

T. Christensen

1 Structural Estimation

Over the course of the next five lectures we shall discuss estimation and inference methods for “structural” econometric models. Structural models attempt to formally model the behavior of economic agents. Once we estimate deep policy-invariant parameters that govern agents’ behavior, we may then perform policy experiments (or *counterfactuals*) to predict how agents’ behavior and economic outcomes would change if we changed policy settings or other features of the economic environment. This ex-ante approach to policy evaluation is needed when experimentation with alternative policies is prohibitively costly or is otherwise infeasible. Next term we shall study a complementary set of “reduced form” methods for ex-post policy evaluation.

2 Some Motivating Examples

Rust (1987). Time is discrete, indexed by $t \in \mathbb{N}_0 := \{0, 1, 2, \dots\}$. At each date t , an individual chooses an action a from a finite set $\mathcal{A} = \{0, 1, 2, \dots, A\}$ to solve

$$V_\theta(s_t) = \mathbb{E} \left[\max_{a \in \mathcal{A}} (u(s_t, a; \theta_1) + \varepsilon_{t,a} + \beta \mathbb{E}[V_\theta(s_{t+1}) | s_t, a; \theta_2]) \right], \quad (1)$$

where θ_1 is a vector of parameters indexing utilities, s_t is an observable (to the econometrician) Markov state taking values in a finite state-space $\mathcal{S} = \{0, 1, 2, \dots, S\}$ with transition distribution $f(s_{t+1} | s_t, a; \theta_2)$ which is known up to the parameter vector θ_2 , $\varepsilon_t = (\varepsilon_{t,0}, \varepsilon_{t,1}, \dots, \varepsilon_{t,A})$ is a vector of unobservable (to the econometrician) utility shocks which are typically assumed to be IID Gumbel,¹ and $\beta \in [0, 1)$ is a discount parameter. The outer expectation is an expectation over ε_t holding s_t fixed while the inner expectation is over s_{t+1} conditional on s_t, a . We wish to estimate model parameters $\theta = (\theta_1, \theta_2)$ so that we may perform policy experiments.

We begin by noting two simplifying implications of the Gumbel assumption. First, using a closed-form expression for the expectation of the maximum of independent Gumbel random variables, we have

$$V_\theta(s_t) = \gamma + \log \left(\sum_{a' \in \mathcal{A}} e^{u(s_t, a'; \theta_1) + \beta \mathbb{E}[V_\theta(s_{t+1}) | s_t, a'; \theta_2]} \right), \quad (2)$$

¹The Gumbel distribution has cdf $F(x) = \exp(-\exp(-x))$.

where $\gamma \approx 0.52277$ is the Euler-Mascheroni constant. The conditional choice probability (CCP) of choosing a in state s_t is also available in closed-form:

$$P(a|s_t; \theta) = \Pr \left(a = \arg \max_{a' \in \mathcal{A}} (u(s_t, a'; \theta_1) + \varepsilon_t(a') + \beta E[V_\theta(s_{t+1})|s_t, a'; \theta_2]) \middle| s_t \right) \quad (3)$$

$$= \frac{e^{u(s_t, a; \theta_1) + \beta E[V_\theta(s_{t+1})|s_t, a; \theta_2]}}{\sum_{a' \in \mathcal{A}} e^{u(s_t, a'; \theta_1) + \beta E[V_\theta(s_{t+1})|s_t, a'; \theta_2]}}. \quad (4)$$

Note that when the agent is myopic ($\beta = 0$) the model reduces to a standard multinomial logit model and the choice probabilities have the usual multinomial logit form for a static problem. Finally, note that according to the model, (s_t, a_t) is Markovian and the transition distribution factorizes as

$$\Pr(a_t, s_t | a_{t-1}, s_{t-1}) = \Pr(a_t | s_t) \Pr(s_t | a_{t-1}, s_{t-1}) \quad (5)$$

$$= P(a_t | s_t; \theta) f(s_t | s_{t-1}, a_{t-1}; \theta_2) \quad (6)$$

which allows us to write the log-likelihood as the sum of the log-likelihood of the CCPs and the log-likelihood for the law of motion of the observable component of the state s_t .

Two sampling schemes are typically used. The first is to assume that we observe a single time-series of the state s_t and the agent's action a_t for a sequence of n periods, so our data are $(s_0, a_0), \dots, (s_n, a_n)$. In this case, we have the likelihood function

$$Q_n(\theta) = \frac{1}{n} \sum_{t=1}^n (\log P(a_t | s_t; \theta) + \log f(s_t | s_{t-1}, a_{t-1}; \theta_2)). \quad (7)$$

Alternatively, we might observe a panel of data on independent, identical agents, $(s_{it}, a_{it})_{t=0}^T$ for $i = 1, \dots, n$ and fixed T , in which case we have the log-likelihood

$$Q_n(\theta) = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T (\log P(a_{it} | s_{it}; \theta) + \log f(s_{it} | s_{it-1}, a_{it-1}; \theta_2)). \quad (8)$$

These likelihoods are “approximate” since we condition on (s_0, a_0) . To work out the full likelihood we could infer the stationary distribution of (s_0, a_0) and then add on the date-0 contribution. Ignoring the date-0 contribution will not matter in large samples for time-series data, but could bias estimates for panel data with fixed T .

To compute $Q_n(\theta)$ we proceed in two steps. First, solve the recursion (2) for V_θ . Then, given V_θ , form the likelihood using expression (4). We maximize Q_n with respect to θ to obtain the *maximum likelihood estimate* (MLE) $\hat{\theta}$ of θ . This estimation method is sometimes referred to as Rust's nested fixed point algorithm. There is an active literature on estimating single and multiple-agent dynamic discrete choice models and dynamic discrete games.

Hansen and Singleton (1982). Consider a representative household that chooses its consumption and investment plan to maximize lifetime utility

$$\sum_{t=0}^{\infty} \delta^t \frac{C_t^{1-\gamma} - 1}{1-\gamma} \quad (9)$$

subject to a budget constraint. Standard arguments deliver the Euler equation

$$\mathbb{E} \left[\delta \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{i,t+1} \middle| \mathcal{I}_t \right] = 1, \quad (10)$$

where \mathcal{I}_t is the information set of the agent at date t and $R_{i,t+1}$ is the gross return on asset i from date t to date $t+1$.

We will use (10) to estimate the preference parameters δ and γ from time-series data on consumption growth C_{t+1}/C_t , a vector of asset returns $R_{t+1} = (R_{1,t+1}, \dots, R_{l,t+1})'$, and a vector of conditioning variables z_t that belong to \mathcal{I}_t (i.e., $\sigma(z_t) \subseteq \mathcal{I}_t$). Then by the law of iterated expectations and (10), at the true values of (δ, γ) we have

$$\mathbb{E} \left[\left(\delta \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{t+1} - 1 \right) \otimes z_t \right] = \mathbb{E} \left[\mathbb{E} \left[\left(\delta \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{t+1} - 1 \right) \otimes z_t \middle| \mathcal{I}_t \right] \right] \quad (11)$$

$$= \mathbb{E} \left[\mathbb{E} \left[\left(\delta \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{t+1} - 1 \right) \middle| \mathcal{I}_t \right] \otimes z_t \right] \quad (12)$$

$$= 0, \quad (13)$$

where $\mathbf{1}$ denotes a conformable vector of ones. Let $\theta = (\delta, \gamma)$ and

$$g(X_t; \theta) = \left(\delta \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{t+1} - 1 \right) \otimes z_t, \quad (14)$$

where $X_t = (C_{t+1}/C_t, R'_{t+1}, z'_t)'$. We know from (13) that the population moment condition

$$\mathbb{E}[g(X_t; \theta)] = 0 \quad (15)$$

must hold at the true value of θ . The idea of the *generalized method of moments* (GMM) is that we choose an estimate $\hat{\theta}$ of θ such that the sample average of $g(X_t; \theta)$ is as close to zero as possible. Let \widehat{W} be a positive definite symmetric matrix. The sample criterion function is

$$Q_n(\theta) = -\frac{1}{2} \left(\frac{1}{n} \sum_{t=1}^n g(X_t; \theta) \right)' \widehat{W} \left(\frac{1}{n} \sum_{t=1}^n g(X_t; \theta) \right). \quad (16)$$

Maximizing Q_n with respect to θ yields the GMM estimate $\hat{\theta}$.

Benhabib, Bisin, and Luo (2019). To identify what drives wealth dynamics in the US, the authors build a heterogeneous agent life-cycle model which exhibits various wealth accumulation factors. Solving the model leads to a stochastic process for wealth across generations and other attributes of heterogeneity. Their description of their empirical framework is as follows:

We estimate the parameters of the model described in the previous section using a [Simulated Method of Moments (SMM)] estimator: i) we fix (or externally calibrate) several parameters of the model; ii) we select some relevant moments of the wealth process as target in the estimation; and iii) we estimate the remaining parameters by matching the targeted moments generated by the stationary distribution induced by the model and those in the data...

We target as moments: the bottom 20%, 20–39%, 40–59%, 60–79%, 80–89%, 90–94%, 95–99%, and the top 1% wealth percentiles; and the diagonal of the social mobility Markov chain transition matrix defined over the same percentile ranges as states. iii) We estimate: the preference parameters μ , A ; and a parameterization of the stochastic process for r ...

In total, therefore, we target 15 moments and we estimate 12 parameters...

We use bootstrapping to generate the standard errors for the statistics related to the return process, e.g. its mean, standard deviation, and autocorrelation coefficient. The procedure is standard. We take the parameter values for generating the return process as given, i.e. the values for the five Markov states and the diagonal matrix of the transition matrix (hence the whole Markov transition matrix), then generate the return process a sufficiently large number of times. We then calculate the mean, standard errors and the autocorrelation coefficient directly using these series of the return processes.

Abstractly, we can write their estimation procedure as maximizing

$$Q_n(\theta) = -\frac{1}{2}(g_n - \gamma_m(\theta))' \widehat{W}(g_n - \gamma_m(\theta)) \quad (17)$$

where g_n are the moments being targeted, $\gamma_m(\theta)$ are the counterparts that are simulated from the model (averaged over m simulations) at parameter value θ , and \widehat{W} is a weight matrix. The resulting estimator $\hat{\theta}$ is the *simulated method of moments* (SMM) estimator. Unlike GMM estimation where the only source of randomness is sampling variation, here there is a second source of randomness that comes from using simulation to compute the model-implied moments. Both sources of randomness play a role in determining the large-sample properties of simulation-based estimators.

3 Preliminaries

3.1 Data

We will work throughout on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ where Ω is the underlying uncertainty space, \mathcal{F} is the collection of all events whose likelihood of occurrence the model will restrict, and \mathbb{P} is the probability measure that assigns the likelihood of events in \mathcal{F} .

The data we observe are a sequence X_1, \dots, X_n where n is the sample size and each of the n observations, denoted X_t , is a random vector (i.e. a vector of random variables). The sample is a map from Ω to $\mathbb{R}^{n \times \dim(X)}$, so we may sometimes use the map $\omega \mapsto X_1(\omega), \dots, X_n(\omega)$ to reflect this. We will always assume this map is *measurable*, in the sense that statements we will make about X_1, \dots, X_n correspond to events in \mathcal{F} .

3.2 Sampling Schemes

We will also allow for two sampling schemes, namely data that are *independent and identically distributed* (IID) or weakly dependent. IID data are most common in cross-sectional econometrics or micro-econometrics when the data are from cross-sectional surveys at the individual, household or firm level (though there may of course be clustering; we ignore this for now). Under IID sampling, the random variables X_1, \dots, X_n are independent and each observation X_t is a draw from a fixed probability distribution (and hence is “identically distributed”). With IID data, the strong law of large numbers (SLLN) asserts

$$\frac{1}{n} \sum_{t=1}^n f(X_t) \rightarrow_{a.s.} \mathbb{E}[f(X_t)] \quad (18)$$

for each function f for which $\mathbb{E}[f(X_t)]$ is finite.

As structural models often have a dynamic component, we also require a sampling scheme for time-series data. In this case we view X_1, \dots, X_n as a segment of the realization of a stochastic processes. Here we shall require that the data are *strictly stationary and ergodic*, which is a notion of “weak” dependence. There are two parts to this definition: strict stationarity requires that the distribution of $(X_{t_1+h}, \dots, X_{t_k+h})$ be independent of h , for each t_1, \dots, t_k . Ergodicity is a technical condition ensuring that time-series averages converge to spatial averages, i.e.:

$$\frac{1}{n} \sum_{t=1}^n f(X_t) \rightarrow_{a.s.} \mathbb{E}[f(X_t)] \quad (19)$$

for each function f for which $\mathbb{E}[f(X_t)]$ is finite. This convergence (19) is the result of what is known formally as the *Ergodic Theorem*, which generalizes the SLLN to weakly dependent data.

For instance, any stationary Markov process is ergodic under very mild conditions. This is useful as structural models with a dynamic component typically impose Markovianity for tractability.

3.3 Convergence Concepts and Mann–Wald Notation

Let Z_1, Z_2, \dots be a sequence of random vectors. Say that Z_1, Z_2, \dots *converges in probability* to Z_0 if

$$\Pr(\|Z_n - Z_0\| > \epsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty \quad (20)$$

for each $\epsilon > 0$. If so, we write $Z_n \rightarrow_p Z_0$. Say that $Z_n = o_p(1)$ if $Z_n \rightarrow_p 0$. Notice that $Z_n \rightarrow_p Z_0$ if and only if $Z_n - Z_0 \rightarrow_p 0$. We use $Z_n \rightarrow_p Z_0$ and $Z_n = Z_0 + o_p(1)$ interchangeably. We can think of $o_p(\cdot)$ as being like a stochastic version of the $o(\cdot)$ notation from analysis.

Say that Z_1, Z_2, \dots *converges almost surely* to Z_0 if

$$\mathbb{P}\left(\left\{\omega : \lim_{n \rightarrow \infty} Z_n(\omega) = Z_0(\omega)\right\}\right) = 1. \quad (21)$$

If so, we write $Z_n \rightarrow_{a.s.} Z_0$. Recall that almost sure convergence implies convergence in probability, i.e. $Z_n \rightarrow_{a.s.} Z_0$ implies $Z_n \rightarrow_p Z_0$. Say that $Z_n = o_{a.s.}(1)$ if $Z_n \rightarrow_{a.s.} 0$. Notice that $Z_n \rightarrow_{a.s.} z_0$ if and only if $Z_n - z_0 \rightarrow_{a.s.} 0$. We use $Z_n \rightarrow_{a.s.} Z_0$ and $Z_n = Z_0 + o_{a.s.}(1)$ interchangeably.

Let $\alpha_1, \alpha_2, \dots$ be a sequence of positive constants. Say that $Z_n = o_p(\alpha_n^{-1})$ (resp. $Z_n = o_{a.s.}(\alpha_n^{-1})$) if $\alpha_n Z_n = o_p(1)$ (resp. $\alpha_n Z_n = o_{a.s.}(1)$).

Say that Z_1, Z_2, \dots is *bounded in probability* or *tight* if for each $\epsilon > 0$ there exists a $M_\epsilon < \infty$ such that

$$\limsup_{n \rightarrow \infty} \Pr(\|Z_n\| > M_\epsilon) \leq \epsilon. \quad (22)$$

If so, we write $Z_n = O_p(1)$. Similarly, say that $Z_n = O_p(\alpha_n)$ if $\alpha_n Z_n = O_p(1)$. We can think of $O_p(\cdot)$ as being like a stochastic version of the $O(\cdot)$ notation from analysis. We have the following relations:

$$o_p(a_n) \times o_p(b_n) = o_p(a_n b_n), \quad (23)$$

$$O_p(a_n) \times o_p(b_n) = o_p(a_n b_n), \text{ and} \quad (24)$$

$$O_p(a_n) \times O_p(b_n) = O_p(a_n b_n). \quad (25)$$

The same relations hold if we replace the o_p and O_p terms by $o_{a.s.}$ and $O_{a.s.}$ terms.

