mathbf{x}_i; \boldsymbol{\theta})$  is invertible in  $\boldsymbol{\theta}$ .

$$\hat{\boldsymbol{\theta}}_{NLLS} = \arg \max_{\boldsymbol{\theta} \in \Theta} - \frac{1}{N} \sum_{i=1}^N [y_i - h(\mathbf{x}_i; \boldsymbol{\theta})]^2$$

- In this case it is  $q(Y_i, \mathbf{x}_i; \boldsymbol{\beta}) = -[Y_i - h(\mathbf{x}_i; \boldsymbol{\theta})]^2$ .

# Score and Hessian

It is useful to define the **score** as the criterion's gradient for  $\boldsymbol{\theta}$ :

$$\mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{\partial q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_2} \\ \vdots \\ \frac{\partial q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_K} \end{bmatrix}$$

as well as the **Hessian** as... the criterion's Hessian for  $\boldsymbol{\theta}$ .

$$\begin{aligned} \mathbf{H}_i(\mathbf{x}_i; \boldsymbol{\theta}) &= \frac{\partial \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} = \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \\ &= \begin{bmatrix} \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_1} & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_K} \\ \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_2} & \cdots & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_1} & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_2} & \cdots & \frac{\partial^2 q(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_K} \end{bmatrix} \end{aligned}$$

Both objects are observation-specific (for  $i = 1, \dots, N$ ).

# Identification of M-Estimators

## Theorem 1

**Identification of M-Estimators.** *In a M-Estimation environment, the “true” parameter set  $\theta_0$  is locally point identified if the following limiting average Hessian matrix evaluated at  $\theta_0$  has full  $K$  rank.*

$$\mathbf{Q}_0 \equiv \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{H}_i(\mathbf{x}_i; \theta_0)]$$

## Proof.

This is indeed a simple application of the Implicit Function Theorem. In a well-defined M-Estimator, the true parameter vector  $\theta_0$  sets the  $K$  First Order Conditions of the empirical criterion function  $\hat{Q}_N(\theta_0)$  equal to zero *at the probability limit*, at some *limiting average score*.

$$\frac{\partial}{\partial \theta} \hat{Q}_N(\theta_0) = \frac{1}{N} \sum_{i=1}^N \left[ \frac{\partial}{\partial \theta} q(\mathbf{x}_i; \theta_0) \right] \xrightarrow{P} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{s}_i(\mathbf{x}_i; \theta_0)] = \mathbf{0}$$

If the Jacobian  $\mathbf{Q}_0$  has full rank, the local solution  $\theta_0$  is unique. □

## Examples: Identification of OLS and NLLS (1/2)

- In OLS, the score and Hessian are:

$$\mathbf{s}_i(y_i, \mathbf{x}_i; \boldsymbol{\beta}) = 2\mathbf{x}_i\varepsilon_i$$

$$\mathbf{H}_i(y_i, \mathbf{x}_i; \boldsymbol{\beta}) = -2\mathbf{x}_i\mathbf{x}_i^T$$

hence, identification requires full rank of a familiar matrix.

$$\mathbf{Q}_0 = -2\mathbf{K}_0 = \lim_{N \rightarrow \infty} -\frac{2}{N} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{x}_i\mathbf{x}_i^T \right]$$

- In NLLS, define the *error term* by  $\varepsilon_i \equiv y_i - h(\mathbf{x}_i; \boldsymbol{\theta})$ ; thus the score and the Hessian are as follows.

$$\mathbf{s}_i(y_i, \mathbf{x}_i; \boldsymbol{\theta}) = 2 \frac{\partial}{\partial \boldsymbol{\theta}} h(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \varepsilon_i$$

$$\begin{aligned} \mathbf{H}_i(y_i, \mathbf{x}_i; \boldsymbol{\theta}) &= -2 \frac{\partial}{\partial \boldsymbol{\theta}} h(\mathbf{x}_i; \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\theta}^T} h(\mathbf{x}_i; \boldsymbol{\theta}) \\ &\quad + 2 \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} h(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \varepsilon_i \end{aligned}$$



## Examples: Identification of OLS and NLLS (2/2)

- Note that, by the Law of Iterated Expectations:

$$\begin{aligned}\mathbb{E} \left[ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} h(\mathbf{x}_i; \boldsymbol{\theta}_0) \cdot \varepsilon_i \right] &= \mathbb{E}_{\mathbf{x}} \left[ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} h(\mathbf{x}_i; \boldsymbol{\theta}_0) \cdot \mathbb{E}[\varepsilon_i | \mathbf{x}_i] \right] \\ &= \mathbf{0}\end{aligned}$$

as  $\mathbb{E}[\varepsilon_i | \mathbf{x}_i] = \mathbb{E}[Y_i | \mathbf{x}_i] - h(\mathbf{x}_i; \boldsymbol{\theta}_0) = 0$ .

- Thus, any expected Hessian matrix of NLLS is just:

$$\mathbb{E}[\mathbf{H}_i(\mathbf{x}_i; \boldsymbol{\theta}_0)] = -2 \mathbb{E}[\mathbf{h}_{0i} \mathbf{h}_{0i}^T]$$

where:

$$\mathbf{h}_{0i} \equiv \frac{\partial}{\partial \boldsymbol{\theta}} h(\mathbf{x}_i; \boldsymbol{\theta}_0)$$

is the derivative of the CEF evaluated at  $\mathbf{x}_i$  and at the true parameters  $\boldsymbol{\theta}_0$ .

- The identification of NLLS is evaluated through  $\mathbb{E}[\mathbf{h}_{0i} \mathbf{h}_{0i}^T]$ .

## Remarks on identification of M-Estimators

- In practical applications, it is customary to verify that the *sample mean* of the Hessian has full rank, as an indication that the model is identified ( $N^{-1}\mathbf{X}^T\mathbf{X}$  for OLS).
- Furthermore, it is useful to check that the rows-columns of the Hessian's sample mean are not *too correlated*; otherwise, identification is said to be *weak*, and the estimates are thus very imprecise with large standard errors.
- This problem is called **quasi-multicollinearity**, and it is intuitively due to the statistical difficulty of distinguishing between two “factors” (like different explanatory variables, columns in  $\mathbf{X}$ ) if they are very similar.
- In the IV and 2SLS cases, this corresponds to the problem of *weak instruments* which occurs if the inverse of  $\mathbf{X}^T\mathbf{P}_Z\mathbf{X}$  is too large.

# Maximum (Likelihood) Estimation (1/3)

- Consider one special case of M-Estimation: MLE.
- Start from  $P$  **structural** relationships  $\mathbf{y}_i = \mathbf{s}(\mathbf{y}_i, \mathbf{z}_i, \boldsymbol{\varepsilon}_i; \boldsymbol{\theta})$ , as well as a joint distribution  $F_{\mathbf{z}, \boldsymbol{\varepsilon}}(\mathbf{z}_i, \boldsymbol{\varepsilon}_i)$  of the exogenous variables (both observable and unobservable).
- One can write the distribution of  $\mathbf{x}_i = (\mathbf{y}_i, \mathbf{z}_i)$  as:

$$f_{\mathbf{z}, \boldsymbol{\varepsilon}}(\mathbf{z}_i, \boldsymbol{\varepsilon}_i; \boldsymbol{\theta}) = f_{\mathbf{z}, \boldsymbol{\varepsilon}}\left(\mathbf{z}_i; \mathbf{s}_{\boldsymbol{\varepsilon}}^{-1}(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta}); \boldsymbol{\theta}\right)$$

where:  $f_{\mathbf{z}, \boldsymbol{\varepsilon}}(\cdot)$  is the probability mass/density function that is associated with  $F_{\mathbf{z}, \boldsymbol{\varepsilon}}(\cdot)$ , whereas  $\mathbf{s}_{\boldsymbol{\varepsilon}}^{-1}(\cdot)$  is the solution of  $\mathbf{s}(\cdot)$  with respect to the unobservables  $\boldsymbol{\varepsilon}_i$  (if it exists).

- Hence, MLE is the M-Estimator with criterion equaling the logarithm of  $f_{\mathbf{z}, \boldsymbol{\varepsilon}}(\cdot)$ ; write it as  $\ell(\mathbf{x}_i; \boldsymbol{\theta})$  where  $\mathbf{x}_i = (\mathbf{y}_i, \mathbf{z}_i)$ .

$$q(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta}) = \log f_{\mathbf{z}, \boldsymbol{\varepsilon}}\left(\mathbf{z}_i; \mathbf{s}_{\boldsymbol{\varepsilon}}^{-1}(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta}); \boldsymbol{\theta}\right) \equiv \ell(\mathbf{x}_i; \boldsymbol{\theta})$$

## Maximum (Likelihood) Estimation (2/3)

- This complies as follows with the definition of M-Estimator.

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta} \in \Theta} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\ell(\mathbf{x}_i; \boldsymbol{\theta})]$$

- This holds since for  $i = 1, \dots, N$  one can maintain that:

$$\begin{aligned} \mathbb{E} [\ell(\mathbf{x}_i; \boldsymbol{\theta})] - \mathbb{E} [\ell(\mathbf{x}_i; \boldsymbol{\theta}_0)] &= \mathbb{E} \left[ \log \frac{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta})}{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}_0)} \right] \\ &\leq \log \mathbb{E} \left[ \frac{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta})}{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}_0)} \right] \\ &= \log \int_{\mathbb{X}} \frac{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta})}{f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}_0)} f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}_0) d\mathbf{x}_i \\ &= 0 \end{aligned}$$

per an application of Jensen's inequality in the second line;  
 $\mathbb{X}$  is here the joint support of  $\mathbf{x}_i = (\mathbf{y}_i, \mathbf{z}_i)$ .

## Maximum (Likelihood) Estimation (3/3)

- It follows that  $\mathbb{E}[\ell(\mathbf{x}_i; \boldsymbol{\theta}_0)] \geq \mathbb{E}[\ell(\mathbf{x}_i; \boldsymbol{\theta})]$  for  $i = 1, \dots, N$ : this shows the argued compliance with M-Estimators.
- A variant of such an approach is to specify the **conditional** distribution of the unobserved factors  $\boldsymbol{\varepsilon}_i$  **only**, while taking the realizations  $\mathbf{y}_i$  of the exogenous variables **as given**.
- The criterion function would then be as follows.

$$q(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta}) = \log f_{\mathbf{z}, \boldsymbol{\varepsilon}} \left( \mathbf{s}_{\boldsymbol{\varepsilon}}^{-1}(\mathbf{y}_i, \mathbf{z}_i; \boldsymbol{\theta}) \mid \mathbf{z}_i \right)$$

- This is referred as: **Conditional Maximum Likelihood Estimation** (CMLE).
- Due to its flexibility, CMLE is prevalent in econometrics – as the next example helps clarify.

## Example: MLE and linear regression (1/4)

- Consider a linear regression model like  $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$ .
- Suppose that the error term is homoscedastic, independent across observations and **normally distributed**:

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

where  $\sigma^2$  is an unknown parameter of the model.

- Also assume that the right-hand side variables  $\mathbf{x}_i$  are **fixed** (*non-stochastic*): one specific realization  $\mathbf{x}_i$  occurs *always*.
- The probability density function of  $\varepsilon_i$  is then as follows.

$$\begin{aligned} f_{\varepsilon}(\varepsilon_i | \boldsymbol{\beta}, \sigma^2) &= f_{\varepsilon}(y_i - \mathbf{x}_i^T \boldsymbol{\beta} | \boldsymbol{\beta}, \sigma^2) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}\right) \end{aligned}$$

## Example: MLE and linear regression (2/4)

- The likelihood function for the parameters  $\theta = (\beta; \sigma^2)$  is:

$$\begin{aligned}\mathcal{L}(\theta | \{y_i, \mathbf{x}_i\}_{i=1}^N) &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \beta)^2}{2\sigma^2}\right) \\ &= \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{N}{2}} \exp\left(-\frac{\sum_{i=1}^N (y_i - \mathbf{x}_i^T \beta)^2}{2\sigma^2}\right)\end{aligned}$$

- ...and the log-likelihood function is as follows.

$$\begin{aligned}\log \mathcal{L}(\theta | \{y_i, \mathbf{x}_i\}_{i=1}^N) &= -\frac{N}{2} (\log 2\pi + \log \sigma^2) \\ &\quad - \frac{\sum_{i=1}^N (y_i - \mathbf{x}_i^T \beta)^2}{2\sigma^2}\end{aligned}$$

- Both functions have a unique maximum.

