

CONDITIONING BY RARE SOURCES

M. GRENDÁR

To Mar, in memoriam.

To George Judge, on the occasion of his eightieth birthday.

ABSTRACT. In this paper we study the exponential decay of posterior probability of a set of sources and conditioning by rare sources for both uniform and general prior distributions of sources. The decay rate is determined by L -divergence and rare sources from a convex, closed set asymptotically conditionally concentrate on an L -projection. L -projection on a linear family of sources belongs to A -family of distributions. The results parallel those of Large Deviations for Empirical Measures (Sanov's Theorem and Conditional Limit Theorem).

1. INTRODUCTION

Information divergence minimization, which is also known as Relative Entropy Maximization or MaxEnt method, has – thanks to Large Deviations Theorems for Empirical Measures – gained a firm probabilistic footing, which justifies its application in the area of the convex Boltzmann Jaynes Inverse Problem (the α -problem, for short). For the β -problem – an 'antipode' of the α -problem – Large Deviations Theorems for Sources single out the L -divergence minimization method.

The paper is organized as follows: First, necessary terminology and notation are introduced. A brief survey of Large Deviations Theorems for Empirical Measures that includes Sanov's Theorem and a Conditional Limit Theorem is given next. Then, a set-up for a study of conditioning by rare sources is formulated and Sanov's Theorem and the Conditional Limit Theorem for Sources are stated; under various assumptions. Next, Theorems are proven for the continuous case and the results are applied to a criterion choice problem associated with the β -problem. An End-Notes section points to relevant literature and contains further discussion.

2000 Mathematics Subject Classification. Primary 60F10; Secondary 94A15.

Key words and phrases. Conditional Limit Theorems, Information projection, Reverse information projection, L -projection, Kerridge's inaccuracy, Watanabe's confirmability, Large Deviations for Sources, Maximum Non-parametric Likelihood, Ill-posed inverse problems, Criterion Choice Problem, Bayesian nonparametric consistency.

Supported in part by Slovak VEGA Grant 1/0424/03 and Slovak Grant APVT-20-004104.

I have greatly benefited from valuable discussions with George Judge. Without implicating him, I am indebted to Ľubomír Snoha for clarifications of a couple of important technical questions. Substantive comments and suggestions from three anonymous referees are gratefully acknowledged. Supported by VEGA 1/7295/20 and 1/3016/06 grants.

Submitted: August 31, 2005

2. TERMINOLOGY AND NOTATION

Let $\mathcal{P}(\mathcal{X})$ be a set of all probability mass functions on a finite alphabet $\mathcal{X} \triangleq \{x_1, x_2, \dots, x_m\}$ of m letters. The support of $p \in \mathcal{P}(\mathcal{X})$ is a set $S(p) \triangleq \{x : p(x) > 0\}$.

A probability mass function (pmf) from $\mathcal{P}(\mathcal{X})$ is rational if it belongs to the set $\mathcal{R} \triangleq \mathcal{P}(\mathcal{X}) \cap \mathbb{Q}^m$. A rational pmf is n -rational, if denominators of all its m elements are n . The set of all n -rational pmf's will be denoted by \mathcal{R}_n .

Let x_1, x_2, \dots, x_n be a sequence of n letters, that is identically and independently drawn from a source $q \in \mathcal{P}(\mathcal{X})$. Type and n -type are other names for empirical measures induced by a sequence of the length n . Formally, type $\nu^n \triangleq [n_1, n_2, \dots, n_m]/n$, where n_i is the number of occurrences of i -th letter of the alphabet in the sequence. Note that there are $\Gamma(\nu^n) \triangleq n!(\prod_{i=1}^m n_i!)^{-1}$ different sequences of length n , which induce the same type ν^n . $\Gamma(\nu^n)$ is called the multiplicity of type. Finally, observe that ν^n is n -rational; $\nu^n \in \mathcal{R}_n$.

Let $\Pi \subseteq \mathcal{P}(\mathcal{X})$. $\Pi_n \triangleq \Pi \cap \mathcal{R}_n$.

The information divergence (\pm -relative entropy, Kullback-Leibler distance etc.) $I(p||q)$ of p with respect to q (both from $\mathcal{P}(\mathcal{X})$) is $I(p||q) \triangleq \sum_{\mathcal{X}} p \log \frac{p}{q}$, with conventions that $0 \log 0 = 0$, $\log b/0 = +\infty$. The information projection \hat{p} of q on Π is $\hat{p} \triangleq \arg \inf_{p \in \Pi} I(p||q)$. The value of the I -divergence at an I -projection of q on Π is denoted by $I(\Pi||q)$.

On $\mathcal{P}(\mathcal{X})$ topology induced by the standard topology on \mathbb{R}^m is assumed.

The support $S(\mathcal{C})$ of a convex set $\mathcal{C} \subset \mathcal{P}(\mathcal{X})$ is just the support of the member of \mathcal{C} for which $S(\cdot)$ contains the support of any other member of the set.

The following families of distributions will be needed:

1) Linear family $\mathcal{L}(u, a) \triangleq \{p : \sum_{\mathcal{X}} p(x)u_j(x) = a_j, j = 1, 2, \dots, k\}$, where u_j is a real-valued function on \mathcal{X} and $a_j \in \mathbb{R}$.

