

GMM and SMM notes and comments for Tyler Ransom

Richard W. Evans

October 2, 2020

1 GMM

1.0 GMM General Specification

My most general definition of GMM is the following. Let \mathbf{x} be an $N \times M$ matrix of data with N observations (rows) and M variables (columns) and typical element $x_{n,m}$. Let $\boldsymbol{\theta} \in \mathbb{R}^K$ be the K -dimensional parameter vector to be estimated, and let $\hat{\boldsymbol{\theta}}_{GMM}$ be the GMM estimator of that parameter vector. The GMM estimator is the one that minimizes some general norm (distance measure) on an R -dimensional vector of model moments $\mathbf{m}(\mathbf{x}|\boldsymbol{\theta})$ and a corresponding R -dimensional vector of data moments $\mathbf{m}(\mathbf{x})$ such that $R \geq K$.

$$\hat{\boldsymbol{\theta}}_{GMM} : \min_{\boldsymbol{\theta}} \|\mathbf{m}(\mathbf{x}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x})\| \quad (1)$$

The most common GMM criterion function or choice of norm is a weighted sum of squared moment errors. We can first define an $R \times 1$ error vector $\mathbf{e}(\mathbf{x}, \boldsymbol{\theta})$ as the following function of the data \mathbf{x} and the parameter vector $\boldsymbol{\theta}$ which can either be given in levels or in percentage differences (there are advantages to each).

$$\mathbf{e}(\mathbf{x}, \boldsymbol{\theta}) \equiv \begin{cases} \mathbf{m}(\mathbf{x}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x}) \\ \text{or} \\ \frac{\mathbf{m}(\mathbf{x}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x})}{\mathbf{m}(\mathbf{x})} \end{cases} \quad (2)$$

The most common weighted sum of squared moment errors GMM specification takes the following form,

$$\hat{\boldsymbol{\theta}}_{GMM} : \min_{\boldsymbol{\theta}} \mathbf{e}(\mathbf{x}, \boldsymbol{\theta})^T \mathbf{W} \mathbf{e}(\mathbf{x}, \boldsymbol{\theta}) \quad (3)$$

where \mathbf{W} is an $R \times R$ weighting matrix. The standard arguments regarding asymptotic biasedness and efficiency apply to the weighting matrix. An optimal weighting matrix is most efficient (inverse variance estimators, take care of heteroscedasticity and autocorrelation), but identity matrix is unbiased.

1.1 GMM Linear Regression: Method 1

Define a specific linear regression model as the following equation with the following assumptions on the error terms,

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots \beta_K x_{K,i} + \varepsilon_i$$

$$\text{where } \varepsilon_i \sim f(\varepsilon), \quad E[\varepsilon_i] = 0, \quad \text{and} \quad E[x_{k,i}\varepsilon_i] = 0 \quad \forall k \quad (4)$$

where the data are the matrix of $i = 1, 2, \dots, N$ observations on $K + 1$ variables $(y_i, x_{1,i}, x_{2,i}, \dots, x_{K,i})$, the parameter vector to be estimated is the $K + 1$ vector $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \dots, \beta_K)$, and the error terms ε_i are i.i.d. according to the general pdf $f(\varepsilon)$ with the moment conditions $E[\varepsilon_i] = 0$ and $E[x_{k,i}\varepsilon_i] = 0$ for all k .

The most standard way of teaching GMM using linear regression is to use the $K + 1$ moments implied by the unconditional expected value condition $E[\varepsilon_i] = 0$ and the K conditional expected value conditions $E[x_{k,i}\varepsilon_i] = 0$ for all k . To map this into the notation in Section 1.0, let the vector of data moments be the vector of $K + 1$ zeros assumed by the $K + 1$ zero-expected values on ε_i .

$$\mathbf{m}(\mathbf{y}, \mathbf{x}) = \begin{bmatrix} E[\varepsilon_i] \\ E[x_{1,i}\varepsilon_i] \\ E[x_{2,i}\varepsilon_i] \\ \vdots \\ E[x_{K,i}\varepsilon_i] \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (5)$$

These $K \times 1$ data moments are zero by assumption and do not actually use the data \mathbf{y} or \mathbf{x} in this case.

