

# Estimadores M

Gabriel V. Montes-Rojas

# Marco general

- Sea  $m(\mathbf{x}, \theta)$  un modelo paramétrico para  $E(y|\mathbf{x})$  (o cualquier momento de interés) donde
  - $m(\cdot)$  es una función conocida de  $\mathbf{x}$  y  $\theta$ .
  - $\mathbf{x}$  es un vector de  $K$  variables explicativas.
  - $\theta$  es un vector  $P \times 1$  de parámetros.  $\theta \in \Theta \subseteq \mathbb{R}^P$ .
- Un modelo está **correctamente especificado** para la media condicional  $E(y|\mathbf{x})$  si para algún  $\theta_0 \in \Theta$ ,

$$E(y|\mathbf{x}) = m(\mathbf{x}, \theta_0)$$

# Marco general

- Tomemos  $q(\mathbf{w}, \theta)$  como una función del vector aleatorio  $\mathbf{w}$  y el vector de parámetros  $\theta \in \Theta$ .
  - En general tenemos  $\mathbf{w} = (y, \mathbf{x})$  con elementos típicos  $\mathbf{w}_i$  de una muestra  $\{\mathbf{w}_i : i = 1, 2, \dots, N\}$ .
  - $q(\cdot)$  es una función conocida de  $\mathbf{w}$  y  $\theta$ .
  - $\theta$  es un vector  $P \times 1$ .  $\theta \in \Theta \subseteq \mathbb{R}^P$ .
- Un **estimador M** (M-estimator en inglés) de  $\theta_0$  resuelve el problema

$$\min_{\theta \in \Theta} N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta),$$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta)$$

- Definamos la función  $\theta \mapsto M_N(\theta)$  donde  $M_N(\theta) \equiv N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta)$  y  $\theta \mapsto M(\theta)$  donde  $M(\theta) \equiv E[q(\mathbf{w}, \theta)]$ .
- Notemos que  $\hat{\theta}$  es una función de la muestra aleatoria  $\{\mathbf{w}_i : i = 1, 2, \dots, N\}$ . Algunas veces usaremos la notación  $\hat{\theta}_N$  para enfatizar que depende de la muestra.

# Marco general

- En general asumimos que

$$\theta_0 = \arg \min_{\theta \in \Theta} E[q(\mathbf{w}, \theta)].$$

- $E[q(\mathbf{w}, \theta)]$  se puede ver en términos estadísticos como una **función de pérdida** (loss function en inglés).

Ejemplos:

- cuadrática:  $(\cdot)^2$ , que es la base de los estimadores MCO
- valor absoluto:  $|\cdot|$ , que es la base de la regresión en la mediana.
- máxima verosimilitud:  $q = -\ell$  el log de la función de verosimilitud (log-likelihood).

# Identificación

*La identificación requiere que la solución a la minimización sea única.*

$$E[q(\mathbf{w}, \theta_0)] < E[q(\mathbf{w}, \theta)], \text{ para todo } \theta \in \Theta, \theta \neq \theta_0.$$

En otros contextos se usan otras definiciones de identificación. Por ejemplo, en modelos lineales de regresión se requiere que los parámetros se puedan obtener como funciones de esperanzas poblacionales.

# Consistencia

Dado que  $\hat{\theta}$  depende de la función  $\theta \mapsto M_N(\theta)$  necesitamos “convergencia funcional”.

*Convergencia uniforme en probabilidad (uniform convergence in probability):*

$$\max_{\theta \in \Theta} \left| N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta) - E[q(\mathbf{w}, \theta)] \right| \xrightarrow{P} 0.$$

# Ley uniforme y débil de los grandes números

**Ley uniforme y débil de los grandes números:** Supongamos que  $\mathbf{w}$  es un vector aleatorio que toma valores en  $\mathcal{W} \subset \mathbb{R}^M$ ,  $\Theta \subset \mathbb{R}^P$ , y  $q: \mathcal{W} \times \Theta \rightarrow \mathbb{R}$  es una función real. Asumamos que

(a)  $\Theta$  es compacto;

(b) para cada  $\theta \in \Theta$ ,  $q(\cdot, \theta)$  es medible en términos de Borel (Borel measurable) sobre  $\mathcal{W}$ ;

(c) para cada  $\mathbf{w} \in \mathcal{W}$ ,  $q(\mathbf{w}, \cdot)$  es continua en  $\Theta$ ; y

(d)  $|q(\mathbf{w}, \theta)| < b(\mathbf{w})$  para todo  $\theta \in \Theta$ , donde  $b$  es una función no negativa de  $\mathcal{W}$  tal que  $E[b(\mathbf{w})] < \infty$ .

