

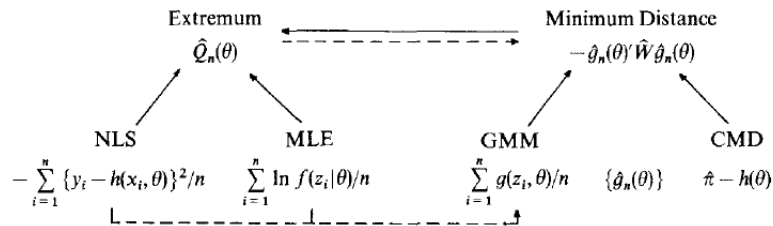
Estimation Principles

Raffaele Saggio

UBC

November 3, 2021

A Map of the World



How to estimate stuff?

- Common framework to analyzing many problems in econometrics summarized in Newey and McFadden (1994)
- Define quantity of interest θ_0 as maximizer of population criterion fn $Q_0(\theta)$
- Define estimator as maximizer of sample criterion fn $\hat{Q}_N(\theta)$
 - “Extremum” estimator (Amemiya, 1973)
 - θ of fixed dimension!
 - Need different tools when $\dim(\theta)$ grows with N

How to estimate stuff?

- Common framework to analyzing many problems in econometrics summarized in Newey and McFadden (1994)
- Define quantity of interest θ_0 as maximizer of population criterion fn $Q_0(\theta)$
- Define estimator as maximizer of sample criterion fn $\hat{Q}_N(\theta)$
 - “Extremum” estimator (Amemiya, 1973)
 - θ of fixed dimension!
 - Need different tools when $\dim(\theta)$ grows with N
- Basic results that hold subject to “usual” regularity conditions:
 - Consistency: $\hat{\theta} \xrightarrow{P} \theta_0$
 - Normality: $\sqrt{N}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \mathbf{H}^{-1} \mathbf{V} \mathbf{H}^{-1})$

Preliminaries

The data are an *iid* sample $\{\mathbf{Z}_i\}_{i=1}^N$ from some d.f. $F_Z(\cdot)$

Parameter of interest is:

$$\theta_0 = \arg \max_{\theta \in \Theta} Q_0(\theta)$$

where $Q_0(\theta)$ is some population criterion function.

Assume θ_0 is a singleton (point-identification)

Guess your estimator!

Example 1?:

$$Q_0(\theta) = E \left[(Y_i - X_i' \theta)^2 \right]$$

Guess your estimator!

Example I?:

$$Q_0(\theta) = E \left[(Y_i - X_i' \theta)^2 \right]$$

Example II?:

$$Q_0(\theta) = E \left[Y_i \ln \Phi(X_i' \theta) + (1 - Y_i) \ln (1 - \Phi(X_i' \theta)) \right]$$

Guess your estimator!

Example I?:

$$Q_0(\theta) = E \left[(Y_i - X_i' \theta)^2 \right]$$

Example II?:

$$Q_0(\theta) = E \left[Y_i \ln \Phi(X_i' \theta) + (1 - Y_i) \ln (1 - \Phi(X_i' \theta)) \right]$$

Example III?:

$$Q_0(\theta) = E \left[(Y_i - X_i' \theta) Z_i \right]' E \left[Z_i Z_i' \right]^{-1} E \left[(Y_i - X_i' \theta) Z_i \right]$$

Estimator

Estimate θ by max'ing sample criterion fn:

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

Estimator

Estimate θ by max'ing sample criterion fn:

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

Example I:

$$\hat{Q}_N(\theta) = \frac{1}{N} \sum_{i=1}^N (Y_i - X_i' \theta)^2$$

Example II:

$$\hat{Q}_N(\theta) = \frac{1}{N} \sum_{i=1}^N Y_i \ln \Phi(X_i' \theta) + (1 - Y_i) \ln (1 - \Phi(X_i' \theta))$$

Example III:

$$\hat{Q}_N(\theta) = \left[\frac{1}{N} \sum_{i=1}^N (Y_i - X_i' \theta) Z_i \right]' \left[\frac{1}{N} \sum_{i=1}^N Z_i Z_i' \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N (Y_i - X_i' \theta) Z_i \right]$$

Consistency

When does $\hat{\theta} \xrightarrow{P} \theta_0$?

Intuition:

- Need for θ_0 to be “well separated”
- Need $\hat{Q}_N(\theta)$ “close” to $Q_0(\theta)$ for maximizer to be close

Standard sufficient conditions:

- $Q_0(\cdot)$ continuous
- Θ compact
- $\sup_{\theta \in \Theta} \left| \hat{Q}_N(\theta) - Q_0(\theta) \right| \xrightarrow{P} 0$ (uniform convergence)

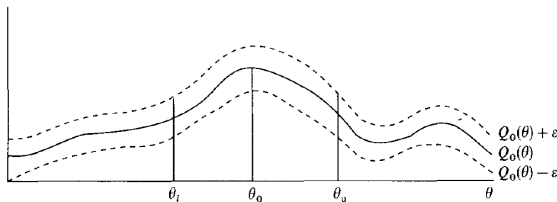
Uniform Convergence

Pointwise convergence usually easy to establish via LLN:

$$\left| \hat{Q}_N(\theta) - Q_0(\theta) \right| \xrightarrow{P} 0 \quad \forall \theta \in \Theta$$

Uniform convergence is stronger, need to trap “worst-case” deviation of fn from limit:

$$\sup_{\theta \in \Theta} \left| \hat{Q}_N(\theta) - Q_0(\theta) \right| \xrightarrow{P} 0$$



On finite grid, pointwise \rightarrow uniform convergence. What can go wrong w/ continuous Θ ?