2) Exponential family $\mathcal{E}(\rho, u, \theta) \triangleq \{p : p(x) = z\rho(x) \exp(\sum_{j=1}^k \theta_j u_j(x)), x \in \mathcal{X}\}$, where a normalizing factor $z \triangleq \sum_{\mathcal{X}} \rho(x) \exp(\sum_{j=1}^k \theta_j u_j(x))$ and ρ belongs to $\mathcal{P}(\mathcal{X})$; $\theta_j \in \mathbb{R}$.

3) Λ -family $\Lambda(\rho, u, \theta, a) \triangleq \{p : p(x) = \rho(x)[1 - \sum_{j=1}^k \theta_k (u_j(x) - a_j)]^{-1}, x \in \mathcal{X}\}$.

The definitions of the families can be extended to continuous \mathcal{X} in a straightforward way.

In what follows, $r \in \mathcal{P}(\mathcal{X})$ will be the 'true' source of sequences and hence types.

3. CONDITIONING BY RARE TYPES

It is convenient to begin with a brief survey of the Large Deviations Theorems for Empirical Measures (Sanov's Theorem and a Conditional Limit Theorem).

First, it is necessary to introduce the probability $\pi(\nu^n; r)$ that the source r generates an n -type ν^n . The probability that r generates a sequence of n letters x_1, x_2, \dots, x_n which induces a type ν^n is $\prod_{i=1}^m (r_i)^{n\nu_i^n}$. As it was already mentioned, there is a number $\Gamma(\nu^n)$ of sequences of length n , which induce the same type ν^n . The probability $\pi(\nu^n; r)$ that r generates type ν^n is thus $\pi(\nu^n; r) \triangleq \Gamma(\nu^n) \prod_{i=1}^m (r_i)^{n\nu_i^n}$. Consequently, for $A \subseteq B \subseteq \mathcal{P}(\mathcal{X})$, $\pi(\nu^n \in A | \nu^n \in B; r) = \frac{\pi(\nu^n \in A; r)}{\pi(\nu^n \in B; r)}$; provided that $\pi(\nu^n \in B; r) \neq 0$.

Π is rare if it does not contain r . Given that the source r produced an n -type from rare Π , it is of interest to know how the conditional probability/measure

spreads among the rare n -types from Π ; especially as n grows beyond any limit. For the rare set of a particular form, this issue is answered by Conditional Limit Theorem (CoLT) which is also known as Conditional Weak Law of Large Numbers.

CoLT can be established by means of Sanov's Theorem (ST).

ST. ([6] *Thm 3*) *Let Π be a set such that its closure is equal to the closure of its interior. Let r be such that $S(r) = \mathcal{X}$. Then,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(\nu^n \in \Pi; r) = -I(\Pi || r).$$

Sanov's Theorem (ST) states that the probability $\pi(\nu^n \in \Pi; r)$ decays exponentially fast, with the decay rate given by the value of the information divergence at an I -projection of the source r on Π .

CoLT. ([8] *Thm 4.1*, [2] *Thm 12.6.2*) *Let Π be a convex, closed rare set. Let $B(\hat{p}, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at I -projection \hat{p} of r on Π . Then for any $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} \pi(\nu^n \in B(\hat{p}, \epsilon) | \nu^n \in \Pi; r) = 1.$$

Informally, CoLT states that if a dense rare set admits a unique I -projection, then asymptotically types conditionally concentrate just on it. Thus, provided that for sufficiently large n a type from rare Π occurred, with probability close to 1 it is just a type close to \hat{p} . Numeric examples of ST and CoLT can be found at [2].

This suggests that, conditionally upon the rare Π , it is the I -projection \hat{p} rather than r , which should be considered as the true *iid* source of data. Gibbs' Conditioning Principle (GCP) - an important strengthening of CoLT - captures this 'intuition'; cf [3], [7].

If $S(\mathcal{L}) = \mathcal{X}$ then the I -projection \hat{p} of r on $\Pi \equiv \mathcal{L}$ is unique and belongs to the exponential family of distributions $\mathcal{E}(r, u, \theta)$; i.e., $\mathcal{L}(u, a) \cap \mathcal{E}(r, u, \theta) = \{\hat{p}\}$.

4. CONDITIONING BY RARE SOURCES

In the above setting there is a fixed source r and a rare set Π_n of n -types. We now consider an opposite setting where the n -type is unique, and there is a set $\mathcal{Q}_n \triangleq \mathcal{Q} \cap \mathcal{R}_n$, where $\mathcal{Q} \subseteq \mathcal{P}(\mathcal{X})$, of rare n -sources of the type.

Furthermore, n -sources q^n are assumed to have prior distribution $\pi(q^n)$. If from \mathcal{R}_n n -source q^n occurs, then the source generates n -type ν^n with the probability $\pi(\nu^n | q^n) \triangleq \Gamma(\nu^n) \prod_{i=1}^m (q_i^n)^{n\nu_i^n}$.

We are interested in the asymptotic behavior of the probability $\pi(q^n \in B | (q^n \in \mathcal{Q}) \wedge \nu^n)$ that if the n -type ν^n and an n -source q^n from a rare set \mathcal{Q} occurred, then the n -source belongs to a subset B of \mathcal{Q} . Note that $\pi(q^n \in B | (q^n \in \mathcal{Q}) \wedge \nu^n) = \frac{\pi(q^n \in B | \nu^n)}{\pi(q^n \in \mathcal{Q} | \nu^n)}$; provided that $\pi(q^n \in \mathcal{Q} | \nu^n) > 0$. The posterior probability $\pi(q^n | \nu^n)$ is related to the defined probabilities $\pi(\nu^n | q^n)$ and $\pi(q^n)$ via Bayes's Theorem.

