# Large Deviations for the Posterior Distributions under Conjugate Prior Distributions

## Takuhisa Shikimi

### Abstract

　　This paper takes up three parametric cases　the normal, Poisson, exponential cases　in order to study a large deviation upper bound for some posterior probabilitiy of the unknown parameter when in each case the prior distribution is assumed to be in a conjugate family. The upper bound will be given explicitly in each case.

**Keywords:** large deviations; posterior distributions; exchangeability.

## Introduction

　　Let $X_1, X_2,...$ be i.i.d. random variables with unknown distribution that belongs to a statistical model $(P : \quad)$　where　is a parameter space. In this paper, we focus on exponential rates of convergence of the posterior distributions in three parametric models　the normal, Poisson and exponential statistical models　when in each case the prior distribution is assumed to be in a conjugate family. There is comparatively little literature on the exponential rate of convergence of posterior distribution. 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.

We will give a large deviation upper bound in an explicit form for posterior probabilities of the event $[\ ,\ )$ given $X_1,...,X_n$ in each of the three parametric cases. In all cases, the basic tool to derive the results is the law of large numbers for exchangeable random variables (Theorem A. 3) together with the conditional Markov inequality.

## Constructing the model

Let $(\ ,\mathscr{U})$ be a measurable space. A stochastic kernel from $(\ ,\mathscr{U})$ to $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ where $\mathscr{B}(\mathbb{R}^n)$ is the Borel -algebra of $\mathbb{R}^n$ ($n$ 1, 2,..., ) is a family $(P:\quad)$ of probability measures on $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ indexed by such that for each $A\quad\mathscr{A}\qquad\mapsto P(A)\quad[0, 1]$ is measurable. As is usual $(P:\quad)$ is referred to as a statistical model. If $P^{(n)}$ is the $n$ dimensional product measure $P\times\quad\times P$ the infinite product probability measure $P^{(\ )}$
$P\times P\times\qquad$ is the unique probability measure on $(\mathbb{R}\ ,\mathscr{B}(\mathbb{R}\ ))$ such that

$$P^{(\ )}(A_1\times\quad\times A_n\times\mathbb{R}\times\mathbb{R}\times\quad)\quad P(A_1)\quad P(A_n)$$
$$P^{(n)}(A_1\times\quad\times A_n)$$

for all $n\geqslant 1$ and $A_1,..., A_n\quad\mathscr{B}(\mathbb{R})$.

**Lemma** 1 *For each $n$ 1, 2,..., the family $(P^{(n)}:\quad)$ is a stochastic kernel from $(\ ,\mathscr{U})$ to $(\mathbb{R}^n,\mathscr{B}(\mathbb{R}^n))$.*

*Proof.* We only show that $(P^{(\ )}:\quad)$ is a stochastic kernel, since $(P^{(n)}:$ $)\quad 1\leqslant n$ will be shown to be stochastic kernels in the same manner.

If we define

$$\mathscr{L} = \{ B \in \mathscr{B}(\mathbb{R}^\infty): \; \theta \mapsto P^{(\theta)}(B) \text{ is measurable} \}$$

then $\mathscr{L}$ is a $\sigma$-class containing the $\pi$-class

$$\mathscr{D} = \{ A_1 \times \cdots \times A_n \times \mathbb{R} \times \mathbb{R} \times \cdots : n \geqslant 1, A_1, \ldots, A_n \in \mathscr{B}(\mathbb{R}) \}$$

It follows that $\mathscr{B}(\mathbb{R}^\infty) = \sigma(\mathscr{D}) \subseteq \mathscr{L}$.　□

For a prior distribution $\mu$ on $(\Theta, \mathscr{U})$ define $\mathbb{P}$ to be the probability measure on $(\Omega, \mathscr{F}) = (\Theta \times \mathbb{R}^\infty, \mathscr{U} \times \mathscr{B}(\mathbb{R}^\infty))$ satisfying

$$\mathbb{P}(U \times B) = \int_U P^{(\theta)}(B)\,\mu(d\theta). \tag{1}$$

for every $U \in \mathscr{U}$ and $B \in \mathscr{B}(\mathbb{R}^\infty)$. It is not difficult to show the existence and uniqueness of $\mathbb{P}$. Now let us introduce the coordinate mappings $\vartheta$, $X$ and $\xi_i$ defined by

$$\vartheta(\omega) = \vartheta(\theta, x) = \theta \tag{2}$$

$$X(\omega) = X(\theta, x) = x \tag{3}$$

$$\xi_i(x) = x_i \quad (i \geqslant 1)$$

for $\omega = (\theta, x) \in \Omega$ and $x = (x_i) \in \mathbb{R}^\infty$. A random element $X$ is a sequence of random variables $X_1, X_2, \ldots$, where $X_i = \xi_i(X)$. We think of $\vartheta$ as the unknown parameter, $X = (X_1, X_2 \ldots)$ a date, where the distribution of $X_i$ is specified by $\vartheta$. By (2), (3) and (1), the parameter $\vartheta$ has $\mu$ as its distribution:

$$\mathbb{P}(\vartheta \in U) = \mu(U);$$

the distribution $\mathbb{P}(X \in dx)$ of $X$ is given by the mixture

$$\mathbb{P}\ (B)\ (d\ )\quad B\quad \mathscr{B}(\mathbb{R}\ );\tag{4}$$

the distribution $\mathbb{P}((X_1,..., X_n)\quad (dx_1,..., dx_n))$ is given by the mixture

$$P^{(n)}(B_n)\ (d\ )\quad B_n\quad \mathscr{B}(\mathbb{R}^n)\tag{5}$$

and the distribution $\mathbb{P}(X_i\quad dx_i)$ of $X_i$ is given by the mixture

$$P\ (A)\ (d\ )\quad A\quad \mathscr{B}(\mathbb{R})\tag{6}$$

In particular, $X_1, X_2...$ are identically distributed (but not independent in general) under $\mathbb{P}$ Distributions defined by (4) (5) and (6) are called prior predictive distributions of $X$, $(X_1,..., X_n)$ and $X_i$, respectively.