Theorem

If the following conditions hold:

- i) Θ is compact*
- ii) $Q_0(\theta)$ is uniquely maximized at θ_0*
- iii) $Q_0(\theta)$ is continuous*
- iv) $\hat{Q}_N(\theta)$ converges uniformly to $Q_0(\theta)$*

Then,

$$\hat{\theta} \xrightarrow{P} \theta_0$$

Proof

Asymptotic Distribution

Theorem

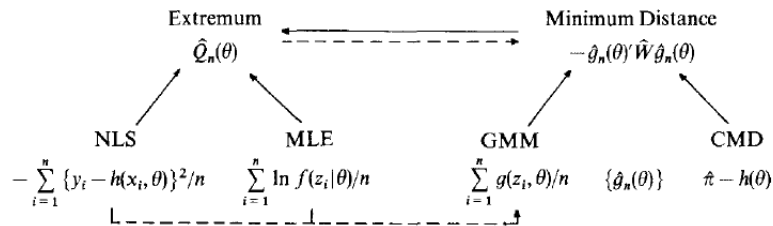
If $\hat{\theta} \xrightarrow{P} \theta_0$ and the following conditions hold:

- i) θ_0 is in the interior of Θ
- ii) $\hat{Q}_N(\theta)$ is twice differentiable in a neighborhood \mathcal{N} of θ_0
- iii) $\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \xrightarrow{d} N(0, \mathbf{V})$
- iv) there is an $\mathbf{H}(\theta)$ that is continuous at θ_0 such that $\sup_{\theta} \left\| \nabla_{\theta\theta} \hat{Q}_N(\theta) - \mathbf{H}(\theta) \right\| \xrightarrow{P} 0$
- v) $\mathbf{H} \equiv \mathbf{H}(\theta_0)$ is nonsingular

Then,

$$\sqrt{N}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \mathbf{H}^{-1} \mathbf{V} \mathbf{H}^{-1})$$

Special Cases



Classical Minimum Distance (CMD)

You have an economic model that says that certain observed moments in the data obeys a certain structure $\mathbf{g}(\theta_0)$.

$$\boldsymbol{\pi} - \mathbf{g}(\theta_0) = 0 \quad (1)$$

Classical Minimum Distance (CMD)

You have an economic model that says that certain observed moments in the data obeys a certain structure $\mathbf{g}(\theta_0)$.

$$\boldsymbol{\pi} - \mathbf{g}(\theta_0) = 0 \quad (1)$$

Objective fn (minimize instead of max):

$$\hat{Q}_N(\theta) = [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]' \hat{\mathbf{W}} [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]$$

- $\hat{\boldsymbol{\pi}} \xrightarrow{P} \boldsymbol{\pi}$ is a vector of reduced form sample moments
- $\mathbf{g}(\theta)$ is a structural function
- $\hat{\mathbf{W}} \xrightarrow{P} \mathbf{W}$ is a symmetric weighting matrix

Classical Minimum Distance (CMD)

You have an economic model that says that certain observed moments in the data obeys a certain structure $\mathbf{g}(\theta_0)$.

$$\boldsymbol{\pi} - \mathbf{g}(\theta_0) = 0 \quad (1)$$

Objective fn (minimize instead of max):

$$\hat{Q}_N(\theta) = [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]' \hat{\mathbf{W}} [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]$$

- $\hat{\boldsymbol{\pi}} \xrightarrow{P} \boldsymbol{\pi}$ is a vector of reduced form sample moments
- $\mathbf{g}(\theta)$ is a structural function
- $\hat{\mathbf{W}} \xrightarrow{P} \mathbf{W}$ is a symmetric weighting matrix

Classical Minimum Distance (CMD)

You have an economic model that says that certain observed moments in the data obeys a certain structure $\mathbf{g}(\theta_0)$.

$$\boldsymbol{\pi} - \mathbf{g}(\theta_0) = 0 \quad (1)$$

Objective fn (minimize instead of max):

$$\hat{Q}_N(\theta) = [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]' \hat{\mathbf{W}} [\hat{\boldsymbol{\pi}} - \mathbf{g}(\theta)]$$

- $\hat{\boldsymbol{\pi}} \xrightarrow{P} \boldsymbol{\pi}$ is a vector of reduced form sample moments
- $\mathbf{g}(\theta)$ is a structural function
- $\hat{\mathbf{W}} \xrightarrow{P} \mathbf{W}$ is a symmetric weighting matrix

Optimal weighting matrix: $\mathbf{W} = \text{Var}(\hat{\boldsymbol{\pi}})^{-1}$ (OMD)

Derivation of asymptotic variance

FOC:

$$\nabla_{\theta} \mathbf{g}(\hat{\theta})' \widehat{\mathbf{W}} [\hat{\pi} - \mathbf{g}(\hat{\theta})] = 0$$

Derivation of asymptotic variance

FOC:

$$\nabla_{\theta} \mathbf{g}(\hat{\theta})' \widehat{\mathbf{W}} [\hat{\pi} - \mathbf{g}(\hat{\theta})] = 0$$