Asymptotic investigations will be first carried on under the assumption of uniform prior distribution of n -sources (Sect. 4.1). The assumption will be relaxed in Section 4.2. Within each of the sections, two cases of convergence will be considered: a static and a dynamic case. For the static case asymptotic investigations are carried over a subsequence of types, which are k -equivalent to ν^{n_0} . A type

$\nu^{kn_0} \triangleq [kn_1, \dots, kn_m]/kn_0$, $k \in \mathbb{N}$, is called k -equivalent to ν^{n_0} . The dynamic case assumes that there is a sequence of n -types which converges in the total variation to some $p \in \mathcal{P}(\mathcal{X})$. For each case what is meant by rare source will be defined separately.

For $p, q \in \mathcal{P}(\mathcal{X})$, the L -divergence $L(q||p)$ of q with respect to p is the map $L : \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R} \cup \{\infty\}$, $L(q||p) \triangleq -\sum_{\mathcal{X}} p \log q$. The L -projection \hat{q} of p on set of sources \mathcal{Q} is: $\hat{q} \triangleq \arg \inf_{q \in \mathcal{Q}} L(q||p)$. The value of L -divergence at an L -projection of p on \mathcal{Q} (i.e., $\inf_{q \in \mathcal{Q}} L(q||p)$) is denoted by $L(\mathcal{Q}||p)$.

4.1 Uniform prior.

Within this section it is assumed that n -sources have a uniform prior distribution. Since there is total $N = \binom{n+m-1}{m-1}$ n -sources (cf. [4]), the uniform prior probability $\pi(q_n) = 1/N$, for all $q^n \in \mathcal{R}_n$.

4.1.1 Static case.

Let there be an n_0 -type ν^{n_0} . A set \mathcal{Q} of sources is rare if it does not contain ν^{n_0} .

Sanov's Theorem for Sources (abbreviated LST) is a counterpart of the Sanov's Theorem for Types.

Static LST. *Let ν^{n_0} be a type. Let \mathcal{Q} be an open set of sources. Then, for $n \rightarrow \infty$ over a subsequence $n = kn_0$, $k \in \mathbb{N}$,*

$$\frac{1}{n} \log \pi(q^n \in \mathcal{Q}|\nu^n) = -\{L(\mathcal{Q}|\nu^{n_0}) - L(\mathcal{P}|\nu^{n_0})\}.$$

Proof. Under the assumption of uniform prior distribution of n -sources

$$\log \pi(q^n \in \mathcal{Q}|\nu^n) = \log \sum_{q^n \in \mathcal{Q}} \prod_{\mathcal{X}} (q^n)^{n\nu^n} - \log \sum_{q^n \in \mathcal{P}} \prod_{\mathcal{X}} (q^n)^{n\nu^n}.$$

Since $N < (n+1)^m$ (cf. Lemma 2.1.2 of [7]), $\frac{1}{n_0} \log \pi(q^{n_0} \in \mathcal{Q}|\nu^{n_0})$ can be bounded from above and below as:

$$\begin{aligned} -L(\mathcal{Q}_{n_0}|\nu^{n_0}) + L(\mathcal{R}_{n_0}|\nu^{n_0}) - \frac{m}{n_0} \log(n_0 + 1) &\leq \frac{1}{n_0} \log \pi(q^{n_0} \in \mathcal{Q}|\nu^{n_0}) \leq \\ &\leq -L(\mathcal{Q}_{n_0}|\nu^{n_0}) + L(\mathcal{R}_{n_0}|\nu^{n_0}) + \frac{m}{n_0} \log(n_0 + 1). \end{aligned}$$

Fix $p \in \mathcal{P}(\mathcal{X})$. Equip $\mathbb{R} \cup \{\infty\}$ with the standard topology (i.e., the topology induced by the total order). As for each open subset A of $\mathbb{R} \cup \{\infty\}$, $L^{-1}(A)$ is an open subset of $\mathcal{P}(\mathcal{X})$, the L -divergence is continuous in q .

\mathcal{Q} is open by the assumption.

Thus, $L(\mathcal{Q}_{n_0}|\nu^{n_0})$ converges to $L(\mathcal{Q}|\nu^{n_0})$ as $n \rightarrow \infty$, $n = kn_0$, $k \in \mathbb{N}$. Also, $L(\mathcal{R}_{n_0}|\nu^{n_0})$ converges to $L(\mathcal{P}|\nu^{n_0})$ for $n \rightarrow \infty$, $n = kn_0$, $k \in \mathbb{N}$. \square

The Law of Large Numbers for Sources (LLLN) is a direct consequence of LST.

Static LLLN. *Let ν^{n_0} be a type. Let \hat{q} be L -projection of ν^{n_0} on $\mathcal{P}(\mathcal{X})$. And let $B(\hat{q}, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at \hat{q} . Then, for $\epsilon > 0$ and $n \rightarrow \infty$ over the types which are k -equivalent with ν^{n_0} ,*

$$\pi(q^n \in B(\hat{q}, \epsilon) | (q^n \in \mathcal{P}) \wedge \nu^n) = 1.$$

Proof. Let $B^C(\hat{q}, \epsilon) \triangleq \mathcal{P}(\mathcal{X}) \setminus B(\hat{q}, \epsilon)$. Since $B^C(\hat{q}, \epsilon)$ is open by the assumption, LST can be applied to it. Since $B^C \subset \mathcal{P}$, $L(B^C|\nu^{n_0}) - L(\mathcal{P}|\nu^{n_0}) > 0$. Thus, $\pi(q^n \in B^C(\hat{q}, \epsilon)|\nu^n)$ converges to 0, as $n \rightarrow \infty$ over a subsequence of $n = kn_0$, $k \in \mathbb{N}$. \square

Obviously, the L -projection \hat{q} of ν^{n_0} on $\mathcal{P}(\mathcal{X})$ is $\hat{q} \equiv \nu^{n_0}$.

