Applied Topics for PhD Monte Carlo & Bootstrap

Michael E. Kummer Theoretical Slide Set, based on Stephen Kastoryano

OTIM, Nova School of Business and Economics

Table of Contents

SIMULATION TO UNDERSTAND ESTIMATORS BETTER

- Generally we always had (y, \mathbf{x}) are and i.i.d. draw from some F
- Then all "statistics", including estimators etc.

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), \theta))$$

SIMULATION TO UNDERSTAND ESTIMATORS BETTER

- Generally we always had (y, \mathbf{x}) are and i.i.d. draw from some F
- Then all "statistics", including estimators etc.

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), \theta))$$

- e.g. t-ratio $(\hat{\theta} \theta)/s(\hat{\theta})$
- or some estimator $\hat{\theta}$ (or $\hat{\beta}$)

SIMULATION TO UNDERSTAND ESTIMATORS BETTER

- Generally we always had (y, \mathbf{x}) are and i.i.d. draw from some F
- Then all "statistics", including estimators etc.

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), \theta))$$

- e.g. t-ratio $(\hat{\theta} \theta)/s(\hat{\theta})$
- or some estimator $\hat{\theta}$ (or $\hat{\beta}$)
- depends on "the data" and *F*, it's distribution.

- We know T_n (estimators, test statistics) are just another random variable
- With it's own expectation and distribution.
- So we model the c.d.f. as:

$$G_n(u,F) = Pr(T_n < u | F)$$

• but while the asymptotic behavior of T_n may be known, the one of G_n generally is not

LETS GO TO THE CASINO

- We can learn about the statistic, by simulation.
- Here's how:
 - We choose some F and then simulate data.
 - This nails down a true θ
 - Now we pass that data to our estimator...
 - …and see if it "get's it right"

LETS GO TO THE CASINO

- We can learn about the statistic, by simulation.
- Here's how:
 - We choose some F and then simulate data.
 - This nails down a true θ
 - Now we pass that data to our estimator...
 - …and see if it "get's it right"
- it'e like simulating it in a sandbox

HERE IS WHAT TO DO:

- Define some DGP (i.e. F) and draw n observations.
- we get *n* observations $(y_i^*, \mathbf{x_i}^*)$ i = 1, ..., n
- compute $T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), \theta))$ on the pseudo data.
- This generates **one** random draw of the unknown $G_n(u, F)$

NOTATION - BTH EXPERIMENT

- repeat over and over.
- *B* times. (like 1000 or 5000 times)
- see how the estimator performs by varying n and F

NOTATION - BTH EXPERIMENT

- repeat over and over.
- B times. (like 1000 or 5000 times)
- see how the estimator performs by varying n and F
- Notation:
 - T_{nb} ... is the outcome of the b^{th} experiment (b = 1, ..., B)
 - So Repeating B times gives:
 - a "random sample of size B"
 - coming from $G_n(u, F) = Pr(T_{nb} < u) = Pr(T_n < u | F)$

RESULTING TEST-STATISTICS

• e.g. Bias

$$\widehat{Bias(\hat{\theta})} = \frac{1}{B} \sum_{b=1}^{B} T_{nb} = \frac{1}{B} \sum_{b=1}^{B} \hat{\theta}_b - \theta$$

• MSE:

$$\widehat{MSE(\hat{\theta})} = \frac{1}{B} \sum_{b=1}^{B} (T_{nb})^2 = \frac{1}{B} \sum_{b=1}^{B} (\hat{\theta}_b - \theta)^2$$

• or even the type I error of a two-sided t-test with $T_n=(\hat{ heta}- heta)/s(\hat{ heta})$

$$\hat{P} = \frac{1}{B} \sum_{b=1}^{B} \mathbb{1}(T_{nb}) \le 1.96)$$

- ...with 1() the indicator function, takes value 1 if condition true
- in other words simply counting.

M Kummer

- Simulate some data for an interesting DGP (with a known Problem), absent other ideas, revisit the omitted variables bias problem
- **②** Now run this not only once but B = 1000 times.
- Sech time run the estimation procedure and store the result.
 - Note: You can either save 1000 datasets, or keep only the estimator.
- Intermediate of the the State of the stimution.
- Ompare with the CDF you know.
- **(**) You can vary *n* and see whether the estimation gets better.

Table of Contents

INTRODUCTION

- Previous lectures were using well defined and understood estimator.
- Asymptotic behavior has been researched by earlier research.
- BUT: On high level research (as PhD-level applied research) the "straight forward" methods are not available

INTRODUCTION

- Previous lectures were using well defined and understood estimator.
- Asymptotic behavior has been researched by earlier research.
- BUT: On high level research (as PhD-level applied research) the "straight forward" methods are not available
- A useful non-parametric method is "the Bootstrap"
 - works in many settings.
 - frequently also for data or residuals from unknown distributions.

• Start from "some test-statistic."

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), F))$$

• Start from "some test-statistic."

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), F))$$

- e.g. t-ratio $(\hat{\theta} \theta)/s(\hat{\theta})$
- or some estimator $\hat{\theta}$ (or $\hat{\beta}$)

• Start from "some test-statistic."

$$T_n = T_n((y_1, x_1), (y_2, x_2), ..., (y_n, x_n), F))$$

- e.g. t-ratio $(\hat{\theta} \theta)/s(\hat{\theta})$
- or some estimator $\hat{\theta}$ (or $\hat{\beta}$)
- depends on "the data" and *F*, it's distribution.

