

Handout 17: Asymptotic behavior of the posterior distribution

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim: We examine the properties of the posterior distribution $\Pi(\theta|y)$, under different sets of conditions, as the number of observations n increases $n \rightarrow \infty$, as well as their implications in inference.

References:

- Ferguson, T. S. (1996, Section 21). A course in large sample theory. Chapman and Hall/CRC.
- Chen, C. F. (1985). On asymptotic normality of limiting density functions with Bayesian implications. Journal of the Royal Statistical Society: Series B (Methodological), 47(3), 540-546.
- Van der Vaart, A. W. (2000, Chapter 10). Asymptotic statistics. Cambridge series in statistical and probabilistic mathematics.

Web-applets

- https://georgios-stats-1.shinyapps.io/demo_conjugatepriors/
- https://georgios-stats-1.shinyapps.io/demo_conjugatejeffreyslplacepriors/
- https://georgios-stats-1.shinyapps.io/demo_mixturepriors/

What is about?

Notation 1. Consider the Bayesian model $(F(x_{1:n}|\theta), \Pi(\theta))$ as

$$\begin{cases} x_{1:n}|\theta & \sim F(\cdot|\theta) \\ \theta & \sim \Pi(\cdot) \end{cases} \quad (1)$$

where a sequence of observables $x_{1:n} = (x_1, \dots, x_n)$ are drawn from the parametric model $F(\cdot|\theta)$ admitting a pdf/pmf $f(\cdot|\theta)$ with unknown parameter $\theta \in \Theta$. The prior $\Pi(\cdot)$ of θ admits pdf/pmf $\pi(\cdot)$.

Question 2. We study the behavior of the posterior distribution $\Pi(\theta|x_{1:n})$ with respect to the number of observables n , first when θ is a discrete parameter, and then when θ is a continuous one.

Note 3. All the theorems in this chapter are frequentist in character, namely we study the posterior laws under the assumption that the observables $x_{1:n}$ is a random sample from the sampling distribution $F(\cdot|\theta^*)$ for some fixed non-random true value $\theta^* \in \Theta$.

Notation 4. We denote the likelihood function as $L_n(\theta) := f(x_{1:n}|\theta)$ and the posterior distribution $\Pi_n(\theta)$ with pdf/pmf $\pi_n(\theta) := \pi(\theta|x_{1:n}) \propto L_n(\theta)\pi(\theta)$, to ease the notation and make clear Note 3.

1 Discrete θ : Asymptotic consistency

Note 5. Given the Bayesian model (1), we consider cases where $\theta \in \Theta$ is a discrete parameter, and Θ is a countable space.

Note 6. The theorem below implies that, if Θ is countable, under conditions, the posterior distribution function of $\theta \in \Theta$ ultimately degenerates to a step function with a single (unit) step at $\theta = \theta^*$, where θ^* is the true value of the unknown discrete parameter θ .

Theorem 7. Assume the Bayesian model (1), let $x_{1:n} = (x_1, \dots, x_n)$ be a sequence of IID observables, $\theta \in \Theta$ be the unknown parameter with prior distribution mass $\pi(\theta)$, and posterior distribution mass $\pi_n(\theta)$, where Θ is a countable parametric space. Suppose $\theta^* \in \Theta$ is the (only) true value of θ such that $\pi(\theta^*) > 0$, and $-KL(f(\cdot|\theta^*)||f(\cdot|\theta)) := \int \log \frac{f(x|\theta)}{f(x|\theta^*)} dF(x|\theta^*) < 0$ for all $\theta \neq \theta^*$. Then

$$\lim_{n \rightarrow \infty} \pi_n(\theta) = \begin{cases} 1 & , \theta = \theta^* \\ 0 & , \theta \neq \theta^* \end{cases}.$$

Proof. Due to exchangeability of $x_{1:n}$, it is

$$\pi_n(\theta) = \frac{\frac{L_n(\theta)}{L_n(\theta^*)} \pi(\theta)}{\sum_{\forall \theta \in \Theta} \frac{L_n(\theta)}{L_n(\theta^*)} \pi(\theta)} = \frac{\exp\left(\sum_{i=1}^n \log \frac{f(x_i|\theta)}{f(x_i|\theta^*)}\right) \pi(\theta)}{\sum_{\forall \theta \in \Theta} \exp\left(\sum_{i=1}^n \log \frac{f(x_i|\theta)}{f(x_i|\theta^*)}\right) \pi(\theta)} = \frac{\exp(S_n(\theta)) \pi(\theta)}{\sum_{\forall \theta \in \Theta} \exp(S_n(\theta)) \pi(\theta)}$$

where $S_n(\theta) = \sum_{i=1}^n \log \frac{f(x_i|\theta)}{f(x_i|\theta^*)}$. From the SLLN, as $n \rightarrow \infty$, it is

$$\lim_{n \rightarrow \infty} \frac{1}{n} S_n(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i|\theta)}{f(x_i|\theta^*)} = E_F \left(\log \frac{f(x|\theta)}{f(x|\theta^*)} \middle| \theta^* \right), \quad \text{a.s.} \quad (2)$$