**Lemma** 2 *The function* $P_{\vartheta(\ )}^{(\ )}(B)$ *defined on* $\times \mathscr{B}(\mathbb{R}\ )$ *is a regular conditional distribution for* $X\quad (X_1, X_2...)$ *given* $\vartheta$ *For each* $n$ *the function* $P_{\vartheta(\ )}^{(n)}(B_n)$ *defined for* $(\ , B_n)\quad \times \mathscr{B}(\mathbb{R}^n)$ *is a regular conditional distribution of* $(X_1,..., X_n)$ *given* $\vartheta$ *Moreover,* $P_{\vartheta(\ )}(A)$ *defined for* $(\ , A)$ $\times \mathscr{B}(\mathbb{R})$ *is a regular conditional distribution of* $X_i$ *given* $\vartheta$ *for every* $i \geqslant 1$

*Proof.* For each        $P_{\vartheta(\ )}^{(\ )}$ is a probability measure on $(\mathbb{R}\quad \mathscr{B}(\mathbb{R}\ ))$ If $B\quad \mathscr{B}(\mathbb{R}\ )$

$$\int_{\vartheta\ U} P_{\vartheta(\ )}^{(\ )}(B)\ \mathbb{P}(d\ )\quad \int_{U} P^{(\ )}(B)\ (d\ )$$
$$\mathbb{P}(U\times\ B)$$
$$\mathbb{P}(\vartheta\quad U, X\quad B)$$

Thus, $P_{\vartheta(\ )}^{(\ )}(B)$ is a version of $\mathbb{P}(X\quad B\quad )(\ )$ because $P_{\vartheta(\ )}^{(\ )}(B)$ is $(\vartheta)$-measurable as a function of    for each $B$.

Likewise, $P_{\vartheta(\ )}^{(n)}(B_n)$ and $P_{\vartheta(\ )}(A)$ are regular conditional distributions for $(X_1,..., X_n)$ and $X_i(i\quad 1, 2...)$    respectively given $\vartheta$   since they are $(\vartheta)$-measurable and almost surely

$$\mathbb{P}((X_1,\dots, X_n) \in B_n \mid \vartheta)(\cdot) = \mathbb{P}(X \in B_n\times \mathbb{R}\times \mathbb{R}\times \cdots \mid \vartheta)(\cdot)$$

$$= P^{(\infty)}_{\vartheta(\cdot)}(B_n\times \mathbb{R}\times \mathbb{R}\times \cdots)$$

$$= P^{(n)}_{\vartheta(\cdot)}(B_n), \quad B_n \in \mathscr{B}(\mathbb{R}^n)$$

$$\mathbb{P}(X_i \in A \mid \vartheta)(\cdot) = \mathbb{P}(X \in \mathbb{R}\times \cdots \times \mathbb{R}\times A\times \mathbb{R}\times \cdots \mid \vartheta)(\cdot)$$

$$= P^{(\infty)}_{\vartheta(\cdot)}(\mathbb{R}\times \cdots \times \mathbb{R}\times A\times \mathbb{R}\times \cdots)$$

$$= P_{\vartheta(\cdot)}(A), \quad A \in \mathscr{B}(\mathbb{R})$$

<div align="right">□</div>

**Lemma** 3 *The random variables* $X_1, X_2\dots$ *are conditionally i.i.d. given* $\vartheta$

*Proof.* For all $n\geqslant 1$ and all $A_1,\dots, A_n \in \mathscr{B}(\mathbb{R})$

$$\mathbb{P}(X_1 \in A_1,\dots, X_n \in A_n \mid \vartheta)(\cdot) = P^{(n)}_{\vartheta(\cdot)}(A_1\times \cdots \times A_n)$$

$$= P_{\vartheta(\cdot)}(A_1)\cdots P_{\vartheta(\cdot)}(A_n)$$

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

where the first and third equalities follow from Lemma 2 Thus, $X_1, X_2,\dots$
are conditionally independent given $\vartheta$ Since $\mathbb{P}(X_i \in A \mid \vartheta)(\cdot) = P_{\vartheta(\cdot)}(A) = \mathbb{P}(X_1 \in A \mid \vartheta)(\cdot)$ a.s. for all $i\geqslant 1$ $X_1, X_2\dots$ are conditionally identically distributed.

Rea-valued random variables $Y_1, Y_2\dots$ are exchangeable if for all $n\geqslant 1$ and all permutations $\sigma$ of $1,\dots, n$

$$(Y_1,\dots, Y_n) \stackrel{d}{=} (Y_{\sigma(1)},\dots, Y_{\sigma(n)}) \tag{7}$$

Here $\stackrel{d}{=}$ stands for equality in distribution. de Finetti's theorem claims that random variables $Y_1, Y_2\dots$ are conditionally i.i.d. given some sub $\sigma$-algebra if and only if they are exchangeable. Lemma 3 tells us that $X_1, X_2,\dots$ are exchangeable random variables. See Aldous (1982) for an abstract version of de Finetti's theorem.

In what follows, we assume that $\Theta$ is a complete seperable metric space,

which is referred to as a Polish space. Accordingly, there exists a regular conditional distribution of $\vartheta$ given $X_1,..., X_n$ for all $n \geqslant 1$ which is termed a posterior distribution of $\vartheta$ given $X_1,..., X_n$ and denoted by $\pi_n(U)$ $(\ , U)$ $\times \mathscr{U}$ More precisely, there exists a function $\pi_n(U)$ on $\ \times \mathscr{U}$ such that

(a) for each $\quad\pi_n(\ )$ is a probability measure on $(\ , U)$;
(b) for each $U \ \mathscr{U},\ \pi_n(U)$ is a variant of $\mathbb{P}(\vartheta\ U\ X_1,..., X_n)(\ )$.