- We know T_n (estimators, test statistics) are just another random variable
- With it's own expectation and distribution.
- So we model the c.d.f. as:

$$G_n(u,F) = Pr(T_n < u | F)$$

- But we do not know F
- So what can we do?

- But we do not know F
- So what can we do?
- We do not know F, but we have tons of data, that were generated by F
- Let's pull ourselves out of the swamp on our own Bootstraps

- But we do not know F
- So what can we do?
- We do not know F, but we have tons of data, that were generated by F
- Let's pull ourselves out of the swamp on our own Bootstraps
 - by learning about the distribution from our data.

 $G_n^*(u) = G_n(u, F_n)$

EMPIRICAL DISTRIBUTION FUNCTION 1

Recall

$$\boldsymbol{F}(\boldsymbol{y}, \boldsymbol{x}) = Pr(y_i < y, \boldsymbol{x}_i < \boldsymbol{x}) = \mathbb{E}(1(y_i < y)1(\boldsymbol{x}_i < \boldsymbol{x}))$$

EMPIRICAL DISTRIBUTION FUNCTION 1

Recall

$$\boldsymbol{F}(\boldsymbol{y}, \boldsymbol{x}) = Pr(y_i < y, \boldsymbol{x}_i < \boldsymbol{x}) = \mathbb{E}(1(y_i < y)1(\boldsymbol{x}_i < \boldsymbol{x}))$$

• The corresponding method of moments estimator simply plugs in the empirical counterpart:

- (...as always.)
- Hence,

$$\boldsymbol{F_n(\mathbf{y}, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i < y)\mathbb{1}(\mathbf{x}_i < \mathbf{x})}$$

• This is called the "Empirical Distribution Function" (EDF)

- EDF consistently estimates CDF. Proof is based on
 - ▶ the insight that for any (y, \mathbf{x}) , $1(y_i < y)1(\mathbf{x}_i < \mathbf{x})$ is an i.i.d. RV drawn from $F(y, \mathbf{x})$
 - WLLN
 - CLT

- EDF consistently estimates CDF. Proof is based on
 - the insight that for any (y, \mathbf{x}) , $1(y_i < y)1(\mathbf{x}_i < \mathbf{x})$ is an i.i.d. RV drawn from $F(y, \mathbf{x})$
 - WLLN
 - CLT

• Hence,

$$\sqrt{n}(\boldsymbol{F_n(y,x)} - \boldsymbol{F(y,x)}) \stackrel{d}{\rightarrow} \boldsymbol{N}(0, \boldsymbol{F(y,x)}(1 - \boldsymbol{F(y,x)}))$$

- And the EDF is a valid **discrete** prob. distribution.
- Probability mass is 1/n for each data-point $(y_i, \mathbf{x_i})$
- i.o.w. think of a random pair $(y_i^*, \mathbf{x_i}^*)$ as drawn from F_n

 $Pr(y_i^* < y, \mathbf{x}_i^* < \mathbf{x}) = F_n(y, \mathbf{x})$

• Sample moments of (y_i^*, \mathbf{x}_i^*)

$$\mathbf{E}h(\mathbf{y}_{i}^{*}, \mathbf{x}_{i}^{*}) = \int (h(y_{i}, \mathbf{x}_{i}) \mathbf{F}_{n}(\mathbf{y}, \mathbf{x}))$$
$$= \sum_{i=1}^{n} h(y_{i}, \mathbf{x}_{i}) \mathbf{P}r(y_{i}^{*} = y_{i}) \mathbf{1}(\mathbf{x}_{i}^{*} = \mathbf{x}_{i})$$
$$= \frac{1}{n} \sum_{i=1}^{n} h(y_{i}, \mathbf{x}_{i})$$



Figure 10.1: Empirical Distribution Functions

• Sample moments of (y_i^*, \mathbf{x}_i^*)

$$\mathbf{E}h(\mathbf{y}_{i}^{*}, \mathbf{x}_{i}^{*}) = \int (h(y_{i}, \mathbf{x}_{i}) \mathbf{F}_{n}(\mathbf{y}, \mathbf{x}))$$
$$= \sum_{i=1}^{n} h(y_{i}, \mathbf{x}_{i}) \mathbf{P}r(y_{i}^{*} = y_{i}) \mathbf{1}(\mathbf{x}_{i}^{*} = \mathbf{x}_{i})$$
$$= \frac{1}{n} \sum_{i=1}^{n} h(y_{i}, \mathbf{x}_{i})$$

Table of Contents

NONPARAMETRIC BOOTSTRAP

- Use EDF as the estimate F_n of F
- Use Monte-Carlo simulations to approximate G_n^*
- Note:
 - The sample size in the simulation should be the same as the sample size
 - (y_i^*, \mathbf{x}_i^*) are drawn from the data, i.e.
 - ▶ You randomly sample (WITH REPLACEMENT!) *n* observations (*y*, **x**)

NONPARAMETRIC BOOTSTRAP

- Use EDF as the estimate F_n of F
- Use Monte-Carlo simulations to approximate G_n^*
- Note:
 - The sample size in the simulation should be the same as the sample size
 - (y_i^*, \mathbf{x}_i^*) are drawn from the data, i.e.
 - You randomly sample (WITH REPLACEMENT!) n observations (y,x)
- And you get your first instance of T_n^* :

$$T_n^* = T_n((y_1^*, x_1^*), (y_2^*, x_2^*), ..., (y_n^*, x_n^*), F_n))$$

• ...that's simply the statistic computed with the simulated sample you have just drawn.

