# Proofs for Large Sample Properties of Generalized Method of Moments Estimators[*]

Lars Peter Hansen
University of Chicago

March 8, 2012

## 1 Introduction

*Econometrica* did not publish many of the proofs in my paper Hansen (1982). These notes provide the missing proofs about consistency of GMM (generalized method of moments) estimators. Section 2 provides the required assumptions and proofs for a Law of Large Numbers for Random Functions to hold. I present a theorem guaranteeing the almost sure uniform convergence of the sample mean of a random function to its population counterpart. In section 3 I establish the consistency of a GMM estimator as a direct application of the theorem. By imposing a special structure on the form of the random function, in Section 4 I relax the assumption of a compact parameter space. I fill in some results on random sets that emerge as the outcome of minimization in section 5. Finally in section 6, I study the consistency of two-step estimation, where the first-step delivers a set of possible estimators and the second step reduces the set to be used in estimating an identified parameter vector.

With minor editorial modification and elaboration, the proofs are very close to their original form. I have made no attempt to connect them to the literature on nonlinear and GMM estimation subsequent to the publication of Hansen (1982).

## 2 Law of Large Numbers for Random Functions

We begin giving some definitions and assumptions.

Let $\mathcal{L}$ be the space of continuous functions defined on a parameter space $P$, and let $\lambda$ be the *sup* norm on $\mathcal{L}$. It is well known that the Law of Large Numbers applies to stochastic processes that are stationary and ergodic. Recall that a stochastic process can be viewed as an indexed (by calendar time) family of random variables. In this chapter we extend the Law

of Large Numbers to stationary and ergodic indexed families (by calendar time) of random functions. A random function is a measurable mapping from the collection of states of the world $\Omega$ into the space $\mathcal{L}$ of continuous functions. It is the natural extension of a random variable.

We assume that the metric space of potential parameter values, $\beta$, is compact.

**Assumption 2.1.** $(P, \pi)$ *is a compact metric space.*

Let $(\Omega, \mathcal{F}, Pr)$ denote the underlying probability space and $\phi$ the random function under consideration.

**Definition 2.1.** *A* **random function** $\phi$ *is a measurable function mapping* $\Omega \to \mathcal{L}$.

**Assumption 2.2.** $\phi$ *is a random function.*

Sometimes we will suppress the dependence on $\omega \in \Omega$ and just write $\phi(\beta)$. Under Assumption 2.2, $\phi(\beta)$ is a random variable for each $\beta$ in $P$. In fact, the measurability restriction in Assumption 2.2 is equivalent to requiring that $\phi(\beta)$ be a random variable for each $\beta$ in $P$ [see Billingsley (1978)].

We construct stochastic processes by using a transformation $S : \Omega \to \Omega$. Our development follows Breiman (1968).

**Definition 2.2.** $S$ *is* **measure preserving** *if for any event $\Lambda$ in $\mathcal{F}$, $Pr(\Lambda) = Pr[S^{-1}(\Lambda)]$.*

**Definition 2.3.** $S$ *is* **ergodic** *if for any event such that $S^{-1}(\Lambda) = \Lambda$, $Pr(\Lambda)$ is equal to zero or one.*[1]

**Assumption 2.3.** $S$ *is measure-preserving and ergodic.*

The process of random functions is constructed using $S$.

$$\phi_t(\omega, \beta) = \phi\left[S^t(\omega), \beta\right]$$

The stochastic process $\{\phi_t(\cdot, \beta) : t \geq 0\}$ is stationary and ergodic for each $\beta$ [see Proposition 6.9 and Proposition 6.18 in Breiman (1968)]. The notation $\sum_T(\phi)$ denotes the random function given by

$$\sum_T\left[\phi(\omega, \beta)\right] = \phi\left[S(\omega), \beta\right] + \phi\left[S^2(\omega), \beta\right] + ... + \phi\left[S^T(\omega), \beta\right].$$

To obtain a Law of Large Numbers for $\phi$ we require that $\phi$ has finite first moments and be continuous in a particular sense.

**Definition 2.4.** *A random function $\phi$ has* finite first moments *if $E|\phi(\beta)| < \infty$ for all $\beta \in P$.*

**Assumption 2.4.** $\phi$ *has finite first moments.*

---

[1]Event for which $S^{-1}(\Lambda) = \Lambda$ are referred to as invariant events.