LLN is a special case of the Conditional Limit Theorem for Sources (LCoLT), which is a consequence of LST, as well.

Static LCoLT. Let ν^{n_0} be a type. Let \mathcal{Q} be a convex, closed rare set of sources. Let \hat{q} be the L -projection of ν^{n_0} on \mathcal{Q} and let $B(\hat{q}, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at \hat{q} . Then, for $\epsilon > 0$ and $n \rightarrow \infty$ over a subsequence $n = kn_0$, $k \in \mathbb{N}$,

$$\pi(q^n \in B(\hat{q}, \epsilon) | (q^n \in \mathcal{Q}) \wedge \nu^n) = 1.$$

Proof. Let $B^C(\hat{q}, \epsilon) \triangleq \mathcal{P}(\mathcal{X}) \setminus B(\hat{q}, \epsilon)$. Clearly,

$$\log \pi(q^{n_0} \in B^C(\hat{q}, \epsilon) | (q^{n_0} \in \mathcal{Q}) \wedge \nu^{n_0}) = \log \pi(q^{n_0} \in B^C | \nu^{n_0}) - \log \pi(q^{n_0} \in \mathcal{Q} | \nu^{n_0}).$$

Since both $B^C(\hat{q}, \epsilon)$ and \mathcal{Q} are open, LST can be applied. As $B^C(\hat{q}, \epsilon) \subset \mathcal{Q}$, $L(B^C|\nu^{n_0}) - L(\mathcal{Q}|\nu^{n_0}) > 0$. Hence $\pi(q^n \in B^C | (q^n \in \mathcal{Q}) \wedge \nu^n)$ converges to 0, as $n \rightarrow \infty$ over a subsequence of $n = kn_0$, $k \in \mathbb{N}$. Since under the assumptions on \mathcal{Q} the L -projection of ν^{n_0} on \mathcal{Q} is unique, the claim of the Theorem follows. \square

Example. Let $\mathcal{X} = \{1, 2, 3, 4\}$. Let $\mathcal{Q} = \{q : \sum_{x \in \mathcal{X}} q(x)x = 1.7\}$. Let $n_0 = 10$ and $\nu^{n_0} = [1, 1, 1, 7]/10$. The L -projection of ν^{n_0} on \mathcal{Q} is $\hat{q} = [0.705, 0.073, 0.039, 0.183]$. Let $\epsilon = 0.1$. The concentration of n -sources on the L -projection, which is captured by the Static LCoLT, is for types k -equivalent to ν^{n_0} ($k = 5, 10, 20, 30$) illustrated in Table 1.

TABLE 1. Values of $\pi(q^n \in B(\hat{q}, \epsilon) | (q^n \in \mathcal{Q}) \wedge \nu^n)$ for $n = kn_0$, $k = 5, 10, 20, 30$.

n	$\pi(\cdot \cdot)$
50	0.868
100	0.948
200	0.994
300	0.999

The L -projection at the above Example can be found by means of the following Proposition.

Proposition. Let $\mathcal{Q} \equiv \mathcal{L}(u, a)$. Let $p \in \mathcal{P}(\mathcal{X})$ be such that $S(p) = S(\mathcal{L})$. Then the L -projection \hat{q} of p on \mathcal{Q} is unique and belongs to $\Lambda(p, u, \theta, a)$ family; i.e., $\mathcal{L}(u, a) \cap \Lambda(p, u, \theta, a) = \{\hat{q}\}$.

Proof. In light of Theorem 9 of [6] it suffices to check that $\hat{q} = p[1 - \sum_{j=1}^k \theta_k(u_j(x) - a_j)]^{-1}$, with θ such that $\hat{q} \in \mathcal{L}(u, a)$, satisfies:

$$\sum_{S(p)} p \left(1 - \frac{q'}{\hat{q}} \right) = 0,$$

for all $q' \in \mathcal{Q}$, which is indeed the case. \square

4.1.2 Dynamic case.

Let there be a sequence of n -types which converges in the total variation to a pmf $p \in \mathcal{P}(\mathcal{X})$, denoted as $\nu^n \rightarrow p$. In this case, a set \mathcal{Q} of sources is rare if it does not contain p .

Dynamic LST. Let $\nu^n \rightarrow p$. Let \mathcal{Q} be an open set of sources. Then,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{Q} | \nu^n) = -\{L(\mathcal{Q} || p) - L(\mathcal{P} || p)\}.$$

Dynamic LLLN. Let $\nu^n \rightarrow p$. Let \hat{q} be L -projection of p on \mathcal{P} . And let $B(\hat{q}, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at \hat{q} . Then, for $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \pi(q^n \in B(\hat{q}, \epsilon) | (q^n \in \mathcal{P}) \wedge \nu^n) = 1.$$

Dynamic LCoLT. Let $\nu^n \rightarrow p$. Let \mathcal{Q} be a convex, closed rare set of sources. Let \hat{q} be the L -projection of p on \mathcal{Q} and let $B(\hat{q}, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at \hat{q} . Then, for $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \pi(q^n \in B(\hat{q}, \epsilon) | (q^n \in \mathcal{Q}) \wedge \nu^n) = 1.$$

Proofs can be constructed along the lines for the static case.

