

Topics in Econometrics

Generalized Method of Moments: Extensions

Instructor: Ma, Jun

Renmin University of China

May 12, 2022

Moment equation models

- ▶ Let $g_i(\beta)$ be a known $l \times 1$ function of the i -th observation W_i ($g_i(\beta) = g(W_i, \beta)$) and the parameter $\beta \in \mathbb{R}^k$. A moment equation model is

$$E[g_i(\beta)] = 0.$$

We know that the true parameter β satisfies the system of equations.

- ▶ For example, in the instrumental variables model $g_i(\beta) = Z_i(Y_i - X_i^\top \beta)$ ($W_i = (Y_i, X_i, Z_i)$).
- ▶ We say the parameter is identified if there is unique β solves the equations. A necessary condition for identification is $l \geq k$.
- ▶ $l = k$: just identified;
- ▶ $l > k$: over-identified.

Method of moments

- ▶ We consider the just identified case: $l = k$
- ▶ The sample analogue of $E[g_i(\beta)]$:

$$\bar{g}_n(\beta) = \frac{1}{n} \sum_{i=1}^n g_i(\beta).$$

- ▶ The method of moments estimator (MME) $\hat{\beta}_{\text{mm}}$ for β is the solution to

$$\frac{1}{n} \sum_{i=1}^n g_i(\hat{\beta}_{\text{mm}}) = 0.$$

Overidentified moment equations

- Define

$$\bar{g}_n(b) = \frac{1}{n} \sum_{i=1}^n g_i(b).$$

- We defined the MME $\hat{\beta}$ for β to be the solution to $\bar{g}_n(\hat{\beta}) = 0$. However, if the model is over-identified, there are more equations than parameters. The MME is not defined.
- We cannot find an estimator $\hat{\beta}$ which sets $\bar{g}_n(\hat{\beta}) = 0$ but we can try to find an estimator $\hat{\beta}$ which makes $\bar{g}_n(\hat{\beta})$ as close to zero as possible.

- ▶ Let W be an $l \times l$ positive definite weight matrix. The GMM criterion function is

$$J(b) = n \cdot \bar{g}_n(b)^\top W \bar{g}_n(b) .$$

- ▶ When $W = I_l$ (l -dimensional identity matrix),
 $J(b) = n \cdot \bar{g}_n(b)^\top \bar{g}_n(b) = n \cdot \|\bar{g}_n(b)\|^2 .$
- ▶ The Generalized method of moments (GMM) estimator is
 $\hat{\beta}_{\text{gmm}} = \operatorname{argmin}_b J_n(b) .$

Asymptotic distribution

- Asymptotic distribution of the GMM estimator

$$\sqrt{n} \left(\hat{\beta}_{\text{gmm}} - \beta \right) \rightarrow_d N(0, V_W).$$

where

$$V_W = \left(Q^\top W Q \right)^{-1} \left(Q^\top W \Omega W Q \right) \left(Q^\top W Q \right)^{-1}$$

with

$$\Omega = E \left[g_i(\beta) g_i(\beta)^\top \right] \text{ and } Q = E \left[\frac{\partial}{\partial b^\top} g_i(b) \Big|_{b=\beta} \right]$$

- If the efficient weight matrix $W = \Omega^{-1}$ is used then

$$V_\beta = \left(Q^\top \Omega^{-1} Q \right)^{-1}.$$

Efficient GMM

- The efficient GMM estimator can be constructed by using

$$\hat{\Omega} = \frac{1}{n} \sum_{i=1}^n g_i(\tilde{\beta}) g_i(\tilde{\beta})^\top - \bar{g}_n(\tilde{\beta}) \bar{g}_n(\tilde{\beta})^\top,$$

with a preliminary consistent estimator $\tilde{\beta}$.

- The asymptotic covariance matrix can be estimated by sample counterparts of the population matrices.