The following notation is used to define a modulus of continuity for a random function:

$$Mod_\phi(\delta, \beta) = \sup \{|\phi(\beta) - \phi(\beta^*)| : \beta^* \in P \text{ and } \pi(\beta, \beta^*) < \delta\}.$$

Since $P$ is compact, $P$ is separable. Consequently, a dense sequence $\{\beta_j : j \geq 1\}$ can be used in place of $P$ in evaluating the *supremum*. In this case, $Mod_\phi(\delta, \beta)$ is a random variable for each positive value of $\delta$ and each $\beta$ in $P$. Also, $Mod_\phi(\delta, \beta) \geq Mod_\phi(\delta^*, \beta)$ if $\delta$ is greater than $\delta^*$. Since $\phi$ maps into $\mathcal{L}$,

$$\lim_{\delta \downarrow 0} Mod_\phi(\delta, \beta) = 0 \text{ for all } \omega \in \Omega \text{ and all } \beta \in P. \tag{1}$$

**Definition 2.5.** *A random function $\phi$ is* first-moment continuous *if for each $\beta \in P$,*

$$\lim_{\delta \downarrow 0} E[Mod_\phi(\delta, \beta)] = 0.$$

**Assumption 2.5.** $\phi$ *is first-moment continuous.*

Since $Mod_\phi(\delta, \beta)$ is decreasing in $\delta$, we have the following result.

**Lemma 2.1** (DeGroot (1970) page 206)**.** *Under Assumptions 2.1 and 2.2, $\phi$ is first-moment continuous if, and only if, for each $\beta \in P$ there is $\delta_\beta > 0$ such that*

$$E[Mod_\phi(\delta_\beta, \beta)] < \infty.$$

*Proof.* The *if* part of the proof follows from the Dominated Convergence Theorem and (1). The *only if* part is immediate. $\square$

Since $S$ is ergodic, a natural candidate for the limit of time series averages of $\phi$ is $E(\phi)$. To establish the Law of Large Numbers for Random Functions, we use: i) the pointwise continuity of $E(\phi)$, ii) a pointwise Law of Large Numbers for $\{(1/T) \sum_T (\phi)(\beta)\}$ for each $\beta$ in $P$, and iii) a pointwise Law of Large Numbers for $\{(1/T) \sum_T [Mod_\phi(\delta, \beta)]\}$ for each $\beta$ in $P$ and positive $\delta$. We establish these approximation results in three lemmas prior to our proof of the Law of Large Numbers for Random Functions. We then demonstrate our main result by showing that the assumption of a compact parameter space $(P, \pi)$ can be used to obtain an approximation that is uniform.

Lemma 2.2 establishes the continuity of $E(\phi)$.

**Lemma 2.2.** *Suppose Assumptions 2.1, 2.2, 2.4 and 2.5 are satisfied. Then there is positive-valued function $\delta^*(\beta, j)$ satisfying*

$$|E[\phi(\beta^*)] - E[\phi(\beta)]| < 1/j$$

*for all $\beta^* \in P$ such that $\pi(\beta^*, \beta) < \delta^*(\beta, j)$ and all integer $j \geq 1$.*

*Proof.* Since $\phi$ is first-moment continuous, there is a function $\delta^*(\beta, j)$ such that

$$E\left(Mod_\phi\left[\beta, \delta^*(\beta, j)\right]\right) < 1/j.$$

Note, however, that

$$\begin{aligned}
|E\phi(\beta^*) - E\phi(\beta)| &\leq E\left|\phi(\beta^*) - \phi(\beta)\right| \\
&\leq E\left[Mod_\phi(\beta, \delta^*)(\beta, j)\right] \\
&< 1/j
\end{aligned}$$

for all $\beta^* \in P$ such that $\pi(\beta, \beta^*) < \delta^*(\beta, j)$. $\qquad \square$

Since $P$ is compact, it can be shown that $\delta^*(\beta, j)$ in Lemma 2.2 can be chosen independent of $\beta$. In other words, $E(\phi)$ is uniformly continuous.

For each element $\beta$ in $P$, $\phi(\beta)$ is a random variable with a finite absolute first moment. Thus the Law of Large Numbers applies pointwise as stated in the following lemma.

**Lemma 2.3.** *Suppose Assumptions 2.1, 2.2, 2.3 and 2.4 are satisfied. Then there is an integer-valued function $T^*(\omega, \beta, j)$ and an indexed set $\Lambda^*(\beta) \in \mathcal{F}$ such that $Pr\{\Lambda^*(\beta)\} = 1$ and*

$$\left|(1/T)\sum_T\left[\phi(\beta)\right] - E\left[\phi(\beta)\right]\right| < 1/j \tag{2}$$

*for all $\beta \in P$, $T \geq T^*(\omega, \beta, j)$, $\omega \in \Lambda^*(\beta)$, and $j \geq 1$.*

*Proof.* Since $S$ is measure-preserving and ergodic, $\{(1/T)\sum_T\left[\phi(\beta)\right] : T \geq 1\}$ converges to $E\left[\phi(\beta)\right]$ on a set $\Lambda^*(\beta) \in \mathcal{F}$ satisfying $Pr\{\Lambda^*(\beta)\} = 1$. $\qquad \square$

The Law of Large Numbers also applies to time series averages of $Mod_\phi(\beta, \delta)$. Since the mean of $Mod_\phi(\beta, \delta)$ can be made arbitrarily small by choosing $\delta$ to be small, we can control the local variation of time series averages of the random function $\phi$.