For example, let X_1, \dots, X_n be IID random vectors with mean μ and finite covariance matrix Σ . Then the sample mean $\bar{X}_n \rightarrow_{a.s.} \mu$ by the SLLN (and hence $\bar{X}_n \rightarrow_p \mu$). By the central limit theorem we have that $\sqrt{n}(\bar{X}_n - \mu) \rightarrow_d N(0, \Sigma)$ hence $(\bar{X}_n - \mu) = O_p(n^{-1/2})$ and $\bar{X}_n = \mu + O_p(n^{-1/2})$.

4 Extremum Estimators and Examples

Consider a class of economic models indexed by a parameter of interest, say θ . The set Θ is the *parameter space*, which is the set of plausible values of θ . We say $\hat{\theta}$ is an *extremum estimator* if it is a (measurable) function of the data X_1, \dots, X_n taking values in Θ and it satisfies

$$Q_n(\hat{\theta}) > \sup_{\theta \in \Theta} Q_n(\theta) - \eta_n, \quad (26)$$

where η_n is a sequence of positive random variables such that $\eta_n \rightarrow_p 0$ and $Q_n : \Theta \rightarrow \mathbb{R}$ is a *sample criterion function* which depends on both θ and the data X_1, \dots, X_n . The function Q_n is a random function: different realizations of the sample X_1, \dots, X_n may produce a different sample criterion function, and therefore a different maximizer $\hat{\theta}$. We include the $-\eta_n$ term because in practice we typically use numerical methods for maximizing Q_n and therefore we are not guaranteed that the optimization routine will converge to the exact maximum.

Throughout we denote the *true* value of θ by θ_0 . The true value $\theta_0 \in \Theta$ solves

$$Q(\theta_0) \geq Q(\theta) \quad \text{for all } \theta \in \Theta \text{ with } \theta \neq \theta_0, \quad (27)$$

where $Q : \Theta \rightarrow \mathbb{R}$ is the population criterion function. Loosely speaking, this is the function to which Q_n would converge if we saw an infinite amount of data.

Say θ_0 is *(point) identified* if

$$Q(\theta_0) > Q(\theta) \quad \text{for all } \theta \in \Theta \text{ with } \theta \neq \theta_0. \quad (28)$$

The first five lectures defined identification in terms of the full set of restrictions that the model contains about θ . The two definitions agree when Q exploits all such restrictions. Some care may need to be taken to do this. In the Hansen and Singleton (1982) example above, the model gives a conditional moment restriction against an information set \mathcal{I}_t but we use only a subset of this information for estimation. It may still be that θ_0 is identified through the GMM population criterion function using only a subset of the information. But in this case the resulting estimate may not be as efficient as it could have been had we exploited the full set of model restrictions.

4.1 M-Estimators

M-estimators are based on the sample criterion function of the form

$$Q_n(\theta) = \frac{1}{n} \sum_{t=1}^n m(X_t; \theta). \quad (29)$$

M-estimators differ in their choice of m . With IID or SSE data the sample average $\frac{1}{n} \sum_{t=1}^n m(X_t; \theta)$ converges a.s. to $E[m(X_t; \theta)]$ by the SLLN or Ergodic Theorem. The *population criterion function* $Q : \Theta \rightarrow \mathbb{R}$ is therefore

$$Q(\theta) = E[m(X_t; \theta)]. \quad (30)$$

Example 4.1. Linear Regression.

Suppose $(X_t, Y_t)_{t=1}^n$ are IID or SSE with Y_t a scalar and X_t a vector of dimension p . Suppose that both Y_t and X_t have finite second moments. The best linear predictor for Y_t given X_t is $X_t' \theta_0$, where θ_0 solves

$$\min_{\theta \in \mathbb{R}^p} E[(Y_t - X_t' \theta)^2]. \quad (31)$$

Equivalently, θ_0 maximizes the population criterion function

$$Q(\theta) = -\frac{1}{2} E[(Y_t - X_t' \theta)^2]. \quad (32)$$

Expanding the quadratic, we see that θ_0 uniquely maximizes Q whenever $E[X_t X_t']$ has full rank. In regression jargon, this is the no-multicollinearity condition.

The sample criterion function is

$$Q_n(\theta) = -\frac{1}{2} \frac{1}{n} \sum_{t=1}^n (Y_t - X_t' \theta)^2. \quad (33)$$

By the usual argument, we see that the least-squares estimator

$$\hat{\theta} = \left(\frac{1}{n} \sum_{t=1}^n X_t X_t' \right)^{-1} \left(\frac{1}{n} \sum_{t=1}^n X_t Y_t \right) \quad (34)$$

maximizes Q_n . This is a rare example of a situation in which we can solve for the estimator in closed form. \square

Example 4.2. Maximum Likelihood (Unconditional).

Suppose X_1, \dots, X_n are IID draws from a distribution with density $f(x; \theta_0)$ with respect to a common dominating measure μ . Continuous distributions have densities with respect to Lebesgue measure; discrete distributions have densities with respect to counting measure. The joint density of X_1, \dots, X_n is

$$f(X_1, \dots, X_n; \theta_0) = \prod_{t=1}^n f(X_t; \theta_0). \quad (35)$$

Of course, we do not know the true θ_0 . The principle of maximum likelihood asserts that we should choose θ to maximize the “probability” of having observed X_1, \dots, X_n . Maximization of the

probability

$$\prod_{t=1}^n f(X_t; \theta) \quad (36)$$

with respect to θ is equivalent to maximizing

$$Q_n(\theta) = \frac{1}{n} \sum_{t=1}^n \log f(X_t; \theta) \quad (37)$$

with respect to θ . (Unconditional) maximum likelihood is therefore a M-estimator whose m function is the log density: $m(X_t; \theta) = \log f(X_t; \theta)$.

We now verify the identification condition (28). The population criterion function is

$$Q(\theta) = E[\log f(X_t; \theta)]. \quad (38)$$

Suppose that the model is correctly specified. That is, X_1, \dots, X_n are IID draws from a distribution with density $f(X_t; \theta_0)$. This means when we take expectations over functions of X we integrating respect to $f(\cdot; \theta_0)$. Then

$$Q(\theta) = E[\log f(X_t; \theta)] - E[\log f(X_t; \theta_0)] + E[\log f(X_t; \theta_0)] \quad (39)$$

$$= -E\left[\log\left(\frac{f(X_t; \theta_0)}{f(X_t; \theta)}\right)\right] + E[\log f(X_t; \theta_0)] \quad (40)$$

$$= -K(f(\cdot; \theta_0) \| f(\cdot; \theta)) + E[\log f(X_t; \theta_0)]. \quad (41)$$

where

$$K(f \| g) = \int f \log(f/g) d\mu \quad (42)$$

is the *Kullback-Leibler divergence* (or *relative entropy*) of g from f . This divergence is a measure of the discrepancy between densities. It has the properties that (i) $K(f \| g) \geq 0$ and (ii) $K(f \| g) = 0$ if and only if $f = g$ holds f -almost everywhere. To see why property (i) holds, by convexity of $-\log(x)$ and Jensen's inequality, we have

$$K(f \| g) = \int -\log(g/f) f d\mu \geq -\log\left(\int (g/f) f d\mu\right) = -\log\left(\int g d\mu\right) = -\log(1) = 0.$$

For property (ii), we note that if $f = g$ holds f -almost everywhere then $\log(f/g) = 0$ holds f -almost everywhere and hence

$$K(f \| g) = \int \log(f/g) f d\mu = \int 0 d\mu = 0.$$

Conversely, suppose $f \neq g$ with positive f -measure. Then by Jensen's inequality, which holds strictly here because $-\log x$ is strictly convex, we have

$$K(f \| g) = \int -\log(g/f) f d\mu > -\log\left(\int g d\mu\right) = 0.$$

Consider (41). As $E[\log f(X_t; \theta_0)]$ does not vary with θ , we see that maximizing Q is equivalent to minimizing $K(f(\cdot; \theta_0) \| f(\cdot; \theta))$. By properties (i) and (ii), a sufficient condition for the identification condition (28) is therefore

$$\Pr(f(X_t; \theta) \neq f(X_t; \theta_0)) > 0 \quad \text{for all } \theta \neq \theta_0, \quad (43)$$

where the probability is with respect to the true distribution (under θ_0). This is a type of “rank” condition that says different θ must induce noticeably different distributions. We can usually verify this condition directly from the model. \square

Example 4.3. Maximum Likelihood (Markov Processes).

Suppose that Y_0, \dots, Y_n are generated by a SSE Markov process with transition density $f(Y_t | Y_{t-1}; \theta_0)$. Once again, “density” means with respect to a dominating measure μ . The Rust (1987) model falls into this framework with $Y_t = (s_t, a_t)$ and some special structure on the transition density.

Suppose that the functional form f is known (i.e., computable) but θ_0 is unknown. Let $X_t = (Y_t, Y_{t-1})$ for $t = 1, \dots, n$. If θ was the true parameter, the joint density of X_1, \dots, X_n would be

$$f(X_1, \dots, X_n; \theta) = \prod_{t=1}^n f(Y_t | Y_{t-1}; \theta) \times f_0(Y_0; \theta) \quad (44)$$

where f_0 is the unconditional (or stationary) density of Y_t . The average log likelihood is

$$\frac{1}{n} \sum_{t=1}^n \log f(Y_t | Y_{t-1}; \theta) + \frac{1}{n} \log f(Y_0; \theta). \quad (45)$$

As n gets large the second term becomes negligible. Therefore, maximizing the average log likelihood is approximately the same as maximizing

$$Q_n(\theta) = \frac{1}{n} \sum_{t=1}^n \log f(Y_t | Y_{t-1}; \theta). \quad (46)$$

Hence we have a M-estimator with $m(X_t; \theta) = \log f(Y_t | Y_{t-1}; \theta)$. The population criterion function is now

$$Q(\theta) = E[\log f(Y_t | Y_{t-1}; \theta)], \quad (47)$$

where the expectation is with respect to the joint (stationary) distribution of (Y_t, Y_{t-1}) .

We may use similar arguments to write the maximum likelihood for higher-order Markov processes (i.e. Markov processes for which the distribution of Y_t depends on Y_{t-1}, \dots, Y_{t-k} but not also on Y_{t-j} for any $j < k$) as M-estimators.

We now verify the identification condition (28). Suppose that the model is correctly specified. Notice that

$$Q(\theta) = -\mathbb{E}\left[\mathbb{E}\left[\log\left(\frac{f(Y_t|Y_{t-1};\theta_0)}{f(Y_t|Y_{t-1};\theta)}\right)\middle|Y_{t-1}\right]\right] + \mathbb{E}[\log f(Y_t|Y_{t-1};\theta_0)]. \quad (48)$$

In the first expression on the right-hand side, the outer expectation is with respect to the stationary distribution of Y_{t-1} implied by $f(\cdot|\cdot;\theta_0)$ while the inner expectation is taken with respect to $Y_t \sim f(\cdot|Y_{t-1};\theta_0)$. By similar arguments to the unconditional case, we see that θ_0 is identified if for all $\theta \neq \theta_0$ there is a set of Y_{t-1} with positive probability (under the stationary distribution) upon which

$$\Pr(f(Y_t|Y_{t-1};\theta) \neq f(Y_t|Y_{t-1};\theta_0)|Y_{t-1}) > 0, \quad (49)$$

where the \Pr is taken with respect to $f(\cdot|Y_{t-1};\theta_0)$. In words, this says that for each $\theta \neq \theta_0$ there exist states Y_{t-1} that we visit with positive probability and in which the transition distribution $f(\cdot|Y_{t-1};\theta)$ is noticeably different from $f(\cdot|Y_{t-1};\theta_0)$.

As a simple example, consider the AR(1) model

$$Y_t = (1 - \rho)\mu + \rho Y_{t-1} + u_t, \quad (50)$$

where the u_t are IID $N(0, \sigma^2)$ and $\rho \in (-1, 1)$. The conditional distribution of Y_t given Y_{t-1} is

$$N((1 - \rho)\mu + \rho Y_{t-1}, \sigma^2). \quad (51)$$

Each of the three parameters (μ, ρ, σ^2) is identified: varying σ^2 changes the dispersion of the distribution, varying μ changes the mean in a common way across all states, and varying ρ changes the mean in a state-dependent way.

The stationary (or unconditional) distribution of Y_t may be shown to be

$$Y_t \sim N(\mu, (1 - \rho^2)^{-1}\sigma^2). \quad (52)$$

A naive estimation strategy would be to ignore the dynamics and attempt to estimate (μ, ρ, σ^2) based on the likelihood of the unconditional distribution. While μ is identified from the stationary distribution, the parameters (ρ, σ^2) are not, as we can choose different (ρ, σ^2) pairs to produce the same value of the unconditional variance $(1 - \rho^2)^{-1}\sigma^2$. \square

4.2 GMM

The M-estimators we just looked at are based on the idea that there is some “true” parameter θ_0 that uniquely maximizes $\mathbb{E}[m(X_t; \theta)]$ over the parameter space $\Theta \subseteq \mathbb{R}^p$. GMM is based upon the idea that there exists a “true” parameter θ_0 in the parameter space Θ that uniquely sets a

vector of population moments to zero, i.e., $E[g(X_t; \theta_0)] = 0$ where $g : \mathbb{R}^{\dim(X)} \times \Theta \rightarrow \mathbb{R}^K$ is jointly measurable in the data X_t and θ . We assume there $K \geq p := \dim(\theta)$ moment conditions.

The Hansen and Singleton (1982) example motivated a particular choice of g from the first-order condition for optimality of a consumption/investment plan. The idea of GMM applies more broadly. For instance, we may g so that it consists of the moments that we really want our model to explain. For instance, in structural models of the labor market we may want to explain quantities like mean wages, mean wage changes among employed individuals, and durations of unemployment, etc, and would choose g accordingly.

Formally, $\hat{\theta}$ is a *GMM estimator* if it is an extremum estimator in the sense of display (26) where

$$Q_n(\theta) = -\frac{1}{2}g_n(\theta)' \widehat{W} g_n(\theta) \quad \text{with} \quad g_n(\theta) = \frac{1}{n} \sum_{t=1}^n g(X_t; \theta), \quad (53)$$

and where \widehat{W} is a $K \times K$ positive-definite symmetric weight matrix that may be a function of the data. Letting W (another positive-definite and symmetric weight matrix) denote the probability limit of \widehat{W} and noting that $g_n(\theta) \rightarrow_{a.s.} E[g(X_t; \theta)]$ by the SLLN or Ergodic Theorem, we see that the population criterion function $Q : \Theta \rightarrow \mathbb{R}$ is

$$Q(\theta) = -\frac{1}{2}E[g(X_t; \theta)]' W E[g(X_t; \theta)] \quad (54)$$

$$= -\frac{1}{2}g(\theta)' W g(\theta) \quad \text{where} \quad g(\theta) = E[g(X_t; \theta)]. \quad (55)$$

We say the GMM model is *correctly specified* if there exists some θ_0 in Θ for which $E[g(X_t; \theta_0)] = 0$. This is a weaker notion of correct specification than with maximum likelihood (ML). With ML the model was correctly specified if there was a value θ_0 that gave us the true distribution of the data; here the model is correctly specified if there is a θ_0 that matches the moments we care about (i.e., those we put in the g function), but not possibly the full distribution of the data.

The model is *identified* if there is only one such θ_0 that sets the population moments to zero:

$$E[g(X_t; \theta)] = 0 \quad \Longleftrightarrow \quad \theta = \theta_0. \quad (56)$$

If the GMM model is identified, then $Q(\theta)$ is uniquely maximized at θ_0 . This is a weaker notion of identification for maximum likelihood: it may be that $E[g(X_t; \theta)] = 0$ for two distinct values of θ , but these values generate different distributions over observables. Hence, maximum likelihood would distinguish the two but GMM would not. In this sense, GMM is a “limited information” estimation procedure as it may not exhaust the full amount of information that the model contains about the parameters of interest.