Mean value expansion:

$$\mathbf{g}(\hat{\theta}) = \mathbf{g}(\theta_0) + \mathbf{G}(\bar{\theta})(\hat{\theta} - \theta_0)$$

where $\mathbf{G}(\theta) = \nabla_{\theta} \mathbf{g}(\theta)$. Substitute in to get:

$$\mathbf{G}(\hat{\theta})' \widehat{\mathbf{W}} [\hat{\pi} - \mathbf{g}(\theta_0) - \mathbf{G}(\bar{\theta})(\hat{\theta} - \theta_0)] = 0$$

Derivation of asymptotic variance

FOC:

$$\nabla_{\theta} \mathbf{g}(\hat{\theta})' \widehat{\mathbf{W}} [\hat{\pi} - \mathbf{g}(\hat{\theta})] = 0$$

Mean value expansion:

$$\mathbf{g}(\hat{\theta}) = \mathbf{g}(\theta_0) + \mathbf{G}(\bar{\theta})(\hat{\theta} - \theta_0)$$

where $\mathbf{G}(\theta) = \nabla_{\theta} \mathbf{g}(\theta)$. Substitute in to get:

$$\mathbf{G}(\hat{\theta})' \widehat{\mathbf{W}} [\hat{\pi} - \mathbf{g}(\theta_0) - \mathbf{G}(\bar{\theta})(\hat{\theta} - \theta_0)] = 0$$

Rearranging, normalizing, and taking limits yields:

$$\mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N} [\hat{\pi} - \mathbf{g}(\theta_0)] = \mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0) \sqrt{N} (\hat{\theta} - \theta_0) + o_p(1)$$

Standard Errors

Solve for $\sqrt{N}(\hat{\theta} - \theta_0)$ to get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = \left[\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0) \right]^{-1} \mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] + o_p(1)$$

By assumption:

$$\sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] \xrightarrow{d} N(0, \mathbf{V}_\pi)$$

Standard Errors

Solve for $\sqrt{N}(\hat{\theta} - \theta_0)$ to get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = \left[\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0) \right]^{-1} \mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] + o_p(1)$$

By assumption:

$$\sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] \xrightarrow{d} N(0, \mathbf{V}_\pi)$$

Slutsky: $\sqrt{N}(\hat{\theta} - \theta_0)$ is (asymptotically) a linear combination of normals w/ variance taking usual sandwich form:

$$\left[\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right]^{-1} \underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{V}_\pi \mathbf{W} \mathbf{G}(\theta_0)}_V \left[\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right]^{-1}$$

Standard Errors

Solve for $\sqrt{N}(\hat{\theta} - \theta_0)$ to get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = \left[\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0) \right]^{-1} \mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] + o_p(1)$$

By assumption:

$$\sqrt{N}[\hat{\pi} - \mathbf{g}(\theta_0)] \xrightarrow{d} N(0, \mathbf{V}_\pi)$$

Slutsky: $\sqrt{N}(\hat{\theta} - \theta_0)$ is (asymptotically) a linear combination of normals w/ variance taking usual sandwich form:

$$\left[\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right]^{-1} \underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{V}_\pi \mathbf{W} \mathbf{G}(\theta_0)}_V \left[\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right]^{-1}$$

Standard errors: replace unknown quantities with sample analogues (i.e. $\mathbf{G}(\hat{\theta})$ for $\mathbf{G}(\theta_0)$, $\hat{\mathbf{V}}_\pi$ for \mathbf{V}_π)

Optimal Weighting

With optimal weights $\mathbf{W} = \mathbf{V}_\pi^{-1}$, asymptotic variance reduces to:

$$\left[\mathbf{G}(\theta_0)' \mathbf{V}_\pi^{-1} \mathbf{G}(\theta_0) \right]^{-1}$$

Optimal Weighting

With optimal weights $\mathbf{W} = \mathbf{V}_\pi^{-1}$, asymptotic variance reduces to:

$$\left[\mathbf{G}(\theta_0)' \mathbf{V}_\pi^{-1} \mathbf{G}(\theta_0) \right]^{-1}$$

Specification testing: optimal weighting yields null distribution for minimized value of criterion function

$$\hat{Q}_N^{OMD}(\hat{\theta}) \sim \chi^2(J - K)$$

Answers question: does my model explain the data up to sampling error? (i.e. could population $R^2 = 1$?)

Practical Problems w/ OMD

Optimal weighting behaves poorly when \mathbf{W} is large relative to sample size (Altonji and Segal, 1996)

- Variant of “Kakwani bias” in FGLS
- Problem emerges from correlation between weights and moments

Practical Problems w/ OMD

Optimal weighting behaves poorly when \mathbf{W} is large relative to sample size (Altonji and Segal, 1996)

- Variant of “Kakwani bias” in FGLS
- Problem emerges from correlation between weights and moments

Potential solutions:

- Use inefficient (but tractable) \mathbf{W} : equal (or diagonal) weight matrix (e.g., Abowd and Card, 1992)
- Jackknife / split sample estimation of \mathbf{W} (Kezdi, Solon, and Hahn, 2002)
- Parameterize $\mathbf{W} = \mathbf{W}(\delta)$ and estimate δ along with θ via “CUGMM” (Hansen, Heaton, and Yaron, 1996; Hausman et al, 2011)
- Generalized Empirical Likelihood (Newey and Smith, 2004)

Spec test without optimal weighting

- Is it possible to test our model if we don't have (want) to use optimal weighting matrix?