**Lemma 2.4.** *Suppose Assumptions 2.1, 2.2, 2.3 and 2.5 are satisfied. There exists an integer-valued function $T^+(\omega, \beta, j)$, a positive function $\delta^+(\beta, j)$ and an indexed set $\Lambda^+(\beta) \in \mathcal{F}$ such that $Pr\{\Lambda^+(\beta)\} = 1$ and*

$$\left|(1/T)\sum_T\left[\phi(\beta) - \phi(\beta^*)\right]\right| < 1/j \tag{3}$$

*for all $\beta^* \in P$ such that $\pi(\beta, \beta^*) < \delta^+(\beta, j)$, $T \geq T^+(\omega, \beta, j)$, $\omega \in \Lambda^+(\beta)$ and $j \geq 1$.*

*Proof.* Since $\phi$ is first-moment continuous, $Mod_\phi(\beta, 1/n)$ has a finite first moment for some positive integer $n$. Since $S$ is measure-preserving and ergodic, $\{(1/T)\sum_T\left[Mod_\phi(\beta, 1/j)\right] : T \geq 1\}$ converges to $E\left[Mod_\phi(\beta, 1/j)\right]$ on a set $\Lambda^-(\beta, j)$ satisfying $Pr\{\Lambda^-(\beta, j)\} = 1$ for $j \geq n$. Let

$$\Lambda^+(\beta) = \bigcap_{j \geq n} \Lambda^-(\beta, j).$$

Then $\Lambda^+(\beta)$ is measurable and $Pr\{\Lambda^+(\beta)\} = 1$.

For each $j$, choose $[1/\delta^+(\beta,j)]$ to equal some integer greater than or equal to $n$ such that

$$E\left((Mod_\phi\left[\beta,\delta^+(\beta,j)\right]\right) < 1/(2j).$$

Since $\{(1/T)\sum_T\{Mod_\phi\left[\beta,\delta^+(\beta,j)\right] : T \geq 1\}$ converges almost surely to $E\left(Mod_\phi\left[\beta,\delta^+(\beta,j)\right]\right)$ on $\Lambda^+(\beta)$, there exists an integer-valued function $T^+(\omega,\beta,j)$ such that

$$\left|(1/T)\sum_T\left(Mod_\phi\left[\beta,\delta^+(\beta,j)\right]\right) - E\left(Mod_\phi\left[\beta,\delta^+(\beta,j)\right]\right)\right| < 1/2j$$

for $T \geq T^+(\omega,\beta,j)$. Therefore,

$$(1/T)\left|\sum_T\left[\phi(\beta) - \phi(\beta^*)\right]\right| < (1/T)\sum_T\left(Mod_\phi\left[\beta,\delta^+(\beta,j)\right]\right) < 1/j$$

for all $\beta^* \in P$ such that $\pi(\beta,\beta^*) < \delta^+(\beta,j)$, $T \geq T^+(\omega,\beta,j)$, $\omega \in \Lambda^+(\beta)$, and $j \geq 1$. $\square$

Our main result establishes the almost sure convergence of time series averages of random functions. Suppose $\phi_1$ and $\phi_2$ are two random functions. Then

$$\lambda(\phi_1,\phi_2) = \sup_{\beta \in P}|\phi_1(\beta) - \phi_2(\beta)|$$

is a measure of distance between these functions that depends on the sample point. Since $P$ is separable, it suffices to take the supremum over a countable dense sequence. Hence $\lambda(\phi_1,\phi_2)$ is a random variable (measurable function). Almost sure convergence of sequences of random functions is defined using the metric $\lambda$.

**Definition 2.6.** *A sequence $\{\phi_j : j \geq 1\}$ of random functions converges almost surely to a random function $\phi_0$ if $\{\lambda(\phi_j,\phi_0) : j \geq 1\}$ converges almost surely to zero.*

We now combine the conclusions from Lemmas 2.2, 2.3 and 2.4 to obtain a Law of Large Numbers for random functions. The idea is to exploit that fact that $P$ is compact to move from pointwise to uniform convergence. Notice that in these three lemmas, $\Lambda^+, \Lambda^*, T^+$ and $T^*$ all depend on $\beta$. In proving this Law of Large Numbers, we will use compactness to show how the dependence on the parameter value can be eliminated.