Entonces

$$\max_{\theta \in \Theta} \left| N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta) - E[q(\mathbf{w}, \theta)] \right| \xrightarrow{P} 0.$$

# Consistencia

**Consistencia de los estimadores M:** *Bajo los supuestos de la convergencia uniforme en probabilidad, y asumiendo identificación, entonces el estimador M  $\hat{\theta}$  (o sea,  $\hat{\theta} = \arg \min_{\theta \in \Theta} N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta)$ ) satisface  $\hat{\theta} \xrightarrow{P} \theta$ .*

*Supongamos que  $\hat{\theta} \xrightarrow{P} \theta$ , y asumamos que  $r(\mathbf{w}, \theta)$  satisface los mismos supuestos de  $q(\mathbf{w}, \theta)$ . Entonces,  $N^{-1} \sum_{i=1}^N r(\mathbf{w}, \hat{\theta}) \xrightarrow{P} E[r(\mathbf{w}, \theta_0)]$ .*



## Ejemplo MCO

- En los modelos MCO  $m(\mathbf{x}, \beta) = E(y|\mathbf{x}) = \mathbf{x}\beta$  y  $q(y, \mathbf{x}, \beta) = [y - m(\mathbf{x}, \beta)]^2$ .
- Para la identificación,

$$\begin{aligned} E[q(y, \mathbf{x}, \beta_0)] &= E[y^2 + m(\mathbf{x}, \beta_0)^2 - 2ym(\mathbf{x}, \beta_0)^2] \\ &= E[(\mathbf{x}\beta_0)^2 + u^2 + (\mathbf{x}\beta_0)^2 - 2(\mathbf{x}\beta_0)^2] = \sigma^2 \end{aligned}$$

$$\begin{aligned} E[q(y, \mathbf{x}, \beta)]_{\beta \neq \beta_0} &= E[y^2 + m(\mathbf{x}, \beta)^2 - 2ym(\mathbf{x}, \beta)^2] \\ &= E[(\mathbf{x}\beta_0)^2 + u^2 + (\mathbf{x}\beta)^2 - 2(\mathbf{x}\beta_0)(\mathbf{x}\beta)] \\ &= \sigma^2 + E(\mathbf{x}\beta_0 - \mathbf{x}\beta)^2 \end{aligned}$$

Entonces,

$$E[q(\mathbf{x}, \beta_0)] < E[q(\mathbf{x}, \beta)] \text{ para todo } \beta \in \mathbf{B}, \beta \neq \beta_0.$$

# Score de la función objetivo

Si  $q(\mathbf{w}, \cdot)$  es continuamente diferenciable en el interior de  $\Theta$ , entonces con probabilidad tendiendo a uno,  $\hat{\theta}$  resuelve la condición de primer orden

$$\sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \hat{\theta}) = \mathbf{0},$$

donde

$$\mathbf{s}(\mathbf{w}, \theta)' = \nabla_{\theta} q(\mathbf{w}, \theta) = \left( \frac{\partial q(\mathbf{w}, \theta)}{\partial \theta_1}, \frac{\partial q(\mathbf{w}, \theta)}{\partial \theta_2}, \dots, \frac{\partial q(\mathbf{w}, \theta)}{\partial \theta_P} \right)'$$

es el **score de la función objetivo**, un vector  $P \times 1$ .

Esta condición se satisface en general con igualdad (pero no siempre, ver regresión por cuantiles). Para MCO siempre es con igualdad.

# Score de la función objetivo

- Si  $E[q(\mathbf{w}, \theta)]$  es continuamente diferenciable en el interior de  $\Theta$ , entonces una condición para la minimización es

$$\nabla_{\theta} E[q(\mathbf{w}, \theta)]_{\theta=\theta_0} = \mathbf{0}.$$

- Si las derivadas y las esperanzas se pueden intercambiar (Teorema de Convergencia Dominada), entonces

$$E[\nabla_{\theta} q(\mathbf{w}, \theta_0)] = E[\mathbf{s}(\mathbf{w}, \theta_0)] = \mathbf{0}.$$

- Además si  $\hat{\theta} \xrightarrow{P} \theta$ , entonces bajo las condiciones estándar de regularidad

$$N^{-1} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \hat{\theta}) \xrightarrow{P} E[\mathbf{s}(\mathbf{w}, \theta_0)].$$

