

4. Preparations for the Asymptotic

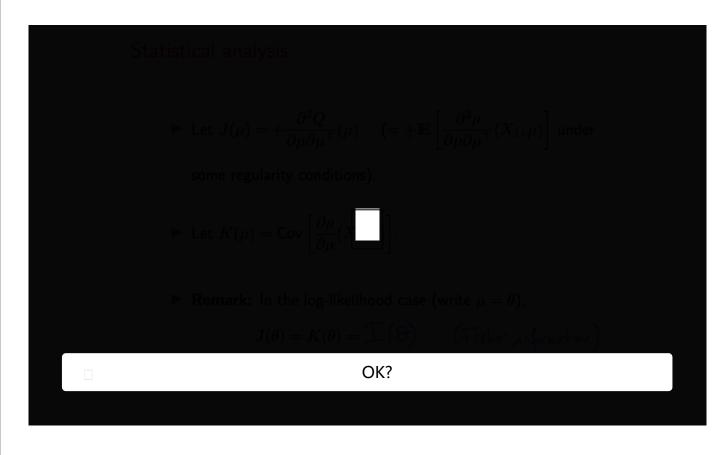
课程 □ Unit 3 Methods of Estimation □ Lecture 12: M-Estimation □ Normality of M-estimators

4. Preparations for the Asymptotic Normality of M-estimators Video note:

In the video below and on the slides, we have used another common notation for the gradient and the Hessian. Below the video, we repeat the definitions of the matrices ${f J}$ and ${f K}$ in the familiar notation we had used earlier.

Also about 4:20, Prof Rigollet said "What is the expectation of the second derivative of the log likelihood minus the expectation of the second derivative of the log likelihood". What he meant was in fact "minus the expectation of the second derivative of the log liklihood", not the difference between this quantity and itself.

Preparation for Asymptotic Normality



One says that it's negative the expectation of the second derivative, and we also know that it's equal to the variance of the first derivative, which is precisely how we describe

K to be.

OK, so in this case, those two guys are actually equal to I of theta. So this is my Fisher information.

OK?

5:01 / 5:01 □ 1.0x End of transcript. Skip to the start.

下载视频文件

下载 SubRip (.srt) file

下载 Text (.txt) file

The **J** and **K** matrices:

Let $\mathbf{X}_1,\ldots,\mathbf{X}_n$ be i.i.d. random vector in \mathbb{R}^k with some unknown distribution \mathbf{P} with some associated parameter $\vec{\mu}^*\in\mathbb{R}^d$ on some sample space E. Let $\mathcal{Q}(\vec{\mu}) = \mathbb{E}\left[\rho\left(\mathbf{X},\vec{\mu}\right)\right]$ for some function $\rho: E \times \mathcal{M} \to \mathbb{R}$, where \mathcal{M} is the set of all possible values of the unknown true parameter $\vec{\mu}^*$.

Then the matrices ${\bf J}$ and ${\bf K}$ are defined as

$$\mathbf{J} = \mathbb{E}\left[\mathbf{H}
ho
ight] \;\; = \;\; \mathbb{E}\left[egin{pmatrix} rac{\partial^2
ho}{\partial\mu_1\partial\mu_1}(\mathbf{X}_1,ec{\mu}) & \ldots & rac{\partial^2
ho}{\partial\mu_1\partial\mu_d}(\mathbf{X}_1,ec{\mu}) \ dots & \ddots & dots \ rac{\partial^2
ho}{\partial\mu_d\partial\mu_1}(\mathbf{X}_1,ec{\mu}) & \ldots & rac{\partial^2
ho}{\partial\mu_d\partial\mu_d}(\mathbf{X}_1,ec{\mu}) \end{pmatrix}
ight] \qquad (d imes d)$$

$$\mathbf{K} = \mathsf{Cov}\left[
abla
ho\left(\mathbf{X}_1,ec{\mu}
ight)
ight] \quad = \quad \mathsf{Cov}\left[\left(egin{array}{c} rac{\partial
ho}{\partial\mu_1}(\mathbf{X}_1,ec{\mu}) \ dots \ rac{\partial
ho}{\partial\mu_d}(\mathbf{X}_1,ec{\mu}) \end{array}
ight)
ight] \qquad (d imes d)\,.$$

