# Large deviations of divergence measures on partitions

Jan Beirlant[a,*], Luc Devroye[b,1], László Györfi[c], Igor Vajda[d,2]

[a] *Department of Mathematics, Katholieke Universiteit Leuven, Celestijnenlaan 200B, B-3000 Leuven,*
*Belgium*
[b] *School of Computer Sciences, McGill University, Montreal, Canada H3A 2K6*
[c] *Department of Computer Science and Information Theory, Technical University of Budapest,*
*1521 Stoczek u. 2, Budapest, Hungary*
[d] *Institute of Information Theory and Automation, Academy of Sciences of the Czech Republic,*
*CZ-182 08 Prague, Czech Republic*

## Abstract

We discuss Chernoff-type large deviation results for the total variation, the I-divergence errors, and the $\chi^2$-divergence errors on partitions. In contrast to the total variation and the I-divergence, the $\chi^2$-divergence has an unconventional large deviation rate. Applications to Bahadur efficiencies of goodness-of-fit tests based on these divergence measures for multivariate observations are given. © 2001 Elsevier Science B.V. All rights reserved.

## 1. Introduction

We consider the problem of testing an unknown probability density function. The test statistics are derived from dissimilarity measures of probability measures, like $\phi$-divergences introduced by Csiszár (1967). The three most important $\phi$-divergences in mathematical statistics and information theory are the total variation distance, the information divergence and the $\chi^2$-divergence. For recent accounts of the theory of $\phi$-divergences, see Liese and Vajda (1987) and Vajda (1989).

* Corresponding author. Tel.: +32-16-322789; fax: +32-16-322831.
*E-mail addresses:* jan.beirlant@wis.kuleuven.ac.be (J. Beirlant), luc@kriek.cs.mcgill.ca (L. Devroye), gyorfi@inf.bme.hu (L. Györfi), vajda@utia.cas.cz (I. Vajda).

If $\mu$ and $v$ are probability measures on $\mathbb{R}^d$ ($d \geqslant 1$), then the *total variation distance* between $\mu$ and $v$ is defined by

$$V(\mu, v) = \sup_A |\mu(A) - v(A)|,$$

where the supremum is taken over all Borel sets $A$.

The *information divergence* (also called I-divergence, Kullback–Leibler number, relative entropy) of $\mu$ and $v$ is defined by

$$I(\mu, v) = \sup_{\{A_j\}} \sum_j \mu(A_j) \log \frac{\mu(A_j)}{v(A_j)},$$

where the supremum is taken over all finite Borel measurable partitions $\{A_j\}$. The term "information divergence" is due to Csiszár (1967), who introduced it to identify the Kullback–Leibler mean information for discrimination between two densities $f$ and $g$ (Kullback and Leibler, 1951) within the more general class of $\phi$-divergences of $f$ and $g$ (Csiszár, 1967, Ali and Silvey, 1966).

The following inequality, termed Pinsker's inequality, gives an upper bound to the total variation in terms of I-divergence (see e.g. Csiszár, 1967; Kullback, 1967; Kemperman, 1969):

$$2\{V(\mu, v)\}^2 \leqslant I(\mu, v). \tag{1}$$

The $\chi^2$-*divergence measure* between $\mu$ and $v$ is defined by

$$\chi^2(\mu, v) = \sup_{\{A_j\}} \sum_j \frac{(\mu(A_j) - v(A_j))^2}{v(A_j)},$$

where again the supremum is taken over all finite Borel measurable partitions $\{A_j\}$.

By using the inequality $\log t \leqslant t - 1$ one easily obtains that

$$\sum_j \mu(A_j) \log \frac{\mu(A_j)}{v(A_j)} \leqslant \sum_j \frac{\mu(A_j)^2}{v(A_j)} - 1 = \sum_j \frac{(\mu(A_j) - v(A_j))^2}{v(A_j)},$$

from which

$$I(\mu, v) \leqslant \chi^2(\mu, v). \tag{2}$$

On the other hand, there are examples where $I(\mu, v_n) \to 0$ and $\chi^2(\mu, v_n) \to \infty$. Therefore the $\chi^2$-divergence is strictly topologically stronger, i.e. the convergence in $\chi^2$-divergence implies convergence in information divergence, but the converse is not true.

Applications of large deviation results in statistical analysis mainly concern the comparison of test procedures using Bahadur efficiencies. We consider the problem of testing hypotheses

   $H_0$: $v = \mu$  versus  $H_1$: $v \neq \mu$

by means of test statistics $T_n = T_n(X_1, \ldots, X_n)$ where $X_1, X_2, \ldots$ are independent and identically distributed along $v$. Considering two tests rejecting $H_0$ for large values of the statistics $T_{n,1}$ and $T_{n,2}$, then (see e.g. Bahadur, 1971; Groeneboom and Shorack,

1981) the efficiency $e_{1,2}$ of $T_{n,1}$ with respect to $T_{n,2}$ is calculated through the Bahadur exact slopes $2b_1(v)$ and $2b_2(v)$ for testing against an alternative $v$: $e_{1,2} = b_1(v)/b_2(v)$. The functions $b_1$ and $b_2$ are then given by

$$b_j(v) = g_j(\psi_j(v)), \quad j = 1, 2 \tag{3}$$

when $v$ is absolutely continuous, provided

$$T_{n,j} \to \psi_j(v) \quad \text{a.s. as } n \to \infty \text{ under } H_1, \tag{4}$$

$$\lim_{n \to \infty} \boldsymbol{E} T_{n,j} = 0 \quad \text{as } n \to \infty \text{ under } H_0 \tag{5}$$

and

$$\lim_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}(T_{n,j} > \varepsilon) = -g_j(\varepsilon) \text{ under } H_0, \ \varepsilon > 0 \tag{6}$$

for $j = 1, 2$.