Suppose that the statistical model $(P:\ \ )$ is dominated by a $\ $-finite measure $\ $ on $(\mathbb{R}\ \mathscr{B}(\mathbb{R}))$ with density function $f(x\ )$ $x\ \mathbb{R}$. We assume that $f(x\ )$ is measurable as a function of $(\ x)\ \ \times \mathbb{R}$. The marginal distribution $\mathbb{P}((X_1,..., X_n)\ (dx_1,..., dx_n))$ of $(X_1,..., X_n)$ has the marginal density function

$$f_n(x_1,..., x_n)\ \ \prod_{i\ 1}^{n} f(x_i\ )\ (d\ )$$

with respect to $\ ^{(n)}$ (the $n$-fold measure of $\ $) i.e.,

$$\mathbb{P}((X_1,..., X_n)\ B_n)\ \ \int_{B_n} f_n(x_1,..., x_n)\ ^{(n)}(d(x_1,..., dx_n))$$

This can be seen from

$$\mathbb{P}(X_1\ A_1,..., X_n\ A_n)\ \ P^{(n)}(A_1 \times\ \ \times A_n)\ (d\ )$$

$$P\ (A_1)\ \ P\ (A_n)\ (d\ )$$

$$\left[\ \int_{A_1} f(x_1\ )\ (dx_1)\ \ \int_{A_1} f(x_n\ )\ (dx_n)\right]\ (d\ )$$

$$\int_{A_1 \times\ \times A_n}\ \prod_{i\ 1}^{n} f(x_i\ )\ ^{(n)}(d(x_1,..., x_n))\ (d\ )$$

$$\int_{A_1 \times\ \times A_n}\ \prod_{i\ 1}^{n} f(x_i\ )\ (d\ )\ ^{(n)}(d(x_1,..., x_n))$$

$$\int_{A_1 \times\ \times A_n} f_n(x_1,..., x_n)\ ^{(n)}(d(x_1,..., x_n))$$

Note that $\mathbb{P}(f_n(X_1, \ldots, X_n) \, 0) \, 0$.

**Lemma** 4   *If the statistical model* $(P : \quad)$ *is dominated by a* -*finite measure* $\quad$ *on* $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ *with density* $f(x \quad)$ *a measurable function on* $\times \mathbb{R}$, *then*

$$_n(U) \quad \left[ \quad_U \frac{\prod_{i\ 1}^{n} f(x_i \quad)}{f_n(X_1, \ldots, X_n)} \ (d \ ) \right] 1_{f_n\ 0} (X_1, \ldots, X_n)$$

$$( U) \, 1_{\{f_n\ 0\}}(X_1, \ldots, X_n)$$

*is a posterior distribution of* $\vartheta$ *given* $X_1, \ldots, X_n$.

*Proof.* It is easily seen that for each $\quad$ $_n( \ )$ is a probability measure on $( \ ,$ $\mathscr{U})$ and that for each $U \quad \mathscr{U}$ $_n(U)$ *is* $(X_1, \ldots, X_n)$-measurable. Thus it suffices to show that $_n(U) \quad \mathbb{P}(\vartheta \quad U \ X_1, \ldots, X_n)( \ )$ a.s. and this can be shown in the following way:

$$_{(X_1, \ldots, X_n)\ B_n} \, _n(U) \, d\mathbb{P} \quad _{\substack{(X_1, \ldots, X_n)\ B_n \\ f_n(X_1, \ldots, X_n)\ 0}} \left[ \quad_U \frac{\prod_{i\ 1}^{n} f(X_i \quad)}{f_n(X_1, \ldots, X_n)} \ (d \ ) \right] d\mathbb{P}$$

$$_U \left[ \quad_{\substack{(X_1, \ldots, X_n)\ B_n \\ f_n(X_1, \ldots, X_n)\ 0}} \frac{\prod_{i\ 1}^{n} f(X_i \quad)}{f_n(X_1, \ldots, X_n)} d\mathbb{P} \right] (d \ )$$

$$_U \left[ \quad_{B_n\ f_n\ 0} \frac{\prod_{i\ 1}^{n} f(x_i \quad)}{f_n(x_1, \ldots, x_n)} f_n(x_1, \ldots, x_n) \ ^{(n)}(d(x_1, \ldots, x_n)) \right] (d \ )$$

$$_U \left[ \quad_{B_n\ f_n\ 0} \prod_{i\ 1}^{n} f(x_i \quad) \ ^{(n)}(d(x_1, \ldots, x_n)) \right] (d \ )$$

$$_U P^{(n)}(B_n \quad f_n\ 0) \ (d \ )$$

$$\mathbb{P}(\vartheta \quad U \quad (X_1, \ldots, X_n) \quad B_n \quad f_n(X_1, \ldots, X_n) \quad 0)$$

$$\mathbb{P}(\vartheta \quad U \quad (X_1, \ldots, X_n) \quad B_n \quad f_n(X_1, \ldots, X_n) \quad 0)$$

$$\mathbb{P}(\vartheta \quad U \quad (X_1, \ldots, X_n) \quad B_n \quad f_n(X_1, \ldots, X_n) \quad 0)$$

$$\mathbb{P}(\vartheta \quad U \quad (X_1, \ldots, X_n) \quad B_n)$$

## The large deviation principle

Let $S$ be a Polish space equipped with the Borel $\sigma$-algebra $\mathscr{B}(S)$. A function $I: S \to [0, \infty]$ is a rate function if for each $M$ the level set $\{x \in S: I(x) \leqslant M\}$ is a compact subset of $S$. A rate function is necessarily a lower semicontinuous function, a function with closed level sets. A family $(Q_n)$ of probability measures on $S$ is defined to satisfy the large deviation principle with rate function $I$ if for each closed $F \subset S$

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

and for each open $G \subset S$

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

Large deviation theory focuses on probability measures $Q_n$ for which $Q_n(A)$ converges to 0 exponentially fast for a class of events $A$. The exponential decay of $Q_n(A)$ is characterized in terms of a rate function defined above. General treatments of the theory of large deviations and a wide variety of applications may be found in Dembo and Zeitouni (1998), Deuschel and Stroock (2000).