GMM generalizes IV estimation. In the special case of an IV model

$$Y_{1t} = Y_{2t}'\theta + u_t, \quad E[u_t Z_t] = 0,$$

where Y_2 is an endogenous regressor and Z is a vector of instruments, we have

$$g(X_t; \theta) = Z_t(Y_{1t} - Y_{2t}'\theta)$$

with $X_t = (Y_{1t}, Y_{2t}, Z_t)$. Here the dimension K of g is the dimension of the IV Z_t . By analogy with linear IV estimation, clearly we need the number of moments K to be at least as large as the number of parameters p for identification. If $K = p$ and the model is identified then we say it is *just identified*. If $K > p$ and the model is identified we say it is *over identified*. In this case, the additional $K - p$ restrictions are testable.

Common choices of \widehat{W} are the $K \times K$ identity matrix or the two-step “optimal” weight matrix

$$\widehat{W} = \left(\frac{1}{n} \sum_{t=1}^n g(X_t; \tilde{\theta}) g(X_t; \tilde{\theta})' \right)^{-1} \quad (57)$$

where $\tilde{\theta}$ is a consistent estimator of θ_0 (such as a GMM estimator we’ve computed using the identity weight matrix). Later we will discuss when this approach is optimal, and in what sense. There are also other estimators for which \widehat{W} itself depends on θ . These fall into the class of “generalized empirical likelihood” estimators, which are outside the scope of this course.

4.3 SMM

A special case of GMM arises when the moment conditions are *separable* in parameters and data:

$$g(X_t; \theta) = g(X_t) - \gamma(\theta) \quad (58)$$

where $g : \mathbb{R}^{\dim(X)} \rightarrow \mathbb{R}^K$ and $\gamma : \Theta \rightarrow \mathbb{R}^K$. Although this looks a bit restrictive, this situation will arise when $E[g(X_t)]$ are “targeted” population moments that we want to match and $\gamma(\theta)$ are their model-implied counterparts. As we don’t observe the true population moments, we estimate them using the sample average:

$$g_n = \frac{1}{n} \sum_{t=1}^n g(X_t) \quad (59)$$

If we know the functional form for γ , then we can use GMM to estimate θ by maximizing

$$Q_n(\theta) = -\frac{1}{2} (g_n - \gamma(\theta))' \widehat{W} (g_n - \gamma(\theta)) \quad (60)$$

for some $K \times K$ weight matrix \widehat{W} . This is sometimes referred to as estimating the model by “calibration”. The matrix \widehat{W} may be chosen to prioritize fitting some moments more than others.

Suppose the model is sufficiently complicated that there is no closed-form expression for $\gamma(\theta)$. For example, your model might have heterogeneous agents or firms and specify the equilibrium cross-sectional distribution of various quantities as

$$f(X_t|\theta). \quad (61)$$

The model-implied aggregate moments would then be

$$\gamma(\theta) = \int \gamma(X_t; \theta) f(X_t|\theta) d\mu(X_t). \quad (62)$$

This integral might be difficult to calculate numerically, even if X is reasonably low-dimensional. The idea of SMM is to use simulation to approximate the integral in the above display. That is, generate a large number of observations $X_1^\theta, \dots, X_m^\theta$ from $f(X|\theta)$ then set

$$\gamma_m(\theta) = \frac{1}{m} \sum_{s=1}^m \gamma(X_s^\theta; \theta). \quad (63)$$

To generate X_s^θ , one typically simulates a vector of primitive shocks ε_s , then solves the model to obtain X_s^θ . SMM can be computationally burdensome to implement when solving the model is a non-trivial exercise. The usual caveats about numerical integration also apply.

Formally, $\hat{\theta}$ is a *SMM estimator* if it is an extremum estimator in the sense of display (26) for which the sample criterion function is

$$Q_n(\theta) = -\frac{1}{2}(g_n - \gamma_m(\theta))' \widehat{W}(g_n - \gamma_m(\theta)) \quad (64)$$

for some $K \times K$ weight matrix \widehat{W} . Note here there are two sources of uncertainty that must be accounted for when computing standard errors, namely (i) sampling uncertainty from the data X_1, \dots, X_n used to compute g_n and (ii) additional uncertainty introduced by replacing the true $\gamma(\theta)$ by a simulation-based estimate $\gamma_m(\theta)$. The population version of the criterion function is

$$Q(\theta) = -\frac{1}{2}(g_0 - \gamma(\theta))' W(g_0 - \gamma(\theta)) \quad (65)$$

where g_0 and W denote the probability limits of g_n and \widehat{W} , respectively.

Note: when implementing SMM, it is important to use the same random seed for generating $X_1^\theta, \dots, X_m^\theta$ for each θ . Otherwise, you can get different values of $Q_n(\theta)$ for the same θ , so whatever numerical procedure you use to maximize $Q_n(\theta)$ will not converge.

4.4 SMD

Simulated minimum distance (SMD) estimation generalizes SMM to a setting where g_n are “sample statistics” that are not necessarily expressible as “moments”. For example, g_n may be quantiles of a wealth distribution. The idea of *minimum distance* (MD) estimation is to estimate θ by minimizing the distance between g_n and its true model-implied counterpart $\gamma(\theta)$. Say $\hat{\theta}$ is a *MD estimator* if it is an extremum estimator in the sense of display (26) for which the sample criterion function is

$$Q_n(\theta) = -\frac{1}{2}(g_n - \gamma(\theta))' \widehat{W} (g_n - \gamma(\theta)). \quad (66)$$

for some $K \times K$ weight matrix \widehat{W} . When the expressions for $\gamma(\theta)$ are not available in closed form but we can simulate from the model, we can compute estimates $\gamma_m(\theta)$ based on simulation. We say $\hat{\theta}$ is a *simulated minimum distance* (SMD) estimator if the sample criterion function is

$$Q_n(\theta) = -\frac{1}{2}(g_n - \gamma_m(\theta))' \widehat{W} (g_n - \gamma_m(\theta)). \quad (67)$$

The population version of these criterion functions is the same as in display (65).

5 Misspecification

When the model is misspecified, we interpret the maximizer θ_0 of the population criterion function Q as the *pseudo-true* parameter. That is, θ_0 is no longer the true data-generating parameter but rather a parameter that brings the (misspecified) model “closest”, in a certain sense, to the true data-generating process.

5.1 Maximum Likelihood

To fix ideas, consider the simplest case of unconditional maximum likelihood. Suppose the data X_1, \dots, X_n are IID from a distribution with true density $g(x)$. Repeating the argument from Example 4.2 above, we have

$$Q(\theta) = E[\log f(X_t; \theta)] - E[\log g(X_t)] + E[\log g(X_t)] \quad (68)$$

$$= -E\left[\log\left(\frac{g(X_t)}{f(X_t; \theta)}\right)\right] + E[\log g(X_t)] \quad (69)$$

$$= -K(g\|f(\cdot; \theta)) + E[\log g(X_t)], \quad (70)$$

where the expectation is taken with respect to $X \sim g$. Evidently, maximizing Q is equivalent to minimizing

$$K(g\|f(\cdot; \theta)) \quad (71)$$

with respect to θ . The pseudo-true parameter θ_0 is therefore that which brings the model-implied distribution $f(\cdot; \theta)$ as close as possible to the true distribution g , where we measure closeness by Kullback–Leibler divergence.

There are infinitely many other ways of measuring “distance” between probability measures: take any non-negative convex function ϕ with $\phi(1) = 0$ and set

$$D_\phi(f\|g) = \int \phi(f/g)g \, d\mu;$$

Kullback–Leibler divergence corresponds to the special case of $\phi(x) = x \log x - x + 1$. We could in principle construct different criterion functions based on different ϕ functions or different notions of distance. Under correct specification these would all have the true parameter as their population maximizer. However, under misspecification the pseudo-true parameters may depend on the choice of criterion function. It is not clear why the parameter that minimizes Kullback–Leibler divergence, as opposed to some other divergence, is the interesting one in this case.

5.2 GMM

A GMM model is misspecified if $E[g(X_t; \theta)] \neq 0$ for all $\theta \in \Theta$. The pseudo-true parameter θ_0 is still well defined as the maximizer of the population criterion function. However, here θ_0 will depend implicitly on the asymptotic weight matrix W . Different choices of \widehat{W} and hence different W correspond to different estimands under misspecification. We can see this most clearly in the case of linear IV:

$$Y_{1t} = Y_{2t}'\theta + u_t, \quad E[u_t Z_t] = 0.$$

Suppose that the model is over identified (i.e., the number of instruments K exceeds the number of regressors p). The model is misspecified if there is no $\theta \in \Theta$ for which $E[Z_t(Y_{1t} - Y_{2t}'\theta)] = 0$.² The population criterion function is

$$Q(\theta) = -\frac{1}{2}E[Z_t(Y_{1t} - Y_{2t}'\theta)]'WE[Z_t(Y_{1t} - Y_{2t}'\theta)].$$

This is a quadratic function of θ so we can solve for the maximizer in closed-form to obtain

$$\theta_0 = (E[Y_{2t}Z_t']WE[Z_tY_{2t}'])^{-1}E[Y_{2t}Z_t']WE[Z_tY_{1t}],$$

²If the model is just identified, we can always solve the moment condition by setting $\theta_0 = E[Z_tY_{2t}']^{-1}E[Z_tY_{1t}]$.

which is well defined whenever $E[Y_{2t}Z_t']$ has full rank p . As $K > p$, we cannot simplify this expression further. The pseudo-true parameter θ_0 therefore depends on W . That is, $\theta_0 = \theta_0(W)$ and so different choices of weight matrix lead to different estimands.

Additional References

- Benhabib, J., A. Bisin, and M. Luo (2019). Wealth Distribution and Social Mobility in the US: A Quantitative Approach. *American Economic Review* 109(5), 1623–1647.
- Hansen, L. P. and K. Singleton (1982). Generalized Instrumental Variables Estimation of Nonlinear Rational Expectations Models. *Econometrica* 50(5), 1269–1286.
- McFadden, D. (1989). A Method of Simulated Moments for Estimation of Discrete Response Models Without Numerical Integration. *Econometrica*, 57(5), 995–1026.
- Rust, J. (1987). Optimal Replacement of GMC Bus Engines: An Empirical Model of Harold Zurcher. *Econometrica* 55(5), 999–1033.
- White, H. (1982). Maximum Likelihood Estimation of Misspecified Models. *Econometrica* 50(1), 1–25.

T. Christensen

1 Consistency of Extremum Estimators

Recall that an estimator $\hat{\theta}$ of θ_0 is *consistent* if

$$\hat{\theta} \rightarrow_p \theta_0 \quad \text{as } n \rightarrow \infty. \quad (1)$$

Consistency is a useful property. It says that as we observe more data, the probability of our estimator $\hat{\theta}$ being close to the estimand θ_0 should approach 1.

The following result is our master consistency result. The result can be used for M-estimation as well as GMM, SMM and MD. Let $\|\cdot\|$ denote a norm on Θ .

Theorem 1 (Consistency of extremum estimators). *Let the following hold:*

(i) (clean maximum) for any $\delta > 0$ we have $\sup_{\theta \in \Theta: \|\theta - \theta_0\| \geq \delta} Q(\theta) < Q(\theta_0)$

(ii) (uniform convergence) $\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| = o_p(1)$.

Then: any estimator $\hat{\theta}$ that satisfies (2) is consistent, i.e. $\hat{\theta} \rightarrow_p \theta$ as $n \rightarrow \infty$.

The intuition is as follows (also see Figure 1). The estimator $\hat{\theta}$ is obtained by maximizing Q_n :

$$Q_n(\hat{\theta}) > \sup_{\theta \in \Theta} Q_n(\theta) - \eta_n, \quad (2)$$

where $\eta_n \geq 0$ is $o_p(1)$. If θ_0 is identified, then the population objective function Q is uniquely maximized at θ_0 . As $\hat{\theta}$ is obtained by maximizing Q_n and we know that Q_n becomes closer to Q as we observe more data, the maximum of Q_n should become closer to θ_0 .

“Clean maximum” means $Q(\theta)$ can only approach $Q(\theta_0)$ as $\theta \rightarrow \theta_0$. This is needed to rule out situations in which $Q(\theta)$ may asymptote to $Q(\theta_0)$ as θ moves along certain directions (see Figure 2).

“Uniform convergence” means Q_n converges to Q in probability uniformly over the parameter space. This rules out, e.g., Q_n having a bump that moves around as n gets large.

For instance, suppose $\Theta = [-1, 1]$ and $Q : \Theta \rightarrow \mathbb{R}$ is continuous, with a unique maximum at $\theta_0 \neq 0$. Suppose also that

$$Q_n(\theta) = \begin{cases} Q(\theta) & \text{if } \theta \neq \frac{1}{n} \\ Q(\theta_0) + 1 & \text{if } \theta = \frac{1}{n} \end{cases} \quad (3)$$

Then $Q_n(\theta)$ converges pointwise to $Q(\theta)$ but not uniformly, because $\sup_{\theta} |Q_n(\theta) - Q(\theta)| \geq 1$. But also note that for each $n \geq 1$ the argmax of $Q_n(\theta)$ is $\hat{\theta} = \frac{1}{n}$, which converges to 0 $\neq \theta_0$.

Proof of Theorem 1. We want to show that $\Pr(\|\hat{\theta} - \theta_0\| > \delta) \rightarrow 0$ (as $n \rightarrow \infty$) for each $\delta > 0$.

Fix any $\delta > 0$. Let $\epsilon = Q(\theta_0) - \sup_{\theta \in \Theta: \|\theta - \theta_0\| \geq \delta} Q(\theta)$. Note $\epsilon > 0$ by (i).

As $\eta_n = o_p(1)$ and $\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| = o_p(1)$, we have with probability approaching one (wpa1) that

$$|\eta_n| < \frac{\epsilon}{3}, \quad \sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| < \frac{\epsilon}{3}. \quad (4)$$

Whenever these inequalities hold, we therefore have that

$$Q(\hat{\theta}) > Q(\theta_0) - \epsilon = \sup_{\theta \in \Theta: \|\theta - \theta_0\| \geq \delta} Q(\theta), \quad (5)$$

where the second equality is by definition of ϵ . It follows that $\|\hat{\theta} - \theta_0\| \leq \delta$ must hold whenever inequality (4) holds. But as (4) holds wpa1, we have therefore shown

$$\Pr(\|\hat{\theta} - \theta_0\| \leq \delta) \rightarrow 1, \quad (6)$$

as required. ■

Remark 1. If we assume $\eta_n = o_{a.s.}(1)$ and replace (ii) with $\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| = o_{a.s.}(1)$, then we can show that $\hat{\theta} \rightarrow_{a.s.} \theta_0$ as $n \rightarrow \infty$. The proof is left as an exercise.

2 Verifying Clean Maximum

There are many sufficient conditions for clean maximum. Here is one set:

Lemma 1 (Verifying “clean maximum”). *Let the following hold:*

(i) Θ is compact

(ii) $Q : \Theta \rightarrow \mathbb{R}$ is continuous

(iii) $Q(\theta_0) > Q(\theta)$ for each $\theta \in \Theta$ with $\theta \neq \theta_0$.

Then: “clean maximum” holds.

Proof. Fix any $\delta > 0$. The set $\{\theta \in \Theta : \|\theta - \theta_0\| \geq \delta\}$ is compact by (i). Then by (ii), we know that there is some $\theta^* \in \{\theta \in \Theta : \|\theta - \theta_0\| \geq \delta\}$ such that $\sup_{\theta \in \Theta: \|\theta - \theta_0\| \geq \delta} Q(\theta) = Q(\theta^*)$ and by (iii) we must have $Q(\theta^*) < Q(\theta_0)$. ■

3 Verifying Uniform Convergence

This is done differently for M-estimators, GMM, SMM, and MD.