## Example: MLE and linear regression (3/4)

- The First Order Conditions are:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \log \mathcal{L} \left( \boldsymbol{\theta} \mid \{y_i, \mathbf{x}_i\}_{i=1}^N \right) &= \\ &= \begin{bmatrix} \sum_{i=1}^N \mathbf{x}_i \left( y_i - \mathbf{x}_i^T \boldsymbol{\beta} \right) \sigma^{-2} \\ -\frac{N}{2} \sigma^{-2} + \frac{1}{2} \sum_{i=1}^N \left( y_i - \mathbf{x}_i^T \boldsymbol{\beta} \right)^2 \sigma^{-4} \end{bmatrix} = \mathbf{0}\end{aligned}$$

- ...and the ML estimator  $\hat{\boldsymbol{\theta}}_{MLE} = \left( \hat{\boldsymbol{\beta}}_{MLE}; \hat{\sigma}_{MLE}^2 \right)$  is:

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{MLE} &= \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \sum_{i=1}^N \mathbf{x}_i y_i \\ \hat{\sigma}_{MLE}^2 &= \frac{\sum_{i=1}^N \left( y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right)^2}{N}\end{aligned}$$

- ...and it is shown to satisfy the Second Order Conditions for a maximum.



## Example: MLE and linear regression (4/4)

- The ML estimator for  $\boldsymbol{\beta}$  is identical to OLS, while the ML estimator for  $\sigma^2$  is smaller than the unbiased estimator  $\widehat{\sigma}^2$  from the small sample analysis of OLS.
- The latter is larger by a factor  $\frac{N}{N-K}$ : hence  $\widehat{\sigma}_{MLE}^2$  is more efficient even if biased.
- Suppose that  $\mathbf{x}_i$  is not fixed; only assume that *conditional on any realization*  $\mathbf{x}_i$ , the error term is normal and that it has constant variance:  $\varepsilon_i | \mathbf{x}_i \sim \mathcal{N}(0, \sigma^2)$ .
- As  $\varepsilon_i = y_i - \mathbf{x}_i^T \boldsymbol{\beta}$ , the following conditional density applies.

$$f_{Y|\mathbf{x}}(y_i | \mathbf{x}_i; \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}\right)$$

- Hence, the more flexible CMLE has the same solution as in the unrealistic “fixed regressors” case!

## Consistency of M-Estimators (1/5)

- It is necessary to study the asymptotics of M-Estimators.
- To prove consistency, one must show that the maximizer of the sample average criterion converges in probability to the maximizer of the *population expected* criterion.

$$\hat{\boldsymbol{\theta}}_M = \max_{\boldsymbol{\theta} \in \Theta} \hat{Q}_N(\boldsymbol{\theta}) \xrightarrow{p} \max_{\boldsymbol{\theta} \in \Theta} Q_0(\boldsymbol{\theta}) = \boldsymbol{\theta}_0$$

- Intuitively, **pointwise convergence** of  $\hat{Q}_N(\boldsymbol{\theta})$  to  $Q_0(\boldsymbol{\theta})$  – for all  $\boldsymbol{\theta} \in \Theta$  – is necessary.

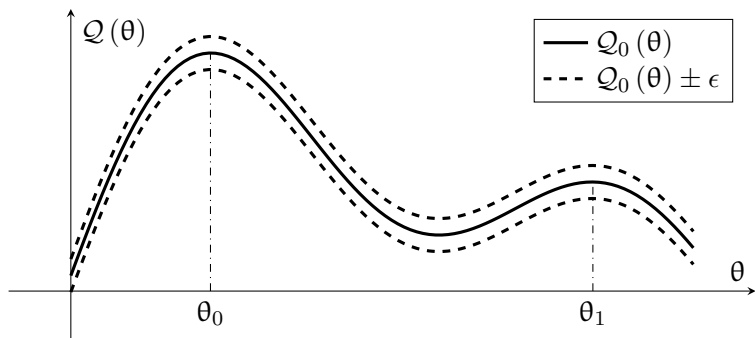
$$\left| \frac{1}{N} \sum_{i=1}^N \{q(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbb{E}[q(\mathbf{x}_i; \boldsymbol{\theta})]\} \right| \xrightarrow{p} 0$$

- However, it is not sufficient – unlike the stronger **uniform convergence** condition.

$$\sup_{\boldsymbol{\theta} \in \Theta} \left| \frac{1}{N} \sum_{i=1}^N \{q(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbb{E}[q(\mathbf{x}_i; \boldsymbol{\theta})]\} \right| \xrightarrow{p} 0$$

## Consistency of M-Estimators (2/5)

- Uniform convergence requires in fact that  $\hat{Q}_N(\theta)$  converges in probability towards  $Q_0(\theta)$  “at the same speed” over the entire parameter space  $\Theta$ .
- The intuition is graphically represented in the figure below:  
 $\hat{Q}_N(\theta) \in [Q_0(\theta) - \epsilon, Q_0(\theta) + \epsilon]$  for any  $\epsilon > 0$  and  $\forall \theta \in \Theta$   
( $\hat{Q}_N$  always stays within a “sleeve” of  $Q(\theta)$ ).



## Consistency of M-Estimators (3/5)

Uniform convergence is ensured if these four conditions hold:

- i.*  $q(\mathbf{x}_i; \boldsymbol{\theta})$  is continuous;
- ii.*  $\Theta$  is a compact set;
- iii.*  $\mathbb{E}[|q(\mathbf{x}_i; \boldsymbol{\theta})|] < \infty$ : that is,  $q(\mathbf{x}_i; \boldsymbol{\theta})$  has a bounded first absolute moment;
- iv.*  $q(\mathbf{x}_i; \boldsymbol{\theta})$  is Borel-measurable on its support.

These conditions – together – allow to invoke a result known as **Uniform Weak Law of Large Numbers**, which then implies uniform convergence.

These conditions are technical; note that *i.* and *ii.* relate to the fact that M-Estimators are maxima; furthermore *iii.* and *iv.* are analogous to conditions from other Laws of Large Numbers.

The full-fledged proof by Newey and McFadden (1994) follows.

# Consistency of M-Estimators (4/5)

## Theorem 2

**Consistency of M-Estimators.** *If i.  $Q_0(\theta)$  is uniquely maximized at  $\theta_0$ , ii.  $\Theta$  is a compact set, iii.  $Q_0(\theta)$  is a continuous function, and iv.  $\hat{Q}_N(\theta)$  uniformly converges in probability to  $Q_0(\theta)$ , it follows that M-Estimators are consistent.*

### Proof.

For any  $\epsilon > 0$ , with probability approaching 1 (w.p.a. 1);

$$\text{by } i.: \quad \hat{Q}_N(\hat{\theta}_M) > \hat{Q}_N(\theta_0) - \frac{\epsilon}{3} \quad (a)$$

$$\text{by } iv.: \quad Q_0(\hat{\theta}_M) > \hat{Q}_N(\hat{\theta}_M) - \frac{\epsilon}{3} \quad (b)$$

$$\text{by } iv.: \quad \hat{Q}_N(\theta_0) > Q_0(\theta_0) - \frac{\epsilon}{3} \quad (c)$$

which allows to show that w.p.a. 1,  $Q_0(\hat{\theta}_M) > Q_0(\theta_0) - \epsilon$ .

(Continues...)

# Consistency of M-Estimators (5/5)

## Theorem 2

### Proof.

(Continued.) Such a property is shown via a chain of inequalities.

$$\mathcal{Q}_0(\widehat{\boldsymbol{\theta}}_M) \stackrel{(b)}{>} \widehat{\mathcal{Q}}_N(\widehat{\boldsymbol{\theta}}_M) - \frac{\epsilon}{3} \stackrel{(a)}{>} \widehat{\mathcal{Q}}_N(\boldsymbol{\theta}_0) - \frac{2\epsilon}{3} \stackrel{(c)}{>} \mathcal{Q}_0(\boldsymbol{\theta}_0) - \epsilon$$

Now, denote by  $\mathbb{U}$  any given open neighborhood of  $\boldsymbol{\theta}_0$  and by  $\mathbb{U}^c$  its complement in  $\boldsymbol{\Theta}$ . Also define,

$$\mathcal{Q}_0(\boldsymbol{\theta}^*) = \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta} \cap \mathbb{U}^c} \mathcal{Q}_0(\boldsymbol{\theta})$$

for some  $\boldsymbol{\theta}^*$ , and notice that  $\mathcal{Q}_0(\boldsymbol{\theta}^*) < \mathcal{Q}_0(\boldsymbol{\theta}_0)$  by *i.-ii.-iii.*: thus, by setting  $\epsilon = \mathcal{Q}_0(\boldsymbol{\theta}_0) - \mathcal{Q}_0(\boldsymbol{\theta}^*)$  it follows that:

$$\mathcal{Q}_0(\widehat{\boldsymbol{\theta}}_M) > \mathcal{Q}_0(\boldsymbol{\theta}_0) - \epsilon \Rightarrow \mathcal{Q}_0(\widehat{\boldsymbol{\theta}}_M) > \mathcal{Q}_0(\boldsymbol{\theta}^*)$$

implying  $\widehat{\boldsymbol{\theta}}_M \in \mathbb{U}$  for any open neighborhood  $\mathbb{U}$ , and so  $\widehat{\boldsymbol{\theta}}_M \xrightarrow{p} \boldsymbol{\theta}_0$ .  $\square$

# Asymptotic normality of M-Estimators (1/5)

The next step is the analysis of M-Estimators' limiting distribution.

