

# A REMARK ON THE CONVERGENCE OF KULLBACK-LEIBLER'S MEAN INFORMATION

By SADA0 IKEDA

(Received June 10, 1960)

## Summary

Under fairly general assumptions two sufficient conditions for the convergence of the Kullback-Leibler mean information are obtained, which are generalizations of the conditions given in Lemma 2.1 and Theorem 2.2 of Chapter 4 in S. Kullback [1].

## 1. A convergence theorem

Let  $(R, S, m)$  be a  $\sigma$ -finite measure space, and  $\{\mu, \nu\}$  a set of two measures on the measurable space  $(R, S)$  dominated by the measure  $m$ . Let  $f(x)$  and  $g(x)$  be the Radon-Nikodym derivatives, i.e., the density functions (with respect to measure  $m$ ) of measures  $\mu$  and  $\nu$ , respectively.  $D(f)$  and  $D(g)$  denote the carriers of  $f(x)$  and  $g(x)$ .

Assume measures  $\mu$  and  $\nu$  are absolutely continuous with respect to each other. Then, it follows that  $D(f)=D(g)(m)$ , and vice versa. Under this assumption, the "mean information for discrimination in favor of  $\mu$  against  $\nu$ " is defined by

$$(1.1) \quad I(f: g) = \int_R f(x) \log \frac{f(x)}{g(x)} dm(x).$$

For brevity, we call the above expression the "Kullback-Leibler mean information".

If both measures  $\mu$  and  $\nu$  are finite, then the inequality

$$(1.2) \quad \mu(R) \log \frac{\mu(R)}{\nu(R)} \leq I(f: g)$$

will be obtained, where the equality holds when and only when the ratio of  $f(x)$  to  $g(x)$  is a constant for almost all ( $m$ )  $x$  in  $R$ . In particular, if both  $\mu$  and  $\nu$  are the probability measures, then (1.2) becomes

$$(1.3) \quad 0 \leq I(f: g),$$

where the equality holds if and only if  $f(x)=g(x)(m)$  on  $R$ .

The following theorem gives a sufficient condition for the convergence

of the mean information (1.1) for finite measures.

**THEOREM 1.1.** *Assume that*

(i)  $f(x)$  and  $g(x)$  are the density functions with respect to  $m$ , of two finite measures  $\mu$  and  $\nu$  dominated by each other, i.e.,  $f(x)$  and  $g(x)$  are nonnegative almost everywhere ( $m$ ) on  $R$ ,  $\mu(R) = \int_R f dm < \infty$ , and  $\nu(R) = \int_R g dm < \infty$ , and  $D(f) = D(g) = E(m)$ ,

(ii)  $\{f_i, g_i\}$  ( $i=1, 2, \dots$ ) is a sequence of pairs of the density functions of finite measures  $\mu_i$  and  $\nu_i$  dominated by each other, or more precisely,  $f_i(x)$  and  $g_i(x)$  are nonnegative almost everywhere ( $m$ ) on  $R$ ,  $\mu_i(R) = \int_R f_i dm < K$ , and  $\nu_i(R) = \int_R g_i dm < K$ , and  $D(f_i) = D(g_i) = E_i(m)$  for  $i=1, 2, \dots$ , where  $K$  is a positive constant independent of  $i$ , and (iii)  $E_i \subset E(m)$  ( $i=1, 2, \dots$ ) and  $\nu(E - E_i) \rightarrow 0$  as  $i \rightarrow \infty$ .

