

Large Deviation Principles for Posterior Distributions of the Normal Parameters

Takuhisa Shikimi

Abstract

Suppose that X_1, X_2, \dots are conditionally i.i.d. random variables with distribution P given $\vartheta = \theta$, where ϑ is an unknown parameter. If P is a normal distribution with mean μ and known variance σ^2 , and if the prior of ϑ is chosen from the conjugate family $N(\mu, \nu^2)$ or proportional to the Lebesgue measure, then it follows that the posterior distributions given X_1, \dots, X_n obey a large deviation principle with a rate function. If P is a normal distribution with known mean and unknown precision τ , and if as a prior we choose the gamma distribution with parameters α and β or the improper distribution $(1/\tau) d\tau$, the Jeffereys' prior, then the posterior distributions of ϑ given X_1, \dots, X_n are shown to satisfy a large deviation principle. The Gärtner-Ellis theorem plays the key role to prove these large deviation principles for the posterior distributions.

Keywords: large deviations; posterior distributions; the Gärtner-Ellis theorem.

1 Introduction

Let ϑ be an unknown parameter with prior π and X_1, X_2, \dots are conditionally i.i.d. random variables with conditional distribution P given $\vartheta = \theta$. In this paper, we will show the large deviation principles for the posterior dis-

tribution of ϑ given X_1, \dots, X_n when P is either $N(\mu, \sigma^2)$ with known σ^2 or $N(\mu, 1/\lambda)$ with known μ . There is comparatively little literature on the exponential rate of convergence of posterior distributions. Fu and Kass (1988) studies the rate of convergence of posterior distributions in the neighborhood of the mode. In the nonparametric Bayesian framework, Shen and Wasserman (2001) studies the rate at which the posterior distribution concentrates around the true parameter, and Ganesh and O'Connell (1999) proves the large deviation principle for posterior distributions given i.i.d. random variables taking values in a finite set. Moderate deviation asymptotic results in a Bayesian setting are studied by Eichelsbacher and Ganesh (2002).

We begin with constructing the Bayesian framework appropriate for formulating and proving our results. Let $(\mathcal{X}, \mathcal{U})$ and $(\mathcal{Y}, \mathcal{A})$ be Polish spaces (complete separable metric spaces) endowed with the Borel σ -algebras \mathcal{U} and \mathcal{A} . A stochastic kernel from $(\mathcal{X}, \mathcal{U})$ to $(\mathcal{Y}, \mathcal{A})$ is a family $(P : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y}))$ of probability measures on $(\mathcal{Y}, \mathcal{A})$ indexed by \mathcal{X} , namely a mapping $\mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$, where $\mathcal{P}(\mathcal{Y})$ is the set of probability measures on $(\mathcal{Y}, \mathcal{A})$, such that for each $A \in \mathcal{A}$, $\mathcal{X} \rightarrow \mathbb{P}(A) \in [0, 1]$ is measurable. The space $\mathcal{P}(\mathcal{Y})$ equipped with the weak topology is metrized by the Lévy-Prohorov metric, and with respect to the metric $\mathcal{P}(\mathcal{Y})$ is Polish. Further, a family $(P : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y}))$ of probability measures is a stochastic kernel from $(\mathcal{X}, \mathcal{U})$ to $(\mathcal{Y}, \mathcal{A})$ if and only if $\mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$ is measurable with respect to \mathcal{U} and the Borel σ -algebra generated by the weak topology of $\mathcal{P}(\mathcal{Y})$. In the literature of statistics, $(\mathcal{X}, \mathcal{U})$ is called a parameter space and a stochastic kernel from $(\mathcal{X}, \mathcal{U})$ to $(\mathcal{Y}, \mathcal{A})$ is called a statistical model on $(\mathcal{Y}, \mathcal{A})$. Let $(\mathcal{X}^{(n)}, \mathcal{A}^{(n)}) = (\mathcal{X}, \mathcal{A}) \times \cdots \times (\mathcal{X}, \mathcal{A})$ be the n -dimensional product measurable space and let $P^{(n)} = P \times \cdots \times P$ be the n -dimensional product measure of P with itself. It is easy to show