In analogous way, let us define the large deviation principle for regular conditional distributions. Let $(\Omega, \mathscr{F}, \mathbb{P})$ be a probability space, $(\mathscr{F}_n)$ a filtration of sub $\sigma$-algebras. We define a function $I: \Omega \times S \to [0, \infty]$ to be a rate function if for each $\omega$, $I(\omega, \cdot)$ is a rate function on $S$.

**Definition** 5 Suppose that $\{Q_n(B) \mid n \geqslant 1\}$ is a family of regular conditional distributions for a random variable taking values in $S$ given $\mathscr{F}_n$. We say that $\{Q_n(B) \mid n \geqslant 1\}$ satisfies the large deviation principle if for each closed set $F$ of $S$

$$\limsup_n \frac{1}{n} \log Q_n(F) \leqslant \inf_{x \in F} I(\ , x) \quad \text{a.s.} \tag{8}$$

and for each open set $G$ of $S$

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

In this paper we restrict ourselves to the analysis on the large deviation upper bound (8) for the posterior distributions of $\vartheta$ given $X_1, \ldots, X_n$. We will examine the posterior distributions $_n$ given $X_1, \ldots, X_n$ in the normal, Poisson and exponential cases and give a large deviation upper bound (8) explicitly for the posterior probability of the closed set $[\ ,\ )$ in each case.

## The normal case

Suppose that

$$P(dx)\quad f(x\ )dx\quad \frac{1}{\sqrt{2}}\exp\left(\ \frac{(x\ )^2}{2}\right)dx\qquad \mathbb{R}$$

and assume that the prior distribution for the normal mean $\vartheta$ is a conjugate distribution

$$(d\ )\quad \frac{1}{\sqrt{2}}\exp\left(\ \frac{(\ \mu)^2}{2\ ^2}\right)d\qquad 0\ \mu\ \mathbb{R}$$

It follows from Lemma 4 that the posterior distribution of $\vartheta$ given $X_1, \ldots, X_n$ is given by

$$_n(d\ )\quad \frac{\prod_{i\ 1}^n f(X_i\ )}{f_n(X_1, \ldots, X_n)}\ (d\ )$$

$$\frac{1}{\sqrt{2}\ _n}\exp\left[\ \frac{(\mu_n(X_1, \ldots, X_n)\ \mu)^2}{2\ _n^2}\right]d$$

where $\mu_n \quad \mu_n(x_1, \ldots, x_n)$ and $_n^2$ are defined by

$$\mu_n(x_1, \ldots, x_n)\quad \left(\frac{1}{1\ n\ ^2}\right)\mu\quad \left(\frac{n\ ^2}{1\ n\ ^2}\right)\bar{x}_n\quad \bar{x}_n\quad \frac{x_1\quad x_n}{n}$$

$$\frac{2}{n} \quad \frac{2}{1 \quad n \quad 2}$$

**Theorem** 6 *For each*

$$\limsup_{n} \frac{1}{n} \log {}_{n}[\ ,\ ) \leqslant \frac{(\quad \vartheta(\ ))^2}{2} \quad on \quad : \quad \vartheta(\ ) \quad a.s.$$

*Proof.* By Markov's inequality for conditional expectations, for all $t$　$0$

$$_{n}[\ ,\ ) \quad \mathbb{P}(\quad : \vartheta(\ ) \geqslant \quad X_1,..., X_n)(\ )$$

$$\mathbb{P}(e^{nt\vartheta(\ )}: \geqslant e^{nt} \quad X_1,..., X_n)(\ )$$

$$\leqslant e^{-nt} \quad \mathbb{E}(e^{nt\vartheta} \quad X_1,..., X_n)(\ )$$

$$e^{-nt} \quad \exp\left[\mu_n(X_1,..., X_n)nt \quad \frac{{}_{n}^{2}n^2t^2}{2}\right] \quad a.s.,$$

so that

$$\frac{1}{n} \log {}_{n}[\ ,\ ) \leqslant \quad t \quad \mu_n(X_1,..., X_n)t \quad \frac{{}_{n}^{2}nt^2}{2}$$

Since $\mu_n(X_1,..., X_n)$　$\mathbb{E}(X_1 \quad \vartheta)$　$\vartheta$ a.s. by Theorem A. 3and Lemma A. 1

we have

$$\limsup_{n} \frac{1}{n} \log {}_{n}[\ ,\ ) \leqslant \quad t \quad \vartheta(\ )t \quad \frac{t^2}{2}$$

Since $t$　$0$is arbitrary

$$\limsup_{n} \frac{1}{n} \log {}_{n}[\ ,\ ) \leqslant \inf_{t\ 0}\left[\quad t \quad \vartheta(\ )t \quad \frac{t^2}{2}\right]$$

$$\frac{(\quad \vartheta(\ ))^2}{2} \quad on \quad : \quad \vartheta(\ ) \quad a.s. \quad (9)$$

$$\square$$

In the same manner, it follows that

$$\limsup_{n} \frac{1}{n} \log {}_{n}(\quad ,\ ] \leqslant \frac{(\quad \vartheta(\ ))^2}{2} \quad on \quad : \quad \vartheta(\ ) \quad a.s.$$

In Theorem 6the rate function $I(\ ,\ )$ $(\ ,\ )$　×　is

$$I(\ ,\ )\quad \frac{(\ -\vartheta(\ ))^2}{2}\quad K(\vartheta(\ ),\ )$$

where $K(\ _1,\ _2)$ is the Kullback-Leibler distance

$$K(\ _1,\ _2)\qquad \log\frac{f(x\ _1)}{f(x\ _2)}f(x\ _1)\,dx\quad \frac{(\ _1\ _2)^2}{2}$$