By using Jensen's inequality and the fact that log is concave, it is

$$E_{x \sim F(\cdot|\theta^*)} \left(\log \frac{f(x|\theta)}{f(x|\theta^*)} \right) \leq \log E_{x \sim F(\cdot|\theta^*)} \left(\frac{f(x|\theta)}{f(x|\theta^*)} \right) = \log(1) = 0 \implies E_{x \sim F(\cdot|\theta^*)} \left(\log \frac{f(x|\theta)}{f(x|\theta^*)} \right) \leq 0 \quad (3)$$

In (3), the equality holds for $\theta = \theta^*$ a.s., and the inequality holds for $\theta \neq \theta^*$ a.s., since Θ is a countable space and $\theta^* \in \Theta$, θ^* is "distinguishable" from the others, according to Theorem 39. Notice that, for any $\theta \neq \theta^*$, (2) and (3) imply that

$$\lim_{n \rightarrow \infty} \frac{1}{n} S_n(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i|\theta)}{f(x_i|\theta^*)} < 0, \quad \text{a.s.}$$

which implies that

$$\lim_{n \rightarrow \infty} S_n(\theta) = \lim_{n \rightarrow \infty} n \frac{1}{n} S_n(\theta) = -\infty, \quad \text{as}$$

Therefore,

- for any $\theta \neq \theta^*$, it is

$$\lim_{n \rightarrow \infty} \pi_n(\theta) = \lim_{n \rightarrow \infty} \frac{\exp(S_n(\theta)) \pi(\theta)}{\sum_{\forall \theta \in \Theta} \exp(S_n(\theta)) \pi(\theta)} = 0, \quad \text{a.s.}$$

- for $\theta = \theta^*$, it is

$$\lim_{n \rightarrow \infty} \pi_n(\theta^*) = 1 - \sum_{\forall \theta \neq \theta^*} \underbrace{\lim_{n \rightarrow \infty} \pi_n(\theta)}_{=0, \text{ for } \theta \neq \theta^*} = 1, \quad \text{a.s.}$$

□

Remark 8. Theorem 7 relies on the condition that the true parameter value θ^* is unique. If there was another θ^{**} such that $f(x|\theta^{**}) = f(x|\theta^*)$, we would observe IID data when θ equaled θ^* or θ^{**} , and hence the data could not discriminate between the two values.

Fact 9. It can be shown that if $\theta^* \notin \Theta$, the posterior degenerates onto the value in Θ which gives the parametric model closest θ^* .

2 Continuous θ : Asymptotic consistency and normality under Cramer's conditions

Note 10. Given the Bayesian model (1), we consider cases that $\theta \in \Theta$ is a continuous parameter, that $\Theta \subset \mathbb{R}^k$ is compact with $k \geq 1$, and that observables $\{x_i\}$ are IID.

Note 11. We show that when θ is continuous and under regularity conditions :

1. The posterior PDF of θ becomes more and more concentrated above an area around the true value θ^* as data size increases beyond a number $n \rightarrow \infty$.
2. The limiting posterior distribution of θ is close to a normal density $N\left(\theta|\hat{\theta}_n, \frac{1}{n}\mathcal{J}(\theta^*)^{-1}\right)$ centered at $\hat{\theta}_n$ (the MLE of (19)), with variance $\frac{1}{n}\mathcal{J}(\theta^*)^{-1}$. Here, $\mathcal{J}(\cdot)$ is the Fisher information, where

$$\mathcal{J}(\theta) = E_{x \sim F(\cdot|\theta)} \left((\nabla_{\theta} \log f(x|\theta))^{\top} (\nabla_{\theta} \log f(x|\theta)) \right) = -E_{x \sim F(\cdot|\theta)} \left(\nabla_{\theta}^2 \log f(x|\theta) \right).$$

3. These conclusions do not depend on the choice of the prior distribution provided that $\pi(\theta^*) > 0$.

Remark 12. The following version of the theorem, by Le Cam (1953), equivalently states that the posterior PDF of (the linear transformation) $\vartheta = \sqrt{n}(\theta - \hat{\theta}_n)$

$$\pi_n(\vartheta) = \frac{L_n(\vartheta \frac{1}{\sqrt{n}} + \hat{\theta}_n) \pi(\vartheta \frac{1}{\sqrt{n}} + \hat{\theta}_n)}{\int L_n(\vartheta \frac{1}{\sqrt{n}} + \hat{\theta}_n) \pi(\vartheta \frac{1}{\sqrt{n}} + \hat{\theta}_n) d\vartheta}$$

approaches the PDF of $N(0, \mathcal{J}(\theta^*)^{-1})$ as $n \rightarrow \infty$.