Continuously-updated GMM

- ▶ An alternative to the two-step GMM estimator can be constructed by letting the weight matrix be an explicit function of b :

$$J(b) = n \cdot \bar{g}_n(b)^\top \left(\frac{1}{n} \sum_{i=1}^n g_i(b) g_i(b)^\top \right)^{-1} \bar{g}_n(b)$$

or

$$J(b) = \bar{g}_n(b)^\top \left(\frac{1}{n} \sum_{i=1}^n g_i(b) g_i(b)^\top - \bar{g}_n(b) \bar{g}_n(b)^\top \right)^{-1} \bar{g}_n(b).$$

- ▶ The $\hat{\beta}$ which minimizes this function is the CU-GMM estimator. The minimization requires numerical methods.
- ▶ We have:

$$\sqrt{n} \left(\hat{\beta}_{\text{cu-gmm}} - \beta \right) \rightarrow_d N(0, V_\beta).$$

Wald statistic

- ▶ The parameter of interest θ is a function of the coefficients, $\theta = r(\beta)$ for some function $r : \mathbb{R}^k \rightarrow \mathbb{R}^q$. The estimator of θ is given by $\hat{\theta} = r(\hat{\beta})$.
- ▶ If $r(\cdot)$ is continuous at the true value of β , then $\hat{\theta} \rightarrow_p \theta$. Suppose that $r : \mathbb{R}^k \rightarrow \mathbb{R}^q$ is continuously differentiable at the true value of β and $R = \partial r(b)^\top / \partial b \big|_{b=\beta}$ has rank q . Then, $\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d N(0, V_\theta)$ where $V_\theta = R^\top V_\beta R$.
- ▶ Consider the Wald statistic

$$W(\theta) = n(\hat{\theta} - \theta)^\top \hat{V}_\theta^{-1}(\hat{\theta} - \theta),$$

where \hat{V}_θ is a consistent estimator of V_θ . Then, $W(\theta) \rightarrow_d \chi_q^2$.

Confidence set

- ▶ A confidence region \hat{C} is a set estimator for $\theta \in \mathbb{R}^q$. A natural confidence region is

$$\hat{C} = \{\theta \in \mathbb{R}^q : W(\theta) \leq c_{1-\alpha}\},$$

with $c_{1-\alpha}$ being the $1 - \alpha$ quantile of the χ_q^2 distribution:
 $F_{\chi_q^2}(c_{1-\alpha}) = 1 - \alpha$.

- ▶ Then,

$$\Pr[\theta \in \hat{C}] \rightarrow \Pr[\chi_q^2 \leq c_{1-\alpha}] = 1 - \alpha.$$

- ▶ Note that the shape of the confidence set \hat{C} is predetermined (i.e., ellipse).

OverIdentification test

- Consider the linear IV model:

$$\begin{aligned}Y_i &= X_i^\top \beta + e_i \\ \mathbb{E}[e_i Z_i] &= 0,\end{aligned}$$

where $X_i \in \mathbb{R}^k$ and $Z_i \in \mathbb{R}^l$. The model is over-identified:
 $l > k$.

- The model specifies

$$\mathbb{E}[e_i Z_i] = 0 \iff \mathbb{E}[Z_i Y_i] = \mathbb{E}[Z_i X_i^\top] \beta.$$

- This is equivalent to saying that $\mathbb{E}[Z_i Y_i]$ is in the column space of $\mathbb{E}[Z_i X_i^\top]$. The model imposes a restriction on the distribution of the observed variables (Y_i, X_i, Z_i) .
- Since β is of dimension $k < l$, it is not certain if such a vector exists. In such a case, we say that the model is misspecified.

- Suppose that $X_i \in \mathbb{R}^1$ and $Z_i = \left(Z_i^{(1)}, Z_i^{(2)} \right)^\top \in \mathbb{R}^2$. Then the model specifies

$$\begin{aligned} \mathrm{E} \left[Z_i^{(1)} Y_i \right] &= \mathrm{E} \left[Z_i^{(1)} X_i \right] \beta \\ \mathrm{E} \left[Z_i^{(2)} Y_i \right] &= \mathrm{E} \left[Z_i^{(2)} X_i \right] \beta, \end{aligned}$$

which requires

$$\frac{\mathrm{E} \left[Z_i^{(1)} Y_i \right]}{\mathrm{E} \left[Z_i^{(1)} X_i \right]} = \frac{\mathrm{E} \left[Z_i^{(2)} Y_i \right]}{\mathrm{E} \left[Z_i^{(2)} X_i \right]}.$$

- The true distribution of (Y_i, X_i, Z_i) may violate this condition.
- We can do a hypothesis test of the model specification. This is known as the overidentification test:

$$H_0 : \text{There exists } \beta \in \mathbb{R}^k \text{ such that } \mathrm{E} \left[Z_i \left(Y_i - X_i^\top \beta \right) \right] = 0.$$

- For the more general model, the null hypothesis of correct model specification is

$$H_0 : \text{There exists } \beta \in \mathbb{R}^k \text{ such that } E[g_i(\beta)] = 0.$$

- H_0 is true if and only if

$$\min_b n \cdot E[g_i(b)]^\top \Omega^{-1} E[g_i(b)] = 0.$$

- We estimate $\min_b n \cdot E[g_i(b)]^\top \Omega^{-1} E[g_i(b)]$ by

$$\min_b n \cdot \bar{g}_n(b)^\top \hat{\Omega}^{-1} \bar{g}_n(b).$$

and if it is large, we reject H_0 .

