

# Notes on Generalized Method of Moments

Chris Conlon

October 28, 2017

## Introduction

These notes are meant as a non-technical introduction/review to the generalized method of moments (GMM). These are intended for my PhD IO class, and are not a serious econometric reference. For a more technical starting point please see <http://www.ssc.wisc.edu/~bhansen/econometrics/>

## Setup and Definitions

In the most basic setup we begin with some data  $w_i$  where  $i = 1, \dots, N$ . Our data  $w_i$  might contain all kinds of things such as dependent variables  $y_i$ , regressors  $x_i$  and excluded instruments  $z_i$ . The main idea is that our economic model provides the following restriction on our data:

$$E[g(w_i, \theta_0)] = 0$$

The idea is that at the true parameter value  $\theta_0 \in \mathbb{R}^k$  our moment conditions  $g(w_i, \theta)$  are on average equal to zero. What does “on average” mean? In theory, we are making an asymptotic statement about what happens as  $N \rightarrow \infty$ . This is what we mean when we write  $E[\cdot]$ . In practice, it is helpful to consider the sample analogue, which we abbreviate with the shorthand  $g_N(\theta) \in \mathbb{R}^q$ , where  $g_N(\theta)$  is a  $q$ -dimensional vector of moment conditions.

$$E[g(w_i, \theta)] \approx \frac{1}{N} \sum_{i=1}^N g(w_i, \theta) \equiv g_N(\theta)$$

We define the Jacobian:  $D(\theta) \equiv E[\frac{\partial g(w_i, \theta)}{\partial \theta}]$ , which is a  $q \times k$  matrix.

Evaluated at the optimum,  $\frac{1}{\sqrt{N}} \sum_{i=1}^N g(w_i, \theta_0) \xrightarrow{d} N(0, S)$  where  $S = E[g(w_i, \theta_0)g(w_i, \theta_0)']$  is a  $q \times q$  matrix.<sup>1</sup> In other words, the moment conditions which are 0 in expectation at  $\theta_0$  are normally distributed with some covariance  $S$ .

Later, we will refer to a weighting matrix  $W_N$  which is a  $q \times q$  positive semi-definite matrix. It tells us how much to penalize the violations of one moment condition relative to another (in quadratic distance).

## Examples

It is easy to see some very simple examples:

**OLS** Here  $y_i = x_i\beta + \epsilon_i$ . Exogeneity implies that  $E[x_i'\epsilon_i] = 0$ . We can write this in terms of just observables and parameters as  $E[x_i'(y_i - x_i\beta)] = 0$  so that  $g(y_i, x_i, \beta) = x_i'(y_i - x_i\beta)$ .

**IV** Again  $y_i = x_i\beta + \epsilon_i$ . Now, endogeneity implies that  $E[x_i'\epsilon_i] \neq 0$ . However there are some instruments  $z_i$  which may be partly contained in  $x_i$  and partly excluded from  $y_i$ , so that  $E[z_i'\epsilon_i] = 0$ .  $E[z_i'(y_i - x_i\beta)] = 0$  so that  $g(y_i, x_i, z_i, \beta) = z_i'(y_i - x_i\beta)$ .

---

<sup>1</sup>Technical conditions to establish this are written down later.

**Maximum Likelihood**  $g(w_i, \theta) = \frac{\partial \log f(w_i, \theta)}{\partial \theta}$  where  $f(w_i, \theta)$  is the density function so that  $\log f(w_i, \theta)$  is the contribution of observation  $i$  to the log-likelihood. Here we set the expected (average) derivative of the log-likelihood (score) function to zero.

**Euler Equation** Assume we have a CRRA utility function  $u(c) = \frac{c^{1-\gamma}-1}{1-\gamma}$  and an agent who maximizes the expected discounted value of their stream of consumption. This leads to an Euler Equation:

$$E \left[ \beta \left( \frac{c_{t+1}}{c_t} \right)^{-\gamma} R_{t+1} - 1 | \Omega_t \right] = 0$$

where  $\Omega_t$  is the “Information Set” (sigma algebra) of everything known to the agent up until time  $t$  (include full histories). We can write a moment restriction of the form for any measurable  $z_t \in \Omega_t$ .

$$E \left[ z_t \left( \beta \left( \frac{c_{t+1}}{c_t} \right)^{-\gamma} R_{t+1} - 1 \right) \right] = 0$$

In the original work by Hansen (1982) on GMM, this  $g(c_t, c_{t+1}, R_{t+1}, \beta, \gamma)$  was used to estimate  $(\beta, \gamma)$ .

