

YouTube lecture 10/10/16

Lecture notes 9 - Asymptotic theory (cont.)

* Test material for test II covered here *

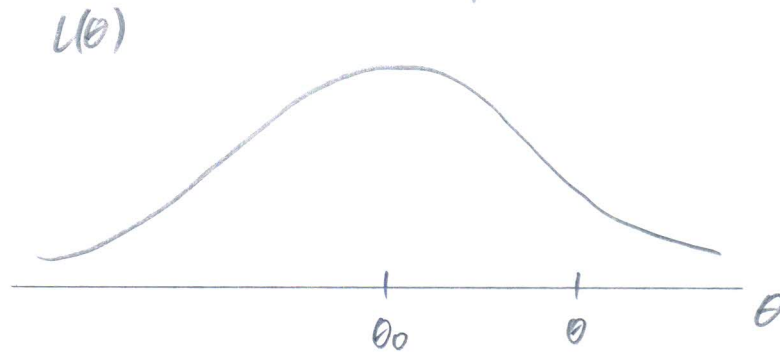
- can bring cheatsheets (1 side)

3 question test

convergence, delta
point estimation
sufficiencyconsistency
of MLE: $\hat{\theta}_{n, MLE} \xrightarrow{P} \theta_0$

(*) Brief review (*)

$$P\left(\frac{L(\theta_0)}{L(\theta)} > 1\right) \rightarrow 1$$

as $n \rightarrow \infty$ θ_0 - true param
value.

- consistency: estimator converges in probability to true parameter
- in general, prove via definition (convergence in prob); or via quadratic mean convergence; or on a case-by-case basis.
- Above result was an attempt to establish whether MLE is, in general, consistent.
- Above result shows that $P(L(\theta_0) > L(\theta)) \rightarrow 1$ i.e. probability of likelihood evaluated at true parameter value θ_0 being greater than likelihood evaluated at an arbitrary param value approaches 1 in the limit as $n \rightarrow \infty$ (ie. more samples collected.)

(*) That is close to; but not quite a proof of the consistency of MLE; again due to the distinction between pointwise and uniform convergence

(*) Consistency of MLE would involve showing that the maximum of the likelihood at points other than θ_0 is always smaller than that of $L(\theta_0)$ (?) - tighter expl. (A1)
in terms of pointwise and uniform convergence.

(*) Proof of MLE consistency in general \rightarrow beyond course scope.

(*) (A1): In general MLE is consistent under regularity conditions:

- 1) dimension of θ is fixed i.e. $\theta \in \mathbb{R}^d$ with a fixed
- 2) $p(x; \theta)$ is a smooth function of θ .
- 3) Identifiability of parameters. - intuitively; can't have two different values of parameter refer to the same distribution

i.e. $\theta_1 \neq \theta_2$

$p(x; \theta_1) \neq p(x; \theta_2)$ where "not equal" means "not equal" in distribution

(*) Practical takeaway: MLE is an effective estimation procedure when there is a large amount of data (large n) and small no. of parameters

0

(*) nonparametric / high dimensional settings \rightarrow this is not the case.

(*) MLE in general is also asymptotically normal

6. Asymptotic Normality of MLE

- we will prove that MLE satisfies:-

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$$

- or informally,

$$\hat{\theta}_n \hat{=} N\left(\theta, \frac{1}{nI(\theta)}\right)$$

- 2 definitions:- score function and Fisher information.

score function: $S_n(\theta) \equiv S_n(\theta, x_1, x_2, \dots, x_n) = t'(\theta) = \frac{\partial \log p(x_1, \dots, x_n; \theta)}{\partial \theta}$

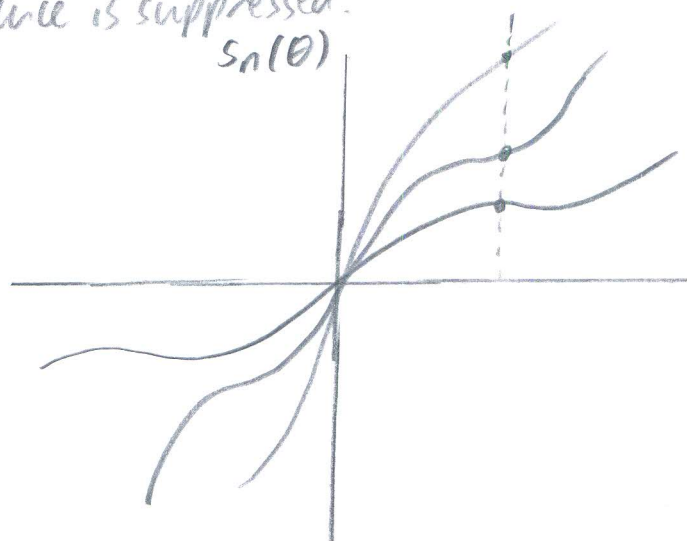
- IID data:- $S_n(\theta) = \sum_{i=1}^n \frac{\partial \log p(x_i; \theta)}{\partial \theta}$

- $S_n(\theta) = S_n(\theta, x_1, \dots, x_n)$ is a function of data and parameters; although formal dependence is suppressed.