If $\quad \vartheta(\ )\quad$ then

$$\frac{(\ -\vartheta(\ ))^2}{2}\quad \inf_{\geqslant} I(\ ,\ )$$

and so the large deviation upper bound inequality (9) is rewritten by using the rate function $I(\ ,\ )$ as

$$\limsup_{n}\frac{1}{n}\log\ _n[\ ,\ )\leqslant \inf_{\geqslant} I(\ ,\ )\quad \text{on}\quad :\quad \vartheta(\ )\quad \text{a.s.}$$

We now turn to the case where the samples are observed from the normal distribution with mean 0 and unknown precision. A precision is the reciprocal of the variance. Accordingly, we assume that

$$P(dx)\quad \left(\frac{}{2}\right)^{1/2}\exp\left(\ -\frac{x^2}{2}\right)dx\qquad (0\ )$$

If the prior distribution $\quad$ is specified by

$$(d\ )\quad \frac{}{(\ )}\quad {}^{1}e\quad 1_{(0\ )}d\qquad 0\quad 0$$

which is a gamma distribution with parameters $\quad$ and $\quad(\quad 0,\quad 0)$ then the posterior distribution of $\vartheta$ given $X_1,..., X_n$ is a gamma distribution with parameters

$$_n\qquad \frac{n}{2}\quad \text{and}\quad _n\quad _n(X_1,..., X_n)\qquad \frac{1}{2}\sum_{i\ 1}^{n} X_i^2$$

Theorem A. 3 together with Lemma A. 1 entails the convergence

$$\frac{\ _n}{n}\quad \frac{1}{2}\mathbb{E}(X_1^2\ )\quad \frac{1}{2\vartheta}\quad \text{a. s.}$$

**Theorem** 7  *For each*        1

$$\limsup_n \frac{1}{n}\log \ _n[\ ,\ ) \leqslant \ \frac{1}{2\vartheta(\ )}(\ \ 1\ \ \vartheta(\ )\log\ )$$

*on*    :    $\vartheta(\ )$   *a.s.*

*Proof*. For almost all            $\vartheta$    and $t$   $(0,\ 1/2\vartheta(\ ))$     there is an $n_0$ such

that   $_n/(\ _n\ \ nt)$    0 for all $n \geqslant n_0$   since

$$\frac{_n}{_n\ \ nt}\quad \frac{_n/n}{_n/n\ \ t}\quad \frac{1/(2\vartheta(\ ))}{1/(2\vartheta(\ ))\ \ t}\quad \frac{1}{1\ \ 2\vartheta(\ )t}$$

By Markov's inequality

$$\frac{1}{n}\log\ _n[\ ,\ )\ \leqslant\ \ t\quad \log \mathbb{E}\,(e^{nt}\ \ X_1,...,\ X_n)\,(\ )$$

$$t\quad \frac{n}{n}\log\left(\frac{_n(X_1,...,\ X_n)}{_n(X_1,...,\ X_n)\ \ nt}\right)$$

It follows that

$$\limsup_n \frac{1}{n}\log\ _n[\ ,\ )\ \leqslant\ \ t\quad \frac{1}{2}\log\left(\frac{1}{1\ \ 2\vartheta(\ )t}\right)$$

for almost all            $\vartheta$   and $t$   $(0,\ 1/2\vartheta(\ ))$    Now we obtain

$$\limsup_n \frac{1}{n}\log\ _n[\ ,\ )\ \leqslant\ \inf_{0\ \ t\ \ 1/2\vartheta(\ )}\left[\ \ t\quad \frac{1}{2}\log\left(\frac{1}{1\ \ 2\vartheta(\ )t}\right)\right]$$

$$\frac{1}{2\vartheta(\ )}(\ \ 1\ \ \vartheta(\ )\ \log\ )$$

on        $\vartheta$   a.s.


## The Poisson case

Let   $_0$ be the counting measure on   $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ and define   $(A)$        $_0(A\ \ $ 0,

$1,...\ )$   $A\ \ \mathscr{B}(\mathbb{R})$    Then    is a   -finite measure on $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$    If

$$P_\vartheta(dx) \quad f(x\mid\vartheta)\,\mu(dx) \quad \frac{e^{-\vartheta}\vartheta^x}{x!}\,\mu(dx) \qquad (\vartheta \in \Theta)$$

and the prior distribution  is a gamma distribution with parameters  and , then the posterior distribution of $\vartheta$ given $X_1,\dots,X_n$ is given by a gamma distribution with parameters  $_n \quad _n(X_1,\dots,X_n)$, $_n$. Here we define

$$_n \quad _n(x_1,\dots,x_n) \quad \sum_{i=1}^n x_i \quad _n \qquad _n$$

**Theorem 8** *For each*

$$\limsup_n \frac{1}{n}\log \quad _n[\ ,\ ) \leqslant (\quad \vartheta(\ )) \quad \vartheta(\ )\log\frac{}{\vartheta(\ )}$$

*on* $\quad:\quad \vartheta(\ ) \quad a.s.$

*Proof.* For all $t \in (0, 1)$  Markov's inequality yields

$$_n[\ ,\ ) \quad \mathbb{P}(\{\ :\vartheta(\ )\geqslant \ \mid X_1,\dots,X_n)(\ )$$
$$\leqslant e^{-nt}\ \mathbb{E}(e^{nt\vartheta}\mid X_1,\dots,X_n)(\ )$$
$$\leqslant e^{-nt}\ \left(\frac{_n}{_n\ nt}\right)^{_n(X_1,\dots,X_n)}\text{a.s.,}$$

and hence for all $t \in (0, 1)$

$$\limsup_n \frac{1}{n}\log \quad _n[\ ,\ ) \leqslant \quad t \quad \lim_n \frac{_n(X_1,\dots,X_n)}{n}\log\left(\frac{_n}{_n\ nt}\right)$$
$$t \quad \mathbb{E}(X_1\mid\vartheta)(\ )\log\left(\frac{1}{1\ t}\right)$$
$$t \quad \vartheta(\ )\log\left(\frac{1}{1\ t}\right)\quad\text{a.s.}$$