- The test statistic is just $J(\hat{\beta}_{\text{gmm}})$. This is known as the J -statistic. The overidentification test is referred to as the Sargan test.

- Under H_0 , $J(\hat{\beta}_{\text{gmm}}) \rightarrow_d \chi_{l-k}^2$. We reject H_0 if

$$J(\hat{\beta}_{\text{gmm}}) > c_{1-\alpha} \text{ with } c_{1-\alpha} \text{ being the } 1 - \alpha \text{ quantile of the } \chi_{l-k}^2 \text{ distribution: } F_{\chi_{l-k}^2}(c_{1-\alpha}) = 1 - \alpha.$$

Maximum likelihood

- ▶ Let (X_1, \dots, X_n) be a random (i.i.d.) sample on a continuous with a density function $f(\cdot; \theta)$, $\theta \in \Theta \subseteq \mathbb{R}^k$. Let x_i be the observed value of X_i . Then we call

$$L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta)$$

the likelihood function of θ given (x_1, x_2, \dots, x_n) , and we call the value of θ that maximizes $L(\theta; X_1, \dots, X_n)$ the maximum likelihood (ML) estimator.

- ▶ The log-likelihood function:

$$\ell(\theta; x_1, \dots, x_n) = \sum_{i=1}^n \log f(x_i; \theta).$$

- ▶ The ML estimator: $\hat{\theta}_{\text{ml}} = \operatorname{argmax}_{\theta \in \Theta} \ell(\theta; X_1, \dots, X_n)$.
- ▶ The model $\{f(\cdot; \theta) : \theta \in \Theta\}$ is correctly specified if there exists $\theta_* \in \Theta$ so that $f(\cdot; \theta_*) = f_X$, where f_X denotes the true density of X_i .

Kullback–Leibler divergence

- The Kullback–Leibler (KL) divergence from a density f to another density g :

$$\mathbb{D}_{\text{kl}}(f \mid g) = \int \log \left(\frac{f(x)}{g(x)} \right) f(x) \, dx.$$

- $\mathbb{D}_{\text{kl}}(f \mid g) \geq 0$ and $\mathbb{D}_{\text{kl}}(f \mid g) = 0$ if and only if $f = g$.
- Jensen's inequality: Let X be a random variable and h be a strictly concave function. That is,

$$h(\lambda a + (1 - \lambda)b) > \lambda h(a) + (1 - \lambda)h(b)$$

for any $a < b$ and $0 < \lambda < 1$. Then $\mathbb{E}[h(X)] < h(\mathbb{E}[X])$.

- If $f \neq g$,

$$\begin{aligned} \int \log \left(\frac{f(x)}{g(x)} \right) f(x) \, dx &= - \int \log \left(\frac{g(x)}{f(x)} \right) f(x) \, dx \\ &> -\log \left(\int g(x) \, dx \right) = 0. \end{aligned}$$

- $\mathbb{D}_{\text{kl}}(f_X \mid f(\cdot; \theta)) \geq 0$ and $\mathbb{D}_{\text{kl}}(f_X \mid f(\cdot; \theta_*)) = 0$. This is equivalent to

$$\begin{aligned}\theta_* &= \operatorname{argmin}_{\theta \in \Theta} \mathbb{D}_{\text{kl}}(f_X \mid f(\cdot; \theta)) \\ &= \operatorname{argmax}_{\theta \in \Theta} - \int \log \left(\frac{f_X(x)}{f(x; \theta)} \right) f_X(x) \, dx \\ &= \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}[\log f(X; \theta)].\end{aligned}$$

- A natural estimator from the perspective of KL divergence is given by $\operatorname{argmax}_{\theta \in \Theta} Q_n(\theta)$ with $Q_n(\theta) = n^{-1} \sum_{i=1}^n \log f(X_i; \theta)$, which is just the ML estimator.
- By LLN, we know that for each θ , $Q_n(\theta) \rightarrow_p Q(\theta)$ where $Q(\theta) = \mathbb{E}[\log f(X; \theta)]$. The ML estimator is defined to be the maximizer of $Q_n(\theta)$. We expect the maximizer should converge to the maximizer of its limit $Q(\theta)$ in probability.