Condition 13. (Cramer conditions) Consider the following regular conditions.

d1 Θ is an open subset of \mathbb{R}^k

d2 second partial derivatives of $f(x|\theta)$ with respect to θ exist and are continuous for all x , and may be passed under the integral operator in $\int f(x|\theta) dx$

d3 there is a function $K(x)$ such that $E_{x \sim F(x|\theta^*)}(K(x)) < \infty$ and each component of $\nabla_{\theta}^2 \log(f(x|\theta))$ is bounded in absolute value by $K(x)$ uniformly in some neighborhood of θ^*

d4 $\mathcal{J}(\theta^*) = -E_{x \sim F(\cdot|\theta^*)} (\nabla_{\theta^*}^2 \log f(x|\theta^*))$ is positive definite

d5 (identifiability) $f(x|\theta) = f(x|\theta^*)$ a.s. then $\theta = \theta^*$

Theorem 14. (Bernstein-von Mises) Let x_1, x_2, \dots be IID random variables drawn from a sample distribution with density $f(x|\theta)$, $\theta \in \Theta$, and let $\theta^* \in \Theta$ denote the true value of θ . Let $L_n(\theta) = f(x_{1:n}|\theta)$ denote the likelihood. Assume that the prior density $\pi(\theta)$ is continuous and $\pi(\theta) > 0$ for all $\theta \in \Theta$. Assume Conditions 13 hold. Then it is

$$\frac{L_n\left(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}\right)}{L_n(\hat{\theta}_n)} \pi\left(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}\right) \xrightarrow{a.s.} \exp\left(-\frac{1}{2} \vartheta^{\top} \mathcal{J}(\theta_0) \vartheta\right) \pi(\theta^*), \quad (4)$$

74 where $\hat{\theta}_n$ is the strongly consistent sequence of roots of the likelihood equation (19) of Theorem 41. If, additionally,

$$75 \int_{\Theta} \frac{L_n(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}})}{L_n(\hat{\theta}_n)} \pi(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}) d\vartheta \xrightarrow{a.s.} \int_{\Theta} \exp(-\frac{1}{2} \vartheta^\top \mathcal{J}(\theta^*) \vartheta) \pi(\theta^*) d\vartheta \quad (5)$$

76 then

$$77 \int_{\Theta} |\pi_n(\vartheta) - N(\vartheta|0, \mathcal{J}(\theta^*)^{-1})| d\vartheta \xrightarrow{a.s.} 0. \quad (6)$$

78 *Proof.* We prove: fist the existence of the the MLE, then (4), and finally (6).

79 **Existence of consistent roots:** I gonna use Theorem 41 (in Appendix) to prove that there exists a consistent sequence Can be
80 $\hat{\theta}_n$ of roots of (Eq. 19 in Appendix), and hence I need to show that its conditions are satisfied. Let $S_\rho =$ skipped
81 $\{\theta : |\theta - \theta^*| \leq \rho\}$, with $\rho > 0$, be a neighborhood of θ^* on which (d3) is satisfied. So for $\Theta = S_\rho$
82 (in Theorem 41). Conditions (c1), (c2), (c5) of Theorem 41 (in Appendix) are automatic! Condition (c4)
83 follows from continuity of $f(x|\theta)$ at θ . Condition (c3), ok.... By Taylor's theorem, I expand $D(x, \theta) =$
84 $\log(f(x|\theta)) - \log(f(x|\theta^*))$ around θ^* as

$$85 D(x, \theta) = D(x, \theta^*) + \nabla_\theta \log(f(x|\theta^*))(\theta - \theta^*) \\ 86 + (\theta - \theta^*) \int_0^1 \int_0^1 v \nabla_\theta^2 \log(f(x|\theta_0 + uv(\theta - \theta^*))) du dv (\theta - \theta^*)$$

87 So because $D(x, \theta^*) = 0$, $\nabla_\theta \log(f(x|\theta^*))$ is integrable, and the components of $\nabla_\theta^2 \log(f(x|\theta))$ are bounded
88 by $K(x)$ uniformly on S_ρ , we get that $D(x, \theta)$ is bounded on S_ρ . So (c3) holds.

89 **Asymptotic Normality:** Let

$$90 \ell_n(\theta) = \log(L_n(\theta)); \quad \dot{\ell}_n(\theta) = \nabla_\theta \log(L_n(\theta)); \quad \ddot{\ell}_n(\theta) = \nabla_\theta^2 \log(L_n(\theta))$$

91 By Taylor's Theorem 38, we expand $\ell_n(\theta)$ around $\hat{\theta}_n$ as

$$92 \ell_n(\theta) = \ell_n(\hat{\theta}_n) + \dot{\ell}_n(\hat{\theta}_n)(\theta - \hat{\theta}_n) + (\theta - \hat{\theta}_n)^\top I_n(\theta)(\theta - \hat{\theta}_n)$$

93 where

$$94 I_n(\theta) = -\frac{1}{n} \int_0^1 \int_0^1 v \ddot{\ell}_n(\hat{\theta}_n + uv(\theta - \hat{\theta}_n)) du dv \quad (7)$$

95 Because, it is $\dot{\ell}_n(\hat{\theta}_n) = 0$ a.s., we get:

$$96 \ell_n(\theta) = \ell_n(\hat{\theta}_n) + (\theta - \hat{\theta}_n)^\top I_n(\theta)(\theta - \hat{\theta}_n) \iff \\ 97 \frac{L_n(\theta)}{L_n(\hat{\theta}_n)} = \exp(-(\theta - \hat{\theta}_n)^\top I_n(\theta)(\theta - \hat{\theta}_n)), \quad \text{a.s.}$$

