

# Choosing the Number of Moments in Conditional Moment Restriction Models

Stephen G. Donald  
Department of Economics  
University of Texas

Guido Imbens  
Department of Economics  
UC-Berkeley

Whitney Newey  
Department of Economics  
MIT

First Draft: January 2002  
This draft: October 2008

## Abstract

Properties of GMM estimators are sensitive to the choice of instruments. Using many instruments leads to high asymptotic efficiency but can cause high bias and/or variance in small samples. In this paper we develop and implement asymptotic mean square error (MSE) based criteria for instrumental variables to use for estimation of conditional moment restriction models. The models we consider include various nonlinear simultaneous equations models with unknown heteroskedasticity. We develop moment selection criteria for the familiar two-step optimal GMM estimator (GMM), a bias corrected version, and generalized empirical likelihood estimators (GEL), that include the continuous updating estimator (CUE) as a special case. We also find that the CUE has lower higher-order variance than the bias-corrected GMM estimator, and that the higher-order efficiency of other GEL estimators depends on conditional kurtosis of the moments.

**JEL Classification:** C13, C30

**Keywords:** Conditional Moment Restrictions, Generalized Method of Moments, Generalized Empirical Likelihood, Mean Squared Error.

# 1 Introduction

It is important to choose carefully the instrumental variables for estimating conditional moment restriction models. Adding instruments increases asymptotic efficiency but also increases small sample bias and/or variance. We account for this trade-off by using a higher-order asymptotic mean-square error (MSE) of the estimator to choose the instrument set. We derive the higher-order MSE for GMM, a bias corrected version of GMM (BGMM), and generalized empirical likelihood (GEL). For simplicity we impose a conditional symmetry assumption, that third conditional moments of disturbances are zero, and use a large number of instruments approximation. We also consider the effect of allowing identification to shrink with the sample size  $n$  at a rate slower than  $1/\sqrt{n}$ . The resulting MSE expressions are quite simple and straightforward to apply in practice to choose the instrument set.

The MSE criteria given here also provide higher order efficiency comparisons. We find that continuously updated GMM estimator (CUE) is higher-order efficient relative to BGMM. We also find that the higher order efficiency of the GEL estimators depends on conditional kurtosis, with all GEL estimators having the same higher-order variance when disturbances are Gaussian. With Gaussian disturbances and homoskedasticity, Rothenberg (1996) showed that empirical likelihood (EL) is higher order efficient relative to BGMM. Our efficiency comparisons generalize those of Rothenberg (1996) to other GEL estimators and heteroskedastic, non Gaussian disturbances. These efficiency results are different than the higher order efficiency result for EL that was shown by Newey and Smith (2004) because Newey and Smith (2004) do not require that conditional third moments are zero. Without that symmetry condition all of the estimators except for EL have additional bias terms that are not corrected for here.

Our MSE criteria is like that of Nagar (1959) and Donald and Newey (2001), being the MSE of leading terms in a stochastic expansion of the estimator. This approach is well known to give the same answer as the MSE of leading terms in an Edgeworth expansion, under suitable regularity conditions (e.g. Rothenberg, 1984). The many instrument and

shrinking identification simplifications seems appropriate for many applications where there is a large number of potential instrumental variables and identification is not very strong. We also assume symmetry, in the sense that conditional third moments of the disturbances are zero. This symmetry assumption greatly simplifies calculations. Also, relaxing it may not change the results much, e.g. because the bias from asymmetry tends to be smaller than other bias sources for large numbers of moment conditions, see Newey and Smith (2004).

Choosing moments to minimize MSE may help reduce misleading inferences that can occur with many moments. For GMM, the MSE explicitly accounts for an important bias term (e.g. see Hansen et. al., 1996, and Newey and Smith, 2004), so choosing moments to minimize MSE avoids cases where asymptotic inferences are poor due to the bias being large relative to the standard deviation. For GEL, the MSE explicitly accounts for higher order variance terms, so that choosing instruments to minimize MSE helps avoid underestimated variances. However, the criteria we consider may not be optimal for reducing misleading inferences. That would lead to a different criteria, as recently pointed out by Jin, Phillips, and Sun (2007) in another context.

The problem addressed in this paper is different than considered by Andrews (1996). Here the problem is how to choose among moments known to be valid while Andrews (1996) is about searching for the largest set of valid moments. Choosing among valid moments is important when there are many thought to be equally valid. Examples include various natural experiment studies, where multiple instruments are often available, as well as intertemporal optimization models, where all lags may serve as instruments.

In Section 2 we describe the estimators we consider and present the criteria we develop for choosing the moments. We also compare the criteria for different estimators, which corresponds to the MSE comparison for the estimators, finding that the CUE has smaller MSE than bias corrected GMM. In Section 3 we give the regularity conditions used to develop the approximate MSE and give the formal results. Section 4 shows optimality of the criteria we propose. A small scale Monte Carlo experiment is conducted in Section 5. Concluding remarks are offered in Section 6.

## 2 The Model and Estimators

We consider a model of conditional moment restrictions like Chamberlain (1987). To describe the model let  $z$  denote a single observation from an i.i.d. sequence  $(z_1, z_2, \dots)$ ,  $\beta$  a  $p \times 1$  parameter vector, and  $\rho(z, \beta)$  a scalar that can often be thought of as a residual<sup>1</sup>. The model specifies a subvector of  $x$ , acting as conditioning variables, such that for a value  $\beta_0$  of the parameters

$$E[\rho(z, \beta_0)|x] = 0,$$

where  $E[\cdot]$  the expectation taken with respect to the distribution of  $z_i$ .

To form GMM estimators we construct unconditional moment restrictions using a vector of  $K$  conditioning variables  $q^K(x) = (q_{1K}(x), \dots, q_{KK}(x))'$ . Let  $g(z, \beta) = \rho(z, \beta)q^K(x)$ . Then the unconditional moment restrictions

$$E[g(z, \beta_0)] = 0$$

are satisfied. Let  $g_i(\beta) \equiv g(z_i, \beta)$ ,  $\bar{g}_n(\beta) \equiv n^{-1} \sum_{i=1}^n g_i(\beta)$ , and  $\hat{\Upsilon}(\beta) \equiv n^{-1} \sum_{i=1}^n g_i(\beta)g_i(\beta)'$ . A two-step GMM estimator is one that satisfies, for some preliminary consistent estimator  $\tilde{\beta}$  for  $\beta_0$ ,

$$\hat{\beta}^H = \arg \min_{\beta \in \mathcal{B}} \bar{g}_n(\beta)' \hat{\Upsilon}(\tilde{\beta})^{-1} \bar{g}_n(\beta), \quad (2.1)$$

where  $\mathcal{B}$  denotes the parameter space. For our purposes  $\tilde{\beta}$  could be some other GMM estimator, obtained as the solution to an analogous minimization problem with  $\hat{\Upsilon}(\tilde{\beta})^{-1}$  replaced by a different weighting matrix, such as  $\tilde{W}_0 = [\sum_{i=1}^n q^K(x_i)q^K(x_i)'/n]^{-1}$ .

The MSE of the estimators will depend not only on the number of instruments but also on their form. In particular, instrumental variables that better predict the optimal instruments will help to lower the asymptotic variance of the estimator for a given  $K$ . Thus, for each  $K$  it is good to choose  $q^K(x)$  that are the best predictors. Often it will be evident in an application how to choose the instruments in this way. For instance, lower order approximating functions (e.g. linear and quadratic) often provide the most

---

<sup>1</sup>The extension to the vector of residuals case is straightforward.

information, and so should be used first. Also, main terms may often be more important than interactions.

The instruments need not form a nested sequence. Letting  $q_{kK}(x)$  depend on  $K$  allows different groups of instrumental variables to be used for different values of  $K$ . Indeed,  $K$  fills a double role here, as the index of the instrument set as well as the number of instruments. We could separate these roles by having a separate index for the instrument set. Instead here we allow for  $K$  to not be selected from all the integers, and let  $K$  fulfill both roles. This restricts the sets of instruments to each have a different number of instruments, but is often true in practice. Also, by imposing upper bounds on  $K$  we also restrict the number of instrument sets we can select among, as seems important for the asymptotic theory.

As demonstrated by Newey and Smith (2004), the correlation of the residual with the derivative of the moment function leads to an asymptotic bias that increases linearly with  $K$ . They suggested an approach that removes this bias (as well as other sources of bias that we will ignore for the moment). This estimator can be obtained by subtracting an estimate of the bias from the GMM estimator and gives rise to what we refer to as the bias adjusted GMM estimator (BGMM). To describe it, let  $q_i = q^K(x_i)$  and

$$\begin{aligned}\hat{\rho}_i &= \rho_i(\hat{\beta}^H), \hat{y}_i = [\partial \rho_i(\hat{\beta}^H) / \partial \beta]', \hat{y} = [\hat{y}_1, \dots, \hat{y}_n]', \\ \hat{\Gamma} &= \sum_{i=1}^n q_i \hat{y}_i' / n, \hat{\Sigma} = \hat{\Upsilon}(\hat{\beta}^H)^{-1} - \hat{\Upsilon}(\hat{\beta}^H)^{-1} \hat{\Gamma} (\hat{\Gamma}' \hat{\Upsilon}(\hat{\beta}^H)^{-1} \hat{\Gamma})^{-1} \hat{\Gamma}' \hat{\Upsilon}(\hat{\beta}^H)^{-1}.\end{aligned}$$

The BGMM estimator is

$$\hat{\beta}^B = \hat{\beta}^H + (\hat{\Gamma}' \hat{\Upsilon}(\hat{\beta}^H)^{-1} \hat{\Gamma})^{-1} \sum_{i=1}^n \hat{y}_i \hat{\rho}_i q_i' \hat{\Sigma} q_i.$$

Also as shown in Newey and Smith (2004) the class of Generalized Empirical Likelihood (GEL) estimators have less bias than GMM. We follow the description of these estimators given in that paper. Let  $s(v)$  be a concave function with domain that is an open interval  $\mathcal{V}$  containing 0,  $s_j(v) = \partial^j s(v) / \partial v^j$ , and  $s_j = s_j(0)$ . We impose the normalizations  $s_1 = s_2 = -1$ . Define the GEL estimator as

$$\hat{\beta}^{GEL} = \arg \min_{\beta \in B} \max_{\lambda \in \hat{\Lambda}_n(\beta)} \sum_{i=1}^n s(\lambda' g_i(\beta))$$

where,  $\hat{\Lambda}_n(\beta) = \{\lambda : \lambda' g_i(\beta) \in \mathcal{V}, i = 1, \dots, n\}$ . This estimator includes as a special cases: empirical likelihood (EL, Qin and Lawless, 1997, and Owen, 1988), where  $s(v) = \ln(1-v)$ , exponential tilting (ET, Johnson, and Spady, 1998, and Kitamura and Stutzer, 1997), where  $s(v) = -\exp(v)$ , and the continuous updating estimator (CUE, Hansen, Heaton, and Yaron 1996), where  $s(v) = -(1+v)^2/2$ . As we will see the MSE comparisons between these estimators depend on  $s_3$ , the third derivative of the  $s$  function, where

$$CUE : s_3 = 0, ET : s_3 = -1, EL : s_3 = -2.$$

## 2.1 Instrument Selection Criteria

The instrument selection is based on minimizing the approximate mean squared error (MSE) of a linear combination  $\hat{t}'\hat{\beta}$  of a GMM estimator or GEL estimator  $\hat{\beta}$ , where  $\hat{t}$  is some vector of (estimated) linear combination coefficients. To describe the criteria, some additional notation is required. Let  $\tilde{\beta}$  be some preliminary estimator,  $\tilde{\rho}_i = \rho_i(\tilde{\beta})$ ,  $\tilde{y}_i = \partial\rho_i(\tilde{\beta})/\partial\beta$ , and

$$\begin{aligned} \hat{\Upsilon} &= \sum_{i=1}^n \tilde{\rho}_i^2 q_i q_i' / n, \hat{\Gamma} = \sum_{i=1}^n q_i \tilde{y}_i' / n, \hat{\Omega} = (\hat{\Gamma}' \hat{\Upsilon}^{-1} \hat{\Gamma}), \hat{\tau} = \hat{\Omega}^{-1} \hat{t}, \\ \tilde{d}_i &= \hat{\Gamma}' \left( \sum_{j=1}^n q_j q_j' / n \right)^{-1} q_i, \tilde{\eta}_i = \tilde{y}_i - \tilde{d}_i, \hat{\xi}_{ij} = q_i' \hat{\Upsilon}^{-1} q_j / n, \hat{D}_i^* = \hat{\Gamma}' \hat{\Upsilon}^{-1} q_i, \\ \hat{\Lambda}(K) &= \sum_{i=1}^n \hat{\xi}_{ii} (\hat{\tau}' \tilde{\rho}_{\beta i})^2, \hat{\Pi}(K) = \sum_{i=1}^n \hat{\xi}_{ii} \tilde{\rho}_i (\hat{\tau}' \tilde{\eta}_i), \\ \hat{\Phi}(K) &= \sum_{i=1}^n \hat{\xi}_{ii} \left\{ \hat{\tau}' (\hat{D}_i^* \tilde{\rho}_i^2 - \tilde{\rho}_{\beta i}) \right\}^2 - \hat{\tau}' \hat{\Gamma}' \hat{\Upsilon}^{-1} \hat{\Gamma} \hat{\tau} \end{aligned}$$

The criteria for the GMM estimator, without a bias correction, is

$$S_{GMM}(K) = \hat{\Pi}(K)^2 / n + \hat{\Phi}(K).$$

Also, let

$$\begin{aligned} \hat{\Pi}_B(K) &= \sum_{i,j=1}^n \tilde{\rho}_i \tilde{\rho}_j (\hat{\tau}' \tilde{\eta}_i) (\hat{\tau}' \tilde{\eta}_j) \hat{\xi}_{ij}^2 = \text{tr}(\tilde{Q} \hat{\Upsilon}^{-1} \tilde{Q} \hat{\Upsilon}^{-1}), \\ \hat{\Xi}(K) &= \sum_{i=1}^n \{5(\hat{\tau}' \hat{d}_i)^2 - \tilde{\rho}_i^4 (\hat{\tau}' \hat{D}_i^*)^2\} \hat{\xi}_{ii}, \hat{\Xi}_{GEL}(K) = \sum_{i=1}^n \{3(\hat{\tau}' \hat{d}_i)^2 - \tilde{\rho}_i^4 (\hat{\tau}' \hat{D}_i^*)^2\} \hat{\xi}_{ii}, \end{aligned}$$

where  $\tilde{Q} = \sum_{i=1}^n \tilde{\rho}_i(\hat{\tau}'\tilde{\eta}_i)q_iq_i'$ . The criteria for the BGMM and GEL estimators are

$$\begin{aligned} S_{BGMM}(K) &= \left[ \hat{\Lambda}(K) + \hat{\Pi}_B(K) + \hat{\Xi}(K) \right] / n + \hat{\Phi}(K), \\ S_{GEL}(K) &= \left[ \hat{\Lambda}(K) - \hat{\Pi}_B(K) + \hat{\Xi}(K) + s_3\hat{\Xi}_{GEL}(K) \right] / n + \hat{\Phi}(K). \end{aligned}$$

For each of the estimators, our proposed instrument selection procedure is to choose  $K$  to minimize  $S(K)$ . As we will show this will correspond to choosing  $K$  to minimize the higher-order MSE of the estimator.

Each of the terms in the criteria have an interpretation. For GMM,  $\hat{\Pi}(K)^2/n$  is an estimate of a squared bias term from Newey and Smith (2004). Because  $\hat{\xi}_{ii}$  is of order  $K$  this squared bias term has order  $K^2/n$ . The  $\hat{\Phi}(K)$  term in the GMM criteria is an asymptotic variance term. Its size is related to the asymptotic efficiency of a GMM estimator with instruments  $q^K(x)$ . As  $K$  grows this terms will tend to shrink, reflecting the reduction in asymptotic variance that accompanies using more instruments. The form of  $\hat{\Phi}(K)$  is analogous to a Mallows criterion, in that it is a variance estimator plus a term that removes bias in the variance estimator.

The terms that appear in  $S(K)$  for BGMM and GEL are all variance terms. No bias terms are present because, as discussed in Newey and Smith (2004), under symmetry GEL removes the GMM bias that grows with  $K$ . As with GMM, the  $\hat{\Phi}(K)$  term accounts for the reduction in asymptotic variance that occurs from adding instruments. The other terms are higher-order variance terms, that will be of order  $K/n$ , because  $\hat{\xi}_{ii}$  is of order  $K$ . The sum of these terms will generally increase with  $K$ , although this need not happen if  $\hat{\Xi}(K)$  is too large relative to the other terms. As we will discuss below,  $\hat{\Xi}(K)$  is an estimator of

$$\Xi(K) = \sum_{i=1}^n \xi_{ii} (\tau' d_i)^2 \{5 - E(\rho_i^4|x_i)/\sigma_i^4\}.$$

As a result if the kurtosis of  $\rho_i$  is too high the higher-order variance of the BGMM and GEL estimators would actually decrease as  $K$  increases. This phenomenon is similar to that noted by Koenker et. al. (1994) for the exogenous linear case. In this case the criteria could fail to be useful as a means of choosing the number of moment conditions, because they would monotonically decrease with  $K$ .

It is interesting to compare the size of the criteria for different estimators, which comparison parallels that of the MSE. As previously noted, the squared bias term for GMM, which is  $\hat{\Pi}(K)^2$ , has the same order as  $K^2/n$ . In contrast the higher-order variance terms in the BGMM and GEL estimators generally have order  $K/n$ , because that is the order of  $\xi_{ii}$ . Consequently, for large  $K$  the MSE criteria for GMM will be larger than the MSE criteria for BGMM and GEL, meaning the BGMM and GEL estimators are preferred over GMM. This comparison parallels that in Newey and Smith (2004) and in Imbens and Spady (2002).

One interesting result is that for the CUE, where  $s_3 = 0$ , the MSE criteria is smaller than it is for BGMM, because  $\hat{\Pi}_B(K)$  is positive. Thus we find that the CUE dominates the BGMM estimator, in terms of higher-order MSE, i.e. CUE is higher-order efficient relative to BGMM. This result is analogous to the higher-order efficiency of the limited information maximum likelihood estimator relative to the bias corrected two-stage least squares estimator that was found by Rothenberg (1983).