NONPARAMETRIC BOOTSTRAP

- Now we need to get the distribution (confidence interval)
- Hence repeat this *B* times.
- Note:
 - B is the number of bootstrap replications.
 - Large B is better, but takes longer to compute
 - A theory about that was developed (Andrews and Buchinsky, 2000), but
 - B = 1000 is typically sufficient and reasonable to compute.

EXAMPLE: T-RATIO

- Now we need to get the distribution (confidence interval)
- Hence repeat this B times.
- Note:
 - B is the number of bootstrap replications.
 - Large B is better, but takes longer to compute
 - A theory about that was developed (Andrews and Buchinsky, 2000), but
 - B = 1000 is typically sufficient and reasonable to compute.
- Finally: The general test statistic T_n usually is a function of F
 - e.g. t-ratio $(\hat{\theta} \theta)/s(\hat{\theta})$ depends on θ :
- Bootstrap replaces F with F_n and similarly θ with θ_n
 - …usually that's simply the newly estimated $\hat{\theta}$
ESTIMATING BIAS

- Denote the Bias $E(\hat{\theta} \theta_0) = \tau_n$
- Let $T_n(\theta) = (\hat{\theta} \theta)$ (*...for any theta*), hence $\tau_n = E(T_n(\theta_0))$
- The Bootstrap estimate of au_n is given by

ESTIMATING BIAS

- Denote the Bias $E(\hat{\theta} \theta_0) = \tau_n$
- Let $T_n(\theta) = (\hat{\theta} \theta)$ (*...for any theta*), hence $\tau_n = E(T_n(\theta_0))$
- The Bootstrap estimate of au_n is given by

 $\tau_n^* = \mathrm{E}(T_n^*(\theta_0))$

• ...it's all in the stars

COMPUTATION OF THE ESTIMATE

- Denote the Bias $E(\hat{\theta} \theta_0) = \tau_n$
- Let $T_n(\theta) = (\hat{\theta} \theta)$ (*...for any theta*), hence $\tau_n = E(T_n(\theta_0))$
- Using the simulation, the estimate of au_n is given by

COMPUTATION OF THE ESTIMATE

- Denote the Bias $E(\hat{\theta} \theta_0) = \tau_n$
- Let $T_n(\theta) = (\hat{\theta} \theta)$ (*...for any theta*), hence $\tau_n = E(T_n(\theta_0))$
- Using the simulation, the estimate of au_n is given by

$$\hat{\tau_n^*} = \frac{1}{B} \sum_{b=1}^{B} T_{nb}^*$$
$$= \frac{1}{B} \sum_{b=1}^{B} \hat{\theta}^* - \hat{\theta}$$
$$= \overline{\hat{\theta}^*} - \hat{\theta}$$

COMPUTATION OF THE ESTIMATE

- Denote the Bias $E(\hat{\theta} \theta_0) = \tau_n$
- Let $T_n(\theta) = (\hat{\theta} \theta)$ (*...for any theta*), hence $\tau_n = E(T_n(\theta_0))$
- Using the simulation, the estimate of au_n is given by

$$\hat{\tau_n^*} = \frac{1}{B} \sum_{b=1}^{B} T_{nb}^*$$
$$= \frac{1}{B} \sum_{b=1}^{B} \hat{\theta}^* - \hat{\theta}$$
$$= \overline{\hat{\theta}^*} - \hat{\theta}$$

• ...that is, subtract your "real" estimate from the BS average estimate. Similarly the variance:

$$\hat{V_n^*} = \frac{1}{B} \sum_{b=1}^{B} (\hat{\theta}_b^* - \overline{\hat{\theta}^*})^2$$

 This procedure also can afford a bias-correction and Percentile Interval Estimation, Symmetric Percentile-t Intervals, Asymptotic Expansion etc. (cf. Hansen, ch. 10)

M Kummer

Table of Contents

NON-LINEAR IN ERRORS

- Models which can not be transformed into linear models. For example, models with limited dependent variables
 - Binary variables: employed vs. not employed outcome.
 - Categorical variables: Choose between 5 political candidates.
 - Nonnegative variables: wages, prices, interest rates.
 - Nonnegative variables with excess zeros: labor supply, doctor visits.
 - Count variables: the number of cigarettes smoked per day.
 - Censored variables: unemployment durations.

NON-LINEAR IN ERRORS

- Models which can not be transformed into linear models. For example, models with limited dependent variables
 - Binary variables: employed vs. not employed outcome.
 - Categorical variables: Choose between 5 political candidates.
 - Nonnegative variables: wages, prices, interest rates.
 - Nonnegative variables with excess zeros: labor supply, doctor visits.
 - Count variables: the number of cigarettes smoked per day.
 - Censored variables: unemployment durations.
- The partial effects no longer straightforward

$$\frac{\partial \mathbf{E}(y_i | \mathbf{x}_i)}{\partial \mathbf{x}_i} \neq \beta$$

- Now the effects depend upon the level of the variables x_i .
- The average partial effects (APE): $E \frac{\partial E(y_i | \mathbf{x}_i)}{\partial \mathbf{x}_i}$ estimated by $\frac{1}{n} \sum_i^n \frac{\partial E(y_i | \mathbf{x}_i)}{\partial \mathbf{x}_i}$.
- The partial effects at the average (PEA): Partial effects $\frac{\partial E(y_i|\mathbf{x}_i)}{\partial \mathbf{x}_i}$ when fixing $\mathbf{x}_i = E(\mathbf{x}_i)$ estimated by $\overline{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i$.