The model moments from Section 1.0 were specified as a function of data and parameters $\mathbf{m}(\mathbf{x}|\boldsymbol{\theta})$. In this linear regression case, the $K \times 1$ model moments m_k are defined by solving the linear regression model (4) for ε_i and taking the sample averages that correspond to the assumed unconditional and conditional expectations on ε_i .

$$\mathbf{m}(\mathbf{y}, \mathbf{x}|\boldsymbol{\beta}) = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i} - \dots \beta_K x_{K,i}) \\ \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i} - \dots \beta_K x_{K,i}) x_{1,i} \\ \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i} - \dots \beta_K x_{K,i}) x_{2,i} \\ \vdots \\ \frac{1}{N} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i} - \dots \beta_K x_{K,i}) x_{K,i} \end{bmatrix} \quad (6)$$

Define the moment error vector, analogous to (2) above, as the following function of data and parameters.

$$\mathbf{e}(\mathbf{y}, \mathbf{x}, \boldsymbol{\beta}) \equiv \mathbf{m}(\mathbf{y}, \mathbf{x}|\boldsymbol{\beta}) - \mathbf{m}(\mathbf{y}, \mathbf{x}) = \mathbf{m}(\mathbf{y}, \mathbf{x}|\boldsymbol{\beta}) \quad (7)$$

Note that we cannot use the percent deviation form of the moment errors shown in the second specification of (2) because the denominator is zero.

The GMM estimator in this case $\hat{\beta}_{GMM1}$ is the following exactly identified problem of choosing $K+1$ parameters to minimize the weighted sum of squared errors in $K+1$ moment conditions,

$$\hat{\beta}_{GMM1} : \min_{\beta} e(\mathbf{y}, \mathbf{x}, \beta)^T \mathbf{W} e(\mathbf{y}, \mathbf{x}, \beta) \quad (8)$$

where \mathbf{W} is a $(K+1 \times K+1)$ weighting matrix.

1.2 GMM Linear Regression: Method 2

The second method for estimating the parameter vector β by GMM is just weighted linear least squares. For nonlinear model cases, this is just weighted nonlinear least squares. This ends up being an overidentified GMM estimation in which each error term is treated as a model moment. This approach does not really use any of the expected value of ε_i conditions in (4). It just uses the idea that we want to choose the parameter vector β to maximize the fit of the predicted values of the endogenous variable \hat{y}_i to the data values of the endogenous variable y_i . This is equivalent to minimizing the absolute value of the error terms ε_i .

Similar to Method 1 in Section 1.1, we assume the data moments are a vector of zeros. However, in this case, the data moment vector is length N , a zero for each observation of ε_i .

$$\mathbf{m}(\mathbf{y}, \mathbf{x}) = \begin{bmatrix} \min |\varepsilon_1| \\ \min |\varepsilon_2| \\ \vdots \\ \min |\varepsilon_N| \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (9)$$

Similar to (5), these $N \times 1$ data moments are zero by assumption and do not actually use the data \mathbf{y} or \mathbf{x} in this case.

The model moments are simply the model error terms found by solving linear regression equation (4) for ε_i for each of the $i = 1, 2, \dots, N$ observations. Note that each model error term is a function of data \mathbf{y} and \mathbf{x} and parameters β .

$$\mathbf{m}(\mathbf{y}, \mathbf{x}|\beta) = \begin{bmatrix} y_1 - \beta_0 - \beta_1 x_{1,1} - \beta_2 x_{2,1} - \dots \beta_K x_{K,1} \\ y_2 - \beta_0 - \beta_1 x_{1,2} - \beta_2 x_{2,2} - \dots \beta_K x_{K,2} \\ \vdots \\ y_N - \beta_0 - \beta_1 x_{1,N} - \beta_2 x_{2,N} - \dots \beta_K x_{K,N} \end{bmatrix} \quad (10)$$

Define the $N \times 1$ moment error vector, analogous to (2), as the following function of data and parameters.

$$e(\mathbf{y}, \mathbf{x}, \beta) \equiv \mathbf{m}(\mathbf{y}, \mathbf{x}|\beta) - \mathbf{m}(\mathbf{y}, \mathbf{x}) = \mathbf{m}(\mathbf{y}, \mathbf{x}|\beta) \quad (11)$$

As with (7), we cannot use the percent deviation form of the moment errors shown in the second specification of (2) because the denominator is zero.

The GMM estimator in this case $\hat{\beta}_{GMM2}$ is the following overidentified problem of choosing $K + 1$ parameters to minimize the weighted sum of squared errors in N moment conditions,

$$\hat{\beta}_{GMM2} : \min_{\beta} e(\mathbf{y}, \mathbf{x}, \beta)^T \mathbf{W} e(\mathbf{y}, \mathbf{x}, \beta) \quad (12)$$

where \mathbf{W} is a $(N \times N)$ weighting matrix.