### Technical Conditions and Asymptotics

Here is the GMM estimator:

$$\hat{\theta} = \arg \min_{\theta} Q_N(\theta) \quad Q_N(\theta) = g_N(\theta)' W_N g_N(\theta)$$

*These are a set of sufficient conditions to establish consistency and asymptotic normality of the GMM estimator. These conditions are stronger than necessary, but they establish the requisite LLN and CLT.*

1.  $\theta \in \Theta$  is compact.
2.  $W_N \xrightarrow{P} W$ .
3.  $g_N(\theta) \xrightarrow{P} E[g(z_i, \theta)]$  (uniformly)
4.  $E[g(z_i, \theta)]$  is continuous.
5. We need that  $E[g(z_i, \theta_0)] = 0$  and  $W_N E[g(z_i, \theta)] \neq 0$  for  $\theta \neq \theta_0$  (global identification condition).
6.  $g_N(\theta)$  is twice continuously differentiable about  $\theta_0$ .
7.  $\theta_0$  is not on the boundary of  $\Theta$ .
8.  $D(\theta_0) W D(\theta_0)'$  is invertible (non-singular).
9.  $g(z_i, \theta)$  has at least two moments finite and finite derivatives at all  $\theta \in \Theta$ .

The first five conditions give us consistency  $\hat{\theta} \xrightarrow{P} \theta_0$  as  $N \rightarrow \infty$ . All nine conditions give us asymptotic normality.

$$\begin{aligned} \sqrt{N}(\hat{\theta} - \theta) &\xrightarrow{d} N(0, V_{\theta}) \\ V_{\theta} &= \underbrace{(D W D')^{-1}}_{\text{bread}} \underbrace{(D W S W' D')}_{\text{filling}} \underbrace{(D W D')^{-1}}_{\text{bread}} \end{aligned}$$

It is common to refer parts of the variance as the *bread* and the *filling* or *meat*, together this is referred to as the *sandwich* estimator of the variance.

## Identification and Examples

The global identification condition is difficult to understand, for the linear model we can replace it with a (local) condition on Jacobian of the moment conditions. Recall the Jacobian:  $D \equiv \frac{\partial g(w_i, \theta)}{\partial \theta}$ , which is a  $q \times k$  matrix. We call the problem under-identified if  $\text{rank}(D) < k$ , just-identified if  $\text{rank}(D) = k$  and over-identified if  $\text{rank}(D) > k$ . In the under-identified case, there may be many such  $\hat{\theta}$  where  $g(w_i, \hat{\theta}) = 0$ . In the just-identified case, it should be possible to find a  $\hat{\theta}$  where  $g_N(\hat{\theta}) = 0$ . We are primarily interested in the over-identified case where we will generally not find  $\hat{\theta}$  which satisfies the moment conditions  $g_N(\hat{\theta}) \neq 0$ . Instead, we search for  $\hat{\theta}$  which minimizes the violations of the moment conditions. We write this as a quadratic form for some positive definite matrix  $W_N$  which is  $q \times q$ .

$$\hat{\theta} = \arg \min_{\theta} Q_N(\theta) \quad Q_N(\theta) = g_N(\theta)' W_N g_N(\theta)$$

For the linear IV problem this becomes:

$$\begin{aligned} g_N(\theta)' W_N g_N(\theta) &= \frac{1}{N^2} \cdot (\mathbf{Z}'(\mathbf{Y} - \mathbf{X}\beta))' W_N (\mathbf{Z}'(\mathbf{Y} - \mathbf{X}\beta)) \\ &= \frac{1}{N^2} \cdot [Y' Z W_N Z' Y - 2\beta' X' Z W_N Z' Y + \beta' X' Z W_N Z' X \beta] \end{aligned}$$