## Theorem 3

**Asymptotic Normality of M-Estimators.** *A given M-Estimator  $\widehat{\theta}_M$  follows an asymptotically normal distribution if the following five conditions hold simultaneously:*

- i.  $\widehat{\theta}_M$  is a consistent estimator of  $\theta_0$ ;*
- ii.  $q(\mathbf{x}_i; \theta)$  is a concave and twice continuously differentiable function in an open neighborhood of  $\theta_0$ ;*
- iii.  $\frac{\partial}{\partial \theta} \mathbb{E}[q(\mathbf{x}_i; \theta)] = \mathbb{E}\left[\frac{\partial}{\partial \theta} q(\mathbf{x}_i; \theta)\right]$ : the derivative can pass through the expectation integral;*
- iv. the data meet the requirements for the application of a Central Limit Theorem (the data are “well behaved”);*
- v. the Hessian matrix is nonsingular, it is continuous in  $\theta$  and it has a bounded absolute first moment.*

*(Continues...)*

# Asymptotic normality of M-Estimators (2/5)

## Theorem 3

*(Continued.) The limiting distribution is:*

$$\sqrt{N} \left( \hat{\boldsymbol{\theta}}_M - \boldsymbol{\theta}_0 \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \mathbf{Q}_0^{-1} \boldsymbol{\Upsilon}_0 \mathbf{Q}_0^{-1} \right)$$

where  $\mathbf{Q}_0$  and  $\boldsymbol{\Upsilon}_0$  are defined as the following probability limits:

$$\lim_{N \rightarrow \infty} \text{Var} \left[ \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) \right] \xrightarrow{p} \boldsymbol{\Upsilon}_0$$
$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{H}_i(\mathbf{x}_i; \boldsymbol{\theta}_0)] \xrightarrow{p} \mathbf{Q}_0$$

which implies the following asymptotic distribution, for a fixed  $N$ .

$$\hat{\boldsymbol{\theta}}_M \overset{A}{\sim} \mathcal{N} \left( \boldsymbol{\theta}_0, \frac{1}{N} \mathbf{Q}_0^{-1} \boldsymbol{\Upsilon}_0 \mathbf{Q}_0^{-1} \right)$$



# Asymptotic normality of M-Estimators (3/5)

## Theorem 3

### Proof.

The derivation is analogous to those given in Lecture 6. By the mean value theorem:

$$\mathbf{s}_i(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_M) = \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) + \mathbf{H}_i(\mathbf{x}_i; \tilde{\boldsymbol{\theta}}_N) (\hat{\boldsymbol{\theta}}_M - \boldsymbol{\theta}_0)$$

where  $\tilde{\boldsymbol{\theta}}_N$  is some convex combination of  $\hat{\boldsymbol{\theta}}_M$  and  $\boldsymbol{\theta}_0$ . Summing over the  $N$  observations and dividing by  $\sqrt{N}$ , one gets:

$$\begin{aligned} \mathbf{0} &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_M) \\ &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) + \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{H}_i(\mathbf{x}_i; \tilde{\boldsymbol{\theta}}_N) \right] \sqrt{N} (\hat{\boldsymbol{\theta}}_M - \boldsymbol{\theta}_0) \end{aligned}$$

by recalling that the sample score evaluated at the solution is equal to zero by definition of M-Estimators. (**Continues...**)

# Asymptotic normality of M-Estimators (4/5)

## Theorem 3

### Proof.

(Continued.) The expression above can be rewritten as:

$$\sqrt{N} \left( \hat{\boldsymbol{\theta}}_M - \boldsymbol{\theta}_0 \right) = - \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{H}_i \left( \mathbf{x}_i; \tilde{\boldsymbol{\theta}}_N \right) \right]^{-1} \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i \left( \mathbf{x}_i; \boldsymbol{\theta}_0 \right)$$

as condition  $v$ . lets invert the average Hessian matrix. Consider that:

1. By  $i$ . and  $v$ . one can apply some suitable Law of Large Numbers to the “sample-averaged” Hessian matrix, showing that:

$$\frac{1}{N} \sum_{i=1}^N \mathbf{H}_i \left( \mathbf{x}_i; \tilde{\boldsymbol{\theta}}_N \right) \xrightarrow{p} \mathbf{Q}_0$$

which follows from the Continuous Mapping Theorem given that, at the same time,  $\tilde{\boldsymbol{\theta}}_N \xrightarrow{p} \boldsymbol{\theta}_0$ . (Continues...)

# Asymptotic normality of M-Estimators (5/5)

## Theorem 3

### Proof.

(Continued.)

2. condition *iii.* yields  $\frac{\partial \mathcal{Q}(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0)] = \mathbf{0}$ ,  
thus:

$$\frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) \xrightarrow{p} \mathbf{0}$$

so by condition *iv.* and the Continuous Mapping Theorem, it is as follows.

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \boldsymbol{\Upsilon}_0)$$

All the intermediate results are – as usual – recombined via Slutskij's Theorem and the Cramér-Wold Device to show the desired result.  $\square$

## Inference for M-Estimators (1/2)

- To perform inference on M-Estimators,  $\Upsilon_0$  and  $\mathbf{Q}_0$  must be estimated: as usual, the analogy principle is applied here.
- The “bread” matrix  $\mathbf{Q}_0$  is estimated easily.

$$\hat{\mathbf{Q}}_N \equiv \frac{1}{N} \sum_{i=1}^N \mathbf{H}_i(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_M) \xrightarrow{p} \mathbf{Q}_0$$

- With the “meat” matrix  $\Upsilon_0$ , distributional assumptions on the sample matter. If the observations are independent but possibly not identically distributed, it is:

$$\Upsilon_0 = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) \mathbf{s}_i^T(\mathbf{x}_i, \boldsymbol{\theta}_0) \right]$$

with  $\Upsilon_0 = \mathbb{E} \left[ \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_0) \mathbf{s}_i^T(\mathbf{x}_i, \boldsymbol{\theta}_0) \right]$  in the i.i.d. case.

## Inference for M-Estimators (2/2)

- With independent observations  $\boldsymbol{\Upsilon}_0$  is estimated as follows.

$$\widehat{\boldsymbol{\Upsilon}}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \widehat{\boldsymbol{\theta}}_M) \mathbf{s}_i^T(\mathbf{x}_i, \widehat{\boldsymbol{\theta}}_M) \xrightarrow{p} \boldsymbol{\Upsilon}_0$$

- Dependent observations complicate matters. The HAC case is feasible, but often complicated. The CCE case yields:

$$\widehat{\boldsymbol{\Upsilon}}_{CCE} = \frac{1}{N} \sum_{c=1}^C \sum_{i=1}^{N_c} \sum_{j=1}^{N_c} \mathbf{s}_{ic}(\mathbf{x}_{ic}; \widehat{\boldsymbol{\theta}}_M) \mathbf{s}_{jc}^T(\mathbf{x}_{jc}; \widehat{\boldsymbol{\theta}}_M) \xrightarrow{p} \boldsymbol{\Upsilon}_0$$

where both observations and scores are also indexed by the group or cluster  $c = 1, \dots, C$  they belong to.

- For a suitable estimator  $\widehat{\boldsymbol{\Upsilon}}_N \xrightarrow{p} \boldsymbol{\Upsilon}_0$ , the variance-covariance matrix of M-Estimators is ultimately estimated as follows.

$$\widehat{\text{Avar}}(\widehat{\boldsymbol{\theta}}_M) = \frac{1}{N} \widehat{\mathbf{Q}}_N^{-1} \widehat{\boldsymbol{\Upsilon}}_N \widehat{\mathbf{Q}}_N^{-1}$$

## Examples: Asymptotics of OLS and NLLS (1/3)

- This analysis is reconciled with that about OLS (Lecture 8) by noting that there consistency follows from  $\mathbb{E}[\varepsilon_i | \mathbf{x}_i] = 0$ , and that in that case it is  $\mathbf{Y}_0 = 4\mathbf{\Xi}_0$  and  $\mathbf{Q}_0 = -2\mathbf{K}_0$ .
- To discuss NNLS, it is useful to define a  $K \times 1$  vector: that is, the derivative of the CEF evaluated at  $\mathbf{x}_i$  and at  $\hat{\boldsymbol{\theta}}_{NLLS}$ .

$$\hat{\mathbf{h}}_i \equiv \frac{\partial}{\partial \boldsymbol{\theta}} h(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_{NLLS})$$

- Recall that by construction, NLLS sets the average score at zero for every value of  $N$ .

$$\mathbf{0} = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(y_i, \mathbf{x}_i; \hat{\boldsymbol{\theta}}_{NLLS}) = \frac{1}{N} \sum_{i=1}^N 2\hat{\mathbf{h}}_i \varepsilon_i \xrightarrow{p} \mathbf{0}$$

- Refer to it as the “condition about the probability limit of the score.”

## Examples: Asymptotics of OLS and NLLS (2/3)

- Also recall the CEF motivation of NLLS.

$$\mathbb{E}[\varepsilon_i | \mathbf{x}_i] = \mathbb{E}[y_i | \mathbf{x}_i] - h(\mathbf{x}_i; \boldsymbol{\theta}_0) = 0$$

- The immediate implication of such a CEF condition is that:

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \left( \frac{\partial}{\partial \boldsymbol{\theta}} h(\mathbf{x}_i; \boldsymbol{\theta}_0) \right) \varepsilon_i \right] &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{h}_{0i} \varepsilon_i] \\ &= \mathbf{0} \end{aligned}$$

which can be reconciled with the above condition about the probability limit of the score if  $\hat{\mathbf{h}}_i \xrightarrow{p} \mathbf{h}_{0i}$  for  $i = 1, \dots, N$  as per some applicable Law of Large Numbers.

- The Continuous Mapping Theorem finally implies that:

$$\hat{\boldsymbol{\theta}}_{NLLS} \xrightarrow{p} \boldsymbol{\theta}_0$$

that is, the NLLS estimator is indeed consistent.

## Examples: Asymptotics of OLS and NLLS (3/3)

- Regarding the asymptotic distribution, note that under the hypothesis of *independent observations* the following holds.

$$\mathbf{Q}_0 = \lim_{N \rightarrow \infty} -\frac{1}{N} \sum_{i=1}^N 2 \cdot \mathbb{E} \left[ \mathbf{h}_{0i} \mathbf{h}_{0i}^T \right]$$

$$\mathbf{R}_0 = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N 4 \cdot \mathbb{E} \left[ \varepsilon_i^2 \mathbf{h}_{0i} \mathbf{h}_{0i}^T \right]$$

- Therefore, the residual  $e_i \equiv y_i - h(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_{NLLS})$  delivers:

$$\widehat{\text{Avar}} \left( \hat{\boldsymbol{\theta}}_{NLLS} \right) = \left[ \sum_{i=1}^N \hat{\mathbf{h}}_i \hat{\mathbf{h}}_i^T \right]^{-1} \left[ \sum_{i=1}^N e_i^2 \hat{\mathbf{h}}_i \hat{\mathbf{h}}_i^T \right] \left[ \sum_{i=1}^N \hat{\mathbf{h}}_i \hat{\mathbf{h}}_i^T \right]^{-1}$$

the **heteroscedasticity-consistent** estimator of the NLLS variance-covariance, akin to the OLS “robust” formula.

- The homoscedastic, CCE and HAC versions are akin too.



# The Trinity of Asymptotic Tests

After estimating an econometric model, a researcher is typically interested into performing some **tests of hypothesis**; these are possibly non-linear:

$$H_0 : \mathbf{v}(\boldsymbol{\theta}_0) = \mathbf{0}$$

$$H_1 : \mathbf{v}(\boldsymbol{\theta}_0) \neq \mathbf{0}$$

where  $\mathbf{v}(\cdot)$ , a vector-valued function, has length  $L$  (for multiple hypotheses).

There are three alternative methods to perform such tests; they are known together as the “Trinity.”

1. Generalized Wald Statistics.
2. Distance, or Likelihood Ratio test.
3. Score – or Lagrange multiplier – test.

These methods are briefly review under the unifying framework of M-Estimation.

# The Generalized Wald Statistics (1/2)

The **generalized** Wald Statistics is:

$$\widetilde{W}_{H_0} = \mathbf{v}^T(\widehat{\boldsymbol{\theta}}_M) \cdot \left[ \widehat{\mathbf{V}} \cdot \widehat{\text{Avar}}(\widehat{\boldsymbol{\theta}}_M) \cdot \widehat{\mathbf{V}}^T \right]^{-1} \cdot \mathbf{v}(\widehat{\boldsymbol{\theta}}_M)$$

where  $\widehat{\mathbf{V}} \equiv \frac{\partial}{\partial \boldsymbol{\theta}^T} \mathbf{v}(\widehat{\boldsymbol{\theta}}_M)$ . Its limiting distribution is:

$$\widetilde{W}_{H_0} \xrightarrow{d} \chi_L^2$$

that is, **under the null hypothesis**  $H_0$  the test statistic has a limiting  $\chi_L^2$  distribution with  $L$  degrees of freedom.

To gain intuition recall the Wald Statistic from the linear model (where  $\mathbf{v}(\boldsymbol{\beta}) = \mathbf{R}\boldsymbol{\beta} - \mathbf{c} = \mathbf{0}$  is some linear function). There, the Wald Statistic is a special case of Hotelling's  $t$ -squared statistic, which is asymptotically chi-squared distributed. This non-linear case is analogously derived by exploiting the Delta Method.