- Likelihood principle (R. A. Fisher, 1922): Choose as estimator of parameter vector θ_0 the value of θ which maximizes likelihood of observing the actual data in sample within proposed model.
 - ▶ In discrete case: likelihood is the probability obtained from the probability mass function.
 - In continuous case: likelihood is the is the density (since probability of observing single realization of continuous variable is always 0).
- Example in discrete case: If one value of θ implies that probability of the observed data occurring is 0.0010, and other value of θ gives a higher probability of 0.0015, then second value better estimator.

Table of Contents

LIKELIHOOD AND LOG LIKELIHOOD FUNCTION

- Joint probability mass function or density f(y, X|θ) is viewed here as a function of θ given the data (y, X).
- This is called the likelihood function and is denoted by $= L(\theta|\mathbf{y}, \mathbf{X})$.
- Maximizing $L(\theta)$ is equivalent to maximizing the log-likelihood function

$$\mathcal{L}(\theta) = lnL(\theta)$$

 ML estimator special property: it is most efficient estimator among consistent and asymptotically normal estimators.

LIKELIHOOD FUNCTION AND CONDITIONAL LIKELIHOOD

- Given data (y_i, x_i), the likelihood L(θ) = f(y, X|θ) = f(y|X, θ)f(X|θ) requires specification of conditional density of y given X and the marginal density of X.
- Often we assume data generating process (DGP) of y given X and of X depend on mutually exclusive set of parameters, f(y|X,θ,η) = f(y|X,θ) and f(X|θ,η) = f(X|η).
- Since our primary interest is in the relationship between y_i and x_i we usually focus on conditional likelihood function: L(θ) = f(y|X,θ)
- Can you think about situation where this would not be appropriate?

LIKELIHOOD FUNCTION AND CONDITIONAL LIKELIHOOD

- Given data (y_i, x_i), the likelihood L(θ) = f(y, X|θ) = f(y|X, θ)f(X|θ) requires specification of conditional density of y given X and the marginal density of X.
- Often we assume data generating process (DGP) of y given X and of X depend on mutually exclusive set of parameters, f(y|X,θ,η) = f(y|X,θ) and f(X|θ,η) = f(X|η).
- Since our primary interest is in the relationship between y_i and x_i we usually focus on conditional likelihood function: L(θ) = f(y|X,θ)
- Can you think about situation where this would not be appropriate?
- For endogenous sampling, consistent estimation requires full density.
- For time series case (leaving out **X** for simplicity) we might have for t = 0, ..., T,

$$f(\mathbf{y}|\boldsymbol{\theta}) = f(y_T|y_{T-1},...,y_1,\boldsymbol{\theta}) \cdot f(y_{T-1}|y_{T-2},...,y_1,\boldsymbol{\theta}) \cdots f(y_2|y_1,\boldsymbol{\theta}) \cdot f(y_1|\boldsymbol{\theta})$$

= $\prod_{t=1}^T f(y_t|y_{t-1},...,y_1,\boldsymbol{\theta})$

MAXIMUM LIKELIHOOD ESTIMATOR

• Joint density of IID sample (y_1, \ldots, y_n) given $(\mathbf{x}_1, \ldots, \mathbf{x}_n)$ is,

$$L(\boldsymbol{\theta}) = \prod_{i=1}^{n} f(y_i | \mathbf{x}_i, \boldsymbol{\theta})$$

The log-likelihood in turn can be written,

$$\mathcal{L}(\theta) = \sum_{i=1}^{n} lnf(y_i | \mathbf{x}_i, \theta)$$

- Note that likelihood-based models should not be specified by making assumptions on the distribution of an error term.
- Assuming linear CEF model and $y_i | \mathbf{x}_i \sim \mathcal{N}(\mathbf{x}'_i \beta, \sigma^2)$, and independence of u_i and x_i implies $u_i | \mathbf{x}_i \sim \mathcal{N}(0, \sigma^2)$.
- So it is simpler to say errors are normally distributed u_i |x_i ~ N(0, σ²), but this is because of assumed additivity of errors.

MAXIMUM LIKELIHOOD ESTIMATOR

• The maximum likelihood estimator (MLE) $\hat{\theta}_{ML}$ is parameter value which maximizes the log-likelihood,

$$\hat{ heta}_{ML} = rg\max_{oldsymbol{ heta}} \mathcal{L}(oldsymbol{ heta})$$

- We can view $\hat{\theta}_{ML}$ as estimator within class of extremum estimators.
- Notice that $\frac{1}{n}\mathcal{L}(\theta) = \frac{1}{n}\sum_{i=1}^{n} lnf(y_i|\mathbf{x}_i, \theta) \xrightarrow{p} E(lnf(y_i|\mathbf{x}_i, \theta)).$
- The f.o.c. imply that $\hat{\theta}_{ML}$ maximises population function at $\theta = \theta_0$:

$$\left. \frac{\partial}{\partial \theta} \mathsf{E}(Inf(y_i | \mathsf{x}_i, \theta)) \right|_{\theta = \theta_0} = \mathbf{0}$$

• Sample counterpart f.o.c. imply that $\hat{ heta}_{ML}$ solves,

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\partial lnf(y_{i}|\mathbf{x}_{i},\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}=\mathbf{0}$$

• In some simple cases, explicit expression exists for $\hat{\theta}_{ML}$ as a function of data, but typically $\hat{\theta}_{ML}$ must be estimated using numerical methods.