98 Let's work on the asymptotics of (7); it is:

$$99 \frac{1}{n} \ddot{\ell}_n(\theta) = \frac{1}{n} \nabla_\theta^2 \log(L_n(\theta)) = \frac{1}{n} \nabla_\theta^2 \log\left(\prod_{i=1}^n f(x_i|\theta)\right) = \frac{1}{n} \sum_{i=1}^n \nabla_\theta^2 \log(f(x_i|\theta)) \\ 100 \xrightarrow{a.s.} E_F(\nabla_\theta^2 \log f(x|\theta)|\theta_0) \quad (8)$$

101 as $n \rightarrow \infty$ by SLLN. Also, it is

$$102 E_F(\nabla_\theta^2 \log f(x|\theta)|\theta_0) = -\mathcal{J}(\theta_0) \quad (9)$$

Hence, from (8) and (9), I get

$$\frac{1}{n} \ddot{\ell}_n(\theta) \xrightarrow{\text{a.s.}} -\mathcal{J}(\theta_0) \quad (10)$$

Therefore,

$$I_n(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}) = -\frac{1}{n} \int_0^1 \int_0^1 v \ddot{\ell}_n(\hat{\theta}_n + uv(\theta - \hat{\theta}_n)) du dv \xrightarrow{\text{a.s.}} \frac{1}{2} \mathcal{J}(\theta^*) \quad (11)$$

because of (10) and because of $\hat{\theta}_n \xrightarrow{\text{a.s.}} \theta^*$ from Theorem 41 (in Appendix).

So back to what we wish to prove, and putting all these together, it is

$$\begin{aligned} \frac{L_n(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}})}{L_n(\hat{\theta}_n)} \pi(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}) &= \exp(-\vartheta^\top I_n(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}) \vartheta) \pi(\hat{\theta}_n + \vartheta \frac{1}{\sqrt{n}}) \\ &\xrightarrow{\text{a.s.}} \exp(-\frac{1}{2} \vartheta^\top \mathcal{J}(\theta^*) \vartheta) \pi(\theta^*) \end{aligned}$$

because of (11) and $\hat{\theta}_n \xrightarrow{\text{a.s.}} \theta^*$.

Now, about the second part of the proof. If (5) then by dividing (4) and (5), I get

$$\pi_n(\vartheta) \xrightarrow{\text{a.s.}} \mathcal{N}(\vartheta|0, \mathcal{J}(\theta^*)^{-1})$$

for all $\theta \in \Theta$. Hence By Scheffe's Theorem 35 (in Appendix) we get (6).

□

Remark 15. Note that Bernstein-von Mises Theorem 14 implies that the posterior distribution of $\vartheta = \sqrt{n}(\theta - \hat{\theta})$ given the data converges to the Normal distribution $\mathcal{N}(0, \mathcal{J}(\theta_0)^{-1})$ in Total Variation Norm, namely

$$\sup_{A \subset \Theta} |\pi_n(\vartheta \in A) - \mathcal{N}(\vartheta \in A|0, \mathcal{J}(\theta^*)^{-1})| dx \rightarrow 0, \quad \text{as } n \rightarrow \infty$$

Corollary. If the conditions of (Bernstein-von Mises) Theorem 14 hold, and if $\mathcal{J}(\theta)$ is continuous at Θ , then

$$\sqrt{n} \mathcal{J}(\hat{\theta}_n)^{-1/2} (\theta - \hat{\theta}_n) \xrightarrow{D} z, \quad \text{where } z \sim \mathcal{N}(0, I_k) \quad (12)$$

this is the result stated in Stat Concepts II notes (Term 2, 2017).

Proof. Bernstein-von Mises Theorem implies $\sqrt{n}(\theta - \hat{\theta}_n) \xrightarrow{D} \mathcal{N}(0, \mathcal{J}(\theta^*)^{-1})$ or equiv.

$$Y_n = \sqrt{n} \mathcal{J}(\theta^*)^{1/2} (\theta - \hat{\theta}_n) \xrightarrow{D} Z, \quad (13)$$

with $Z \sim \mathcal{N}(0, I_k)$. From Theorem 41 (in Appendix) I get $\hat{\theta}_n \rightarrow \theta^*$ a.s.. Due to continuity of $\mathcal{J}(\theta)$, it is

$$X_n = \mathcal{J}(\hat{\theta}_n)^{1/2} \mathcal{J}(\theta^*)^{-1/2} \xrightarrow{\text{a.s.}} I_k \quad (14)$$

According to Slutsky's theorem¹ by multiplying (13), (14), I get $X_n Y_n \xrightarrow{D} Z$, i.e., $\sqrt{n} \mathcal{J}(\hat{\theta}_n)^{1/2} (\theta - \hat{\theta}_n) \xrightarrow{D} \mathcal{N}(0, I_k)$.

□

Example 16. Consider a Bayesian model

$$\begin{cases} x_i | \theta & \stackrel{\text{iid}}{\sim} \text{Bn}(\theta), & i = 1, \dots, n \\ \theta & \sim \text{Be}(a, b) \end{cases}$$

where $a > 0$, $b > 0$, and $n > 2$. Find the asymptotic posterior distribution of θ as $n \rightarrow \infty$, given Cramer's conditions.