## The Generalized Wald Statistics (2/2)

**Example.** Suppose that interest lies in one specific hypothesis about the linear model:

$$H_0 : \sum_{k=1}^K \beta_k = 1 \qquad H_1 : \sum_{k=1}^K \beta_k \neq 1$$

like say the hypothesis of constant return to scale in production functions (with the constant parameter being  $\beta_0$ ). In this case:

$$\widetilde{W}_{H_0} = \frac{\left(\sum_{k=1}^K \widehat{\beta}_{k,OLS} - 1\right)^2}{\sum_{k=1}^K \sum_{q=1}^K \widehat{\sigma}_{\beta_k q}} \xrightarrow{d} \chi_1^2$$

where the  $kq$ -th element of the estimated variance-covariance of the OLS estimates is denoted as  $\widehat{\sigma}_{\beta_k q}$ .

The Wald Test is very easy to implement but it suffers from two issues. First, it performs poorly in small samples; and second, it is not transformation-invariant – for example, different statistics are calculated whether  $H_0 : \beta_k = 0$  or  $H_0 : \exp(\beta_k) = 1$ .

## The Distance, or Likelihood Ratio test (1/2)

The “Distance Test” or “Likelihood Ratio Test,” originates with MLE. With respect to the simpler Generalized Wald Test, it is transformation-invariant and deals non-linear hypotheses quite well. However, it is more computationally demanding.

In fact, to perform it one must compute the main estimate of  $\theta$  as well as an additional “restricted” estimate

$$\hat{\theta}_V = \arg \max_{\theta \in \Theta_V} \hat{Q}_N(\theta)$$

with a “restricted parameter space”  $\Theta_V = \{\theta \in \Theta : v(\theta) = \mathbf{0}\}$ . Then, the “Distance” Statistic in all cases but MLE is:

$$D_{H_0} = N \left[ \hat{Q}_N(\hat{\theta}_M) - \hat{Q}_N(\hat{\theta}_V) \right] \xrightarrow{d} \chi_L^2$$

while:

$$LR_{H_0} = 2 \left[ \log \hat{Q}_N(\hat{\theta}_M) - \log \hat{Q}_N(\hat{\theta}_V) \right] \xrightarrow{d} \chi_L^2$$

is the “Likelihood Ratio” from MLE, where  $\hat{Q}_N(\theta) = \hat{\mathcal{L}}_N(\theta)$ .

## The Distance, or Likelihood Ratio test (2/2)

Intuitively, the test compares the advantage in terms of **fitting** when letting the model be estimated without the restriction.

**Example.** Expanding on the toy test that used the Generalized Wald Statistic, one can estimate a “restricted” model, such as:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots \\ + \beta_{K-1} X_{(K-1)i} + \left(1 - \sum_{k=1}^{K-1} \beta_k\right) X_{Ki} + \epsilon_i$$

which can also be written, for  $\ddot{X}_{ki} \equiv X_{ki} - X_{Ki}$  ( $k = 1, \dots, K$ ):

$$Y_i - X_{Ki} = \beta_0 + \beta_1 \ddot{X}_{1i} + \beta_2 \ddot{X}_{2i} + \dots + \beta_{K-1} \ddot{X}_{(K-1)i} + \epsilon_i$$

restricting  $\beta_K$  or possibly another coefficient (except  $\beta_0$ ). Thus:

$$D_{H_0} = \left[ \sum_{i=1}^N \left( y_i - x_{Ki} - \ddot{\mathbf{x}}_i^T \hat{\boldsymbol{\beta}}_V \right)^2 - \sum_{i=1}^N \left( y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{OLS} \right)^2 \right] \xrightarrow{d} \chi^2$$

is the Distance Test, and  $\hat{\boldsymbol{\beta}}_V$  are the “restricted” estimates.

## The score, or Lagrange multiplier, test (1/2)

This last test has the same advantages as the previous one, but it does not require the “unrestricted” model to be estimated: it is computationally more parsimonious. The test is based on the properties of the sample average score function evaluated at **one specific parameter value**  $\theta_v$  implied by the null hypothesis.

By construction, in all M-Estimators the average score is always equal to zero when evaluated at the unrestricted estimate  $\hat{\theta}_M$ .

$$\frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \hat{\theta}_M) = \mathbf{0}$$

Consider a *restricted* parameter value  $\theta_v$  such that  $\mathbf{v}(\theta_v) = \mathbf{0}$ . It follows that:

$$\frac{1}{N} \sum_{i=1}^N |\mathbf{s}_i(\mathbf{x}_i; \theta_v)| > 0$$

the  $K$ -dimensional sample score vector deviates from zero when evaluated at any “suboptimal” parameter choice.

## The score, or Lagrange multiplier, test. (2/2)

The Lagrange Multiplier statistic is:

$$\text{LM}_{H_0} = N \left[ \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_v) \right]^T \widehat{\boldsymbol{\Upsilon}}_v^{-1} \left[ \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_v) \right] \xrightarrow{d} \chi_K^2$$

where:

$$\widehat{\boldsymbol{\Upsilon}}_v = \widehat{\text{Avar}} \left[ \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_i; \boldsymbol{\theta}_v) \right].$$

It indicates how “statistically relevant” are the deviations of the restricted score from the zero benchmark.

**Example.** Return again to the running examples: consider the “restricted” estimates  $\boldsymbol{\beta}_v$  obtained by forcing all key coefficients to sum up to one. The Lagrange Multiplier Statistic is:

$$\text{LM}_{H_0} = \left[ \sum_{i=1}^N \mathbf{x}_i e_i(\boldsymbol{\beta}_v) \right]^T \left[ \sum_{i=1}^N e_i^2(\boldsymbol{\beta}_v) \mathbf{x}_i \mathbf{x}_i^T \right]^{-1} \left[ \sum_{i=1}^N \mathbf{x}_i e_i(\boldsymbol{\beta}_v) \right]$$

with  $\text{LM}_{H_0} \xrightarrow{d} \chi_K^2$ , given the new residuals  $e_i(\boldsymbol{\beta}_v) = y_i - \mathbf{x}_i^T \boldsymbol{\beta}_v$ .

## The Maximum Likelihood special case

- As seen in Lectures 5 and 6, with MLE if the data are i.i.d. the Information Matrix Equality applies.

$$\Upsilon_0 = -\mathbf{Q}_0$$

- Hence, the asymptotic distribution is:

$$\hat{\boldsymbol{\theta}}_{MLE} \stackrel{A}{\sim} \mathcal{N}(\boldsymbol{\theta}_0, [\mathbf{I}_N(\boldsymbol{\theta}_0)]^{-1})$$

where  $\mathbf{I}_N(\boldsymbol{\theta}_0)$  is the information matrix.

- Importantly, the information matrix hits the **Cramér-Rao bound**, making MLE the most efficient estimator under its distributional assumptions.
- The information matrix can be **estimated** either by  $\hat{\boldsymbol{\Upsilon}}_N$  or via  $-\hat{\mathbf{Q}}_N$ ; the former approach is called **outer product of the gradients**, typically more computationally convenient.



## Example: MLE linear regression expanded (1/2)

- Return to the MLE framework for linear regression. There, the score is as follows.

$$\mathbf{s}_i(y_i, \mathbf{x}_i; \boldsymbol{\theta}) = \begin{bmatrix} \mathbf{x}_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \sigma^{-2} \\ -\frac{1}{2} \sigma^{-2} + \frac{1}{2} (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \sigma^{-4} \end{bmatrix}$$

- The (symmetric) Hessian matrix is instead as follows.

$$\begin{aligned} \mathbf{H}_i(y_i, \mathbf{x}_i; \boldsymbol{\theta}) &= \\ &= \begin{bmatrix} -\mathbf{x}_i \mathbf{x}_i^T \sigma^{-2} & -\mathbf{x}_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \sigma^{-4} \\ -\mathbf{x}_i^T (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \sigma^{-4} & \frac{1}{2} \sigma^{-4} - (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \sigma^{-6} \end{bmatrix} \end{aligned}$$

- By plugging the MLE estimators  $\hat{\boldsymbol{\beta}}_{MLE}$  and  $\hat{\sigma}_{MLE}^2$  into the above expressions one can construct **consistent** estimators for the information matrix.

## Example: MLE linear regression expanded (2/2)

- The estimator based on the Hessian is as follows.

$$\begin{aligned}\frac{\hat{\mathbf{Q}}_N^{-1}}{N} &= \left[ \sum_{i=1}^N \mathbf{H}_i \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\theta}}_{MLE} \right) \right]^{-1} \\ &= - \begin{bmatrix} \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \hat{\sigma}_{MLE}^2 & \mathbf{0} \\ \mathbf{0}^T & \frac{2}{N} \hat{\sigma}_{MLE}^4 \end{bmatrix}\end{aligned}$$

- The next is based on the outer product of the gradients.

$$\begin{aligned}\frac{\hat{\mathbf{Y}}_N^{-1}}{N} &= \left[ \sum_{i=1}^N \mathbf{s}_i \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\theta}}_{MLE} \right) \mathbf{s}_i^T \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\theta}}_{MLE} \right) \right]^{-1} \\ &= \begin{bmatrix} \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \hat{\sigma}_{MLE}^2 & \mathbf{0} \\ \mathbf{0}^T & \frac{2}{N} \hat{\sigma}_{MLE}^4 \end{bmatrix}\end{aligned}$$

- Note the different sign, the zero vectors at the border, and the bottom-right element for the variance of  $\hat{\sigma}_{MLE}^2$ .

## Quasi-Maximum Likelihood (1/2)

- The nice properties of MLE break apart when one deviates from the i.i.d. benchmark. For example, in case of clustered dependence the likelihood function reads:

$$\mathcal{L}(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{c=1}^C f_{\mathbf{x}_1, \dots, \mathbf{x}_{N_c}}(\mathbf{x}_{1c}, \dots, \mathbf{x}_{N_c}; \boldsymbol{\theta})$$

but it cannot be factorized further. Hence, the Information Matrix Equality no longer holds and  $\hat{\boldsymbol{\Upsilon}}_{CCE}$  must be used.

- Things are even worse when MLE is built around incorrect distributional assumptions (even if the data are i.i.d.).
- In *misspecification* cases, the estimator in question is called **Quasi-Maximum Likelihood Estimator**  $\hat{\boldsymbol{\theta}}_{QMLE}$  and it has a probability limit called **pseudo-true value**  $\boldsymbol{\theta}^*$ .

$$\hat{\boldsymbol{\theta}}_{QMLE} = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_N) \xrightarrow{p} \boldsymbol{\theta}^*$$

## Quasi-Maximum Likelihood (2/2)

- A key question is: when is the Quasi-Maximum Likelihood Estimator (QMLE) **consistent**, that is  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$ ?
- In the MLE **linear** regression example,  $\widehat{\boldsymbol{\beta}}_{MLE}$  is consistent under a linear CEF even if the errors are not normal!
- If the assumed distribution is  $f_{Y,\mathbf{x}}(y_i, \mathbf{x}_i; \boldsymbol{\theta})$  for some scalar endogenous variable  $Y_i$ , and it belongs to the **exponential macro-family** – it can be decomposed in terms of primitive scalar functions  $a[\cdot]$ ,  $b[\cdot]$  and  $c[\cdot]$  as:

$$\begin{aligned} f_{Y,\mathbf{x}}(y_i, \mathbf{x}_i; \boldsymbol{\theta}) &= \\ &= \exp \left\{ a \left[ \mu_{Y|\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}) \right] + b[y_i] + c \left[ \mu_{Y|\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}) \right] y_i \right\} \end{aligned}$$

where:  $\mu_{Y|\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta}) \equiv \mathbb{E}[Y_i | \mathbf{x}_i; \boldsymbol{\theta}]$  is a **correctly specified** parametric expression of the CEF of  $Y_i$  given  $\mathbf{x}_i$ , in general (Q)MLE consistently estimates the function  $\mu_{Y|\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\theta})$ .