1.3 GMM Comparison of Methods

2 SMM

2.0 SMM General Specification

Simulated method of moments (SMM) is analogous to the generalized method of moments (GMM) estimator. SMM could really be thought of as a particular type of GMM estimator. The SMM estimator chooses model parameters θ to make simulated model moments match data moments. Seminal papers developing SMM are [McFadden \(1989\)](#), [Lee and Ingram \(1991\)](#), and [Duffie and Singleton \(1993\)](#). Good textbook treatments of SMM are found in [Adda and Cooper \(2003, pp.87-100\)](#) and [Davidson and MacKinnon \(2004, pp. 383-394\)](#).

In the general specification of GMM estimation from Section 1.0, we use data \mathbf{x} and model parameters θ to minimize the distance between model moments $\mathbf{m}(\mathbf{x}|\theta)$ and data moments $\mathbf{m}(\mathbf{x})$.

$$\hat{\theta}_{GMM} : \min_{\theta} \|\mathbf{m}(\mathbf{x}|\theta) - \mathbf{m}(\mathbf{x})\| \quad (1)$$

The following difficulties can arise with GMM making it not possible or very difficult.

- The model moment function $\mathbf{m}(\mathbf{x}|\theta)$ is not known analytically.
- The data moments you are trying to match come from another model (indirect inference).
- The model moments $\mathbf{m}(\mathbf{x}|\theta)$ are derived from latent variables that are not observed by the modeler. You only have moments, not the underlying data. See [Laroque and Salanie \(1993\)](#).
- The model moments $\mathbf{m}(\mathbf{x}|\theta)$ are derived from censored variables that are only partially observed by the modeler.
- The model moments $\mathbf{m}(\mathbf{x}|\theta)$ are just difficult to derive analytically. Examples include moments that include multiple integrals over nonlinear functions as in [McFadden \(1989\)](#).

Let \mathbf{x} be the $N \times M$ vector of data as in Section 1.0, and let $\tilde{\mathbf{x}}_s$ be the s th simulation of the dataset \mathbf{x} . SMM estimation is simply to simulate the model data S times, and use the average values of the moments from the simulated data as the estimator for the model moments. Let $\tilde{\mathbf{x}} = \{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_s, \dots, \tilde{\mathbf{x}}_S\}$ be the S simulations of the model data.

It is important to note a major difference between SMM and GMM. In the GMM examples in Section 1, the parameter vector $\boldsymbol{\theta}$ could be estimated without strong distributional assumptions about the random variable error term ε_i . In the first linear regression example in Section 1.1 with the $K + 1$ moments, all that was needed was that the errors be distributed i.i.d. according to general pdf $\varepsilon_i \sim f(\varepsilon)$ and the two sets of moment conditions on the distribution of the errors $E[\varepsilon_i] = 0$ and $E[\varepsilon_i x_{k,i}] = 0$. In SMM, we need to know the exact distribution $f(\varepsilon)$ in order to simulate the data. For example, if we assume the errors are distributed normally, we need to know the parameters of the distribution (variance σ^2 in the case of the normal distribution). We need to be able

The $R \times 1$ data moment vector $\mathbf{m}(\mathbf{x})$ is the same as in the GMM specification of Section 1.0. We can compute a set of R model moments $\mathbf{m}(\tilde{\mathbf{x}}_s|\boldsymbol{\theta})$ from each of the simulated datasets $\tilde{\mathbf{x}}_s$. One conceptual advantage of SMM relative to GMM is that the model moments are computed in the same way as the data moments because the model moments are computed from simulated datasets $\tilde{\mathbf{x}}_s$ that are supposed to be representative of the real world dataset \mathbf{x} . The final model moment vector $\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta})$ to be used in SMM is an average across each of the moment vectors for each simulation $\mathbf{m}(\tilde{\mathbf{x}}_s|\boldsymbol{\theta})$.

$$\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta}) \equiv \frac{1}{S} \sum_{s=1}^S \mathbf{m}(\tilde{\mathbf{x}}_s|\boldsymbol{\theta}) \quad (13)$$

Once we have an estimate of the model moments $\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta})$ from our S simulations, SMM estimation is very similar to our presentation of GMM. The SMM approach of estimating the parameter vector $\hat{\boldsymbol{\theta}}_{SMM}$ is to choose $\boldsymbol{\theta}$ to minimize some distance measure of the simulated model moments $\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta})$ from the data moments $\mathbf{m}(\mathbf{x})$.

$$\hat{\boldsymbol{\theta}}_{SMM} : \min_{\boldsymbol{\theta}} \|\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x})\| \quad (14)$$