- Newey (1985): yes! Consider estimation of θ that is based upon

$$[\hat{\pi} - g(\theta)]' \mathbf{A} [\hat{\pi} - g(\theta)] \quad (2)$$

where \mathbf{A} is a possibly stochastic matrix.

- Newey showed that

$$N[\hat{\pi} - g(\theta)]' \mathbf{R}^- [\hat{\pi} - g(\theta)] \sim \chi^2_{J-K} \quad (3)$$

where \mathbf{R}^- is the generalized inverse of the matrix $\mathbf{R} = \mathbf{M} \mathbf{V}_\pi \mathbf{M}'$ with

$$\mathbf{M} = \mathbf{I} - \mathbf{G}(\hat{\theta})(\mathbf{G}(\hat{\theta})' \mathbf{A} \mathbf{G}(\hat{\theta}))^{-1} \mathbf{G}(\hat{\theta})' \mathbf{A} \quad (4)$$

- Need generalized inverse because \mathbf{R} does not have full rank (presence of \mathbf{M}).

On the frontier

- Important assumption is “correct” specification.

$$\sqrt{N}(\hat{\pi} - g(\theta_0)) \sim^a \mathcal{N}(0, \mathbf{V}_\pi) \quad (5)$$

- Important assumption is “correct” specification.

$$\sqrt{N}(\hat{\pi} - g(\theta_0)) \sim^d \mathcal{N}(0, \mathbf{V}_\pi) \quad (5)$$

- What happens under misspecification?
- Chamberlain (1994): be careful with weights! Randomness of weights increases noise in resulting estimator.
- Active frontier in econometrics today: what if my model is wrong?
Can my standard errors account for that?
 - Armstrong and Kolesar (2018); Bonhomme and Weidner (2019).

Criterion function (min once again):

$$\hat{Q}_{GMM}(\theta) = \hat{\mathbf{g}}(\theta)' \hat{\mathbf{W}} \hat{\mathbf{g}}(\theta)$$

- $\hat{\mathbf{g}}(\theta) = \frac{1}{N} \sum_i \mathbf{f}(\mathbf{Z}_i, \theta)$ is a $J \times 1$ vector of moment conditions
- CMD nested by separable case where $\mathbf{f}(\mathbf{Z}_i, \theta) = \boldsymbol{\pi}(\mathbf{Z}_i) + \tilde{\mathbf{f}}(\theta)$
- MM nested by just-id case where $J = \dim(\theta)$

Identification

Population moment restrictions:

$$E [\mathbf{f} (\mathbf{Z}_i, \theta)] = 0 \text{ iff } \theta = \theta_0$$

Examples:

- Instrumental variables (orthogonality condition):

$$E [(Y_i - X_i' \beta) Z_i] = 0$$

- Rational expectations restrictions (e.g., Hall, 1978) (real interest rate = 0)

$$u'(C_t) = \beta E[u'(C_{t+1}) | \Omega_t]$$

Derivation of Asymptotic Variance

FOC:

$$\widehat{\mathbf{G}}(\widehat{\theta})' \widehat{\mathbf{W}} \widehat{\mathbf{g}}(\widehat{\theta}) = 0$$

Mean value expansion:

$$\widehat{\mathbf{g}}(\widehat{\theta}) = \widehat{\mathbf{g}}(\theta_0) + \widehat{\mathbf{G}}(\bar{\theta})(\widehat{\theta} - \theta_0)$$

Substitute in to get:

$$\widehat{\mathbf{G}}(\widehat{\theta})' \widehat{\mathbf{W}} \widehat{\mathbf{g}}(\theta_0) + \widehat{\mathbf{G}}(\widehat{\theta})' \widehat{\mathbf{W}} \widehat{\mathbf{G}}(\bar{\theta})(\widehat{\theta} - \theta_0) = 0$$

Therefore

$$\sqrt{N}(\widehat{\theta} - \theta_0) = (\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0))^{-1} \mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N} \widehat{\mathbf{g}}(\theta_0) + o_p(1)$$

Derivation of Asymptotic Variance

By assumption:

$$E \left[\sqrt{N} \hat{\mathbf{g}}(\theta_0) \right] = 0$$

Hence, for $\mathbf{V}_f = E \left[\mathbf{f}(\mathbf{Z}_i, \theta_0) \mathbf{f}(\mathbf{Z}_i, \theta_0)' \right]$, we have

$$AVAR \left(\mathbf{G}(\theta_0)' \mathbf{W} \sqrt{N} \hat{\mathbf{g}}(\theta_0) \right) = \mathbf{G}(\theta_0)' \mathbf{W} \mathbf{V}_f \mathbf{W} \mathbf{G}(\theta_0)$$

Consequently $AVAR \left(\sqrt{N} (\hat{\theta} - \theta_0) \right)$ is:

$$\left(\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right)^{-1} \underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{V}_f \mathbf{W} \mathbf{G}(\theta_0)}_V \left(\underbrace{\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)}_H \right)^{-1}$$

Plugin sample analogues to get standard errors

Optimal Weighting

When $\mathbf{W} = \mathbf{V}_f^{-1}$ (optimal weighting) AVAR simplifies to:

$$(\mathbf{G}(\theta_0)' \mathbf{V}_f^{-1} \mathbf{G}(\theta_0))^{-1}$$

2-step estimation approach (ala 3SLS):