## Example: Poisson regression (1/4)

- A **count data model** is a model for explaining variables  $Y_i$  that assume non-negative integer values  $Y_i = 0, 1, 2, \dots$  *only* and where smaller values occur more frequently.
- An example is the **Poisson regression**, that is:

$$\mathbb{P}(Y_i | \mathbf{x}_i) = \frac{\lambda_i(\mathbf{x}_i)^{Y_i} \exp(-\lambda_i(\mathbf{x}_i))}{Y_i!}$$

where the Poisson parameter  $\lambda_i(\mathbf{x}_i)$  is treated as a function of the individual characteristics  $\mathbf{x}_i$ .

- By the properties of the Poisson distribution:

$$\lambda_i(\mathbf{x}_i) = \mathbb{E}[Y_i | \mathbf{x}_i] = \text{Var}[Y_i | \mathbf{x}_i]$$

that is,  $\lambda_i(\mathbf{x}_i)$  equals both the **conditional mean** and the **conditional variance** of  $Y_i$  given  $\mathbf{x}_i$ .

## Example: Poisson regression (2/4)

- The most common choice is  $\lambda_i(\mathbf{x}_i) = \exp(\mathbf{x}_i^T \boldsymbol{\beta}_0)$ , yielding:

$$\mathbb{P}(Y_i | \mathbf{x}_i) = \frac{\exp(Y_i \cdot \mathbf{x}_i^T \boldsymbol{\beta}_0) \exp[-\exp(\mathbf{x}_i^T \boldsymbol{\beta}_0)]}{Y_i!}$$

- ... an implication of which is:

$$\begin{aligned} \frac{\partial \exp(\mathbf{x}_i^T \boldsymbol{\beta}_0)}{\partial \mathbf{x}_i} &= \exp(\mathbf{x}_i^T \boldsymbol{\beta}_0) \boldsymbol{\beta}_0 \\ \Rightarrow \boldsymbol{\beta}_0 &= \frac{\partial \mathbb{E}[Y_i | \mathbf{x}_i]}{\partial \mathbf{x}_i} \frac{1}{\mathbb{E}[Y_i | \mathbf{x}_i]} \end{aligned}$$

hence,  $\boldsymbol{\beta}_0$  can be interpreted as a **semi-elasticity** like in a *log-lin* model that can be estimated via OLS.

- Unlike a *log-lin* model the Poisson regression allows  $Y_i = 0$ , which can be often useful or convenient!

## Example: Poisson regression (3/4)

- The **log-likelihood** function, given a sample  $\{(y_i, \mathbf{x}_i)\}_{i=1}^N$ , writes as follows for this model.

$$\begin{aligned}\log \mathcal{L} \left( \boldsymbol{\beta} \mid \{(y_i, \mathbf{x}_i)\}_{i=1}^N \right) &= \\ &= \sum_{i=1}^N \left[ y_i \cdot \mathbf{x}_i^T \boldsymbol{\beta} - \exp \left( \mathbf{x}_i^T \boldsymbol{\beta} \right) - \log y_i! \right]\end{aligned}$$

- The First Order Conditions of the MLE problem, expressed as the sum of the individual scores, consequently are:

$$\sum_{i=1}^N \mathbf{s}_i \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\beta}}_{MLE} \right) = \sum_{i=1}^N \mathbf{x}_i \left[ y_i - \exp \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right) \right] = \mathbf{0}$$

- ... and they lack a closed form solution; therefore, the MLE estimator must be obtained by numerical methods.

## Example: Poisson regression (4/4)

- The empirical Hessian matrix is:

$$\begin{aligned}\hat{\mathbf{Q}}_N &= \frac{1}{N} \sum_{i=1}^N \mathbf{H}_i \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\beta}}_{MLE} \right) \\ &= -\frac{1}{N} \sum_{i=1}^N \exp \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right) \cdot \mathbf{x}_i \mathbf{x}_i^T\end{aligned}$$

the opposite of which – if divided by  $N$  – is an appropriate estimator of the information matrix.

- The Poisson distribution belongs to the exponential family: therefore, even if the likelihood is misspecified, the **CEF** of  $Y_i$  given  $\mathbf{x}_i$  is consistently estimated if it is well specified.
- With the “exponential” specification for  $\lambda_i(\mathbf{x}_i)$  this implies that the MLE is interpretable in terms of semi-elasticities.



## Quasi-Maximum Likelihood: more remarks (1/3)

- The Poisson regression is an extreme case where QMLE is safe, whereas in the linear model  $\hat{\sigma}_{MLE}^2$  has little meaning without the normality assumption.
- This bears implications for inference. In general:

$$\sqrt{N} \left( \hat{\boldsymbol{\theta}}_{QMLE} - \boldsymbol{\theta}^* \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, [\mathbf{Q}^*]^{-1} \boldsymbol{\Upsilon}^* [\mathbf{Q}^*]^{-1} \right)$$

where  $\mathbf{Q}^*$  and  $\boldsymbol{\Upsilon}^*$  are analogues of  $\mathbf{Q}_0$  and  $\boldsymbol{\Upsilon}_0$  respectively, but evaluated at  $\boldsymbol{\theta}^*$  instead of  $\boldsymbol{\theta}_0$ .

- It is thus safer to conduct inference using the “sandwiched” formula of M-Estimation – the extensive one!
- What if QMLE is actually inconsistent:  $\boldsymbol{\theta}^* \neq \boldsymbol{\theta}_0$ ? Not all is lost: QMLE can still be given an interpretation in terms of “best approximation,” like OLS with a non-linear CEF.

## Quasi-Maximum Likelihood: more remarks (2/3)

- To see this, consider the **Kullback-Leibler Information Criterion** (KLIC).

$$\begin{aligned}\mathcal{K}_x(\boldsymbol{\theta}) &\equiv \mathbb{E}_g \left[ \log \left( \frac{g_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta}_0)}{f_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta})} \right) \right] = \\ &= \int_{\mathbb{X}} \log \left( \frac{g_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta}_0)}{f_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta})} \right) g_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta}_0) d\mathbf{x}_1 \dots d\mathbf{x}_N\end{aligned}$$

- Here,  $f_x(\cdot)$  is the **assumed** joint mass or density function that is thought to generate the data.
- Instead  $g_x(\cdot)$  is the **true** function, which is taken as given; the expectation is taken with respect to  $g_x(\cdot)$ .
- Note that by construction,  $\mathcal{K}_x(\boldsymbol{\theta}) \geq 0$  for all  $\boldsymbol{\theta} \in \Theta$ .
- If the distribution is correctly specified, it is  $f_x(\cdot) = g_x(\cdot)$ , and  $\mathcal{K}_x(\boldsymbol{\theta}_0) = 0$ : the KLIC would attain its minimum.

## Quasi-Maximum Likelihood: more remarks (3/3)

- In addition:

$$\begin{aligned}\mathcal{K}_x(\boldsymbol{\theta}) &= \mathbb{E}_g [\log g_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta}_0) - \log f_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta})] \\ &= \mathbb{E}_g [\log g_x(\mathbf{x}_1, \dots, \mathbf{x}_N; \boldsymbol{\theta}_0)] - \log \mathcal{L}_0^g(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_N)\end{aligned}$$

where  $\log \mathcal{L}_0^g(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_N)$  is here a **pseudo-population likelihood** that results if the assumed distribution is  $f_x(\cdot)$ , but the true one is  $g_x(\cdot)$ .

- So, under general assumptions:

$$\hat{\boldsymbol{\theta}}_{QMLE} \xrightarrow{p} \boldsymbol{\theta}^* = \max_{\boldsymbol{\theta} \in \Theta} \log \mathcal{L}_0^g(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_N) = \min_{\boldsymbol{\theta} \in \Theta} \mathcal{K}_x(\boldsymbol{\theta})$$

the pseudo-true value  $\boldsymbol{\theta}^*$ , to which the QMLE converges in probability, is both the maximizer of the pseudo-population likelihood and the minimizer of the KLIC!

- Trying to “augment” an MLE model, however misspecified, aids this KLIC-motivated “approximation” interpretation.

# Introduction to Binary Outcome Models (1/4)

A subclass of econometric models is characterized by a **limited dependent variable** (LDV):  $Y_i$  has a *discrete, finite* support.

Many questions in economics revolve in fact around a **binary** dependent variable:  $Y_i \in \{0, 1\}$ . Here are some examples.

- What are the determinants of enrollment in college?
- Which factors influence a firm's probability of default?
- What are the causes of civil wars in a country?

Other LDV models take multiple outcomes.

- What means of transportations do individuals choose?
- What determines a country's political regime?
- What type of insurance contract is preferred by people?
- Which individual characteristics predict survey responses?

## Introduction to Binary Outcome Models (2/4)

It is useful to overview models for **binary outcomes**  $Y_i \in \{0, 1\}$  so as to appreciate the usefulness of MLE in this setting.

- If one treats  $Y_i$  as random (Bernoulli) event, it is natural to think about its realization *probability* as **conditional** upon some observable variables  $\mathbf{x}_i$ :

$$\mathbb{P}(Y_i = 1 | \mathbf{x}_i) = G(\mathbf{x}_i, \boldsymbol{\beta}_0)$$

$$\mathbb{P}(Y_i = 0 | \mathbf{x}_i) = 1 - G(\mathbf{x}_i, \boldsymbol{\beta}_0)$$

where  $G(\cdot)$  is some function of  $\mathbf{x}_i$  parametrized by  $\boldsymbol{\beta}_0$ .

- As the problem is binary, the probability of either outcome can be treated residually with respect to the other's. Also, one can write the CEF of  $Y_i$  given  $\mathbf{x}_i$  as follows.

$$\begin{aligned}\mathbb{E}[Y_i | \mathbf{x}_i] &= 1 \cdot [G(\mathbf{x}_i, \boldsymbol{\beta}_0)] + 0 \cdot [1 - G(\mathbf{x}_i, \boldsymbol{\beta}_0)] \\ &= G(\mathbf{x}_i, \boldsymbol{\beta}_0)\end{aligned}$$

## Introduction to Binary Outcome Models (3/4)

- Can such a model be estimated via linear regression? This is called the **linear probability model** (LPM):

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}_0 + \epsilon_i, \quad y_i \in \{0, 1\}$$

which is equivalent to assuming  $G(\mathbf{x}_i, \boldsymbol{\beta}_0) = \mathbf{x}_i^T \boldsymbol{\beta}_0$ .

- By definition of regression:  $\epsilon_i = Y_i - \mathbb{E}[Y_i | \mathbf{x}_i] = Y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0$ ; this has a “natural heteroscedasticity” implication.

$$\mathbb{P}(\epsilon_i = 1 - \mathbf{x}_i^T \boldsymbol{\beta}_0 | \mathbf{x}_i) = \mathbf{x}_i^T \boldsymbol{\beta}_0$$

$$\mathbb{P}(\epsilon_i = -\mathbf{x}_i^T \boldsymbol{\beta}_0 | \mathbf{x}_i) = 1 - \mathbf{x}_i^T \boldsymbol{\beta}_0$$

- A consequent implication is:

$$\mathbb{E}[\epsilon_i | \mathbf{x}_i] = (1 - \mathbf{x}_i^T \boldsymbol{\beta}_0) \mathbf{x}_i^T \boldsymbol{\beta}_0 - \mathbf{x}_i^T \boldsymbol{\beta}_0 (1 - \mathbf{x}_i^T \boldsymbol{\beta}_0) = 0$$

hence, OLS still delivers unbiased and consistent estimates of  $\boldsymbol{\beta}_0$  even if the problem is naturally heteroscedastic.