# Score de la función objetivo

*Un estimador que se define por*

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \left\| \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \theta) \right\|_M$$

*se llama estimador Z (del inglés Z-estimator, por “zero”). Aquí  $\|\cdot\|_M$  es una norma pre-especificada (en general la normal euclidiana).*

*Para este caso deberíamos usar una condición de identificación diferente:*

$$E[\mathbf{s}(\mathbf{w}, \theta)] = \mathbf{0} \iff \theta = \theta_0.$$

# Hessiano de la función objetivo

- Si  $q(\mathbf{w}, \cdot)$  es diferenciable al menos dos veces en el interior de  $\Theta$ , entonces definamos el hessiano de la función objetivo como,

$$\mathbf{H}(\mathbf{w}, \theta) = \nabla_{\theta}^2 q(\mathbf{w}, \theta) = \frac{\partial^2 q(\mathbf{w}, \theta)}{\partial \theta \partial \theta'} = \frac{\partial \mathbf{s}(\mathbf{w}, \theta)}{\partial \theta'}$$

- Por el teorema del valor medio aplicado a  $q(\mathbf{w}, \theta)$  en el valor  $\theta_0$ ,

$$\sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \hat{\theta}) = \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \theta_0) + \left( \sum_{i=1}^N \dot{\mathbf{H}}_i \right) (\hat{\theta} - \theta_0),$$

donde  $\dot{\mathbf{H}}_i \equiv \mathbf{H}(\mathbf{w}_i, \tilde{\theta})$  y  $\tilde{\theta}$  es un vector donde cada elemento está en el segmento entre  $\hat{\theta}$  y  $\theta_0$ .

- Notar que la solución implica  $\sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \hat{\theta}) = \mathbf{0}$ .
- Por otro lado notar que  $N^{-1} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \hat{\theta})$  es un promedio, que converge en probabilidad a  $E[\mathbf{s}(\mathbf{w}, \theta_0)]$  dado que  $\hat{\theta} \xrightarrow{P} \theta_0$  cuando  $N \rightarrow \infty$ , usando LGN.

# Expansiones asintóticas

- Ahora con  $N^{-1} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \theta_0) \xrightarrow{P} E[\mathbf{s}(\mathbf{w}, \theta_0)] = \mathbf{0}$  también podemos aplicar la LGN.
- También vamos a usar el TCL,  $N^{-1/2} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{B}_0)$ , donde  $\mathbf{B}_0$  se define más abajo, la varianza de  $\mathbf{s}(\mathbf{w}, \theta_0)$ .
- Usando argumentos similares para el score  $\hat{\theta} \xrightarrow{P} \theta_0$ , entonces por teoría asintótica (LGN y teorema de Slutsky)

$$N^{-1} \sum_{i=1}^N \mathbf{H}(\mathbf{w}_i, \tilde{\theta}) \xrightarrow{P} E[\mathbf{H}(\mathbf{w}, \theta_0)].$$

- Para estimadores M, si  $\theta_0$  está identificado entonces  $E[\mathbf{H}(\mathbf{w}, \theta_0)]$  es definida positiva.

# Expansiones asintóticas

- La expansión en el valor real de los parámetros  $\theta_0$  juega un rol fundamental en el análisis asintótico.
- Tomemos la igualdad anterior y multipliquemos ambos lados de la igualdad por  $1/\sqrt{N}$ ,

$$0 = 1/\sqrt{N} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \theta_0) + 1/N \left( \sum_{i=1}^N \ddot{\mathbf{H}}_i \right) \sqrt{N}(\hat{\theta} - \theta_0).$$

- Reordenando los términos obtenemos la expansión asintótica de primer orden de  $\hat{\theta}$ :

$$\sqrt{N}(\hat{\theta} - \theta_0) = \left( N^{-1} \sum_{i=1}^N \ddot{\mathbf{H}}_i \right)^{-1} \left[ -N^{-1/2} \sum_{i=1}^N \mathbf{s}_i(\theta_0) \right]$$

donde  $\mathbf{s}_i(\theta_0) = \mathbf{s}(\mathbf{w}_i, \theta_0)$ .