- 1 get inefficient starting estimates $\hat{\theta}^{(1)}$ from choosing $\hat{\mathbf{W}} = \mathbf{I}$.
- 2 minimize $Q_{GMM}(\theta)$ using $\hat{\mathbf{W}} = \mathbf{W}(\hat{\theta}^{(1)})$

Keep going?: further steps provide no (first order) asymptotic advantage but often perform better in finite samples

CUGMM (Hansen, Heaton, and Yaron, 1996)

As w/ OMD, problems emerge when weight matrix has too many unknown parameters

One solution: “continuously update” weight matrix $\mathbf{W}(\theta)$ and minimize

$$\hat{Q}_{CUGMM}(\theta) = \hat{\mathbf{g}}(\theta)' \hat{\mathbf{W}}(\theta) \hat{\mathbf{g}}(\theta)$$

Better asymptotic performance than GMM but can be difficult to find minimum.

Empirical Likelihood

GMM a bit like a quadratic approx to a log likelihood

EL: maximize non-parametric likelihood fn subject to moment restrictions

$$\begin{aligned} & \max_{\{\pi_i\}, \theta} \prod_{i=1}^N \pi_i \\ \text{s.t.} \quad & \sum_{i=1}^N \pi_i \mathbf{f}(\mathbf{Z}_i, \theta) = 0 \\ & \sum_{i=1}^N \pi_i = 1 \end{aligned}$$

Difficult optimization problem:

- w/o constraints $\hat{\pi}_i = \frac{1}{N}$ (returns the EDF)
- w/ constraints $\{\hat{\pi}_i\}$ give efficient estimates of joint distribution of data and $\hat{\theta}$ higher order unbiased (Newey and Smith, 2004)
- Generally difficult to compute (extension to CUGMM).

Maximum Likelihood

You are willing to specify the *entire* conditional distribution of the outcome given covariates.

Maximum Likelihood

You are willing to specify the *entire* conditional distribution of the outcome given covariates. Objective fn:

$$\hat{Q}_{ML}(\theta) = \frac{1}{N} \sum_i l(\mathbf{z}_i, \theta)$$

FOC:

$$\frac{1}{N} \sum_i \mathbf{s}(\mathbf{z}_i, \hat{\theta}) = 0$$

- ML: if model is correct $\rightarrow E[\mathbf{s}(\mathbf{Z}_i, \theta_0)] = 0$
- No weighting problems!
- Likelihood chooses the right moments for you! (embedded in the score function)
 - “Efficient Method of Moments” on FOCs (Gallant and Tauchen, 1996)
- ML is prone to misspecification. Even getting a minor feature of the entire distribution wrong can make the ML estimator inconsistent.
- Can you think of a relevant counter-example?

Asymptotic variance

“Influence function” representation – 1st order effect of adding an obs on estimator:

$$\sqrt{N} \left(\hat{\theta} - \theta_0 \right) = \hat{\mathbf{H}}(\theta_0)^{-1} \frac{1}{\sqrt{N}} \sum_i \mathbf{s}(\mathbf{Z}_i, \theta_0) + o_p(1)$$

Variance of IF gives asymptotic variance:

$$\mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1}$$

where $\mathbf{V} = \mathbf{E} [\mathbf{s}(\mathbf{Z}_i, \theta_0) \mathbf{s}(\mathbf{Z}_i, \theta_0)']$

Asymptotic variance

“Influence function” representation – 1st order effect of adding an obs on estimator:

$$\sqrt{N}(\hat{\theta} - \theta_0) = \hat{\mathbf{H}}(\theta_0)^{-1} \frac{1}{\sqrt{N}} \sum_i \mathbf{s}(\mathbf{Z}_i, \theta_0) + o_p(1)$$

Variance of IF gives asymptotic variance:

$$\mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1}$$

where $\mathbf{V} = \mathbf{E} [\mathbf{s}(\mathbf{Z}_i, \theta_0) \mathbf{s}(\mathbf{Z}_i, \theta_0)']$

Information matrix equality:

$$\mathbf{V} = \mathbf{H}(\theta_0)$$

When this holds, asymptotic variance reduces to Cramer-Rao lower bound $\mathbf{H}(\theta_0)^{-1}$.

Nothing beats the CRLB → nothing beats a properly specified likelihood model.

Simulation Estimators

Often difficult to write down the likelihood

- Particularly difficult for multinomial choice problems where regions of integration can quickly become intractable
- Or for parametric models with lots of unobserved heterogeneity

Simulation Estimators

Often difficult to write down the likelihood

- Particularly difficult for multinomial choice problems where regions of integration can quickly become intractable
- Or for parametric models with lots of unobserved heterogeneity

Simulation methods

- Can use computer to compute likelihood via simulation (MSL)
- With continuous variables often easier to simulate moment conditions (MSM)
- Great (free) introduction by Train (2009) or Chapter 12 of Cameron and Trivedi.

Simulation Estimators

Often difficult to write down the likelihood

- Particularly difficult for multinomial choice problems where regions of integration can quickly become intractable
- Or for parametric models with lots of unobserved heterogeneity

Simulation methods

- Can use computer to compute likelihood via simulation (MSL)
- With continuous variables often easier to simulate moment conditions (MSM)
- Great (free) introduction by Train (2009) or Chapter 12 of Cameron and Trivedi.

Warning: can easily waste a lot of time on this stuff

- Tendency to use these methods when you don't fully understand the model / identification
- Important to “warm up” by estimating simpler models that you understand!