3.1 Consistency of M-Estimators

For M-estimators, it suffices to show that the following *uniform* law of large numbers holds:

$$\sup_{\theta \in \Theta} \left| \underbrace{\frac{1}{n} \sum_{t=1}^n m(X_t, \theta)}_{Q_n(\theta)} - \underbrace{E[m(X_t, \theta)]}_{Q(\theta)} \right| = o_p(1). \quad (7)$$

Note that this is a stronger notion than the law of large numbers which asserts the pointwise result

$$\frac{1}{n} \sum_{t=1}^n m(X_t, \theta) - E[m(X_t, \theta)] = o_p(1) \quad (8)$$

for each θ .

We establish uniform convergence using a notion of the “size” or “complexity” of the class of functions whose average we are taking. It will turn out that uniform convergence holds whenever the class of functions $\mathcal{M} = \{m(\cdot; \theta) : \theta \in \Theta\}$ is small enough that it has *finite bracketing numbers*. Later in the course, we will see that similar notions of size or complexity are used to establish convergence results for nonparametric and modern machine learning methods.

Let $L^1 = \{f(X_t) : E[|f(X_t)|] < \infty\}$ and let $\mathcal{F} \subset L^1$ be a collection of functions of interest. The list of pairs of functions

$$l_{\varepsilon,1}, u_{\varepsilon,1}, l_{\varepsilon,2}, u_{\varepsilon,2}, \dots, l_{\varepsilon,N}, u_{\varepsilon,N} \subset L^1 \quad (9)$$

is said to *bracket* \mathcal{F} at level ε if for each $f \in \mathcal{F}$ we can choose a pair $l_{\varepsilon,i}$ and $u_{\varepsilon,i}$ such that $l_{\varepsilon,i} \leq f \leq u_{\varepsilon,i}$ and $E[u_{\varepsilon,i} - l_{\varepsilon,i}] \leq \varepsilon$ for each i . The ε -*bracketing number* of \mathcal{F} , denoted $N_{[\cdot]}(\mathcal{F}, \varepsilon)$, is the minimal number pairs required to bracket \mathcal{F} at level ε . If $N_{[\cdot]}(\mathcal{F}, \varepsilon) < \infty$ for all $\varepsilon > 0$ then we say that \mathcal{F} has *finite bracketing numbers*.

Lemma 2 (Uniform Strong Law of Large Numbers (ULLN)). *Let the following hold:*

- (i) X_1, \dots, X_n are IID or SSE
- (ii) $N_{[\cdot]}(\mathcal{M}, \varepsilon) < \infty$ for each $\varepsilon > 0$.

Then:

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{t=1}^n m(X_t, \theta) - E[m(X_t, \theta)] \right| = o_{a.s.}(1). \quad (10)$$

Proof. Take any rational $\varepsilon > 0$. For each $\theta \in \Theta$ there is a pair $l_{\varepsilon,i(\theta)}(X_t), u_{\varepsilon,i(\theta)}(X_t)$ with

$$l_{\varepsilon,i(\theta)}(X_t) \leq m(X_t; \theta) \leq u_{\varepsilon,i(\theta)}(X_t) \quad (11)$$

for all X_t and

$$\mathbb{E}[u_{\varepsilon,i(\theta)}(X_t) - l_{\varepsilon,i(\theta)}(X_t)] \leq \varepsilon \quad (12)$$

and $i(\theta) \in \{1, \dots, N_{[\cdot]}(\mathcal{M}, \varepsilon)\}$. Therefore, for each $\theta \in \Theta$ we have

$$\begin{aligned} Q_n(\theta) - Q(\theta) &= \frac{1}{n} \sum_{t=1}^n m(X_t; \theta) - \mathbb{E}[m(X_t; \theta)] \\ &\leq \frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i(\theta)}(X_t) - \mathbb{E}[m(X_t; \theta)] \end{aligned} \quad (13)$$

$$= \frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i(\theta)}(X_t) - \mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] + (\mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] - \mathbb{E}[m(X_t; \theta)]) \quad (14)$$

$$\leq \frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i(\theta)}(X_t) - \mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] + \varepsilon \quad (15)$$

because $\mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] - \mathbb{E}[m(X_t; \theta)] \leq \mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] - \mathbb{E}[l_{\varepsilon,i(\theta)}(X_t)] \leq \varepsilon$.

Taking the sup over $\theta \in \Theta$:

$$\sup_{\theta \in \Theta} (Q_n(\theta) - Q(\theta)) \leq \sup_{\theta \in \Theta} \left(\frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i(\theta)}(X_t) - \mathbb{E}[u_{\varepsilon,i(\theta)}(X_t)] \right) + \varepsilon \quad (16)$$

$$\leq \max_{1 \leq i \leq N} \left(\frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i}(X_t) - \mathbb{E}[u_{\varepsilon,i}(X_t)] \right) + \varepsilon \quad (17)$$

where $N = N_{[\cdot]}(\mathcal{M}, \varepsilon)$. Applying the SLLN or Ergodic Theorem yields

$$\frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i}(X_t) - \mathbb{E}[u_{\varepsilon,i}(X_t)] \rightarrow_{a.s.} 0 \quad (18)$$

for each $1 \leq i \leq N_{[\cdot]}(\mathcal{M}, \varepsilon)$, and so

$$\max_{1 \leq i \leq N} \left(\frac{1}{n} \sum_{t=1}^n u_{\varepsilon,i}(X_t) - \mathbb{E}[u_{\varepsilon,i}(X_t)] \right) \rightarrow_{a.s.} 0. \quad (19)$$

Therefore,

$$\sup_{\theta \in \Theta} (Q_n(\theta) - Q(\theta)) \leq \varepsilon + o_{a.s.}(1). \quad (20)$$

A similar argument with the lower bracket gives us

$$\inf_{\theta \in \Theta} (Q_n(\theta) - Q(\theta)) \geq -\varepsilon + o_{a.s.}(1). \quad (21)$$

Combining the preceding two inequalities, we obtain

$$\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| \leq \varepsilon + o_{a.s.}(1). \quad (22)$$

By definition of almost sure convergence, this means that there exists a set $S_\varepsilon \in \mathcal{F}$ with $\mathbb{P}(S_\varepsilon) = 1$ such that:

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} |Q_n(\theta; \omega) - Q(\theta)| \leq \varepsilon \quad (23)$$

for all $\omega \in S_\varepsilon$. Take $S = \cap_{\varepsilon \in \mathbb{Q}_+} S_\varepsilon$ where \mathbb{Q}_+ is the set of positive rational numbers. Then $\mathbb{P}(S) = 1$ and for each rational $\varepsilon > 0$ we have

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} |Q_n(\theta; \omega) - Q(\theta)| \leq \varepsilon \quad (24)$$

for all $\omega \in S$. As $\varepsilon \in \mathbb{Q}_+$ is arbitrary, we have shown that:

$$\lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} |Q_n(\theta; \omega) - Q(\theta)| = 0 \quad (25)$$

for all $\omega \in S$. ■

How do we show a collection of functions has finite bracketing numbers? The following result uses compactness of Θ , continuity, and dominance assumptions.

Lemma 3. *Let the following hold:*

(i) Θ is compact

(ii) $m(X_t; \theta)$ is continuous in θ for all X_t

(iii) $\mathbb{E}[\sup_{\theta \in \Theta} |m(X_t; \theta)|] < \infty$.

Then: $\mathcal{M} = \{m(X_t, \theta) : \theta \in \Theta\}$ has finite bracketing numbers. If, in addition,

(iv) X_1, \dots, X_n are IID or SSE, then: $\sup_{\theta \in \Theta} |\frac{1}{n} \sum_{t=1}^n m(X_t, \theta) - \mathbb{E}[m(X_t, \theta)]| \rightarrow_{a.s.} 0$.

Proof. Fix any $\delta > 0$. As Θ is compact, we can cover Θ with finitely many open balls of radius δ centered at $\theta_1, \dots, \theta_J$. For each $j = 1, \dots, J$, define

$$l_{\delta,j}(\cdot) = \inf_{\theta \in \Theta: \|\theta - \theta_j\| \leq \delta} m(\cdot; \theta) \quad \text{and} \quad u_{\delta,j}(\cdot) = \sup_{\theta \in \Theta: \|\theta - \theta_j\| \leq \delta} m(\cdot; \theta), \quad (26)$$

so that $l_{\delta,j}(\cdot) \leq m(\cdot; \theta) \leq u_{\delta,j}(\cdot)$ holds for each θ with $\|\theta - \theta_j\| \leq \delta$. Note that the inf and sup are always finite by (i) and (ii).

Let $\varepsilon(\delta) = \max_{1 \leq j \leq J} \mathbb{E}[u_{\delta,j}(X_t) - l_{\delta,j}(X_t)]$. We have shown that $N_{[\cdot]}(\varepsilon(\delta), \mathcal{M}) \leq J < \infty$. It remains to show that $\varepsilon(\delta) \rightarrow 0$ as $\delta \rightarrow 0$. This will ensure that for any $\epsilon > 0$ we can choose a δ such that $\varepsilon(\delta) \leq \epsilon$, and hence that $N_{[\cdot]}(\epsilon, \mathcal{M}) < \infty$ for each $\epsilon > 0$.

Let $M_\delta(\cdot) = \max_{1 \leq j \leq J} (u_{\delta,j}(\cdot) - l_{\delta,j}(\cdot))$. We may use (i) and (ii) to deduce that $M_\delta(X_t) \rightarrow 0$ as $\delta \rightarrow 0$ for each X_t (Exercise: use the fact that a continuous function on a compact set is uniformly continuous to show this formally). Also notice that $|M_\delta(X_t)| \leq 2 \sup_{\theta \in \Theta} |m(X_t; \theta)|$. Then by (iii) we may apply the dominated convergence theorem to obtain:

$$\lim_{\delta \rightarrow 0} \varepsilon(\delta) \leq \lim_{\delta \rightarrow 0} \mathbb{E}[M_\delta(X_t)] = \mathbb{E}[\lim_{\delta \rightarrow 0} M_\delta(X_t)] = 0 \quad (27)$$

as required. ■

Combining Lemmas 1, 2 and 3 gives the following consistency result.

Theorem 2 (Consistency of M-estimators). *Let the following hold:*

- (i) X_1, \dots, X_n are IID or SSE
- (ii) Θ is compact
- (iii) $m(X_t; \theta)$ is continuous in θ for all X_t
- (iv) $\mathbb{E}[\sup_{\theta \in \Theta} |m(X_t; \theta)|] < \infty$
- (v) $Q(\theta_0) > Q(\theta)$ for all $\theta \in \Theta$ with $\theta \neq \theta_0$.

Then: $\hat{\theta} \rightarrow_p \theta_0$ as $n \rightarrow \infty$.

Proof. By Theorem 1 we just need to verify “clean maximum” and “uniform convergence”.

We use Lemma 1 to verify “clean maximum”. By conditions (ii) and (v), it is enough to show that Q is continuous under the stated conditions. To verify continuity of Q , take any $\theta^* \in \Theta$ and let $(\theta_n)_{n \in \mathbb{N}} \subset \Theta$ be a sequence such that $\|\theta_n - \theta^*\| \rightarrow 0$ as $n \rightarrow \infty$. By condition (iii) we know that $\lim_{n \rightarrow \infty} m(X_t; \theta_n) = m(X_t; \theta^*)$ for all X_t . Then by condition (iv) we may apply the dominated convergence theorem to deduce

$$\lim_{n \rightarrow \infty} Q(\theta_n) = \lim_{n \rightarrow \infty} \mathbb{E}[m(X_t; \theta_n)] = \mathbb{E}[\lim_{n \rightarrow \infty} m(X_t; \theta_n)] = \mathbb{E}[m(X_t; \theta^*)] = Q(\theta^*), \quad (28)$$

which verifies continuity of Q . Therefore “clean maximum” holds.

Conditions (ii)–(iv) give finite bracketing numbers by Lemma 3. Moreover, $\mathcal{M} \subset L^1$ by (iv). This, together with (i), gives “uniform convergence” by Lemma 2. ■

3.2 Consistency of GMM Estimators

We're going to apply Theorem 1 to establish consistency of the GMM estimator. This requires verifying “clean maximum” and “uniform convergence”.

To apply Lemma 2, we need some notation. Write

$$g(X_t; \theta) = \begin{pmatrix} g_1(X_t; \theta) \\ g_2(X_t; \theta) \\ \vdots \\ g_K(X_t; \theta) \end{pmatrix}. \quad (29)$$

Then with this notation,

$$g_n(\theta) - g(\theta) = \begin{pmatrix} \frac{1}{n} \sum_{t=1}^n g_1(X_t; \theta) - E[g_1(X_t; \theta)] \\ \frac{1}{n} \sum_{t=1}^n g_2(X_t; \theta) - E[g_2(X_t; \theta)] \\ \vdots \\ \frac{1}{n} \sum_{t=1}^n g_K(X_t; \theta) - E[g_K(X_t; \theta)] \end{pmatrix}. \quad (30)$$

We will use Lemma 2 to ensure that each entry of $g_n(\theta) - g(\theta)$ converges in probability to zero (uniformly in θ), and hence $\sup_{\theta \in \Theta} \|g_n(\theta) - g(\theta)\| \rightarrow_p 0$. Let $\mathcal{G} = \{g_k(\cdot; \theta) : \theta \in \Theta, 1 \leq k \leq K\}$.

Theorem 3 (Consistency of GMM estimators). *Let the following hold:*

- (i) X_1, \dots, X_n are IID or SSE
 - (ii) Θ is compact
 - (iii) $g(\theta)$ is continuous
 - (iv) $\widehat{W} \rightarrow_p W$ where W is positive definite and symmetric
 - (v) $g(\theta) = 0$ if and only if $\theta = \theta_0$
 - (vi) \mathcal{G} has finite bracketing numbers.
- Then: $\hat{\theta} \rightarrow_p \theta_0$ as $n \rightarrow \infty$.

Proof. We verify the conditions of Theorem 1.

Continuity of $Q(\theta)$ follows from continuity of $g(\theta)$ and positive-definiteness of W . Therefore “clean maximum” holds by Lemma 1 (under conditions (ii)–(v)).