Such a limit assumption on the tail of the distribution of $T_{n,j}$ means that

$$\boldsymbol{P}(T_{n,j} > \varepsilon) = e^{-n[g_j(\varepsilon) + o(1)]}.$$

Each of the divergence measures defined above, when applied to the distance between the null-hypothesis distribution $\mu$ and the empirical distribution $\mu_n$, both restricted to a partition (see Tusnády, 1977; Barron, 1989), does lead to a test procedure. Hence, the large deviation results given in the subsequent sections offer the possibility to calculate exact Bahadur slopes for these tests. Quine and Robinson (1985) derived large deviation results for test statistics for uniformity based on $I$ and $\chi^2$, i.e. for the classical likelihood ratio and chi-square goodness-of-fit tests. These authors derived these results establishing probability inequalities. Here, we show how such results can be obtained using a classical result of Sanov (1957). Also, we extend these results to a more general setting and derive large deviation results for other test statistics such as the $L_1$-test statistic. We also refer to Nikitin (1995) for a survey on Bahadur efficiencies of different well-known tests.

## 2. The $L_1$ error

We now consider some goodness-of-fit tests for $H_0$ given in the Introduction. Suppose that $\mu$ is nonatomic. Assume a sample of independent random vectors $X_1, \ldots, X_n$, distributed according to a probability measure $\mu$, and let $\mu_n$ denote the empirical measure.

Györfi and van der Meulen (1991) introduced the test statistic

$$J_n = \sum_{j=1}^{m_n} |\mu(A_{n,j}) - \mu_n(A_{n,j})|,$$

based on a finite partition $\mathscr{P}_n = \{A_{n,1}, \ldots, A_{n,m_n}\}$, $(n \geqslant 2)$, of $\mathbb{R}^d$. These authors also showed that under $H_0$

$$\boldsymbol{P}(J_n \geqslant \varepsilon) \leqslant e^{-n(\varepsilon^2/8 + o(1))}.$$

Moreover, the asymptotic normality of this test statistic under the null hypothesis in case $\mu$ has a density was discussed in Beirlant et al. (1995). Let $\bar{\mu}_n$ and $\mu_n^*$ be the restrictions of $\mu$ and $\mu_n$ to the partition $\mathscr{P}_n$, then

$$J_n = 2V(\bar{\mu}_n, \mu_n^*).$$

**Theorem 1.** *Assume that*

$$\lim_{n \to \infty} \max_j \mu(A_{n,j}) = 0 \tag{7}$$

*and*

$$\lim_{n \to \infty} \frac{m_n \log n}{n} = 0. \tag{8}$$

*Then for all $0 < \varepsilon < 2$*

$$\lim_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\{J_n > \varepsilon\} = -g(\varepsilon), \tag{9}$$

*where*

$$g(\varepsilon) = \inf_{0 \,<\, p \,<\, 1-\varepsilon/2} \left( p \log \frac{p}{p + \varepsilon/2} + (1 - p)\log \frac{1 - p}{1 - p - \varepsilon/2} \right). \tag{10}$$

In the proofs of our theorems we shall use the function

$$D(\alpha \| \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha)\log \frac{1 - \alpha}{1 - \beta} \tag{11}$$

and the following lemma.

**Lemma 1** (Sanov, 1957; see p. 16 in Dembo and Zeitouni, 1992; or Problem 1.2.11 in Csiszár and Körner, 1981)**.** *Let $\Sigma$ be a finite set of measurable sets (alphabet), $\mathscr{L}_n$ be a set of types (possible empirical distributions) on $\Sigma$, and let $\Gamma$ be a set of distributions on $\Sigma$. Then*

$$\left| \frac{1}{n} \log \boldsymbol{P}\{\mu_n^* \in \Gamma\} + \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n) \right| \leqslant \frac{|\Sigma| \log(n + 1)}{n}, \tag{12}$$

*where $|\Sigma|$ denotes the cardinality of $\Sigma$.*

**Proof of Theorem 1.** We apply (12) for

$$\Sigma = \{A_{n,1}, \ldots, A_{n,m_n}\},$$

such that

$$\Gamma = \{\tau: 2V(\bar{\mu}_n, \tau) \geqslant \varepsilon\}.$$

Then, according to (12),

$$\left| \frac{1}{n} \log \boldsymbol{P}\{J_n \geqslant \varepsilon\} + \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n) \right| \leqslant \frac{m_n \log(n + 1)}{n}$$

and therefore, under (8),

$$\lim_{n\to\infty} \frac{1}{n} \log \boldsymbol{P}\{J_n > \varepsilon\} = -\lim_{n\to\infty} \inf_{\tau\in\Gamma\cap\mathscr{L}_n} I(\tau, \bar{\mu}_n).$$

It now remains to show that

$$\lim_{n\to\infty} \inf_{\tau\in\Gamma\cap\mathscr{L}_n} I(\tau, \bar{\mu}_n) = g(\varepsilon).$$

The distributions in $\mathscr{L}_n$ are possible empirical distributions, having components of the form $r/n$, where $r$ is integer. Because of (7) we have that

$$m_n \to \infty,$$

therefore because of the continuity of $V(\tau, \bar{\mu}_n)$ and $I(\tau, \bar{\mu}_n)$

$$\lim_{n\to\infty} \inf_{\tau\in\Gamma\cap\mathscr{L}_n} I(\tau, \bar{\mu}_n) = \lim_{n\to\infty} \inf_{2V(\tau,\bar{\mu}_n)\geqslant\varepsilon} I(\tau, \bar{\mu}_n).$$