# Expansiones asintóticas

- Además, usando  $\mathbf{A}_0 \equiv E[\mathbf{H}(\mathbf{w}, \theta_0)]$ ,

$$\begin{aligned}
 \sqrt{N}(\hat{\theta} - \theta_0) &= \mathbf{A}_0^{-1} \left[ -N^{-1/2} \sum_{i=1}^N \mathbf{s}_i(\theta_0) \right] \\
 &+ \left( \left( N^{-1} \sum_{i=1}^N \ddot{\mathbf{H}}_i \right)^{-1} - \mathbf{A}_0^{-1} \right) \left[ -N^{-1/2} \sum_{i=1}^N \mathbf{s}_i(\theta_0) \right] \\
 &= \mathbf{A}_0^{-1} \left[ -N^{-1/2} \sum_{i=1}^N \mathbf{s}_i(\theta_0) \right] + o_p(1) \cdot O_p(1) \\
 &= \mathbf{A}_0^{-1} \left[ -N^{-1/2} \sum_{i=1}^N \mathbf{s}_i(\theta_0) \right] + o_p(1)
 \end{aligned}$$

- Definamos la función de influencia de  $\hat{\theta}$  como  $\mathbf{e}(\mathbf{w}_i, \theta_0) \equiv \mathbf{e}_i(\theta_0) \equiv \mathbf{A}_0^{-1} \mathbf{s}_i(\theta_0)$ . Esta mide la “influencia” de cada observación particular  $i$  en el estimador.



# Normalidad asintótica

**Normalidad asintótica:** Además de los supuestos necesarios para consistencia, asumamos que

(a)  $\theta_0 \in \text{int}\Theta$ ;

(b)  $\mathbf{s}(\mathbf{w}, \cdot)$  es continuamente diferenciable en el interior de  $\Theta$  para todo  $\mathbf{w} \in \mathcal{W}$ ;

(c) Cada elemento en  $\mathbf{H}(\mathbf{w}, \theta)$  esta acotado en valor absoluto por una función  $b(\mathbf{w})$ , donde  $E[b(\mathbf{w})] < \infty$ ;

(d)  $\mathbf{A}_0 \equiv E[\mathbf{H}(\mathbf{w}, \theta_0)]$  es definida positiva;

(e)  $E[\mathbf{s}(\mathbf{w}, \theta_0)] = \mathbf{0}$ ; y

(f) Cada elemento en  $\mathbf{s}(\mathbf{w}, \theta_0)$  tiene momento segundo finito (varianza finita).

Entonces,

$$\sqrt{N}(\hat{\theta} - \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1}),$$

donde  $\mathbf{B}_0 \equiv E[\mathbf{s}(\mathbf{w}, \theta_0)\mathbf{s}(\mathbf{w}, \theta_0)'] = \text{Var}[\mathbf{s}(\mathbf{w}, \theta_0)]$ . Así,

$$\text{Avar}(\hat{\theta}) = \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1} / N.$$

Notemos la fórmula sandwich:  $\mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1}$ . Esta aparecerá muchas veces en el futuro.

# Ejemplo MCO

- Resolver este caso para MCO.
- Notar que

$$\mathbf{s}(\mathbf{w}, \beta) = \nabla_{\theta} q(\mathbf{w}, \beta) = -2\mathbf{x}'(y - \mathbf{x}\beta),$$

que nos da la clásica condición  $E[\mathbf{x}'u] = \mathbf{0}$ .

- El hessiano es

$$\mathbf{A}_0 = E[\mathbf{H}(\mathbf{w}, \beta_0)] = 2E[\mathbf{x}'\mathbf{x}].$$

- Si asumimos homoscedasticidad ( $E[\mathbf{u}\mathbf{u}'|\mathbf{x}] = \sigma^2\mathbf{I}$ ),

$$\mathbf{B}_0 = E[\mathbf{s}(\mathbf{w}, \beta_0)\mathbf{s}(\mathbf{w}, \beta_0)'] = 4\sigma^2 E[\mathbf{x}'\mathbf{x}]$$

- Entonces,

$$\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, \sigma^2 E[\mathbf{x}'\mathbf{x}]^{-1}).$$

# Ejemplo MCO

- Si tenemos heteroscedasticidad:

$$\mathbf{B}_0 = E[\mathbf{s}(\mathbf{w}, \beta_0)\mathbf{s}(\mathbf{w}, \beta_0)'] = 4E[\mathbf{x}'u'ux]$$

y entonces  $\mathbf{B}_0 \neq \mathbf{A}_0$ .