Nonparametric likelihood

- ▶ The moment equation model is nonparametric in the sense that we do not fully specify the distribution of the observed variables.
- ▶ Rather than specifying a parametric model for X_i , we assume the variables follow a discrete distribution supported on the observations X_1, \dots, X_n .
- ▶ The parameters corresponding to this “model” is p_1, \dots, p_n with $(p_1, \dots, p_n) \in \Delta$, where

$$\Delta = \left\{ (p_1, \dots, p_n) : \sum_{i=1}^n p_i = 1, p_i \geq 0, i = 1, \dots, n \right\}.$$

- ▶ The nonparametric log-likelihood is

$$\ell(p_1, \dots, p_n; X_1, \dots, X_n) = \sum_{i=1}^n \log(n \cdot p_i), (p_1, \dots, p_n) \in \Delta.$$

- ▶ The maximum of the above log-likelihood function is attained at $p_i = 1/n, \forall i$, which is the empirical distribution.

- ▶ Maximizing the nonparametric log-likelihood is equivalent to minimizing the KL divergence from the empirical distribution $(1/n, \dots, 1/n)$ to (p_1, \dots, p_n) : $\sum_{i=1}^n n^{-1} \log (n^{-1}/p_i)$.
- ▶ Consider the moment equation model:

$$\mathbb{E} [g_i (\beta)] = \int g (w, \beta) f_W (w) \mathrm{d}w = 0,$$

where f_W denotes the true density of W_i .

- ▶ The model imposes a restriction on f_W .
- ▶ The empirical likelihood method is a constrained nonparametric likelihood with the constraint $\sum_{i=1}^n p_i g_i (b) = 0$ imposed.

Empirical likelihood (EL)

- ▶ The EL criterion function:

$$\ell_{\text{el}}(b) := \max_{p_1, \dots, p_n} 2 \sum_{i=1}^n \log(n \cdot p_i)$$

subject to $\sum_{i=1}^n p_i g_i(b) = 0, (p_1, \dots, p_n) \in \Delta.$

- ▶ The EL estimator $\hat{\beta}_{\text{el}}$ is the maximizer of $\ell_{\text{el}}(b)$.
- ▶ $\sqrt{n}(\hat{\beta}_{\text{el}} - \beta) \rightarrow_d N(0, V_\beta)$, where V_β is the asymptotic variance of the efficient GMM (Qin and Lawless, 1994).
- ▶ The EL estimator is efficient and avoids estimating the optimal weighting matrix in the first step.

EL ratio inference

- ▶ The EL ratio statistic:

$$LR(\theta) = \max_b \ell_{\text{el}}(b) - \max_{r(b)=\theta} \ell_{\text{el}}(b).$$

Then, $LR(\theta) \rightarrow_d \chi_q^2$.

- ▶ Estimation of the asymptotic variance is not needed.
- ▶ The EL confidence set:

$$\hat{C} = \{\theta \in \mathbb{R}^q : LR(\theta) \leq c_{1-\alpha}\}.$$

- ▶ The shape of the EL confidence set is data-driven.
- ▶ The EL method has many other favorable properties relative to efficient GMM. See Kitamura (2006) for a review.

Duality

- ▶ It seems that the high dimensionality of the parameter space makes the maximization problem infeasible in practice.
- ▶ Instead of directly solving it, we fix b first and use the Lagrange multiplier method to solve

$$\begin{aligned}\ell_{\text{el}}(b) &:= \max_{p_1, \dots, p_n} 2 \sum_{i=1}^n \log(n \cdot p_i) \\ &\text{subject to } \sum_{i=1}^n p_i g_i(b) = 0, (p_1, \dots, p_n) \in \Delta.\end{aligned}$$

- ▶ The Lagrangian associated with the constrained optimization problem is

$$\mathcal{L}(p_1, \dots, p_n, \lambda) = \sum_{i=1}^n \log(p_i) + \gamma \left(1 - \sum_{i=1}^n p_i \right) - n \cdot \lambda^\top \sum_{i=1}^n p_i g_i(b),$$

where $\gamma \in \mathbb{R}$ and $\lambda \in \mathbb{R}^l$ are Lagrange multipliers.