Maximum Simulated Likelihood

$$\hat{\theta}_{MSL} \equiv \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \ln \hat{L}_M(\mathbf{Z}_i, \theta)$$

- $\hat{L}_M(\mathbf{Z}_i, \theta) = \hat{L}(\mathbf{Z}_i, \theta, u_{i1}, u_{i2}, \dots, u_{iM})$ is simulated likelihood
- Fully parametric: drawing $\{u_{im}\}_{m=1}^M$ from known distribution
- Hold $\{\mathbf{u}_i\}_{i=1}^N$ fixed when searching over θ to avoid “chatter”
- Need large M to ensure consistency:
$$\lim_{M \rightarrow \infty} P(|\hat{L}_M(\mathbf{Z}_i, \theta) - L(\mathbf{Z}_i, \theta)| > \varepsilon) = 0$$

Example (Probit)

Probit likelihood:

$$L(Y_i, X_i, \beta) = \Phi(X_i' \beta)^{Y_i} [1 - \Phi(X_i' \beta)]^{1-Y_i}$$

Simulated analogue:

$$\begin{aligned} \hat{L}_M(Y_i, X_i, \beta) &= \left(\frac{1}{M} \sum_{m=1}^M 1[X_i' \beta - u_{im} > 0] \right)^{Y_i} \\ &\quad \times \left(\frac{1}{M} \sum_{m=1}^M 1[X_i' \beta - u_{im} < 0] \right)^{1-Y_i} \end{aligned}$$

where $\left\{ u_{im} \stackrel{iid}{\sim} N(0, 1) \right\}_{i=1, m=1}^{N, M}$

- Unbiased simulator: $E[\hat{L}_M(Y_i, X_i, \beta)] = L(Y_i, X_i, \beta)$
- But $\hat{L}_M(Y_i, X_i, \beta)$ nondifferentiable for finite M (step function)

Example (Random Coefficient Logit)

Suppose

$$Y_i = 1 [b_i X_i + \varepsilon_i > 0]$$

where $\varepsilon_i \sim \text{Logistic}(0, 1)$ and $b_i \sim N(\mu, \sigma^2)$

- Could simulate likelihood by taking draws from (ε_i, b_i) .
- Better to integrate out ε_i analytically to obtain smooth objective fn:

$$\begin{aligned} \hat{L}_M(Y_i, X_i, \mu, \sigma) &\equiv \left(\frac{1}{M} \sum_{m=1}^M \Lambda((\mu + \sigma u_{im}) X_i) \right)^{Y_i} \\ &\times \left(1 - \frac{1}{M} \sum_{m=1}^M \Lambda((\mu + \sigma u_{im}) X_i) \right)^{1-Y_i} \end{aligned}$$

Asymptotics

Pakes and Pollard (1989): general asymptotic framework for simulation estimators (“stochastic equicontinuity”)

Gourieroux and Monfort (1991): if $M, N \rightarrow \infty$ and $\sqrt{N}/M \rightarrow 0$, then $\hat{\theta}_{MSL} \xrightarrow{P} \theta_0$ and

$$\sqrt{N} \left(\hat{\theta}_{MSL} - \theta_0 \right) \xrightarrow{d} N \left(0, \mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1} \right)$$

Asymptotics

Pakes and Pollard (1989): general asymptotic framework for simulation estimators (“stochastic equicontinuity”)

Gourieroux and Monfort (1991): if $M, N \rightarrow \infty$ and $\sqrt{N}/M \rightarrow 0$, then $\hat{\theta}_{MSL} \xrightarrow{P} \theta_0$ and

$$\sqrt{N} \left(\hat{\theta}_{MSL} - \theta_0 \right) \xrightarrow{d} N \left(0, \mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1} \right)$$

Practical advice: start with $M = 100$ to get initial estimates then check stability as M increases

Asymptotics

Pakes and Pollard (1989): general asymptotic framework for simulation estimators (“stochastic equicontinuity”)

Gourieroux and Monfort (1991): if $M, N \rightarrow \infty$ and $\sqrt{N}/M \rightarrow 0$, then $\hat{\theta}_{MSL} \xrightarrow{P} \theta_0$ and

$$\sqrt{N} \left(\hat{\theta}_{MSL} - \theta_0 \right) \xrightarrow{d} N \left(0, \mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1} \right)$$

Practical advice: start with $M = 100$ to get initial estimates then check stability as M increases

Std errors: plug in simulation estimators of $\mathbf{H}(\hat{\theta}_{MSL})$, $\mathbf{V}(\hat{\theta}_{MSL})$

Method of Simulated Moments

$$\hat{\theta}_{MSM} \equiv \arg \min_{\theta} \hat{\mathbf{g}}_M(\theta)' \hat{\mathbf{W}} \hat{\mathbf{g}}_M(\theta)$$

- $\hat{\mathbf{g}}_M(\theta) = \frac{1}{N} \sum_i \hat{\mathbf{f}}_M(\mathbf{Z}_i, \theta)$ is simulated moment condition
- McFadden (1989): for fixed M , as $N \rightarrow \infty$, $\hat{\theta}_{MSM} \xrightarrow{p} \theta_0$ and $AVAR\left(\sqrt{N}\left(\hat{\theta}_{MSM} - \theta_0\right)\right) =$

$$\left(\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)\right)^{-1} \mathbf{G}(\theta_0)' \mathbf{W} \mathbf{V}_{f,M} \mathbf{W} \mathbf{G}(\theta_0) \left(\mathbf{G}(\theta_0)' \mathbf{W} \mathbf{G}(\theta_0)\right)^{-1}$$

Which moments to match?