We now verify “uniform convergence”. Step 1: we show $\sup_{\theta \in \Theta} \|g_n(\theta) - g(\theta)\| \rightarrow_p 0$. Assumption (vi) implies that each of the K component functions in $g(X_t; \theta)$ has finite bracketing numbers. We may then apply Lemma 2 to deduce that each entry of $g_n(\theta) - g(\theta)$ converges in probability to zero (uniformly in θ), and hence

$$\sup_{\theta \in \Theta} \|g_n(\theta) - g(\theta)\| \rightarrow_p 0. \quad (31)$$

Before proceeding, we note that as g is continuous and Θ is compact, we also have:

$$\sup_{\theta \in \Theta} \|g(\theta)\| < \infty. \quad (32)$$

Combining (31) and (32) gives:

$$\sup_{\theta \in \Theta} \|g_n(\theta)\| \leq \sup_{\theta \in \Theta} \|g_n(\theta) - g(\theta)\| + \sup_{\theta \in \Theta} \|g(\theta)\| = o_p(1) + \sup_{\theta \in \Theta} \|g(\theta)\| = O_p(1). \quad (33)$$

Step 2: we show $\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| \rightarrow_p 0$. Adding and subtracting terms:

$$2Q(\theta) - 2Q_n(\theta) = g_n(\theta)' \widehat{W} g_n(\theta) - g(\theta)' W g(\theta) \quad (34)$$

$$= g_n(\theta)' (\widehat{W} - W) g_n(\theta) + g_n(\theta)' W g_n(\theta) - g(\theta)' W g(\theta) \quad (35)$$

$$= g_n(\theta)' (\widehat{W} - W) g_n(\theta) + (g_n(\theta) - g(\theta))' W (g_n(\theta) + g(\theta)). \quad (36)$$

Notice that for any K -vectors x, y and $K \times K$ matrix A we have

$$|x' A y| \leq \|x\| \|y\| \|A\| \quad (37)$$

where $\|x\|$ and $\|y\|$ are the Euclidean norms of x and y and $\|A\|$ is the spectral norm (largest singular value) of A . Applying the triangle inequality then inequality (37) to (36) yields:

$$\begin{aligned} \sup_{\theta \in \Theta} 2|Q_n(\theta) - Q(\theta)| &\leq \sup_{\theta \in \Theta} |g_n(\theta)' (\widehat{W} - W) g_n(\theta)| \\ &\quad + \sup_{\theta \in \Theta} |(g_n(\theta) - g(\theta))' W (g_n(\theta) + g(\theta))| \end{aligned} \quad (38)$$

$$\begin{aligned} &\leq \underbrace{\left(\sup_{\theta \in \Theta} \|g_n(\theta)\| \right)^2}_{=O_p(1) \text{ by (33)}} \times \underbrace{\|\widehat{W} - W\|}_{=o_p(1) \text{ by (iv)}} \\ &\quad + \sup_{\theta \in \Theta} \underbrace{\|(g_n(\theta) - g(\theta))\|}_{=o_p(1) \text{ by (31)}} \times \underbrace{\sup_{\theta \in \Theta} \|g_n(\theta) + g(\theta)\|}_{=O_p(1) \text{ by (32) and (33)}} \times \|W\| \end{aligned} \quad (39)$$

$$= O_p(1) \times o_p(1) + o_p(1) \times O_p(1) \times \text{constant} = o_p(1), \quad (40)$$

which verifies “uniform convergence”. ■

3.3 Consistency of SMM Estimators

Consistency for SMM requires special treatment because of the additional noise introduced by the simulation draws. Let's suppose that the simulated data $X_1^\theta, \dots, X_m^\theta$ are generated as functions of

i.i.d. draws $\varepsilon_1, \dots, \varepsilon_m$ which represent the “shocks” used to simulate the data. That is,

$$X_s^\theta = a(\varepsilon_s, \theta) \quad (41)$$

for each $1 \leq s \leq m$ and each $\theta \in \Theta$. We expand the probability space to jointly accommodate the true data X_1, \dots, X_n and the simulated draws $\varepsilon_1, \dots, \varepsilon_m$.¹ All probability statements we make in reference to SMM are to be understood with respect to the joint probability law of the data and simulated draws. As the sample size n gets large, we will be taking $m \rightarrow \infty$ also. If we don’t, the simulation error will eventually dominate and the SMM estimator will not converge.

We again establish consistency by verifying “clean maximum” and “uniform convergence”.

Theorem 4 (Consistency of SMM estimators). *Let the following hold:*

- (i) Θ is compact
 - (ii) $\gamma(\theta)$ is continuous in θ
 - (iii) $\sup_{\theta \in \Theta} \|\gamma_m(\theta) - \gamma(\theta)\| = o_p(1)$
 - (iv) $g_n \rightarrow_p g_0$ and $\widehat{W} \rightarrow_p W$ where W is positive definite and symmetric
 - (v) $\gamma(\theta) = g_0$ if and only if $\theta = \theta_0$.
- Then: $\hat{\theta} \rightarrow_p \theta_0$ as $n \rightarrow \infty$.

Note that in (iii) we explicitly assume the simulated moments converge (uniformly) to the moment function $\gamma(\theta)$ as the number of simulations increases. This can be verified under more primitive conditions by applying Lemma 2, substituting $\varepsilon_1, \dots, \varepsilon_m$ for X_1, \dots, X_n and $\gamma(a(\varepsilon_s, \theta); \theta)$ for $m(X_t, \theta)$.

Proof. By Theorem 1 we just need to verify “clean maximum” and “uniform convergence”.

We use Lemma 1 to verify “clean maximum”, noting Θ is compact (by (i)), $Q(\theta)$ is continuous (by (ii) and finiteness of W), and $Q(\theta_0) > Q(\theta)$ for any $\theta \neq \theta_0$ (by (v) and positive-definiteness of W).

We verify “uniform convergence” by similar arguments to the proof of Lemma 3. Adding and subtracting terms:

$$2Q(\theta) - 2Q_n(\theta) = (g_n - \gamma_m(\theta))' \widehat{W} (g_n - \gamma_m(\theta)) - (g_0 - \gamma(\theta))' W (g_0 - \gamma(\theta)) \quad (42)$$

$$\begin{aligned} &= (g_n - \gamma_m(\theta))' (\widehat{W} - W) (g_n - \gamma_m(\theta)) \\ &\quad + (g_n - g_0 + \gamma(\theta) - \gamma_m(\theta))' W (g_n - \gamma_m(\theta) + g_0 - \gamma(\theta)). \end{aligned} \quad (43)$$

Conditions (iii) and (iv) imply that

$$\sup_{\theta \in \Theta} \|g_n - g_0 + \gamma(\theta) - \gamma_m(\theta)\| \leq \|g_n - g_0\| + \sup_{\theta \in \Theta} \|\gamma_m(\theta) - \gamma(\theta)\| = o_p(1) \quad (44)$$

¹This is achieved by joining the σ -fields of the two and using the fact that the simulation draws are totally independent of the data.

and, moreover,

$$\sup_{\theta \in \Theta} \|g_n - \gamma_m(\theta)\| \leq \sup_{\theta \in \Theta} \|g_0 - \gamma(\theta)\| + \sup_{\theta \in \Theta} \|g_n - g_0 + \gamma(\theta) - \gamma_m(\theta)\| = O_p(1) \quad (45)$$

because $\sup_{\theta \in \Theta} \|g_0 - \gamma(\theta)\| < \infty$ by conditions (i) and (ii). Applying the triangle inequality then inequality (37) to (43), we obtain

$$\begin{aligned} \sup_{\theta \in \Theta} 2|Q_n(\theta) - Q(\theta)| &\leq \underbrace{\left(\sup_{\theta \in \Theta} \|g_n - \gamma_m(\theta)\| \right)^2}_{=O_p(1) \text{ by (45)}} \times \underbrace{\|\widehat{W} - W\|}_{=o_p(1) \text{ by (iv)}} \\ &\quad + \underbrace{\sup_{\theta \in \Theta} \|g_n - g_0 + \gamma(\theta) - \gamma_m(\theta)\|}_{=o_p(1) \text{ by (44)}} \\ &\quad \times \underbrace{\sup_{\theta \in \Theta} \|g_n - \gamma_m(\theta) + g_0 - \gamma(\theta)\| \times \|W\|}_{=O_p(1) \text{ by (45)}} \end{aligned} \quad (46)$$

$$= O_p(1) \times o_p(1) + o_p(1) \times O_p(1) \times \text{constant} = o_p(1), \quad (47)$$

which verifies “uniform convergence”. ■

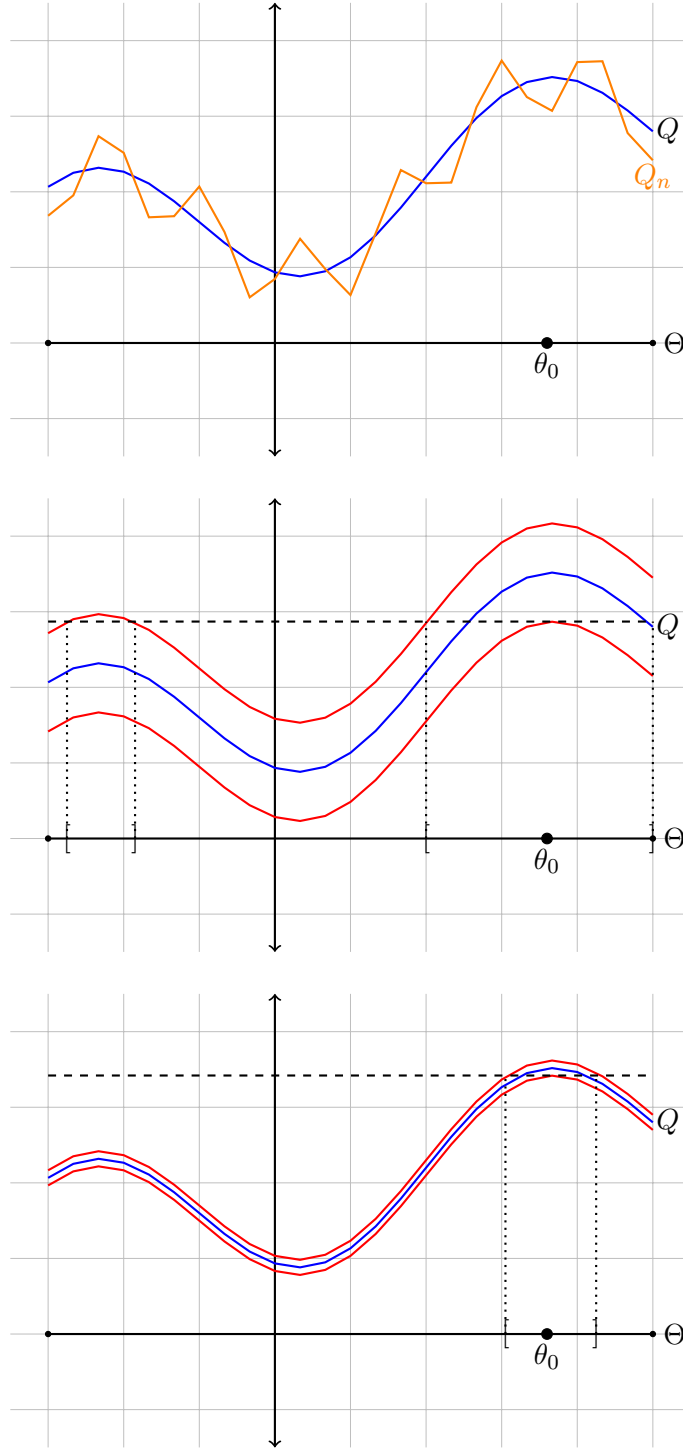


Figure 1: Consistency of Extremum Estimators. When $Q_n(\theta)$ lies uniformly in $[Q(\theta) - \epsilon, Q(\theta) + \epsilon]$ we know that $\hat{\theta}$ must be in the set $\{\theta : Q(\theta) \geq Q(\theta_0) - \epsilon\}$. Provided clean maximum holds, this set becomes a shrinking interval around θ_0 as ϵ decreases.

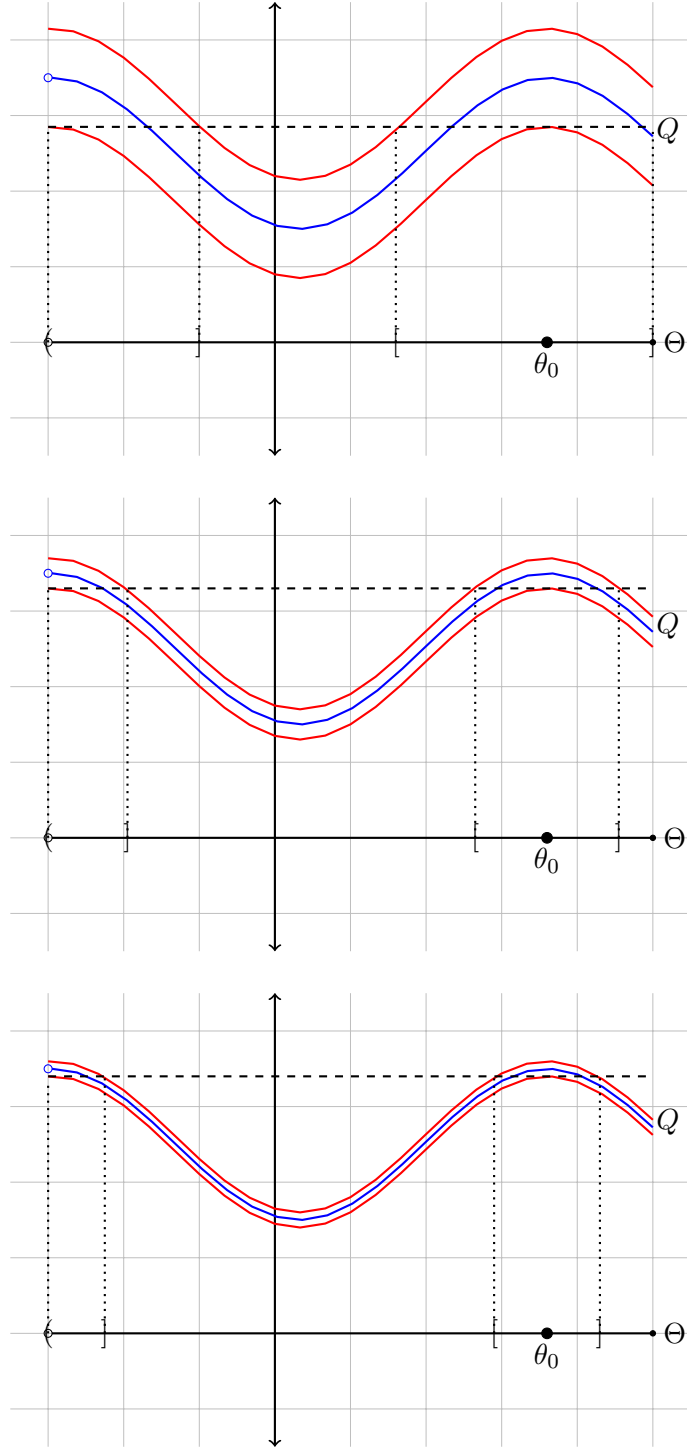


Figure 2: Necessity of Clean Maximum.

T. Christensen

1 Asymptotic Normality of Extremum Estimators

1.1 Preliminaries

The following three preliminary results will be used repeatedly:

Theorem 1 (Slutsky). *Let $X_n \rightarrow_p c$ where c is a constant and let $f(\cdot)$ be a function that is continuous at c . Then $f(X_n) \rightarrow_p f(c)$.*

Theorem 2 (Continuous mapping). *Let $X_n \rightarrow_d X$ where X_n and X are random vectors of length $k \geq 1$ and let $f : \mathbb{R}^k \rightarrow \mathbb{R}^q$ be continuous. Then $f(X_n) \rightarrow_d f(X)$.*

Theorem 3 (Mean-value). *Let $f : \mathbb{R}^k \rightarrow \mathbb{R}^q$ be continuous and differentiable on an open, convex neighborhood N of x_0 . Then for any $x \in N$:*

$$f(x) = f(x_0) + \frac{\partial f(\tilde{x})}{\partial x'}(x - x_0) \quad (1)$$

where the \tilde{x} lies in the segment between x and x_0 and may be different for each row of $\frac{\partial f(\tilde{x})}{\partial x'}$.

This lecture will establish asymptotic normality of extremum estimators. Formally, for a M, GMM, or SMM estimator $\hat{\theta}$, we will look at conditions under which

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, \Omega) \quad (2)$$

in which case we say that $\hat{\theta}$ is *root- n consistent* and *asymptotically normal* and its *asymptotic variance* is Ω . We will also discuss how to construct estimators $\hat{\Omega}$ of Ω . Results (2) should be thought of as an approximation. In practice we have a finite data set X_1, \dots, X_n consisting of n observations. If (2) is a good approximation, then the sampling distribution of $\hat{\theta}$ is approximately $N(\theta_0, \Omega/n)$. If $\hat{\Omega}$ is consistent for Ω , then the 95% confidence interval for the i th element $\theta_{0,i}$ of θ_0 , given by

$$\hat{\theta}_i \pm 1.96 \sqrt{\frac{\hat{\Omega}_{ii}}{n}}, \quad (3)$$

where $\hat{\Omega}_{ii}$ is the i th diagonal entry of $\hat{\Omega}$, should contain $\theta_{0,i}$ with probability approximately 0.95 across repeated samples.