- En este caso,

$$\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, E[\mathbf{x}'\mathbf{x}]^{-1}E[\mathbf{x}'u'ux]E[\mathbf{x}'\mathbf{x}]^{-1}).$$

# Método delta

- La expansión asintótica se puede usar para funciones de parámetros. Consideremos la función  $\theta \mapsto g(\theta)$ . Nos gustaría saber la distribución de  $\sqrt{N}(g(\hat{\theta}) - g(\theta_0))$ .
- Asumamos que
  - $\sqrt{N}(\hat{\theta} - \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\hat{\theta}})$ ,
  - $\theta \mapsto g(\theta)$  es diferenciable en  $\theta_0$  (o sea, la primera derivada existe y es continua en  $\theta_0$ ).

Entonces,

$$\sqrt{N}(g(\hat{\theta}) - g(\theta_0)) \xrightarrow{d} N(\mathbf{0}, \nabla_{\theta}g(\theta_0) \mathbf{V}_{\hat{\theta}} \nabla_{\theta}g(\theta_0)')$$

- Por el teorema de Slutsky:

$$\nabla_{\theta}g(\hat{\theta}) \hat{\mathbf{V}}_{\hat{\theta}} \nabla_{\theta}g(\hat{\theta})' \xrightarrow{p} \nabla_{\theta}g(\theta_0) \mathbf{V}_{\hat{\theta}} \nabla_{\theta}g(\theta_0)',$$

donde  $\hat{\mathbf{V}}_{\hat{\theta}}$  es un estimador consistente de  $\mathbf{V}_{\hat{\theta}}$ .

# Ejemplo: Elasticidades de largo plazo

- Consideremos el modelo  $y_t = \beta_0 + x_t\beta_1 + y_{t-1}\alpha_1 + u_t$  donde  $(y, x)$  están en logaritmos.
- La elasticidad de corto plazo viene dado por el efecto de  $x$  en  $y$  viene dado por  $\beta_1$ .
- La elasticidad de largo plazo viene dado por el efecto de  $x$  en  $y$  viene dado por  $\frac{\beta_1}{1-\alpha_1}$ .
- Obtener la distribución de  $\frac{\hat{\beta}_1}{1-\hat{\alpha}_1}$ .

# Método de los momentos generalizados

- Supongamos el modelo de regresión lineal  $y = \mathbf{x}\beta + u$ . En este modelo podemos tener que  $\mathbf{x}$  es endógena (o algunas  $x$ s), por lo que  $\beta$  no puede ser estimado con MCO. Supongamos que  $\mathbf{x}$  es un vector  $1 \times K$  y  $\beta$  es  $K \times 1$ .  $\mathbf{X}$  es la matriz de datos  $N \times K$ .  $\mathbf{y}$  y  $\mathbf{u}$  son vectores  $N \times 1$ .
- Supongamos un conjunto de  $J$  variables exógenas,  $\mathbf{z}$  que dan lugar a las condiciones de momento,  $E(u|\mathbf{z}) = \mathbf{0}_J$ .
- Si  $J = K$  el modelo está exactamente identificado. Ejemplo:
  - $E(u|\mathbf{z}) = \mathbf{0}_J$  da lugar a las condiciones de momento  $E(\mathbf{z}'u) = \mathbf{0}_J$  (probar).
  - Entonces podemos plantear las siguientes condiciones empíricas  $\mathbf{X}'\mathbf{u}(\hat{\beta}) = \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{0}_J$ . La solución es MCO.
- Si  $J > K$  el modelo está sobre-identificado (*over-identified*). Notemos que en este caso la solución no es única. Para cada combinación de  $K$  condiciones de momento tenemos un estimador diferente...

# Método de los momentos generalizados

- Tomemos  $\mathbf{g}(\mathbf{w}, \theta)$  un vector  $J \times 1$ ,  $\mathbf{g}(\mathbf{w}, \theta) = (g_1(\mathbf{w}, \theta), \dots, g_J(\mathbf{w}, \theta))'$  que mapea  $\mathcal{W} \times \Theta \mapsto \mathbf{g}(\mathbf{w}, \theta)$ , y también  $\mathbf{g}_i(\theta) \equiv \mathbf{g}(\mathbf{w}_i, \theta)$ .
- Asumir que  $E[\mathbf{g}(\mathbf{w}_i, \theta_0)] = \mathbf{0}$ .
- Una condición necesaria mínima para la identificación es  $\theta_0$  is  $J \geq P$  (que haya al menos tantos momentos como parámetros a estimar).
- Cuando  $J = P$ , entonces  $\theta_0$  se estima por la contraparte muestral  $N^{-1} \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta_0) = \mathbf{0}$ . Este es el caso de MCO donde  $\mathbf{g}(\mathbf{w}, \theta) = \mathbf{x}'(y - \mathbf{x}\beta)$ , el modelo lineal que teníamos antes.
- Cuando  $J > P$ , entonces el Método de los momentos generalizados (GMM en inglés) usa una métrica cuadrada para minimizar  $\sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta)$ :