If model is fully parametric, best to match the scores: “method of simulated scores” Hajivassiliou and McFadden (1998)

- Problem: typically hard to get unbiased simulator of scores
- Example (binary choice):

$$E \left[Y_i \frac{\frac{\partial}{\partial \theta} P(Y_i = 1 | X_i)}{P(Y_i = 1 | X_i)} - (1 - Y_i) \frac{\frac{\partial}{\partial \theta} P(Y_i = 1 | X_i)}{1 - P(Y_i = 1 | X_i)} \right] = 0$$

Which moments to match?

If model is fully parametric, best to match the scores: “method of simulated scores” Hajivassiliou and McFadden (1998)

- Problem: typically hard to get unbiased simulator of scores
- Example (binary choice):

$$E \left[Y_i \frac{\frac{\partial}{\partial \theta} P(Y_i = 1 | X_i)}{P(Y_i = 1 | X_i)} - (1 - Y_i) \frac{\frac{\partial}{\partial \theta} P(Y_i = 1 | X_i)}{1 - P(Y_i = 1 | X_i)} \right] = 0$$

- Can you see a problem with this?

MSL vs MSM

Rules of thumb:

- parametric discrete choice model \Rightarrow MSL
 - MSL chooses the “right” moments
 - Works well for moderate M
- semi-parametric model or model w/ mixed continuous-discrete choices \Rightarrow MSM
 - don't want to match moments you don't believe
 - too hard to simulate multidimensional density of continuous variables

MSL vs MSM

Rules of thumb:

- parametric discrete choice model \Rightarrow MSL
 - MSL chooses the “right” moments
 - Works well for moderate M
- semi-parametric model or model w/ mixed continuous-discrete choices \Rightarrow MSM
 - don't want to match moments you don't believe
 - too hard to simulate multidimensional density of continuous variables

Eisenhauer, Heckman, and Mosso (2015, IER):

- Empirical horserace of MSL vs MSM for modern dynamic discrete choice model
- Conclusion: $MSL > MSM$ (assuming model is correct)

Appendix

Proof

For any ε we have with probability approaching 1:

- a) $\hat{Q}_N(\hat{\theta}) > \hat{Q}_N(\theta_0) - \varepsilon$ (by virtue of maximizer)
- b) $Q_0(\hat{\theta}) > \hat{Q}_N(\hat{\theta}) - \varepsilon$ (by uniform convergence)
- c) $\hat{Q}_N(\theta_0) > Q_0(\theta_0) - \varepsilon$ (by uniform convergence)

Therefore:

$$Q_0(\hat{\theta}) \stackrel{b)}{>} \hat{Q}_N(\hat{\theta}) - \varepsilon \stackrel{a)}{>} \hat{Q}_N(\theta_0) - 2\varepsilon \stackrel{c)}{>} Q_0(\theta_0) - 3\varepsilon.$$

Proof

For any ε we have with probability approaching 1:

- a) $\hat{Q}_N(\hat{\theta}) > \hat{Q}_N(\theta_0) - \varepsilon$ (by virtue of maximizer)
- b) $Q_0(\hat{\theta}) > \hat{Q}_N(\hat{\theta}) - \varepsilon$ (by uniform convergence)
- c) $\hat{Q}_N(\theta_0) > Q_0(\theta_0) - \varepsilon$ (by uniform convergence)

Therefore:

$$Q_0(\hat{\theta}) \stackrel{b)}{>} \hat{Q}_N(\hat{\theta}) - \varepsilon \stackrel{a)}{>} \hat{Q}_N(\theta_0) - 2\varepsilon \stackrel{c)}{>} Q_0(\theta_0) - 3\varepsilon.$$

Let \mathcal{N} be an open neighborhood of θ_0 and \mathcal{N}^c its complement.

Note that $\sup_{\theta \in \Theta \cap \mathcal{N}^c} Q_0(\theta) = Q_0(\theta^*) < Q_0(\theta_0)$.

- Choose $3\varepsilon = Q_0(\theta_0) - \sup_{\theta \in \Theta \cap \mathcal{N}^c} Q_0(\theta)$.
- It follows that $Q_0(\hat{\theta}) > \sup_{\theta \in \Theta \cap \mathcal{N}^c} Q_0(\theta)$ w.p.a.1.

Hence $\hat{\theta} \in \mathcal{N}$ \square

Proof Sketch

Start w/ FOC:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = 0$$

- Mean value expansion:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = \nabla_{\theta} \hat{Q}_N(\theta_0) + \hat{\mathbf{H}}(\bar{\theta})(\hat{\theta} - \theta_0)$$

where $\bar{\theta} \in [\hat{\theta}, \theta_0]$ and $\hat{\mathbf{H}}(\theta) \equiv \nabla_{\theta\theta} \hat{Q}_N(\theta)$.

- Rearranging we get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\hat{\mathbf{H}}(\bar{\theta})^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right)$$

Proof Sketch

Start w/ FOC:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = 0$$

- Mean value expansion:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = \nabla_{\theta} \hat{Q}_N(\theta_0) + \hat{\mathbf{H}}(\bar{\theta})(\hat{\theta} - \theta_0)$$

where $\bar{\theta} \in [\hat{\theta}, \theta_0]$ and $\hat{\mathbf{H}}(\theta) \equiv \nabla_{\theta\theta} \hat{Q}_N(\theta)$.