We can ignore the  $\frac{1}{N^2}$  and take the first-order condition:

$$\begin{aligned} 2X' Z W_N Z' Y &= 2X' Z W_N Z' X \beta \\ \hat{\beta}_{GMM} &= (X' Z W_N Z' X)^{-1} X' Z W_N Z' Y \end{aligned}$$

### OLS Estimator

Suppose that we do not have any excluded instruments so that  $Z = X$  (and thus  $q = k$ ). Also suppose that  $W_N = \mathbf{I}_q$  (the identity matrix). Then we can see that:

$$\begin{aligned} \hat{\beta}_{GMM} &= (X' X \mathbf{I}_q X' X)^{-1} X' X \mathbf{I}_q X' Y \\ &= (X' X X' X)^{-1} X' X X' Y \\ &= (X' X)^{-1} (X' X)^{-1} (X' X) X' Y = (X' X)^{-1} X' Y = \hat{\beta}_{OLS} \end{aligned}$$

In other words, OLS is a special case of the GMM estimator. Also, the identification condition  $D = \frac{\partial g(w_i, \theta)}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N z'_i x_i = \frac{1}{N} \sum_{i=1}^N x'_i x_i$  becomes that  $\text{rank}(X' X) = k$  the well-known OLS rank condition.

### 2SLS Estimator

Suppose that we do have excluded instruments so that  $\dim(Z) = q > \dim(X) = k$  and that  $W_N = (Z' Z)^{-1}$ . It immediately follows that:

$$\hat{\beta}_{GMM} = (X' Z (Z' Z)^{-1} Z' X)^{-1} X' Z (Z' Z)^{-1} Z' Y = \hat{\beta}_{2SLS}$$

If  $\dim(Z) = q = \dim(X) = k$  then  $(X' Z)$  is square (and invertible). This expression further simplifies:

$$\begin{aligned} (X' Z (Z' Z)^{-1} Z' X)^{-1} &= (Z' X)^{-1} (Z' Z) (X' Z)^{-1} \\ \rightarrow \hat{\beta}_{GMM} &= (Z' X)^{-1} (Z' Z) (X' Z)^{-1} X' Z (Z' Z)^{-1} Z' Y = (Z' X)^{-1} Z' Y = \hat{\beta}_{IV} \end{aligned}$$

## Optimal Weighting Matrix

An important question remains how one should choose the weighting matrix  $W_N$ . We've already seen two options: (1) The identity matrix  $\mathbf{I}_q$  equally penalizes violations of all  $q$  moments; and (2) the TSLS weighting matrix  $(Z' Z)^{-1}$  which can be thought about as the inverse of the covariance of the instruments.

We are interested in *efficient GMM* which is the GMM estimator with the lowest variance. In order to find the  $W_N$  which minimizes the variance of  $\hat{\theta}_{GMM}$  we recall the asymptotic variance of the GMM estimator:

$$V_{\theta} = (D W D')^{-1} (D W S W' D') (D W D')^{-1}$$

It turns out that the best choice of  $W_N = S^{-1}$  (which sets *filling* = *bread*). This is easy to see, because  $W_N$  is positive semi-definite.

$$(DS^{-1}D')^{-1}(DS^{-1}SS^{-1'}D')(DS^{-1}D')^{-1} = (DS^{-1}D')^{-1}(DS^{-1}D')(DS^{-1}D')^{-1} = (DS^{-1}D')^{-1}$$

This gives us some insight into what we are looking for from moment conditions. In general we want  $S$  to be small (we want the sampling variation/noise of our moments to be as small as possible). We also want  $D$  (the Jacobian of the moments) to be large. This means that small violations in moment conditions lead to large changes in the objective function. In practical terms, the problem is well identified when the objective function is steep around  $\theta_0$ . When the problem becomes flat, it becomes hard to distinguish one  $\theta$  in favor of another.