¹Sluky's theorem: If $Y_n \xrightarrow{D} Z$ and $X_n \xrightarrow{\text{a.s.}} c$, where $c \in \mathbb{R}^k$ is a constant, then $X_n Y_n \xrightarrow{D} cZ$

Solution. I will find the MLE $\hat{\theta}_n$ of θ . The likelihood is $L_n(\theta) = \prod_{i=1}^n \text{Bn}(x_i|\theta)$. Then

$$\frac{d}{d\theta} \log f(x_{1:n}|\theta) = \frac{d}{d\theta} \sum_{i=1}^n \log(\text{Bn}(x_i|\theta)) = \frac{n\theta - \sum_{i=1}^n x_i}{\theta(1-\theta)} \implies$$

$$\frac{d}{d\theta} \log f(x_{1:n}|\theta)|_{\theta=\hat{\theta}_n} = 0 \implies \hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

I will find the Fisher Information; it is

$$\frac{d^2}{d\theta^2} \log(f(x|\theta)) = \frac{d^2}{d\theta^2} \log(\text{Bn}(x|\theta)) = -\frac{x}{\theta^2} - \frac{1-x}{(1-\theta)^2} \implies$$

$$\mathcal{J}(\theta) = -\mathbb{E}_{\text{Bn}(\theta)} \left(-\frac{x}{\theta^2} - \frac{1-x}{(1-\theta)^2} \right) = \frac{1}{\theta(1-\theta)}$$

According to Bernstein-von Mises Theorem 14, it is $\theta|x_{1:n} \sim \mathcal{N}\left(\hat{\theta}_n, \frac{1}{n} \mathcal{J}(\theta^*)^{-1}\right)$, where θ^* is the true value of θ .

According to Corollary 2, it is $\theta|x_{1:n} \sim \mathcal{N}\left(\hat{\theta}_n, \frac{1}{n} \mathcal{J}(\hat{\theta}_n)^{-1}\right)$ as well.

3 Continuous θ : Asymptotic distribution under Chen (1985) conditions

Notation 17. Let $U_n(\theta) = \log(\pi_n(\theta))$. Let $|\theta| = \sqrt{\theta^\top \theta}$. Let $B_\delta(\theta^*) = \{\theta \in \Theta; |\theta - \theta^*| < \delta\}$.

Assumption 18. Assume that the posterior density has maximum at $\theta = m_n$ such that $\dot{U}_n(m_n) = 0$, and $\Sigma_n = -(\ddot{U}_n(m_n))^{-1}$ where $\Sigma_n > 0$ is positive-definite.

Note 19. We consider weaker conditions which guarantee that the posterior density can be approximated by a Normal distribution around a small neighborhood of a posterior density maximum m_n as $n \rightarrow \infty$.

Remark 20. Assumption 18 is so general that (i.) $\{m_n\}$ is not assumed to converge, (ii.) $\pi_n(\theta)$ can be multimodal for each n , (iii.) m_n need not be the global maximum point of $\pi_n(\theta)$ for each n .

Condition 21. Consider the following regularity conditions

e1 (Steepness) $\bar{\sigma}_n^2 \rightarrow 0$ as $n \rightarrow \infty$ where $\bar{\sigma}_n^2$ is the largest eigenvalue of Σ_n

e2 (Smoothness) For any $\epsilon > 0$ there exists N and $\delta > 0$ such that, for any $n > N$ and $\theta \in B_\delta(m_n)$, $\ddot{U}_n(\theta)$ exists and satisfies

$$I - A(\epsilon) \leq \ddot{U}_n(\theta) (\ddot{U}_n(m_n))^{-1} \leq I + A(\epsilon),$$

where I is the $k \times k$ identity matrix, and $A(\epsilon)$ is a $k \times k$ symmetric positive-semidefinite matrix whose largest eigenvalue tends to zero as $\epsilon \rightarrow 0$

e3 (Concentration) For any $\delta > 0$, as $n \rightarrow \infty$.

$$Q_n := \int_{B_\delta(m_n)} \pi_n(\theta) d\theta \rightarrow 1. \quad (15)$$

Remark 22. Conditions 21 are weaker than Cramer Conditions 13. Conditions (e1 & e2) imply that $\pi_n(\theta)$ becomes pick around m_n and behave like a normal kernel inside a neighborhood of m_n . Condition (e3) ensures that the mass outside that neighborhood is negligible. No IID sampling is assumed.

Lemma 23. If conditions (e1) and (e2) hold then

$$\lim_{n \rightarrow \infty} \pi_n(m_n) |\Sigma_n|^{1/2} \leq (2\pi)^{-k/2}. \quad (16)$$

The equality holds when condition (e3) is satisfied.