that $(P^{(n)}: \mathcal{P} \rightarrow \mathcal{P}^{(n)})$ is a stochastic kernel from $(\mathcal{X}, \mathcal{U})$ to $(\mathcal{X}^{(n)}, \mathcal{A}^{(n)})$. We assume that each P is absolutely continuous with respect to a σ -finite measure ν on $(\mathcal{X}, \mathcal{A})$, i.e., $P \ll \nu$ for all n . If $f(x|y)$ is a density of P with respect to ν

$$P^{(n)}(A) = \int_A \prod_{i=1}^n f(x_i|y) d(\nu \times \cdots \times \nu), A \in \mathcal{A}^{(n)}.$$

We further assume that $f(x|y)$ is measurable as a function of $(y, x) \in \mathcal{X} \times \mathcal{X}$.

Let \mathbb{P} be a prior distribution on $(\mathcal{X}, \mathcal{U})$. Define probability measures $\mathbb{P}^{(n)}$ on $(\mathcal{X} \times \mathcal{X}^{(n)}, \mathcal{U} \times \mathcal{A}^{(n)})$ by

$$\mathbb{P}^{(n)}(C) = \int_C \prod_{i=1}^n f(x_i|y) d(\mathbb{P} \times \nu \times \cdots \times \nu), C \in \mathcal{U} \times \mathcal{A}^{(n)}, n \geq 1.$$

Since the sequence $\mathbb{P}^{(1)}, \mathbb{P}^{(2)}, \dots$ is consistent, that is,

$$\mathbb{P}^{(n+1)}(C \times \mathcal{X}) = \mathbb{P}^{(n)}(C)$$

for all $n \geq 1$ and $C \in \mathcal{U} \times \mathcal{A}^{(n)}$, it follows by Kolmogorov's consistency theorem that there is a unique probability measure $\mathbb{P} = \mathbb{P}^\infty$ on the infinite product space $(\mathcal{X}^\infty, \mathcal{F}) = (\mathcal{X} \times \mathcal{X}^{(\infty)}, \mathcal{U} \times \mathcal{A}^{(\infty)})$ such that

$$\mathbb{P}(C \times \mathcal{X} \times \mathcal{X} \times \cdots) = \mathbb{P}^{(n)}(C), n \geq 1, C \in \mathcal{U} \times \mathcal{A}^{(n)}.$$

Now let us introduce the coordinate functions. For $\vartheta = (y, (x_i))$ define

$$\begin{aligned} \vartheta(\vartheta) &= \vartheta, \\ X_i(\vartheta) &= x_i \in \mathcal{X}, i = 1, 2, \dots \end{aligned}$$

We think of ϑ as the unknown parameter and X_1, X_2, \dots as data. If $U \in \mathcal{U}$ and $A \in \mathcal{A}^{(n)}$, then

$$\begin{aligned}
\mathbb{P}(\vartheta \in U, (X_1, \dots, X_n) \in A) &= \mathbb{P}(U \times A \times \mathcal{X} \times \mathcal{X} \times \dots) \\
&= \mathbb{P}^{(n)}(U \times A) = \int_{U \times A} \prod_{i=1}^n f(x_i | \cdot) d(\mathbf{x} \times \dots \times \cdot) \\
&= \int_U \prod_{i=1}^n f(x_i | \cdot) d(\mathbf{x} \times \dots \times \cdot) d \\
&= \int_U P^{(n)}(A) (d) .
\end{aligned}$$