4.2 General prior.

Let $\pi(q)$ be a prior pmf on \mathcal{R} . From this pmf, a prior distribution $\pi^{\mathcal{A}}(q^n)$ on \mathcal{R}_n is constructed by a quantization $\mathcal{A} \triangleq \{A_1, A_2, \dots, A_N\}$ of \mathcal{R} into disjoint sets, such that each $A \in \mathcal{A}$ contains just one q_n from \mathcal{R}_n . Then $\pi^{\mathcal{A}}(q^n) \triangleq \pi(\{A_j : q^n \in A_j, j = 1, 2, \dots, N\})$.

Let $\mathcal{S} \triangleq S(\pi(\cdot))$. Let $\mathcal{Q}^\pi \triangleq \mathcal{Q} \cap \mathcal{S}$, $\mathcal{P}^\pi \triangleq \mathcal{P} \cap \mathcal{S}$.

As the static case is subsumed under the dynamic one, only the latter limit theorems will be presented.

General prior LST. Let $\nu^n \rightarrow p$. Let \mathcal{Q}^π be an open set of sources. Then,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi^{\mathcal{A}}(q^n \in \mathcal{Q}^\pi | \nu^n) = -\{L(\mathcal{Q}^\pi || p) - L(\mathcal{P}^\pi || p)\}.$$

Proof. For a zero-prior-probability n -source, the posterior probability is zero as well; so such sources can be excluded from considerations. Let $\mathcal{S}_n \triangleq S(\pi^{\mathcal{A}}(q_n))$, $\mathcal{Q}_n^\pi \triangleq \mathcal{Q} \cap \mathcal{S}_n$, $\mathcal{P}_n^\pi \triangleq \mathcal{P} \cap \mathcal{S}_n$.

$$\log \pi^{\mathcal{A}}(q^n \in \mathcal{Q} | \nu^n) = \log \sum_{q^n \in \mathcal{Q}_n^\pi} \pi^{\mathcal{A}}(q^n) \prod_{\mathcal{X}} (q^n)^{n\nu^n} - \log \sum_{q^n \in \mathcal{P}_n^\pi} \pi^{\mathcal{A}}(q^n) \prod_{\mathcal{X}} (q^n)^{n\nu^n}.$$

Denote by $\lambda(\mathcal{Q}_n^\pi || \nu^n) \triangleq \inf_{q^n \in \mathcal{Q}_n^\pi} \lambda(q^n || \nu^n)$, where $\lambda(q^n || \nu^n) \triangleq L(q^n || \nu^n) - \frac{1}{n} \log \pi^{\mathcal{A}}(q^n)$. Using this notation and invoking the same argument as in the proof

of LST for uniform prior, $\frac{1}{n} \log \pi^{\mathcal{A}}(q^n \in \mathcal{Q} | \nu^n)$ can be bounded from above and below as:

$$\begin{aligned} -\lambda(\mathcal{Q}_n^\pi | \nu^n) + \lambda(\mathcal{P}_n^\pi | \nu^n) - \frac{m}{n} \log(n+1) &\leq \frac{1}{n} \log \pi^{\mathcal{A}}(q^n \in \mathcal{Q} | \nu^n) \leq \\ &\leq -\lambda(\mathcal{Q}_n^\pi | \nu^n) + \lambda(\mathcal{P}_n^\pi | \nu^n) + \frac{m}{n} \log(n+1). \end{aligned}$$

Since for $n \rightarrow \infty$, $\mathcal{S}_n = \mathcal{S}$, and $\nu^n \rightarrow p$, and \mathcal{Q}^π is open, it taken together, implies that $\lambda(\mathcal{Q}_n^\pi | \nu^n)$ converges to $L(\mathcal{Q}^\pi | p)$. Similarly, $\lambda(\mathcal{P}_n^\pi | \nu^n)$ converges to $L(\mathcal{P}^\pi | p)$. \square

Let $\nu^n \rightarrow p$. A set of sources is rare if it does not contain p . Then, from the General prior LST, follows

General prior LCoLT. *Let $\nu^n \rightarrow p$. Let \mathcal{Q}^π be a convex, closed rare set of sources. Let \hat{q}^π be the L-projection of p on \mathcal{Q}^π . Let $B(\hat{q}^\pi, \epsilon)$ be a closed ϵ -ball defined by the total variation metric, centered at \hat{q}^π . Then, for $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} \pi^{\mathcal{A}}(q^n \in B(\hat{q}^\pi, \epsilon) | (q^n \in \mathcal{Q}^\pi) \wedge \nu^n) = 1.$$

4.3 CONDITIONING BY RARE SOURCES: CONTINUOUS ALPHABET

Sanov's Theorem for continuous alphabet can be established either via 'the method of types + discrete approximation' approach (cf. [4]) or by means of the large deviations theory (cf. [7]). The former approach will be used here to formulate continuous alphabet version of LST.

Let $(\mathcal{Y}, \mathcal{F})$ be a measurable space. Let \mathcal{T}^m be a partition of the alphabet \mathcal{Y} into finite number m of sets $\mathcal{T}^m \triangleq (T_1, T_2, \dots, T_m)$; $T_i \in \mathcal{F}$. The \mathcal{T}^m -quantized P , denoted by $P^{\mathcal{T}}$, is defined as the distribution $P(T_1), P(T_2), \dots, P(T_m)$ on the finite set $\mathcal{X} \triangleq \{1, 2, \dots, m\}$.