Proof. Omitted but provided in the Exercise sheet. \square

Theorem 24. Assume posterior density $\pi_n(\theta)$ has maximum at $\theta = m_n$ such that $\dot{U}_n(m_n) = 0$, $\Sigma_n > 0$ where $\Sigma_n = -(\ddot{U}_n(m_n))^{-1}$, and $U_n(\theta) = \log(\pi_n(\theta))$. Let $\phi_n = \Sigma_n^{-1/2}(\theta - m_n)$, where $\Sigma_n^{-1/2}$ is the inverse of the lower matrix from the Cholesky decomposition of Σ_n . Given (e1), and (e2), (e3) is necessary and sufficient condition so that

$$\Sigma_n^{-1/2}(\theta - m_n) \xrightarrow{D} Z; \quad \text{where } Z \sim N(0, 1).$$

Proof. Define $Z_n = \Sigma_n^{-1/2}(\theta - m_n)$. Assume $a, b \in \Theta$, such that $a \leq b$. It is sufficient to show that for any $a \leq 0$ and $b \geq 0$, it is $\lim_{n \rightarrow \infty} P_n(a, b) = P(a, b)$, where $P_n(a, b) = P(a \leq Z_n \leq b)$ and $P(a, b) = P(a \leq Z \leq b)$, if and only if C-3 holds.

Write

$$P_n(a, b) = \int_{R_n} \pi_n(\theta) d\theta,$$

where $R_n = \{\theta \mid \Sigma_n^{1/2}a \leq (\theta - m_n) \leq \Sigma_n^{1/2}b\} \subseteq B_\delta(m_n)$, for any $\delta > 0$ and sufficiently large n , by (e1).

for every $\epsilon > 0$, $P_n(a, b) \in [P_n^-(a, b, \epsilon), P_n^+(a, b, \epsilon)]$, where

$$P_n^+(a, b, \epsilon) = \pi_n(m_n) |\Sigma_n|^{1/2} |I - A(\epsilon)|^{-1/2} \int_{R(\epsilon)} \exp(-\frac{1}{2} z^\top z) dz;$$

$$P_n^-(a, b, \epsilon) = \pi_n(m_n) |\Sigma_n|^{1/2} |I + A(\epsilon)|^{-1/2} \int_{R(\epsilon)} \exp(-\frac{1}{2} z^\top z) dz,$$

where $R(\epsilon) = \{z \mid [I - A(\epsilon)]^{-1/2}a \leq z \leq [I - A(\epsilon)]^{-1/2}b\}$.

By letting $\epsilon \rightarrow 0$, and under (e1), (e2), we get

$$\lim_{n \rightarrow \infty} P_n(a, b) = \lim_{n \rightarrow \infty} \pi_n(m_n) |\Sigma|^{1/2} \int_R \exp(-\frac{1}{2} z^\top z) dz,$$

where $R = \{z \mid a \leq z \leq b\}$. According to Lemma 23, $\lim_{n \rightarrow \infty} P_n(a, b) = P(a, b)$ if and only if (e3) holds. \square

Remark 25. Conditions (e1) and (e2) in Theorem 24 are relatively easy to check in practice. Condition (e3) maybe be a bit tricky, hence, two alternative conditions (e3.1) and (e3.2) for the tail behaviors of $\pi_n(\theta)$ are provided. They are especially useful when m_n is the global maximum point of $\pi_n(\theta)$ for all n , such as in the unimodal case.

Proposition 26. Assume that (e1) and (e2) hold. Then, either (e3.1) or (e3.2) implies (e3).

e3.1 For any $\delta > 0$, there exists an integer N , and true numbers $c > 0$, $p > 0$ such that, for any $n > N$ and $\theta \notin B_\delta(m_n)$,

$$U_n(\theta) - U_n(m_n) < -c((\theta - m_n)^\top \Sigma_n^{-1}(\theta - m_n))^p$$

e3.2 For any $\delta > 0$, there exists an integer N , and real numbers $c > 0$, $p > 0$ such that, for any $n > N$ and $\theta \notin B_\delta(m_n)$,

$$U_n(\theta) - U_n(m_n) < -c/|\Sigma_n|^p + \log(g(\theta)),$$

for some integrable function $g(\theta)$, i.e. $\int g(\theta) d\theta < \infty$.

Proof. Under (e3.1)

$$Q_n < \pi_n(m_n) |\Sigma_n|^{1/2} \int_{|z| > \delta/\sigma_n} \exp(-c(z^\top z)^d) dz \quad (17)$$

Under (e3.2)

$$Q_n < \pi_n(m_n) |\Sigma_n|^{1/2} |\Sigma_n|^{-1/2} \int_{|z| > \delta/\sigma_n} \exp(-c|\Sigma_n|^{-d}) dz \quad (18)$$

where Q_n in (15). From Lemma 23 $\lim_{n \rightarrow \infty} \pi_n(\theta) |\Sigma_n|^{1/2}$ is bounded by $(2\pi)^{-k/2}$. Also rest terms in (17) and (18) tend to zero as $n \rightarrow \infty$. Then $\lim_{n \rightarrow \infty} Q_n = 0$, and e3 is implied. \square

Remark 27. Assumptions (e3.1) and (e3.2) do not require the computation of the, often unknown, normalizing constant because it is simplified,

$$U_n(\theta) - U_n(m_n) = \log(f(x_{1:n}|\theta)) - \log(f(x_{1:n}|m_n)) + \log(\pi(\theta)) - \log(\pi(m_n)).$$

Remark 28. When the sample size is large enough, most priors will lead to the same inference and this inference will be equivalent to the one based only on the likelihood function.