Here

$$I(\tau, \bar{\mu}_n) = \sum_{j=1}^{m_n} \tau(A_{n,j}) \log \frac{\tau(A_{n,j})}{\mu(A_{n,j})}.$$

Put

$$L = \{j : \mu(A_{n,j}) > \tau(A_{n,j})\}$$

and

$$A_n = \bigcup_{j\in L} A_{n,j}.$$

Then

$$2V(\tau, \bar{\mu}_n) = 2(\mu(A_n) - \tau(A_n))$$

and, by the Information Processing Theorem of Csiszár (1967),

$$I(\tau, \bar{\mu}_n) \geqslant D(\tau(A_n) \| \mu(A_n)),$$

where the equality holds iff $\tau(A_{n,j})/\mu(A_{n,j})$ is constant both on $L$ and $L^c$. Thus,

$$\lim_{n\to\infty} \inf_{2V(\tau,\bar{\mu}_n)\geqslant\varepsilon} I(\tau, \bar{\mu}_n)$$
$$= \inf_{0 < p < 1-\varepsilon/2 : \tau(A_n)=p, \mu(A_n)=p+\varepsilon/2} D(\tau(A_n) \| \mu(A_n)),$$
$$= \inf_{0 < p < 1-\varepsilon/2} \left( p \log \frac{p}{p+\varepsilon/2} + (1-p)\log \frac{1-p}{1-p-\varepsilon/2} \right)$$
$$= g(\varepsilon),$$

and Theorem 1 is proved. $\square$

**Remark 1.** Note that a lower bound for $g$ follows from Pinsker's inequality (1) since

$$\inf_{2V(\tau,\bar{\mu}_n)\geqslant\varepsilon} I(\tau, \bar{\mu}_n) \geqslant \inf_{I(\tau,\bar{\mu}_n)\geqslant\varepsilon^2/2} I(\tau, \bar{\mu}_n),$$

and therefore

$$g(\varepsilon) \geqslant \varepsilon^2/2.$$

The best-known lower bound is due to Toussaint (1975):

$$g(\varepsilon) \geqslant \varepsilon^2/2 + \varepsilon^4/36 + \varepsilon^6/280.$$

An upper bound $\hat{g}(\varepsilon)$ of $g(\varepsilon)$ can be obtained substituting $p$ by $(1 - \varepsilon/2)/2$ in the definition of $g(\varepsilon)$. Then

$$\hat{g}(\varepsilon) = \frac{\varepsilon}{2} \log \frac{2 + \varepsilon}{2 - \varepsilon} \geqslant g(\varepsilon)$$

(Vajda, 1970). Further bounds can be found on pp. 294–295 in Vajda (1989). Remark that also in Lemma 5.1 in Bahadur (1971) it was observed that

$$g(\varepsilon) = \frac{\varepsilon^2}{2}(1 + o(1))$$

as $\varepsilon \to 0$.

Local Bahadur efficiencies of the test $J_n$ with respect to other goodness-of-fit tests can now be computed on the basis of the above theorem.

**Example 1** (*Testing for uniformity*). Let $\mu$ be the restriction of the Lebesgue measure on (0,1). For $\theta \neq 0$ consider alternatives $v$ with density $f_\theta(x) = 1 + \theta h(x)$ ($x \in (0, 1)$) where $\int_0^1 h = 0$. Then, using (3), (9), the asymptotic formula for $g(\varepsilon)$ in Remark 1, and the fact that the $\psi$-function for this $L_1$-test equals $\int_0^1 |f_\theta(x) - 1| \, \mathrm{d}x$, we find that the exact Bahadur slope for the $J_n$-statistic behaves as

$$\theta^2 \left( \int_0^1 |h(x)| \, \mathrm{d}x \right)^2 \quad \text{when } \theta \to 0.$$

Remark that condition (5) can be verified here using Theorem 2.1 in Beirlant and Györfi (1998).

From Groeneboom and Shorack (1981) it follows that the exact Bahadur slope of the Kolmogorov–Smirnov test, respectively the Anderson–Darling test, behaves as $4\theta^2 (\sup_{x \in (0,1)} |\int_0^x h(t) \, \mathrm{d}t|)^2$, respectively $2\theta^2 \int_0^1 (\int_0^x h)^2 [x(1 - x)]^{-1} \, \mathrm{d}x$, when $\theta \to 0$. Remark that

$$\left| \int_0^x h(t) \, \mathrm{d}t \right| = \left| \int_0^1 h(t) \left( 1_{[0,x]}(t) - \frac{1}{2} \right) \mathrm{d}t \right|$$

$$\leqslant \int_0^1 |h(t)| \left| 1_{[0,x]}(t) - \frac{1}{2} \right| \mathrm{d}t = \frac{1}{2} \left| \int_0^1 h(t) \, \mathrm{d}t \right|.$$

Hence the $J_n$-test is more Bahadur efficient than the Kolmogorov–Smirnov test in the given setting. It is also more efficient than the Anderson–Darling test (AD) for several examples of $h$. For example when $h(x) = \operatorname{sgn}(\frac{1}{2} - x)$ ($x \in (0, 1)$), then $e_{J,\mathrm{AD}} = 1/(4 \ln(2) - 2) \sim 1.29$.

**Remark 2.** The above result obviously can be generalized to a large deviation result for the $L_1$-inaccuracy rate of the histogram density estimator $f_n$ based on a sample of independent random vectors $X_1, \ldots, X_n$, distributed according to a probability measure $\mu$ with density $f$ with respect to the Lebesgue measure $\lambda$.