**Theorem 2.1.** *Suppose Assumptions 2.1, 2.2, 2.3, 2.4 and 2.5 are satisfied. Then $\{(1/T)\sum_T(\phi) : T \geq 1\}$ converges almost surely to $E(\phi)$.*

*Proof.* Let

$$Q(\beta,j) = \{\beta^* \in P : \pi(\beta,\beta^*) < \min\left\{\delta^*(\beta,j),\delta^+(\beta,j)\right\}\}.$$

Then for each $j \geq 1$,

$$P = \bigcup_{\beta \in P}Q(\beta,j).$$

Since $P$ is compact a finite number of $\beta_i$ can be selected so that

$$P = \bigcup_{i \geq 1}^{N(j)} Q(\beta_i, j)$$

where $N(j)$ is integer-valued and $\{\beta_i : i \geq 1\}$ is a sequence in $P$. Construct

$$\Lambda = \bigcap_{i \geq 1} \left[ \Lambda^*(\beta_i) \bigcap \Lambda^+(\beta_i) \right].$$

Then $\Lambda \in \mathcal{F}$ and $Pr(\Lambda) = 1$. Let

$$
\begin{aligned}
T(\omega, j) &= \max\{ T^*(\omega, \beta_1, j), T^*(\omega, \beta_2, j), ..., T^*\left[\omega, \beta_{N(j)}, j\right], \\
&\quad T^+(\omega, \beta_1, j), T^+(\omega, \beta_2, j), ..., T^+\left[\omega, \beta_{N(j)}, j\right] \}.
\end{aligned}
$$

For $T \geq T(\omega, j)$, Lemmas 2.2, 2.3 and 2.4 imply that

$$
\left| (1/T) \sum_T [\phi(\beta)] - E[\phi(\beta)] \right| \leq (1/T) \left| \sum_T [\phi(\beta)] - \sum_T [\phi(\beta_i)] \right|
$$

$$
+ \left| (1/T) \sum_T [\phi(\beta_i)] - E[\phi(\beta_i)] \right| + |E[\phi(\beta_i)] - E[\phi(\beta)]|
$$

$$
< \quad 3/j
$$

where $\beta_i$ is chosen so that $\beta \in Q(\beta_i, j)$ for some $1 \leq i \leq N(j)$. Hence

$$\lambda \left[ (1/T) \sum_T (\phi), E\phi \right] \leq 3/j$$

for $T \geq T(\omega, j)$ and $\omega \in \Lambda$. Therefore, $\{\lambda[(1/T) \sum_T (\phi), E\phi] : T \geq 1\}$ converges to zero on $\Lambda$. □

# 3   Consistency of the GMM Estimator

We apply Theorem 2.1 to obtain an approximation result for a GMM estimator as defined in Hansen (1982). So far, the results obtained in the previous section have been for the case of random functions that map into $\mathbb{R}$. The GMM estimator works by making sample analogues of population orthogonality conditions close to zero. We will map the assumptions in this note to those in Hansen (1982) and show that the consistency of the GMM estimator (Theorem 2.1 in Hansen (1982)) is an application of Theorem 2.1 above.

First, construct a stochastic process by using the transformation $S : \Omega \to \Omega$. Let,

$$x_t(\omega) = x\left[S^t(\omega)\right]$$

where $x_t(\omega)$ has $p$ components as in Hansen (1982).

Assumption 2.3 from the previous section ensures that this process is stationary and ergodic, which is Assumption 2.1 in Hansen (1982). Hansen (1982) restricts the parameter space to be a separable metric space (Assumption 2.2). This is implied by assumption 2.1 above, since a compact space is separable.

To represent the population orthogonality conditions we will consider a function $f : \mathbb{R}^p \times P \to \mathbb{R}^r$. The random function $\phi$ in section 2 is given by:

$$\phi(\omega, \beta) = f[x(\omega), \beta]. \tag{4}$$

Hansen (1982) requires that $f(\cdot, \beta)$ is Borel measurable for each $\beta$ in $P$ and $f(x, \cdot)$ is continuous on $P$ for each $x$ in $\mathbb{R}^p$ (Assumption 2.3). Given construction (4), these restrictions imply Assumption 2.2.

Assumption 2.4 in the GMM paper requires that $Ef(x, \beta)$ exists and is finite for all $\beta \in P$ and $Ef(x, \beta_0) = 0$. The first part of this assumption is equivalent to Assumption 2.4 given construction (4) of the random function $\phi$. To match the GMM paper, we impose the identification restriction:

**Assumption 3.1.** *The equation*
$$E\phi(\beta) = 0 \tag{5}$$
*is satisfied on $P$ if, and only if $\beta = \beta_0$.*

Finally, as in Assumption 2.5 of Hansen (1982) we introduce a sequence of weighting matrices:

**Assumption 3.2.** *The sequence of random matrices $\{A_T : T \geq 1\}$ converges almost surely to a constant matrix $A_0$ where $(A_0)'A_0$ has full rank.*

We now have all the ingredients to state the following corollary to Theorem 2.1.