- $S_n^{(1)}(\theta)$

- $S_n^{(2)}(\theta)$

- $S_n^{(3)}(\theta)$



- score function $S_n(\theta)$ is a function of θ .

- But it is a random function that depends on dataset x_1, \dots, x_n .

- Graphic shows score function for 3 datasets of size n .

- $S_n^{(i)}(\theta) \quad i=1, \dots, 3$

(e.g. through simulation)

- we can then also conceive of mean and variance of the score function.

Fisher Information

(*) The Fisher Information is defined to be:-

$$I_n(\theta) = \text{Var}_\theta(S_n(\theta))$$

(information is defined as
variance of the score function)

(*) In context of graphics, it can help to visualise this generatively. In that we can think of computing the variance of the score function over the 3 datasets we have at a particular value of θ . Doing this for all values of θ yields the Fisher information (i.e. a slice of graph)

• Also note that:-

$$\text{Var}(\hat{\theta}_{\text{MLE}}) \approx \frac{1}{I_n(\theta)}$$

- Under MLE; variance of estimator is approx equal to Fisher info.

lw: called 'information' \rightarrow if $I_n(\theta)$ is large, there is a lot of information; and $\text{Var}(\hat{\theta}_{\text{MLE}})$ is small, i.e. our estimator is precise.

lw: Provides theoretical understanding of MLE; but also a means of computing standard errors (standard dev of estimators).

Theorem 7

- visualise by looking at graph

- under regularity conditions; expected value of score fn is 0

$$\mathbb{E}_\theta[S_n(\theta)] = 0$$

$$\mathbb{E}_\theta[S_n(\theta)] = \int \dots \int \left(\frac{\partial \log p(x_1, \dots, x_n; \theta)}{\partial \theta} \right) p(x_1, \dots, x_n; \theta) dx_1 \dots dx_n = 0$$

(*) Integral/exp. is wrt to the joint distribution of the data $p(x^n; \theta)$

- noting that $p(x^n; \theta) = \prod_{n=1}^N p(x_n; \theta)$

(*) We assume that the model is correct and that the data is generated from a joint distribution under a particular value of θ .

Proof:

$$\begin{aligned} \mathbb{E}_\theta[S_n(\theta)] &= \int \dots \int \frac{\partial \log p(x_1, \dots, x_n; \theta)}{\partial \theta} p(x_1, \dots, x_n; \theta) dx_1 \dots dx_n \\ &= \int \dots \int \frac{\partial}{\partial \theta} \log p \cdot \frac{\partial p}{\partial \theta} \cdot p dx_1, \dots, dx_n \\ &= \int \dots \int \frac{1}{p(x_1, \dots, x_n; \theta)} \cdot \frac{\partial p(x_1, \dots, x_n; \theta)}{\partial \theta} \cdot p(x_1, \dots, x_n; \theta) dx_1 \dots dx_n \\ &= \int \dots \int \frac{\partial}{\partial \theta} p(x_1, \dots, x_n; \theta) dx_1 \dots dx_n \end{aligned}$$

(*) regularity conditions allow us to exchange order of operators.

(*) we require that the set over which density is defined to not depend on

θ .
e.g. Normal is defined over \mathbb{R} ; independent of params μ, σ^2
(univ.)

$X \sim \text{unif}(0, \theta) \rightarrow$ range of random variable X depends on the parameter θ ; violates condition.

(*) Notes:- If support of p depends on θ ; then $\int \dots \int$ and $\frac{\partial}{\partial \theta}$ cannot be switched.