In particular, the marginal distributions of ϑ , (X_1, \dots, X_n) and X_i are given as follows:

$$\mathbb{P}(\vartheta \in U) = \int_U d, U \in \mathcal{U},$$

$$\mathbb{P}((X_1, \dots, X_n) \in A) = \int_A P^{(n)}(A) (d), A \in \mathcal{A}^{(n)},$$

$$\mathbb{P}(X_i \in A) = \int_A P(A) (d), A \in \mathcal{A}, i = 1, 2, \dots$$

For each ϑ , $P_{\vartheta}^{(n)}$ is a probability measure on $(\mathcal{X}^{(n)}, \mathcal{A}^{(n)})$. If $A \in \mathcal{A}^{(n)}$, $P_{\vartheta}^{(n)}(A)$ is (ϑ) -measurable as a function of ϑ and

$$\int_{\vartheta} P_{\vartheta}^{(n)}(A) d\mathbb{P} = \int_U P^{(n)}(A) (d) = \mathbb{P}(\vartheta \in U, (X_1, \dots, X_n) \in A),$$

which shows that $P_{\vartheta}^{(n)}(A)$ is a conditional probability of $(X_1, \dots, X_n) \in A$ given ϑ . Thus $P_{\vartheta}^{(n)}(A)$ is a regular conditional distribution of (X_1, \dots, X_n) given ϑ . In the same way, $P_{\vartheta}(A)$, $A \in \mathcal{A}$ is shown to be a regular conditional distribution of X_i given ϑ . In particular, $\mathbb{P}(X_i \in A | \vartheta) = P_{\vartheta}(A)$ a.s., $A \in \mathcal{A}$. Since

$$\begin{aligned}
\mathbb{P}((X_1 \in A_1, \dots, X_n \in A_n) | \vartheta) &= P_{\vartheta}^{(n)}(A_1 \times \dots \times A_n) \\
&= P_{\vartheta}(A_1) \cdots P_{\vartheta}(A_n)
\end{aligned}$$

$$= \mathbb{P}(X_1 \in A_1 | \vartheta) \cdots \mathbb{P}(X_n \in A_n | \vartheta) \text{ a.s.,}$$

it follows that X_1, X_2, \dots are conditionally i.i.d. given ϑ .

Since ϑ is a random variable with values in a Polish space $(\mathcal{D}, \mathcal{B})$, there is a regular conditional distribution of ϑ given X_1, \dots, X_n , which we will call the posterior distribution of ϑ given X_1, \dots, X_n and denote it by $\hat{\mathbb{P}}(\vartheta \in d | X_1, \dots, X_n)$ or $\pi_n(d)$.

2 The large deviation principle

Let (Q_n) be a sequence of probability measures on a Polish space $(S, \mathcal{B}(S))$, where $\mathcal{B}(S)$ is the Borel σ -algebra. A function $I : S \rightarrow [0, \infty]$ is called a rate function on S if for each $M < \infty$ the level set $\{x \in S : I(x) \leq M\}$ is a compact subset of S . It follows necessarily that I is a lower semicontinuous function. In most applications, I is a convex function. The large deviation principle focuses on the asymptotic behavior of the sequence (Q_n) in terms of a rate function I . The precise formulation is as follows.

Definition 1. Let (Q_n) be a sequence of probability measures on S and I a rate function on S . We say that (Q_n) satisfies the large deviation principle with rate function I if the following two conditions hold.

- (i) For each closed subset $F \subset S$

$$\limsup_n \frac{1}{n} \log Q_n(F) \leq - \inf_{x \in F} I(x).$$

- (ii) For each open subset $G \subset S$

$$\liminf_n \frac{1}{n} \log Q_n(G) \geq - \inf_{x \in G} I(x).$$

It is well known that if (Q_n) satisfies the large deviation principle with a

rate function, then the rate function is uniquely determined. We refer the reader to Dembo and Zeitouni (1998) or Deuschel and Strook (2000) for basic results and a wide variety of large deviation techniques.