Thus on $\quad:\quad \vartheta(\ )$

$$\limsup_n \frac{1}{n}\log \quad _n[\ ,\ ) \leqslant \inf_{0\ t\ 1}\left[\quad t \quad \vartheta(\ )\log\left(\frac{1}{1\ t}\right)\right]$$
$$(\quad \vartheta(\ )) \quad \vartheta(\ )\log\frac{}{\vartheta(\ )}\quad\text{a.s.}$$

□

## The exponential case

Suppose that $(0, )$ and that for each

$$P(dx) \quad e^{-x}1_{(0, )}dx.$$

If the prior distribution is a gamma distribution with parameters and then the posterior distribution given $X_1, ..., X_n$ is a gamma distribution with parameters $_n$ and $_n$ $_n(X_1, ..., X_n)$ where

$$_n \quad n \quad _n \quad _n(x_1, ..., x_n) \quad \sum_{i=1}^{n} x_i$$

**Theorem** 9 *For each*

$$\limsup_n \frac{1}{n} {}_n[ \ , \ ) \leqslant 1 \quad \vartheta( ) \quad \log( \ \vartheta( ))$$

*on* $:$ $\vartheta( )$ *a.s.*

*Proof.* For almost all $\{ \ \vartheta\}$ and $t \ (0, \vartheta( ))$ there is an $n_0$ such that

$$\frac{_n(X_1, ..., X_n)}{_n(X_1, ..., X_n)} \quad 0 \text{ for all } n \geqslant n_0 \quad \text{since}$$

$$\frac{_n(X_1, ..., X_n)}{_n(X_1, ..., X_n)} \quad nt \quad \frac{\mathbb{E}(X_1 \ \vartheta)( )}{\mathbb{E}(X_1 \ \vartheta)( )} \quad \frac{\vartheta( )}{\vartheta( )} \quad 0$$

Thus for almost all $\{ \ \vartheta\}$ and all $t \ (0, \vartheta( ))$

$$\frac{1}{n}\log \ _n[ \ , \ ) \leqslant \quad t \quad \frac{-n}{n}\log\left(\frac{_n(X_1, ..., X_n)}{_n(X_1, ..., X_n) \quad nt}\right)$$

for all $n \geqslant n_0$ so that for $\{ \ \vartheta\}$ and $t \ (0, \vartheta( ))$

$$\limsup_n \frac{1}{n}\log \ _n[ \ , \ ) \leqslant \quad t \quad \log\left(\frac{\vartheta( )}{\vartheta( ) \quad t}\right)$$

Consequently

$$\limsup_n \frac{1}{n}\log \ _n[ \ , \ ) \leqslant \inf_{0 \ t \ \vartheta( )}\left[ \quad t \quad \log\left(\frac{\vartheta( )}{\vartheta( ) \quad t}\right)\right]$$

$$1 \quad (\ ) \quad \log (\ (\ ))$$

□

## Appendix

**Lemma A**. 1   Let $Y_1$ and $Y_2$ be random variables on $(\ , \mathscr{F}, \mathbb{P})$ with values in measurable spaces $(E_1, \mathscr{E}_1)$ and $(E_2, \mathscr{E}_2)$ respectively, and $\mathscr{G}$ a sub- -algebra with respect to which $Y_2$ is measurable. If $\mu$ is a regular conditional distribution for $Y_1$ given $\mathscr{G}$ then for every measurable function $f$ $E_1 \times E_2$ $\mathbb{R}$ such that $h(Y_1, Y_2)$ $L^1(\ , \mathscr{F}, \mathbb{P})$,

$$\int_{E_1} h(y_1, Y_2(\ ))\mu(\ , dy_1) \tag{A. 1}$$

is $\mathscr{G}$-measurable and

$$\mathbb{E}(h(Y_1, Y_2) \mathscr{G})(\ ) \quad \int_{E_1} h(y_1, Y_2(\ ))\mu(\ , dy_1) \quad \text{a.s.} \tag{A. 2}$$

In other words   (A. 1) is a version of $\mathbb{E}(h(Y_1, Y_2) \mathscr{G})$.

*Proof.* If $h$ $1_{A_1 \times A_2}$ $A_i$ $\mathscr{E}_i$, then (A. 1) is $\mathscr{G}$-measurable and (A. 2) holds. Since

$$\mathscr{H} \quad \left\{ A \quad \mathscr{E}_1 \times \mathscr{E}_2 \quad \int_{E_1} 1_A(y_1, Y_2(\ ))\mu(\ , dy_1) \text{ is a version of } \mathbb{E}(1_A(Y_1; Y_2) \mathscr{G})(\ ) \right\}$$

is a -class and $\mathscr{H}$ contains the -class

$$\mathscr{D} \quad \{A_1 \times A_2 \quad A_i \quad \mathscr{E}_i, i \quad 1, 2\}$$

$\mathscr{E}_1 \times \mathscr{E}_2$ $\mathscr{H}$ Thus (A. 1) is a version of $\mathbb{E}(h(Y_1, Y_2) \mathscr{G})$ whenever $h$ is an indicator function. By linearity   (A. 1) is a version of $\mathbb{E}(h(Y_1, Y_2) \mathscr{G})$ for all simple functions $h$, and hence for all nonnegative functions by the monotone convergence theorem. For the general case, the result follows by splitting the function into positive and negative parts.     □

Let $Y_1$, $Y_2$... be real-valued random variables defined on a probability space $(\ ,\mathscr{F}, \mathbb{P})$ and $\mathscr{G}$ a sub -algebra. If for all $n \geqslant 1$ and $A_1,..., A_n \quad \mathscr{B}(\mathbb{R})$

$$\mathbb{P}(Y_1 \quad A_1,..., Y_n \quad A_n \mid \mathscr{G}) \quad \prod_{i \ 1}^{n} \mathbb{P}(X_i \quad A_i \mid \mathscr{G}) \quad \text{a.s.,}$$