The comparison of the higher-order MSE for the CUE and the other GEL estimators depends on the kurtosis of the residual. Let  $\rho_i = \rho(z_i, \beta_0)$  and  $\sigma_i^2 = E[\rho_i^2|x_i]$ . For conditionally normal  $\rho_i$  we have  $E[\rho_i^4|x_i] = 3\sigma_i^4$  and consequently  $\hat{\Xi}_{GEL}(K)$  will converge to zero for each  $K$ , that all the GEL estimators have the same higher-order MSE. When there is excess kurtosis, with  $E[\rho_i^4|x_i] > 3\sigma_i^4$ , ET will have larger MSE than CUE, and EL will have larger MSE than ET, with these rankings being reversed when  $E[\rho_i^4|x_i] < 3\sigma_i^4$ . These comparisons parallel those of Newey and Smith (2004) for a heteroskedastic linear model with exogeneity.

The case with no endogeneity has some independent interest. In this setting the GMM estimator can often be interpreted as using "extra" moment conditions to improve efficiency in the presence of heteroskedasticity of unknown functional form. Here the MSE criteria will give a method for choosing the number of moments used for this purpose. Dropping the bias terms, which are not present in exogenous cases, leads to criteria of



the form

$$\begin{aligned}
S_{GMM}(K) &= \hat{\Xi}(K)/n + \hat{\Phi}(K) \\
S_{GEL}(K) &= \left[ \hat{\Xi}(K) + s_3 \hat{\Xi}_{GEL}(K) \right] /n + \hat{\Phi}(K)
\end{aligned}$$

Here GMM and CUE have the same higher-order variance, as was found by Newey and Smith (2002). Also, as in the general case, these criteria can fail to be useful if there is too much kurtosis.

### 3 Assumptions and MSE Results

#### 3.1 Basic Expansion

As in Donald and Newey (2001), the MSE approximations are based on a decomposition of the form,

$$\begin{aligned}
nt'(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)'t &= \hat{Q}(K) + \hat{R}(K), & (3.2) \\
E(\hat{Q}(K)|X) &= t'\Omega^{*-1}t + S(K) + T(K), \\
[\hat{R}(K) + T(K)]/S(K) &= o_p(1), K \rightarrow \infty, n \rightarrow \infty.
\end{aligned}$$

where  $X = [x_1, \dots, x_n]'$ ,  $t = plim(\hat{t})$ ,  $\Omega^* = \sum_{i=1}^n \sigma_i^{-2} d_i d_i' / n$ ,  $\sigma_i^2 = E[\rho_i^2 | x_i]$ , and  $d_i = E[\partial \rho_i(\beta_0) / \partial \beta | x_i]$ . Here  $S(K)$  is part of conditional MSE of  $\hat{Q}$  that depends on  $K$  and  $\hat{R}(K)$  and  $T(K)$  are remainder terms that goes to zero faster than  $S(K)$ . Thus,  $S(K)$  is the MSE of the dominant terms for the estimator. All calculations are done assuming that  $K$  increases with  $n$ . The largest terms increasing and decreasing with  $K$  are retained. Compared to Donald and Newey (2001) we have the additional complication that none of our estimators has a closed form solution. Thus, we use the first order condition that defines the estimator to develop approximations to the difference  $\sqrt{nt}'(\hat{\beta} - \beta_0)$  where remainders are controlled using the smoothness of the relevant functions and the fact that under our assumptions the estimators are all root-n consistent.

To describe the results, let

$$\begin{aligned}
\rho_i &= \rho(z_i, \beta_0), \rho_{\beta i} = \partial \rho_i(\beta_0) / \partial \beta, \eta_i = \rho_{\beta i} - d_i, q_i = q^K(x_i), \kappa_i = E[\rho_i^4 | x_i] / \sigma_i^4, \\
\Upsilon &= \sum_{i=1}^n \sigma_i^2 q_i q_i' / n, \Gamma = \sum_i q_i d_i' / n, \tau = \Omega^{*-1} t, \xi_{ij} = q_i \Upsilon^{-1} q_j' / n, E[\tau' \eta_i \rho_i | x_i] = \sigma_i^{\rho \eta}, \\
\Pi &= \sum_{i=1}^n \xi_{ii} \sigma_i^{\rho \eta}, \Pi_B = \sum_{i,j=1}^n \sigma_i^{\rho \eta} \sigma_j^{\rho \eta} \xi_{ij}^2, \Lambda = \sum_{i=1}^n \xi_{ii} E[(\tau' \eta_i)^2 | x_i], \\
\Xi &= \sum_{i=1}^n \xi_{ii} (\tau' d_i)^2 (5 - \kappa_i), \Xi_{GEL} = \sum_{i=1}^n \xi_{ii} (\tau' d_i)^2 (3 - \kappa_i),
\end{aligned}$$

where we suppress the  $K$  argument for notational convenience. The terms involving fourth moments of the residuals are due to estimation of the weight matrix  $\Upsilon^{-1}$  for the optimal GMM estimator. This feature did not arise in the homoskedastic case considered in Donald and Newey (2001) where an optimal weight matrix depends only on the instruments.

### 3.2 Assumptions and Results

We impose the following fundamental condition on the data, the approximating functions  $q^K(x)$  and the distribution of  $x$ :

**Assumption 1 (Moments):** *Assume that  $z_i$  are i.i.d., and*

- (i)  $\beta_0$  is unique value of  $\beta$  in  $\mathcal{B}$  (a compact subset of  $\mathbb{R}^p$ ) satisfying  $E[\rho(z_i, \beta) | x_i] = 0$ ;
- (ii)  $\sum_{i=1}^n \sigma_i^{-2} d_i d_i' / n$  is uniformly positive definite and finite (w.p.1.).
- (iii)  $\sigma_i^2$  is bounded and bounded away from zero.
- (iv)  $E(\eta_{ji}^{\iota_1} \rho_i^{\iota_2} | x_i) = 0$  for any non-negative integers  $\iota_1$  and  $\iota_2$  such that  $\iota_1 + \iota_2 = 3$ .
- (v)  $E(\|\eta_i\|^\iota + |\rho_i|^\iota | x_i)$  is bounded for  $\iota = 6$  for GMM and BGMM and  $\iota = 8$  for GEL.

For identification, this condition only requires that  $E[\rho(z_i, \beta) | x_i] = 0$  has a unique solution at  $\beta = \beta_0$ . Estimators will be consistent under this condition because  $K$  is allowed to grow with  $n$ , as in Donald, Imbens, and Newey (2003). Part of this assumption is a

restriction that the third moments are zero. This greatly simplifies the MSE calculations. The last condition is a restriction on the moments that is used to control the remainder terms in the MSE expansion. The condition is more restrictive for GEL which has a more complicated expansion involving more terms and higher moments. The next assumption concerns the properties of the derivatives of the moment functions. Specifically, in order to control the remainder terms we will require certain smoothness conditions so that Taylor series expansions can be used and so that we can bound the remainder terms in such expansions.

**Assumption 2 (Expansion):** *Assume that  $\rho(z, \beta)$  is at least five times continuously differentiable in a neighborhood  $\mathcal{N}$  of  $\beta_0$ , with derivatives that are all dominated in absolute value by the random variable  $b_i$  with  $E(b_i^2) < \infty$  for GMM and BGMM and  $E(b_i^5) < \infty$  for GEL.*

This assumption is used to control remainder terms and has as an implication that for instance,

$$\sup_{\beta \in \mathcal{N}} \|(\partial/\partial\beta') \rho(z, \beta)\| < b_i$$

It should be noted that in the linear case only the first derivative needs to be bounded since all other derivatives would be zero. It is also interesting to note that although we allow for nonlinearities in the MSE calculations, they do not have an impact on the dominant terms in the MSE. The condition is stronger for GEL reflecting the more complicated remainder term. Our next assumption concerns the “instruments” represented by the vector  $q^K(x_i)$ .

**Assumption 3 (Approximation):** *(i) There is  $\zeta(K)$  such that for each  $K$  there is a nonsingular constant matrix  $B$  such that  $\tilde{q}^K(x) = Bp^K(x)$  for all  $x$  in the support of  $x_i$  and  $\sup_{x \in X} \|\tilde{q}^K(x)\| \leq \zeta(K)$  and  $E[\tilde{q}^K(x)\tilde{q}^K(x)']$  has smallest eigenvalue that is bounded away from zero, and  $\sqrt{K} \leq \zeta(K) \leq CK$  for some finite constant  $C$ . (ii) For each  $K$  there exists a sequence of constants  $\pi_K$  and  $\pi_K^*$  such that  $E(\|d_i - q_i'\pi_K\|^2) \rightarrow 0$  and  $\zeta(K)^2 E(\|d_i/\sigma_i^2 - q_i'\pi_K^*\|^2) \rightarrow 0$  as  $K \rightarrow \infty$ .*

The first part of the assumption gives a bound on the norm of the basis functions, and is used extensively in the MSE derivations to bound remainder terms. The second part of the assumption implies that  $d_i$  and  $d_i/\sigma_i^2$  be approximated by linear combinations of  $q_i$ . Because  $\sigma_i^2$  is bounded and bounded away from zero, it is easily seen that for the same coefficients  $\pi_K$ ,  $\|d_i/\sigma_i - \sigma_i q_i \pi_K^*\|^2 \leq \sigma_i^2 \|d_i/\sigma_i^2 - q_i' \pi_K\|^2$  so that  $d_i/\sigma_i$  can be approximated by a linear combination of  $\sigma_i q_i$ . Indeed the variance part of the MSE measures the mean squared error in the fit of this regression. Since  $\zeta(K) \rightarrow \infty$  the approximation condition for  $d_i/\sigma_i^2$  is slightly stronger than for  $d_i$ . This is to control various remainder terms where  $d_i/\sigma_i$  needs to be approximated in uniform manner. Since in many cases one can show that the expectations in (ii) are bounded by  $K^{-2\alpha}$  where  $\alpha$  depends on the smoothness of the function  $d_i/\sigma_i^2$ , the condition can be met by assuming that  $d_i/\sigma_i^2$  is a sufficiently smooth function of  $x_i$ .

We will assume that the preliminary estimator  $\tilde{\beta}$  used to construct the weight matrix is a GMM estimator is itself a GMM estimator with weighting matrix that may not be optimal. where we do not require either optimal weighting or that the number of moments increase. In other words we let  $\tilde{\beta}$  solve,

$$\min_{\beta} \tilde{g}_n(\beta)' \tilde{W}_0 \tilde{g}_n(\beta), \quad \tilde{g}_n(\beta) = (1/n) \sum_{i=1}^n \tilde{q}(x_i) \rho_i(\beta)$$

for some  $\tilde{K}$  vector of functions  $\tilde{q}(x_i)$  and some  $\tilde{K} \times \tilde{K}$  matrix  $\tilde{W}_0$  which potentially could be  $I_{\tilde{K}}$  or it could be random as would be the case if more than one iteration were used to obtain the GMM estimator. We make the following assumption regarding this preliminary estimator.

**Assumption 4 (Preliminary Estimator):** : Assume (i)  $\tilde{\beta} \xrightarrow{p} \beta_0$  (ii) there exist some non-stochastic matrix  $W_0$  such that  $\|\tilde{W}_0 - W_0\| \xrightarrow{p} 0$  and we can write  $\tilde{\beta} = \beta_0 + \frac{1}{n} \sum_{i=1}^n \tilde{\phi}_i \rho_i + o_p(n^{-1/2})$ ,  $\tilde{\phi}_i = -(\tilde{\Gamma}' W_0 \tilde{\Gamma})^{-1} \tilde{\Gamma}' W_0 \tilde{q}_i$  with  $\tilde{\Gamma} = \sum_{i=1}^n \tilde{q}(x_i) d_i / n$  and  $E(\|\rho_i^2 \tilde{\phi}_i \tilde{\phi}_i'\|) < \infty$

Note that the assumption requires that we just use some root-n consistent and asymptotically normally distributed estimator. The asymptotic variance of the preliminary

estimator will be,

$$p \lim((\tilde{\Gamma}'W_0\tilde{\Gamma})^{-1}\tilde{\Gamma}'W_0\tilde{Y}W_0\tilde{\Gamma}((\tilde{\Gamma}'W_0\tilde{\Gamma})^{-1}), \tilde{Y} = \sum_{i=1}^n \tilde{q}(x_i)\tilde{q}(x_i)'\sigma_i^2/n$$

and if the preliminary estimator uses optimal weighting we can show that this coincides with  $p \lim \Omega^*$  provided that  $\tilde{K}$  increase with  $n$  in a way that the assumptions of Donald, Imbens and Newey (2003) are satisfied. Also note that for the GMM estimator we can write,

$$\hat{\beta} = \beta_0 + \frac{1}{n} \sum_{i=1}^n \phi_i^* \rho_i + o_p(n^{-1/2}), \phi_i^* = -\Omega^{*-1} d_i / \sigma_i^2$$

The covariance between the (normalized) preliminary estimator and the GMM estimator is then,

$$\frac{1}{n} \sum_i \tilde{\phi}_i \phi_i^* \sigma_i^2 = \Omega^{*-1}$$

a fact that will be used in the MSE derivations to show that the MSE for BGMM does not depend on the preliminary estimator. Finally we use Assumption 6 of Donald, Imbens and Newey (2003) for the GEL class of estimators.

**Assumption 5 (GEL):**  $s(v)$  is at least five times continuously differentiable and concave on its domain, which is an open interval containing the origin,  $s_1(0) = -1$ ,  $s_2(v) = -1$  and  $s_j(v)$  is bounded in a neighborhood of  $v = 0$  for  $j = 1, \dots, 5$ .

The following three propositions give the MSE results for the three estimators considered in this paper. The results are proved in Appendix A and use an expansion that is provided in Appendix B.

**Proposition 1:** For GMM under Assumptions 1 - 4, if w.p.a.1 as  $n \rightarrow \infty$ ,  $|\Pi| \geq cK$  for some  $c > 0$ ,  $K \rightarrow \infty$ , and  $\zeta(K)\sqrt{K/n} \rightarrow 0$  then the approximate MSE for  $t'\sqrt{n}(\hat{\beta}^H - \beta_0)$  is given by,

$$S^H(K) = \Pi^2/n + \tau'(\Omega^* - \Gamma'\Upsilon^{-1}\Gamma)\tau$$

**Proposition 2:** For BGMM under Assumptions 1 -4, the condition that w.p.a.1 as  $n \rightarrow \infty$ ,

$$\Lambda + \Pi_B + \Xi_1 \geq cK$$

for some  $c > 0$ , and assuming that  $K \rightarrow \infty$  with  $\zeta(K)^2 \sqrt{K/n} \rightarrow 0$  the approximate MSE for  $t' \sqrt{n}(\hat{\beta}^B - \beta_0)$  is given by,

$$S^B(K) = [\Lambda + \Pi_B + \Xi] / n + \tau'(\Omega^* - \Gamma' \Upsilon^{-1} \Gamma) \tau$$

**Proposition 3:** For GEL, if Assumptions 1 - 3, 5 are satisfied, w.p.a.1 as  $n \rightarrow \infty$ ,

$$\{\Lambda - \Pi_B + \Xi + s_3 \Xi_{GEL}\} \geq cK,$$

$K \rightarrow \infty$ , and  $\zeta(K)^2 K^2 / \sqrt{n} \rightarrow 0$  the approximate MSE for  $t' \sqrt{n}(\hat{\beta}^{GEL} - \beta_0)$  is given by,

$$S^{GEL}(K) = [\Lambda - \Pi_B + \Xi + s_3 \Xi_{GEL}] / n + \tau'(\Omega^* - \Gamma' \Upsilon^{-1} \Gamma) \tau$$

For comparison purposes, and to help interpret the formulas, it is useful to consider the homoskedastic case. Let

$$\begin{aligned} \sigma^2 &= E[\rho_i^2], \sigma_{\eta\rho} = E[\tau' \eta_i \rho_i], \sigma_{\eta\eta} = E[(\tau' \eta_i)^2], \kappa = E[\rho_i^4] / \sigma^4, Q_{ii} = q_i' \left( \sum_{j=1}^n q_j q_j' \right)^{-1} q_i \\ \Delta(K) &= \sigma^{-2} \left\{ \sum_{i=1}^n (\tau' d_i)^2 - \sum_{i=1}^n \tau' d_i q_i' \left( \sum_{i=1}^n q_i q_i' \right)^{-1} \sum_{i=1}^n \tau' d_i q_i \right\} / n, \end{aligned}$$

Then we have the following expressions under homoskedasticity,

$$\begin{aligned} S^H(K) &= (\sigma_{\rho\eta} / \sigma^2)^2 K^2 / n + \Delta(K), \\ S^B(K) &= (\sigma_{\eta\eta} / \sigma^2 + \sigma_{\eta\rho}^2 / \sigma^4) K / n + \sigma^{-2} (5 - \kappa) \sum_i (\tau' d_i)^2 Q_{ii} / n + \Delta(K), \\ S^{GEL}(K) &= (\sigma_{\eta\eta} / \sigma^2 - \sigma_{\eta\rho}^2 / \sigma^4) K / n + \sigma^{-2} [(5 - \kappa) + s_3 (3 - \kappa)] \sum_i (\tau' d_i)^2 Q_{ii} / n + \Delta(K). \end{aligned}$$

For GMM, the MSE is the same as that presented in Donald and Newey (2001) for 2SLS, which is the same as Nagar (1959) for large numbers of moments. The leading  $K/n$  term in the MSE of BGMM is the same as the MSE of the bias-corrected 2SLS estimator, but

there is also an additional term, where  $(5 - \kappa)$  appears, that is due to the presence of the estimated weighting matrix. This term is also present for GMM, but is dominated by the  $K^2/n$  bias term, and so does not appear in our large  $K$  approximate MSE. As long as  $\kappa < 5$ , this additional term adds to the MSE of the estimator, representing a penalty for using a heteroskedasticity robust weighting matrix. When  $\kappa > 5$ , using the heteroskedasticity robust weighting matrix lowers the MSE, a phenomenon that was considered in Koenker et. al. (1994).