Hence, we have from above:-

$$\mathbb{E}_\theta[S_n(\theta)] = \frac{\partial}{\partial \theta} \underbrace{\int \dots \int p(x_1, \dots, x_n; \theta) dx_1 \dots dx_n}_{=1}$$

$$= 0.$$



(*) some nuances from notes:-

- If the expectation value is taken at the same θ as we evaluate $S_n(\theta)$; then $\mathbb{E}_\theta[S_n(\theta)] = 0$

- When the θ s mismatch, i.e. $\theta_1 \neq \theta_2$; then

$$\mathbb{E}_{\theta_1}[S_n(\theta_2)] \neq 0.$$

(*) Example 8

- Let $X_1, \dots, X_n \sim N(\theta, 1)$
- $L(\theta) \propto e^{-\frac{1}{2} \sum_{i=1}^n (X_i - \theta)^2}$
- $\ell(\theta) = -\frac{1}{2} \sum_{i=1}^n (X_i - \theta)^2$

$$S_n(\theta) = \sum_{i=1}^n (X_i - \theta) = n(\bar{X} - \theta) \quad (\text{multiplying and dividing by } n)$$

- Note $\mathbb{E}_{\theta}[S_n(\theta)] = 0$

- $S_n(\theta)$ is a random function.
- $S_n(\theta) = f(\theta, x_1, \dots, x_n)$
- different datasets (i.e. different values of \bar{X}) will shift $S_n(\theta)$

(*) Further properties of Fisher information

- As Fisher info := $I_n = \text{Var}(S_n(\theta))$ and $\mathbb{E}_{\theta}[S_n(\theta)] = 0$; we have

$$I_n(\theta) = \mathbb{E}_{\theta}[S_n^2(\theta)]$$

Lemma 9

- iid case
- $I_n(\theta) = nI(\theta)$ where $I(\theta)$ is F.I for $n=1$ (via definition of $\ell(\theta)$ as sum for iid).

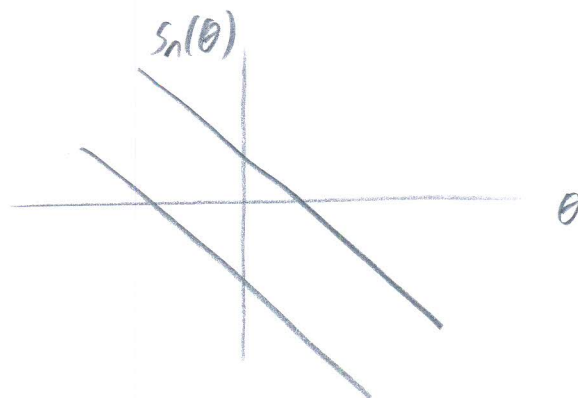
- LW: Helps simplify calculations (additive charact. of information).

- LW: follows since log-likelihood, and hence score, is the sum of n -independent terms.

(*) Lemma 10

- under reg. conditions:-

$$I_n(\theta) = -\mathbb{E}_{\theta}\left[\frac{\partial^2}{\partial \theta^2} \ell_n(\theta)\right]$$



(*) - Take log likelihood, 2nd deriv, negative; provides another means of computing $I_n(\theta)$; in particular when combined with lemma 9.

Proof Take $n=1$

- Note: - $\int p = 1 \Rightarrow \int p' = 0 \Rightarrow \int p'' = 0 \Rightarrow \int \frac{p''}{p} p = 0$

$$\Rightarrow \mathbb{E}\left(\frac{p''}{p}\right) = 0$$

Let $\ell = \log p$ and $S = \ell' = \frac{p'}{p}$; then $\ell'' = \left(\frac{p''}{p}\right) - \left(\frac{p'}{p}\right)^2$

(*) Above identities ~~from~~^{rely on} reg. conditions.

$$\begin{aligned} \text{Var}(S) &= \mathbb{E}[S^2] - \underbrace{(\mathbb{E}[S])^2}_{=0} = \mathbb{E}[S^2] = \mathbb{E}\left[\left(\frac{p'}{p}\right)^2\right] = \mathbb{E}\left[\left(\frac{p'}{p}\right)^2\right] - \underbrace{\mathbb{E}\left[\left(\frac{p''}{p}\right)\right]}_{=0} \\ &= -\mathbb{E}\left[\left(\frac{p''}{p}\right) - \left(\frac{p'}{p}\right)^2\right] = -\mathbb{E}(\ell'') \end{aligned}$$

W: Score fn is derivative of likelihood

Fisher information is variance of score function; or above form.

(*) extension to vectors (cosmetic)

$$\theta = (\theta_1, \dots, \theta_K)$$

- $\ln(\theta)$ and $\ell_n(\theta)$ do not change

- But score $S_n(\theta)$ is affected \rightarrow it is derivative of $\ell(\theta)$ for univariate case.