We assume that the extremum estimator $\hat{\theta}$ solves the approximate first-order condition

$$o_p(n^{-1/2}) = \frac{\partial Q_n(\hat{\theta})}{\partial \theta}. \quad (4)$$

The interpretation of the left-hand side is that there may be a small approximation error but this does not affect the asymptotic behavior of our estimator to first order. Multiplying by \sqrt{n} and taking a mean-value expansion around θ_0 yields

$$o_p(1) = H_n \sqrt{n}(\hat{\theta} - \theta_0) - Z_n \quad (5)$$

where H_n is a symmetric $p \times p$ matrix whose probability limit, denoted H , is invertible, and

$$Z_n \rightarrow_d N(0, \Sigma). \quad (6)$$

It follows by the continuous mapping theorem (CMT) that

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, H^{-1} \Sigma H^{-1}). \quad (7)$$

The main regularity conditions we need are (i) that $\hat{\theta}$ is consistent (and hence within a neighborhood of θ_0 with probability approaching one, allowing us to perform a mean value expansion), (ii) that θ_0 is in the interior of the parameter space (again, so we can perform a mean value expansion) and (iii) that the criterions are “smooth” in θ and enough moments exist.

The expressions for H_n , Z_n , H , and Σ will differ depending on whether we’re using M, GMM, or SMM estimators. We discuss these in more detail in the following sections. What is generally true is that H will represent the *curvature* of population the objective function Q at θ_0 , which is related to how cleanly the population criterion identifies θ_0 . The term Σ is a measure of *sampling uncertainty*. A unifying theme that emerges is that *efficient* estimators balance curvature against sampling uncertainty in the sense that $H = \Sigma$.

1.2 M-Estimators

Define

$$s(X_t; \theta) = \frac{\partial m(X_t; \theta)}{\partial \theta}, \quad (8)$$

$$D(X_t; \theta) = \frac{\partial^2 m(X_t; \theta)}{\partial \theta \partial \theta'}, \quad (9)$$

$$Z_n = \frac{1}{\sqrt{n}} \sum_{t=1}^n s(X_t; \theta_0), \quad (10)$$

$$H = -E[D(X_t; \theta_0)], \quad (11)$$

where $s(X_t; \theta)$ and Z_n are $p \times 1$ vectors and $D(X_t; \theta)$ and S are $p \times p$. The s is for *score* and the H is for *Hessian*.

Theorem 4. *Let the following hold:*

- (i) X_1, \dots, X_n are IID or SSE
- (ii) $\hat{\theta} \rightarrow_p \theta_0$
- (iii) θ_0 is in the interior of Θ
- (iv) $m(X_t; \theta)$ is twice continuously differentiable in θ for any X_t
- (v) $Z_n \rightarrow_d N(0, \Sigma)$ with Σ positive definite
- (vi) $E[\sup_{\theta \in N} \|D(X_t; \theta)\|] < \infty$ for some convex, compact neighborhood N of θ_0
- (vii) H is positive definite.

Then: $\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, H^{-1}\Sigma H^{-1})$ as $n \rightarrow \infty$.

Proof. A mean-value expansion (see Theorem 3) of (4) yields

$$o_p(1) = Z_n + \underbrace{\frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'}}_{=:-H_n} \sqrt{n}(\hat{\theta} - \theta_0) \quad (12)$$

where

$$(H_n)_{ij} = \frac{1}{n} \sum_{t=1}^n D_{ij}(X_t; \tilde{\theta}) \quad (13)$$

for some $\tilde{\theta}$ in the segment between θ and θ_0 (this is formally justified by (ii), (iii) and (iv)).

The function $\theta \mapsto D_{ij}(X_t; \theta)$ is continuous and dominated by assumptions (iv) and (vi). Therefore $\mathcal{D}_{ij} = \{D_{ij}(X_t; \theta) : \theta \in N\} \subset L^1$ has finite bracketing numbers by Lecture 7, Lemma 3, so we have by the ULLN that

$$\sup_{\theta \in N} \left| \frac{1}{n} \sum_{t=1}^n D_{ij}(X_t; \theta) - E[D_{ij}(X_t; \theta)] \right| \rightarrow_p 0 \quad (14)$$

for each $1 \leq i, j \leq p$, where D_{ij} denotes the (i, j) element of D . Also notice that $E[D_{ij}(X_t; \theta)]$ is continuous at θ_0 , which can be shown using the dominated convergence theorem. Hence,

$$\sup_{\theta: \|\theta - \theta_0\| \leq \epsilon} |E[D_{ij}(X_t; \theta)] - E[D_{ij}(X_t; \theta_0)]| \rightarrow 0 \quad \text{as } \epsilon \rightarrow 0 \quad (15)$$

for each $1 \leq i, j \leq p$.

Because $\|\hat{\theta} - \theta_0\| = o_p(1)$ (by (ii)), we combine (14) and (15) to obtain:

$$\frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} \rightarrow_p -H. \quad (16)$$

The result now follows by the CMT. ■

The matrix Σ will depend on the sampling scheme and structure of the model, as we discuss below.

Example 1.1. Linear Regression.

Suppose that (Y_t, X_t) , $t = 1, \dots, n$ is IID. You wish to estimate the coefficients θ_0 for the best linear predictor $X_t' \theta_0$ of Y_t given X_t , where θ_0 solves

$$\min_{\theta} E[(Y_t - X_t' \theta)^2]. \quad (17)$$

The least-squares estimator $\hat{\theta}$ of θ_0 is a M-estimator with

$$m((X_t, Y_t); \theta) = -\frac{1}{2}(Y_t - X_t' \theta)^2. \quad (18)$$

Let $u_t = Y_t - X_t' \theta_0$. Note by the FOC for θ_0 that $E[X_t u_t] = 0$. Here we have

$$s((X_t, Y_t); \theta) = X_t(Y_t - X_t' \theta), \quad (19)$$

$$D((X_t, Y_t); \theta) = -X_t X_t', \quad (20)$$

$$Z_n = \frac{1}{\sqrt{n}} \sum_{t=1}^n X_t u_t, \quad (21)$$

$$H = E[X_t X_t']. \quad (22)$$

Provided $E[\|X_t u_t\|^2] < \infty$, then by the CLT we have $Z_n \rightarrow_d N(0, \Sigma)$ with $\Sigma = E[u_t^2 X_t X_t']$. Finally, if H is finite and positive definite (i.e., there is “no multicollinearity”) then we have

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, E[X_t X_t']^{-1} E[u_t^2 X_t X_t'] E[X_t X_t']^{-1}). \quad (23)$$

Note the asymptotic covariance matrix is the usual Eicker–Huber–White “heteroskedasticity-robust” asymptotic covariance matrix. □

Example 1.2. Maximum Likelihood.

Suppose that X_1, \dots, X_n are IID with density $f(\cdot; \theta_0)$. You know the parametric family f but not θ_0 . You estimate θ_0 using maximum likelihood. We have

$$m(X_t; \theta) = \log f(X_t; \theta), \quad s(X_t; \theta) = \frac{\partial \log f(X_t; \theta)}{\partial \theta}, \quad D(X_t; \theta) = \frac{\partial^2 \log f(X_t; \theta)}{\partial \theta \partial \theta'} \quad (24)$$

and

$$Z_n = \frac{1}{\sqrt{n}} \sum_{t=1}^n \frac{\partial \log f(X_t; \theta)}{\partial \theta}, \quad H = -E \left[\frac{\partial^2 \log f(X_t; \theta_0)}{\partial \theta \partial \theta'} \right]. \quad (25)$$

Note that $s(X_t; \theta_0)$ has mean zero by virtue of the FOC for θ_0 , irrespective of whether or not the model is correctly specified. As the data are IID, if $E[\|s(X_t; \theta_0)\|^2] < \infty$ then by CLT we have $Z_n \rightarrow_d N(0, \Sigma)$ with

$$\Sigma = E \left[\left(\frac{\partial \log f(X_t; \theta)}{\partial \theta} \right) \left(\frac{\partial \log f(X_t; \theta)}{\partial \theta} \right)' \right]. \quad (26)$$

Hence, the general “sandwich” asymptotic covariance $H^{-1} \Sigma H^{-1}$ holds under correct specification and misspecification.

Differentiating

$$\int f(u; \theta) du = 1 \quad (27)$$

twice respect to θ and exchanging differentiation and integration, we get

$$\int \frac{\partial^2 f(u; \theta)}{\partial \theta \partial \theta'} du = \mathbf{0}_{p \times p}. \quad (28)$$

But notice

$$\frac{\partial \log f(u; \theta)}{\partial \theta} = \frac{1}{f(u; \theta)} \frac{\partial f(u; \theta)}{\partial \theta} \quad (29)$$

$$\frac{\partial^2 \log f(u; \theta)}{\partial \theta \partial \theta'} = \frac{1}{f(u; \theta)} \frac{\partial^2 f(u; \theta)}{\partial \theta \partial \theta'} - \frac{1}{f(u; \theta)^2} \frac{\partial f(u; \theta)}{\partial \theta} \frac{\partial f(u; \theta)}{\partial \theta'} \quad (30)$$

$$= \frac{1}{f(u; \theta)} \frac{\partial^2 f(u; \theta)}{\partial \theta \partial \theta'} - \frac{\partial \log f(u; \theta)}{\partial \theta} \frac{\partial \log f(u; \theta)}{\partial \theta'}, \quad (31)$$

hence

$$\frac{\partial^2 f(u; \theta)}{\partial \theta \partial \theta'} = \left(\frac{\partial^2 \log f(u; \theta)}{\partial \theta \partial \theta'} + \frac{\partial \log f(u; \theta)}{\partial \theta} \frac{\partial \log f(u; \theta)}{\partial \theta'} \right) f(u; \theta). \quad (32)$$

Substituting into (28) and setting $\theta = \theta_0$ gives:

$$\int \left(\frac{\partial^2 \log f(u; \theta_0)}{\partial \theta \partial \theta'} + \frac{\partial \log f(u; \theta_0)}{\partial \theta} \frac{\partial \log f(u; \theta_0)}{\partial \theta'} \right) f(u; \theta_0) du = \mathbf{0}_{p \times p}. \quad (33)$$

If the model is correctly specified, then integrating with respect to $f(\cdot; \theta_0)$ is the same as taking

expectation. So in this case we have the *information equality*:

$$\mathbb{E} \left[\frac{\partial^2 \log f(X_t; \theta_0)}{\partial \theta \partial \theta'} + \frac{\partial \log f(X_t; \theta_0)}{\partial \theta} \frac{\partial \log f(X_t; \theta_0)}{\partial \theta'} \right] = \mathbf{0}_{p \times p}. \quad (34)$$

Equivalently,

$$H = \Sigma \quad (35)$$

in our earlier notation. If this equality holds, then we call the common value of H and Σ the *information matrix*, and denote it $\mathbb{I}(\theta_0)$. So in this case,

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, \mathbb{I}(\theta_0)^{-1}). \quad (36)$$

Intuitively, $\mathbb{I}(\theta_0)$ is larger when the criterion function is more curved at θ_0 , reflecting the fact that the model contains more “information” about θ_0 . If $\mathbb{I}(\theta_0)$ is large, then $\mathbb{I}(\theta_0)^{-1}$ is small, hence the asymptotic variance is small. \square

1.2.1 IID Data

When X_1, \dots, X_n are IID and $\mathbb{E}[\|s(X_t; \theta_0)\|^2] < \infty$, the central limit theorem tells us

$$Z_n \rightarrow_d N(0, \Sigma) \quad \text{with} \quad \Sigma = \mathbb{E}[s(X_t; \theta_0)s(X_t; \theta_0)']. \quad (37)$$

1.2.2 SSE Data and MDS Scores

Let \mathcal{F}_t denote the σ -algebra generated by (X_t, X_{t-1}, \dots) where $(X_t)_{t \in \mathbb{Z}}$ are strictly stationary and ergodic.¹ We say that $(s(X_t; \theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ is a *martingale difference sequence* (MDS) if

$$\mathbb{E}[s(X_{t+1}; \theta_0) | \mathcal{F}_t] = 0 \quad (38)$$

with probability 1 for each $t \in \mathbb{Z}$. Hence, $s(X_{t+1}; \theta_0)$ has conditional mean zero for almost every realization of (X_t, X_{t-1}, \dots) .

Let X_1, \dots, X_n be SSE, let $(s(X_t; \theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ be a MDS, and let $\mathbb{E}[\|s(X_t; \theta_0)\|^2] < \infty$. By the

¹We interpret X_t as the date- t observation of a stochastic process that was started in the infinite past. Even though our data are X_1, \dots, X_n , we can still condition on the σ -algebra \mathcal{F}_t representing all “information” in (X_t, X_{t-1}, \dots) .

martingale property, the covariance of $s(X_t; \theta_0)$ and $s(X_{t+k}; \theta_0)$ is necessarily zero for each $k \geq 1$:

$$\mathbb{E}[s(X_{t+k}; \theta_0)s(X_t; \theta_0)] \stackrel{(\text{LIE})}{=} \mathbb{E}[\mathbb{E}[s(X_{t+k}; \theta_0)s(X_t; \theta_0)' | \mathcal{F}_{t+k-1}]] \quad (39)$$

$$= \mathbb{E}[\underbrace{\mathbb{E}[s(X_{t+k}; \theta_0) | \mathcal{F}_{t+k-1}]}_{=0 \text{ by (38)}} s(X_t; \theta_0)'] = 0 \quad (40)$$

The central limit theorem for stationary and ergodic MDSs (Billingsley, 1961) yields

$$Z_n \rightarrow_d N(0, \Sigma) \quad \text{with} \quad \Sigma = \mathbb{E}[s(X_t; \theta_0)s(X_t; \theta_0)'] \quad (41)$$

just as in the IID case.

Example 1.3. Maximum Likelihood for Markov Processes.

Suppose we have a model for the transition distribution $f(Y_{t+1}|Y_t; \theta)$ of a first-order Markov process. Here we have

$$s(X_{t+1}; \theta_0) = \frac{\partial \log f(Y_{t+1}|Y_t; \theta_0)}{\partial \theta} \quad (42)$$

with $X_t = (Y_t, Y_{t-1})$. As $f(Y_{t+1}|Y_t; \theta)$ is a density for each θ , we must have

$$1 = \int f(u|Y_t; \theta) du \quad (43)$$

for each $\theta \in \Theta$. Differentiating both sides with respect to θ and interchanging integration and differentiation:

$$0 = \int \frac{\partial f(u|Y_t; \theta)}{\partial \theta} du \quad (44)$$

taking $\theta = \theta_0$ and multiplying and dividing by $f(\cdot|Y_t; \theta_0)$:

$$0 = \int \frac{\partial f(u|Y_t; \theta_0)}{\partial \theta} \frac{1}{f(u|Y_t; \theta_0)} f(u|Y_t; \theta_0) du = \int \frac{\partial \log f(u|Y_t; \theta_0)}{\partial \theta} f(u|Y_t; \theta_0) du. \quad (45)$$

When the model is correctly specified, the term on the right is the conditional expectation of $s(X_{t+1}; \theta_0)$ given Y_t , hence

$$\mathbb{E}[s(X_{t+1}; \theta_0) | Y_t] = 0. \quad (46)$$

Because $X_t = (Y_t, Y_{t-1})$ and Y is a first-order Markov process, conditioning on Y_t is the same as conditioning on \mathcal{F}_t . Therefore $\mathbb{E}[s(X_{t+1}; \theta_0) | \mathcal{F}_t] = 0$.