For GEL the leading  $K/n$  term is the same as for LIML, and is smaller than the corresponding term for BGMM. This comparison is identical to that for 2SLS and LIML, and represents an efficiency improvement from using GEL. For the CUE (or any other estimator where  $s_3 = 0$ ) the additional term is the same for BGMM and CUE, so that CUE has smaller MSE. The comparison among GEL estimators depends on the kurtosis  $\kappa$ . For Gaussian  $\rho(z, \beta_0)$ ,  $\kappa = 3$ , and the MSE of all the GEL estimators is the same. For  $\kappa > 3$ , the MSE of EL is greater than ET which is greater than CUE, with the order reversed for  $\kappa < 3$ . For Gaussian disturbances the relationships between the asymptotic MSE of LIML, BGMM, and EL were reported by Rothenberg (1996), though expressions were not given.

When there is heteroskedasticity, the comparison between estimators is exactly analogous to that for homoskedasticity, except that the results for LIML and B2SLS no longer apply. In particular, CUE has smaller MSE than BGMM, and BGMM and all GEL estimators have smaller MSE than GMM for large enough  $K$ . Since the comparisons are so similar, and since many of them were also discussed in the last Section, we omit them for brevity.

## 4 Monte Carlo Experiments

In this section we examine the performance of the different estimators and moment selection criteria in the context of a small scale Monte Carlo experiment based on the setup in Hahn and Hausman (2002) that was also used in Donald and Newey (2001).

The basic model used is of the form,

$$\begin{aligned} y_i &= \gamma Y_i + \rho_i \\ Y_i &= X_i' \pi + \eta_i \end{aligned} \tag{4.3}$$

for  $i = 1, \dots, n$  and the moment functions take the form (for  $K$  instruments),

$$g_i(\gamma) = \begin{pmatrix} X_{1i} \\ X_{2i} \\ \vdots \\ X_{Ki} \end{pmatrix} (y_i - \gamma Y_i)$$

where we are interested in methods for determining how many of the  $X_{ji}$  should be used to construct the moment functions. Because of the invariance of the estimators to the value of  $\gamma$  we set  $\gamma = 0$  and for different specifications of  $\pi$  we generate artificial random samples under the assumptions that

$$E \left( \begin{pmatrix} \rho_i \\ \eta_i \end{pmatrix} \begin{pmatrix} \rho_i & \eta_i \end{pmatrix} \right) = \Sigma = \begin{pmatrix} 1 & c \\ c & 1 \end{pmatrix}$$

and  $X_i \sim N(0, I_{\bar{K}})$  where  $\bar{K}$  is the maximal number of instruments considered. As shown in Hahn and Hausman (2002) the specification implies a theoretical first stage  $R$ -squared that is of the form,

$$R_f^2 = \frac{\pi' \pi}{\pi' \pi + 1} \tag{4.4}$$

We consider one of the models that was considered in Donald and Newey (2001) where,

$$\pi_k^2 = c(\bar{K}) \left( 1 - \frac{k}{\bar{K} + 1} \right)^4 \text{ for } k = 1, \dots, \bar{K}$$

where the constant  $c(\bar{K})$  is chosen so that  $\pi' \pi = R_f^2 / (1 - R_f^2)$ . In this model all the instruments are relevant but they have coefficients that are declining. This represents a situation where one has prior information that suggests that certain instruments are more important than others and the instruments have been ranked accordingly. In this model all of the potential  $\bar{K}$  moment conditions should be used for the estimators to be asymptotically efficient. Note also, that in our setup LIML and 2SLS are also asymptotically efficient estimators provided that we eventually use all of the instruments  $X_{ji}$ .



Indeed in the experiments we compute not only GMM, BGMM, ET, EL and CUE (the last three being members of the GEL class) but we also examine the performance of 2SLS and LIML along with the instrument selection methods proposed in Donald and Newey (2001). This allows us to gauge the small sample cost of not imposing heteroskedasticity. As in Donald and Newey (2001) we report for each of the seven different estimators, summary statistics for the version that uses all available instruments or moment conditions plus the summary statistics for the estimators based on a set of moment conditions or instruments that were chosen using the respective moment or instrument selection criterion.

For each model experiments were conducted with the specifications for sample sizes of  $n = 200$  and  $n = 1000$ . When the sample size is 200 we set  $R_f^2 = 0.1$ ,  $\bar{K} = 10$  and performed 500 replications, while in the larger sample size we set  $R_f^2 = 0.1$ ,  $\bar{K} = 20$  and we performed 200 replications (due to time constraints). Both of these choices reflect the fairly common situation where there may be a relatively small amount of correlation between the instruments and the endogenous variable (see Staiger and Stock (1997) and Stock and Wright (2000) as well as the fact that with larger data sets empirical researchers are more willing to use more moment conditions to improve efficiency. For each of these cases we consider  $c \in \{.1, .5, .9\}$ , though for brevity we will only report results for  $c = .5$ . In addition we consider the impact of having excess kurtosis, which as noted above has differential effect on the higher order MSE across the different estimators. The distribution we consider is that of

$$\begin{pmatrix} \rho_i \\ \eta_i \end{pmatrix} = |e_i| \begin{pmatrix} \rho_i^* \\ \eta_i^* \end{pmatrix}, \begin{pmatrix} \rho_i^* \\ \eta_i^* \end{pmatrix} \sim N(0, \Sigma), e_i \sim \text{logistic}(0,1).$$

where  $e_i$  is independent of  $\rho_i^*$  and  $\eta_i^*$  and is distributed as a logistic random variable with mean zero and variance equal to one. Given this particular setup we will have that  $(\rho_i, \eta_i)$  are jointly distributed with mean zero and a covariance matrix equal to  $\Sigma$ , and a coefficient of kurtosis of approximately  $\kappa = 12.6$ . With two different models, two different distributions for the errors, and three different choices for residual correlations there are a total of 12 specifications for each sample size.

The estimator that uses all moments or instruments is indicated by the suffix “-all” while the estimator that uses a number of moment conditions as chosen by the respective moment or instrument selection criterion is indicated by “-op”. Therefore, for instance, GMM-all and GMM-op are the two step estimator that uses all of the moment conditions and the moment conditions the minimize the estimated MSE criterion respectively. The preliminary estimates of the objects that appear in the criteria were in each case based on a number of moment conditions that was optimal with respect to cross validation in the first stage.

As in Donald and Newey (2001) we present robust measures of central tendency and dispersion. We computed the median bias (Med. Bias) for each estimator, the median of the absolute deviations (MAD) of the estimator from the true value of  $\gamma = 0$  and examined dispersion through the difference between the 0.1 and 0.9 quantile (Dec. Rge) in the distribution of each estimator. We also examined statistical inference by computing the coverage rate for 95% confidence intervals as well as the rejection rate for an overidentification test (in cases where overidentifying restrictions are present) using the test statistic corresponding to the estimator and a significance level of 5%. In addition we report some summary statistics concerning the choices of  $K$  in the experiments, including the modal choice of  $K$  if one used the actual MSE to choose  $K$ . There was very little dispersion in this variable across replications and generally the optimal  $K$  with the true criterion was equal to the same value in most if not all replications. In cases where there was some dispersion it was usually either being some cases on either side of the mode. To indicate such cases we use + and -, so that for instance 3+ means that the mode was 3 but that there were some cases where 4 was optimal. The notation 3++ means that the mode was 3 but that a good proportion of the replications had 4 as being optimal.

Tables I-VI and X-XV contains the summary statistics for the estimators for  $n = 200$  and  $n = 800$  respectively, while Tables VII-IX and XVI-XVIII contain the summary statistics for the chosen number of moments across the replications. In general the results are encouraging for all the estimators. As expected the GEL and LIML estimators are less dispersed when the optimal number of moments is used, while for GMM and 2SLS

the use of the criterion reduces the bias that occurs when there is a high degree of covariance between the residuals. The improvements for the GEL estimators are more marked when there is a low to moderate degree of covariance. It is noteworthy that in such situations there is also a dramatic improvement in the quality of inference as indicated by the coverage rates for the confidence interval. As far as testing the overidentifying restrictions only when there is a high degree of covariance is there any problem with testing these restrictions. This occurs with most of the estimators in the small sample with a high covariance and with GMM and TSLS in the large sample with a high covariance. It also seems that using the criteria does not really help in fixing any of these problems.

There are a number of things to note about the results for  $\hat{K}$ . First, the estimated criteria give values for  $\hat{K}$  that are often near the values that minimize the true criterion, suggesting that the estimated criterion is a good approximation to the true criterion. It is also noteworthy that, as one would expect, the criteria suggest use of a small number of moments for GMM and 2SLS when there is a high error covariance and for the GEL estimators when there is a low covariance. For BGMM the optimal number is quite stable as the covariance increases. In the larger sample the optimal number decreases as the covariance increases, but is slightly larger when the residuals have fat tails compared to the situation where they do not. Among the GEL estimators increasing the covariance and having fat tailed errors has the most dramatic impact on CUE as one would expect given the criteria.

Concerning the effect of excess kurtosis, it does appear that the improvement from using the criteria is more noticeable for EL, which is most sensitive to having fat tailed errors. There also was some evidence that going from normal to fat tailed errors helped CUE more than the other estimators, as suggested in the theory, although this led to a lower improvement from using the moment selection criterion.

## 5 Conclusion

In this paper we have developed approximate MSE criteria for moment selection for a variety of estimators in conditional moment contexts. We found that the CUE has smaller MSE than the bias corrected GMM estimator. In addition we proposed data based methods for estimating the approximate MSE, so that in practice the number of moments can be selected by minimizing these criteria. The criteria seemed to perform adequately in a small scale simulation exercise.

The present paper has considered a restrictive environment in which the data are considered a random sample. It would be useful to extend the results in two directions. The first would be to the dynamic panel data case. In that situation there will typically be different sets of instruments available for each residual coming from sequential moment restrictions. It would also be useful to extend the results to a purely time series context where one would need to deal with serial correlation. Kuersteiner (2002) has derived interesting results in this direction.

## Appendix A: MSE Derivation Results

Throughout the Appendix repeated use of the Cauchy Schwarz (CS), Markov (M) and Triangle inequalities is made. We let  $\|\cdot\|$  denote the usual matrix norm. The following Maximum Eigenvalue (ME) inequality is also used repeatedly,

$$\|A'BC\| \leq \lambda_{\max}(BB) \|A\| \|C\|$$

for a square symmetric matrix  $B$  and conformable matrices  $A$  and  $C$ . For simplicity of notation and without loss of generality we assume that the true value of the coefficients are all zero and we only perform the calculation for the case where there is one parameter (in addition to the auxiliary parameters  $\lambda$ ). Because higher order derivatives are required for the MSE calculations (even though they do not appear in the final result) we use the following notation: for  $j = 0, 1, \dots, 4$  we let,

$$\begin{aligned} \Gamma_j &= \frac{1}{n} \sum_i \Gamma_{ji}, \bar{\Gamma}_j = \frac{1}{n} \sum_i \Gamma_{ji}(0), \hat{\Gamma}_j = \frac{1}{n} \sum_i \Gamma_{ji}(\hat{\beta}) \\ \Gamma_{ji} &= q_i E\left(\frac{\partial^{j+1}}{\partial \beta^{j+1}} \rho_i(0) | x_i\right), \Gamma_{ji}(\beta) = q_i \frac{\partial^{j+1}}{\partial \beta^{j+1}} \rho_i(\beta) \\ \eta_{ji} &= \frac{\partial^{j+1}}{\partial \beta^{j+1}} \rho_i(0) - E\left(\frac{\partial^{j+1}}{\partial \beta^{j+1}} \rho_i(0) | x_i\right) \end{aligned}$$

and  $\bar{\Gamma}_j^*$  denotes  $\bar{\Gamma}_j$  evaluated at some point(s) lying between the respective estimator and its true value. Hence  $\Gamma_0$  corresponds to  $\Gamma$  in the text. In addition we assume as in Donald, Imbens and Newey (2003) (hereafter DIN) that  $q_i$  has been normalized so that  $\|q_i\| < C\zeta(K)$  and  $E(q_i q_i') = I_K$  so that,

$$\begin{aligned} \frac{1}{n} \sum_i \|q_i\|^2 &= O(K), \lambda_{\max}\left(\left(\frac{1}{n} \sum_i q_i q_i' \delta_i\right)^2\right) = O(1), \text{ if } \delta_i < C < \infty \\ \lambda_{\min}\left(\frac{1}{n} \sum_i q_i q_i' \delta_i\right) &> 0, \text{ for } 0 < c < \delta_i < C < \infty \end{aligned}$$

where here and elsewhere we let  $c$  denote a generic small constant and  $C$  a generic large constant. The MSE are based on an expansion that is contained in Appendix B. The remainder term for GEL is dealt with in a technical appendix that is available on request.

Derivatives that are used in the expansion are also available in a technical appendix that is available on request.

**Proof of Proposition 1:**

In deriving the MSE for GMM we simplify notation and refer to the estimator as  $\hat{\beta}$  as distinct from the preliminary estimator. Since we need the results for BGMM we expand the estimator and display all terms that are needed to perform the calculations in Proposition 2. Terms that will not be needed for GMM and BGMM are those that are  $o(K^2/n^{3/2})$  and  $o(K/n^{3/2})$  respectively. Now for GMM we have from Newey and Smith (2004) (hereafter NS) we have that GMM can be written as the solution to the First Order Conditions,

$$\begin{aligned}\frac{1}{n} \sum_i \frac{\partial}{\partial \beta} \rho_i(\hat{\beta}) q_i' \hat{\lambda} &= 0 \\ \frac{1}{n} \sum_i q_i \rho_i(\hat{\beta}) + \hat{\Upsilon}(\tilde{\beta}) \hat{\lambda} &= 0\end{aligned}$$

where by DIN  $\|\hat{\lambda}\| = O(\sqrt{K/n})$  and  $\|\hat{\beta}\| = O(1/\sqrt{n})$ . Using Appendix B and the partitioned inverse formula applied to

$$M^{-1} = \begin{pmatrix} 0 & \Gamma_0' \\ \Gamma_0 & \Upsilon^{-1} \end{pmatrix}^{-1} = \begin{pmatrix} -\Omega^{-1} & \Omega^{-1} \Gamma_0' \Upsilon^{-1} \\ \Upsilon^{-1} \Gamma_0 \Omega^{-1} & \Sigma \end{pmatrix}$$

where  $\Omega = \Gamma_0' \Upsilon^{-1} \Gamma_0$  and  $\Sigma = \Upsilon^{-1} - \Upsilon^{-1} \Gamma_0 \Omega^{-1} \Gamma_0' \Upsilon^{-1}$  we have that,

$$-M^{-1}m = \begin{pmatrix} -\Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g} \\ -\Sigma \bar{g} \end{pmatrix} = \begin{pmatrix} T_1^\beta \\ T_1^\lambda \end{pmatrix}$$

Note that we have  $\|\Omega^{-1}\| = O(1)$ ,  $\|\Upsilon^{-1} \Gamma_0\| = O(1)$  by ME  $\|\Gamma_0\| = O(1)$  and  $\lambda_{\max}(\Upsilon^{-1}) = O(1)$  and finally,  $\lambda_{\max}(\Sigma) \leq \lambda_{\max}(\Upsilon^{-1}) = O(1)$ . Similarly by Appendix B and the technical appendix,

$$\begin{aligned}-M^{-1}(\hat{M} - M)\theta &= \begin{pmatrix} \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} - \Omega^{-1} \Gamma_0' \Upsilon^{-1} \left( (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon) \hat{\lambda} + (\bar{\Gamma}_0 - \Gamma_0) \hat{\beta} \right) \\ -\Upsilon^{-1} \Gamma_0 \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} - \Sigma \left( (\bar{\Gamma}_0 - \Gamma_0) \hat{\beta} + (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon) \hat{\lambda} \right) \end{pmatrix} \\ &= \begin{pmatrix} T_2^\beta \\ T_2^\lambda \end{pmatrix} = \begin{pmatrix} O(K/n) + O(K/n) + O(1/n) \\ O(K/n) + O(\sqrt{K}/n) + O(\zeta(K)K/n) \end{pmatrix} \\ \begin{pmatrix} T_3^\beta \\ T_3^\lambda \end{pmatrix} &= -M^{-1} \sum_j \theta_j A_j \theta / 2 = \begin{pmatrix} \Omega^{-1} \hat{\beta} \Gamma_1' \hat{\lambda} - (1/2) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\beta}^2 \Gamma_1 \\ -\Upsilon^{-1} \Gamma_0 \Omega^{-1} \hat{\beta} \Gamma_1' \hat{\lambda} - (1/2) \Sigma \hat{\beta}^2 \Gamma_1 \end{pmatrix} = \begin{pmatrix} O(1/n) \\ O(1/n) \end{pmatrix}\end{aligned}$$

$$\begin{aligned}
-M^{-1} \sum_j \theta_j (\hat{A}_j - A_j) \theta &= \begin{pmatrix} \Omega^{-1} \hat{\beta} (\bar{\Gamma}_1 - \Gamma_1)' \hat{\lambda} - (1/2) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\beta}^2 (\bar{\Gamma}_1 - \Gamma_1) \\ -\Upsilon^{-1} \Gamma_0 \Omega^{-1} \hat{\beta} (\bar{\Gamma}_1 - \Gamma_1)' \hat{\lambda} - (1/2) \Sigma \hat{\beta}^2 (\bar{\Gamma}_1 - \Gamma_1) \end{pmatrix} \\
&= \begin{pmatrix} T_{41}^\beta + T_{42}^\beta \\ T_{41}^\lambda + T_{42}^\lambda \end{pmatrix} = \begin{pmatrix} O(K/n^{3/2}) + O(\sqrt{K}/n^{3/2}) \\ O(K/n^{3/2}) + O(\sqrt{K}/n^{3/2}) \end{pmatrix}
\end{aligned}$$