The problem is that  $S = E[g(w_i, \theta_0)g(w_i, \theta_0)']$  is not something that we readily observe from our data. In fact, the asymptotic covariance evaluated at  $\theta_0$  is *infeasible*. The best we can hope for is to use some sample analogue  $W_N = \hat{S}^{-1}$  in its place. One way to compute that is the covariance of the moments estimated at some  $\hat{\theta}$  for an initial guess of  $W$ :

$$\hat{W} = \hat{S}^{-1} = \left( \frac{1}{N} \sum_{i=1}^N (g(w_i, \hat{\theta}) - g_N(\hat{\theta})) (g(w_i, \hat{\theta}) - g_N(\hat{\theta}))' \right)^{-1}$$

Because  $E[g(w_i, \theta_0)] = 0$  at  $\theta_0$  there is a tendency to use  $\left( \frac{1}{N} \sum_{i=1}^N g(w_i, \hat{\theta}) g(w_i, \hat{\theta})' \right)^{-1}$  (without de-meaning the moments). In theory this would work fine, but in practice **it is nearly always a bad idea**.

The overall procedure works as follows:

1. Pick some initial weighting matrix  $W_0$ : often  $\mathbf{I}_q$  or  $(Z'Z)^{-1}$ .
2. Solve  $\hat{\theta} = \arg \min_{\theta} g_N(\theta)' W_0 g_N(\theta)$ .
3. Update  $\hat{W} = \left( \frac{1}{N} \sum_{i=1}^N (g(w_i, \hat{\theta}) - g_N(\hat{\theta})) (g(w_i, \hat{\theta}) - g_N(\hat{\theta}))' \right)^{-1}$
4. Solve  $\hat{\theta}_{GMM} = \arg \min_{\theta} g_N(\theta)' \hat{W} g_N(\theta)$ .
5. Compute  $D(\hat{\theta}_{GMM})$  and  $S(\hat{\theta}_{GMM})$  and compute standard errors.

#### Example

For the linear IV estimator when  $i$  is independent then  $g(w_i, \theta) = z_i \epsilon_i$  and  $E[z_i \epsilon_i] = 0$

$$\hat{S} = \frac{1}{N} \sum_{i=1}^N z_i z_i' \epsilon_i^2$$

When there is homoskedastic variance  $E[\epsilon_i^2 | z_i] = \sigma^2$  and the covariance of the moments becomes  $\frac{\sigma^2}{N} \sum_{i=1}^N z_i z_i'$ . Because scaling weighting matrix by a constant has no effect on the maximum this is equivalent to the 2SLS weight matrix:  $\sum_{i=1}^N z_i z_i'$  or  $Z'Z$ . Thus 2SLS is only the efficient estimator when homoskedasticity is a reasonable assumption. Likewise, if all regressors are exogenous then  $X = Z$  and we are left with the GMM formula coincides with the covariance for heteroskedasticity robust standard errors.

Similarly, when appropriate we can consider extensions such as *clustered standard errors* which are robust to weaker forms of independence.

As a practical matter, we should always use the *sandwich* form when calculating the GMM standard errors, rather than the simpler *bread* version which is only correct at  $\theta_0$  under asymptotic optimality conditions.

## Common Questions/Extensions

**Semiparametric Efficiency** In a famous paper, Chamberlain (1987) showed that GMM obtained the semi-parametric efficiency bound asymptotically. That means as our sample gets large and without making additional parametric assumptions (such as placing distributions on error terms) the efficient GMM provides the most efficient estimator.

**My first stage estimates look ok, but my second stage estimates look like garbage** This means there is a problem with your weighting matrix. Make sure you are subtracting off the average  $g_N(\theta)$  when computing the covariance. Another problem appears is when you take the inverse. There is a statistic known as the *condition number* which measures the ratio of the minimum and maximum eigenvalue of a matrix. When these eigenvalues are close together then small errors in  $A + \epsilon$  lead to small errors in  $(A + \epsilon)^{-1}$ . Software sometimes reports either the condition number, or its inverse. In an ideal world this would be  $\approx 1$ . If the condition number is  $10^{\pm 13}$  it approaches the numerical precision of your computer. Then even tiny (sampling) errors can lead to nearly infinite weighting matrices. Usually this happens if the gradient of  $Q_n(\theta)$  is not close to zero in a particular dimension or the average Jacobian  $D(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{\partial g(w_i, \theta)}{\partial \theta}$  is not close to zero in some dimension. [You are not an optimum of  $Q_n(\theta)$ !].