- Rearranging we get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\hat{\mathbf{H}}(\bar{\theta})^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right)$$

- Want to show:

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\mathbf{H}(\theta_0)^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right) + o_p(1)$$

Proof Sketch

Start w/ FOC:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = 0$$

- Mean value expansion:

$$\nabla_{\theta} \hat{Q}_N(\hat{\theta}) = \nabla_{\theta} \hat{Q}_N(\theta_0) + \hat{\mathbf{H}}(\bar{\theta})(\hat{\theta} - \theta_0)$$

where $\bar{\theta} \in [\hat{\theta}, \theta_0]$ and $\hat{\mathbf{H}}(\theta) \equiv \nabla_{\theta\theta} \hat{Q}_N(\theta)$.

- Rearranging we get:

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\hat{\mathbf{H}}(\bar{\theta})^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right)$$

- Want to show:

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\mathbf{H}(\theta_0)^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right) + o_p(1)$$

- Need $\left\| \hat{\mathbf{H}}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\| \xrightarrow{p} 0$

Proof Sketch

Triangle inequality:

$$\begin{aligned}\left\| \hat{\mathbf{H}}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\| &\leq \underbrace{\left\| \hat{\mathbf{H}}(\bar{\theta}) - \mathbf{H}(\bar{\theta}) \right\|}_{\text{Noise in Hessian}} + \underbrace{\left\| \mathbf{H}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\|}_{\text{Noise in parameter}} \\ &\leq \sup_{\theta \in \Theta} \left\| \hat{\mathbf{H}}(\theta) - \mathbf{H}(\theta) \right\| + \left\| \mathbf{H}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\|\end{aligned}$$

Proof Sketch

Triangle inequality:

$$\begin{aligned}\left\| \hat{\mathbf{H}}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\| &\leq \underbrace{\left\| \hat{\mathbf{H}}(\bar{\theta}) - \mathbf{H}(\bar{\theta}) \right\|}_{\text{Noise in Hessian}} + \underbrace{\left\| \mathbf{H}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\|}_{\text{Noise in parameter}} \\ &\leq \sup_{\theta \in \Theta} \left\| \hat{\mathbf{H}}(\theta) - \mathbf{H}(\theta) \right\| + \left\| \mathbf{H}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\|\end{aligned}$$

- By assumption $\sup_{\theta \in \Theta} \left\| \hat{\mathbf{H}}(\theta) - \mathbf{H}(\theta) \right\| \xrightarrow{P} 0$ (kills 1st term)
- Because $\hat{\theta} \xrightarrow{P} \theta_0$, it must also be that $\bar{\theta} \xrightarrow{P} \theta_0$ (squeeze theorem).
- Since $\mathbf{H}(\cdot)$ is continuous, it follows from the continuous mapping theorem that $\left\| \mathbf{H}(\bar{\theta}) - \mathbf{H}(\theta_0) \right\| \xrightarrow{P} 0$ (kills 2nd term)
- Therefore $\hat{\mathbf{H}}(\bar{\theta}) \xrightarrow{P} \mathbf{H}(\theta_0)$

Recall

$$\sqrt{N} \left(\hat{\theta} - \theta_0 \right) = -\hat{\mathbf{H}} \left(\bar{\theta} \right)^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N \left(\theta_0 \right) \right)$$

Recall

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\hat{\mathbf{H}}(\bar{\theta})^{-1} \left(\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \right)$$

- We established $\hat{\mathbf{H}}(\bar{\theta}) \xrightarrow{P} \mathbf{H}(\theta_0)$
- By continuous mapping theorem, $\hat{\mathbf{H}}(\bar{\theta})^{-1} \xrightarrow{P} \mathbf{H}(\theta_0)^{-1}$
- And by assumption: $\sqrt{N} \nabla_{\theta} \hat{Q}_N(\theta_0) \xrightarrow{d} N(0, \mathbf{V})$
- Therefore, by Slutsky,
$$\sqrt{N}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, \mathbf{H}(\theta_0)^{-1} \mathbf{V} \mathbf{H}(\theta_0)^{-1}\right) \quad \square$$

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

- $\Theta = [0, 1]$ is compact

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

- $\Theta = [0, 1]$ is compact
- $Q_N(\cdot)$ is continuous

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

- $\Theta = [0, 1]$ is compact
- $Q_N(\cdot)$ is continuous
- $\lim_{N \rightarrow \infty} Q_N(\theta) = 0 \quad \forall \theta \in [0, 1]$ (pointwise convergence)

Example: “folding tent”

$$Q_N(\theta) = \begin{cases} \max\{4N\theta - 2, 0\} & \text{if } \theta \leq \frac{3}{4N} \\ \max\{-4N\theta + 4, 0\} & \text{if } \theta > \frac{3}{4N} \end{cases}$$

- $\Theta = [0, 1]$ is compact
- $Q_N(\cdot)$ is continuous
- $\lim_{N \rightarrow \infty} Q_N(\theta) = 0 \quad \forall \theta \in [0, 1]$ (pointwise convergence)
- But $\max_{\theta \in \Theta} |Q_N(\theta) - 0| = 1 \quad \forall N \in \mathbb{N}!$