A convex function $\psi : \mathbb{R} \rightarrow (-\infty, \infty]$ is said to be essentially smooth if

- (i) E° is not empty, where E° is the interior of $E = \{t \in \mathbb{R} : \psi(t) < \infty\}$;
- (ii) ψ is differentiable on E° ;
- (iii) ψ is steep, namely

$$\lim_n |\psi'(t_n)| = \infty$$

whenever (t_n) is a sequence in E° converging to a boundary point of E° .

Here we state the Gärtner-Ellis theorem, an extension of Crámer’s theorem to non-i.i.d. random variables. This is a key theorem in this paper for proving the large deviation principles for posterior distributions of the normal parameters.

Theorem 2 . *Let (Q_n) be a sequence of probability measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ such that the cumulant generating function*

$$\psi(t) = \lim_n \frac{1}{n} \log \int_{\mathbb{R}} e^{ntx} Q_n(dx) \tag{1}$$

exists in $[-\infty, \infty]$ for each $t \in \mathbb{R}$. If E° contains the origin, then $\psi(t) > -\infty$, ψ is convex, and the Fenchel-Legendre transform of ψ (the conjugate function of ψ)

$$I(x) = \sup_{t \in \mathbb{R}} [tx - \psi(t)] \tag{2}$$

is a rate function on \mathbb{R} .

Theorem 3 (Gärtner-Ellis) . *Let (Q_n) be a sequence of probability measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ such that (1) exists for all $t \in \mathbb{R}$. If E° contains the origin, then the followings hold with the rate function I defined by (2) .*

(i) For each closed subset $F \subseteq S$

$$\limsup_n \frac{1}{n} \log Q_n(F) \leq - \inf_{x \in F} I(x) .$$

(ii) For each open subset $G \subseteq S$

$$\liminf_n \frac{1}{n} \log Q_n(G) \geq - \inf_{x \in G} I(x) ,$$

where H is the set of exposed points of I :

$$H = \{x \in \mathbb{R} : \exists t \in E \circ y \text{ s.t. } ty - I(y) > tx - I(x) \} .$$

If, in addition, I is essentially smooth and lower semicontinuous, then (Q_n) satisfies the large deviation principle with the rate function I .

For the proofs of Theorem 2 and Theorem 3, see Dembo and Zeitouni (1998). Our purpose is to obtain the large deviation principles for posteriors of the normal parameters given data. Accordingly, we have to slightly modify the definitions of rate function and the large deviation principle in order to cover a sequence of conditional distributions. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, (\mathcal{G}_n) an increasing sequence of sub σ -algebras, and Y random variable taking values in $(S, \mathcal{B}(S))$. Then for each $n \geq 1$ there exists a function $Q_n: (\Omega, \mathcal{G}_n) \times \mathcal{B}(S) \rightarrow [0, 1]$ such that

- (i) for each $B \in \mathcal{B}(S)$ $Q_n(B)$ is a probability measure on $(S, \mathcal{B}(S))$;
- (ii) for each $B \in \mathcal{B}(S)$ $Q_n(B)$ is a version of $\mathbb{P}(Y \in B | \mathcal{G}_n)$.

The function Q_n is called a regular conditional distribution of Y given \mathcal{G}_n .

We define a measurable function $I: \mathbb{R} \times S \rightarrow [0, \infty]$ to be a rate function on S if $I(\cdot, \cdot)$ is a rate function for each \mathcal{G}_n . It is natural to define the large deviation principle for a sequence of conditional distributions (Q_n) as follows.

Definition 4. Let $I: \mathbb{R} \times S \rightarrow [0, \infty]$ be a rate function on S . We say that a

sequence (Q_n) satisfies the large deviation principle with rate function I if the following two conditions hold.