**My point estimates look okay, but my standard errors look crazy** Assuming it is not one of the problems described above, you have probably dropped a  $\frac{1}{N}$  or  $\frac{1}{\sqrt{N}}$  somewhere. You can ignore the  $N$ 's when obtaining point estimates because scaling  $Q_n(\theta)$  or  $W_n$  does not affect the location of  $\hat{\theta}$ , but you need to be more careful in computing  $V_{\hat{\theta}}$ .

**If 2 step GMM is good, is 3-step GMM better?** Asymptotically the answer is no. Even with a sub-optimal choice of weighting matrix, your first stage  $\hat{\theta} \xrightarrow{P} \theta_0$  as  $N \rightarrow \infty$ . This means that in the limit, you always recover the efficient choice of  $\hat{W}$ . In finite sample, anything can happen, but numerous studies have found little to no improvement from considering a  $K$ -step GMM estimator.

**What about Infinite-step GMM?** There isn't such a thing. But there is the Continuously Updating GMM estimator of Hansen Heaton and Yaron (1996). In this case we let

$W(\theta) = \left( \frac{1}{N} \sum_{i=1}^N (g(w_i, \theta) - g_N(\theta))(g(w_i, \theta) - g_N(\theta))' \right)$  so that

$$Q_n(\theta)^{CUE} = g_N(\theta)' \left( \frac{1}{N} \sum_{i=1}^N (g(w_i, \theta) - g_N(\theta))(g(w_i, \theta) - g_N(\theta))' \right) g_N(\theta)$$

Because the weighting matrix now changes with  $\theta$ , even for linear models this becomes very difficult to estimate. The original GMM problem was a quadratic optimization problem. The CUE problem is no longer a convex optimization problem. This means that even numerical Quasi-Newton approaches are no longer guaranteed to work. In practice, CUE is not very popular for this reason.

CUE has some additional advantages over GMM. We see this by examining the GMM FOC.

$$\hat{D}(\theta) W_n \hat{g}_N(\hat{\theta}) = 0$$

If we take an Edgeworth expansion we can see how the *asymptotic bias* term looks:

$$\frac{1}{N^2} \sum_{i=1}^N E[g(w_i, \theta) W_n g(w_i, \theta)]$$

As long as the variance is finite we can see that this bias disappears as  $N \rightarrow \infty$ . However it also tends to grow linearly with  $q$  the number of moment restrictions. This is often referred to as the *too many moments* or *many instruments* problem.

The challenge is that in general when we estimate  $\hat{g}(\hat{\theta})$  and  $\hat{D}(\hat{\theta})$  we construct them so that they are mechanically correlated with one another. It turns out the CUE estimator fixes this in exactly the right way. How is beyond the scope of this note (and course). If you are really interested you should look at Newey and Smith (2004).

### For Demand Estimation

The BLP moments are given by  $g(z_i, \theta) = z_{jt}\xi_{jt}(\theta)$  and the Jacobian  $\frac{\partial g(z_i, \theta)}{\partial \theta} = z_{jt} \frac{\partial \xi_{jt}(\theta)}{\partial \theta}$ . The Jacobian is hard to work out analytically (see Nevo's Practitioner's Guide: Appendix):  $\frac{\partial \xi_{jt}(\theta)}{\partial \theta}$  can be calculated using the implicit function theorem given  $\frac{\partial s_{jt}}{\partial \xi_{kt}} = \frac{\partial s_{jt}}{\partial \delta_{kt}}$  and  $\frac{\partial s_{jt}}{\partial \theta}$  where the former  $A$  is stacked into a  $J \times J$  matrix and the latter  $B$  is stacked into a  $J \times K$  matrix and then applying  $\frac{\partial \xi_{jt}(\theta)}{\partial \theta} = A^{-1}B$ .