$$\min_{\theta \in \Theta} \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right]' \hat{\Sigma} \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right]$$

donde  $\hat{\Sigma}$  es una matriz de pesos  $J \times J$  simétrica y semidefinida positiva.

- $\hat{\Sigma}$  tiene un rol central en el modelo GMM. En particular se encarga de “pesar” cada condición de momento  $g_j$ . Ejemplo: Supongamos que  $g_1$  tiene más varianza que  $g_2$ . Entonces deberíamos usar más información de  $g_2$  que de  $g_1$ .

# Método de los momentos generalizados

- Definamos  $Q_N(\theta) = \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right]' \hat{\Xi} \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right]$ . Bajo condiciones de regularidad estándares  $Q_N(\theta)$  converge uniformemente a  $\{E[\mathbf{g}_i(\theta)]\}' \Xi_0 \{E[\mathbf{g}_i(\theta)]\}$  donde  $\hat{\Xi} \xrightarrow{P} \Xi_0$ .
- Entonces el estimador GMM es

$$\hat{\theta}_{GMM} = \operatorname{argmin}_{\theta \in \Theta} \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right]' \hat{\Xi} \left[ \sum_{i=1}^N \mathbf{g}(\mathbf{w}_i, \theta) \right].$$



# Método de los momentos generalizados

- Bajo el supuesto que  $\mathbf{g}(\mathbf{w}, \cdot)$  es continuamente diferenciable en  $\text{int}(\Theta)$ ,  $\theta_0 \in \Theta$ , entonces la condición de primer orden para  $\hat{\theta}$  es

$$\left[ \sum_{i=1}^N \nabla_{\theta} \mathbf{g}_i(\hat{\theta}) \right]' \hat{\Xi} \left[ \sum_{i=1}^N \mathbf{g}_i(\hat{\theta}) \right] \equiv \mathbf{0}.$$

- Definamos

- matriz  $J \times P$  de rango  $P$ ,  $\mathbf{G}_0 \equiv E[\nabla_{\theta} \mathbf{g}_i(\hat{\theta})]$ ;
- matriz  $P \times P$  de rango  $P$ ,  $\mathbf{A}_0 \equiv \mathbf{G}_0' \Xi_0 \mathbf{G}_0$ ;
- matriz  $P \times P$  de rango  $P$ ,  $\mathbf{B}_0 \equiv \mathbf{G}_0' \Xi_0 \Lambda_0 \Xi_0 \mathbf{G}_0$ ;
- matriz  $J \times J$  de rango  $J$ ,  $\Lambda_0 \equiv E[\mathbf{g}_i(\theta_0) \mathbf{g}_i(\theta_0)'] = \text{Var}[\mathbf{g}_i(\theta_0)]$ .

- Con manipulaciones algebraicas llegamos a

$$\mathbf{0} = \mathbf{G}_0' \Xi_0 N^{-1/2} \sum_{i=1}^N \mathbf{g}_i(\theta_0) + \mathbf{A}_0 \sqrt{N}(\hat{\theta} - \theta_0) + o_p(1)$$

- Entonces,

$$\sqrt{N}(\hat{\theta} - \theta_0) = -\mathbf{A}_0^{-1} \mathbf{G}_0' \Xi_0 N^{-1/2} \sum_{i=1}^N \mathbf{g}_i(\theta_0) + o_p(1) \xrightarrow{d} N(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

# Método de los momentos generalizados

- El estimador GMM más eficiente requiere  $\Xi_0 = \Lambda_0^{-1}$  porque

$$(\mathbf{G}'_0 \Xi_0 \mathbf{G}_0)^{-1} (\mathbf{G}'_0 \Xi_0 \Lambda_0 \Xi_0 \mathbf{G}_0) (\mathbf{G}'_0 \Xi_0 \mathbf{G}_0)^{-1} - (\mathbf{G}'_0 \Lambda_0^{-1} \mathbf{G}_0)^{-1},$$

o sea, la diferencia entre dos estimadores GMM es una matriz semidefinida positiva.