**Corollary 3.1.** *(Theorem 2.1 in Hansen (1982)) Suppose that Assumptions 2.1, 2.2, 2.3, 2.4, 2.5, 3.1 and 3.2 are satisfied. Then*

$$b_T = \arg\min_{\beta \in P} \left[\frac{1}{T}\sum_{t=1}^{T} f(x_t, \beta)\right]' (A_T)' A_T \left[\frac{1}{T}\sum_{t=1}^{T} f(x_t, \beta)\right]$$

*converges almost surely to $\beta_0$.*

*Proof.* Theorem 2.1 implies that the objective function for the minimization problem converges uniformly to the continuous function:

$$g_0(\beta) = Ef(x_t, \beta)' (A_0)' A_0 Ef(x_t, \beta) \tag{6}$$

with probability one. Assumptions 3.1 and 3.2 guarantee that the limiting objective $g_0$ has a unique minimizer at $\beta_0$ with a minimized objective equal to zero. The conclusion follows. $\square$

# 4  Parameter Separation

We now explore an alternative consistency argument that avoids assuming a compact parameter space. Instead we suppose that we have estimation equations with a special separable structure. We take as our starting point:

$$EXh(\beta) = 0. \tag{7}$$

where the random matrix $X$ is a general function of the data and is $r$ by $m$. This matrix may in fact contain nonlinear transformations of the data. The key restriction is that it does not depend on unknown parameters.

**Assumption 4.1.** *$P$ is a locally compact metric space.*

We presume that the parameter vector $\beta_0$ can be identified from the moment condition 7. We will sharpen this assumption in what follows.

**Assumption 4.2.** *$h : P \to \Gamma \subset \mathbb{R}^m$ is a homeomorphism.*[2]

Neither $P$ nor $\Gamma$ is necessarily compact. While we can estimate $EX$ in large samples, approximation errors get magnified if $\Gamma$ is unbounded as $\beta$ ranges over the parameter space.

Prior to studying the estimation problem that interests us, we first consider an auxiliary problem. Define:

$$\gamma_0 = h(\beta_0).$$

Consider estimating $\gamma_0$ by solving:

$$\tilde{c}_T = \arg\min_{\gamma \in \Gamma} \frac{\gamma' \left[ \frac{1}{T} \sum_T (X) \right]' A_T' A_T \left[ \frac{1}{T} \sum_T (X) \right] \gamma}{1 + |\gamma|^2}$$

The scaling by $1 + |\gamma|^2$ limits the magnitude of the approximation error. With this in mind we form $\Gamma^*$ to be the closure of the bounded set:

$$\hat{\Gamma} = \left\{ \frac{\gamma}{\sqrt{1 + |\gamma|^2}} : \gamma \in \Gamma \right\}.$$

Assumption 2.2 and 2.3 guarantees that

$$\frac{1}{T} \sum_T (X) \to EX$$

provided that $X$ has a finite first moment.

**Assumption 4.3.** *The matrix $X$ has a finite first moment and $EX\gamma = 0$ on $\Gamma^*$ if, and only if $\gamma = \frac{\gamma_0}{\sqrt{1 + |\gamma_0|^2}}$.*

---

[2]$h$ is continuous, one-to-one, and has a continuous inverse.

Since we have added in the closure points, this is a stronger condition than the usual identification condition using moment conditions.

**Theorem 4.1.** *Suppose that Assumptions 2.2, 2.3, 3.2, 4.1, 4.2, and 4.3 are satisfied. Then $\tilde{c}_T$ converges to $\gamma_0$ almost surely.*

*Proof.* First transform the problems to be:

$$\hat{c}_T = \arg\min_{\gamma^* \in \Gamma^*} \gamma^{*\prime} \left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right]' \gamma^*.$$

Since the objective function converges uniformly almost surely, the sequence of minimizers $\hat{c}_T$ converges to:

$$\frac{\gamma_0}{\sqrt{1+|\gamma_0|^2}} = \arg\min_{\gamma^* \in \Gamma^*} \gamma^{*\prime}(EX)' A_0{}' A_0 (EX)\gamma^*.$$

Since $\hat{\Gamma}$ is locally compact, for sufficiently large $T$, $\hat{c}_T$ is in the set $\hat{\Gamma}$. Thus:

$$\tilde{c}_T = \frac{\hat{c}_T}{\sqrt{1-|\hat{c}_T|^2}}$$

for sufficiently large $T$. The conclusion follows from the convergence of $\hat{c}_T$. $\square$

While the auxiliary estimator is of interest in its own right, we will now use the auxiliary estimator as a bound for:

$$c_T = \arg\min_{\gamma \in \Gamma} \gamma' \left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right] \gamma.$$

As we will see,

$$|c_T| \le |\tilde{c}_T| \tag{8}$$

because the scaling of the objective function for the auxiliary estimator implicitly rewards magnitude.