## Introduction to Binary Outcome Models (4/4)

“Natural heteroscedasticity” is no issue for the LPM: it can be addressed via “robust” estimation, or FGLS in small samples.

The main issue is that the linear conditional expectation  $\mathbf{x}_i^T \boldsymbol{\beta}_0$  cannot be constrained to lie within the (0, 1) interval. Hence:

1. the conditional variance of the error term  $\epsilon_i$  takes negative values;

$$\text{Var}[\epsilon_i | \mathbf{x}_i] = \mathbf{x}_i^T \boldsymbol{\beta}_0 (1 - \mathbf{x}_i^T \boldsymbol{\beta}_0) \gtrless 0$$

2. the predicted probabilities  $\hat{\mathbb{E}}[Y_i | \mathbf{x}_i = \mathbf{x}_i] = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{LPM} = \hat{y}_i$  take values outside the  $[0, 1]$  interval.

Both facts that make **no probabilistic sense** and call for more parametric models estimated by MLE. Yet, the LPM is however used especially when interest falls on transparent **causal** effects of  $\mathbf{x}_i$  on  $Y_i$ , as it is easier to implement.

# MLE for Binary Outcome Models (1/8)

All MLE approaches to binary outcome models are based on an assumption of the kind that, for each realization  $\mathbf{x}_i$  of  $\mathbf{x}_i$ :

$$G(\mathbf{x}_i, \boldsymbol{\beta}_0) = F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0) = F_{\mathbf{x}}(\lambda_i)$$

where  $F_{\mathbf{x}}(\cdot)$  is some **probability distribution function** with one “free” parameter (usually a location parameter)  $\lambda_i = \mathbf{x}_i^T \boldsymbol{\beta}_0$ .

Most applications use a **symmetric distribution function**, so that  $F(\lambda_i) = 1 - F(-\lambda_i)$ . This solves the problem of predicted probabilities, since conditionally on any realization  $\mathbf{x}_i$ :

$$\lim_{\mathbf{x}_i^T \boldsymbol{\beta}_0 \rightarrow +\infty} \mathbb{P}(Y_i = 1 | \mathbf{x}_i) = \lim_{\mathbf{x}_i^T \boldsymbol{\beta}_0 \rightarrow +\infty} F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0) = 1$$

$$\lim_{\mathbf{x}_i^T \boldsymbol{\beta}_0 \rightarrow -\infty} \mathbb{P}(Y_i = 1 | \mathbf{x}_i) = \lim_{\mathbf{x}_i^T \boldsymbol{\beta}_0 \rightarrow -\infty} F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0) = 0$$

and symmetrically for  $Y_i = 0$ .



## MLE for Binary Outcome Models (2/8)

Another advantage is a model of this sort can be based upon a **structural** interpretation about “choices.” Specifically, write:

$$y_i^* = \mathbf{x}_i^T \boldsymbol{\beta}_0 + \varepsilon_i$$
$$y_i = \begin{cases} 1 & \text{if } y_i^* > \alpha_0 \\ 0 & \text{if } y_i^* \leq \alpha_0 \end{cases}$$

where  $Y_i^*$  is a **latent variable** that represents the **cost-benefit** evaluation of the binary choice by the  $i$ -th individual.

- The latent variable  $Y_i^*$  is unobserved, and it is a theoretical abstraction used for interpretation’s sake (like “utility”).

In such a model, conditionally on any realization  $\mathbf{x}_i$ :

$$\mathbb{P}(Y_i = 1 | \mathbf{x}_i) = \mathbb{P}(Y_i^* > \alpha_0 | \mathbf{x}_i) = \mathbb{P}(\varepsilon_i > -\mathbf{x}_i^T \boldsymbol{\beta}_0 + \alpha_0 | \mathbf{x}_i)$$

and intuitively, if  $\mathbf{x}_i$  contains a constant element the associated “intercept” parameter and  $\alpha_0$  are not separately identified. The **normalization**  $\alpha_0 = 0$  is thus typically imposed.

## MLE for Binary Outcome Models (3/8)

If  $F_{\mathbf{x}}(\cdot)$  is a symmetric distribution, the model reshapes as:

$$\begin{aligned}\mathbb{P}(Y_i = 1 | \mathbf{x}_i) &= \mathbb{P}(Y_i^* > 0 | \mathbf{x}_i) \\ &= \mathbb{P}(\varepsilon_i > -\mathbf{x}_i^T \boldsymbol{\beta}_0 | \mathbf{x}_i) \\ &= 1 - \mathbb{P}(\varepsilon_i \leq -\mathbf{x}_i^T \boldsymbol{\beta}_0 | \mathbf{x}_i) \\ &= \mathbb{P}(\varepsilon_i < \mathbf{x}_i^T \boldsymbol{\beta}_0 | \mathbf{x}_i) \\ &= F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0)\end{aligned}$$

where the fourth line exploits the symmetry of  $F_{\mathbf{x}}(\cdot)$ . This fact reconciles the latent variable model with the specification of the conditional probability for the outcome  $Y_i$ .

- Latent variable models are not specific of binary outcomes: multinomial LDV models are usually motivated by complex variations of this approach. Latent variable models are also used in the structural analysis of empirical strategic games.

## MLE for Binary Outcome Models (4/8)

Note that the distribution  $F_{\mathbf{x}}(\cdot)$  should *not* contain a **variable** scale parameter; if it is say normal, its variance must be known or normalized, e.g.  $\sigma^2 = 1$ ); or else the  $K$  parameters in  $\boldsymbol{\beta}_0$  and the scale parameter would not be separately identified.

To see intuitively why consider the case where  $\alpha_0 = 0$  and  $F_{\mathbf{x}}(\cdot)$  features a scale parameter  $\sigma$ . Here, the two equations:

$$\begin{aligned}y_i^* &= \mathbf{x}_i^T \boldsymbol{\beta}_0 + \varepsilon_i \\ \sigma y_i^* &= \sigma \left( \mathbf{x}_i^T \boldsymbol{\beta}_0 + \varepsilon_i \right)\end{aligned}$$

are observationally equivalent:  $F_{\mathbf{x}}\left(\mathbf{x}_i^T \boldsymbol{\beta}_0\right) = F_{\mathbf{x}}\left(\sigma \cdot \mathbf{x}_i^T \boldsymbol{\beta}_0\right)$ .

Intuitively, one can only observe whether  $Y_i^*$  lies above ( $Y_i = 1$ ) or below ( $Y_i = 0$ ) the implied threshold, and not its variation as a function of the variation of  $\mathbf{x}_i$ . For a similar reason  $\sigma$  *could* be identified *if* both  $\alpha_0 = 0$  *and*  $\mathbf{x}_i$  do not include a constant term.

## MLE for Binary Outcome Models (5/8)

In light of all these considerations it is easier to characterize the MLE problem for binary outcome models.

Given a symmetric probability distribution  $F_{\mathbf{x}}(\cdot)$  and a sample  $\{(y_i, \mathbf{x}_i)\}_{i=1}^N$ , the likelihood function is expressed as:

$$\mathcal{L}(\boldsymbol{\beta} | \{(y_i, \mathbf{x}_i)\}_{i=1}^N) = \prod_{i=1}^N [F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta})]^{y_i} [1 - F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta})]^{1-y_i}$$

which generalizes the likelihood function for a Bernoulli sample.

The corresponding log-likelihood function is as follows.

$$\begin{aligned} \log \mathcal{L}(\boldsymbol{\beta} | \{(y_i, \mathbf{x}_i)\}_{i=1}^N) &= \\ &= \sum_{i=1}^N \left\{ y_i \log F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}) + (1 - y_i) \log [1 - F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta})] \right\} \end{aligned}$$

## MLE for Binary Outcome Models (6/8)

The First Order Conditions are:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \log \mathcal{L} \left( \hat{\boldsymbol{\beta}}_{MLE} \mid \{(y_i, \mathbf{x}_i)\}_{i=1}^N \right) &= \sum_{i=1}^N \mathbf{s}_i \left( y_i, \mathbf{x}_i; \hat{\boldsymbol{\beta}}_{MLE} \right) = \\ &= \sum_{i=1}^N \left[ \frac{y_i f_{\mathbf{x}} \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right)}{F_{\mathbf{x}} \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right)} - \frac{(1 - y_i) f_{\mathbf{x}} \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right)}{1 - F_{\mathbf{x}} \left( \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{MLE} \right)} \right] \mathbf{x}_i = \mathbf{0} \end{aligned}$$

where  $f_{\mathbf{x}} \left( \mathbf{x}_i^T \boldsymbol{\beta} \right)$  is the probability **density function** associated with the – implicitly continuous – distribution  $F_{\mathbf{x}} \left( \mathbf{x}_i^T \boldsymbol{\beta} \right)$ .

- Like in the Poisson model, there is no closed form solution and the estimator must be calculated numerically.
- The variance-covariance is more conveniently estimated via the outer product of the gradients.

## MLE for Binary Outcome Models (7/8)

The most common choices for  $F_{\mathbf{x}}(\cdot)$  are:

- the **probit** model, where  $F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0) = \Phi(\mathbf{x}_i^T \boldsymbol{\beta}_0)$  and  $\Phi(\cdot)$  is a **cumulative standard normal distribution**:

$$\mathbb{P}(Y_i = 1 | \mathbf{x}_i) = \Phi(\mathbf{x}_i^T \boldsymbol{\beta}_0) = \int_{-\infty}^{\mathbf{x}_i^T \boldsymbol{\beta}_0} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt$$

- the **logit** model, where  $F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}_0) = \Lambda(\mathbf{x}_i^T \boldsymbol{\beta}_0)$  and  $\Lambda(\cdot)$  is a **scale-normalized cumulative logistic distribution**.

$$\mathbb{P}(Y_i = 1 | \mathbf{x}_i) = \Lambda(\mathbf{x}_i^T \boldsymbol{\beta}_0) = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta}_0)}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta}_0)}$$

The two models typically yield quite similar results; of the two, the logit is often more convenient to manipulate and compute.

## MLE for Binary Outcome Models (8/8)

- How to interpret the estimates  $\hat{\boldsymbol{\beta}}_{MLE}$  in a binary outcome model? **Not** in terms of causal effects, or probabilities.
- They must be interpreted in terms of **marginal effects**:

$$\frac{\partial \mathbb{P}(Y_i = 1 | \mathbf{x}_i)}{\partial \mathbf{x}_i} = \frac{\partial \mathbb{E}[Y_i | \mathbf{x}_i]}{\partial \mathbf{x}_i} = \frac{\partial F_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta})}{\partial \mathbf{x}_i} = f_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}) \boldsymbol{\beta}$$

which, however, vary as a function of the realizations  $\mathbf{x}_i$ .

- There are two **asymptotically equivalent** ways to derive marginal effects that aid the interpretation of the estimates.
- In one approach one evaluates  $f_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}) \boldsymbol{\beta}$  at  $\hat{\boldsymbol{\beta}}_{MLE}$  and at  $\mathbf{x} = \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ , the average values of  $\mathbf{x}_i$  in the data.
- In the other, one evaluates  $f_{\mathbf{x}}(\mathbf{x}_i^T \boldsymbol{\beta}) \boldsymbol{\beta}$  at  $\hat{\boldsymbol{\beta}}_{MLE}$  and at  $\mathbf{x}_i$  for  $i = 1, \dots, N$ ; the resulting quantities are averaged out.