Let $\mathcal{P}(\mathcal{Y})$ be the set of all probability measures on $(\mathcal{Y}, \mathcal{F})$. Let $\mathcal{Q} \subseteq \mathcal{P}$. For probability measures (pm's) $P, Q \in \mathcal{P}(\mathcal{Y})$, the L^m -divergence $L^m(Q || P)$ of Q with respect to P is defined as

$$L^m(Q || P) \triangleq \sup_{\mathcal{T}^m} L(Q^{\mathcal{T}} || P^{\mathcal{T}}),$$

where the supremum is taken over all m -element partitions. $L^m(\mathcal{Q} || P)$ denotes $\sup_{Q \in \mathcal{Q}} L^m(Q || P)$. Let $\mathcal{Q}^{\mathcal{T}} \triangleq \{Q : Q^{\mathcal{T}} \in \mathcal{Q}\}$, $L^m(\mathcal{Q}^{\mathcal{T}} || P^{\mathcal{T}}) \triangleq \sup_{\mathcal{Q}} L(Q^{\mathcal{T}} || P^{\mathcal{T}})$.

The empirical distribution $\nu^{n,m}$ of an n -sequence of \mathcal{Y} -valued random variables Y with respect to a partition \mathcal{T}^m is defined as

$$\nu_j^{n,m} = \frac{1}{n} \text{Card}\{Y_i : Y_i \in T_j; 1 \leq i \leq n\}, \quad 1 \leq j \leq m.$$

The τ^m -topology of pm's on $(\mathcal{Y}, \mathcal{F})$ is the topology in which a pm belongs to the interior of a set \mathcal{Q} of pm's iff for some partition \mathcal{T}^m and $\epsilon > 0$

$$\{Q' : |Q'(T_j) - Q(T_j)| < \epsilon, j = 1, 2, \dots, m\} \subset \mathcal{Q}.$$

Thus, an n -source $q^n \in \mathcal{R}_n(\mathcal{X})$ belongs to the interior of \mathcal{Q} if there exists \mathcal{T}^m of \mathcal{Y} and $\epsilon > 0$ such that the set $\{Q' : |Q'(T_j) - q_j^n| < \epsilon, j = 1, 2, \dots, m\}$ is a subset of \mathcal{Q} .

Under the assumption of uniform prior distribution of n -sources, a continuous analogue to the Dynamic LST is:

Continuous LST. Let, as $n \rightarrow \infty$, $\nu^{n,m} \rightarrow R$, $R \in \mathcal{R}(\mathcal{X})$. Let \mathcal{Q} be a rare (i.e., $R \notin \mathcal{Q}$) open subset of $\mathcal{P}(\mathcal{Y})$. Then

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{Q} | \nu^{n,m}) = -\{L^m(\mathcal{Q}||R) - L^m(\mathcal{P}||R)\}.$$

Proof. First, an asymptotic lower bound to $\frac{1}{n} \log \pi(q^n \in \mathcal{Q} | \nu^n)$ will be established. Pick up a Q such that for a \mathcal{T}^m , and an $\epsilon > 0$, $Q \in \mathcal{Q}$. Let $\mathcal{M}^{\mathcal{T}}(Q) \triangleq \{q^n : |q_j^n - Q(T_j)| < \epsilon, j = 1, 2, \dots, m\}$. By the Dynamic LST for uniform prior, $\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{M}^{\mathcal{T}}(Q) | \nu^n) = -\{L(\mathcal{M}^{\mathcal{T}}(Q)||R^{\mathcal{T}}) - L(R^{\mathcal{T}}||R^{\mathcal{T}})\}$ which is greater or equal to $-\{L(Q^{\mathcal{T}}||R^{\mathcal{T}}) - L(R^{\mathcal{T}}||R^{\mathcal{T}})\}$, since $Q^{\mathcal{T}} \in \mathcal{M}^{\mathcal{T}}(Q)$. Let $\mathcal{M}(Q) \triangleq \cup_{\mathcal{T}^m} \mathcal{M}^{\mathcal{T}}(Q)$. Then, for $n \rightarrow \infty$, $\frac{1}{n} \log \pi(q^n \in \mathcal{M}(Q) | \nu^n) \geq \sup_{\mathcal{T}^m} -\{L(Q^{\mathcal{T}}||R^{\mathcal{T}}) - L(R^{\mathcal{T}}||R^{\mathcal{T}})\} \equiv -\{L^m(Q||R) - L^m(R||R)\}$. Since $\pi(q^n \in \mathcal{Q} | \nu^n) \geq \sup_{Q \in \mathcal{Q}} \pi(q^n \in \mathcal{M}(Q) | \nu^n)$,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{Q} | \nu^n) &\geq \sup_{Q \in \mathcal{Q}} -\{L^m(Q||R) - L^m(R||R)\} \\ &\equiv -\{L^m(\mathcal{Q}||R) - L^m(\mathcal{P}||R)\}. \end{aligned}$$

Asymptotic upper bound: for \mathcal{T}^m as above, by the Dynamic LST with a uniform prior,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{Q}^{\mathcal{T}} | \nu^n) &= -\{L^m(\mathcal{Q}^{\mathcal{T}}||R^{\mathcal{T}}) - L^m(\mathcal{P}^{\mathcal{T}}||R^{\mathcal{T}})\} \\ &\equiv \sup_{\mathcal{Q}} -\{L(Q^{\mathcal{T}}||\mathcal{P}^{\mathcal{T}}) - L(R^{\mathcal{T}}||R^{\mathcal{T}})\}. \end{aligned}$$