Introducing a partition $\mathscr{P}_n = \{A_{n,j}; \ j \geqslant 1\}$, $(n \geqslant 2)$ of $\mathbb{R}^d$ such that $\sup_{j \geqslant 1} \lambda(A_{n,j}) < \infty$, then the *histogram density estimator* is defined by

$$f_n(x) = \frac{\mu_n(A_n(x))}{\lambda(A_n(x))}. \tag{13}$$

where $A_n(x) = A_{n,i}$ if $x \in A_{n,i}$.

If $\mu$ and $v$ are absolutely continuous with respect to a $\sigma$-finite measure $\lambda$ with densities $f$ and $g$, respectively, then

$$||f - g|| := \int |f(x) - g(x)| \lambda(\mathrm{d}x) = 2 \, V(\mu, v).$$

Hence for any density estimator the $L_1$-consistency implies the consistency in total variation.

Using the assumption that for each sphere $S$ centered at the origin

$$\lim_{n \to \infty} \sup_{A_{n,j} \cap S \neq \emptyset} \operatorname{diam}(A_{n,j}) = 0 \tag{14}$$

and

$$\lim_{n \to \infty} \frac{|\{A_{n,j} \cap S \neq \emptyset\}|}{n} = 0 \tag{15}$$

then

$$\lim_{n \to \infty} ||f_n - f|| = 0 \tag{16}$$

a.s. and

$$\lim_{n \to \infty} \boldsymbol{E}||f_n - f|| = 0 \tag{17}$$

(see Devroye and Györfi, 1985). Combining this result with Theorem 1 then yields the following result.

**Corollary 1.** *Assume* (14). *If there exists a sequence of spheres $S_n$ centered at the origin such that $S_n \uparrow \mathbb{R}^d$ and*

$$\lim_{n \to \infty} \operatorname{card}\{A_{n,j} \cap S_n \neq \emptyset\} \frac{\log n}{n} = 0 \tag{18}$$

*then for all $0 < \varepsilon < 2$*

$$\lim_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\{||f - f_n|| > \varepsilon\} = -g(\varepsilon), \tag{19}$$

*where $g(\varepsilon)$ is defined by* (10).

Note that in this statement there is no condition on $f$, and so (19) holds for all $f$, and the rate function $g(\varepsilon)$ does not depend on $f$.

**Proof of Corollary 1.** Without loss of generality, assume that

$$A_{n,j} \cap S_n \neq \emptyset, \quad \text{for } j = 1, 2, \ldots, m_n - 1$$

and

$$A_{n,j} \cap S_n = \emptyset, \quad \text{for } j = m_n, m_n + 1, \ldots .$$

Put

$$A'_{n,m_n} = \bigcup_{j=m_n}^{\infty} A_{n,j}$$

and

$$\Delta_n = \sum_{j=1}^{m_n-1} |\mu(A_{n,j}) - \mu_n(A_{n,j})| + |\mu(A'_{n,m_n}) - \mu_n(A'_{n,m_n})|.$$

Then

$$\|f_n - f\| = \sum_{j=1}^{\infty} \int_{A_{n,j}} |f_n(x) - f(x)| \lambda(\mathrm{d}x)$$

$$\geqslant \sum_{j=1}^{m_n-1} \int_{A_{n,j}} |f_n(x) - f(x)| \lambda(\mathrm{d}x)$$

$$+ \left| \sum_{j=m_n}^{\infty} \int_{A_{n,j}} f_n(x) \lambda(\mathrm{d}x) - \sum_{j=m_n}^{\infty} \int_{A_{n,j}} f(x) \lambda(\mathrm{d}x) \right|$$

$$\geqslant \sum_{j=1}^{m_n-1} \left| \int_{A_{n,j}} f_n(x) \lambda(\mathrm{d}x) - \int_{A_{n,j}} f(x) \lambda(\mathrm{d}x) \right|$$

$$+ \left| \int_{A'_{n,m_n}} f_n(x) \lambda(\mathrm{d}x) - \int_{A'_{n,m_n}} f(x) \lambda(\mathrm{d}x) \right|$$

$$= \sum_{j=1}^{m_n-1} |\mu(A_{n,j}) - \mu_n(A_{n,j})| + |\mu_n(A'_{n,m_n}) - \mu(A'_{n,m_n})|$$

$$= \Delta_n.$$

On the other hand,

$$\|f_n - f\| \leqslant \|f_n - \boldsymbol{E} f_n\| + \|\boldsymbol{E} f_n - f\|$$

$$= \sum_{j=1}^{\infty} |\mu(A_{n,j}) - \mu_n(A_{n,j})| + \|\boldsymbol{E} f_n - f\|$$

$$\leqslant \sum_{j=1}^{m_n-1} |\mu(A_{n,j}) - \mu_n(A_{n,j})| + |\mu(A'_{n,m_n}) - \mu_n(A'_{n,m_n})|$$

$$+ 2\mu(A'_{n,m_n}) + ||\boldsymbol{E}f_n - f||$$
$$= \Delta_n + 2\mu(A'_{n,m_n}) + ||\boldsymbol{E}f_n - f||.$$

By definition $A'_{n,m_n} \subset S_n^c$, and by assumption $S_n \uparrow \mathbb{R}^d$, and therefore $\mu(A'_{n,m_n}) \to 0$. So together with (17)

$$2\mu(A'_{n,m_n}) + ||\boldsymbol{E}f_n - f|| \to 0$$

and consequently Corollary 1 is proved if

$$\lim_{n\to\infty} \frac{1}{n} \log \boldsymbol{P}\{\Delta_n > \varepsilon\} = -g(\varepsilon), \tag{20}$$

which follows from Theorem 1 with $J_n = \Delta_n$. $\quad\square$

**Remark 3.** Louani (1998) derived large deviation results for a kernel density estimator in terms of the sup-norm distance. Several assumptions are required for this case, such as the boundedness of the density, and the results depend strongly on the bandwidth parameter. Louani (2000) showed a large deviation limit for the $L_1$-error of the kernel estimate, with the same rate function $g$ as given in Theorem 1.