**Theorem 4.2.** *Suppose that Assumptions 2.2, 2.3, 3.2, 4.1, 4.2, and 4.3 are satisfied. Then $c_T$ converges to $\gamma_0$ almost surely.*

*Proof.* There are two inequalities implied by the minimization problems used to construct $c_T$ and $\tilde{c}_T$:

$$(c_T)'\left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right](c_T) \le (\tilde{c}_T)'\left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right](\tilde{c}_T)$$

and

$$\frac{(c_T)'\left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right](c_T)}{1+|c_T|^2} \ge \frac{(\tilde{c}_T)'\left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right](\tilde{c}_T)}{1+|\tilde{c}_T|^2}.$$

9

Bound (8) follows because to achieve the second inequality, the left-hand side must have a smaller denominator.

Since $\hat{c}_T$ converges to $\gamma_0$, it follows that for sufficiently large $T$, $|c_T| \leq |\gamma_0| + 1$ and hence the estimator $c_T$ is eventually in the compact set given by the closure of:

$$\{\gamma \in \Gamma : |\gamma| \leq |\gamma_0| + 1\}. \tag{9}$$

Since the objective converges uniformly almost surely on this compact set and the limiting objective has a unique minimizer at $\gamma_0$, the conclusion follows except that we used the closure of the set (9) and not $\Gamma$. Since $\Gamma$ is locally compact, for $T$ sufficiently large minimizers over the set (9) will eventually be in $\Gamma$. $\qquad\square$

Finally to produce an estimator of $\beta_0$ we solve:

$$c_T = h(b_T).$$

Notice that we may equivalently define $b_T$ as:

$$b_T = \arg\min_{\beta \in P} h(\beta)' \left[\frac{1}{T}\sum\nolimits_T (X)\right]' A_T{}' A_T \left[\frac{1}{T}\sum\nolimits_T (X)\right] h(\beta).$$

**Corollary 4.1.** *Suppose that Assumptions 2.3, 3.2, 4.1, 4.2, and 4.3 are satisfied. Then $b_T$ converges to $\beta_0$ almost surely.*

*Proof.* This follows from the fact that $h$ is one-to-one with a continuous inverse. $\qquad\square$

Theorem 2.2 in Hansen (1982) is a special case of the latter corollary. That paper assumes (Assumption 2.6) that

$$f(x_t, \beta) = \begin{bmatrix} C_0(x_t) & C_1(x_t) \end{bmatrix} \begin{bmatrix} 1 \\ \lambda(\beta) \end{bmatrix}.$$

Thus to see the relation, we may set

$$X_t = \begin{bmatrix} C_0(x_t) & C_1(x_t) \end{bmatrix},$$

and

$$h(\beta) = \begin{bmatrix} 1 \\ \lambda(\beta) \end{bmatrix}.$$

In addition to Assumptions 2.1 to 2.6, for which we have already provided a mapping with the assumptions of this paper, Theorem 2.2 in Hansen (1982) requires that the parameter space be locally compact, that $\lambda$ be a homeomorphism and an identification assumption. These are implied by Assumptions 4.1, 4.2, and 4.3 in this note, respectively.

# 5 Random sets

So far, the discussion has been bit casual about the the minimizers of GMM objectives. We address this gap by viewing the collection of the compact subsets of $(P, \pi)$ as a metric space and appealing to some known results about the corresponding collection of Borel sets.

Let $\mathcal{K}$ be the family of nonempty compact subsets of $P$. The Hausdorff metric $\delta$ on the collection of compact sets in $(P, \pi)$ is:

$$\delta(K_1, K_2) = \max\{\rho(K_1, K_2), \rho(K_2, K_1)\}$$

where

$$\rho(K_1, K_2) = \sup_{\beta \in K_1} \inf_{\alpha \in K_2} \pi(\alpha, \beta)$$

**Definition 5.1.** *A random set $\Theta$ is a Borel measurable mapping from $\Omega \to \mathcal{K}$.*

**Lemma 5.1.** *Suppose Assumption 2.1 is satisfied, $\phi$ is a random function (Assumption 2.2), and $\Theta$ is a random set. Let $\Delta$ denote the set of minimizers of $\phi$ over the set $\Theta$. Then $\Delta$ is a random set.*

*Proof.* Let $\{\beta_n\}$ be a countable dense subset of $P$, and let $O$ be any open subset of $P$. Define

$$\tilde{\phi}(\omega, n, j) = \begin{cases} +\infty & \text{if } \left\{\beta \in P : \pi(\beta, \beta_n) < \frac{1}{j}\right\} \cap O \cap \Theta(\omega) = \varnothing \\ \phi(\omega, \beta_n) & \text{otherwise} \end{cases}.$$

From Debreu (1967) (page 355), the set

$$\mathcal{O} = \left\{K \in \mathcal{K} : \left\{\beta \in P : \pi(\beta, \beta_n) < \frac{1}{j}\right\} \cap O \cap K \neq \varnothing\right\}$$

is open. Since $\Theta$ is Borel measurable,

$$\{\omega \in \Omega : \Theta(\omega) \in \mathcal{O}\} = \left\{\omega \in \Omega : \left\{\beta \in P : \pi(\beta, \beta_n) < \frac{1}{j}\right\} \cap O \cap \Theta(\omega) \neq \varnothing\right\}$$

is measurable, and similarly the complement set

$$\left\{\omega \in \Omega : \left\{\beta \in P : \pi(\beta, \beta_n) < \frac{1}{j}\right\} \cap O \cap \Theta(\omega) = \varnothing\right\}$$

is measurable. Moreover, $\tilde{\phi}(\cdot, n, j)$ is a Borel measurable function as a mapping onto the extended real numbers.