In one dimension, i.e. d=1, the matrices reduce to the following:

$$J\left(\mu
ight) \;\;\; = \mathbb{E}\left[rac{\partial^2
ho}{\partial\mu^2}(X_1,\mu)
ight]$$

$$K\left(\mu
ight) = \mathrm{Var}\left[rac{\partial
ho}{\partial\mu}(X_1,\mu)
ight]$$

Concept Check: M-estimators vs. Maximum Likelihood Estimation

1/1 point (graded)

Let ho denote a loss function, and let $X_1,\ldots,X_n\stackrel{iid}{\sim} \mathbf{P}$. Let $\widehat{\mu}$ denote the M-estimator for some unknown parameter $\mu^*= \mathop{\mathrm{argmin}}_{\mu \in \mathbb{R}} \mathbb{E}\left[\rho\left(X_1,\mu\right) \right] \in \mathbb{R}$ associated with \mathbf{P} . (Here we are assuming that μ^* is a one-dimensional parameter.)

Consider the following functions

$$J\left(\mu
ight) \; = \mathbb{E}\left[rac{\partial^2
ho}{\partial\mu^2}(X_1,\mu)
ight]$$

$$K\left(\mu
ight) \ = \mathrm{Var}\left[rac{\partial
ho}{\partial\mu}(X_1,\mu)
ight]$$

Which of the following statements are true? (Choose all that apply.)

- \square It is always true that $J\left(\mu
 ight)=K\left(\mu
 ight)$.
- $I(\mu)=K(\mu)$ when ho is the negative log-likelihood– in this case, both of these functions are equal to the Fisher information. \Box
- ullet Under some technical conditions, the functions $J(\mu)$ and $K(\mu)$ determine the asymptotic variance of the M-estimator $\widehat{\mu}$. \Box

Solution:

- The response "It is always true that $J(\mu) = K(\mu)$." is incorrect. In general, the functions $J(\mu)$ and $K(\mu)$ will not be equal to each other. For example, if the loss function is given in terms of Huber's loss (as we will see later in this lecture), $J(\mu) \neq K(\mu)$.
- The choice " $J(\mu) = K(\mu)$ when ρ is the negative log-likelihood– in this case, both of these functions are equal to the Fisher information." is correct. In the special case where $\rho(x,\mu)$ is defined to be the negative log-likelihood of the statistical model, then it is true that $J(\mu) = K(\mu)$. This was derived in <u>Lecture 11</u>.
- The choice "Under some technical conditions, the functions $J(\mu)$ and $K(\mu)$ determine the asymptotic variance of the M-estimator $\widehat{\mu}$." is correct. This is content of the theorem on the slide "Asymptotic Normality," which shows that the asymptotic variance of $\widehat{\mu}_n$, assuming some hypotheses, is given by $J(\mu^*)^{-1}K(\mu^*)J(\mu^*)^{-1}$.

Remark on signs:

Let us match the signs in the definition of ${f J}$ and ${f K}$ with those in the definition of Fisher information. For maximum likelihood estimation,

$$\rho_n\left(\theta\right) := \rho\left(\mathbf{X}_1, \ldots, \mathbf{X}_n, \theta\right) = -\ell_n\left(\theta\right) \quad \text{where } \ell_n\left(\theta\right) = \ln L_n\left(\mathbf{X}_1, \ldots, \mathbf{X}_n, \theta\right).$$

For this particular loss function $\, oldsymbol{
ho}, \,$ the $\, {f J} \,$ and $\, {f K} \,$ matrices are

$$\mathbf{J} = \mathbb{E}\left[\mathbf{H}\rho_{1}\left(\theta\right)\right] = -\mathbb{E}\left[\mathbf{H}\ell_{1}\left(\theta\right)\right]$$

$$\mathbf{K} = \operatorname{Cov}\left[\nabla\rho_{1}\left(\theta\right)\right] = \operatorname{Cov}\left[-\nabla\ell_{1}\left(\theta\right)\right] = \operatorname{Cov}\left[\nabla\ell_{1}\left(\theta\right)\right] \qquad \left(\operatorname{Cov}\left[\mathbf{Y}\right] = \operatorname{Cov}\left[-\mathbf{Y}\right] \text{ for any random vector } \mathbf{Y}.$$

Both of these matrices equals the Fisher information matrix.

提交	你已经尝试了1次(总共可以尝试2次)	
□ Answers are displayed within the problem		
讨论		显示讨论
主题: Unit 3 Methods of Estimation:Lecture 12: M-Estimation / 4. Preparations for the Asymptotic Normality of M-estimators		<u> </u>

© 保留所有权利