- Notemos que no podemos tener  $\Lambda_0$  antes de estimar  $\hat{\theta}$ . Sin embargo, solo necesitamos un estimador consistente,  $\tilde{\theta}$ , para estimar primero  $\Lambda_0$ , y luego se construye el estimador GMM.

# Método de los momentos generalizados

- Siguiendo con el modelo lineal sobre-identificado, podemos plantear el problema como encontrar un cero a  $\mathbf{Z}'\mathbf{u}(\beta) = \mathbf{0}_J$  donde  $\mathbf{u}(\beta) = \mathbf{y} - \mathbf{X}\beta$  es un vector  $N \times 1$ . Como  $J > P = K$  tenemos más ecuaciones que parámetros, entonces no hay solución única.
- Para transformar esto en un problema de solución única, lo transformamos en una forma cuadrática. Para cualquier matriz  $\Xi$  positiva semi-definida  $J \times J$ ,

$$\begin{aligned}\hat{\beta}_{GMM}^{\Xi} &= \arg \min_{\beta} (\mathbf{Z}'\mathbf{u}(\beta))' \Xi (\mathbf{Z}'\mathbf{u}(\beta)) = \arg \min_{\beta} (\mathbf{Z}'(\mathbf{y} - \mathbf{X}\beta))' \Xi (\mathbf{Z}'(\mathbf{y} - \mathbf{X}\beta)) \\ &= \arg \min_{\beta} \left[ \sum_{i=1}^N \mathbf{z}'_i u_i(\beta) \right]' \Xi \left[ \sum_{i=1}^N \mathbf{z}'_i u_i(\beta) \right]\end{aligned}$$

De esta manera  $J$  condiciones de momento se transformaron en una función cuadrática a ser minimizada, que tienen solución única.

- También lo podemos escribir en función de lo anterior como

$$\hat{\beta}_{GMM} = \arg \min_{\beta} Q_N(\beta)$$

# Método de los momentos generalizados

- Tomemos ahora la primera derivada de esta función y resolvamos para que evaluada en  $\hat{\beta}$  sea 0:

$$\frac{\partial Q_N(\hat{\beta})}{\partial \beta} = \mathbf{0}_J = \frac{\partial d}{\partial \hat{\mathbf{u}}} (\hat{\mathbf{u}}' (\mathbf{Z} \Xi \mathbf{Z}') \hat{\mathbf{u}}) \frac{\partial d(\mathbf{y} - \mathbf{X} \hat{\beta})}{\partial \hat{\beta}} = \frac{2}{N} \hat{\mathbf{u}}' \mathbf{Z} \Xi \mathbf{Z}' (-\mathbf{X}),$$

donde  $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X} \hat{\beta}$  y usamos  $\partial \mathbf{A} \mathbf{b} / \partial \mathbf{b} = \mathbf{A}$  y  $\partial (\mathbf{b}' \mathbf{A} \mathbf{b}) / \partial \mathbf{b} = 2 \mathbf{b}' \mathbf{A}$ , donde  $\mathbf{b}$  es un vector columna y  $\mathbf{A}$  es una matriz simétrica. Entonces tenemos,

$$\hat{\beta}_{GMM}^{\Xi} = (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{y},$$

donde enfatizamos que depende de  $\Xi$ .

# Método de los momentos generalizados

- Si las condiciones de momento son válidas, el estimador es consistente:

$$\begin{aligned}\hat{\beta}_{GMM}^{\Xi} &= (\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'\mathbf{y} = (\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'(\mathbf{X}\beta + \mathbf{u}) \\ &= \beta + (\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\Xi\mathbf{Z}'\mathbf{u} \xrightarrow{P} \beta, N \rightarrow \infty\end{aligned}$$

usando  $\frac{1}{N}\mathbf{Z}'\mathbf{u} \xrightarrow{P} \mathbf{0}_J$  y que  $\frac{1}{N}\mathbf{Z}'\mathbf{X}$  converge a una matriz no nula cuando  $N \rightarrow \infty$ .

- Este estimador no es factible porque no tenemos  $\Xi$ . Para ello necesitamos calcular la varianza asintótica...