ML IN NORMAL REGRESSION MODEL

- Normal regression model assumes linear CEF model and $y_i | \mathbf{x}_i \sim \mathcal{N}(\mathbf{x}'_i \beta, \sigma^2)$, and independence of u_i and x_i (homoskedasticity)
- The log-likelihood of the normal regression model is,

$$\mathcal{L}(\beta,\sigma^2) = \sum_{i=1}^n \ln\left(\frac{1}{(2\pi\sigma^2)^{1/2}}\exp(-\frac{1}{2\sigma^2}(y_i - \mathbf{x}'_i\beta)^2)\right)$$
$$= -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln(\sigma^2) - \frac{n}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)'$$

ML IN NORMAL REGRESSION MODEL

• Maximization with respect to β yields f.o.c.

$$\begin{aligned} &\frac{\partial}{\partial\beta}\mathcal{L}(\beta,\sigma^2) = \frac{1}{\sigma^2}(\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\beta}_{ML}) = \mathbf{0} \\ &\Leftrightarrow \qquad \hat{\beta}_{ML} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \end{aligned}$$

• The equivalence of $\hat{\beta}_{ML} = \hat{\beta}_{OLS}$ should not be surprising since maximizing $\mathcal{L}(\beta, \sigma^2)$ with respect to β is equivalent to minimizing sum of squared error criterion function for OLS.

ML IN NORMAL REGRESSION MODEL

• Maximization with respect to β yields f.o.c.

$$\begin{aligned} \frac{\partial}{\partial \beta} \mathcal{L}(\beta, \sigma^2) &= \frac{1}{\sigma^2} (\mathbf{X}' \mathbf{y} - \mathbf{X}' \mathbf{X} \hat{\beta}_{ML}) = \mathbf{0} \\ \Leftrightarrow \qquad \hat{\beta}_{ML} &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \end{aligned}$$

- The equivalence of $\hat{\beta}_{ML} = \hat{\beta}_{OLS}$ should not be surprising since maximizing $\mathcal{L}(\beta, \sigma^2)$ with respect to β is equivalent to minimizing sum of squared error criterion function for OLS.
- Maximization with respect to σ^2 yields f.o.c.

$$\frac{\partial}{\partial \sigma^2} \mathcal{L}(\beta, \sigma^2) = -\frac{n}{2\hat{\sigma}_{ML}^2} + \frac{1}{2(\hat{\sigma}_{ML}^2)^2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) = 0$$

• Plugging in $\hat{\beta}_{ML}$ for β and solving for $\hat{\sigma}^2_{ML}$ results in

$$\hat{\sigma}_{ML}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\hat{\beta}_{ML})' (\mathbf{y} - \mathbf{X}\hat{\beta}_{ML})$$

Table of Contents

SCORE AND HESSIAN

• The likelihood score is the gradient (vector of partial derivatives) of the log-likelihood evaluated at true parameter θ_0 ,

$$\mathbf{S}_{i} = \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{0}) = \frac{\partial}{\partial \boldsymbol{\theta}} lnf(y_{i} | \mathbf{x}_{i}, \boldsymbol{\theta}_{0})$$

• The Hessian is the matrix of partial derivatives of the score,

$$\mathbf{H}_{i}(\boldsymbol{\theta}) = \frac{\partial^{2}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \mathcal{L}_{i}(\boldsymbol{\theta}_{0}) = \frac{\partial^{2}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ln f(y_{i} | \mathbf{x}_{i}, \boldsymbol{\theta}_{0})$$

• We use here following shorthand notation,

$$\begin{aligned} \mathbf{S}_{i}(\tilde{\theta}) &= \frac{\partial}{\partial \theta} lnf(y_{i}|\mathbf{x}_{i},\theta) \bigg|_{\theta=\tilde{\theta}} = \frac{\partial}{\partial \theta} lnf(y_{i}|\mathbf{x}_{i},\tilde{\theta}) \\ \mathbf{H}_{i}(\tilde{\theta}) &= \frac{\partial^{2}}{\partial \theta \partial \theta'} lnf(y_{i}|\mathbf{x}_{i},\theta) \bigg|_{\theta=\tilde{\theta}} = \frac{\partial^{2}}{\partial \theta \partial \theta'} lnf(y_{i}|\mathbf{x}_{i},\tilde{\theta}) \end{aligned}$$

PROPERTIES OF LIKELIHOOD

$$\left. \frac{\partial}{\partial \theta} \mathsf{E}_{\theta_0}(\mathit{Inf}(y_i | \mathsf{x}_i, \theta)) \right|_{\theta = \theta_0} = \mathbf{0}$$

.

where $E_{\theta_0}(\cdot)$ means we are taking expectation under the DGP characterized by θ_0 .