If the model is misspecified, then we still have that $\mathbb{E}[s(X_{t+1}; \theta_0)] = 0$ from the FOC for θ_0 , but we no longer have the stronger condition that $\mathbb{E}[s(X_{t+1}; \theta_0) | Y_t] = 0$. \square

1.2.3 SSE but not MDS

If $s(X_t; \theta_0)$ is not a MDS, then the asymptotic variance Σ will depend on the autocovariances of $s(X_t; \theta_0)$. Let

$$C_j = E[s(X_t; \theta_0)s(X_{t-j}; \theta_0)'] \quad (47)$$

for each $j \in \mathbb{Z}$. Suppose that the sum $\sum_{j=-\infty}^{\infty} C_j$ converges. Then under some additional conditions, we have

$$Z_n \rightarrow_d N(0, \Sigma) \quad \text{with} \quad \Sigma = \sum_{j=-\infty}^{\infty} C_j = C_0 + \sum_{j=1}^{\infty} (C_j + C_j'). \quad (48)$$

The quantity Σ is referred to as the *long-run variance* of the stochastic process $(s(X_t; \theta_0))_{t \in \mathbb{Z}}$. Estimation of Σ is nontrivial in this case as we have to estimate infinitely many covariances with a finite amount of data. We'll explore this issue in greater detail when we study GMM estimation. The IID/MDS results correspond to the special case in which $C_j = 0$ for all $j \neq 0$.

1.3 GMM Estimators

Define

$$d(X_t; \theta) = \frac{\partial g(X_t; \theta)}{\partial \theta'}, \quad (49)$$

$$G_n(\theta) = \frac{1}{n} \sum_{t=1}^n d(X_t; \theta), \quad (50)$$

$$G(\theta) = E[d(X_t; \theta)], \quad (51)$$

$$G = G(\theta_0). \quad (52)$$

notice that $d(X_t; \theta)$, $G_n(\theta)$ and G are all $K \times p$. We also use the notation

$$\sqrt{n}g_n(\theta_0) = \frac{1}{\sqrt{n}} \sum_{t=1}^n g(X_t; \theta_0). \quad (53)$$

Notice $E[g(X_t; \theta_0)] = 0$, so we should expect that $\sqrt{n}g_n(\theta_0) \rightarrow_d N(0, S)$ by a central limit theorem.

The following result presents conditions for asymptotic normality of a GMM estimator $\hat{\theta}$.

Theorem 5. *Let the following hold:*

- (i) X_1, \dots, X_n are IID or SSE
- (ii) $\hat{\theta} \rightarrow_p \theta_0$
- (iii) $\widehat{W} \rightarrow_p W$ where W is positive definite and symmetric
- (iv) θ_0 is in the interior of Θ
- (v) $g(X_t; \theta)$ is continuously differentiable in θ for any X_t

(vi) $(\sqrt{n}g_n(\theta_0)) \rightarrow_d N(0, S)$ with S positive definite
(vii) $E[\sup_{\theta \in N} \|d(X_t; \theta)\|] < \infty$ for some convex, compact neighborhood N of θ_0
(viii) G has full column rank p — (a “local” identification condition).
Then: $\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, (G'WG)^{-1}G'WSWG(G'WG)^{-1})$.

Note that the asymptotic variance is of the form $\Omega = H^{-1}\Sigma H$ where now

$$H = G'WG, \quad (54)$$

$$\Sigma = G'WSWG. \quad (55)$$

Proof. A mean-value expansion (see Theorem 3) of (4) yields

$$o_p(n^{-1/2}) = \frac{\partial Q_n(\hat{\theta})}{\partial \theta} \quad (56)$$

$$= -G_n(\hat{\theta})' \widehat{W} g_n(\hat{\theta}) \quad (57)$$

$$= -G_n(\hat{\theta})' \widehat{W} \left(g_n(\theta_0) + G_n(\tilde{\theta})(\hat{\theta} - \theta_0) \right). \quad (58)$$

The mean-value expansion is formally justified whenever $\hat{\theta}$ is in an open neighborhood of θ_0 , which it is wpa1 by conditions (ii) and (iv). Again, $\tilde{\theta}$ is in the segment between $\hat{\theta}$ and θ_0 and is possibly different for each row of $G_n(\cdot)$.

We can therefore rewrite the above display:

$$o_p(1) = (G_n(\hat{\theta})' \widehat{W} G_n(\tilde{\theta})) \left(\sqrt{n}(\hat{\theta} - \theta_0) \right) - \left(-G_n(\hat{\theta})' \widehat{W} (\sqrt{n}g_n(\theta_0)) \right). \quad (59)$$

We show below that $G_n(\hat{\theta}) \rightarrow_p G$ and $G_n(\tilde{\theta}) \rightarrow_p G$. This, together with the fact that $\widehat{W} \rightarrow_p W$ (condition (iii)), implies that

$$(G_n(\hat{\theta})' \widehat{W} G_n(\tilde{\theta})) \rightarrow_p G'WG, \quad (60)$$

$$G_n(\hat{\theta})' \widehat{W} (\sqrt{n}g_n(\theta_0)) \rightarrow_d N(0, G'WSWG), \quad (61)$$

where the second line is by part (vi) and the CMT. The result now follows by the CMT.

It remains to show that $G_n(\tilde{\theta}) \rightarrow_p G$ and $G_n(\hat{\theta}) \rightarrow_p G$. We just prove the first; the second is similar. Notice that $\theta \mapsto d(X_t; \theta)$ is continuous by assumption (v) and dominated on N by assumption (vii). We therefore apply Lecture 7, Lemma 3 to obtain

$$\sup_{\theta \in N} |[G_n(\theta)]_{ij} - [G(\theta)]_{ij}| \rightarrow_p 0. \quad (62)$$

Also notice that $G(\theta)$ is continuous at θ_0 under assumptions (v) and (vii). Hence,

$$[G(\tilde{\theta})]_{ij} - [G(\theta_0)]_{ij} = o_p(1) \quad (63)$$

by Slutsky's theorem. So, whenever $\tilde{\theta} \in N$ (which it is wpa1), we have

$$|[G_n(\tilde{\theta})]_{ij} - [G(\theta_0)]_{ij}| \leq |[G_n(\tilde{\theta})]_{ij} - [G(\tilde{\theta})]_{ij}| + |[G(\tilde{\theta})]_{ij} - [G(\theta_0)]_{ij}| \quad (64)$$

$$\leq \underbrace{\sup_{\theta \in N} |[G_n(\theta)]_{ij} - [G(\theta)]_{ij}|}_{=o_p(1) \text{ by (62)}} + \underbrace{|[G(\tilde{\theta})]_{ij} - [G(\theta_0)]_{ij}|}_{=o_p(1) \text{ by (63)}} = o_p(1). \quad (65)$$

■

In *just identified models* ($K = p$) the choice of weighting matrix will not influence the asymptotic variance. This is because G is square and has full rank, and hence is invertible. In this case, the asymptotic variance simplifies:

$$(G'WG)^{-1}G'WSWG(G'WG)^{-1} = G^{-1}SG'^{-1} = (G'S^{-1}G)^{-1}. \quad (66)$$

In *over identified models* ($K > p$) the matrix G is not invertible so we do not have this simplification. So, in over-identified models the choice of weighting matrix will affect the asymptotic distribution (and hence the efficiency) of our estimators.

If we choose \widehat{W} so that $\widehat{W} \rightarrow_p S^{-1}$ then the asymptotic variance of $\hat{\theta}$ collapses to $(G'S^{-1}G)^{-1}$. This weighting is “optimal” because it leads to the smallest asymptotic variance:

$$(G'WG)^{-1}G'WSWG(G'WG)^{-1} \geq (G'S^{-1}G)^{-1} \quad (67)$$

for any positive definite matrix W , where the inequality is in the sense of nonnegative-definite matrices.

To see why this inequality must hold, note that we can write the left-hand side as $A'SA$ with $A' = (G'WG)^{-1}G'W$ and the right-hand side as $B'SB$ with $B' = (G'S^{-1}G)^{-1}G'S^{-1}$. But then $B'SA = B'SB$ so $B'S(A - B) = 0$. With this in mind, we may write the left-hand side as

$$A'SA = (B + A - B)'S(B + A - B) \quad (68)$$

$$= B'SB + (A - B)'S(A - B) + B'S(A - B) + (A - B)'SB \quad (69)$$

$$= B'SB + (A - B)'S(A - B) \quad (70)$$

$$\geq B'SB. \quad (71)$$

GMM estimators whose weighting matrix \hat{W} satisfies $\hat{W} \rightarrow_p S^{-1}$ are called *optimally-weighted* GMM estimators. Implementing optimally-weighted GMM typically consists of two steps:

1. Compute an estimator $\tilde{\theta}$ of θ_0 using a deterministic weighting matrix, such as the identity.
2. Estimate θ_0 again by maximizing $Q_n(\theta)$ with respect to θ , but this time using

$$\widehat{W} = \left(\frac{1}{n} \sum_{t=1}^n g(X_t; \tilde{\theta}) g(X_t; \tilde{\theta})' \right)^{-1} \quad (72)$$

for IID data or when $(g(X_t; \theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ is a MDS; otherwise, use the Newey-West estimator of the long-run variance of $(g(X_t; \theta_0))_{t \in \mathbb{Z}}$. The maximizing value $\hat{\theta}$ is our two-step optimally weighted estimator of θ_0 . Its asymptotic variance is $(G' S^{-1} G)^{-1}$.

Example 1.4. Linear IV Model.

Consider the linear IV model

$$Y_{1t} = Y_{2t}' \theta_0 + u_t, \quad E[u_t Z_t] = 0, \quad (73)$$

where Y_2 is an endogenous regressor and Z is a vector of instruments. Recall that the linear IV estimator is a GMM estimator with

$$g(X_t; \theta) = Z_t(Y_{1t} - Y_{2t}' \theta) \quad (74)$$

with $X_t = (Y_{1t}, Y_{2t}, Z_t)$. Suppose the data are IID. We have

$$G = Z_t Y_{2t}', \quad S = E[u_t^2 Z_t Z_t']. \quad (75)$$

In just-identified models the asymptotic variance of the linear IV estimator is

$$(E[Y_{2t} Z_t'] E[u_t^2 Z_t Z_t']^{-1} E[Z_t Y_{2t}'])^{-1}. \quad (76)$$

In over-identified models with asymptotic weight matrix W (e.g. $W = E[Z_t Z_t']^{-1}$ for 2SLS), the asymptotic variance is

$$(E[Y_{2t} Z_t'] W E[Z_t Y_{2t}'])^{-1} (E[Y_{2t} Z_t'] W E[u_t^2 Z_t Z_t'] W E[Z_t Y_{2t}']) (E[Y_{2t} Z_t'] W E[Z_t Y_{2t}'])^{-1}. \quad (77)$$

The smallest asymptotic variance is

$$(E[Y_{2t} Z_t'] E[u_t^2 Z_t Z_t']^{-1} E[Z_t Y_{2t}'])^{-1}, \quad (78)$$

which can be achieved using the weight matrix

$$\hat{W} = \left(\frac{1}{n} \sum_{t=1}^n \hat{u}_t^2 Z_t Z_t' \right)^{-1}, \quad (79)$$

where $\hat{u}_t = Y_{1t} - Y_{2t}' \tilde{\theta}$ for a consistent estimator $\tilde{\theta}$ of θ_0 . Note that the 2SLS estimator will only be efficient if $E[u_t^2 Z_t Z_t']^{-1} = c E[Z_t Z_t']^{-1}$ for some constant $c \in (0, \infty)$, such as when the errors u_t are conditionally homoskedastic in the sense that $E[u_t^2 | Z_t = z] = \sigma^2$. \square

1.4 SMM Estimators

Derivation of asymptotic normality for SMM is more involved. We just sketch the derivation here. Rigorous derivations are provided in Pakes and Pollard (1989) and Duffie and Singleton (1993). Notice that for any $\theta \in \Theta$:

$$Q_n(\theta) = -\frac{1}{2}(g_n - \gamma_m(\theta_0) + \gamma_m(\theta_0) - \gamma_m(\theta))' \widehat{W} (g_n - \gamma_m(\theta_0) + \gamma_m(\theta_0) - \gamma_m(\theta)) \quad (80)$$

$$\begin{aligned} &= Q_n(\theta_0) + (g_n - \gamma_m(\theta_0))' \widehat{W} (\gamma_m(\theta) - \gamma_m(\theta_0)) \\ &\quad - \frac{1}{2}(\gamma_m(\theta) - \gamma_m(\theta_0))' \widehat{W} (\gamma_m(\theta) - \gamma_m(\theta_0)). \end{aligned} \quad (81)$$

Suppose that the model is correctly specified, so

$$E[g(X_t)] = \int \gamma(X_t; \theta_0) f(X_t | \theta_0) dX_t =: \gamma(\theta_0). \quad (82)$$

We then have

$$\sqrt{n}(g_n - \gamma(\theta_0)) = \frac{1}{\sqrt{n}} \sum_{t=1}^n g(X_t) - E[g(X_t)] \rightarrow_d N(0, S) \quad (83)$$

where

$$S = E[(g(X_t) - \gamma(\theta_0))(g(X_t) - \gamma(\theta_0))'] \quad (84)$$

if X_1, \dots, X_n are IID or if $(g(X_t) - \gamma(\theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ is a MDS. Otherwise, S is the long-run variance of $(g(X_t) - \gamma(\theta_0))_{t \in \mathbb{Z}}$.

Also suppose that

$$\frac{n}{m} \rightarrow \tau \geq 0 \quad \text{as } n \rightarrow \infty. \quad (85)$$

If $\tau = 0$ then the number of simulations is growing faster than sample size. If $\tau \in (0, \infty)$ then the number of simulations is proportional to sample size. We do not allow for the possibility that $n/m \rightarrow \infty$, as this will mean the number of simulations is growing slower than sample size. In this case the asymptotics behavior of our estimator will be determined by the simulation draws and not by the data.

Then we have

$$\sqrt{n}(\gamma_m(\theta_0) - \gamma(\theta_0)) = \underbrace{\sqrt{\frac{n}{m}}}_{\rightarrow \sqrt{\tau}} \times \underbrace{\frac{1}{\sqrt{m}} \sum_{s=1}^m (\gamma(X_s^{\theta_0}; \theta_0) - \gamma(\theta_0))}_{\rightarrow_d N(0, S)} \quad (86)$$

$$\rightarrow_d N(0, \tau S). \quad (87)$$

As the actual data X_1, \dots, X_n and the simulated data $X_1^{\theta_0}, \dots, X_m^{\theta_0}$ are independent, we have the joint convergence

$$\sqrt{n} \begin{pmatrix} g_n - \gamma(\theta_0) \\ \gamma_m(\theta_0) - \gamma(\theta_0) \end{pmatrix} \rightarrow_d N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} S & 0 \\ 0 & \tau S \end{pmatrix} \right). \quad (88)$$

In view of the joint convergence, we may apply the CMT to obtain:

$$\sqrt{n}(g_n - \gamma_m(\theta_0)) \rightarrow_d N(0, (1 + \tau)S). \quad (89)$$

The factor τ represents the efficiency loss due to using simulation. Ideally $n/m \rightarrow 0$ in order to make estimates as efficient as possible, but it may be computationally costly to make m large.