(i) For each closed subset $F \subseteq S$

$$\limsup_n \frac{1}{n} \log Q_n(F) \leq - \inf_{x \in F} I(\mu, x) \quad \text{a.s.}$$

(ii) For each open subset $G \subseteq S$

$$\liminf_n \frac{1}{n} \log Q_n(G) \geq - \inf_{x \in G} I(\mu, x) \quad \text{a.s.}$$

3 The normal model with unknown mean and known variance

Suppose that $(\mu, \nu) = (\mathcal{X}, \mathcal{A}) = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$, where $\mathcal{B}(\mathbb{R})$ is the Borel σ -algebra of \mathbb{R} . Let $P = N(\mu, \nu^2)$, the normal distribution with mean μ and known variance ν^2 . If we choose the normal distribution $N(\mu, \nu^2)$ for the prior distribution of ϑ , then the posterior distribution of ϑ given X_1, \dots, X_n is the normal distribution $N(\mu_n, \nu_n^2)$, where

$$\mu_n = \frac{\nu^2/n}{\nu^2/n + \nu^2} \mu + \frac{\nu^2}{\nu^2/n + \nu^2} \bar{X}_n, \quad \bar{X}_n = \frac{X_1 + \dots + X_n}{n},$$

$$\frac{1}{\nu_n^2} = \frac{1}{\nu^2} + \frac{1}{\nu^2/n}.$$

Since X_1, X_2, \dots are conditionally i.i.d. given ϑ and X_1 is integrable under \mathbb{P} , the law of large numbers for conditionally i.i.d. random variables entails the convergence $\bar{X}_n \xrightarrow{\mathbb{P}} \mathbb{E}(X_1 | \vartheta) = \vartheta$ a.s., and hence we have

$$\mu_n \xrightarrow{\mathbb{P}} \vartheta \quad \text{a.s.}$$

Since

$$e^{nt} \int_{\mathbb{R}} p_n(d) = \exp\left(\mu_n nt + \frac{v_n^2 n^2 t^2}{2}\right)$$

we have the cumulant generating function

$$\psi(t) = \lim_n \frac{1}{n} \log \int_{\mathbb{R}} e^{t \cdot} p_n(d) = \vartheta t + \frac{v^2 t^2}{2} \tag{3}$$

which is that of $N(\vartheta, v^2)$. If we select the Lebesgue measure on \mathbb{R} as the prior for ϑ , then the posterior distribution is given by $N(\bar{X}_n, v^2/n)$, which is the limiting distribution of $N(\mu_n, v_n^2)$ as v^2 . It follows that the cumulant generating function of the posterior distribution is the same as (3).

It is easy to see that $E^{\vartheta} = \mathbb{R}$, is differentiable and steep for each ϑ , so that $I(\vartheta, \cdot)$ is essentially smooth and continuous for each ϑ . It is not difficult to show that the Fenchel-Legendre transform of ψ for each ϑ is given by

$$I(\vartheta, x) = \frac{(x - \vartheta)^2}{2} \tag{4}$$

Now we obtain the following theorem.

Theorem 5. *Suppose that X_1, X_2, \dots are conditionally i.i.d. with distribution $N(\cdot, v^2)$ given $\vartheta = \mu \in \mathbb{R}$. If $N(\mu, v^2), \mu \in \mathbb{R}, v^2 > 0$ is the prior for ϑ , then the posterior distributions satisfy the large deviation principle with rate function (4). If we choose the Lebesgue measure as the prior for ϑ , the same result follows.*

4 The normal model with known mean and unknown precision

Let X_1, X_2, \dots be conditionally i.i.d. random variables sampled from $P = N(\mu, 1/\lambda)$, $\lambda = (0, \infty)$. Assume that the prior distribution of the precision λ is a gamma distribution $G(\alpha, \beta)$, $\alpha > 0, \beta > 0$. Then the posterior dis-

tribution of ϑ given X_1, \dots, X_n is a gamma $G(n, n)$, where

$$\begin{aligned} n &= \frac{n}{2}, \\ n &= \frac{n}{2} S_n^2, S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2. \end{aligned}$$

Note that by the law of large numbers for conditionally i.i.d random variables

$$S_n^2 \xrightarrow{\text{a.s.}} \mathbb{E}[(X_1 - \mu)^2 | \vartheta] = \frac{1}{\vartheta} \text{ a.s.}$$