$Y_1$, $Y_2$... are declared conditionally independent given $\mathscr{G}$ If $\mathscr{G}$ ( ) for some random element $\quad Y_1$, $Y_2$... are called conditionally independent given $\quad$ In addition to the conditional independence, if for all $i \geqslant 1$ $\mathbb{P}(Y_i \quad A$

$\mid \mathscr{G}) \quad \mathbb{P}(Y_1 \quad A \mid \mathscr{G})$ a.s., $Y_1$, $Y_2$... are defined to be conditionally independent and identically distributed (abbreviated to conditionally i.i.d.) given $\mathscr{G}$. If $Y_1$, $Y_2$... are conditionally i.i.d. and $\varphi$ is a measurable function, then $\varphi$ $(Y_1)$, $\varphi(Y_2)$,... are conditionally i.i.d.

**Lemma A. 2** If $Y_1$, $Y_2$... are conditionally i.i.d. given $\mathscr{G}$, there exists a regular conditional distribution $\mu(\ , B), (\ , B) \quad \times \mathscr{B}(\mathbb{R})$ for $Y \quad (Y_1, Y_2...)$ given $\mathscr{G}$ such that for each $\quad$ the coordinate functions $_1$, $_2$... on $(\mathbb{R}\ , \mathscr{B}(\mathbb{R}\ ), \mu(\ , \ ))$ are i.i.d. Moreover, if $Y_1$ is integrable, then $_1$, $_2$... are integrable with respect to $\mu(\ , \ )$ for almost all

*Proof.* Since $\mathbb{R}\ $ is a Borel space, there is a regular conditional distribution $_0$ $(\ , B)$ for $Y \quad (Y_1, Y_2...)$ given $\mathscr{G}$. For each $i \geqslant 1$ and each $r \quad \mathbb{Q}$ there is a null set $N_{i,r} \quad \mathscr{G}$ such that for each $\quad / N_{i,r}$

$$_0(\ , \ _i \leqslant r) \qquad _0(\ , \mathbb{R} \times \quad \times \mathbb{R} \times (\quad , r] \times \mathbb{R} \times \quad )$$
$$\mathbb{P}(Y \quad \mathbb{R} \times \quad \times \mathbb{R} \times (\quad , r] \times \mathbb{R} \times \quad \mid \mathscr{G})(\ )$$
$$\mathbb{P}(Y_i \leqslant r \mid \mathscr{G})(\ ) \quad \mathbb{P}(Y_1 \leqslant r \mid \mathscr{G})(\ )$$
$$_0(\ , \ _1 \leqslant r)$$

and hence for all $\quad / N \quad \bigcup_{i \geqslant 1, r \ \mathbb{Q}} N_{i,r}$ and for all $i \geqslant 1, r \quad \mathbb{Q}$, we have

$$_0(\ , \ _i \leqslant r) \qquad _0(\ , \ _1 \leqslant r).$$

Since the sets of the form ( , $r$] , $r$   ℚ form a    -class generating $\mathscr{B}(\mathbb{R})$ it follows that for each    / $N$    ₀( , ᵢ ) and ₀( , ₁ ) agree as probability measures on $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$    For each    define a measure    by

$$
(\;) \quad \begin{cases} _0(\; , \;_1\;) & / N \\ (\;) & N \end{cases}
$$

where    is any probability measure on $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$    Now we define a probability measure

$$
\mu (\; , \;) \quad (\; \times \quad \times \quad )(\;)
$$

for each    on $(\mathbb{R}\; , \mathscr{B}(\mathbb{R}\; ))$    We will show that $\mu$ is a regular conditional distribution given $\mathscr{S}$ that satisfies the requirement of the theorem. Since $\mu$ ( , ) is the infinite-dimensional product measure of    with itself, the coordinate functions    ₁, ₂... are necessarily i.i.d. random variables on $(\mathbb{R}\; , \mathscr{B}(\mathbb{R}\; ), \mu (\; , \;))$ for each    with distribution

$$
\mu (\; , \;_i \; A) \quad (A) \quad \begin{cases} _0(\; , \;_1 \; A) & / N \\ (A) & N \end{cases} \quad A \quad \mathscr{B}(\mathbb{R})
$$

To show that $\mu$ ( , B) is a regular conditional distribution for $Y$  ( $Y_1$, $Y_2$...) given $\mathscr{S}$ it suffices to verify that $\mu$ ( , B) is a version of $\mathbb{P}(Y \; B \; \mathscr{S})$ for each $B$   $\mathscr{B}(\mathbb{R}\; )$    since $\mu$ ( , ) is a probability measure by definition. If $A_1$,..., $A_n$   $\mathscr{B}(\mathbb{R})$ , $n \geqslant 1$   then

$$
\mu (\; , A_1 \times \quad \times A_n \times \mathbb{R} \times \quad) \quad (A_1) \quad (A_n)\, 1_{N^c} \quad (A_1) \quad (A_n)\, 1_N
$$
$$
_0(\; , \;_1 \; A_1) \quad _0(\; , \;_1 \; A_n)\, 1_{N^c} \quad (A_1) \quad (A_n)\, 1_N,
$$

and therefore $\mu$ ( , $A_1 \times$    $\times A_n \times \mathbb{R} \times$    ) is $\mathscr{S}$ measurable. Besides outside

the $\mathscr{G}$-null set $N$

$$\mu(\ ,A_1\times\ \cdots\times A_n\times\mathbb{R}\times\ )\quad {}_0(\ ,\ {}_1\ A_1)\qquad {}_0(\ ,\ {}_1\ A_n)$$

$$_0(\ ,\ {}_1\ A_1)\qquad {}_0(\ ,\ {}_n\ A_n)$$

$$\mathbb{P}(Y_1\ A_1\ \mathscr{A})(\ )\quad \mathbb{P}(Y_n\ A_n\ \mathscr{A})(\ )$$

$$\mathbb{P}(Y_1\ A_1,...,\ Y_n\ A_n\ \mathscr{A})(\ )$$

$$\mathbb{P}(Y\ A_1\times\quad A_n\times\mathbb{R}\times\qquad \mathscr{A})(\ )\ \text{a. s.}$$