## 3. The information divergence

In the literature on goodness-of-fit testing two statistics are related to the information divergence, namely the *I-divergence statistic*

$$\tilde{I}_n = I(\bar{\mu}_n, \mu_n^*) = \sum_{j=1}^{m_n} \mu(A_{n,j}) \log \frac{\mu(A_{n,j})}{\mu_n(A_{n,j})}$$

and the *reversed I-divergence statistic*

$$I_n = I(\mu_n^*, \bar{\mu}_n) = \sum_{j=1}^{m_n} \mu_n(A_{n,j}) \log \frac{\mu_n(A_{n,j})}{\mu(A_{n,j})}.$$

If there are empty cells ($\mu_n(A_{n,j}) = 0$) then $\tilde{I}_n = \infty$, and hence $I_n$ is more common. We show that

$$\boldsymbol{P}\{I_n > \varepsilon\} = e^{-n(\varepsilon + o(1))}$$

and

$$\boldsymbol{P}\{\tilde{I}_n > \varepsilon\} = e^{-(n/m_n)(1 + o(1))}.$$

Here we can refer to Tusnády (1977) and Barron (1989) who first discussed the exponential character of the tails of $I_n$.

Again a large deviation result concerning $I_n$ can be derived proceeding similarly as in the proof of Theorem 1.

**Theorem 2** (Corollary 2.4 in Kallenberg, 1985; Theorem 2 in Quine and Robinson, 1985)**.** *Under* (7) *and* (8), *for all* $\varepsilon > 0$

$$\lim_{n\to\infty} \frac{1}{n} \log \boldsymbol{P}\{I_n > \varepsilon\} = -\varepsilon.$$

This leads to the following observation concerning the goodness-of-fit tests based on $I_n$.

**Example 2.** Again considering alternatives with density $f_\theta(x) = 1 + \theta h(x)$ $(x \in (0,1))$ where $\int_0^1 h = 0$, when testing for uniformity, the above result implies that the test based on $I_n$ has an exact Bahadur slope $2 \int_0^1 (1 + \theta h(x)) \log(1 + \theta h(x)) \, dx \sim \theta^2 \int_0^1 h^2$ for $\theta \to 0$. Hence the Cauchy–Schwarz inequality implies that this test is more efficient than the test based on $J_n$.

From this point on we suppose that for all $n$

$$\mu(A_{n,j}) = \frac{1}{m_n}, \quad j = 1, \ldots, m_n. \tag{21}$$

This condition can be satisfied for any nonatomic distribution $\mu$ on $\mathbb{R}^d$, e.g. by using tree partitions $\Sigma$ described in Chapter 20 in Devroye et al. (1996).

**Theorem 3.** *Assume that* (21) *holds. If, in addition,*

$$\lim_{n \to \infty} \frac{m_n^2 \log n}{n} = 0, \tag{22}$$

*then we have for all $\varepsilon > 0$*

$$\lim_{n \to \infty} \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{I}_n > \varepsilon\} = -1.$$

**Proof.** Again apply (12) for $\Sigma = \{A_{n,1}, \ldots, A_{n,m_n}\}$ such that

$$\Gamma = \{\tau \colon I(\bar{\mu}_n, \tau) \geqslant \varepsilon\}.$$

Then, according to (12),

$$\left| \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{I}_n \geqslant \varepsilon\} + m_n \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n) \right| \leqslant \frac{m_n^2 \log(n+1)}{n}.$$

Therefore, because of (21) and (22),

$$\lim_{n \to \infty} \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{I}_n > \varepsilon\} = - \lim_{n \to \infty} m_n \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n)$$
$$= - \lim_{n \to \infty} m_n \inf_{I(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n).$$

Using the notation

$$\tau = (\tau_1, \ldots, \tau_{m_n}),$$

we obtain

$$I(\tau, \bar{\mu}_n) = \sum_{j=1}^{m_n} \tau_j \log(m_n \tau_j) \tag{23}$$

and

$$I(\bar{\mu}_n, \tau) = -\frac{1}{m_n} \sum_{j=1}^{m_n} \log(m_n \tau_j).$$

Considering the distribution

$$\tau^* = \left( 0, \frac{1}{m_n - 1}, \dots, \frac{1}{m_n - 1} \right) \tag{24}$$

we find that

$$m_n \inf_{I(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n) \leqslant m_n I(\tau^*, \bar{\mu}_n) = m_n \log(m_n/(m_n - 1)) \to 1.$$

Concerning the lower bound we show that for all but finitely many $n$, the minimizing distribution is of the form

$$\tilde{\tau} = \left( \tau_1, \frac{1 - \tau_1}{m_n - 1}, \dots, \frac{1 - \tau_1}{m_n - 1} \right), \tag{25}$$

where $\tau_1$ (possibly depending on $n$) is bounded by

$$0 < \tau_1 \leqslant \frac{e^{-\varepsilon m_n}}{m_n}.$$

To this end fix $n$ and suppose that $\tau$ is the minimizing distribution, with coordinates ordered in the sense $\tau_1 \leqslant \tau_2 \leqslant \cdots \leqslant \tau_{m_n}$.