- The first-order conditions:

$$0 = \frac{1}{p_i} - \gamma - n \left(\lambda^\top g_i(b) \right)$$

$$0 = 1 - \sum_{i=1}^n p_i$$

$$0 = n \sum_{i=1}^n p_i g_i(b).$$

- The first-order conditions are solved by $\gamma = n$ and $(p_1, \dots, p_n, \lambda)$ are given by the solution to

$$p_i = \frac{1}{n(1 + \lambda^\top g_i(b))}$$
$$0 = \sum_{i=1}^n \frac{g_i(b)}{1 + \lambda^\top g_i(b)}.$$

- The l equations $0 = \sum_{i=1}^n g_i(b) / (1 + \lambda^\top g_i(b))$ are the first-order conditions of the convex minimization problem $\min_{\lambda} - \sum_{i=1}^n \log(1 + \lambda^\top g_i(b))$.

- The EL estimator is therefore

$$\hat{\beta}_{\text{el}} = \operatorname{argmax}_b \min_{\lambda} - \sum_{i=1}^n \log \left(1 + \lambda^{\top} g_i(b) \right).$$

- For fixed b , $\min_{\lambda} - \sum_{i=1}^n \log \left(1 + \lambda^{\top} g_i(b) \right)$ is a convex minimization problem, for which a simple Newton algorithm works.
- The maximization of $\min_{\lambda} - \sum_{i=1}^n \log \left(1 + \lambda^{\top} g_i(b) \right)$ with respect to b is harder to solve. It is solved by a nonlinear optimization algorithm.

Implied probabilities

- Once $\hat{\beta}_{\text{el}}$ is calculated, we get the implied probabilities

$$\hat{p}_i = \frac{1}{n \left(1 + \hat{\lambda}^\top g_i \left(\hat{\beta}_{\text{el}} \right) \right)},$$

where $\hat{\lambda}$ is the solution of the equations

$$0 = \sum_{i=1}^n \frac{g_i \left(\hat{\beta}_{\text{el}} \right)}{1 + \hat{\lambda}^\top g_i \left(\hat{\beta}_{\text{el}} \right)}.$$

- Suppose that we are interested in estimating $\text{E} [h (W_i)]$, where $h (\cdot)$ is a known function.
- $\sum_i \hat{p}_i h (W_i)$ is an efficient estimator of $\text{E} [h (W_i)]$ relative to the sample mean $n^{-1} \sum_i h (W_i)$ (Brown and Newey, 1998).
- $(\hat{p}_1, \dots, \hat{p}_n)$ is also a more efficient estimator than the empirical distribution, from which we do bootstrap resampling.

- ▶ We have a small dataset on $W_i = (Y_i, X_i)$, $i = 1, \dots, n$, but X_i includes a rich set of variables so that the regression model is not suffering from the omitted variable bias.
- ▶ Suppose that M_i is the vector collecting a small subset of variables in W_i . We have another auxiliary dataset on M_i . Such a dataset has a very large sample size N .
- ▶ We can calculate the implied probabilities

$$\hat{p}_i = \frac{1}{n \left(1 + \hat{\lambda}^\top (M_i - \overline{M}) \right)}$$

where $\hat{\lambda}$ is the solution of the equations

$$0 = \sum_{i=1}^n \frac{M_i - \overline{M}}{1 + \hat{\lambda}^\top (M_i - \overline{M})},$$

where \overline{M} is the sample mean of M_i computed by using the auxiliary dataset.

- ▶ The reweighted estimator $\left(\sum_{i=1}^n \hat{p}_i X_i X_i^\top \right)^{-1} \left(\sum_{i=1}^n \hat{p}_i X_i Y_i \right)$ is more efficient than the OLS (Hellerstein and Imbens, 1999).

Cressie-Read divergence

- ▶ EL can be thought of as minimizing the KL divergence (distance) of the empirical distribution and the discrete distribution supported on the sample with a constraint.
- ▶ We can consider other distance. E.g., $\sum_{i=1}^n p_i \log(n^{-1}/p_i)$ (reverse KL divergence, exponential tilting, Kitamura and Stutzer, 1998) and $\sum_{i=1}^n (n \cdot p_i - 1)^2$ (Euclidean distance, continuously-updated GMM/Euclidean likelihood).
- ▶ Cressie-Read divergence:

$$\frac{1}{\gamma(\gamma+1)} \sum_{i=1}^n [(n \cdot p_i)^{-\gamma} - 1], \gamma \in \mathbb{R}.$$

Special cases: $\gamma = -2$, continuously-updated GMM; $\gamma = -1$, exponential tilting; $\gamma = 0$, EL among many others.

- ▶ In the literature, various papers show that some method has certain advantages over other methods, from different perspectives.