# Introduction to Simulated Maximum Estimation

- Sometimes, the **numerical** evaluation of the M-Estimation criterion function  $q(\mathbf{x}_i; \boldsymbol{\theta})$  is so **complicated** that practical applications of the theory discussed thus far is unfeasible.
- In such cases, the typical solutions are estimators based on **simulation methods** that help circumvent the problem.
- The leading techniques based on simulations lie in the MLE domain, and are especially applied in some of LDV models.
- In econometrics this applies *especially but not exclusively* to models of Industrial Organization, and related ones.
- For convenience, the overview that follows next starts from simulation-based estimators that extend MLE.
- It then moves to general simulation-based M-Estimators.



# Unsolvable criterion functions

- Suppose that the probability mass or density function of all observable variables  $\mathbf{x}_i$  is expressed given the parameters  $\boldsymbol{\theta}$  as follows.

$$f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta}) = \int_{\mathbb{U}} f_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_i; \boldsymbol{\theta}) dH_{\mathbf{u}}(\mathbf{u}_i)$$

- Here,  $\mathbf{u}_i$  is a **random vector** with cumulative distribution  $H_{\mathbf{u}}(\mathbf{u}_i)$ . In the above expression,  $\mathbf{u}_i$  is integrated out over its support  $\mathbb{U}$ .
- Consider the case where there is **no closed form solution** for the above integral, even if by itself the conditional mass or density function  $f_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_i; \boldsymbol{\theta})$  is tractable.
- It is obvious that the MLE problem associated to  $f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$  **cannot be solved**, at least not easily.

## Example: Random coefficients logit (1/2)

- Consider the following **bivariate** logit model:

$$\mathbb{P}[Y_i = 1 | X_i] = \Lambda(\beta_0 + \beta_{1i}X_i)$$

a LDV with one *binary* dependent variable  $Y_i \in \{0, 1\}$  and one *possibly* continuous independent variable  $X_i$ .

- Note that the parameter  $\beta_{1i}$  is **observation specific**.
- As there are only  $N$  observations, one cannot estimate **all** the  $N$  parameters  $\beta_{1i}$  for  $i = 1, \dots, N$  and  $\beta_0$ .
- Yet there are actual settings where it is important to assess **heterogeneity** in the individual response of  $Y_i$  to  $X_i$ .
- These settings are framed as **random coefficients** models – as the logit model expressed above; they are more difficult to handle in LDV cases than in, say, linear environments.

## Example: Random coefficients logit (2/2)

- One could make some **distributional assumptions** about the individual coefficients, such as normality.

$$\beta_{1i} \sim \mathcal{N}(\beta_1, \sigma^2)$$

- The parameter set of interest is thus  $\theta = (\beta_0, \beta_1, \sigma^2)$ .
- Define  $u_i \equiv (\beta_{1i} - \beta_1) / \sigma$  and recall that  $\phi(\cdot)$  is the density function of the standard normal distribution.
- Under these hypotheses the conditional mass function of  $Y_i$  is expressed as an integral without closed form solution.

$$\begin{aligned} f_{Y_i|X_i}(y_i|x_i; \beta_0, \beta_1, \sigma^2) &= \\ &= \int_{\mathbb{R}} \Lambda[\beta_0 + (\beta_1 + \sigma u_i)x_i]^{y_i} \cdot \\ &\quad \cdot \{1 - \Lambda[\beta_0 + (\beta_1 + \sigma u_i)x_i]\}^{1-y_i} \phi(u_i) du_i \end{aligned}$$

# Simulation-based likelihood evaluations

- How to solve this problem and similar ones? A *brute force* approach to the integral is generally impractical.
- Common approaches are based upon the **simulation** of the likelihood components  $f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$ .
- The **Direct Monte Carlo Sampling** method is based on a sample  $\{\mathbf{u}_s\}_{s=1}^S$  of  $S$  random draws of  $\mathbf{u}_i$  from  $H_{\mathbf{u}}(\mathbf{u}_i)$  – these are used to calculate a *Monte Carlo estimate*:

$$\hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \tilde{f}_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_s; \boldsymbol{\theta})$$

where function  $\tilde{f}_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_s; \boldsymbol{\theta})$  is called a **subsimulator**.

- The subsimulator is typically an unbiased predictor:

$$\mathbb{E} \left[ \tilde{f}_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_s; \boldsymbol{\theta}) \right] = f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$$

where the expectation is taken over the support of  $\mathbf{u}_i$ .

# Maximum Simulated Likelihood

- With an unbiased subsimulator, by standard asymptotics:

$$\hat{f}_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta}) \xrightarrow{P} f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$$

that is, the simulator is **consistent** insofar as  $S \rightarrow \infty$ .

- Therefore the **Maximum Simulated Likelihood (MSL)** can be easily motivated and constructed as follows.

$$\hat{\boldsymbol{\theta}}_{MSL} = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N \log \hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta})$$

- MSL simulators need to be differentiable with respect to  $\boldsymbol{\theta}$ .  
**Example:** this is verifiable in the random coefficient logit.

$$\begin{aligned} \tilde{f}_{Y_i, X_i | U_s}(y_i, x_i | u_s; \beta_0, \beta_1, \sigma^2) &= \\ &= \Lambda[\beta_0 + (\beta_1 + \sigma u_s) x_i]^{y_i} \{1 - \Lambda[\beta_0 + (\beta_1 + \sigma u_s) x_i]\}^{1-y_i} \end{aligned}$$

- Elements  $\hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta})$  are conveniently calculated using the same draw  $\{\mathbf{u}_s\}_{s=1}^S$  for  $i = 1, \dots, N$ .

# Maximum Simulated Likelihood: Asymptotics

## Theorem 4

### Asymptotic Efficiency of Maximum Simulated Likelihood.

*Suppose that the mass or density function  $f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$  that describes the model's data generation process meets the requirements of Theorem 18 of Lecture 6, hence the corresponding "theoretical" MLE has a limiting distribution as per the statement of that Theorem. A SML estimator based on an unbiased subsimulator is asymptotically equivalent to the "theoretical" MLE and it has the same limiting distribution:*

$$\sqrt{N} \left( \hat{\boldsymbol{\theta}}_{SML} - \boldsymbol{\theta}_0 \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, [\mathbf{I}(\boldsymbol{\theta}_0)]^{-1} \right)$$

*if  $S, N \rightarrow \infty$  (which is sufficient for consistency) and  $\sqrt{N}/S \rightarrow 0$ .*

### Proof.

*(Outline.)* Gouriéroux and Monfort (1991) develop the standard Taylor expansion of the First Order Conditions and elaborate how it depends on two sources of noise: the one coming from the data  $\{\mathbf{x}_i\}_{i=1}^N$  and the one due to the simulation draws  $\{\mathbf{u}_s\}_{s=1}^S$ . The latter vanishes asymptotically if  $S$  grows at a rate higher than that of  $N$ .  $\square$

# Asymptotic bias-corrected MSL

- This result, while powerful, requires  $S$  to be very large, or else MSL is inconsistent: the problem arises as the log of the subsimulator is *not* unbiased, even if the level is.
- A second-order Taylor expansion of  $\log \hat{f}_{\mathbf{x}|S}(\mathbf{x}_i | \boldsymbol{\theta})$  around  $\log f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})$  suggested by Gouriéroux and Monfort (1991):

$$\log f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta}) \simeq \mathbb{E} \left[ \log \hat{f}_{\mathbf{x}|S}(\mathbf{x}_i | \boldsymbol{\theta}) \right] + \frac{\text{Var} \left[ \left( \hat{f}_{\mathbf{x}|S}(\mathbf{x}_i | \boldsymbol{\theta}) - f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta}) \right)^2 \right]}{2 [f_{\mathbf{x}}(\mathbf{x}_i | \boldsymbol{\theta})]^2}$$

- ...motivates the **asymptotic bias-corrected MSL**.

$$\hat{\boldsymbol{\theta}}_{BCMSL} = \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^N \left[ \log \hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta}) + \sum_{s=1}^S \frac{\left[ \tilde{f}_{\mathbf{x}|u}(\mathbf{x}_i | \mathbf{u}_s; \boldsymbol{\theta}) - \hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta}) \right]^2}{2S \left[ \hat{f}_{\mathbf{x},S}(\mathbf{x}_i | \boldsymbol{\theta}) \right]^2} \right]$$

# Simulated Maximum Estimation

- The analysis further extends to all M-Estimators where the criterion  $q(\mathbf{x}_i; \boldsymbol{\theta})$  cannot be expressed in closed form.
- A **Simulated M-Estimator** (SM), of which MSL is but a special case, is defined as:

$$\hat{\boldsymbol{\theta}}_{SM} = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N \hat{q}_S(\mathbf{x}_i; \boldsymbol{\theta})$$

and  $\hat{q}_S(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \tilde{q}_s(\mathbf{x}_i, \mathbf{u}_s; \boldsymbol{\theta})$  is typically an average of subsimulators that are written as  $\tilde{q}_s(\mathbf{x}_i, \mathbf{u}_s; \boldsymbol{\theta})$ ; these are based upon pseudo-random draws of  $\mathbf{u}_i$  like in SML case.

- The SM estimator is **consistent** if  $S, N \rightarrow \infty$ ; furthermore it is as efficient as the paired non-simulated M-Estimator if  $\sqrt{N}/S \rightarrow 0$ ; if  $S$  is too small, like in SML a bias correction may be necessary.



# Inference for Simulated Maximum Estimators

- In SM, a consistent estimator for  $\Upsilon_0$  – if the observations are independent (extensions to CCE and HAC exist) – is:

$$\hat{\Upsilon}_{M,S} = \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{s}}_{Si}(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_M) \hat{\mathbf{s}}_{Si}^T(\mathbf{x}_i, \hat{\boldsymbol{\theta}}_M) \xrightarrow{p} \Upsilon_0$$

where  $\hat{\mathbf{s}}_{Si}(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{\partial \hat{q}_S(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ .

- That for  $\mathbf{Q}_0$  is based on  $\hat{\mathbf{H}}_{Si}(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{\partial \hat{\mathbf{s}}_{Si}(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} = \frac{\partial^2 \hat{q}_S(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}$ .

$$\hat{\mathbf{Q}}_N \equiv \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{H}}_{Si}(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_M) \xrightarrow{p} \mathbf{Q}_0$$

- With i.i.d. data, the MSL information matrix is estimated as  $\hat{\Upsilon}_{M,S} \xrightarrow{p} [\mathbf{I}(\boldsymbol{\theta}_0)]^{-1}$ , and the score specializes as follows.

$$\hat{\mathbf{s}}_{Si}(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{\sum_{s=1}^S \frac{\partial \tilde{f}_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_s; \hat{\boldsymbol{\theta}}_{SML})}{\partial \boldsymbol{\theta}}}{\sum_{s=1}^S \tilde{f}_{\mathbf{x}|\mathbf{u}}(\mathbf{x}_i | \mathbf{u}_s; \hat{\boldsymbol{\theta}}_{SML})}$$

# Applications of Maximum Estimation: overview

This Lecture is concluded by showcasing three applications of M-Estimation in actual econometric problems.

The following applications are covered.

1. The Constant Elasticity of Substitution (CES) production functions, which naturally call for the use of NLLS.
2. Sample selection models: peculiar binary outcome models with a clear structural, economic interpretation.
3. The Bresnahan model for detecting collusion in a market: a classical application of MLE in Industrial Organization.

## CES production functions and NLLS (1/2)

- **Constant Elasticity of Substitution** (CES) production functions are key ingredients of many economic models.
- In its simplest form a CES production function writes as:

$$Y_i = [\alpha_K K_i + \alpha_L L_i]^{\frac{1}{\rho}} + \varepsilon_i$$

where  $Y_i$ ,  $K_i$ ,  $L_i$  are output, capital, labor; while  $\alpha_K > 0$  and  $\alpha_L > 0$  are the **saliency** parameters that determine the relative importance of each input. Furthermore,  $\varepsilon_i$  is an econometric error term.