Since $\pi(q^n \in \mathcal{Q} | \nu^n) \leq \sup_{\mathcal{T}^m} \pi(q^n \in \mathcal{Q}^{\mathcal{T}} | \nu^n)$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \pi(q^n \in \mathcal{Q} | \nu^n) \leq -\{L^m(\mathcal{Q}||R) - L^m(\mathcal{P}||R)\}.$$

As the asymptotic lower and upper bounds coincide, the claim follows. \square

5. APPLICATION TO CRITERION CHOICE PROBLEM

1. Let there be an alphabet \mathcal{X} (finite, for simplicity) and prior distribution $\pi(q)$ of sources. From the prior $\pi(q)$ a source is drawn, and the source then generates an n -type ν^n . We are not given the actual source, but rather a set \mathcal{Q} to which the source belongs. Given the alphabet \mathcal{X} , the n -type ν^n , the prior distribution of sources $\pi(\cdot)$ and the set $\mathcal{Q} \subseteq \mathcal{P}(\mathcal{X})$ the objective is to select a source $q \in \mathcal{Q}$. This constitutes the β -problem. Since \mathcal{Q} in general contains more than one source the problem is under-determined and in this sense ill-posed. This paper is concerned with the special case of the β -problem where rational sources (i.e., n -sources) are considered.

If $\mathcal{Q} \equiv \mathcal{P}(\mathcal{X})$, then under the assumption of uniform prior distribution of n -sources, Static LLLN shows that asymptotically (along the types k -equivalent with ν^n) it is just $\hat{q} \equiv \nu^n$ which is the 'only-possible' source of ν^n (i.e., of itself)¹.

¹Note that in the case of unrestricted \mathcal{Q} , ν^n is known to be Non-parametric Maximum Likelihood Estimator of the source. Here, ν^n is the Maximum A-posteriori Probability source.

Dynamic LLLN, assuming that $\nu^n \rightarrow r$, implies that the n -sources concentrate on the true source. However, they do not, if a general prior is assumed, such that it puts zero probability on the true source. In the dynamic case ($\nu^n \rightarrow r$) with general prior, n -sources concentrate on the L -projection of r on \mathcal{P}^π .

What if \mathcal{Q} does not contain ν^n ? How should an n -source be selected in this case of static rare \mathcal{Q} ? One possibility is to select q^n from \mathcal{Q} by minimization of a distance or a convex statistical distance measure [16] between ν^n and \mathcal{Q}_n . In this way, the original β -problem of selecting $q^n \in \mathcal{Q}$ is transformed into an associated Criterion Choice Problem (CCP).

If the rare \mathcal{Q} is convex and closed, Static LCoLT shows that - at least for n sufficiently large - the CCP associated with this instance of the β -problem should be solved by minimization of the L -divergence over \mathcal{Q} . A major qualifier has to be added to this statement: it holds provided that uniform prior distribution of n -sources is assumed. If a general prior, strictly positive on the entire set of rational sources is assumed, then the statement still holds. Prior matters only if it is not strictly positive on the entire \mathcal{R} . Then, it is the L -projection of ν^n on \mathcal{Q}^π that should be selected (recall the General prior LCoLT).

2. Confront the β -problem with the following α -problem (also known as Boltzmann Jaynes Inverse Problem): let there be a source q that emits letters from an alphabet \mathcal{X} . From the source q an n -type was drawn. We are not given the actual n -type, but rather a set Π to which the n -type belongs. Given the alphabet \mathcal{X} , the source q and the set Π the objective is to select an n -type $\nu^n \in \Pi$.

The CCP associated with the α -problem is solved by CoLT and GCP provided that Π is a convex, closed rare set. The Theorems imply that at least for sufficiently large n , the I -projection of q on Π should be selected.

6. ENDNOTES

0) While preparing the final form of the paper, the author learned about the work [10] by Ayalvadi Ganesh and Neil O'Connell, where an inverse of Sanov's Theorem has already been studied and established. The authors also discuss relevance of the Theorem for Bayesian nonparametric consistency. Some differences: the authors work with prior distribution on $\mathcal{P}(\mathcal{X})$ where the alphabet \mathcal{X} is finite. Here the concept of n -source is used and continuous \mathcal{X} is also considered. The rate function of General prior LST of the present work appears to be more general than that of Theorem 1 of [10], as the latter was established for the case of $\pi(p) > 0$. The Conditional Limit Theorem (LCoLT) was not explicitly considered at [10].

1) The terminology and notation of this paper follow more or less closely [2], [4], [6], [7]. The brief survey of Large Deviations Theorems for Empirical Measures (Sect. 3) draws from the same sources. Reader interested in tracking evolution of the Theorems is directed to [1], [3], [7], [8], [12], [16], [18], [19], [22], [23], [24], [25], [27], among others; see [13] for new developments. In relation to the Proposition of Sect. 4.1 see also [9]. The continuous case of conditioning by rare sources (Sect. 5) is built parallel with [12] and [4].

2) This work is motivated by [11], where a problem of selecting between Empirical Likelihood and Maximum Entropy Empirical Likelihood (cf. [20], [21]) has been addressed on probabilistic, rather than statistical, grounds. Further discussion, relevant also to the CCP associated with the α and β -problems, can be found there.