• Proof:

M Kummer

$$\begin{aligned} \frac{\partial}{\partial \theta} \mathbf{E}_{\theta_0}(lnf(y_i|\mathbf{x}_i,\theta)) \bigg|_{\theta=\theta_0} &= \frac{\partial}{\partial \theta} \int lnf(y_i|\mathbf{x}_i,\theta) f(y_i|\mathbf{x}_i,\theta_0) dy_i \bigg|_{\theta=\theta_0} \\ &= \int \frac{\partial}{\partial \theta} f(y_i|\mathbf{x}_i,\theta) \frac{f(y_i|\mathbf{x}_i,\theta_0)}{f(y_i|\mathbf{x}_i,\theta)} dy_i \bigg|_{\theta=\theta_0} \\ &= \int \frac{\partial}{\partial \theta} f(y_i|\mathbf{x}_i,\theta_0) \frac{f(y_i|\mathbf{x}_i,\theta_0)}{f(y_i|\mathbf{x}_i,\theta_0)} dy_i \\ &= \int \frac{\partial}{\partial \theta} f(y_i|\mathbf{x}_i,\theta) dy_i \bigg|_{\theta=\theta_0} \\ &= \frac{\partial}{\partial \theta} \int f(y_i|\mathbf{x}_i,\theta) dy_i \bigg|_{\theta=\theta_0} = \frac{\partial}{\partial \theta} 1 \bigg|_{\theta=\theta_0} = \mathbf{0} \end{aligned}$$

PROPERTIES OF LIKELIHOOD

 $\mathbf{E}_{\boldsymbol{\theta}_0}(\mathbf{S}_i(\boldsymbol{\theta}_0)) = \mathbf{0}$

• Proof:

$$E_{\theta_0}(\mathbf{S}_i(\theta_0)) = E_{\theta_0}(\frac{\partial}{\partial \theta} lnf(y_i | \mathbf{x}_i, \theta_0))$$
$$= \frac{\partial}{\partial \theta} E_{\theta_0}(lnf(y_i | \mathbf{x}_i, \theta)) \bigg|_{\theta = \theta_0}$$
$$= \mathbf{0}$$

• and similarly,

$$\mathbf{E}_{\boldsymbol{\theta}_0}\left(\frac{\frac{\partial^2}{\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}'}f(y_i|\mathbf{x}_i,\boldsymbol{\theta}_0)}{f(y_i|\mathbf{x}_i,\boldsymbol{\theta}_0)}\right) = \mathbf{0}$$

M Kummer

THE HESSIAN AND THE EXPECTED OUTER PRODUCT SCORE

$$-\mathrm{E}_{\boldsymbol{\theta}_0}(\mathbf{H}_i(\boldsymbol{\theta}_0)) = \mathrm{E}_{\boldsymbol{\theta}_0}(\mathbf{S}_i(\boldsymbol{\theta}_0)\mathbf{S}_i'(\boldsymbol{\theta}_0))$$

$$\frac{\partial^2}{\partial\theta\partial\theta'} lnf(y_i|\mathbf{x}_i,\theta_0) = \frac{\frac{\partial^2}{\partial\theta\partial\theta'}f(y_i|\mathbf{x}_i,\theta_0)}{f(y_i|\mathbf{x}_i,\theta_0)} - \frac{\frac{\partial}{\partial\theta}f(y_i|\mathbf{x}_i,\theta_0)\frac{\partial}{\partial\theta}f(y_i|\mathbf{x}_i,\theta_0)'}{f(y_i|\mathbf{x}_i,\theta_0)^2}$$

• Take expectations on both sides under the DGP characterized by θ_0 ,

$$-\mathrm{E}_{\boldsymbol{\theta}_0}(\mathbf{H}_i(\boldsymbol{\theta}_0)) = \mathrm{E}_{\boldsymbol{\theta}_0}(\mathbf{S}_i(\boldsymbol{\theta}_0)\mathbf{S}_i(\boldsymbol{\theta}_0)')$$

Table of Contents

INFORMATION MATRIX EQUALITY

• Define $\mathcal{H} = -E_{\theta_0}(\mathbf{H}_i(\theta_0))$ and $\dot{} = E_{\theta_0}(\mathbf{S}_i(\theta_0)\mathbf{S}_i(\theta_0)')$ then

$$\mathcal{H} = \dot{=} \mathcal{I}$$

- *I* is the Fisher information matrix, it is a way of measuring the amount of information that
 observable random variables (y_i, x_i) carry about θ₀.
- The relation $\mathcal{H} = := \mathcal{I}$ is called the information matrix equality.

ML ASYMPTOTIC DISTRIBUTION

- We will not go extensively through the proof since it is very similar to that of a GMM estimator.
- Check Newey and McFadden (1994) for thorough proof and B. Hansen (appendix B.11) for proof sketch.
- Under regularity conditions,

$$\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \stackrel{d}{\rightarrow} \mathcal{H}^{-1}\mathcal{N}(\mathbf{0}, \cdot) = \mathcal{N}(\mathbf{0}, \mathcal{H}^{-1} \cdot \mathcal{H}^{-1}) = \mathcal{N}(\mathbf{0}, \mathcal{I}^{-1})$$

• This is a beautiful result since it shows that the asymptotic variance of the MLE estimator is

$$\mathbf{V}_{\boldsymbol{\theta}_{ML}} = \mathcal{H}^{-1} = \boldsymbol{\cdot}^{-1} = \mathcal{I}^{-1}$$

INTUITION BEHIND $\mathbf{V}_{\theta_{ML}} = \mathcal{H}^{-1} = \mathbf{}^{-1} = \mathcal{I}^{-1}$

- Intuitively, think of asymptotic variance $V_{\theta_{ML}}$ as measure of precision of θ_{ML} . It is expected value of (scaled) square of 'average' *spread* of θ_{ML} when we have infinite data.
- Right side of information matrix equality says that the precision of $\hat{\theta}_{ML}$ depends on shape of log-likelihood function near $\hat{\theta}_{ML}$.
- If log-likelihood very *curved* or *steep* around $\hat{\theta}_{ML}$, then θ_0 will be precisely estimated. In this case, we say that we have a lot of information about θ_0 .