- Instead,  $\rho > 0$  is a parameter related to the **elasticity of substitution** between inputs, which is clearly a constant and writes as  $\sigma = (1 + \rho)^{-1} \in (0, 1)$ .
- The CES production function becomes a Cobb-Douglas as in Lecture 7, with  $\alpha_K = \beta_K$  and  $\alpha_L = \beta_L$ , for  $\rho \rightarrow 0$ .

## CES production functions and NLLS (2/2)

- Clearly, this model must be estimated via NLLS. No closed form solution exists, and numerical methods must be used.
- A typical estimation algorithm splits the problem as:

$$(\hat{\rho}, \hat{\alpha}_K, \hat{\alpha}_L)_{NLLS} \in \arg \min_{\rho \in \mathbb{R}_{++}} \left[ \arg \min_{(\alpha_K, \alpha_L) \in \mathbb{R}_{++}^2} \sum_{i=1}^N \left( y_i - [\alpha_K k_i + \alpha_L l_i]^{\frac{1}{\rho}} \right)^2 \right]$$

where  $(y_i, k_i, l_i)$  denote observations of  $(Y_i, K_i, L_i)$ .

- Hence, numerical algorithms feature an *inner maximizer* of  $(\alpha_K, \alpha_L)$  given a value of  $\rho$ , and an *outer maximizer* for  $\rho$ .
- Unfortunately, the solution is typically unstable: especially with more inputs. Hence, practitioners usually prefer linear approximations or alternative production functions.

# The Heckit Sample Selection Model (1/4)

- Labor economists are interested in **labor supply** choices, especially those by women (who traditionally participated less than men in the labor market).
- Consider a labor supply **intensity** equation such as:

$$h_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$$

where  $h_i$  is the amount of time spent by woman  $i$  in wage labor, and  $\mathbf{x}_i$  are her characteristics.

- However,  $h_i$  is only observed for women who *do work*:

$$z_i^* = \mathbf{w}_i^T \boldsymbol{\gamma} + v_i$$
$$h_i \begin{cases} > 0 & \text{if } z_i^* > 0 \\ = 0 & \text{if } z_i^* \leq 0 \end{cases}$$

as governed by some latent variable  $z_i^*$ , a function of some set of characteristics  $\mathbf{w}_i$ , with a structural interpretation.

## The Heckit Sample Selection Model (2/4)

- Although the binary outcome model for **participation** can be estimated (say via MLE) interest falls on the **intensity** equation for  $h_i$ , or other labor market outcomes like wages.
- OLS cannot estimate the intensity equation consistently: if  $H_i$  is the random variable whence  $h_i$  is drawn, it is:

$$\begin{aligned}\mathbb{E}[H_i | \mathbf{x}_i, h_i > 0] &= \mathbb{E}[H_i | \mathbf{x}_i, z_i^* > 0] \\ &= \mathbf{x}_i^T \boldsymbol{\beta} + \mathbb{E}[\varepsilon_i | \mathbf{x}_i, z_i^* > 0] \\ &= \mathbf{x}_i^T \boldsymbol{\beta} + \mathbb{E}[\varepsilon_i | \mathbf{x}_i, v_i > -\mathbf{w}_i^T \boldsymbol{\gamma}]\end{aligned}$$

and  $\lambda(\mathbf{w}_i) \equiv \mathbb{E}[\varepsilon_i | v_i > -\mathbf{w}_i^T \boldsymbol{\gamma}] \neq 0$  if the two error terms  $\varepsilon_i$  and  $v_i$  are correlated (which is likely).

- The **omitted variable**  $\lambda(\mathbf{w}_i)$  thus leads to inconsistency. For example, wealthy women may be **both** inclined not to participate, and to participate with low  $h_i$  if they do.

## The Heckit Sample Selection Model (3/4)

- Heckman (1979) devised a solution to this problem, which awarded him the Nobel prize in economics.
- His “**heckit**” model posits a parametric assumption about the two error terms  $(\varepsilon_i, v_i)$ , like the bivariate normal.

$$\begin{pmatrix} \varepsilon_i \\ v_i \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

where  $\rho$  is the correlation coefficient between  $\varepsilon_i$  and  $v_i$ ,  $\sigma^2$  is the variance of  $\varepsilon_i$ , while the variance of  $v_i$  is *normalized* to 1 like in probit and logit.

- The resulting model could be estimated jointly by MLE.
- Yet the full-fledged likelihood function may be complicated to handle. An alternative two-step procedure, which is also consistent, is more popular.

# The Heckit Sample Selection Model (4/4)

The two-step procedure is as follows.

1. Run a probit on the participation equation; obtain  $\hat{\boldsymbol{\gamma}}_{MLE}$ .
2. For each observation, calculate the **inverse Mills ratio**:

$$\lambda_i = \left[ \frac{\phi(\mathbf{w}_i^T \hat{\boldsymbol{\gamma}}_{MLE})}{\Phi(\mathbf{w}_i^T \hat{\boldsymbol{\gamma}}_{MLE})} \right]$$

where  $\phi(\cdot)$  and  $\Phi(\cdot)$  are, respectively, the p.d.f. and c.d.f. of the standard normal distribution.

3. Run OLS on a *modified* intensity equation:

$$h_i = \mathbf{x}_i^T \boldsymbol{\beta} + \rho \lambda_i + \varepsilon_i$$

where  $\rho$  is now also the OLS coefficient for  $\lambda_i$ .

While consistent, this procedure delivers higher standard errors.



# Detecting collusion in oligopolies (1/6)

- Industrial Organization is currently the “structural” field of economics *par excellence*.
- Decades ago however it saw little empirical work. Classical questions, like the one regarding the sudden 45% increase in the US automobile production and sales in 1955 (that some attributed to a temporary price war), were left unanswered.
- Paul Samuelson himself had allegedly once said that he... *would flunk any econometrics paper that claimed to provide an explanation of 1955 auto sales.*
- Yet Bresnahan (1987) defied Samuelson as he developed an original model estimated via MLE that can be used to *test for collusion*. The results suggest that a price war occurred.
- Although this model is now obsolete, it is still a classic and has great instruction value.

## Detecting collusion in oligopolies (2/6)

- Consider  $N$  types of cars each with **quality**  $X_i = X(z_i, \beta)$  being a function of one car's characteristics  $z_i$  (given some parameters  $\beta$ ). Qualities can be ordered from best to worst – without loss of generality,  $X_i > X_h$  if  $i > h$ .
- Also consider some well-microfounded **demand functions** of each car for each year  $t = 1, \dots, T$ :

$$Q_{it}^D = D(P_{ht}, P_{it}, P_{jt}, X_{ht}, X_{it}, X_{jt}, \gamma)$$

where  $Q_{it}$  is the **quantity** of product  $i$ ,  $P_{it}$  its **price**,  $h, i, j$  are three **consecutive** products in the **order** of “qualities” and  $\gamma$  are some parameters.

- This specification makes prices and quantities dependent in equilibrium only on those of the “neighbors” of one product in the product space – it results from a unique specification of consumers' utility as given by Bresnahan.

## Detecting collusion in oligopolies (3/6)

Supply is standard: the profits from the sale of product  $i$  are:

$$\pi_{it} = P_{it}Q_{it} - c(X_{it})Q_{it}$$

with  $c(X_{it}) = \mu \exp(X_{it})$ . Bresnahan analyzes two scenarios.

1. **Competition:** each firm sets its own price  $P_{it}$  taking the price of neighbors  $h$  and  $j$  as given, with FOCs:

$$\frac{\partial \pi_{it}}{\partial P_{it}} = Q_{it} + (P_{it} - c(X_{it})) \frac{\partial Q_{it}(\cdot)}{\partial P_{it}} = 0$$

where  $Q_{it}$  is a function of  $P_{it}$  as in the demand function.

2. **Cooperation:** the firm(s) selling two products  $i$  and  $j$  set prices  $P_{it}$  and  $P_{jt}$  so as to maximize the joint profits, with FOCs for  $P_{it}$  price as follows (and symmetrically for  $P_{jt}$ ).

$$\begin{aligned} \frac{\partial [\pi_{it} + \pi_{jt}]}{\partial P_{it}} &= Q_{it} + (P_{it} - c(X_{it})) \frac{\partial Q_{it}(\cdot)}{\partial P_{it}} \\ &\quad + (P_{jt} - c(X_{jt})) \frac{\partial Q_{jt}(\cdot)}{\partial P_{it}} = 0 \end{aligned}$$

## Detecting collusion in oligopolies (4/6)

- Bresnahan then defines several **matrices**  $\mathbf{H}_t$  such that, in each year,

$$h_{(ij)t} = \begin{cases} 1 & \text{cooperation between products } i \text{ and } j \\ 0 & \text{competition between products } i \text{ and } j \end{cases}$$

and he characterizes several hypothetical scenarios for 1955 and surrounding years with corresponding matrices  $\mathbf{H}_t$ .

- Thus, for a **given choice** of matrix  $\mathbf{H}_t$  the supply function can be written as

$$q_{it}^S = S(P_{ht}, P_{it}, P_{jt}, X_{ht}, X_{it}, X_{jt}, \mathbf{H}_t, \boldsymbol{\gamma}, \mu)$$

where the demand parameters  $\boldsymbol{\gamma}$  enter via the derivative of the demand functions implied in the supply FOCs.

- This sets up alternative **counterfactuals** about collusion.

## Detecting collusion in oligopolies (5/6)

- By setting the equilibrium condition  $Q_{it}^D = Q_{it}^S = Q_{it}^*$  (and similarly for prices) for each product  $i = 1, \dots, N$  in every year  $t = 1, \dots, T$ , the **reduced form** is derived.

$$P_{it} = P^*(X_{ht}, X_{it}, X_{jt}, \mathbf{H}_t, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\mu})$$

$$Q_{it} = Q^*(X_{ht}, X_{it}, X_{jt}, \mathbf{H}_t, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\mu})$$

- This requires to develop the full-fledged solution given the assumptions about demand and supply, and a matrix  $\mathbf{H}_t$ .
- Introduce some error terms that make the model stochastic, and endow them of distributional assumptions as follows.

$$\begin{pmatrix} P_{it} - P^* \\ Q_{it} - Q^* \end{pmatrix} = \begin{pmatrix} \xi_{it}^P \\ \xi_{it}^Q \end{pmatrix} = \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_P^2 & 0 \\ 0 & \sigma_Q^2 \end{pmatrix} \right)$$

- All these error terms are implicit functions of the **reduced form** “predictions” of the model.

## Detecting collusion in oligopolies (6/6)

- For any **given** matrix  $\mathbf{H}_t$ , the model is estimated via MLE using the following likelihood function.

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\mu} | \mathbf{H}_t, \{p_{it}, q_{it}, \mathbf{z}_i\}_{i=1}^N) &= \\ &= \prod_{t=1}^T \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_P^2}} \exp\left(-\frac{(\xi_{it}^P)^2}{2\sigma_P^2}\right) \\ &\quad \times \prod_{t=1}^T \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_Q^2}} \exp\left(-\frac{(\xi_{it}^Q)^2}{2\sigma_Q^2}\right)\end{aligned}$$

- Tests about alternative matrices  $\mathbf{H}_{t0}$  and  $\mathbf{H}_{t1}$  (say, one for “competition” and one for “collusion”) are performed with the likelihood ratio method, as follows.

$$\begin{aligned}C_H &= 2 \left[ \log \mathcal{L}(\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}, \tilde{\boldsymbol{\mu}} | \mathbf{H}_{t1}, \mathbf{z}_1, \dots, \mathbf{z}_N) \right. \\ &\quad \left. - \log \mathcal{L}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\mu}} | \mathbf{H}_{t0}, \mathbf{z}_1, \dots, \mathbf{z}_N) \right]\end{aligned}$$