The distance measure $\|\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x})\|$ can be any kind of norm. But it is important to recognize that your estimates $\hat{\boldsymbol{\theta}}_{SMM}$ will be dependent on what distance measure (norm) you choose. The most widely studied and used distance metric in GMM and SMM estimation is the L^2 norm or the weighted sum of squared errors in moments. Define the moment error function $e(\tilde{\mathbf{x}}, \mathbf{x}|\boldsymbol{\theta})$ as the simple difference or percent difference in the vector of simulated model moments from the data moments.

$$e(\tilde{\mathbf{x}}, \mathbf{x}, \boldsymbol{\theta}) \equiv \begin{cases} \hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x}) \\ \text{or} \\ \frac{\hat{\mathbf{m}}(\tilde{\mathbf{x}}|\boldsymbol{\theta}) - \mathbf{m}(\mathbf{x})}{\mathbf{m}(\mathbf{x})} \end{cases} \quad (15)$$

The weighted sum of squared simulated moment errors SMM specification takes the following form,

$$\hat{\theta}_{SMM} : \min_{\theta} e(\tilde{x}, x, \theta)^T W e(\tilde{x}, x, \theta) \quad (16)$$

where W is a $R \times R$ weighting matrix in the criterion function. We call the quadratic form expression $e(\tilde{x}, x|\theta)^T W e(\tilde{x}, x|\theta)$ the criterion function because it is a strictly positive scalar that is the object of the minimization in the SMM problem statement. The $R \times R$ weighting matrix W in the criterion function allows the econometrician to control how each moment is weighted in the minimization problem. For example, an $R \times R$ identity matrix for W would give each moment equal weighting, and the criterion function would be a simply sum of squared percent deviations (errors). Other weighting strategies can be dictated by the nature of the problem or model.

One last item to emphasize with SMM is that the errors that are drawn for the S simulations of the model must be drawn only once so that the minimization problem for $\hat{\theta}_{SMM}$ does not have the underlying sampling changing for each guess of a value of θ . Put more simply, you want the random draws for all the simulations to be held constant so that the only thing changing in the minimization problem is the value of the vector of parameters θ .

2.1 SMM Linear Regression: Method 1

The SMM approach analogous to the GMM approach in Section 1.1 does not work. In most cases, the parameters are not identified. Loss of identification comes first because we have a maximum of $K + 1$ moments with potentially $K + 2$ parameters to estimate. A second characteristic that takes away a degree of freedom in this case is that one of the moment conditions $E[\varepsilon_i] = 0$ might be true for many parameter values in the assumed random variable distribution

In SMM, we need to modify the regression equation (4) to include a specific parameterized pdf for the i.i.d. random shocks ε_i . We assume here a normal distribution, which could add another parameter σ to our vector of parameters to be estimated θ .

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots \beta_K x_{K,i} + \varepsilon_i \\ \text{where } \varepsilon_i &\sim N(0, \sigma^2), \quad E[\varepsilon_i] = 0, \quad \text{and} \quad E[x_{k,i} \varepsilon_i] = 0 \quad \forall k \end{aligned} \quad (17)$$

The $K + 1$ model moment conditions are either not functions of the parameters

2.2 SMM Linear Regression: Method 2

Only the SMM method analogous to the second GMM method in Section 1.2 will work. The data \mathbf{y} and \mathbf{x} in linear regression model (17) has N observations on one dependent variable y_i and K independent variables $x_{k,i}$. The parameter vector to be estimated is now $\theta = (\beta_0, \beta_1, \dots, \beta_K, \sigma)$, which has length $K + 2$ and includes the standard deviation σ of the mean-zero, normally distributed random shocks ε_i .

In this case, let the vector of data moments be the vector of N dependent variable values y_i .

$$\mathbf{m}(\mathbf{y}, \mathbf{x}) = \mathbf{m}(\mathbf{y}) = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \mathbf{y} \quad (18)$$

Then let the data moments for the s th simulation be the vector of predicted values including the simulated error terms $\varepsilon_{s,i}$,

$$\mathbf{m}(\tilde{\mathbf{y}}_s, \tilde{\mathbf{x}}_s | \boldsymbol{\theta}) = \mathbf{m}(\tilde{\mathbf{x}}_s | \boldsymbol{\theta}) = \begin{bmatrix} \beta_0 + \beta_1 x_{1,1} + \beta_2 x_{2,1} + \dots \beta_K x_{K,1} + \varepsilon_{1,s} \\ \beta_0 + \beta_1 x_{1,2} + \beta_2 x_{2,2} + \dots \beta_K x_{K,2} + \varepsilon_{2,s} \\ \vdots \\ \beta_0 + \beta_1 x_{1,N} + \beta_2 x_{2,N} + \dots \beta_K x_{K,N} + \varepsilon_{N,s} \end{bmatrix} = \begin{bmatrix} \hat{y}_{1,s} \\ \hat{y}_{2,s} \\ \vdots \\ \hat{y}_{N,s} \end{bmatrix} = \hat{\mathbf{y}}_s \quad (19)$$

where $\varepsilon_{i,s}$ is the i th observation error term in the s th simulation of the errors, and $\hat{y}_{i,s}$ is the i th observation predicted value of the dependent variable in the s th simulation of the errors.