Example 29. (Cont. Example 16) Consider the posterior distribution is $\theta|x_{1:n} \sim \text{Be}(a_n, b_n)$, where $a_n = a + n\bar{x}$, and $b_n = b + n - n\bar{x}$. Find the asymptotic distribution of θ as $n \rightarrow \infty$.

Solution. It is

$$U_n(\theta) = \log(\pi(\theta|x_{1:n})) = (a_n - 1) \log(\theta) + (b_n - 1) \log(1 - \theta) - \log f(x_{1:n})$$

So

$$\dot{U}_n(\theta) = \frac{a_n - 1}{\theta} - \frac{b_n - 1}{1 - \theta}; \quad \ddot{U}_n(\theta) = \frac{a_n - 1}{\theta^2} - \frac{b_n - 1}{(1 - \theta)^2};$$

Then

$$m_n := \frac{a_n - 1}{a_n + b_n - 2}; \quad \Sigma_n := (-U_n''(m_n))^{-1} = \frac{(a_n - 1)(b_n - 1)}{(a_n + b_n - 2)^3}.$$

Condition (e1) holds because $\lim_{n \rightarrow \infty} (-\ddot{U}_n(m_n))^{-1} = 0$. Condition (e2) holds because $\dot{U}_n(\theta)$ is a continuous with respect of θ . Condition (e3) holds by using the same arguments as in Theorem ???. Therefore, θ has asymptotic posterior distribution $\theta|x_{1:n} \dot{\sim} N(m_n, \Sigma_n)$.

4 Continuous θ : Asymptotic efficiency of Bayes Estimates

Remark 30. Consider the squared error loss $\ell(\theta, \delta) = (\theta - \delta)^\top (\theta - \delta)$ which implies the posterior expectation $\delta^\pi = E_\Pi(\theta|x_{1:n})$ as Bayes point estimator. Given that we can interchange the limit and the expectation operator of $\vartheta = \sqrt{n}(\theta - \hat{\theta}_n)$, we get $\sqrt{n}(\delta^\pi - \hat{\theta}_n) \xrightarrow{P} 0$ meaning that δ^π and $\hat{\theta}_n$ are asymptotically equivalent; i.e. $\delta^\pi - \hat{\theta}_n \xrightarrow{P} 0$. Hence

$$\sqrt{n}(\delta^\pi - \theta^*) = \sqrt{n}(\delta^\pi - \hat{\theta}_n) + \sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{D} N(0, \mathcal{I}(\theta^*)^{-1})$$

which means that the Bayes estimator is asymptotically efficient.

Example 31. (Cont. Example 29) How the Bayes estimators under the square loss and under the 0-1 loss behave as $n \rightarrow \infty$?

Solution. The exact posterior distr. is $\theta|x_{1:n} \sim \text{Be}(a + n\bar{x}, b + n - n\bar{x})$. The squared loss and the 0 - 1 loss imply the posterior mean $\delta_1(x_{1:n}) = \frac{n\bar{x} + a}{n + a + b}$ and posterior mode $\delta_2(x_{1:n}) = \frac{n\bar{x} + a - 1}{n + a + b - 2}$ as Bayes estimators correspondingly. Both converge to the MLE $\hat{\theta}_n = \bar{x}$ since $\lim_{n \rightarrow \infty} \delta_1(x_{1:n}) = \bar{x}$ and $\lim_{n \rightarrow \infty} \delta_2(x_{1:n}) = \bar{x}$.

5 Exercises

Try the Exercises ??, ?? from the Exercise sheet.

Appendix

A An inventory of definitions

Definition 32. (Types of converge) Assume a probability triplet $\{\Omega, \mathcal{F}, P\}$, and a sequence of random quantities $\{x_n; n = 1, 2, \dots\}$, such that $x_n : \Omega \rightarrow \mathbb{R}^d$, $d > 0$. Then

- $\{x_n\}$ converges almost surely to a random quantify x if and only if

$$P(\lim_{n \rightarrow \infty} x_n = x) = 1.$$

It is demoted as $x_n \xrightarrow{\text{a.s.}} x$.

- $\{x_n\}$ converges in distribution to a random quantify x if and only if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} F_{x_n}(t) = F_x(t)$$

for all continuity points t of F . in \mathbb{R} , where $F_{x_n}(t) = P(x_n \leq t)$ and $F_x(t) = P(x \leq t)$ are CDFs of x_n and x . It is demoted as $x_n \xrightarrow{D} x$.

- $\{x_n\}$ converges in total variation a random quantity x if and only if

$$\lim_{n \rightarrow \infty} \sup_{\forall B \subset \Theta} |P(x_n \in B) - P(x \in B)| = 0,$$

It is demoted as $x_n \xrightarrow{\text{T.V.}} x$.