- score function is now a vector of functions: -

$$S_n(\theta) = \begin{bmatrix} \frac{\partial \ell_n(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell_n(\theta)}{\partial \theta_K} \end{bmatrix} \quad S_n(\theta) \in \mathbb{R}^K \quad (S_n(\theta))_k = \frac{\partial \ell_n(\theta)}{\partial \theta_k}$$

- Fisher information is a variance-covariance matrix of the score vector.

$$I_n(\theta) \in \mathbb{R}^{K \times K}$$

$$I_n(\theta) = \text{cov}(S_n(\theta)) \quad I_n(\theta) = -E_\theta \left[\frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} \right]$$

(*) Example

$$X_1, \dots, X_n \sim \text{Bern}(p)$$

likelihood: $\ell(p) = p^S (1-p)^{n-S}$ $S = \sum_{n=1}^N X_n$

log-likelihood: $S \log p + (n-S) \log(1-p) = \ell(p)$

score: $S_n(p) = \frac{\partial \ell(p)}{\partial p} = \frac{S}{p} - \frac{n-S}{1-p}$

2nd deriv $\ell(p)$: $\ell''(p) = -\frac{S}{p^2} - \frac{n-S}{(1-p)^2}$

Fisher info: $I_n(p) = -E[\ell''(p)]$

$$= \frac{np}{p^2} + \frac{n-np}{(1-p)^2} = \frac{n}{p} + \frac{n}{(1-p)} = \frac{n}{p(1-p)}$$

W: Score function is

a random function (of data)

Fisher info is just a function of the parameter

(*) Example II

- suppose $X_1, \dots, X_n \sim N(\mu, \gamma)$ $\gamma = \sigma^2$ (?)

(A1) verify this yourself as part of review:-

$$S_n(\mu, \gamma) = \begin{bmatrix} \frac{1}{\gamma} \sum (x_i - \mu) \\ -\frac{n}{2\gamma} + \frac{1}{2\gamma^2} \sum (x_i - \mu)^2 \end{bmatrix} \quad I_n(\mu, \gamma) = \begin{bmatrix} \frac{n}{\gamma} & 0 \\ 0 & \frac{n}{2\gamma^2} \end{bmatrix}$$

(*) score and F.I. cannot be computed without regularity conditions.

6. Asymptotic Normality of MLE

we will prove this Theorem 12

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$$

$$\hat{\theta} \approx N\left(\theta, \frac{1}{nI(\theta)}\right) = N\left(\theta, \frac{1}{I_n(\theta)}\right)$$

we will also prove the following for MLE estimators:-

$$\hat{\theta} = \theta + \frac{1}{n} \sum_{i=1}^n \psi^*(x_i) + o_p(n^{-1/2}) \quad \text{where } \psi^*(x) = \frac{s(\theta, x_i)}{I(\theta)} \quad (1)$$

where ψ^* is an influence function

any well-behaved estimator can be written as (1) for some ψ and that $\text{Var}(\psi) \geq \text{Var}(\psi^*)$ later

(*) influence function with smallest variance gives MLE. (later)

the computation of confidence intervals relies on asymptotic normality

Proof of Theorem 12:

$\hat{\theta}_{MLE}$ (MLE est.)

we know:-

$$t'(\hat{\theta}_{MLE}) = 0$$

expand $t'(\hat{\theta}_{MLE})$ about true value of param:-

$$0 = t'(\hat{\theta}_{MLE}) = t'(\theta) + (\hat{\theta} - \theta)t''(\theta) + \dots$$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta) \approx \frac{\frac{1}{\sqrt{n}} t'(\theta)}{-\frac{1}{\sqrt{n}} t''(\theta)} = \frac{A}{B}$$

$$A = \frac{1}{\sqrt{n}} t'(\theta) = \sqrt{n} \times \frac{1}{n} \sum_{i=1}^n s(\theta, x_i) = \sqrt{n}(\bar{s} - 0)$$

$$A \xrightarrow{d} N(0, I(\theta))$$

② ?
don't understand
the reasoning
for this

- And also

$$B \xrightarrow{P} -E(l'') = I(\theta)$$

- via Slutsky:-

$$\frac{A}{B} \xrightarrow{d} \frac{\sqrt{I(\theta)} z}{I(\theta)} = \frac{z}{\sqrt{I(\theta)}} = N\left(0, \frac{1}{I(\theta)}\right)$$

so:-

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{1}{I(\theta)}\right)$$

②
(A2) - don't undigest
this argument; review.
B \xrightarrow{P} (need to
move on).

□