INTUITION BEHIND $\mathbf{V}_{\theta_{ML}} = \mathcal{H}^{-1} = \mathbf{}^{-1} = \mathcal{I}^{-1}$

- Intuitively, think of asymptotic variance $V_{\theta_{ML}}$ as measure of precision of θ_{ML} . It is expected value of (scaled) square of 'average' spread of θ_{ML} when we have infinite data.
- Right side of information matrix equality says that the precision of $\hat{\theta}_{ML}$ depends on shape of log-likelihood function near $\hat{\theta}_{ML}$.
- If log-likelihood very *curved* or *steep* around $\hat{\theta}_{ML}$, then θ_0 will be precisely estimated. In this case, we say that we have a lot of information about θ_0 .
- Steepness of log-likelihood captured by first derivative which relates to $S_i(\theta_0)$.
- Curvature of log-likelihood captured by its second derivative. Information in sample therefore also relates to minus the Hessian $-\mathbf{H}_i(\theta_0)$ (since the Hessian is negative semi-definite).
- Information matrix equality says when sample size 'plays no role' in increasing efficiency, all three measures have a simple equivalence.

$\hat{\beta}_{ML}$ and $\widehat{\mathbf{V}}_{\beta_{ML}}$ estimators for normal regression under homoskedasticity

Under homoskedasticity assumption, asymptotic distribution of MLE estimators is

$$\left[\begin{array}{c} \hat{\beta}_{ML} \\ \hat{\sigma}_{ML}^2 \end{array}\right] \quad \stackrel{d}{\rightarrow} \quad \mathcal{N}\left(\left[\begin{array}{c} \beta_0 \\ \sigma_0^2 \end{array}\right], \left[\begin{array}{c} \sigma^2(\mathbf{X}'\mathbf{X})^{-1} & \mathbf{0} \\ 0 & \frac{2}{n}\sigma^4 \end{array}\right]\right)$$

• Additive separability and normality of the errors implies here that covariance terms are 0.

Table of Contents

CRAMER-RAO LOWER BOUND AND INFORMATION MATRIX

- Cramer-Rao Lower Bound theorem shows that asymptotic variance of the MLE estimator is smallest possible within a certain class of estimators.
- Cramer-Rao Lower Bound (CRLB): If $\tilde{\theta}$ is an unbiased regression estimator of θ , then $V_{\tilde{\theta}} \ge (n\mathcal{I})^{-1}$.
- Implies asymptotic variance of the standardized estimator $\sqrt{n}(\hat{\theta}_{ML} \theta_0)$ is bounded below by \mathcal{I}^{-1} .
- $\hat{\theta}_{ML}$ is therefore asymptotically efficient.
- Efficiency bounds: CRLB is to likelihood framework what Gauss-Markov theorem (or GMM extension) is to second moment context.

CRAMER-RAO LOWER BOUND PROOF

- Define $\mathbf{S}(\theta_0) = \frac{\partial}{\partial \theta} lnf(\mathbf{y}|\mathbf{X}, \theta_0) = \sum_{i=1}^{n} \mathbf{S}_i(\theta_0)$ which from slide 14-15 we know has mean $E(\mathbf{S}(\theta_0)) = \mathbf{0}$ and variance $E(\mathbf{S}(\theta_0)\mathbf{S}'(\theta_0)) = n\mathcal{I}$.
- **(a)** We assume $\tilde{\theta}$ is unbiased estimator for any θ and writing $\tilde{\theta} = \tilde{\theta}(\mathbf{y}, \mathbf{X})$ as function of (\mathbf{y}, \mathbf{X}) we have,

$$oldsymbol{ heta} = \mathrm{E}(ilde{oldsymbol{ heta}}) = \int ilde{oldsymbol{ heta}}(\mathbf{y},\mathbf{X}) f(\mathbf{y}|\mathbf{X},oldsymbol{ heta}) d\mathbf{y}$$

CRAMER-RAO LOWER BOUND PROOF

- Define $\mathbf{S}(\theta_0) = \frac{\partial}{\partial \theta} lnf(\mathbf{y}|\mathbf{X}, \theta_0) = \sum_{i=1}^{n} \mathbf{S}_i(\theta_0)$ which from slide 14-15 we know has mean $E(\mathbf{S}(\theta_0)) = \mathbf{0}$ and variance $E(\mathbf{S}(\theta_0)\mathbf{S}'(\theta_0)) = n\mathcal{I}$.
- **(a)** We assume $\tilde{\theta}$ is unbiased estimator for any θ and writing $\tilde{\theta} = \tilde{\theta}(\mathbf{y}, \mathbf{X})$ as function of (\mathbf{y}, \mathbf{X}) we have,

$$oldsymbol{ heta} = \mathrm{E}(ilde{oldsymbol{ heta}}) = \int ilde{oldsymbol{ heta}}(\mathbf{y},\mathbf{X}) f(\mathbf{y}|\mathbf{X},oldsymbol{ heta}) d\mathbf{y}$$