  (a) It holds $\tau_1 > 0$. Indeed, if $\tau_1 = 0$ then, since the Shannon entropy is maximum for the uniform distribution, it follows that $I(\tau, \bar{\mu}_n) \geqslant I(\tau^*, \bar{\mu}_n) = \log(1/(1 - 1/m_n))$. This however contradicts the fact that the functions

$$I(\tilde{\tau}, \bar{\mu}_n) = D(\tau_1 \| 1/m_n), \quad I(\bar{\mu}_n, \tilde{\tau}) = D(1/m_n \| \tau_1), \tag{26}$$

(cf. (11) and (25)) are continuous, strictly decreasing in the variable $\tau_1 \in (0, 1/m_n]$ with

$$D(0 \| 1/m_n) = \log \frac{1}{1 - 1/m_n}, \quad D(1/m_n \| 0) = \infty \tag{27}$$

and

$$D(1/m_n \| 1/m_n) = 0. \tag{28}$$

  (b) There exists $1 \leqslant r_n < m_n$ such that

$$\tau_1 = \cdots = \tau_{r_n} < \tau_{r_n+1} = \cdots = \tau_{m_n}. \tag{29}$$

If (29) does not hold then either $\tau_1 = \tau_2 = \cdots = \tau_{m_n} = 1/m_n$ or there exist $1 < r_n < s_n \leqslant m_n$ such that $\tau_1 < \tau_{r_n} < \tau_{s_n}$. The first possibility contradicts the assumption $I(\bar{\mu}_n, \tau) \geqslant \varepsilon > 0$. To see that the second possibility contradicts the definition of $\tau$, suppose for simplicity that $r_n = 2$ and $s_n = 3$, i.e. let (cf. (a))

$$0 < \tau_1 < \tau_2 < \tau_3. \tag{30}$$

Consider a new distribution $\bar{\tau}$ with all coordinates coinciding with those of $\tau$ except $\bar{\tau}_1 = \tau_1 - \delta$, $\bar{\tau}_2 = \tau_2 + \delta t$ and $\bar{\tau}_3 = \tau_3 - \delta(t - 1)$, where $0 < \delta < \tau_1$ and $0 < t < 1 + \tau_3/\tau_1$ are chosen such that $I(\bar{\mu}_n, \bar{\tau}) = I(\bar{\mu}_n, \tau)$, i.e. such that

$$\sum_{j=1}^{3} \log \frac{\bar{\tau}_j}{\tau_j} = 0.$$

This means that

$$\left(-\frac{1}{\tau_1} + \frac{t}{\tau_2} - \frac{t-1}{\tau_3}\right)\delta + o(\delta) = 0 \quad \text{as } \delta \to 0,$$

i.e.

$$t = \frac{x-1}{y-1} + o(1) \quad \text{for } x = \frac{\tau_3}{\tau_1}, \quad y = \frac{\tau_3}{\tau_2}, \quad \frac{x-1}{y-1} > 1.$$

It is easy to verify that

$$I(\bar{\tau}, \bar{\mu}_n) = I(\tau, \bar{\mu}_n) + \delta(\ln x - t \ln y) + o(\delta) \quad \text{as } \delta \to 0,$$

where $x > y > 1$ and $(\ln x)/(x-1) < (\ln y)/(y-1)$. Thus $I(\bar{\mu}_n, \bar{\tau}) < I(\bar{\mu}_n, \tau)$ for sufficiently small $\delta > 0$ while $I(\bar{\mu}_n, \bar{\tau}) \geqslant \varepsilon$ for all $\delta > 0$.

(c) In (29), $r_n = 1$ for all but finitely many $n$. To see this, we define

$$\alpha_n = \tau_1 r_n$$

so that $\tau_1 = \alpha_n/r_n$, $\tau_{m_n} = (1 - \alpha_n)/(m_n - r_n)$, and

$$\beta_n = \frac{r_n}{m_n} > \alpha_n.$$

Then, using (11),

$$I(\tau, \bar{\mu}_n) = D(\alpha_n \| \beta_n) \quad \text{and} \quad I(\bar{\mu}_m, \tau) = D(\beta_n \| \alpha_n).$$

Since $\beta_n > \alpha_n$, the assumption $I(\bar{\mu}_n, \tau) \geqslant \varepsilon$ implies

$$\beta_n \log \frac{\beta_n}{\alpha_n} \geqslant D(\beta_n \| \alpha_n) \geqslant \varepsilon,$$

so that $\alpha_n < \beta_n \, e^{-\varepsilon/\beta_n}$. Combining this with the monotonicity of $D(\alpha \| \beta_n)$ in the domain $\alpha \in (0, \beta_n]$, we obtain

$$I(\tau, \bar{\mu}_n) \geqslant D(\beta_n \, e^{-\varepsilon/\beta_n} \| \beta_n)$$
$$= -\varepsilon \, e^{-\varepsilon/\beta_n} + (1 - \beta_n \, e^{-\varepsilon/\beta_n}) \log \frac{1 - \beta_n \, e^{-\varepsilon/\beta_n}}{1 - \beta_n}. \tag{31}$$

If $r_n = \beta_n m_n \geqslant 2$ for infinitely many $n$, then the last inequality contradicts the upper bound

$$I(\tau, \bar{\mu}_n) \leqslant D(0 \| 1/m_n) = \log \frac{1}{1 - 1/m_n}$$

proved above. This completes the proof of (c).