The final model moment vector $\hat{\mathbf{m}}(\tilde{\mathbf{x}} | \boldsymbol{\theta})$ to be used in SMM is an average across each of the moment vectors for each simulation $\mathbf{m}(\tilde{\mathbf{x}}_s | \boldsymbol{\theta})$.

$$\hat{\mathbf{m}}(\tilde{\mathbf{x}} | \boldsymbol{\theta}) \equiv \frac{1}{S} \sum_{s=1}^S \mathbf{m}(\tilde{\mathbf{x}}_s | \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \hat{\mathbf{y}}_s = \hat{\mathbf{y}} \quad (20)$$

Define the moment error function as the following.

$$\mathbf{e}(\tilde{\mathbf{x}}, \mathbf{y}, \boldsymbol{\theta}) = \hat{\mathbf{m}}(\tilde{\mathbf{x}} | \boldsymbol{\theta}) - \mathbf{m}(\mathbf{y}) = \hat{\mathbf{y}} - \mathbf{y} \quad (21)$$

The SMM approach of estimating the parameter vector $\hat{\boldsymbol{\theta}}_{SMM}$ is to choose $\boldsymbol{\theta}$ to minimize some distance measure of the simulated model moments $\hat{\mathbf{m}}(\tilde{\mathbf{x}} | \boldsymbol{\theta})$ from the data moments $\mathbf{m}(\mathbf{x})$.

$$\hat{\boldsymbol{\theta}}_{SMM} : \min_{\boldsymbol{\theta}} \mathbf{e}(\tilde{\mathbf{x}}, \mathbf{y}, \boldsymbol{\theta})^T \mathbf{W} \mathbf{e}(\tilde{\mathbf{x}}, \mathbf{y}, \boldsymbol{\theta}) \quad (22)$$

1. Draw S samples of N uniformly distributed random errors $u_{s,i} \sim U(0, 1)$ such that the sample of errors for the s th simulation is $\mathbf{u}_s = (u_{s,1}, u_{s,2}, \dots, u_{s,i}, \dots, u_{s,N})$.
2. Guess values for the parameter vector to be estimated $\boldsymbol{\theta}_v = (\beta_{0,v}, \beta_{1,v}, \dots, \beta_{K,v}, \sigma_v)$, where v indexes the iteration of the guess.
3. Compute S samples of N normally distributed random shocks $\varepsilon_{i,s}$ that correspond to the uniformly distributed random errors $u_{s,i}$ in step (1) and the current guess of the standard deviation σ_v of the normally distributed errors.
4. Use the simulated values of $\varepsilon_{i,s}$ and the guessed values of the coefficients $(\beta_{0,v}, \beta_{1,v}, \dots, \beta_{K,v})$ to compute average predicted values $\hat{\mathbf{y}}$ using (19) and (20).

5. Compute the value of the criterion function from (22) given the guess for the parameter vector θ_v .
6. Update guess for parameter vector to θ_{v+1} to minimize criterion function from (22). This should be done with a minimizer function.

3 Linking SMM to Cross Validation in Machine Learning

References

- Adda, Jérôme and Russell Cooper**, *Dynamic Economics: Quantitative Methods and Applications*, MIT Press, 2003.
- Davidson, Russell and James G. MacKinnon**, *Econometric Theory and Methods*, Oxford University Press, 2004.
- Duffie, Darrell and Kenneth J. Singleton**, “Simulated Moment Estimation of Markov Models of Asset Prices,” *Econometrica*, July 1993, 61 (4), 929–952.
- Laroque, G. and B. Salanie**, “Simulation Based Estimation Models with Lagged Latent Variables,” *Journal of Applied Econometrics*, December 1993, 8 (Supplement), 119–133.
- Lee, Bong-Soo and Beth Fisher Ingram**, “Simulation Based Estimation of Time-series Models,” *Journal of Econometrics*, February 1991, 47 (2-3), 197–205.
- McFadden, Daniel**, “A Method of Simulated Moments for Estimation of Discrete Response Models without Numerical Integration,” *Econometrica*, September 1989, 57 (5), 995–1026.