By the continuity of $\phi(\omega, \cdot)$,

$$B(\omega, O) = \inf_n \lim_{j \to \infty} \tilde{\phi}(\omega, n, j) = \inf_{O \cap \Theta(\omega)} \phi(\omega, \beta).$$

is Borel measurable since we can express it as a limit of an infimum over a countable collection of Borel measurable functions. Notice that

$$\{\omega \in \Omega : O \cap \Delta(\omega) \neq \varnothing\} = \{\omega \in \Omega : B(\omega, O) = B(\omega, P)\}.$$

Both $B(\cdot, O)$ and $B(\cdot, P)$ are Borel measurable, and as a consequence it follows that the set on left-hand side is measurable. Construct the open set (of compact sets):

$$\widetilde{\mathcal{O}} = \{K \in \mathcal{K} : O \cap K \neq \varnothing\},$$

and note that

$$\{\omega \in \Omega : O \cap \Delta(\omega) \neq \varnothing\} = \{\omega \in \Omega : \Delta(\omega) \in \widetilde{\mathcal{O}}\}.$$

From Debreu (1967) (page 355), the finding that this set is measurable for any open set $O$ implies that $\Delta$ is Borel measurable.[3]  □

So far I have been rather informal about the construction of a GMM estimator. I have not proved that in finite samples there is a unique minimizer. There are two strategies to add formality at the level of generality reflected in Corollary 3.1. One possibility is to study consistency using the Hausdorff metric, even though we know that in the limit the set of minimizers converges to a singleton set $\{\beta_0\}$. Another possibility is to introduce a measurable selection rule. For instance if $P$ is a subset of $\mathbb{R}^q$, measurable selection rules include coordinate-by-coordinate lexicographic minimization or maximization over a compact subset of $\mathbb{R}^q$.[4]

To see the value of allowing for a random set to be used in a minimization problem, suppose that the parameter space $P$ is a subset of $\mathbb{R}^q$. Partition an element of $\mathbb{R}^q$: $\beta' = (\beta_1', \beta_2')$ and suppose that

$$\pi(\beta, \alpha) = \max\left\{\pi_1(\alpha_1, \beta_1), \pi_2(\alpha_2, \beta_2)\right\}.$$

Using $P$ and the partitioning we form two new parameter spaces:

$$P_1 = \left\{\beta_1 : \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \in P \text{ for some } \beta_2\right\}$$

$$P_2 = \left\{\beta_2 : \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \in P \text{ for some } \beta_1\right\} \tag{10}$$

When $P$ is compact, so are $P_1$ and $P_2$.

Let $b_1$ be a random vector that takes on values in $P_1$, and form the random set

$$\Theta(\omega) = \left\{\beta \in P : \beta = \begin{bmatrix} b_1(\omega) \\ \beta_2 \end{bmatrix} \text{ for some } \beta_2 \in P_2\right\}.$$

---

[3]Debreu provides a formal analysis, but attributes the result to an unpublished theorem of Dubins and Ornstein.

[4]See also Jennrich (1969), Lemma 2. Jennrich, however, does not consider the role of random parameter spaces in establishing the existence of estimators that are Borel measurable.

It is convenient to view the construction of $b_1$ as the first-step of a two-step estimation problem. The second step selects an estimator from the random set $\Theta(\omega)$. For a given open set $O$, let $O_1$ denote the subset:

$$O_1 = \left\{ \beta_1 \in P_1 : \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \in O \text{ for some } \beta_2 \in P_2 \right\}.$$

Given $\beta_1 \in O_1$, there exists a $\beta_2$ and a "rectangular neighborhood" containing $\begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$ contained in $O$. We may use this rectangular neighborhood to construct a neighbor of $\beta_1$ in $P_1$ and hence $O_1$ is open in $P_1$. Then

$$\{\omega \in \Omega : b_1(\omega) \in O_1\} = \{\omega \in \Omega : \Theta(\omega) \cap O \neq \varnothing\}$$

is measurable implying that the constructed random set is Borel measurable.