If (c) holds, then $\tau$ under consideration coincides with $\tilde{\tau}$ given by (25) and (3). The desired result $I(\tau, \bar{\mu}_n) \geqslant 1/m_n + o(1/m_n)$ thus follows from (31) applied to $\alpha_n = \tau_1$ and $\beta_n = 1/m_n$. □

**Remark 4.** Remark that formula (3) cannot be applied directly for the test based on $\tilde{I}_n$ because (5) does not hold in this case, in contrast to the test based on the $I$-divergence statistic (see Remark 5 below). The same holds for the tests based on the reversed Pearson statistic discussed in the next section.

## 4. The $\chi^2$-divergence

Of course, the best-known goodness-of-fit test is based on the $\chi^2$ statistic. We refer to Neyman (1949), Watson (1958) for some important historic references. See also Kallenberg et al. (1985) for a more recent discussion. The *Pearson statistic* is given by

$$\chi_n^2 = \chi^2(\mu_n^*, \bar{\mu}_n) = \sum_{j=1}^{m_n} \frac{(\mu(A_{n,j}) - \mu_n(A_{n,j}))^2}{\mu(A_{n,j})},$$

while the reversed Pearson statistic, also known as *Neyman* or *Neyman-modified Pearson statistic* is defined as

$$\tilde{\chi}_n^2 = \chi^2(\bar{\mu}_n, \mu_n^*) = \sum_{j=1}^{m_n} \frac{(\mu(A_{n,j}) - \mu_n(A_{n,j}))^2}{\mu_n(A_{n,j})}.$$

We show that

$$P\{\chi_n^2 > \varepsilon\} = e^{-(n \log m_n / \sqrt{m_n})(\sqrt{\varepsilon}/2 + o(1))}$$

and

$$P\{\tilde{\chi}_n^2 > \varepsilon\} = e^{-(n/m_n)(1 + o(1))}.$$

**Theorem 4** (Theorem 1 in Quine and Robinson, 1985). *Suppose that* (21) *holds. If, in addition,*

$$\frac{m_n^{3/2} \log n}{n \log m_n} \to 0, \tag{32}$$

*then for all $\varepsilon > 0$*

$$\lim_{n \to \infty} \frac{\sqrt{m_n}}{n \log m_n} \log P\{\chi_n^2 > \varepsilon\} = -\sqrt{\varepsilon}/2.$$

**Proof.** Again apply (12) for $\Sigma = \{A_{n,1}, \ldots, A_{n,m_n}\}$ such that

$$\Gamma = \{\tau: \chi^2(\tau, \bar{\mu}_n) \geqslant \varepsilon\}.$$

Then, according to (12),

$$\left| \frac{\sqrt{m_n}}{n \log m_n} \log P\{\chi_n^2 \geqslant \varepsilon\} + \frac{\sqrt{m_n}}{\log m_n} \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n) \right| \leqslant \frac{m_n^{3/2} \log(n+1)}{n \log m_n}.$$

Therefore, because of (32) and (21),

$$\lim_{n \to \infty} \frac{\sqrt{m_n}}{n \log m_n} \log P\{\chi_n^2 > \varepsilon\} = - \lim_{n \to \infty} \frac{\sqrt{m_n}}{\log m_n} \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n)$$

$$= - \lim_{n \to \infty} \frac{\sqrt{m_n}}{\log m_n} \inf_{\chi^2(\tau, \bar{\mu}_n) \geqslant \varepsilon} I(\tau, \bar{\mu}_n).$$

Using the notation

$$\tau = (\tau_1, \ldots, \tau_{m_n})$$

we have (23) and

$$\chi^2(\tau, \bar{\mu}_n) = m_n \sum_{j=1}^{m_n} \tau_j^2 - 1.$$

Concerning the upper bound consider the distribution $\tilde{\tau}$ defined by (25) where $\tau_1$ solves the equation

$$\chi^2(\tilde{\tau}, \bar{\mu}_n) = \varepsilon,$$

or,

$$\tau_1 = \frac{1 + \sqrt{\varepsilon(m_n - 1)}}{m_n}.$$

Then

$$\frac{\sqrt{m_n}}{\log m_n} \inf_{\chi^2(\tau, \bar{\mu}_n) \geqslant \varepsilon} I(\tau, \bar{\mu}_n) \leqslant \frac{\sqrt{m_n}}{\log m_n} I(\tilde{\tau}, \bar{\mu}_n) \to \sqrt{\varepsilon}/2. \tag{33}$$

Concerning the lower bound one shows along similar lines as in the proof of Theorem 3 that for all but finitely many $n$, $\tilde{\tau}$ is the minimizing distribution.

**Remark 5.** Since $\boldsymbol{E}(I_n) \leqslant \boldsymbol{E}(\chi_n^2) = (m_n - 1)/n$, (5) holds for the statistics $I_n$ and $\chi_n^2$ when $m_n/n \to 0$.

**Remark 6.** Theorem 4 means that when $m_n \to \infty$, the Bahadur exact slope of the $\chi^2$-test is identically zero. Another interpretation is that the tail of the $\chi^2$-test statistic is of sub-exponential nature, that is, is heavier than an exponential tail. This is due to cells $A_{n,j}$ with small probabilities, which put too much weight on the squared difference $(\mu(A_{n,j}) - \mu_n(A_{n,j}))^2$. Although there is a widespread believe in literature that the I-divergence test $I_n$ and the $\chi^2$-test have a similar behaviour, their Bahadur slopes are quite different.