Therefore $\mu(\ ,A_1\times\ \cdots\times A_n\times\mathbb{R}\times\ )$ is a version of $\mathbb{P}(Y\ A_1\times\quad A_n\times\mathbb{R}$

$\times\quad \mathscr{A})$    Note that

$$\mathscr{D}\quad\{A_1\times\quad A_n\times\mathbb{R}\times\quad :n\geqslant 1,A_i\ \mathscr{B}(\mathbb{R}),i\ 1,...,\ n\}$$

is a  -class that generates $\mathscr{B}(\mathbb{R})$    Since

$$\mathscr{H}\quad\{B\ \mathscr{B}(\mathbb{R}):\mu(\ ,B)\text{ is a version of }\mathbb{P}(Y\ B\ \mathscr{A})\}$$

is a  -class with $\mathscr{D}\ \mathscr{H}\ \mathscr{B}(\mathbb{R})\quad\mathscr{H}$   This implies that $\mu(\ ,B)$ is a version of $\mathbb{P}(Y\ B\ \mathscr{A})$ for each $B\ \mathscr{B}(\mathbb{R})$.

Finally by Lemma A. 1

$$_{\mathbb{R}}\quad {}_i(y)\ \mu(\ ,dy)\qquad_{\mathbb{R}}\quad {}_1(y)\ \mu(\ ,dy)$$

$$\mathbb{E}(\ {}_1(Y)\ \mathscr{A})(\ )\quad \mathbb{E}(\ Y_1\ \mathscr{A})(\ )\quad\text{a.s.}$$

The integrability of $Y_1$ entails $\mathbb{E}(\ Y_1\ \mathscr{A})(\ )$      a.s., and hence the claims follows. This completes the proof.                                □

**Theorem A**. 3   If $Y_1,\ Y_2,...$ are conditionally i.i.d. random variables given a sub  -algebra $\mathscr{G}$ and if $Y_1$ is integrable, then

$$\mathcal{Y}_n\quad \frac{Y_1\quad Y_n}{n}\quad \mathbb{E}(Y_1\ \mathscr{A})\quad\text{a.s. }(n\quad ).$$

*Proof.* Let $\mu\ (B)\quad\mu(\ ,B),(\ ,B)\quad\times\mathscr{B}(\mathbb{R})$ be a regular conditional dis-

tribution for $Y$ $(Y_1, Y_2...)$ given $\mathscr{G}$ such that the coordinate functions $_1$, $_2$... are i.i.d. random variables on $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \mu)$ for each    We will show that

$$\mathbb{P}\left(\sup_{n \geq m} \, \overline{Y_n} \quad \mathbb{E}(Y_1 \mid \mathscr{G}) \quad \right) \quad 0 \quad (m \quad), \tag{A. 3}$$

which is equivalent to the convergence $\overline{Y_n}$ $\mathbb{E}(Y_1 \mid \mathscr{G})$ a.s. as $n$    For all    0

$$\mathbb{P}\left(\sup_{n \geq m} \, \overline{Y_n} \quad \mathbb{E}(Y_1 \mid \mathscr{G}) \quad \right) \quad \mathbb{E}\left[\mathbb{P}\left(\sup_{n \geq m} \, \overline{Y_n} \quad \mathbb{E}(Y_1 \mid \mathscr{G}) \quad \mid \mathscr{G}\right)\right]$$

$$\mathbb{E}\left[\mathbb{P}\left(\sup_{n \geq m} \, \frac{1}{n} \sum_{i=1}^{n} \, _i(Y) \quad \mathbb{E}(Y_1 \mid \mathscr{G}) \quad \mid \mathscr{G}\right)\right]$$

$$\mathbb{E}\left[\mu \left\{ y \in \mathbb{R} : \sup_{n \geq m} \frac{1}{n} \sum_{i=1}^{n} \, _i(y) \quad \mathbb{E}(Y_1 \mid \mathscr{G})(\,) \quad \right\}\right].$$

The last equation follows from Lemma A. 1 Since $Y_1$ is assumed to be integrable, Lemma A. 2 shows that $_1$, $_2$,... are i.i.d. integrable random variables on $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \mu)$ for almost all    It follows by the strong law of large numbers and Lemma A. 1 that

$$\frac{1}{n} \sum_{i=1}^{n} \, _i \quad \int_{\mathbb{R}} \, _1 d\mu \quad \mathbb{E}(\, _1(Y) \mid \mathscr{G})(\,)$$
$$\mathbb{E}(Y_1 \mid \mathscr{G})(\,) \quad \mu \text{-a.s.}$$

for almost all    It follows that

$$\mu \left\{ y \in \mathbb{R} : \sup_{n \geq m} \frac{1}{n} \sum_{i=1}^{n} \, _i(y) \quad \mathbb{E}(Y_1 \mid \mathscr{G})(\,) \quad \right\} \quad 0$$

for almost all    And now (A. 3) is obtained by the dominated convergence theorem.

### References

[ 1 ] Aldous, D. J ( 1982) On exchangeabilitiy and conditional independence, In Koch, G. and Spizzichino, F., editors, *Exchangeabilitiy in Probability and Statistics* 165-170

North-Holland, Amsterdam.

[ 2 ] Dembo, A. and Zeitouni, O.（1998）*Large Deviations Techniques and Applications*
2nd ed., Springer-Verlag, New York.

[ 3 ] Deuschel, J. D. and Stroock, D. W.（2000）*Large Deviations*, AMS Chelsea Publish-
ing, Amer. Math. Soc.

[ 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.