Next we have that,

$$\begin{aligned}
\begin{pmatrix} T_5^\beta \\ T_5^\lambda \end{pmatrix} &= -M^{-1} \sum_j \sum_k \theta_j \theta_k B_{jk} \theta / 6 = -(1/6) \begin{pmatrix} -\Omega^{-1} 3 \hat{\beta}^2 \Gamma_2' \hat{\lambda} + \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\beta}^3 \Gamma_2 \\ \Upsilon^{-1} \Gamma_0 \Omega^{-1} 3 \hat{\beta}^2 \Gamma_2' \hat{\lambda} + \Sigma \hat{\beta}^3 \Gamma_2 \end{pmatrix} \\
&= \begin{pmatrix} O(\|\hat{\beta}\|^2 \|\Gamma_2\| \|\hat{\lambda}\|) \\ O(\|\hat{\beta}\|^3 \|\Gamma_2\|) \end{pmatrix} = \begin{pmatrix} O(\sqrt{K}/n^{3/2}) \\ O(1/n^{3/2}) \end{pmatrix}
\end{aligned}$$

by CS, the results for  $\|\hat{\beta}\|$  and  $\|\hat{\lambda}\|$ , Assumption 3, and the conditions on the elements of  $M^{-1}$ . Next we have,

$$\begin{aligned}
\sum_j \sum_k \theta_j \theta_k (\hat{B}_{jk} - B_{jk}) \theta / 6 &= (1/6) \begin{pmatrix} 3 \hat{\beta}^2 (\bar{\Gamma}_2 - \Gamma_2)' \hat{\lambda} \\ \hat{\beta}^3 (\bar{\Gamma}_2 - \Gamma_2) \end{pmatrix} \\
&= \begin{pmatrix} O(\|\hat{\beta}\|^2 \|\bar{\Gamma}_2 - \Gamma_2\| \|\hat{\lambda}\|) \\ O(\|\hat{\beta}\|^3 \|\bar{\Gamma}_2 - \Gamma_2\|) \end{pmatrix} = \begin{pmatrix} O(K/n^2) \\ O(\sqrt{K}/n^2) \end{pmatrix}
\end{aligned}$$

The last term is  $(1/24)$  times,

$$\begin{aligned}
\sum_j \sum_k \theta_j \theta_k \sum_l \theta_l \hat{C}_{jkl} \theta &= \begin{pmatrix} \hat{\beta}^4 (\bar{\Gamma}_4^* \hat{\lambda}) + 4 \hat{\beta}^3 \bar{\Gamma}_3^* \hat{\lambda} \\ \bar{\Gamma}_3^* \hat{\beta}^4 \end{pmatrix} \\
&= \begin{pmatrix} O(\|\hat{\beta}\|^3 \|\hat{\lambda}\| (\|\hat{\beta}\| \|\bar{\Gamma}_4^*\| + \|\bar{\Gamma}_3^*\|)) \\ O(\|\bar{\Gamma}_3^*\| \|\hat{\beta}\|^4) \end{pmatrix} = O(\sqrt{K}/n^2)
\end{aligned}$$

Therefore,  $\|R_{n,K}^\beta\| = O(K/n^2)$  and  $\|R_{n,K}^\lambda\| = O(K/n^2)$  by CS and ME and the condition on the elements of  $M^{-1}$ . Here we have used  $\|\Gamma_0' \Upsilon^{-1} (\bar{\Gamma}_0 - \Gamma_0)\| = O(1/\sqrt{n})$  and  $\|\Gamma_0' \Upsilon^{-1} (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon)\| = O(\sqrt{K}/n + \zeta(K)/\sqrt{n}) = O(\zeta(K)/\sqrt{n})$  which follow from DIN Lemma A4.

Note that by Assumption 2 we can write,

$$\begin{aligned}
\hat{\Upsilon}(\tilde{\beta}) - \Upsilon &= \hat{\Upsilon}_{v1} + 2\tilde{\beta} \Upsilon_{\rho\eta} + 2\tilde{\beta} \hat{\Upsilon}_{v2} + 2\tilde{\beta} \hat{\Upsilon}_{v3} + \tilde{\beta}^2 \hat{\Upsilon}_r \\
&= \hat{\Upsilon}_{v1} + 2\tilde{\beta} \Upsilon_{\rho\eta} + \hat{R}_{n,K}^\Upsilon
\end{aligned}$$

where,

$$\begin{aligned}\hat{\Upsilon}_{vj} &= \frac{1}{n} \sum_i q_i q'_i v_{ji} \text{ for } j = 1, 2, 3 \text{ } E(v_{ji}|x_i) = 0 \\ v_{1i} &= \rho_i^2 - \sigma_i^2, v_{2i} = d_i \rho_i, v_{3i} = (\eta_{0i} \rho_i - \sigma_{\rho\eta}(x_i)) \\ \Upsilon_{\rho\eta} &= \frac{1}{n} \sum_i q_i q'_i \sigma_{\rho\eta}(x_i), \hat{\Upsilon}_r = \frac{1}{n} \sum_i q_i q'_i r_i\end{aligned}$$

and  $E(v_{ji}|x_i) = 0$ ,  $E(v_{ji}^2|x_i) < C$  Assumption 1.  $E(\|r_i\| | x_i) < C$  by Assumption 2. Hence, we have that,

$$\begin{aligned}\|\hat{\Upsilon}_{vj}\| &= O(\zeta(K)\sqrt{K/n}), \|\hat{R}_{n,K}^\Upsilon\| = O\left((\zeta(K)\sqrt{K} + K)/n\right) \\ \lambda_{\max}(\Upsilon_{\rho\eta} \Upsilon_{\rho\eta}) &= O(1), \lambda_{\max}(\hat{\Upsilon}_r \hat{\Upsilon}_r) = O(1) \\ \|\hat{\Upsilon}_{vj} \Sigma \bar{g}\| &= O(\zeta(K)\sqrt{K}/n) \text{ for } j = 1, 3 \\ \|\hat{\Upsilon}_{vj} \Sigma (\bar{\Gamma}_0 - \Gamma_0)\| &= O(\zeta(K)\sqrt{K}/n) \text{ for } j = 1, 3 \\ \|\hat{\Upsilon}_{v2} \Sigma \bar{g}\| &= O(\zeta(K)K/n), \|\hat{\Upsilon}_{v2} \Sigma (\bar{\Gamma}_0 - \Gamma_0)\| = O(\zeta(K)K/n) \\ \|\Gamma_0 \Upsilon^{-1} \hat{\Upsilon}_{vj} \Sigma\| &= O(\sqrt{K/n}), \|\Gamma_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Sigma\| = O(1), \|\Gamma_0 \Upsilon^{-1} \hat{\Upsilon}_r \Sigma\| = O(1) \\ \|\bar{\Gamma}_0 - \Gamma_0\| &= O(\sqrt{K/n}), \|(\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0\| = O(1/\sqrt{n})\end{aligned}$$

with the last fact following from M and Assumptions 1 and 3. For the  $\hat{\lambda}$  terms from the above expansion it follows that,

$$\hat{\lambda} = -\Sigma \bar{g} - \Upsilon^{-1} \Gamma_0 \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} - \Sigma (\bar{\Gamma}_0 - \Gamma_0) \hat{\beta} - \Sigma (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon) \hat{\lambda} + R_1^\lambda$$

where  $\|R_1^\lambda\| = O(1/n)$  under the condition on  $K$  for GMM and hence for BGMM also. Repeated substitution and using the facts that by CS and the results for  $\hat{\lambda}$  and  $\|\bar{\Gamma}_0 - \Gamma_0\|$  and the fact that from the above expansion  $\hat{\beta} = -\Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g} + R_1^\beta$  with  $\|R_1^\beta\| = O(K/n)$ ,

$$\begin{aligned}\hat{\lambda} &= -\Sigma \bar{g} + \Upsilon^{-1} \Gamma_0 \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} + \Sigma (\bar{\Gamma}_0 - \Gamma_0) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g} \\ &\quad + \Sigma \hat{\Upsilon}_{vj} \Sigma \bar{g} + 2\tilde{\beta} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} + R_2^\lambda, \\ \|\Upsilon^{-1} \Gamma_0 \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g}\| &= O(K/n) \|\Sigma (\bar{\Gamma}_0 - \Gamma_0) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g}\| = O(\sqrt{K}/n) \\ \|\Sigma \hat{\Upsilon}_{vj} \Sigma \bar{g}\| &= O(\zeta(K)\sqrt{K}/n), \|2\tilde{\beta} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g}\| = O(\sqrt{K}/n), \\ \|R_2^\lambda\| &= O(\zeta(K)\sqrt{K}/n) O(\zeta(K)\sqrt{K/n})\end{aligned}$$



where  $\|R_2^\lambda\| = o(\sqrt{K}/n)$  under the condition on  $K$  for BGMM and  $\|R_2^\lambda\| = O(1/\sqrt{n})o(\zeta(K)K/n) = O(1/\sqrt{n})o(K^2/n)$  under the condition on  $K$  for GMM.

Take the lead term in the expansion for  $\hat{\beta}$ ,  $T_2^\beta$  and substitute in for  $\hat{\lambda}$  and apply CS, and M to get that,

$$\begin{aligned}\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} &= -\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} + \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} \\ &\quad + \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma (\bar{\Gamma}_0 - \Gamma_0) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g} \\ &\quad + \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} + 2\tilde{\beta} \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} + R_2^\beta \\ R_2^\beta &= -\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' R_2^\lambda\end{aligned}$$

Here  $\|R_2^\beta\| = O(\sqrt{K}/n)o(\sqrt{K}/n) = O(1/\sqrt{n})o(K/n)$  under the condition for BGMM, and

$$\begin{aligned}-\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} &= O(K/n) \\ \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma (\bar{\Gamma}_0 - \Gamma_0) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \bar{g} &= O(1/\sqrt{n})O(K/n) \\ \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} &= O(1/\sqrt{n})O(K/n) \\ \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} &= O(1/\sqrt{n})O(\zeta(K)K/n) \\ 2\tilde{\beta} \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} &= O(1/\sqrt{n})O(K/n)\end{aligned}$$

Under the condition on  $K$  for GMM we have,

$$-\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} = -\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} + \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} + R_3^\beta$$

where  $\|R_3^\beta\| = O(1/\sqrt{n})o(K^2/n)$ .

Now consider the second term in  $T_2^\beta$  given by  $-\Omega^{-1} \Gamma_0' \Upsilon^{-1} (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon) \hat{\lambda}$ . Using the facts above we can write,

$$\begin{aligned}-\Omega^{-1} \Gamma_0' \Upsilon^{-1} (\hat{\Upsilon}(\tilde{\beta}) - \Upsilon) &= -\Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v1} - 2\tilde{\beta} \Omega^{-1} \Gamma_0' \Upsilon^{-1} \Upsilon_{\rho\eta} - 2\tilde{\beta} \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v2} + o(\sqrt{K}/n) \\ \left\| -\Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v1} \right\| &= O(\sqrt{K}/n), \quad \left\| 2\tilde{\beta} \Omega^{-1} \Gamma_0' \Upsilon^{-1} \Upsilon_{\rho\eta} \right\| = O(1/\sqrt{n}), \\ \left\| 2\tilde{\beta} \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v2} \right\| &= O(\sqrt{K}/n)\end{aligned}$$

so that using the above expansion for  $\hat{\lambda}$  and repeated use of CS and M,

$$\begin{aligned}
-\Omega^{-1}\Gamma'_0\Upsilon^{-1}\left(\hat{\Upsilon}(\tilde{\beta})-\Upsilon\right)\hat{\lambda} &= \Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\bar{g} + 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\bar{g} \\
&+ 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v2}\Sigma\bar{g} - \Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Upsilon^{-1}\Gamma_0\Omega^{-1}\left(\bar{\Gamma}_0-\Gamma_0\right)'\Sigma\bar{g} \\
&- 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Upsilon^{-1}\Gamma_0\Omega^{-1}\left(\bar{\Gamma}_0-\Gamma_0\right)'\Sigma\bar{g} \\
&- \Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} - 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} \\
&- 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} + R^\beta \\
\|R^\beta\| &= O(\sqrt{K/n})O(\|R_2^\lambda\|)
\end{aligned}$$

using the fact that also under the Assumption 1,

$$\left\|-\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\left(\bar{\Gamma}_0-\Gamma_0\right)\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g}\right\| = O(1/\sqrt{n})O(\sqrt{K}/n)$$

Here we have that,

$$\begin{aligned}
\left\|\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\bar{g}\right\| &= O(1/\sqrt{n})O(\sqrt{K/n}) \\
\left\|2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\bar{g}\right\| &= O(1/\sqrt{n})O(1/\sqrt{n}) \\
\left\|2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v2}\Sigma\bar{g}\right\| &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

$$\begin{aligned}
\left\|\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Upsilon^{-1}\Gamma_0\Omega^{-1}\left(\bar{\Gamma}_0-\Gamma_0\right)'\Sigma\bar{g}\right\| &= O(1/\sqrt{n})O(K/n) \\
\left\|2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Upsilon^{-1}\Gamma_0\Omega^{-1}\left(\bar{\Gamma}_0-\Gamma_0\right)'\Sigma\bar{g}\right\| &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

$$\begin{aligned}
\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} &= O(\sqrt{K/n})O(\zeta(K)\sqrt{K}/n) = O(1/\sqrt{n})O(\zeta(K)K/n) \\
2\tilde{\beta}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} &= O(1/\sqrt{n})O(\zeta(K)\sqrt{K}/n) \\
2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

and under the condition on  $K$  for BGMM  $\|R^\beta\| = O(1/\sqrt{n})o(K/n)$ . So under the conditions for GMM we have that  $\|R^\beta\| = O(1/\sqrt{n})o(K^2/n)$  and that by the above results,

$$\begin{aligned}
-\Omega^{-1}\Gamma'_0\Upsilon^{-1}\left(\hat{\Upsilon}(\tilde{\beta})-\Upsilon\right)\hat{\lambda} &= \Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\bar{g} + 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\bar{g} \\
&- \Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} - 2\tilde{\beta}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\hat{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} + R_3^\beta
\end{aligned}$$

with  $\|R_3^\beta\| = O(1/\sqrt{n})o(K^2/n)$ .

Now take the third term in  $T_2^\beta$  and using the above results and the fact that  $\|-\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\hat{\beta}\| = O(1/n)$  we have that

$$\hat{\beta} = -\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} - \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} + o(K/n)$$

so,

$$\begin{aligned} -\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\hat{\beta} &= \Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} \\ &\quad + \Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} + O(1/\sqrt{n})o(K/n) \end{aligned}$$

under the conditions for BGMM and

$$-\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\hat{\beta} = \Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} + O(1/\sqrt{n})o(K^2/n)$$

under the conditions for GMM where,

$$\begin{aligned} \|\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}\| &= O(1/\sqrt{n})O(1/\sqrt{n}) \\ \|\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g}\| &= O(1/\sqrt{n})O(K/n) \end{aligned}$$

Next we have for  $T_3^\beta$  by the above expansion for  $\hat{\lambda}$ ,

$$\begin{aligned} \Gamma_1'\hat{\lambda} &= -\Gamma_1'\Sigma\bar{g} - \Gamma_1'\Upsilon^{-1}\Gamma_0\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} + \Gamma_1'\Sigma(\bar{\Gamma}_0 - \Gamma_0)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} \\ &\quad + \Gamma_1'\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g} + 2\tilde{\beta}\Gamma_1'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} + \Gamma_1'R_2^\lambda, \end{aligned}$$

By showing that  $\|\Gamma_1'\Sigma\hat{\Upsilon}_{v1}\Sigma\bar{g}\| = O(\sqrt{K}/n)$  and  $\|\Gamma_1'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}\| = O(1/\sqrt{n})$  we can show that under the condition for BGMM

$$\Gamma_1'\hat{\lambda} = -\Gamma_1'\Sigma\bar{g} - \Gamma_1'\Upsilon^{-1}\Gamma_0\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} + o(K/n)$$

so that, after substituting in for  $\hat{\beta}$  we have,

$$\begin{aligned} &\Omega^{-1}\hat{\beta}\Gamma_1'\hat{\lambda} - (1/2)\hat{\beta}^2\Omega^{-1}\Gamma_0'\Upsilon^{-1}\Gamma_1 \\ &= \Omega^{-1}\Gamma_1'\Sigma\bar{g}\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} - \Gamma_1'\Upsilon^{-1}\Gamma_0\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g}\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} \\ &\quad + \Omega^{-1}\Gamma_1'\Sigma\bar{g}\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} - (1/2)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\Gamma_1(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g})^2 \\ &\quad - \Omega^{-1}\Gamma_0'\Upsilon^{-1}\Gamma_1\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} + O(1/\sqrt{n})o(K/n) \end{aligned}$$

where

$$\begin{aligned}
\Omega^{-1}\Gamma_1'\Sigma\bar{g}\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} &= O(1/\sqrt{n})O(1/\sqrt{n}) \\
\Gamma_1'\Upsilon^{-1}\Gamma_0\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g}\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} &= O(1/\sqrt{n})O(K/n) \\
\Omega^{-1}\Gamma_1'\Sigma\bar{g}\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} &= O(1/\sqrt{n})O(K/n) \\
(1/2)\Omega^{-1}\Gamma_0'\Upsilon^{-1}\Gamma_1(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g})^2 &= O(1/n) \\
\Omega^{-1}\Gamma_0'\Upsilon^{-1}\Gamma_1\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)'\Sigma\bar{g} &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

For the term  $T_4^\beta$  we have,

$$\begin{aligned}
T_4^\beta &= (1/2)2\hat{\beta}\Omega^{-1}(\bar{\Gamma}_1 - \Gamma_1)'\hat{\lambda} + O(1/\sqrt{n})o(K/n) \\
&= \Omega^{-1}(\bar{\Gamma}_1 - \Gamma_1)'\Sigma\bar{g}\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} + O(1/\sqrt{n})o(K/n) \\
\|\Omega^{-1}((\bar{\Gamma}_1 - \Gamma_1)'\Sigma\bar{g})\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}\| &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

Finally we note that we have  $\|T_5^\beta\| = O(1/\sqrt{n})o(K/n)$  so that we have now found all terms that are the right order.

For GMM let  $\gamma_{K,n} = K^2/n + \Delta_{K,n}$  with  $\Delta_{K,n} = \Omega^* - \Omega$ . Then by,

$$\begin{aligned}
n\left(\zeta(K)K/n^{3/2} + \sqrt{K}/n + 1/n\right)^2 &= o(K^2/n) \\
n\left(\zeta(K)K/n^{3/2} + \sqrt{K}/n + 1/n\right)K/\sqrt{n} &= O(\zeta(K)/\sqrt{n})O(K^2/n) = o(K^2/n)
\end{aligned}$$

so with  $h = -\Gamma_0'\Upsilon^{-1}\bar{g}$ , we can write,

$$\begin{aligned}
\sqrt{n}\hat{\beta}^H &= \sqrt{n}\Omega^{-1}(h + \sum_{j=1}^4 T_j^H + Z^H) \\
\|T_1^H\| &= O(K/\sqrt{n}), \|T_2^H\| = O(\sqrt{K/n}), \|T_3^H\| = O(1/\sqrt{n}) \\
\|T_4^H\| &= O(\zeta(K)K/n), \|Z^H\| = o(K^2/n)
\end{aligned}$$

Using the expansion,  $\Omega^{-1} = \Omega^{*-1} + \Omega^{*-1}(\Omega^* - \Omega)\Omega^{*-1} + O(\Delta_{K,n}^2)$  and noting that under the condition on  $K$ ,  $T_1^H = o(1)$  then we can write,

$$\left(\sqrt{n}\hat{\beta}^H\right)^2 = n\Omega^{-1}hh'\Omega^{-1} + \Omega^{*-1}\left(T_1^HT_1^H + 2\sum_{j=1}^4 hT_j^H\right)\Omega^{*-1} + o(\gamma_{K,n})$$

$$\Omega^{-1}E(hh')\Omega^{-1} = \Omega^{-1} = \Omega^{*-1} + \Omega^{*-1}(\Omega^* - \Omega)\Omega^{*-1} + O(\Delta_{K,n}^2)$$

with  $E(hT_1^H) = E(hT_2^H) = 0$  by the third moment condition. Next,

$$E(nT_1^H T_1^H) = \frac{1}{n} \left( \sum_i E(\rho_i \eta_{0i}) \xi_{ii} \right)^2 + o(\gamma_{K,n})$$

and by the third moment condition,

$$\begin{aligned} nE(T_1^H h) &= nE(\Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \bar{g} h) = \frac{1}{n} \sum_i (\kappa_i - 1) d_i^2 \xi_{ii} + o(\gamma_{K,n}) \\ &= O(K/n) = o(\gamma_{K,n}) \end{aligned}$$

$$\begin{aligned} nE(T_3^H h) &= -nE((\bar{\Gamma}_0 - \Gamma_0)' \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} h) + nE(\Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} h) \\ &\quad - 2nE(\tilde{\beta} \Gamma_0' \Upsilon^{-1} \Upsilon_{\rho\eta} \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g}) \\ &= -\frac{1}{n} \sum_i (\kappa_i - 1) d_i^2 \xi_{ii} + o(\gamma_{K,n}) \end{aligned}$$

Therefore we have the result,

$$\begin{aligned} E(nt\Omega^{-1} \left( h^2 + 2hT_h^1 + (T_h^1)^2 + 2hT_h^2 \right) \Omega^{-1}t) \\ = \Omega^{*-1} + \Pi^2/n + (\Omega^* - \Gamma_0' \Upsilon^{-1} \Gamma_0) + o(\gamma_{K,n}) \end{aligned}$$

using the fact that  $\Omega^* \Omega^{-1} = I + (\Omega^* - \Omega) \Omega^{*-1} + O(\Delta_{K,n}^2)$  and the fact that  $\tau = \Omega^* t$

**Proof of Proposition 2:** For BGMM we have the additional term,

$$\begin{aligned} (\hat{\Gamma}_0' \hat{\Upsilon}^{-1} \hat{\Gamma}_0)^{-1} \sum_{i=1}^n \hat{\Gamma}_{0i}' \hat{\Sigma} \hat{g}_i / n^2 &= \hat{\Omega}^{-1} \sum_{i=1}^n \hat{\Gamma}_{0i}' \hat{\Sigma} \hat{g}_i / n^2 \\ \hat{\Gamma}_0 &= \sum_{i=1}^n \hat{\Gamma}_{0i} / n, \hat{\Gamma}_{0i} = q_i \hat{y}_i, \hat{y}_i = [\partial \rho_i(\hat{\beta}^H) / \partial \beta]', \hat{g}_i = q_i \rho_i(\hat{\beta}^H) \\ \hat{\Sigma} &= \hat{\Upsilon}^{-1} - \hat{\Upsilon}^{-1} \hat{\Gamma}_0 \hat{\Omega}^{-1} \hat{\Gamma}_0' \hat{\Upsilon}^{-1}, \hat{\Omega} = \hat{\Gamma}_0' \hat{\Upsilon}^{-1} \hat{\Gamma}_0 \end{aligned}$$

First we have by Assumption 3, letting  $\Gamma_{0i} = q_i E(\partial \rho_i(0) / \partial \beta)$ ,  $\Gamma_{0i}(0) = q_i \partial \rho_i(0) / \partial \beta$

$$\begin{aligned} \hat{\Gamma}_{0i} &= \Gamma_{0i} + (\Gamma_{0i}(0) - \Gamma_{0i}) + \hat{\beta} \Gamma_{1i} + \hat{\beta}(\Gamma_{1i}(0) - \Gamma_{1i}) + \left( \hat{\beta} \right)^2 \Gamma_{2i}^* \\ \hat{g}_i &= g_i + \Gamma_{0i} \hat{\beta} + (\Gamma_{0i}(0) - \Gamma_{0i}) \hat{\beta} + \left( \hat{\beta} \right)^2 \Gamma_{1i}^* \end{aligned}$$

Using,

$$\hat{\Omega}^{-1} = \Omega^{-1} + \Omega^{-1} (\Omega - \hat{\Omega}) \Omega^{-1} + \Omega^{-1} (\Omega - \hat{\Omega}) \hat{\Omega}^{-1} (\Omega - \hat{\Omega}) \Omega^{-1}$$

and,

$$\begin{aligned} \hat{\Omega} - \Omega &= \hat{\Gamma}'_0 \hat{\Upsilon}^{-1} \hat{\Gamma}_0 - \hat{\Gamma}'_0 \Upsilon^{-1} \hat{\Gamma}_0 + \hat{\Gamma}'_0 \Upsilon^{-1} \hat{\Gamma}_0 - \Gamma'_0 \Upsilon^{-1} \Gamma_0 \\ \hat{\Gamma}'_0 \hat{\Upsilon}^{-1} \hat{\Gamma}_0 - \hat{\Gamma}'_0 \Upsilon^{-1} \hat{\Gamma}_0 &= \Gamma'_0 (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \Gamma_0 + (\hat{\Gamma}_0 - \Gamma_0)' (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \Gamma_0 \\ &\quad + \Gamma'_0 (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) (\hat{\Gamma}_0 - \Gamma_0) + (\hat{\Gamma}_0 - \Gamma_0)' (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) (\hat{\Gamma}_0 - \Gamma_0) \end{aligned}$$

with

$$\begin{aligned} \hat{\Gamma}_0 - \Gamma_0 &= \bar{\Gamma}_0 - \Gamma_0 + \hat{\beta} \bar{\Gamma}_1 + \hat{\beta}^2 \bar{\Gamma}_1^* \\ &= \bar{\Gamma}_0 - \Gamma_0 + \hat{\beta} \Gamma_1 + \hat{\beta} (\bar{\Gamma}_1 - \Gamma_1) + \hat{\beta}^2 \bar{\Gamma}_1^* \end{aligned}$$

we can write,

$$\begin{aligned} (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) &= \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} + \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \hat{\Upsilon}^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \\ \Gamma'_0 (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \Gamma_0 &= \Gamma'_0 \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \Gamma_0 \\ &\quad + \Gamma'_0 \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \hat{\Upsilon}^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \Gamma_0 \\ \Gamma'_0 \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \hat{\Upsilon}^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \Gamma_0 &= O(K/n) \\ \Gamma'_0 \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \Gamma_0 &= -\Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} \Gamma_0 \\ &\quad - \hat{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 + O(1/n) \\ \Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} \Gamma_0 &= O(1/\sqrt{n}) \\ \hat{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 &= O(1/\sqrt{n}) \\ (\hat{\Gamma}_0 - \Gamma_0)' (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \Gamma_0 &= O(\zeta(K) \sqrt{K}/n) \\ (\hat{\Gamma}_0 - \Gamma_0)' (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) (\hat{\Gamma}_0 - \Gamma_0) &= o(K/n) \\ \hat{\Gamma}'_0 \Upsilon^{-1} \hat{\Gamma}_0 - \Gamma'_0 \Upsilon^{-1} \Gamma_0 &= (\hat{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 + \Gamma'_0 \Upsilon^{-1} (\hat{\Gamma}_0 - \Gamma_0) \\ &\quad + (\hat{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} (\hat{\Gamma}_0 - \Gamma_0) \\ (\hat{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 &= (\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 + \hat{\beta} \Gamma_1 \Upsilon^{-1} \Gamma_0 + O(1/n) \\ (\hat{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} (\hat{\Gamma}_0 - \Gamma_0) &= O(K/n) \end{aligned}$$

so that,

$$\begin{aligned}\hat{\Omega} - \Omega &= \Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} \Gamma_0 + 2\hat{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 - (\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma_0 - \hat{\beta} \Gamma_1 \Upsilon^{-1} \Gamma_0 \\ &\quad - \Gamma'_0 \Upsilon^{-1} (\bar{\Gamma}_0 - \Gamma_0) - \hat{\beta} \Gamma'_0 \Upsilon^{-1} \Gamma_1 + O(\zeta(K)\sqrt{K}/n) \\ \|\hat{\Omega} - \Omega\| &= O(1/\sqrt{n})\end{aligned}$$

Now,

$$\sum_{i=1}^n \hat{\Gamma}'_{0i} \hat{\Sigma} \hat{g}_i / n^2 = \left(\frac{1}{n^2}\right) \sum_{i=1}^n \hat{\Gamma}'_{0i} \Sigma \hat{g}_i / n^2 + \left(\frac{1}{n^2}\right) \sum_{i=1}^n \hat{\Gamma}'_{0i} (\hat{\Sigma} - \Sigma) \hat{g}_i / n^2$$

and,

$$\begin{aligned}\sum_{i=1}^n \hat{\Gamma}'_{0i} \Sigma \hat{g}_i / n^2 &= \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 + \sum_{i=1}^n \Gamma'_{0i} \Sigma \Gamma_{0i} \hat{\beta} / n^2 + \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma g_i / n^2 \\ &\quad + \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma (\Gamma_{0i}(0) - \Gamma_{0i}) \hat{\beta} / n^2 \\ &\quad + \hat{\beta} \sum_{i=1}^n (\Gamma_{1i}(0) - \Gamma_{1i})' \Sigma g_i / n^2 + O(1/\sqrt{n})o(K/n)\end{aligned}$$

with,

$$\begin{aligned}\sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 &= O(1/\sqrt{n})O(\zeta(K)\sqrt{K}/n) \\ \sum_{i=1}^n \Gamma'_{0i} \Sigma \Gamma_{0i} \hat{\beta} / n^2 &= O(1/\sqrt{n})O(K/n) \\ \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma g_i / n^2 &= O(1/\sqrt{n})O(K/\sqrt{n}) \\ \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma (\Gamma_{0i}(0) - \Gamma_{0i}) \hat{\beta} / n^2 &= O(1/\sqrt{n})O(K/n) \\ \hat{\beta} \sum_{i=1}^n (\Gamma_{1i}(0) - \Gamma_{1i})' \Sigma g_i / n^2 &= O(1/\sqrt{n})O(K/n)\end{aligned}$$

Next

$$\begin{aligned}
& \sum_{i=1}^n \hat{\Gamma}'_{0i} (\hat{\Sigma} - \Sigma) \hat{g}_i / n^2 \\
= & \sum_{i=1}^n \hat{\Gamma}'_{0i} (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \hat{g}_i / n^2 + \sum_{i=1}^n \hat{\Gamma}'_{0i} \left( \Upsilon^{-1} \Gamma_0 \hat{\Omega}^{-1} \hat{\Gamma}'_0 \hat{\Upsilon}^{-1} - \hat{\Upsilon}^{-1} \hat{\Gamma}_0 \Omega^{-1} \hat{\Gamma}'_0 \hat{\Upsilon}^{-1} \right) \hat{g}_i / n^2 \\
& + \sum_{i=1}^n \hat{\Gamma}'_{0i} \left( \hat{\Upsilon}^{-1} \hat{\Gamma}_0 \Omega^{-1} \hat{\Gamma}'_0 \hat{\Upsilon}^{-1} - \Upsilon^{-1} \Gamma_0 \Omega^{-1} \Gamma'_0 \Upsilon^{-1} \right) \hat{g}_i / n^2
\end{aligned}$$

For the first term using the above expansions for  $\hat{\Upsilon}^{-1} - \Upsilon^{-1}$  and for  $\Upsilon - \hat{\Upsilon}$  in Proposition 1, we can show that using T, followed by CS and then ME and the conditions on  $K$ , for BGMM we have

$$\begin{aligned}
\sum_{i=1}^n \hat{\Gamma}'_{0i} \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \hat{g}_i / n^2 &= - \sum_{i=1}^n \Gamma'_{0i} \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 - \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 \\
&\quad - 2\hat{\beta} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} g_i / n^2 + O(1/\sqrt{n})o(K/n)
\end{aligned}$$

with,

$$\begin{aligned}
\left\| \sum_{i=1}^n \Gamma'_{0i} \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 \right\| &= (1/\sqrt{n})O(\zeta(K)K/n) \\
\left\| \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 \right\| &= (1/\sqrt{n})O(\zeta(K)K/n) \\
\left\| 2\hat{\beta} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} g_i / n^2 \right\| &= O(1/\sqrt{n})O(K/n)
\end{aligned}$$

where the second part of the first term satisfies (using similar arguments),

$$\left\| \sum_{i=1}^n \hat{\Gamma}'_{0i} \Upsilon^{-1} (\Upsilon - \hat{\Upsilon}) \hat{\Upsilon}^{-1} (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \hat{g}_i / n^2 \right\| = O(1/\sqrt{n})o(K/n)$$

For the remaining terms we can show using similar arguments that they are each  $O(1/\sqrt{n})o(K/n)$ .

So altogether  $\sum_{i=1}^n \hat{\Gamma}'_{0i} \hat{\Sigma} \hat{g}_i / n^2 = O(K/\sqrt{n})$  so that by the condition on  $K$  and,

$$O(K/\sqrt{n}) \left( O(\zeta(K)\sqrt{K}/n) + O(\Delta_{K,n}) \right) = O(\gamma_{K,n})$$



we have that

$$\begin{aligned}
& \hat{\Omega}^{-1} \sum_{i=1}^n \hat{\Gamma}'_{0i} \hat{\Sigma} \hat{g}_i / n^2 = \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 + \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma \Gamma_{0i} \hat{\beta} / n^2 \\
& + \Omega^{-1} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma g_i / n^2 + \Omega^{-1} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma (\Gamma_{0i}(0) - \Gamma_{0i}) \hat{\beta} / n^2 \\
& + \Omega^{-1} \hat{\beta} \sum_{i=1}^n (\Gamma_{1i}(0) - \Gamma_{1i})' \Sigma g_i / n^2 - \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 \\
& \quad - \Omega^{-1} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 \\
& \quad - 2\hat{\beta} \Omega^{-1} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} g_i / n^2 \\
& \quad + \Omega^{-1} \Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} \Gamma_0 \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 \\
& \quad + 2\hat{\beta} \Omega^{-1} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 \\
& \quad - \Omega^{-1} \left( \Gamma'_0 \Upsilon^{-1} (\bar{\Gamma}_0 - \Gamma_0) + (\bar{\Gamma}_0 - \Gamma_0)' \Upsilon^{-1} \Gamma \right) \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 \\
& \quad - \hat{\beta} \Omega^{-1} (\Gamma'_0 \Upsilon^{-1} \Gamma_1 + \Gamma_1 \Upsilon^{-1} \Gamma) \Omega^{-1} \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 + O(1/\sqrt{n})o(K/n)
\end{aligned}$$

Now combining terms here with the terms that are not  $O(1/\sqrt{n})o(K/n)$  from GMM and using the facts that for  $j = 0, 1$ ,

$$- (\bar{\Gamma}_j - \Gamma_j)' \Sigma \bar{g} + \sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 = \sum_{i \neq j}^n \Gamma'_{0i} \Sigma g_j / n^2 = O(1/\sqrt{n})O(\sqrt{K/n})$$

and the fact that  $\hat{\beta} = -\Omega^{-1} \Gamma'_0 \Upsilon^{-1} \bar{g} + o(1/\sqrt{n})$  have that for instance,

$$(\bar{\Gamma}_1 - \Gamma_1)' \Sigma \bar{g} \Omega^{-1} \Gamma'_0 \Upsilon^{-1} \bar{g} + \hat{\beta} \sum_{i=1}^n (\Gamma_{1i}(0) - \Gamma_{1i})' \Sigma g_i / n^2 = O(1/\sqrt{n})O(\sqrt{K/n})$$

This also occurs by combining terms from GMM and BGMM. Then for BGMM we have

that for  $\gamma_{K,n} = K/n + \Delta_{K,n}$ ,

$$\begin{aligned}\sqrt{n}\hat{\beta}^B &= \Omega^{-1}(h + \sum_{j=1}^5 T_j^B + Z^B) \\ \|T_1^B\| &= O(\sqrt{K/n}), \|T_2^B\| = O(\zeta(K)K/n), \|T_3^B\| = O(1/\sqrt{n}) \\ \|T_4^B\| &= O(\zeta(K)\sqrt{K/n}) \quad \|T_5^B\| = O(K/n), \|Z^B\| = o(\gamma_{K,n})\end{aligned}$$

then under the condition on  $K$  as with GMM,

$$n(\hat{\beta}^B)^2 = n\Omega^{-1}hh'\Omega^{-1} + \Omega^{*-1}(T_1^B T_1^B + \sum_{j=1}^5 2T_j^B h)\Omega^{*-1} + o(\gamma_{K,n})$$

Now doing the calculations we get  $E(hT_3^B) = O$  by Assumption 1(iv). Then we have for  $E(T_1^B T_1^B)$  the terms,

$$\begin{aligned}nE\left(\sum_{i \neq j}^n \Gamma'_{0i} \Sigma g_j / n^2\right) &= \frac{1}{n} \sum_i E(\eta_{0i}^2) \xi_{ii} + \sum_{i,j} E(\rho_i \eta_{0i}) E(\rho_j \eta_{0j}) \xi_{ij}^2 + o(\gamma_{K,n}) \\ nE(\Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \bar{g} \sum_{i \neq j}^n \Gamma'_{0i} \Sigma g_j / n^2) &= o(\gamma_{K,n}) \\ nE(\Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \bar{g} \bar{g}' \Sigma \hat{\Upsilon}_{v1} \Upsilon^{-1} \Gamma) &= \frac{1}{n} \sum_i (\kappa_i - 1) d_i^2 \xi_{ii} + o(\gamma_{K,n})\end{aligned}$$

For the term  $2T_1^B h$  we have by Assumption 1(iv)

$$2nE(\Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \bar{g} h) = -2 \frac{1}{n} \sum_i (\kappa_i - 1) d_i^2 \xi_{ii} + o(\gamma_{K,n})$$

Next for  $2T_2^B h$  we have,

$$\begin{aligned}-2nE(\Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} h) &= o(\gamma_{K,n}) \\ -2nE((\bar{\Gamma}_0 - \Gamma_0)' \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} h') &= o(\gamma_{K,n}) \\ 4nE(\tilde{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Sigma \hat{\Upsilon}_{v1} \Sigma \bar{g} h') &= o(\gamma_{K,n}) \\ -2nE\left(\sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 h'\right) &= o(\gamma_{K,n}) \\ -2nE\left(\sum_{i=1}^n \Gamma'_{0i} \Upsilon^{-1} \hat{\Upsilon}_{v1} \Upsilon^{-1} g_i / n^2 h'\right) &= o(\gamma_{K,n})\end{aligned}$$

For the terms in  $2T_4^B h$  we have,

$$2nE\left(\sum_{i=1}^n \Gamma'_{0i} \Sigma g_i / n^2 h'\right) = -2\frac{1}{n} \sum_i d_i^2 \xi_{ii} + o(\gamma_{K,n})$$

Next we must deal with the many terms in  $2T_4^B h$ . We have for the terms coming from Proposition 1,

$$\begin{aligned} 4nE(\tilde{\beta} \Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v2} \Sigma \bar{g} h') &= 4\frac{1}{n} \sum_i d_i^2 \xi_{ii} + o(\gamma_{K,n}) \\ -4nE(\tilde{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} h') &= -4\Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 \frac{1}{n} \sum_i E(\rho_i \eta_{0i}) \xi_{ii} \\ &\quad + o(\gamma_{K,n}) \\ -4nE(\tilde{\beta} \Gamma'_0 \Upsilon^{-1} \hat{\Upsilon}_{v1} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} h') &= o(\gamma_{K,n}) \\ 2nE((\bar{\Gamma}_0 - \Gamma_0)' \Sigma (\bar{\Gamma}_0 - \Gamma_0) \Omega^{-1} \Gamma'_0 \Upsilon^{-1} \bar{g} h') &= -2\frac{1}{n} \sum_i E(\eta_{0i}^2) \xi_{ii} + o(\gamma_{K,n}) \\ 4nE(\tilde{\beta} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} h') &= \frac{1}{n} \sum_{i,j} E(\rho_i \eta_{0i}) E(\rho_j \eta_j) \xi_{ij}^2 + o(\gamma_{K,n}) \\ 2nE(\Gamma'_1 \Sigma \bar{g} \Omega^{-1} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \bar{g} h') &= o(\gamma_{K,n}) \end{aligned}$$

and the terms coming from the expansion of the bias adjustment factor,

$$\begin{aligned} 2nE\left(\sum_{i=1}^n \Gamma'_{0i} \Sigma \Gamma_{0i} \hat{\beta} / n^2 h'\right) &= 2\frac{1}{n} \sum_i d_i^2 \xi_{ii} \\ 2nE\left(\sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i}) \Sigma (\Gamma_{0i}(0) - \Gamma_{0i}) \hat{\beta} / n^2 h'\right) &= 2\frac{1}{n} \sum_i E(\eta_i^2) \xi_{ii} + o(\gamma_{K,n}) \\ -4nE\left(\hat{\beta} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} g_i / n^2 h'\right) &= -4\frac{1}{n} \sum_{i,j} E(\rho_i \eta_{0i}) E(\rho_j \eta_j) \xi_{ij}^2 + o(\gamma_{K,n}) \end{aligned}$$

$$\begin{aligned} &4nE(\hat{\beta} \Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 \Omega^{-1} \sum_{i=1}^n (\Gamma_{0i}(0) - \Gamma_{0i})' \Sigma g_i / n^2 h) \\ &= 4\Gamma'_0 \Upsilon^{-1} \Upsilon_{\rho\eta} \Upsilon^{-1} \Gamma_0 \frac{1}{n} \sum_i E(\rho_i \eta_{0i}) \xi_{ii} + o(\gamma_{K,n}) \end{aligned}$$

Then collect all terms use the definition of  $\tau$  to get the result.

**Proof of Proposition 3:** For ease of notation assume that  $\hat{\beta}$  the GEL estimator has population zero, so that both  $\hat{\beta}$  and  $\hat{\lambda}$  have population value zero. Also note that under the conditions of the proposition  $\gamma_{K,n} \geq C(K/n + \Delta_{K,n})$  so that terms in the expansion can be dropped if they are  $o(K/n)$ . Then the FOC given by,  $\sum_i m_i(\hat{\theta})/n = 0$  where,

$$m_i(\theta) = -s_1(\lambda' g_i(\beta)) \begin{pmatrix} \Gamma_{0i}(\beta)' \lambda \\ g_i(\beta) \end{pmatrix}$$

and we have by Appendix B that  $M$  and  $M^{-1}$  are the same as for GMM. Using Appendix B we have

$$\begin{pmatrix} T_1^\beta \\ T_1^\lambda \end{pmatrix} - M^{-1}m = \begin{pmatrix} -\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} \\ -\Sigma\bar{g} \end{pmatrix} = \begin{pmatrix} O(1/\sqrt{n}) \\ O(\sqrt{K/n}) \end{pmatrix}$$

For  $-M^{-1}(\hat{M} - M)\hat{\theta}$  we have terms,

$$\begin{aligned} T_2^\beta &= \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} - \Omega^{-1}\Gamma_0'\Upsilon^{-1}\tilde{\Upsilon}_v\hat{\lambda} - \Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\hat{\beta} \\ &= O(K/n) + O(\sqrt{K}/n) + O(1/n) \\ T_2^\lambda &= -\Upsilon^{-1}\Gamma_0\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} - \Sigma\left((\bar{\Gamma}_0 - \Gamma_0)\hat{\beta} + \tilde{\Upsilon}_v\hat{\lambda}\right) \\ &= O(K/n) + O(\sqrt{K}/n) + O(\zeta(K)\sqrt{K}/n) \end{aligned}$$

using CS and arguments similar to Proposition 1. For the term  $-(1/2)M^{-1}\sum_j \hat{\theta}_j A_j \hat{\theta}$  we have similarly,

$$\begin{aligned} T_3^\beta &= \Omega^{-1}\left(\hat{\beta}\Gamma_1'\hat{\lambda} + \hat{\lambda}'\Upsilon_{\rho\eta}\hat{\lambda}\right) - \Omega^{-1}\Gamma_0'\Upsilon^{-1}\left((1/2)\hat{\beta}^2\Gamma_1 + 2\hat{\beta}\Upsilon_{\rho\eta}\hat{\lambda}\right) \\ &= O(1/n) + O(K/n) + O(1/n) + O(1/n) \\ T_3^\lambda &= -\Upsilon^{-1}\Gamma_0\Omega^{-1}\left(\hat{\beta}\Gamma_1'\hat{\lambda} + \hat{\lambda}'\Upsilon_{\rho\eta}\hat{\lambda}\right) - \Sigma\left((1/2)\hat{\beta}^2\Gamma_1 + 2\hat{\beta}\Upsilon_{\rho\eta}\hat{\lambda}\right) \\ &= O(1/n) + O(K/n) + O(1/n) + O(\sqrt{K}/n) \end{aligned}$$

and from  $-M^{-1} \sum_j \hat{\theta}_j (\hat{A}_j - A_j) \hat{\theta}/2$  we have terms,

$$\begin{aligned}
T_4^\beta &= \Omega^{-1} \hat{\beta} (\bar{\Gamma}_1 - \Gamma_1)' \hat{\lambda} + \Omega^{-1} \hat{\lambda}' (\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3}) \hat{\lambda} - (1/2) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\beta}^2 (\bar{\Gamma}_1 - \Gamma_1) \\
&\quad - 2\Omega^{-1} \Gamma_0' \Upsilon^{-1} \hat{\beta} (\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3}) \hat{\lambda} + (s_3/2) \Omega^{-1} \Gamma_0' \Upsilon^{-1} \sum_j \hat{\lambda}_j \hat{\Upsilon}_{v4}^j \hat{\lambda} \\
&= O(K/n^{3/2}) + O(K/n) O(\zeta(K) \sqrt{K/n}) + O(1/n^{3/2}) + O(K/n^{3/2}) \\
&\quad + O(K/n) O(\sqrt{K} \zeta(K) / \sqrt{n}) \\
T_4^\lambda &= -\Upsilon^{-1} \Gamma_0 \Omega^{-1} (3/2) \hat{\beta} (\bar{\Gamma}_1 - \Gamma_1)' \hat{\lambda} - +\Upsilon^{-1} \Gamma_0 \Omega^{-1} \hat{\lambda}' (\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3}) \hat{\lambda} - \Sigma (1/2) \hat{\beta}^2 (\bar{\Gamma}_1 - \Gamma_1) \\
&\quad - \Sigma \hat{\beta} (\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3}) \hat{\lambda} - (s_3/2) \Sigma \sum_j \hat{\lambda}_j \hat{\Upsilon}_{v4}^j \hat{\lambda} \\
&= O(K/n^{3/2}) + O(K/n) O(\zeta(K) \sqrt{K/n}) + O(\sqrt{K}/n^{3/2}) + O(\zeta(K) \sqrt{K}/n^{3/2}) \\
&\quad + O(\zeta(K) K/n^{3/2}) + O(K/n) O(\zeta(K)^2 \sqrt{K/n}) \\
&= o(1/n)
\end{aligned}$$

where we use the notation  $\hat{\Upsilon}_{v4}^j = \sum_i q_i q_i \rho_i^3 q_i^j / n$ . Note that for  $T_4^\beta$  the order of the last term follows from,

$$\begin{aligned}
\left\| \frac{1}{n} \sum_j \hat{\lambda}_j \Gamma_0' \Upsilon^{-1} \hat{\Upsilon}_{v4}^j \hat{\lambda} \right\| &\leq \|\hat{\lambda}\| \left\| \hat{\lambda}' \frac{1}{n} \sum_i \bar{D}_i q_i q_i' \rho_i^3 \right\| \leq \|\hat{\lambda}\|^2 \left\| \frac{1}{n} \sum_i \bar{D}_i q_i q_i' \rho_i^3 \right\| \\
&\leq O(K/n) O(\sqrt{K} \zeta(K) / \sqrt{n})
\end{aligned}$$

by,

$$E\left( \left\| \frac{1}{n} \sum_i \bar{D}_i q_i q_i' \rho_i^3 \right\|^2 \right) \leq \frac{1}{n^2} \sum_i \bar{D}_i^2 E(\rho_i^6) \|q_i\|^4 = O(K \zeta(K)^2 / n)$$

while for  $T_4^\lambda$  it follows from use of T, CS, the result for  $\|\hat{\lambda}\|$  and the fact that,

$$\begin{aligned}
\sum_j E\left( \left\| \hat{\Upsilon}_{v4}^j \right\|^2 \right) &\leq \frac{1}{n^2} \lambda_{\max}(\Upsilon) \sum_i \|q_i\|^2 E(\rho_i^6) \sum_j (q_i^j)^2 q_i' \Upsilon^{-1} q_i \\
&\leq \frac{1}{n} C \lambda_{\max}(\Upsilon) \sum_i \|q_i\|^2 (q_i' \Upsilon^{-1} q_i / n) \|q_i\|^2 = O(K \zeta(K)^4 / n)
\end{aligned}$$

Next for  $-M^{-1} \sum_j \sum_k \hat{\theta}_j \hat{\theta}_k B_{jk} \hat{\theta} / 6$  we have terms,

$$\begin{aligned}
T_5^\beta &= \hat{\beta} \Omega^{-1} \hat{\lambda}' \Upsilon_{\eta\eta+dd} \hat{\lambda} + \hat{\beta} \Omega^{-1} \hat{\lambda}' \Upsilon_{\rho 1} \hat{\lambda} - (1/2) s_3 \Omega^{-1} \hat{\lambda} \sum_j \hat{\lambda}_j \Upsilon_{2d}^j \hat{\lambda} \\
&\quad + (3/2) \hat{\beta} s_3 \Omega^{-1} \Gamma_0' \Upsilon^{-1} \sum_j \hat{\lambda}_j \Upsilon_{2d}^j \hat{\lambda} - (1/6) s_4 \Omega^{-1} \Gamma_0' \Upsilon^{-1} \sum_{j,k} \Upsilon_4^{jk} \hat{\lambda}_j \hat{\lambda}_k \hat{\lambda} \\
&\quad - (1/6) \hat{\beta}^3 \Omega^{-1} \Gamma_0' \Upsilon^{-1} \Gamma_2 \hat{\beta} - \hat{\beta}^2 \Omega^{-1} \Gamma_0' \Upsilon^{-1} \Upsilon_{\eta\eta+dd} \hat{\lambda} - \hat{\beta}^2 \Omega^{-1} \Gamma_0' \Upsilon^{-1} \Upsilon_{\rho 1} \hat{\lambda} + \hat{\beta}^2 \Omega^{-1} \Gamma_2' \hat{\lambda} \\
\Upsilon_{\eta\eta+dd} &= \sum_i q_i q_i (E(\eta_i^2 | x_i) + d_i^2), \quad \Upsilon_{\rho 1} = \sum_i q_i q_i E(\rho_i \eta_{1i} | x_i) / n, \\
\Upsilon_{2d}^j &= \sum_i q_i q_i q_i^j d_i E(\rho_i^2 | x_i) / n, \quad \Upsilon_4^{jk} = \sum_i q_i q_i q_i^j q_i^k E(\rho_i^4 | x_i) / n
\end{aligned}$$

Here we can show that the last three terms are all  $o(K/n^{3/2})$ . The first two terms are  $O(K/n^{3/2})$  using CS and M, while for the remaining terms we can show that since,

$$\left\| \sum_j \hat{\lambda}_j \Upsilon_{2d}^j \hat{\lambda} \right\| = O(\|\hat{\lambda}\|^2) O(\zeta(K)K)$$

we can show using CS that each term is  $O(1/\sqrt{n})o(\sqrt{K/n})$  given the condition on  $K$ . Similarly it is possible to show that  $\|T_5^\lambda\| = o(\sqrt{K}/n)$ . The remaining terms in the expansion are dealt with in a technical Appendix that is available on request – there we show that the term in the remainder corresponding to  $\hat{\beta}$  is  $o(K/n^{3/2})$  and that the term for  $\hat{\lambda}$  is  $o(\sqrt{K}/n)$  as will be required in the remainder of the proof. The technical appendix also contains all required derivatives.

Before collecting terms we note that by repeated substitution and using the preliminary bounds for the terms in the expansion for  $\hat{\lambda}$  we can write,

$$\begin{aligned}
\hat{\lambda}' \Upsilon_{\rho\eta} \hat{\lambda} &= \bar{g}' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} - 2\bar{g}' \Sigma \tilde{\Upsilon}_{v1} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} - 4\hat{\beta} \bar{g}' \Sigma \Upsilon_{\rho\eta} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} \\
&\quad - 2\hat{\beta} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} + O(1/\sqrt{n})o(K/n) \\
\|\bar{g}' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g}\| &= O(K/n), \quad \left\| 2\bar{g}' \Sigma \tilde{\Upsilon}_{v1} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} \right\| = O(\zeta(K)K/n^{3/2}) \\
\left\| 4\hat{\beta} \bar{g}' \Sigma \Upsilon_{\rho\eta} \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} \right\| &= O(K/n^{3/2}), \quad \left\| 2\hat{\beta} (\bar{\Gamma}_0 - \Gamma_0)' \Sigma \Upsilon_{\rho\eta} \Sigma \bar{g} \right\| = O(K/n^{3/2})
\end{aligned}$$

Then we can derive an expansion for  $\hat{\lambda}$ ,

$$\begin{aligned}\hat{\lambda} &= -\Sigma\bar{g} + \Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g} + \Upsilon^{-1}\Gamma_0\Omega^{-1} \left( (\bar{\Gamma}_0 - \Gamma_0)' \Sigma\bar{g} - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} \right) + 2\Sigma\hat{\beta}\Upsilon_{\rho\eta}\Sigma\bar{g} + \Sigma R_\lambda^3 \\ &\quad \left\| R_\lambda^3 \right\| = o(\sqrt{K}/n) \\ &\quad \Upsilon^{-1}\Gamma_0\Omega^{-1} \left( (\bar{\Gamma}_0 - \Gamma_0)' \Sigma\bar{g} - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} \right) = O(\sqrt{K}/n) \\ &\quad 2\Sigma\hat{\beta}\Upsilon_{\rho\eta}\Sigma\bar{g} = O(\sqrt{K}/n)\end{aligned}$$

Now take the first  $\hat{\beta}$  term in  $T_2^\beta$  and substitute for  $\hat{\lambda}$  and combine with the second part of  $T_3^\beta$  and substitute in for  $\hat{\beta} = -\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g} + O(1/\sqrt{n})$  to get,

$$\begin{aligned}&\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \hat{\lambda} + \Omega^{-1}\hat{\lambda}'\Upsilon_{\rho\eta}\hat{\lambda} \\ &= -\Omega^{-1}((\bar{\Gamma}_0 - \Gamma_0)' \Sigma\bar{g} - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}) - \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma(\bar{\Gamma}_0 - \Gamma_0) (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\ &\quad + \Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g} - 4\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\ &\quad - 2\Omega^{-1}\bar{g}'\Sigma\tilde{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} + 4(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \Omega^{-1}\bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} \\ &\quad + O(1/\sqrt{n})o(K/n)\end{aligned}$$

where in addition to the terms that appeared in the expansion for GMM we have,

$$\begin{aligned}\left\| \Omega^{-1}((\bar{\Gamma}_0 - \Gamma_0)' \Sigma\bar{g} - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}) \right\| &= O(\sqrt{K}/n) \\ \left\| 2\Omega^{-1}(\bar{\Gamma}_0 - \Gamma_0)' \Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \right\| &= O(K/n^{3/2})\end{aligned}$$

with the other terms having the same order as described earlier. Hence we have that,

$$\hat{\beta} = -(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) + O(\sqrt{K}/n)$$

and so for the second term in  $T_2^\beta$ ,

$$\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0)\hat{\beta} = \Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_0 - \Gamma_0) (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) + O(1/\sqrt{n})o(K/n)$$

For the third term in  $T_2^\beta$  we have,

$$\begin{aligned}-\Omega^{-1}\Gamma_0'\Upsilon^{-1}\tilde{\Upsilon}_v\hat{\lambda} &= \Omega^{-1}\Gamma_0'\Upsilon^{-1}\tilde{\Upsilon}_v\Sigma\bar{g} - \Omega^{-1}\Gamma_0'\Upsilon^{-1}\tilde{\Upsilon}_{v1}\Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g} \\ &\quad + 2(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \Omega^{-1}\Gamma_0'\Upsilon^{-1}\tilde{\Upsilon}_v\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g} + O(1/\sqrt{n})o(K/n)\end{aligned}$$

The additional terms for  $T_3^\beta$  are (multiplied by  $\Omega^{-1}$ ),

$$\begin{aligned}
\hat{\beta}\Gamma_1'\hat{\lambda} &= -\hat{\beta}\Gamma_1'\Sigma\bar{g} + \hat{\beta}\Gamma_1'\Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g} + O(1/\sqrt{n})O(\sqrt{K}/n) \\
&= -\hat{\beta}\Gamma_1'\Sigma\bar{g} + O(1/\sqrt{n})O(\sqrt{K}/n) \\
&= (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g})\Gamma_1'\Sigma\bar{g} + O(1/\sqrt{n})O(\sqrt{K}/n) \\
-(1/2)\Gamma_0'\Upsilon^{-1}\Gamma_1\hat{\beta}^2 &= -(1/2)\Gamma_0'\Upsilon^{-1}\Gamma_1(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g})^2 + O(1/\sqrt{n})O(\sqrt{K}/n) \\
-2\Gamma_0'\Upsilon^{-1}\hat{\beta}\Upsilon_{\rho\eta}\hat{\lambda} &= 2\Gamma_0'\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\bar{g}(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) - 2\Gamma_0'\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g}(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\
&\quad + O(1/\sqrt{n})O(\sqrt{K}/n)
\end{aligned}$$

under the condition that  $\zeta(K)^2K^2/\sqrt{n} \rightarrow 0$ . From  $T_4^\beta$  we have,

$$\begin{aligned}
\Omega^{-1}\hat{\beta}(\bar{\Gamma}_1 - \Gamma_1)'\hat{\lambda} &= \Omega^{-1}(\bar{\Gamma}_1 - \Gamma_1)'\Sigma\bar{g}(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\
&\quad + O(1/\sqrt{n})o(K/n) \\
\Omega^{-1}\hat{\lambda}'(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\hat{\lambda} &= \Omega^{-1}\bar{g}'\Sigma(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\Sigma\bar{g} + O(1/\sqrt{n})o(1/n) \\
-(1/2)\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\bar{\Gamma}_1 - \Gamma_1)\hat{\beta}^2 &= O(1/\sqrt{n})o(1/n) \\
-2\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\hat{\lambda}\hat{\beta} &= -2\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\Sigma\bar{g}(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\
&\quad + O(1/\sqrt{n})o(K/n) \\
-\Omega^{-1}\Gamma_0'\Upsilon^{-1}(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\Sigma\bar{g}(\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) &= O(K/n^{3/2}) \\
\Omega^{-1}\Gamma_0'\Upsilon^{-1}(s_3/2)\sum_j\hat{\lambda}_j\hat{\Upsilon}_{v4}^j\hat{\lambda} &= (s_3/2)\Omega^{-1}\bar{g}'\Sigma\left(\frac{1}{n}\sum_i\bar{D}_iq_iq_i'\rho_i^3\right)\Sigma\bar{g} \\
&\quad + O(1/\sqrt{n})o(K/n) \\
(s_3/2)\Omega^{-1}\bar{g}'\Sigma\left(\frac{1}{n}\sum_i\bar{D}_iq_iq_i'\rho_i^3\right)\Sigma\bar{g} &= O(K\zeta(K)^2/n^{3/2})
\end{aligned}$$



Finally, from  $T_5^\beta$  we get terms,  $\hat{\beta}\Omega^{-1}\hat{\lambda}'\Upsilon_{\eta\eta+dd}\hat{\lambda} + \hat{\beta}\Omega^{-1}\hat{\lambda}'\Upsilon_{\rho 1}\hat{\lambda} - (1/2)s_3\Omega^{-1}\hat{\lambda}'\sum_j\hat{\lambda}_j\Upsilon_{2d}^j\hat{\lambda}$

$$\begin{aligned}
\hat{\beta}\Omega^{-1}\hat{\lambda}'\Upsilon_{\eta\eta+dd}\hat{\lambda} &= -\Omega^{-1}\bar{g}'\Sigma\Upsilon_{\eta\eta+dd}\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) + O(1/\sqrt{n})o(K/n) \\
\hat{\beta}\Omega^{-1}\hat{\lambda}'\Upsilon_{\rho 1}\hat{\lambda} &= -\Omega^{-1}\bar{g}'\Sigma\Upsilon_{\rho 1}\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) + O(1/\sqrt{n})o(K/n) \\
-(1/2)s_3\Omega^{-1}\hat{\lambda}'\sum_j\hat{\lambda}_j\Upsilon_{2d}^j\hat{\lambda} &= (1/2)s_3\Omega^{-1}\bar{g}'\Sigma\sum_j(e'_j\Sigma\bar{g})\Upsilon_{2d}^j\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\
&\quad + O(1/\sqrt{n})o(K/n) \\
(3/2)\hat{\beta}s_3\Omega^{-1}\Gamma_0'\Upsilon^{-1}\sum_j\hat{\lambda}_j\Upsilon_{2d}^j\hat{\lambda} &= -(3/2)s_3\Omega^{-1}\Gamma_0'\Upsilon^{-1}\sum_j(e'_j\Sigma\bar{g})\Upsilon_{2d}^j\Sigma\bar{g} (\Omega^{-1}\Gamma_0'\Upsilon^{-1}\bar{g}) \\
&\quad + O(1/\sqrt{n})o(K/n) \\
-(1/6)s_4\Omega^{-1}\Gamma_0'\Upsilon^{-1}\sum_{j,k}\Upsilon_4^{jk}\hat{\lambda}_j\hat{\lambda}_k\hat{\lambda} &= (1/6)s_4\Omega^{-1}\Gamma_0'\Upsilon^{-1}\sum_{j,k}(e'_j\Sigma\bar{g})(e'_k\Sigma\bar{g})\Upsilon_4^{jk}(e'_j\Sigma\bar{g})\Sigma\bar{g} \\
&\quad + O(1/\sqrt{n})o(K/n)
\end{aligned}$$

with the order of each term on the right the same as the order of the term on the left as given above.

Then we can write,

$$\begin{aligned}
\sqrt{n}\hat{\beta}^{GEL} &= \Omega^{-1}(h + \sum_{j=1}^3 T_j^G + Z^G), \quad \|T_1^G\| = O(\sqrt{K/n}), \\
\|T_2^G\| &= O(1/\sqrt{n}), \quad \|Z^B\| = o(\gamma_{K,n})
\end{aligned}$$

and where the terms in  $T_3^G$  are all  $o(\sqrt{K}/n)$ , but are not necessarily  $o(\gamma_{K,n})$ . Given this it is then the case that  $T_3^G$  is of an order

$$n\left(\hat{\beta}^G\right)^2 = n\Omega^{-1}hh'\Omega^{-1} + \Omega^{*-1}(T_1^B T_1^B + \sum_{j=1}^3 2T_j^B h)\Omega^{*-1} + o(\gamma_{K,n})$$

As in the case of BGMM and GMM we have that,

$$n\Omega^{-1}E(hh')\Omega^{-1} = \Omega^{*-1} + \Omega^{*-1}(\Omega^* - \Omega)\Omega^{*-1} + o(\gamma_{K,n})$$

Then in place of the first term in  $E(T_1^B T_1^B)$  calculated in Proposition 2 we have,

$$\begin{aligned}
&nE((\bar{\Gamma}_0 - \Gamma_0)' \Sigma\bar{g} - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}) (\bar{g}'\Sigma(\bar{\Gamma}_0 - \Gamma_0) - \bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}) \\
&= \frac{1}{n}\sum_i E(\eta_i^2)\xi_{ii} - \frac{1}{n}\sum_{i,j} E(\rho_i\eta_{0i})E(\rho_j\eta_j)\xi_{ij}^2 + o(\gamma_{K,n})
\end{aligned}$$

Next for

$$\begin{aligned}
-4nE(\bar{g}'\Sigma\tilde{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}h') &= o(\gamma_{K,n}) \\
8nE((\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})\bar{g}'\Sigma\Upsilon_{\rho\eta}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}h') &= -8\frac{1}{n}\sum_{i,j}E(\rho_i\eta_{0i})E(\rho_j\eta_j)\xi_{ij}^2 + o(\gamma_{K,n}) \\
4nE(\Gamma'_0\Upsilon^{-1}\tilde{\Upsilon}_{v1}\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g}h') &= o(\gamma_{K,n}) \\
4nE(\Gamma'_0\Upsilon^{-1}\Upsilon_{\rho\eta}\Sigma\tilde{\Upsilon}_{v1}\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= o(\gamma_{K,n}) \\
2nE((\bar{\Gamma}_1 - \Gamma_1)'\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= -2\frac{1}{n}\sum_i E(\rho_i\eta_{1i})\xi_{ii} + o(\gamma_{K,n}) \\
-2nE(\bar{g}'\Sigma\Upsilon_{\rho 1}\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= 2\frac{1}{n}\sum_i E(\rho_i\eta_{1i})\xi_{ii} + o(\gamma_{K,n}) \\
8nE((\bar{\Gamma}_0 - \Gamma_0)'\Sigma\Upsilon_{\rho\eta}\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= 8\frac{1}{n}\sum_{i,j}E(\rho_i\eta_{0i})E(\rho_j\eta_j)\xi_{ij}^2 + o(\gamma_{K,n})
\end{aligned}$$

Then we have,

$$\begin{aligned}
2nE(\bar{g}'\Sigma(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\Sigma\bar{g}h') &= -2\frac{1}{n}\sum_i d_i^2\xi_{ii} + o(\gamma_{K,n}) \\
-4E(\Gamma'_0\Upsilon^{-1}(\hat{\Upsilon}_{v2} + \hat{\Upsilon}_{v3})\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= 4\frac{1}{n}\sum_i d_i^2\xi_{ii} + o(\gamma_{K,n}) \\
-2E(\bar{g}'\Sigma\Upsilon_{\eta\eta+dd}\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= 2\frac{1}{n}\sum_i E(\eta_{0i}^2)\xi_{ii} + 2\frac{1}{n}\sum_i d_i^2\xi_{ii} + o(\gamma_{K,n})
\end{aligned}$$

Finally for the terms that depend on  $s_3$  we have,

$$\begin{aligned}
-2E((3/2)s_3\Omega^{-1}\Gamma'_0\Upsilon^{-1}\sum_j(e'_j\Sigma\bar{g})\Upsilon_{2d}^j\Sigma\bar{g}(\Omega^{-1}\Gamma'_0\Upsilon^{-1}\bar{g})h') &= 3s_3\frac{1}{n}\sum_{i,j}d_i^2\xi_{ii} + o(\gamma_{K,n}) \\
2E((s_3/2)\Omega^{-1}\bar{g}'\Sigma\left(\frac{1}{n}\sum_i\bar{D}_iq_iq_i^3\rho_i^3\right)\Sigma\bar{g}h') &= -s_3\frac{1}{n}\sum_{i,j}\kappa_id_i^2\xi_{ii} + o(\gamma_{K,n})
\end{aligned}$$

Collecting terms with those already found in Proposition 2, we have

$$\begin{aligned}
E(T_1^G T_1^G + \sum_{j=1}^3 2T_j^G h) &= \frac{1}{n}\sum_i E(\eta_i^2)\xi_{ii} - \frac{1}{n}\sum_{i,j} E(\rho_i\eta_{0i})E(\rho_j\eta_j)\xi_{ij}^2 \\
&\quad + 5\frac{1}{n}\sum_i d_i^2\xi_{ii} - \frac{1}{n}\sum_i \kappa_id_i^2\xi_{ii} + s_3\frac{1}{n}\sum_{i,j}(3 - \kappa_i)d_i^2\xi_{ii} + o(\gamma_{K,n})
\end{aligned}$$

Then using the definition of  $\tau$  the result follows.

## Appendix B: Expansion

For ease of notation we assume that  $\theta = (\beta, \lambda)'$  is the estimator with population value 0. Then for an estimator satisfying

$$\hat{m}(\theta) = \frac{1}{n} \sum_{i=1}^n m_i(\theta) = 0$$

, we can write the following expansion,

$$\begin{aligned} \hat{m}(\theta) &= \hat{m}(0) + \hat{M}\theta + (1/2) \sum_j \theta_j \hat{A}_j \theta + (1/6) \sum_j \sum_k \theta_j \theta_k \hat{B}_{jk} \theta + \sum_j \sum_k \sum_l \theta_j \theta_k \theta_l \hat{C}_{jkl}^* \theta \\ &= m + M\theta + (\hat{M} - M)\theta + (1/2) \sum_j \theta_j A_j \theta + (1/2) \sum_j \theta_j (\hat{A}_j - A_j) \theta \\ &\quad + (1/6) \sum_j \sum_k \theta_j \theta_k B_{jk} \theta + (1/6) \sum_j \sum_k \theta_j \theta_k (\hat{B}_{jk} - B_{jk}) \theta + (1/24) \sum_j \sum_k \sum_l \theta_j \theta_k \theta_l \hat{C}_{jkl}^* \theta \end{aligned}$$

where

$$\begin{aligned} \hat{M} &= \frac{1}{n} \sum_{i=1}^n \left( \frac{\partial}{\partial \theta} \right) m_i(0)' = \frac{1}{n} \sum_{i=1}^n M_i(0), \\ \hat{A}_j &= \frac{1}{n} \sum_{i=1}^n \left( \frac{\partial}{\partial \theta_j} \right) M_i(0), \quad \hat{B}_{jk} = \frac{1}{n} \sum_{i=1}^n \left( \frac{\partial^2}{\partial \theta_k \partial \theta_j} \right) M_i(0) \end{aligned}$$

and where the rows of  $\hat{C}_{jkl}^*$  are the third derivatives with respect to  $\theta_j$ ,  $\theta_k$ , and  $\theta_l$  evaluated at a value  $\theta^*$  that lies on the line segment joining  $\theta$  and the limit value which is 0. Consequently  $\|\theta_j^*\| < \|\theta_j\|$  for all  $j$  and  $\|\beta^*\| < \|\beta\|$ ,  $\|\lambda^*\| < \|\lambda\|$ . In the second line the terms  $M$ ,  $A_j$ , and  $B_{jk}$  are the nonstochastic (given  $x_i$ ) population values of  $\hat{M}$ ,  $\hat{A}_j$ , and  $\hat{B}_{jk}$ . Then we have,

$$\begin{aligned} \theta &= -M^{-1}m - M^{-1}(\hat{M} - M)\theta - (1/2) \sum_j \theta_j M^{-1} A_j \theta - (1/2) \sum_j \theta_j M^{-1} (\hat{A}_j - A_j) \theta \\ &\quad + (1/6) \sum_j \sum_k \theta_j \theta_k M^{-1} B_{jk} \theta + R_{n,K} \\ R_{n,K} &= -(1/6) M^{-1} \sum_j \sum_k \theta_j \theta_k (\hat{B}_{jk} - B_{jk}) \theta - (1/24) M^{-1} \sum_j \sum_k \sum_l \theta_j \theta_k \theta_l M^{-1} \bar{C}_{jkl}^* \theta \end{aligned}$$

For both GMM and GEL it is possible to show that the part of  $R_{n,K}$  that corresponds to  $\beta$ , say  $R_{n,K}^\beta$  satisfies,  $\|R_{n,K}^\beta\| = O(1/\sqrt{n})o(\gamma_{K,n})$ .

## References

- Andrews, D.W.K. (1999): "Consistent Moment Selection Procedures for Generalized Method of Moments Estimation," *Econometrica* 67, 543-564.
- Chamberlain, G. (1987): "Asymptotic Efficiency in Estimation with Conditional Moment Restrictions," *Journal of Econometrics* 34, 305-334.
- de Jong, R.M., and H. Bierens (1994): "On the Limit Behavior of a Chi-square Type Test if the Number of Conditional Moments Tested Approaches Infinity," *Econometric Theory* 9, 70-90.
- Donald, S.G. and W.K. Newey (2001): "Choosing the Number of Instruments," *Econometrica* 69, 1161-1191.
- Donald, S. G., G. Imbens and W. K. Newey (2003): "Empirical Likelihood Estimation and Consistent Tests with Conditional Moment Restrictions," *Journal of Econometrics* 117, 55-93.
- Hahn, J. and J.A. Hausman (2002): "A New Specification Test for the Validity of Instrumental Variables," *Econometrica* 70, 163-189.
- Hansen, L. P. (1982): "Large Sample Properties of Generalized Method of Moments Estimators", *Econometrica* 50, 1029-1054.
- Hansen, L.P., J. Heaton and A. Yaron (1996): "Finite-Sample Properties of Some Alternative GMM Estimators", *Journal of Business and Economic Statistics* 14, 262-280.
- Imbens, G.W., R.H. Spady and P. Johnson (1998): "Information Theoretic Approaches to Inference in Moment Condition Models," *Econometrica* 66, 333-357.
- Imbens, G.W. and R.H. Spady (2005): "The Performance of Empirical Likelihood and its Generalizations," in Andrews and Stock (eds.), *Identification and Inference for*

*Econometric Models, Essays in Honor of Thomas Rothenberg*, Cambridge University Press.

Kitamura, Y., and M. Stutzer (1997): "An Information-Theoretic Alternative to Generalized Method of Moments Estimation", *Econometrica* 65, 861-874.

Koenker, R., Machado, J.A.F., Skeels, C., and Welsh, A.H. (1994): "Momentary Lapses: Moment Expansions and the Robustness of Minimum Distance Estimation," *Econometric Theory* 10, 172-190.

Kuersteiner, G. M. (2002), "Selecting the Number of Instruments for GMM Estimators of Linear Time Series Models.", mimeo UC Davis.

Nagar, A.L. (1959): "The Bias and Moment Matrix of the General k-Class Estimators of the Parameters in Simultaneous Equations," *Econometrica* 27, 575-595.

Newey, W.K. (1991): "Uniform Convergence in Probability and Stochastic Equicontinuity," *Econometrica* 59, 1161-1167.

Newey, W.K. (1993): "Efficient Estimation of Models with Conditional Moment Restrictions," in G.S. Maddala, C.R. Rao, and H.D. Vinod, eds., *Handbook of Statistics, Volume 11: Econometrics*. Amsterdam: North-Holland.

Newey, W.K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics* 79, 147-168.

Newey, W.K. and D. McFadden (1994): "Large Sample Estimation and Hypothesis Testing," in Engle, R. and D. McFadden, eds., *Handbook of Econometrics, Vol. 4*, New York: North Holland.

Newey, W.K. and R.J. Smith (2004): "Higher-Order Properties of GMM and Generalized Empirical Likelihood Estimators," *Econometrica* 72, 219-255.

Jin, Sainan, P.C.B. Phillips, and Y. Sun (2007): "Optimal Bandwidth Selection in Heteroscedasticity-Autocorrelation Robust Testing," *Econometrica*, forthcoming.

- Owen, A. (1988): "Empirical Likelihood Ratio Confidence Regions for a Single Functional," *Biometrika* 75, 237-249.
- Qin, J. and Lawless, J. (1994): "Empirical Likelihood and General Estimating Equations", *Annals of Statistics* 22, 300-325.
- Rothenberg, T. J. (1983): "Asymptotic Properties of Some Estimators in Structural Models," in Karlin, S., T. Amemiya, and L. A. Goodman (eds.) *Studies in Econometrics, Time Series and Multivariate Statistics*, edited by New York: Academic Press.
- Rothenberg, T.J. (1984): "Approximating the Distributions of Econometric Estimators and Test Statistics," in Griliches, Z. and M.D. Intriligator, eds., *Handbook of Econometrics, Vol. 2*, New York: North-Holland.
- Rothenberg, T.J. (1996): "Empirical Likelihood Parameter Estimation Under Moment Restrictions," seminar notes, Harvard/M.I.T. and Bristol.
- Smith, R. J. (1997): "Alternative Semi-Parametric Likelihood Approaches to Generalized Method of Moments Estimation", *Economic Journal* 107, 503-519.
- Staiger, D. and J.H. Stock (1997): "Instrumental Variables Regression with Weak Instruments," *Econometrica* 65, 557-586.
- Stock, J.H., and J.H. Wright (2000): "GMM with Weak Identification", *Econometrica* 68, 1055 - 1096.

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.028	0.129	0.489	0.934	0.018
GMM-op	0.019	0.143	0.537	0.942	0.022
BGMM-all	0.013	0.163	0.616	0.864	0.012
BGMM-op	0.011	0.152	0.586	0.936	0.036
EL-all	-0.011	0.190	0.712	0.806	0.054
EL-op	0.011	0.158	0.597	0.934	0.048
ET-all	-0.004	0.195	0.716	0.790	0.048
ET-op	0.010	0.155	0.593	0.936	0.042
CUE-all	0.006	0.192	0.733	0.770	0.010
CUE-op	0.013	0.151	0.596	0.924	0.032
2SLS-all	0.027	0.126	0.447	0.958	0.026
2SLS-op	0.018	0.137	0.509	0.974	0.034
LIML-all	-0.009	0.183	0.649	0.974	0.030
LIML-op	0.009	0.141	0.564	0.980	0.026

Table I:  $n = 200$ , Cov=0.1, Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.018	0.113	0.422	0.932	0.034
GMM-op	0.013	0.125	0.478	0.926	0.044
BGMM-all	0.001	0.135	0.513	0.864	0.032
BGMM-op	0.021	0.137	0.529	0.916	0.040
EL-all	-0.018	0.173	0.646	0.782	0.174
EL-op	0.007	0.149	0.586	0.882	0.104
ET-all	-0.008	0.158	0.601	0.798	0.110
ET-op	0.014	0.148	0.564	0.878	0.088
CUE-all	-0.006	0.160	0.590	0.787	0.024
CUE-op	0.008	0.144	0.562	0.880	0.042
2SLS-all	0.034	0.118	0.443	0.948	0.040
2SLS-op	0.031	0.143	0.516	0.952	0.050
LIML-all	0.001	0.182	0.710	0.972	0.044
LIML-op	0.017	0.152	0.567	0.962	0.052

Table II:  $n = 200$ , Cov=0.1, Logistic

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.149	0.165	0.436	0.782	0.038
GMM-op	0.065	0.153	0.530	0.858	0.036
BGMM-all	0.064	0.169	0.598	0.842	0.032
BGMM-op	0.047	0.154	0.532	0.91	0.036
EL-all	-0.002	0.182	0.761	0.854	0.072
EL-op	0.036	0.162	0.552	0.896	0.052
ET-all	0.003	0.180	0.711	0.860	0.066
ET-op	0.035	0.155	0.533	0.898	0.048
CUE-all	0.002	0.177	0.734	0.840	0.022
CUE-op	0.039	0.153	0.528	0.886	0.038
2SLS-all	0.143	0.161	0.426	0.836	0.066
2SLS-op	0.066	0.152	0.517	0.900	0.046
LIML-all	0.006	0.170	0.680	0.964	0.044
LIML-op	0.041	0.154	0.527	0.946	0.048

Table III:  $n = 200$ , Cov=0.5, Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.131	0.161	0.438	0.768	0.038
GMM-op	0.079	0.154	0.516	0.854	0.044
BGMM-all	0.062	0.160	0.540	0.816	0.032
BGMM-op	0.048	0.148	0.527	0.880	0.038
EL-all	0.016	0.187	0.701	0.796	0.160
EL-op	0.041	0.156	0.578	0.860	0.090
ET-all	0.012	0.178	0.635	0.796	0.108
ET-op	0.039	0.153	0.555	0.868	0.078
CUE-all	-0.004	0.170	0.638	0.776	0.014
CUE-op	0.041	0.154	0.530	0.866	0.036
2SLS-all	0.147	0.172	0.461	0.800	0.076
2SLS-op	0.081	0.160	0.550	0.874	0.058
LIML-all	-0.007	0.175	0.707	0.936	0.06
LIML-op	0.045	0.149	0.581	0.920	0.054

Table IV:  $n = 200$ , Cov=0.5, Logistic



Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.274	0.275	0.368	0.460	0.180
GMM-op	0.124	0.189	0.565	0.798	0.078
BGMM-all	0.128	0.183	0.583	0.738	0.092
BGMM-op	0.091	0.171	0.600	0.814	0.072
EL-all	0.016	0.165	0.688	0.876	0.096
EL-op	0.056	0.168	0.599	0.846	0.126
ET-all	0.020	0.165	0.690	0.874	0.084
ET-op	0.059	0.166	0.603	0.842	0.126
CUE-all	0.024	0.165	0.681	0.880	0.034
CUE-op	0.063	0.169	0.589	0.838	0.078
2SLS-all	0.274	0.275	0.334	0.484	0.198
2SLS-op	0.115	0.186	0.559	0.820	0.062
LIML-all	0.006	0.161	0.648	0.944	0.056
LIML-op	0.041	0.156	0.623	0.900	0.108

Table V:  $n = 200$ , Cov=0.9, Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.213	0.213	0.349	0.568	0.134
GMM-op	0.092	0.136	0.505	0.874	0.076
BGMM-all	0.065	0.146	0.484	0.802	0.078
BGMM-op	0.073	0.134	0.472	0.854	0.084
EL-all	-0.006	0.146	0.620	0.886	0.158
EL-op	0.044	0.133	0.504	0.870	0.160
ET-all	-0.010	0.134	0.551	0.898	0.118
ET-op	0.039	0.126	0.470	0.892	0.144
CUE-all	-0.016	0.129	0.530	0.879	0.034
CUE-op	0.030	0.122	0.472	0.886	0.066
2SLS-all	0.242	0.244	0.347	0.580	0.190
2SLS-op	0.081	0.134	0.485	0.882	0.076
LIML-all	-0.008	0.131	0.595	0.952	0.056
LIML-op	0.032	0.127	0.557	0.934	0.108

Table VI:  $n = 200$ , Cov=0.9, Logistic

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	5	3	3	3	3	5	3+
	Mode	2	2	2	2	2	3	2
	1Q	3	2	2	2	2	3	2
	Med.	5	3	3	3	3	4	3
	3Q	8	4	5	5	4	6	4
Logistic	K	5	4+	2+	3	5+	5	3+
	Mode	10	3	2	2	3	3	3
	1Q	4	3	2	3	3	3	2
	Med.	6	4	4	4	4	4	3
	3Q	9	6	6	7	7	6	4

Table VII: Statistics for  $\hat{K}$ ,  $n = 200$ , cov=0.1

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	3	3-	4-	4-	4-	3	4
	Mode	2	2	2	2	2	2	3
	1Q	2	2	2	2	2	2	2
	Med.	3	3	3	3	3	3	3
	3Q	4	4	5	5	5	4	5
Logistic	K	3	4-	3-	4-	10-	3	4-
	Mode	3	3	3	3	3	2	2
	1Q	2	2	2	3	3	2	2
	Med.	4	4	4	4	5	3	3
	3Q	6	6	7	7	7	4	5

Table VIII: Statistics for  $\hat{K}$ ,  $n = 200$ , cov=0.5

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	2	2+	4+	4+	4+	2	5
	Mode	2	2	3	3	3	2	3
	1Q	2	2	3	3	3	2	3
	Med.	2	3	4	4	4	2	4
	3Q	3	3	7	7	6	3	7
Logistic	K	2	3+	3-	4	10	2	5
	Mode	2	2	10	10	10	2	3
	1Q	2	2	3	4	4	2	3
	Med.	2	3	6	6	6	2	4
	3Q	3	5	9	9	9	3	7

Table IX: Statistics for  $\hat{K}$ ,  $n = 200$ , cov=0.9

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.019	0.07	0.262	0.925	0.035
GMM-op	0.01	0.078	0.265	0.92	0.055
BGMM-all	0.006	0.081	0.302	0.92	0.025
BGMM-op	0.001	0.084	0.298	0.92	0.06
EL-all	0.000	0.085	0.305	0.895	0.065
EL-op	-0.001	0.080	0.296	0.900	0.060
ET-all	0.005	0.084	0.314	0.895	0.070
ET-op	0.001	0.080	0.290	0.905	0.065
CUE-all	0.006	0.082	0.307	0.895	0.025
CUE-op	0.004	0.082	0.298	0.915	0.055
2SLS-all	0.022	0.066	0.24	0.925	0.050
2SLS-op	0.005	0.074	0.275	0.920	0.060
LIML-all	-0.001	0.079	0.305	0.945	0.060
LIML-op	-0.006	0.083	0.296	0.930	0.055

Table X:  $n = 800$ , cov=0.1, Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.007	0.068	0.234	0.93	0.025
GMM-op	0.007	0.073	0.253	0.925	0.025
BGMM-all	-0.007	0.076	0.271	0.89	0.025
BGMM-op	-0.002	0.07	0.269	0.905	0.025
EL-all	-0.012	0.082	0.300	0.850	0.135
EL-op	-0.002	0.082	0.288	0.910	0.085
ET-all	-0.015	0.083	0.286	0.845	0.105
ET-op	-0.003	0.073	0.290	0.900	0.080
CUE-all	-0.005	0.08	0.281	0.856	0.025
CUE-op	-0.001	0.073	0.276	0.887	0.035
2SLS-all	0.005	0.067	0.243	0.965	0.060
2SLS-op	0.007	0.072	0.260	0.975	0.025
LIML-all	-0.012	0.078	0.314	0.975	0.060
LIML-op	-0.010	0.069	0.297	0.98	0.045

Table XI:  $n = 800$ , cov=0.1, Logistic

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.087	0.094	0.237	0.770	0.070
GMM-op	0.034	0.081	0.269	0.910	0.035
BGMM-all	0.022	0.085	0.297	0.860	0.065
BGMM-op	0.016	0.077	0.278	0.940	0.020
EL-all	0.004	0.089	0.322	0.890	0.075
EL-op	0.015	0.084	0.282	0.930	0.065
ET-all	0.005	0.089	0.314	0.880	0.075
ET-op	0.015	0.085	0.282	0.935	0.065
CUE-all	0.009	0.090	0.322	0.870	0.050
CUE-op	0.015	0.082	0.276	0.935	0.040
2SLS-all	0.089	0.090	0.231	0.805	0.065
2SLS-op	0.035	0.077	0.255	0.915	0.035
LIML-all	0.004	0.085	0.319	0.960	0.055
LIML-op	0.018	0.083	0.281	0.955	0.040

Table XII:  $n = 800$ ,  $\text{cov}=0.5$ , Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.082	0.091	0.241	0.790	0.05
GMM-op	0.042	0.078	0.274	0.870	0.045
BGMM-all	0.022	0.079	0.291	0.870	0.025
BGMM-op	0.025	0.081	0.285	0.895	0.055
EL-all	-0.003	0.089	0.312	0.875	0.155
EL-op	0.016	0.083	0.287	0.885	0.100
ET-all	0.000	0.082	0.303	0.870	0.115
ET-op	0.018	0.080	0.275	0.885	0.080
CUE-all	0.001	0.086	0.302	0.838	0.025
CUE-op	0.017	0.077	0.273	0.880	0.045
2SLS-all	0.093	0.093	0.224	0.800	0.055
2SLS-op	0.041	0.076	0.278	0.890	0.060
LIML-all	-0.009	0.076	0.286	0.955	0.055
LIML-op	0.021	0.078	0.274	0.925	0.06

Table XIII:  $n = 800$ ,  $\text{cov}=0.5$ , Logistic

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.176	0.176	0.212	0.415	0.145
GMM-op	0.064	0.093	0.273	0.880	0.055
BGMM-all	0.060	0.100	0.297	0.815	0.070
BGMM-op	0.044	0.083	0.294	0.875	0.060
EL-all	0.010	0.078	0.287	0.915	0.085
EL-op	0.032	0.082	0.272	0.895	0.105
ET-all	0.016	0.078	0.279	0.920	0.115
ET-op	0.025	0.079	0.273	0.920	0.135
CUE-all	0.012	0.080	0.280	0.925	0.075
CUE-op	0.025	0.079	0.276	0.920	0.115
2SLS-all	0.166	0.166	0.217	0.455	0.140
2SLS-op	0.061	0.089	0.268	0.88	0.050
LIML-all	0.020	0.079	0.285	0.95	0.060
LIML-op	0.035	0.080	0.276	0.93	0.100

Table XIV:  $n = 800$ , cov=0.9, Normal

Est.	Med. Bias	Med. AD	Dec. Rge	Cov.	Over.
GMM-all	0.143	0.145	0.179	0.530	0.130
GMM-op	0.046	0.089	0.274	0.885	0.045
BGMM-all	0.033	0.070	0.230	0.885	0.075
BGMM-op	0.039	0.074	0.263	0.875	0.075
EL-all	0.003	0.066	0.289	0.910	0.180
EL-op	0.024	0.072	0.256	0.920	0.165
ET-all	-0.002	0.086	0.281	0.910	0.115
ET-op	0.015	0.079	0.281	0.910	0.125
CUE-all	-0.001	0.083	0.264	0.919	0.035
CUE-op	0.013	0.079	0.277	0.911	0.040
2SLS-all	0.161	0.161	0.199	0.510	0.135
2SLS-op	0.058	0.089	0.304	0.870	0.085
LIML-all	0.004	0.077	0.304	0.975	0.075
LIML-op	0.016	0.076	0.263	0.955	0.115

Table XV:  $n = 800$ , cov=0.9, Logistic

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	10	7	7+	7+	7+	10	8-
	Mode	8	6	6	6	6	8	6
	1Q	8	6	6	6	6	7	6
	Med.	12	7	7	7	7	9	7
	3Q	17	9	9	9	9	13	9
Logistic	K	10	9+	6	7	11-	10	8
	Mode	20	10	7	8	8	8	7
	1Q	10	8	6	7	8	7	6
	Med.	15	11	8	10	12	9	7
	3Q	19	16	12	15	17	12	8

Table XVI: Statistics for  $\hat{K}$ ,  $n = 800$ , cov=0.1

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	6-	7-	8	8	8	6-	9-
	Mode	5	6	7	7	7	5	7
	1Q	5	6	6	6	6	5	6
	Med.	6	7	8	8	8	6	8
	3Q	7	8	10	10	9	7	10
Logistic	K	6-	8	7-	8	20	6-	9-
	Mode	6	6	6	6	20	6	6
	1Q	5	7	6	7	9	5	6
	Med.	6	10	9	11	13	6	8
	3Q	8	15	14	15	18	7	10

Table XVII: Statistics for  $\hat{K}$ ,  $n = 800$ , cov=0.5

		GMM	BGMM	EL	ET	CUE	TSLS	LIML
Normal	K	4	6+	9	9	9	4	11
	Mode	4	5	8	9	9	4	8
	1Q	4	5	8	8	7	4	8
	Med.	4	6	10	10	10	4	10
	3Q	5	7	15	15	14	5	14
Logistic	K	4	7	7	9	20	4	11
	Mode	4	7	20	20	20	4	8
	1Q	4	6	8	10	13	3	7
	Med.	4	8	12	15	17	4	9
	3Q	5	12	18	19	19	4	14

Table XVIII: Statistics for  $\hat{K}$ ,  $n = 800$ , cov=0.9