③ Differentiating both sides with respect to θ' and evaluating at θ_0 ,

$$\begin{split} \mathbf{I} &= \int \tilde{\theta}(\mathbf{y}, \mathbf{X}) \frac{\partial}{\partial \theta'} f(\mathbf{y} | \mathbf{X}, \theta) d\mathbf{y} \\ &= \int \tilde{\theta}(\mathbf{y}, \mathbf{X}) \mathbf{S}(\theta_0)' f(\mathbf{y} | \mathbf{X}, \theta_0) d\mathbf{y} \quad (\text{see slide 13}) \\ &= \mathrm{E}((\tilde{\theta} - \theta_0) \mathbf{S}(\theta_0)') \quad (\text{from } \mathrm{E}(\theta_0 \mathbf{S}(\theta_0)) = \theta_0 \mathrm{E}(\mathbf{S}(\theta_0)) = \mathbf{0}) \end{split}$$

CRAMER-RAO LOWER BOUND PROOF

- **O** Define $\mathbf{S}(\theta_0) = \frac{\partial}{\partial \theta} lnf(\mathbf{y}|\mathbf{X}, \theta_0) = \sum_{i=1}^{n} \mathbf{S}_i(\theta_0)$ which from slide 14-15 we know has mean $E(\mathbf{S}(\theta_0)) = \mathbf{0}$ and variance $E(\mathbf{S}(\theta_0)\mathbf{S}'(\theta_0)) = n\mathcal{I}$.
- **a** We assume $\tilde{\theta}$ is unbiased estimator for any θ and writing $\tilde{\theta} = \tilde{\theta}(\mathbf{y}, \mathbf{X})$ as function of (\mathbf{y}, \mathbf{X}) we have.

$$oldsymbol{ heta} = \mathrm{E}(ilde{oldsymbol{ heta}}) = \int ilde{oldsymbol{ heta}}(\mathbf{y},\mathbf{X}) f(\mathbf{y}|\mathbf{X},oldsymbol{ heta}) d\mathbf{y}$$

Differentiating both sides with respect to θ' and evaluating at θ_0 . 3

$$\begin{split} \mathbf{I} &= \int \tilde{\theta}(\mathbf{y}, \mathbf{X}) \frac{\partial}{\partial \theta'} f(\mathbf{y} | \mathbf{X}, \theta) d\mathbf{y} \\ &= \int \tilde{\theta}(\mathbf{y}, \mathbf{X}) \mathbf{S}(\theta_0)' f(\mathbf{y} | \mathbf{X}, \theta_0) d\mathbf{y} \quad (\text{see slide 13}) \\ &= \mathrm{E}((\tilde{\theta} - \theta_0) \mathbf{S}(\theta_0)') \quad (\text{from } \mathrm{E}(\theta_0 \mathbf{S}(\theta_0)) = \theta_0 \mathrm{E}(\mathbf{S}(\theta_0)) = \mathbf{0}) \end{split}$$

Last, by Cauchy-Schwartz inequality for matrices we have,

$$\begin{split} \mathbf{V}_{\tilde{\boldsymbol{\theta}}} &= \mathrm{E}((\tilde{\boldsymbol{\theta}} - \theta_0)(\tilde{\boldsymbol{\theta}} - \theta_0)') \\ &\geq \mathrm{E}((\tilde{\boldsymbol{\theta}} - \theta_0)\mathbf{S}(\theta_0)')\mathrm{E}(\mathbf{S}(\theta_0)\mathbf{S}(\theta_0)')^{-1}\mathrm{E}(\mathbf{S}(\theta_0)(\tilde{\boldsymbol{\theta}} - \theta_0)') \\ &= \mathrm{E}(\mathbf{S}(\theta_0)\mathbf{S}(\theta_0)')^{-1} \\ &= \mathrm{E}(\mathbf{S}(\theta$$

M Kummer

CRAMER-RAO LOWER BOUND AND INVARIANCE PRINCIPLE

- Finally, consider functions of parameters. If $\psi = \mathbf{g}(\theta)$ then the MLE of $\hat{\psi}_{ML}$ is $\mathbf{g}(\hat{\theta}_{ML})$.
- This is because maximization of objective function is unaffected by parameterization and transformation.

$$\sqrt{n}(\hat{\psi}_{ML} - \psi) \simeq \mathbf{G}_0 \sqrt{n}(\hat{\theta}_{ML} - \theta_0) \quad \stackrel{d}{\rightarrow} \quad \mathcal{H}^{-1} \mathcal{N}(\mathbf{0}, \mathbf{G}_0' \mathcal{I}^{-1} \mathbf{G}_0)$$

with $\mathbf{G}_0 = \frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{g}(\boldsymbol{\theta}_0)$.

• A useful extension is that Cramer-Rao lower bound for $\psi = \mathbf{g}(\hat{\theta})$ is $\mathbf{G}_0' \mathcal{I}^{-1} \mathbf{G}_0$, and the MLE $\hat{\psi}_{ML} = \mathbf{g}(\hat{\theta}_{ML})$ is asymptotically efficient.