# Método de los momentos generalizados

- Calculemos la varianza asintótica,

$$\text{Var}(\hat{\beta}_{GMM}^{\Xi}) = \text{plim}_{N \rightarrow \infty} (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{u} \mathbf{u}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X}) (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X})^{-1}.$$

- Si hacemos  $\hat{\Xi} = \frac{1}{N} \mathbf{Z}' \mathbf{u} \xrightarrow{P} E[z_i' u_i u_i' z_i] = (\Xi^*)^{-1}$ , entonces tenemos la varianza óptima:

$$\text{Var}(\hat{\beta}_{GMM}^*) = \text{plim}_{N \rightarrow \infty} (\mathbf{X}' \mathbf{Z} \Xi^* \mathbf{Z}' \mathbf{X})^{-1}.$$

Prueba: Para ver que este es el estimador más eficiente calculemos  $\text{Var}(\hat{\beta}_{GMM}^*) - \text{Var}(\hat{\beta}_{GMM}^{\Xi}) = \text{plim}_{N \rightarrow \infty} (\mathbf{X}' \mathbf{Z} \mathbf{Z}' \mathbf{u} \mathbf{u}' \mathbf{Z} \mathbf{Z}' \mathbf{X})^{-1} - (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{u} \mathbf{u}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X}) (\mathbf{X}' \mathbf{Z} \Xi \mathbf{Z}' \mathbf{X})^{-1} = (\Sigma_{XZ} \Xi^* \Sigma'_{XZ})^{-1} - (\Sigma_{XZ} \Xi \Sigma'_{XZ})^{-1} (\Sigma_{XZ} \Xi \Xi^* \Xi \Sigma'_{XZ}) (\Sigma_{XZ} \Xi \Sigma'_{XZ})^{-1}$  donde  $\Sigma_{XZ} = \frac{1}{N} \text{plim}_{N \rightarrow \infty} \mathbf{X}' \mathbf{Z}$ . Esto puede escribirse como una forma cuadrática  $= -\mathbf{D}[\mathbf{I} - \mathbf{H}(\mathbf{H}'\mathbf{H})^{-1}\mathbf{H}']\mathbf{D}$  donde  $\mathbf{D} = (\Sigma_{XZ} \Xi \Sigma'_{XZ})^{-1} (\Sigma_{XZ} \Xi (\Xi^*)^{1/2})$  y  $\mathbf{H} = (\Xi^*)^{-1/2} \Sigma_{XZ}$ .

- El estimador GMM se tiene que armar con ponderaciones que son inversamente proporcionales a la varianza de las condiciones de momento.
- Estos modelos reciben el nombre de **mínimos cuadrados generalizados (MCG)**.

# Método de los momentos generalizados

- Supongamos que  $u_i \sim i.i.d.(0, \sigma^2)$  y que  $\mathbf{Z} = \mathbf{X}$ , es decir son exógenas. Entonces tenemos los supuestos de Gauss-Markov. En este caso,  $\Xi^{-1} = \text{Var}(\mathbf{Z}'\mathbf{u}) = E(\mathbf{Z}'\mathbf{u}\mathbf{u}'\mathbf{Z}) = E[\mathbf{Z}E(\mathbf{u}\mathbf{u}'|\mathbf{Z})\mathbf{Z}'] = \sigma^2 E[\mathbf{Z}\mathbf{Z}'] = \sigma^2 E[\mathbf{X}\mathbf{X}']$ , entonces  $\hat{\beta}_{GMM}^* = \hat{\beta}_{MCO}$ .
- Supongamos que  $u_i|x_i \sim i.i.d.(0, h(x_i))$  y que  $\mathbf{Z} = \mathbf{X}$ . Entonces  $\Xi^{-1} = \text{Var}(\mathbf{Z}'\mathbf{u})$ , que no se puede simplificar.
- Supongamos que  $E(\mathbf{u}\mathbf{u}'|\mathbf{X}) = \Omega$  y que  $\mathbf{Z} = \mathbf{X}$ . Modelo de efectos aleatorios, clusters, etc. (ver datos en panel)
- Supongamos que  $\mathbf{Z} \neq \mathbf{X}$ . Variables instrumentales.

# Referencias

*Estas notas se basan en*

- *Capítulos 12, 13 y 14 de Wooldridge.*
- *Newey, W.K., y McFadden, D. (1994), "Large Sample Estimation and Hypothesis Testing," en Handbook of Econometrics, Volumen 4, ed. R.F. Engle y D. McFadden. Amsterdam: North Holland, 2111–2245.*
- *Van der Vaart, A.W. (1998), Asymptotic Statistics. Cambridge University Press.*