**Theorem 5.** *Under the conditions of Theorem* 3

$$\lim_{n \to \infty} \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{\chi}_n^2 > \varepsilon\} = -1.$$

**Proof.** Again apply (12) for $\Sigma = \{A_{n,1}, \ldots, A_{n,m_n}\}$ such that

$$\Gamma = \{\tau \colon \chi^2(\bar{\mu}_n, \tau) \geqslant \varepsilon\}.$$

Then, according to (12),

$$\left| \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{\chi}_n^2 \geqslant \varepsilon\} + m_n \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n) \right| \leqslant \frac{m_n^2 \log(n+1)}{n}.$$

Therefore, because of (21) and (22),

$$\lim_{n \to \infty} \frac{m_n}{n} \log \boldsymbol{P}\{\tilde{\chi}_n^2 > \varepsilon\} = - \lim_{n \to \infty} m_n \inf_{\tau \in \Gamma \cap \mathscr{L}_n} I(\tau, \bar{\mu}_n)$$

$$= - \lim_{n \to \infty} m_n \inf_{\chi^2(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n).$$

Concerning the upper bound consider distribution (24). Then

$$m_n \inf_{\chi^2(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n) \leqslant m_n I(\tau_0, \bar{\mu}_n) \to 1.$$

Concerning the lower bound apply Theorem 3 and (2) to obtain

$$m_n \inf_{\chi^2(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n) \geqslant m_n \inf_{I(\bar{\mu}_n, \tau) \geqslant \varepsilon} I(\tau, \bar{\mu}_n) \to 1.$$

## References

Ali, M.S., Silvey, S.D., 1966. A general class of coefficients of divergence of one distribution from another. J. Roy. Statist. Soc. Ser B 28, 131–140.

Bahadur, R.R., 1971. Some Limit Theorems in Statistics. SIAM, Philadelphia.

Barron, A.R., 1989. Uniformly powerful goodness of fit tests. Ann. Statist. 17, 107–124.

Beirlant, J., Györfi, L., 1998. On the $L_1$-error in histogram density estimation: the multidimensional case. Nonparametric Statist. 9, 197–216.

Beirlant, J., Györfi, L., Lugosi, G., 1995. On the asymptotic normality of the $L_1$- and the $L_2$-errors in histogram density estimation. Canadian J. Statist. 22, 309–318.

Csiszár, I., 1967. Information-type measures of divergence of probability distributions and indirect observations. Studia Sci. Math. Hungar. 2, 299–318.

Csiszár, I., Körner, J., 1981. Information Theory: Coding Theorems for Memoryless Systems. Academic Press, New York.

Dembo, A., Zeitouni, O., 1992. Large Deviations Techniques and Applications. Jones and Bartlett Publishers, Boston.

Devroye, L., Györfi, L., 1985. Nonparametric Density Estimation: the $L_1$ View. Wiley, New York.

Devroye, L., Györfi, L., Lugosi, G., 1996. A Probabilistic Theory of Pattern Recognition. Springer, New York.

Groeneboom, P., Shorack, G.R., 1981. Large deviations of goodness of fit statistics and linear combinations of order statistics. Ann. Probab. 9, 971–987.

Györfi, L., van der Meulen, E.C., 1991. A consistent goodness-of-fit test based on the total variation distance. In: Roussas, G. (Ed.), Nonparametric Functional Estimation and Related Topics. Kluwer, Boston, pp. 631–646.

Kallenberg, W.C.M., 1985. On moderate and large deviations in multinomial distributions. Ann. Statist. 13, 1554–1580.

Kallenberg, W.C.M., Oosterhof, J., Shriever, B.F., 1985. The number of classes in chi-squared goodness-of-fit tests. J. Amer. Statist. Assoc. 80, 959–968.

Kemperman, J.H.B., 1969. An optimum rate of transmitting information. Ann. Math. Statist. 40, 2156–2177.

Kullback, S., 1967. A lower bound for discrimination in terms of variation. IEEE Trans. Inform. Theory 13, 126–127.

Kullback, S., Leibler, R.A., 1951. On information and sufficiency. Ann. Math. Statist. 22, 79–86.

Liese, F., Vajda, I., 1987. Convex Statistical Distances. Teubner, Leipzig.

Louani, D., 1998. Large deviation limit theorems for the kernel density estimator. Scand. J. Statist. 25, 243–253.

Louani, D., 2000. Large deviations for the $L_1$-distance in kernel density estimation, J. Statist. Plann. Inf. 90, 177–182.

Neyman, J., 1949. Contribution to the theory of the $\chi^2$ test. Proceedings of the First Berkeley Symposium on Mathematical Statistics and Probability, Berkeley University Press, Berkeley, CA, pp. 239–273.

Nikitin, Ya., 1995. Asymptotic Efficiency of Nonparametric Tests. Cambridge University Press, Cambridge.

Quine, M.P., Robinson, J., 1985. Efficiencies of chi-square and likelihood ratio goodness-of-fit tests. Ann. Statist. 13, 727–742.

Sanov, I.N., 1957. On the probability of large deviations of random variables. Mat. Sb. 42, 11–44 (English translation in Sel. Transl. Math. Statist. Probab. 1 (1961) 213–244).

Toussaint, G.T., 1975. Sharper lower bounds for information in term of variation. IEEE Trans. Inform. Theory IT-21, 99–103.

Tusnády, G., 1977. On asymptotically optimal tests. Ann. Statist. 5, 385–393.

Vajda, I., 1970. Note on discrimination information and variation. IEEE Trans. Inform. Theory IT-16, 771–773.

Vajda, I., 1989. Theory of Statistical Inference and Information. Kluwer Academic Publishers, Dordrecht.

Watson, G.S., 1958. On chi-squared goodness-of-fit tests for continuous distributions. J. Roy. Statist. Soc. Ser. B 20, 44–61.