Under these assumptions, it holds that

$$(1.4) \quad I(f_i : g_i) \rightarrow I(f : g) \quad (i \rightarrow \infty),$$

if either of the following conditions is satisfied:

$$(A) \quad \begin{cases} (a) & h_i(g, g_i) = \text{ess sup}_{E_i} \left| 1 - \frac{g_i}{g} \right| \rightarrow 0 \quad (i \rightarrow \infty), \text{ and} \\ (b) & d_i(p, p_i) = \text{ess sup}_{E_i} |p - p_i| \rightarrow 0 \quad (i \rightarrow \infty) \end{cases}$$

and

$$(B) \quad \begin{cases} (a)' & h_i(g, g_i) = \text{ess sup}_{E_i} \left| 1 - \frac{g_i}{g} \right| \rightarrow 0 \quad (i \rightarrow \infty), \text{ and} \\ (b)' & h_i(p, p_i) = \text{ess sup}_{E_i} \left| 1 - \frac{p_i}{p} \right| \rightarrow 0 \quad (i \rightarrow \infty) \end{cases}$$

where

$$p(x) = \frac{f(x)}{g(x)}(m), \quad \text{and} \quad p_i(x) = \frac{f_i(x)}{g_i(x)}(m) \quad (i=1, 2, \dots),$$

respectively on  $E$  and  $E_i$ , and  $\text{ess sup}$  is taken with respect to measure  $m$ .

**PROOF.** 1° In the first place, we consider the case when

$$(1.5) \quad |I(f : g)| < \infty.$$

From definition (1.1) we can write as

$$(1.6) \quad \text{and} \quad \begin{aligned} I(f : g) &= \int_E p \log p d\nu, \\ I(f_i : g_i) &= \int_{E_i} p_i \log p_i d\nu_i, \end{aligned} \quad (i=1, 2, \dots),$$

therefore, we have

$$(1.7) \quad |I(f: g) - I(f_i: g_i)| \leq \left| \int_{E_i} p \log p d\nu - \int_{E_i} p_i \log p_i d\nu_i \right| + \left| \int_{E-E_i} p \log p d\nu \right| \quad (i=1, 2, \dots).$$

The second term of the right-hand side of (1.7) converges to zero as  $i \rightarrow \infty$ , by virtue of assumption (iii) and the finiteness of  $I(f: g)$ . The first term of the right-hand side of (1.7) will be evaluated by the following inequality:

$$(1.8) \quad \left| \int_{E_i} p \log p d\nu - \int_{E_i} p_i \log p_i d\nu_i \right| \leq \left| \int_{E_i} p \log p d\nu - \int_{E_i} p \log p d\nu_i \right| + \left| \int_{E_i} p \log p d\nu_i - \int_{E_i} p_i \log p_i d\nu_i \right|.$$

First, we consider the first part of the right-hand side of (1.8). Since

$$(1.9) \quad \left| \int_{E_i} p \log p d\nu - \int_{E_i} p \log p d\nu_i \right| \leq \int_{E_i} |p \log p| |g - g_i| dm \leq h_i(g, g_i) \int_{E_i} |p \log p| d\nu,$$

it follows from assumption (iii), (1.5) and condition (a) of (A) or (B) that

$$(1.10) \quad \left| \int_{E_i} p \log p d\nu - \int_{E_i} p \log p d\nu_i \right| \rightarrow 0 \quad (i \rightarrow \infty).$$

Next, we consider the second part of the right-hand side of (1.8). From Lemma 2.1 (ii) of the author's paper [2], and condition (A), it will easily be seen that, for any  $\varepsilon > 0$ , there exists a positive integer  $N$  such that  $N \leq i$  implies  $h_i(g, g_i) < \varepsilon$  and  $d_i(p, p_i) < \varepsilon$ , and

$$(1.11) \quad \left| \int_{E_i} p \log p d\nu_i - \int_{E_i} p_i \log p_i d\nu_i \right| \leq \int_{E_i} |p \log p - p_i \log p_i| d\nu_i \leq \int_{E_i} \varepsilon(p+1+\varepsilon) d\nu_i = \varepsilon \left\{ \int_{E_i} p g_i dm + (1+\varepsilon) \int_{E_i} g_i dm \right\}.$$

Since  $g_i(x) < (1+\varepsilon)g(x)(m)$  on  $E_i$ , it holds that

$$\int_{E_i} p g_i dm \leq (1+\varepsilon) \int_{E_i} f dm.$$

Therefore, from assumptions (i) and (ii) it will easily be seen that the values of the members within the bracket in the last expression

of (1.11) are bounded. