

## LECTURE 9 NOTES

**1. Asymptotic normality of GMM estimators.** Recall GMM estimators minimize the quadratic form

$$Q_n(\theta) := \frac{1}{2} \left\| \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta) \right\|_{\widehat{W}_n}^2.$$

Its first two derivatives are

$$\begin{aligned} \nabla Q_n(\theta) &= \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\theta) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta) \right) \\ \nabla^2 Q_n(\theta) &= \left( \frac{1}{n} \sum_{i \in [n]} \nabla^2 g_{\mathbf{x}_i}(\theta^*) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta) \right) \\ &\quad + \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\theta) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\theta) \right)^T. \end{aligned}$$

There is a similar asymptotic normality result for GMM estimators. The key ingredients are

1. the consistency of GMM estimators,
2. the asymptotic normality of  $\frac{1}{\sqrt{n}} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta^*)$ ,
3. the validity of the second-order Taylor expansion of  $g_x$  at  $\theta^*$ .

The proof even mimics the proof of Lecture 8, Theorem 2.1.

**THEOREM 1.1.** *Let  $\{\mathbf{x}_i\}$  be i.i.d. random variables distributed according to  $F_{\theta^*}$ . If*

1.  $\frac{1}{\sqrt{n}} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta^*) \xrightarrow{d} \mathcal{N}(0, \mathbf{var}_{\theta^*}[g_{\mathbf{x}_1}(\theta^*)])$
2.  $g_x$  is continuously differentiable,
3.  $\mathbf{E}_{\theta^*}[\nabla^2 Q(\theta^*)]$  is non-singular.

*If the GMM estimator is consistent, then*

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, (GWG^T)^{-1}GWJ(GW)^T(GWG^T)^{-1}),$$

*where  $G := \mathbf{E}_{\theta^*}[\nabla g_{\mathbf{x}_1}(\theta^*)]$  and  $J := \mathbf{var}_{\theta^*}[g_{\mathbf{x}_1}(\theta^*)]$ .*

**PROOF SKETCH.** Since  $\hat{\theta}_n$  minimizes  $Q_n$ ,

$$0 = \nabla Q_n(\hat{\theta}_n) = \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\hat{\theta}_n) \right),$$

which, by Taylor's theorem, is

$$= \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta^*) + \nabla g_{\mathbf{x}_i}(\tilde{\theta}_n)(\hat{\theta}_n - \theta^*) \right).$$

Rearranging,

$$\begin{aligned} & \sqrt{n}(\hat{\theta}_n - \theta^*) \\ &= \left( \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\tilde{\theta}_n) \right)^T \right)^{-1} \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{\sqrt{n}} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta^*) \right). \end{aligned}$$

The technical part of the proof consists of showing

1.  $\left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{\sqrt{n}} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta^*) \right) \xrightarrow{d} \mathcal{N}(0, (GW)J(GW)^T)$ , and
2.  $\left( \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\hat{\theta}_n) \right) \widehat{W}_n \left( \frac{1}{n} \sum_{i \in [n]} \nabla g_{\mathbf{x}_i}(\tilde{\theta}_n) \right)^T \right)^{-1} \xrightarrow{p} (GWG^T)^{-1}$

by appealing to the consistency of  $\hat{\theta}_n$  and the usual barrage of convergence theorems. We skip the technical details here. We combine the two limits by Slutsky's theorem to deduce

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, (GWG^T)^{-1}GWJ(GW)^T(GWG^T)^{-1}),$$

where we plugged in  $GWG^T$  for  $\nabla^2 Q(\theta^*)$ .  $\square$

Although the choice of  $W$  does not affect the consistency of GMM estimators, it affects the asymptotic variance of the estimator. As we shall see, choosing  $W = J^{-1}$  minimizes the asymptotic variance.

LEMMA 1.2. *For any  $W \succ 0$ ,*

$$(GWG^T)^{-1}GWJ(GW)^T(GWG^T)^{-1} \succeq (GJ^{-1}G^T)^{-1}.$$

Of course,  $\mathbf{var}_P[g_{\mathbf{x}}(\theta^*)]$  is typically unknown in practice. To obtain an efficient GMM estimator, i.e. a GMM estimator with minimum asymptotic variance, practitioners usually resort to a two-step procedure:

1. obtain a consistent pilot estimator  $\tilde{\theta}_n$  e.g. by an inefficient GMM estimator.
2. obtain  $\hat{\theta}_n$  by  $\arg \min_{\theta \in \Theta} \left\| \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\theta) \right\|_{\widehat{W}_n}^2$ , where

$$\widehat{W}_n \leftarrow \frac{1}{n} \sum_{i \in [n]} g_{\mathbf{x}_i}(\tilde{\theta}_n) g_{\mathbf{x}_i}(\tilde{\theta}_n)^T.$$

Since  $\widehat{W}_n \xrightarrow{p} \mathbf{var}_P[g_{\mathbf{x}_1}(\theta^*)]$  by the LLN, it is possible to adapt the proof of Theorem 1.1 to show that  $\{\hat{\theta}_n\}$  is asymptotically normal, and its asymptotic variance is  $(GJ^{-1}G^T)^{-1}$ .

**2. Asymptotic efficiency.** In the notes on point estimation, we considered a decision theoretic approach to evaluating point estimators. Here, we consider an asymptotic approach.

The appeal of the asymptotic approach is the possibility of *exactly* characterizing the risk (typically MSE) of a point estimator. For any sequence of consistent estimators, the asymptotic bias vanishes. Thus the asymptotic MSE is the asymptotic variance.

**THEOREM 2.1** (Cramér-Rao inequality). *Assume the density of  $\mathbf{x}$  obeys*

$$\mathbf{E}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})]] = 0.$$

Let  $\tau(\theta) := \mathbf{E}_\theta[\delta(\mathbf{x})]$ . For any estimator  $\delta(\mathbf{x})$  that obeys

$$(2.1) \quad \nabla_j \tau(\theta) = \nabla_j \mathbf{E}_\theta[\delta(\mathbf{x})] = \int_{\mathcal{X}} \delta(\mathbf{x}) \nabla_j f_\theta(x) dx,$$

for any  $j \in [p]$ , we have

$$\mathbf{var}_\theta[\delta(\mathbf{x})] \succeq \nabla \tau(\theta) \mathbf{var}_\theta[\nabla_\theta \log f_\theta(\mathbf{x})]^{-1} \nabla \tau(\theta),$$

where  $A \succeq B$  means  $A - B$  is positive semi-definite.

**PROOF.** Consider

$$\begin{aligned} & \mathbf{var}_\theta \left[ \begin{bmatrix} \delta(\mathbf{x}) \\ \nabla_\theta \log f_\theta(\mathbf{x}) \end{bmatrix} \right] \\ &= \begin{bmatrix} \mathbf{var}_\theta[\delta(\mathbf{x})] & \mathbf{cov}_\theta[\delta(\mathbf{x}) \nabla_\theta[\log f_\theta(\mathbf{x})]^T] \\ \mathbf{cov}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})] \delta(\mathbf{x})^T] & \mathbf{var}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})]] \end{bmatrix}. \end{aligned}$$

It is, by definition, positive semi-definite, which implies the Schur complement of the lower-right block

$$\begin{aligned} & \mathbf{var}_\theta[\delta(\mathbf{x})] \\ & - \mathbf{cov}_\theta[\delta(\mathbf{x}) \nabla_\theta[\log f_\theta(\mathbf{x})]^T] \mathbf{var}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})]]^{-1} \mathbf{cov}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})] \delta(\mathbf{x})^T] \end{aligned}$$

is positive semi-definite. Rearranging, we obtain

$$\begin{aligned} & \mathbf{var}_\theta[\delta(\mathbf{x})] \\ & \succeq \mathbf{cov}_\theta[\delta(\mathbf{x}) \nabla_\theta[\log f_\theta(\mathbf{x})]^T] \mathbf{var}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})]]^{-1} \mathbf{cov}_\theta[\nabla_\theta[\log f_\theta(\mathbf{x})] \delta(\mathbf{x})^T]. \end{aligned}$$

To complete the proof, we must show that

$$\nabla_{\theta} \mathbf{E}_{\theta}[\delta(\mathbf{x})] = \mathbf{cov}_{\theta} \left[ \delta(\mathbf{x}) \nabla_{\theta} [\log f_{\theta}(\mathbf{x})]^T \right].$$

Indeed, by differentiating under the integral sign, we have

$$\begin{aligned} \nabla_j \mathbf{E}_{\theta}[\delta(\mathbf{x})] &= \int_{\mathcal{X}} \delta(\mathbf{x}) \nabla_j f_{\theta}(x) dx \\ &= \int_{\mathcal{X}} \delta(\mathbf{x}) \frac{\nabla_j f_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \\ &= \int_{\mathcal{X}} \delta(\mathbf{x}) \nabla_j [\log f_{\theta}(x)] f_{\theta}(x) dx \\ &= \mathbf{E}_{\theta} \left[ \delta(\mathbf{x}) \nabla_j [\log f_{\theta}(\mathbf{x})] \right], \end{aligned}$$

which, by the assumption  $\mathbf{E}_{\theta}[\nabla_{\theta} [\log f_{\theta}(\mathbf{x})]] = 0$ , is

$$= \mathbf{cov}_{\theta} \left[ \delta(\mathbf{x}) \nabla_j [\log f_{\theta}(\mathbf{x})] \right].$$

for any  $j \in [p]$ . □

As the Fisher Information increases, the Cramér-Rao inequality decreases. When  $\mathbf{x}$  consists of *i.i.d.* random variables  $\mathbf{x}_i$ ,  $i \in [n]$ , the Fisher Information increases linearly as the sample size grows:

$$\begin{aligned} I(\theta) &= \mathbf{var}_{\theta} [\nabla_{\theta} [\log f_{\theta}(\mathbf{x})]] \\ &= \mathbf{var}_{\theta} \left[ \sum_{i \in [n]} \nabla_{\theta} [\log f_{\theta}(\mathbf{x}_i)] \right] \\ &= n \mathbf{var}_{\theta} [\nabla_{\theta} [\log f_{\theta}(\mathbf{x}_1)]], \end{aligned}$$

where  $f_{\theta}(x_1)$  is the density of  $\mathbf{x}_1$ .

We hasten to mention that the Cramér-Rao lower bound (as stated) is valid for any estimator that obeys (2.1). If  $\delta(\mathbf{x})$  is unbiased, then

$$\nabla_{\theta} \mathbf{E}_{\theta}[\delta(\mathbf{x})] = \nabla_{\theta} \theta = I_p,$$

and the Cramér-Rao lower bound simplifies to

$$(2.2) \quad \mathbf{var}_{\theta}[\delta(\mathbf{x})] \succeq \mathbf{var}_{\theta}[\nabla_{\theta} [\log f_{\theta}(\mathbf{x})]]^{-1}.$$

Some sources refer to (2.2) as the Cramér-Rao lower bound.

Finally, we remark that the Cramér-Rao inequality is non-asymptotic: it is valid for finite sample sizes. It is generally not sharp; i.e. there may not

be any estimator that saturates the inequality. Even in the usually favorable setting of a one-dimensional exponential family, it is only known that  $\phi(\mathbf{x})$  is an unbiased estimator of  $\mathbf{E}_\theta[\phi(\mathbf{x})]$  that attains the lower bound. It is not known whether there are estimators of other functions of  $\theta$  that attain the lower bound. However, as we shall see, the MLE saturates the inequality in asymptopia.

DEFINITION 2.2. A sequence of estimators  $\{\hat{\theta}_n\}$  is asymptotically efficient at a parameter  $\tau(\theta) \in \mathbf{R}^p$  if

1.  $\{\hat{\theta}_n\}$  is consistent:  $\hat{\theta}_n \xrightarrow{P} \theta$
2.  $\{\hat{\theta}_n\}$  is asymptotically normal:  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} \mathcal{N}(0, V(\theta))$ ,
3. the asymptotic variance of  $\{\hat{\theta}_n\}$  attains the Cramér-Rao lower bound:

$$V(\theta) = I(\theta^*)^{-1}.$$

By comparing the asymptotic distribution of the MLE with the Cramér-Rao lower bound, we see that the MLE is asymptotically efficient. The ratio between the asymptotic variance of the MLE and that of another estimator is the *asymptotic relative efficiency* (ARE) of the other estimator.

EXAMPLE 2.3. Let  $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \text{Poi}(\lambda)$ . We wish to estimate the probability of observing no events  $\mathbf{P}_\lambda(\mathbf{x} = 0)$ , which, by the Poisson density, is  $e^{-\lambda}$ . A naive estimator is

$$\hat{\theta}_{\text{Ber},n} = \frac{1}{n} \sum_{i \in [n]} \mathbf{1}\{\mathbf{x}_i = 0\}.$$

By the CLT, its asymptotic distribution is

$$\sqrt{n}(\hat{\theta}_{\text{Ber},n} - e^{-\lambda}) \xrightarrow{d} \mathcal{N}(0, e^{-\lambda}(1 - e^{-\lambda})).$$

The MLE of  $e^{-\lambda}$  is  $e^{-\bar{\mathbf{x}}}$ , where  $\bar{\mathbf{x}}$  is the MLE of  $\lambda$ . By the CLT, the asymptotic variance of  $\bar{\mathbf{x}}_n$  is  $\lambda$ . By the delta method,

$$\sqrt{n}(e^{-\bar{\mathbf{x}}_n} - e^{-\lambda}) \xrightarrow{d} \mathcal{N}(0, \lambda e^{-2\lambda})$$

We see that the asymptotic relative efficiency of  $\hat{\theta}_{\text{Ber}}$  to the MLE  $e^{-\bar{\mathbf{x}}_n}$  is

$$\frac{\lambda e^{-2\lambda}}{e^{-\lambda}(1 - e^{-\lambda})} = \frac{\lambda}{e^\lambda - 1},$$

which decreases exponentially as  $\lambda$  grows.