By the moment generating function of $G(n, n)$, we have

$${}_0 e^{nt} {}_n(d) = \begin{cases} \left(\frac{n}{n-nt}\right)^n, & t < \frac{n}{n}, \\ , & t \geq \frac{n}{n}, \end{cases}$$

and hence we obtain the cumulant generating function

$$\begin{aligned} (t) &= \lim_n \frac{1}{n} \log {}_0 e^{nt} {}_n(d) \\ &= \begin{cases} \log \left(\frac{1/(2\vartheta)}{1/(2\vartheta)-t}\right)^{1/2}, & t < \frac{1}{2\vartheta}, \\ , & t \geq \frac{1}{2\vartheta}. \end{cases} \end{aligned} \tag{5}$$

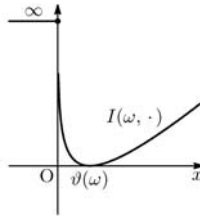
This is the cumulant generating function of $G(1/2, 1/(2\vartheta))$. If we utilize the improper prior

$$(d) \propto \frac{1}{d} \tag{6}$$

which is the Jefferys' prior for ϑ and obtained formally by letting $\rightarrow 0$ and

$\rightarrow 0$ in $G(\cdot, \cdot)$, then the same cumulant generating function as (5) is obtained. The Fenchel-Legendre transform of (5) has the following form.

$$I(\omega, x) = \begin{cases} \frac{1}{2} \left[\frac{x}{\vartheta(\omega)} - \log \frac{x}{\vartheta(\omega)} - 1 \right], & x > 0, \\ 0, & x \leq 0. \end{cases} \tag{7}$$



$I(\omega, \cdot)$ defined by (7).

By Theorem 2, $I(\omega, \cdot)$ is a convex rate function for each ω , and by Theorem 3 we have the following result.

Theorem 6. *Suppose that X_1, X_2, \dots are conditionally i.i.d. with distribution $N(\mu, 1/\vartheta)$ given $\vartheta = (\vartheta, \mu) = (0, \mu)$. If a gamma distribution $G(\alpha, \beta)$, $\alpha > 0$, $\beta > 0$ is the prior for ϑ , then*

$$\limsup_n \frac{1}{n} \log \hat{\mathbb{P}}(\vartheta \in G | X_1, \dots, X_n)(\omega) \leq - \inf_{x \in G} I(\omega, x) \text{ a.s.}$$

for all open sets $G \subset (0, \infty)$, and

$$\liminf_n \frac{1}{n} \log \hat{\mathbb{P}}(\vartheta \in F | X_1, \dots, X_n)(\omega) \geq - \inf_{x \in F \setminus H} I(\omega, x) \text{ a.s.}$$

for all closed sets $F \subset (0, \infty)$, where H is the set of exposed points of $I(\omega, \cdot)$.

If we select the Jeffreys's prior (6) for ϑ , the same result follows.

References

[1] Dembo, A. and Zeitouni, O. (1998). *Large Deviations Techniques and Applications*, 2nd ed., Springer-Verlag, New York.
 [2] Deuschel, J. D. and Stroock, D. W. (2000). *Large Deviations*, AMS Chelsea Publishing, Amer. Math. Soc.

- [3] Eichelsbacher, P. and Ganash, A.J . (2002) . Moderate deviations for Bayes posteriors. *Scand. J. Statist.* , 29 , 153-167 .
- [4] Fu, J. C. and Kass, R. E . (1988) . The exponential rate of convergence of posterior distributions, *Ann. Inst. Statist. Math.* , 40 , 683-691 .
- [5] Ganesh, A. and O'Connell, N . (1999) . An inverse of Sanov's theorem, *Statist. Probab. Lett.* , 42 , 201-206 .
- [6] Shen, X. and Wasserman, L . (1998) . Rates of convergence of posterior distributions, *Ann. Probab.* , 29 , 687-714 .