Going back to the quadratic approximation (81), we now have

$$Q_n(\theta) = Q_n(\theta_0) + \left(\frac{1}{\sqrt{n}} Z_n \right)' \widehat{W}(\gamma_m(\theta) - \gamma_m(\theta_0)) - \frac{1}{2}(\gamma_m(\theta) - \gamma_m(\theta_0))' \widehat{W}(\gamma_m(\theta) - \gamma_m(\theta_0)) \quad (90)$$

where

$$Z_n = \sqrt{n}(g_n - \gamma_m(\theta_0)) \rightarrow_d N(0, (1 + \tau)S). \quad (91)$$

It remains to write the terms $\gamma_m(\theta) - \gamma_m(\theta_0)$ in terms of $(\theta - \theta_0)$. Under some regularity conditions, we treat γ_m and γ as “close” and may take a mean-value expansion of $\gamma(\theta)$ around θ_0 to obtain

$$Q_n(\theta) = Q_n(\theta_0) + \left(\frac{1}{\sqrt{n}} Z_n \right)' W\Gamma(\theta - \theta_0) - \frac{1}{2}(\theta - \theta_0)' \Gamma' W\Gamma(\theta - \theta_0) + o_p(\|\theta - \theta_0\|^2) \quad (92)$$

with $\Gamma = \frac{\partial \gamma(\theta_0)}{\partial \theta'}$. As $\hat{\theta}$ maximizes Q_n this with respect to θ , we obtain

$$\sqrt{n}(\hat{\theta} - \theta_0) = (\Gamma' W\Gamma)^{-1} \Gamma' W Z_n + o_p(1). \quad (93)$$

Therefore,

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, (1 + \tau)(\Gamma' W\Gamma)^{-1} \Gamma' W S W\Gamma (\Gamma' W\Gamma)^{-1}). \quad (94)$$

By analogy with GMM, we know that we can get the smallest variance if we set $W = S^{-1}$ (this is only for an over-identified model; in a just-identified model the W will cancel out). Choosing an

efficient weighting matrix is actually quite easy here: we can take \widehat{W} as the inverse of the sample covariance of $g(X_t)$ in the IID or MDS cases; otherwise, we take \widehat{W} as the inverse of the Newey-West estimator of the long-run variance of $(g(X_t))_{t \in \mathbb{Z}}$ under general weak dependence.

Doing so, we obtain:

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, (1 + \tau)(\Gamma' S^{-1} \Gamma)^{-1}). \quad (95)$$

This is still not efficient unless $\tau = 0$. If $\tau = 0$ then we get

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, (\Gamma' S^{-1} \Gamma)^{-1}) \quad (96)$$

where $(\Gamma' S^{-1} \Gamma)^{-1}$ is the smallest possible asymptotic variance.

Additional References

- Billingsley, P. (1961). The Lindeberg–Lévy Theorem for Martingales. *Proceedings of the American Mathematical Society*, 12(5), 788–792.
- Duffie, D. and K. J. Singleton (1993). Simulated Moments Estimation of Markov Models of Asset Prices. *Econometrica*, 61(4), 929–952.
- Hansen, L. P. (1982). Large Sample Properties of Generalized Method of Moments Estimators. *Econometrica*, 50(4), 1029–1054.
- Pakes, A. and D. Pollard (1989). Simulation and the Asymptotics of Optimization Estimators. *Econometrica*, 57(5), 1027–1057.

T. Christensen

1 Estimating Standard Errors

1.1 M-Estimators

Recall from Lecture 8 that

$$s(X_t; \theta) = \frac{\partial m(X_t; \theta)}{\partial \theta}, \quad (1)$$

$$D(X_t; \theta) = \frac{\partial^2 m(X_t; \theta)}{\partial \theta \partial \theta'}, \quad (2)$$

$$Z_n = \frac{1}{\sqrt{n}} \sum_{t=1}^n s(X_t; \theta_0), \quad (3)$$

$$H = -E[D(X_t; \theta_0)], \quad (4)$$

where $s(X_t; \theta)$ and Z_n are $p \times 1$ vectors and $D(X_t; \theta)$ and H are $p \times p$. We also had

$$Z_n \rightarrow_d N(0, \Sigma), \quad (5)$$

where

$$\Sigma = E[s(X_t; \theta_0)s(X_t; \theta_0)'] \quad (6)$$

with IID data or when $(s(X_t; \theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ is a MDS; otherwise,

$$\Sigma = C_0 + \sum_{j=1}^{\infty} (C_j + C_j'), \quad (7)$$

where

$$C_j = E[s(X_{t+j}; \theta_0)s(X_t; \theta_0)']. \quad (8)$$

In Lecture 8, Theorem 4 we showed that

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, H^{-1}\Sigma H^{-1}). \quad (9)$$

We estimate the asymptotic variance $\Omega := H^{-1}\Sigma H^{-1}$ using the analogue approach, replacing expectations with sample averages and the true parameter θ_0 with our estimator $\hat{\theta}$. So, our estimator

of the asymptotic variance of $\hat{\theta}$ is

$$\hat{\Omega} = (\hat{H})^{-1} \hat{\Sigma} (\hat{H})^{-1}, \quad (10)$$

where

$$\hat{H} = -\frac{1}{n} \sum_{t=1}^n D(X_t; \hat{\theta}) \quad (11)$$

and

$$\hat{\Sigma} = \frac{1}{n} \sum_{t=1}^n s(X_t; \hat{\theta}) s(X_t; \hat{\theta})' \quad (12)$$

in the IID or MDS cases. Outside these cases, we need to estimate the long-run variance (7) of the process $(s(X_t; \theta_0))_{t \in \mathbb{Z}}$. We do this using the approach of Newey and West (1987). The estimator is

$$\hat{\Sigma} = \hat{C}_0 + \sum_{j=1}^{J_n} \left(1 - \frac{j}{J_n + 1}\right) (\hat{C}_j + \hat{C}_j'), \quad (13)$$

$$\hat{C}_j = \frac{1}{n} \sum_{t=j+1}^n s(X_t; \hat{\theta}) s(X_{t-j}; \hat{\theta})', \quad (14)$$

where the number of lags J_n is a tuning parameter chosen by the researcher such that $J_n \rightarrow \infty$ as $n \rightarrow \infty$, but slower than n . The weighting term $(1 - \frac{j}{J_n + 1})$ ensures that the resulting estimator \hat{S} is positive definite. In practice J_n is fixed at something sensible, so we might use 8, 12 or 16 lags with quarterly data. The standard error estimates can be sensitive to the choice of J_n .

In either case, we form 95% CIs for elements θ_{0i} of θ_0 as

$$\hat{\theta}_i \pm 1.96 \sqrt{\frac{((\hat{H})^{-1} \hat{\Sigma} (\hat{H})^{-1})_{ii}}{n}}. \quad (15)$$

1.2 GMM Estimators

Recall from Lecture 8 that

$$d(X_t; \theta) = \frac{\partial g(X_t; \theta)}{\partial \theta'}, \quad (16)$$

$$G_n(\theta) = \frac{1}{n} \sum_{t=1}^n d(X_t; \theta), \quad (17)$$

$$G(\theta) = E[d(X_t; \theta)], \quad (18)$$

$$G = G(\theta_0). \quad (19)$$

We also had that

$$Z_n := \sqrt{n} g_n(\theta_0) = \frac{1}{\sqrt{n}} \sum_{t=1}^n g(X_t; \theta_0) \rightarrow_d N(0, S), \quad (20)$$

where

$$S = E[g(X_t; \theta_0)g(X_t; \theta_0)'] \quad (21)$$

with IID data or if $(g(X_t; \theta_0), \mathcal{F}_t)_{t \in \mathbb{Z}}$ is a MDS. Otherwise, S is of the form on the right-hand side of equation (7) with

$$C_j = E[g(X_{t+j}; \theta_0)g(X_t; \theta_0)']. \quad (22)$$

Lecture 8, Theorem 5 showed that the asymptotic distribution of the GMM estimator is

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, (G'WG)^{-1}G'WSWG(G'WG)^{-1}), \quad (23)$$

where W is the probability limit of the weight matrix \hat{W} . We again use the analogue approach to estimate the asymptotic variance:

$$\hat{\Omega} = (\hat{G}'\hat{W}\hat{G})^{-1}\hat{G}'\hat{W}\hat{S}\hat{W}\hat{G}(\hat{G}'\hat{W}\hat{G})^{-1}, \quad (24)$$

where

$$\hat{G} = G_n(\hat{\theta}) \quad (25)$$

and

$$\hat{S} = \frac{1}{n} \sum_{t=1}^n g(X_t; \hat{\theta})g(X_t; \hat{\theta})', \quad (26)$$

in the IID or MDS cases; otherwise,

$$\hat{S} = \hat{C}_0 + \sum_{j=1}^{J_n} \left(1 - \frac{j}{J_n + 1}\right) (\hat{C}_j + \hat{C}_j'), \quad (27)$$

$$\hat{C}_j = \frac{1}{n} \sum_{t=j+1}^n g(X_t; \hat{\theta})g(X_{t-j}; \hat{\theta})', \quad (28)$$

where J_n again grows slowly with sample size n .

In either case, we form 95% CIs for elements θ_{0i} of θ_0 as

$$\hat{\theta}_i \pm 1.96 \sqrt{\frac{((\hat{G}'\hat{W}\hat{G})^{-1}\hat{G}'\hat{W}\hat{S}\hat{W}\hat{G}(\hat{G}'\hat{W}\hat{G})^{-1})_{ii}}{n}}. \quad (29)$$

The asymptotic variance of SMM estimators is essentially the same, except they are inflated by a factor $(1 + \tau)$ where τ is the ratio of sample size n to number of simulation draws m .

2 Hypothesis Tests about Parameters

First suppose we are interested in $a(\theta_0)$ for some $a : \Theta \rightarrow \mathbb{R}^k$. If a is continuously differentiable, then by the δ -method we have

$$\sqrt{n}(a(\hat{\theta}) - a(\theta_0)) \rightarrow_d N(0, A(\theta_0)\Omega A(\theta_0)'), \quad (30)$$

where

$$A(\theta_0) = \frac{\partial a(\theta_0)}{\partial \theta'}. \quad (31)$$

When $k = 1$ (so a is scalar-valued) we can perform inference on $a(\theta_0)$ by constructing the 95% confidence interval

$$a(\hat{\theta}) \pm 1.96 \sqrt{\frac{A(\hat{\theta})\hat{\Omega}A(\hat{\theta})'}{n}}, \quad (32)$$

where $\hat{\Omega}$ is one of the estimators proposed above. We can test $H_0 : a(\theta_0) = 0$ against $H_1 : a(\theta_0) \neq 0$ by checking whether the confidence interval spans zero. If it doesn't, then we reject H_0 .

When $1 < k \leq p$ (so a is vector-valued) we cannot reduce the problem of testing $H_0 : a(\theta_0) = 0$ against $H_1 : a(\theta_0) \neq 0$ down to a single scalar confidence interval. Instead, we perform a *Wald test* based on the statistic

$$\xi_W := n(a(\hat{\theta}))'(A(\hat{\theta})\hat{\Omega}A(\hat{\theta})')^{-1}a(\hat{\theta}). \quad (33)$$

Under H_0 , this test is asymptotically distributed as χ_k^2 . For this asymptotic approximation to be valid we require that the rank of $A(\theta_0)$ is k , so that none of the hypotheses are redundant. The decision rule is to reject H_0 if ξ_W exceeds the 0.95 quantile of the χ_k^2 distribution.

For the intuition behind the test, suppose H_0 is correct. Then by virtue of (30), we have

$$\sqrt{n}(A(\theta_0)\Omega A(\theta_0)')^{-1/2}a(\hat{\theta}) \rightarrow_d N(0, I_k). \quad (34)$$

Because $A(\hat{\theta})\hat{\Omega}A(\hat{\theta})' \rightarrow_p A(\theta_0)\Omega A(\theta_0)'$, we also have

$$\sqrt{n}(A(\hat{\theta})\hat{\Omega}A(\hat{\theta})')^{-1/2}a(\hat{\theta}) \rightarrow_d N(0, I_k). \quad (35)$$

But note that ξ_W is precisely the squared Euclidean norm $\|\cdot\|^2$ of the left-hand side, and a χ_k^2 distributed random variable can be written as $\|Z\|^2$ where $Z \sim N(0, I_k)$.

There are also *Lagrange Multiplier* and *Quasi-Likelihood Ratio* tests. But for models in which $H \neq \Sigma$ (e.g. GMM without optimal weighting or ML under misspecification) the asymptotic distribution of these test statistics is no longer χ_k^2 , but of a more complicated form involving nuisance parameters.

3 Tests of Over-identifying Restrictions

In over identified GMM models (with number of moments K exceeding number of parameters p) we can use the $K - p$ excess restrictions to test whether the model fits the data. This can be interpreted as a specification test of the model. The null hypothesis that we are testing is that $E[g(X_t; \theta_0)] = 0$ for some $\theta_0 \in \Theta$. The test statistic is

$$J = -2nQ_n(\hat{\theta}) \quad (36)$$

where the weighting matrix \hat{W} must be a consistent estimator of S^{-1} . Under the null hypothesis,

$$J \rightarrow_d \chi^2_{(K-p)} \quad (37)$$

as $n \rightarrow \infty$. This test is variously referred to as the *test of over-identifying restrictions*, the *Hansen–Sargan test*, or the *J test*.

Theorem 1. *Let the Assumptions of Lecture 8, Theorem 5 hold and let the GMM estimator be optimally weighted, i.e., $\hat{W} \rightarrow_p S^{-1}$. Then: $J \rightarrow_d \chi^2_{(K-p)}$.*

Proof. As $\hat{W} \rightarrow_p S^{-1}$ we can write:

$$J = (\sqrt{n}g_n(\hat{\theta}))'(S^{-1} + o_p(1))(\sqrt{n}g_n(\hat{\theta})). \quad (38)$$

Taking a mean-value expansion of $g_n(\hat{\theta})$ about θ_0 and using the fact that $G_n(\tilde{\theta}) \rightarrow_p G$ (see the proof of Lecture 8, Theorem 5), we obtain:

$$J = (\sqrt{n}g_n(\theta_0) + G_n(\tilde{\theta})\sqrt{n}(\hat{\theta} - \theta_0))'(S^{-1} + o_p(1))(\sqrt{n}g_n(\theta_0) + G_n(\tilde{\theta})\sqrt{n}(\hat{\theta} - \theta_0)) \quad (39)$$

$$= (\sqrt{n}g_n(\theta_0) + G\sqrt{n}(\hat{\theta} - \theta_0))'S^{-1}(\sqrt{n}g_n(\theta_0) + G\sqrt{n}(\hat{\theta} - \theta_0)) + o_p(1). \quad (40)$$

Substituting

$$\sqrt{n}(\hat{\theta} - \theta_0) = -(G'WG)^{-1}G'W(\sqrt{n}g_n(\theta_0)) + o_p(1) \quad (41)$$

into the above expression for J and rearranging yields

$$J = (S^{-1/2}\sqrt{n}g_n(\theta_0))' \left\{ I - S^{-1/2}G(G'S^{-1}G)^{-1}G'S^{-1/2} \right\} (S^{-1/2}\sqrt{n}g_n(\theta_0)) \quad (42)$$

where $S^{-1/2}$ denotes the inverse of the positive definite square root of S . Notice that

$$S^{-1/2}\sqrt{n}g_n(\theta_0) \rightarrow_d N(0, I) \quad (43)$$

where I is the $K \times K$ identity matrix. The matrix in braces in display (42) is a $K \times K$ orthogonal projection matrix with rank $K - p$ (Exercise: show that the rank of this matrix is $K - p$).

Finally, using the fact that if $Z \sim N(0, K)$ and Q is a $K \times K$ orthogonal projection matrix, then $Z'QZ \sim \chi^2_{\text{rank}(Q)}$, it follows by display (42) that $J \rightarrow_d \chi^2_{K-p}$. ■

We reject H_0 if the J statistic exceeds the 0.95 quantile of the $\chi^2_{(K-p)}$ critical value. If so, we conclude that some moments are misspecified. Note that the test of over-identifying restrictions tells you whether or not some moment conditions are misspecified, but does not tell you *which* moments are misspecified.

Additional References

- Newey, W. K. and K. D. West (1987). A Simple, Positive Semi-definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix. *Econometrica*, 55(3), 703–708.