3) Any of the results presented here may be stated in terms of reverse I -projections [5]. For instance the right-hand side of the General prior LST could be equivalently expressed as $-(I(p||\mathcal{Q}^\pi) - I(p||\mathcal{P}^\pi))$, where $I(p||C) \triangleq \inf_{q \in C} I(p||q)$ is the value of the I -divergence at a reverse I -projection of p on C . The above mentioned statistical considerations (and 4) below) served as a motivation for stating the results in terms of the newly introduced L -divergence, though the L -projection is formally identical with the reverse I -projection, which is already in use in a parametric context, cf. [5]. The present work leaves open the issue whether it is more advantageous to state the Theorems of conditioning by rare sources in terms of the L -projection or in terms of the reverse I -projection.

4) If p is an n -type then the L -divergence is known as Kerridge's inaccuracy; cf. [14], [15]. Watanabe in a fundamental work [26] which also addresses questions related to that of the present paper, calls negative of Kerridge's inaccuracy confirmability. A reviewer pointed out that the L -divergence can be identified with mean code length.

5) For any prior $\pi(\cdot)$, the L -projection \hat{q}^π of p on \mathcal{Q}^π is the same as the source which has asymptotically supremal over \mathcal{Q}^π value of the posterior probability $\pi(q^n|\nu^n)$. In the case of uniform prior the correspondence holds for any n .

REFERENCES

1. Bártfai, P., *On a conditional limit theorem*, Progress in Statistics **1** (1974), 85–91.
2. Cover, T. and Thomas, J., *Elements of Information Theory*, John Wiley and Sons, NY, 1991.
3. Csiszár, I., *Sanov Property, Generalized I-projection and a Conditional Limit Theorem*, Ann. Probab. **12** (1984), 768–793.
4. ———, *The Method of Types*, IEEE Trans. IT **44** (1998), 2505–2523.
5. Csiszár, I. and Matúš, F., *Information projections revisited*, IEEE Trans. IT **49** (2004), 1474–1490.
6. Csiszár, I. and Shields, P., *Notes on Information Theory and Statistics: A tutorial*, Foundations and Trends in Communications and Information Theory **1** (2004), 1–111.
7. Dembo, A and Zeitouni, O., *Large Deviations Techniques and Applications*, 2-nd ed., Springer, Application of Mathematics, vol. 38, NY, 1998.
8. Ellis, R. S., *The theory of large deviations: from Boltzmann's 1877 calculation to equilibrium macrostates in 2D turbulence*, Physica D (1999), 106–113.
9. Friedlander M. P. and Gupta M. R., *On minimizing distortion and relative entropy*, IEEE Trans. IT (to appear).
10. Ganesh, A. and O'Connell, N., *An inverse of Sanov's theorem*, Stat. & Prob. Letters **42** (1999), 201–206.
11. Grendár, M. and Judge, G., *Probabilistic approach to Criterion Choice Problem: Estimating Equations case* (2005) (working paper).
12. Groeneboom, P., Oosterhoff, J. and Ruymgaart, F. H., *Large deviation theorems for empirical probability measures*, Ann. Probab. **7** (1979), 553–586.
13. Harremoës, P., *Information topologies with applications* (to appear at *Bolyai Studies*).
14. Kerridge, D. F., *Inaccuracy and inference*, J. Roy. Statist. Soc. Ser. B **23** (1961), 284–294.
15. Kulhavý, R., *Recursive Nonlinear Estimation: A Geometric Approach*, vol. 216, Springer-Verlag, London, 1996 (Lecture Notes in Control and Information Sciences).
16. La Cour, B. R. and Schieve, W. C., *Macroscopic determinism in interacting systems using Large Deviations Theory*, Jour. Stat. Phys. **107 3/4** (2002), 729–755.
17. Liese, F. and Vajda, I., *Convex Statistical Distances*, Teubner, Leipzig, 1987.
18. Leonard, Ch. and Najim, J., *An extension of Sanov's Theorem: Application to the Gibbs Conditioning Principle*, Bernoulli **8 (6)** (2002), 721–743.
19. Lewis, J. T., Pfister, C.-E. and Sullivan, W. G., *Entropy, concentration of probability and conditional theorems*, Markov Proc. Rel. Field. **1** (1995), 319–386.
20. Mittelhammer, R., Judge, G. and Miller, D., *Econometric Foundations*, CUP, Cambridge, 2000.

21. Owen, A., *Empirical Likelihood*, Chapman & Hall/CRC, NY, 2001.
22. Sanov, I. N., *On the probability of large deviations of random variables*, Mat. Sbornik **42** (1957), 11–44 (In Russian).
23. van Campenhout, J. M. and Cover, T. M., *Maximum entropy and conditional probability*, IEEE Trans. IT **27** (1981), 483–489.
24. Vasicek, O. A., *A conditional law of large numbers*, Ann. Probab. **8** (1980), 142–147.
25. Vincze, I., *On the maximum probability principle in statistical physics*, Coll. Math. Soc. J. Bolyai **9** (1972), 869–893.
26. Watanabe, S., *Information-theoretic aspects of inductive and deductive inference*, IBM Journal (1960), 208–231.
27. Zabel, S., *Rates of convergence for conditional expectations*, Ann. Probab. **8** (1980), 928–941.

DEPARTMENT OF MATHEMATICS; FPV UMB; TAJOVSKÉHO 40; SK-974 01 BANSKA BYSTRICA; SLOVAKIA.

INSTITUTE OF MATHEMATICS AND COMPUTER SCIENCE; BANSKA BYSTRICA; SLOVAKIA.

INSTITUTE OF MEASUREMENT SCIENCE; BRATISLAVA; SLOVAKIA

E-mail: marian.grendar@savba.sk