# 6 Two-step estimation

To formulate two-step or recursive GMM estimation, partition the function $f$ as:

$$f(x, \beta) = \begin{bmatrix} f_1(x, \beta_1) \\ f_2(x, \beta) \end{bmatrix}.$$

Notice that first component of $f$, $f_1$, depends only $\beta_1$. Use

$$E\left[f_1(X_t, \beta_{1,0})\right] = 0$$

to identify and estimate $\beta_{1,0}$. Then Corollary 3.1 could be used to justify the consistency of the resulting GMM estimator $\{b_{T,1}\}$ for estimating $\beta_{1,0}$ using the parameter space $P_1$ given in (10). The second set of moment restrictions:

$$E\left[f_2(X_t, \beta)\right] = 0$$

may be used to construct an estimator of $\beta_{2,o}$ using the initial estimator for $\beta_{1,0}$. We represent such an estimation approach using a block diagonal construction of the selection matrix $A_T$:

$$A_T = \begin{bmatrix} A_{11,T} & 0 \\ 0 & A_{22,T} \end{bmatrix}$$

where the rows are partitioned in accordance with the parameter vector $\beta$ and columns in accordance with the function $f$.

The following result extends 3.1 to apply to a recursive procedure.

**Theorem 6.1.** *(Theorem 2.3 in Hansen (1982)) Suppose that Assumptions 2.1, 2.2, 2.3, 2.4, 2.5, 3.1 and 3.2 are satisfied. In addition suppose that*

*i) $\{b_{1,T}\}$ converges almost surely to $\beta_{o,1}$.*

*ii) For any $\beta \in P$ and any sequence $\{\beta_{1,j} : j = 1, 2, ...\}$ in $P_1$ that converges to $\beta_1$, there exists a sequence $\{(\beta_{1,j}', \beta_{2,j}')' : j = 1, 2, ...\}$ in $P$ that converges to $\beta$.[5]*

*Then*

$$\Delta_T = \arg\min_{\beta \in \Theta_T} \left[\frac{1}{T}\sum_{t=1}^{T} f(x_t, \beta)\right]' (A_T)' A_T \left[\frac{1}{T}\sum_{t=1}^{T} f(x_t, \beta)\right]$$

*converges almost surely to $\{\beta_0\}$ where*

$$\Theta_T = \left\{\beta \in P : \beta = \begin{bmatrix} b_{1,T} \\ \beta_2 \end{bmatrix} \text{ for some } \beta_2 \in P_2\right\}$$

---

[5]We may be weaken this restriction to apply to some subset of $P$ containing $\beta_0$ when the parameter space is augmented to include closure points required for compactness.

*Proof.* For any $\epsilon > 0$, form an open ball $O$ of $\beta_0$ with radius $\epsilon$. Let

$$\delta = \frac{1}{3} \min_{P-O} g_0(\beta)$$

where $g_0$ is given by (6). Note that the minimizer is attained because $P - N$ is compact, $\delta > 0$, and $g_0$ is continuous and has a unique minimizer $\beta_0$ over the compact set $P$ with $g(\beta_0) = 0$. In particular, the uniqueness follows from Assumptions 3.1 and 3.2. Since is $g_0$ is continuous at $\beta_0$, we may shrink the radius of the ball to $\epsilon^* < \epsilon$ in order that

$$g_0(\beta) \leq \delta$$

for all $\beta$ in this smaller ball $O^*$.

For sufficiently large $T$,

$$\Theta_T \cap O^* \neq \varnothing.$$

Let

$$g_T(\beta) = \left[ \frac{1}{T} \sum_{t=1}^{T} f\left(x_t, \beta\right) \right]' A_T{}' A_T \left[ \frac{1}{T} \sum_{t=1}^{T} f\left(x_t, \beta\right) \right]$$

The sequence of functions $\{g_T\}$ converges uniformly almost surely to $g_0$ as implied by Theorem 2.1 and Assumption 3.2. For sufficiently large $T$, $\beta \in P - O$ and $\alpha \in \Theta_T \cap O^*$,

$$g_T(\beta) > g_0(\beta) - \delta \geq 2\delta \geq g_0(\alpha) + \delta \geq g_T(\alpha).$$

It follows that $\Delta_T \subset O$. Since we are free to make the radius $\epsilon$ as small as possible, the conclusion follows.

$\square$

15

# References

Breiman, Leo. 1968. *Probability*. Addison-Wesley Series in Statistics.

Debreu, G. 1967. Integration of Correspondences. Proc. 5th Berkeley Symp. Math. Stat. Probab. 1965/66, 2, Part 1.

DeGroot, Morris H. 1970. *Optimal Statistical Decisions*. New York: McGraw-Hill.

Hansen, L. P. 1982. Large Sample Properties of Generalized Method of Moments Estimators. *Econometrica* 50:1029–1054.

Jennrich, Robert I. 1969. Asymptotic Properties of Non-Linear Least Squares Estimators. *The Annals of Mathematical Statistics* 40 (2):pp. 633–643.