Definition 33. (Upper semicontinuous) A real-valued function, $f(\theta)$, defined on Θ is said to be upper semicontinuous (u.s.c.) on Θ , if for all $\theta \in \Theta$ and for any sequence θ_n in Θ such that $\theta_n \rightarrow \theta$, we have

$$\lim_{n \rightarrow \infty} \sup f(\theta_n) \leq f(\theta)$$

Proposition 34. If $\pi_n(\cdot)$ and $\pi(\cdot)$ are the PDFs of x_n and x correspondingly, then

$$\sup_{\forall B \subset \Theta} |P(x_n \in B) - P(x \in B)| = \int \frac{1}{2} |\pi_n(t) - \pi(t)| dt$$

Theorem 35. (Scheffe convergence theorem²) If $f_n(\cdot)$ and $g(\cdot)$ are density functions such that for all $x \in \mathcal{X}$ $\lim_{n \rightarrow \infty} f_n(x) = g(x)$, then

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} |f_n(x) - g(x)| dx = 0$$

(that is a point-wise convergence of densities)

Theorem 36. If random variables x_n has density $f_n(x)$ and random variable x has density $g(x)$, and if $\lim_{n \rightarrow \infty} \int_{\mathcal{X}} |f_n(x) - g(x)| dx = 0$ then

$$\sup_{\forall A \subset \mathcal{X}} |P(x_n \in A) - P(x \in A)| \rightarrow 0, \quad \text{as } n \rightarrow \infty$$

that is called convergence in Total Variation.

Theorem 37. (A Strong law of large numbers (SLLN)) Let $\{x_i\}_{i=1}^n$ be a sequence of IID random quantities, with $E(x_i) = \mu < \infty$, and $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$ then $\bar{x}_n \xrightarrow{\text{a.s.}} \mu$.

²This is not the original version, but it is what we need

Theorem 38. (Taylor's theorem) If $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and if $\nabla^2 f(x) = \nabla(\nabla f(x))^\top$ is continuous in the ball $\{x \in \mathcal{X} : |x - x_0| < r\}$, then for $|t| < r$, it is

$$f(x_0 + t) = f(x_0) + \nabla f(x_0)t + t^\top \cdot \int_0^1 \int_0^1 u \nabla^2 f(x_0 + uvt) du dv \cdot t$$

Theorem 39. (Shannon-Kolmogorov Information Inequality) Let $f_0(x)$ and $f_1(x)$ be densities with respect to Lebesgue measure dx . Then

$$KL(f_0 || f_1) = E_{F_0(x)}(\log \frac{f_0(x)}{f_1(x)}) = \int_{\mathcal{X}} \log \frac{f_0(x)}{f_1(x)} f_0(x) dx \geq 0,$$

with equality if and only if $f_1(x) = f_0(x)$ a.s.

Lemma 40. (Passing the derivative under the integral operator) If $(\partial/\partial\theta)g(x, \theta)$ exists and is continuous in θ for all x and all θ in an open interval \mathcal{S} and if $|(\partial/\partial\theta)g(x, \theta)| \leq K(x)$ on \mathcal{S} where $\int K(x)dx < \infty$ and if $\int g(x, \theta)dx$ exists on \mathcal{S} , then

$$\frac{d}{d\theta} \int g(x, \theta)dx = \int \frac{d}{d\theta} g(x, \theta)dx$$

B Strong consistency of Maximum Likelihood Estimates

In frequentist statistics, given that $\nabla_\theta f(x|\theta)$ exists, one may seek to find the MLE $\hat{\theta}_n$ as the solution of the likelihood equation:

$$\hat{\theta}_n : \nabla_\theta \log f(x_{1:n}|\theta)|_{\theta=\hat{\theta}_n} = \sum_{i=1}^n \nabla_\theta \log f(x_i|\theta)|_{\theta=\hat{\theta}_n} = 0 \quad (19)$$

The following theorem states (more or less) that the MLE $\hat{\theta}_n$ in (19) is consistent.

Theorem 41. (Strong consistency of MLE) Let x_1, x_2, \dots be IID random variables with density $f(x|\theta)$ (with respect to measure dx), $\theta \in \Theta$, and let θ^* denote the true value of θ . If the following conditions are satisfied:

c1 Θ is a closed and bounded set in \mathbb{R}^k

c2 $f(x|\theta)$ is u.s.c. in θ for all $x \in \mathcal{X}$ $f(x|\theta)$ is continuous in θ for all x

c3 there is a function $K(x)$ such that $E^{f(x|\theta^*)}(|K(x)|) < \infty$ and

$$\log(f(x|\theta)) - \log(f(x|\theta^*)) \leq K(x), \quad \forall x, \forall \theta$$

c4 for all $\theta \in \Theta$ and sufficiency small $\rho > 0$, $\sup_{|\theta' - \theta| < \rho} f(x|\theta')$ is measurable in x

c5 (identifiability) $f(x|\theta) = f(x|\theta^*)$ a.s. then $\theta = \theta^*$

then, for any sequence of maximum-likelihood estimates $\hat{\theta}_n$ of θ , it is

$$\hat{\theta}_n \xrightarrow{a.s.} \theta^* \quad (20)$$

The following theorem states (more or less) that the MLE is asymptotically normal.

Theorem 42. (Cramer) Let x_1, x_2, \dots be IID random variables density $f(x|\theta)$ (with respect to some distribution $F(x|\theta)$), $\theta \in \Theta$, and let θ^* denote the true value of θ . If the conditions (d1)-(d5) stated in Theorem 14 (check in the next Theorem) are satisfied, then there exists a strongly consistent sequence $\hat{\theta}_n$ of roots of the likelihood equation (19) such that

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{D} N(0, \mathcal{J}(\theta^*)^{-1})$$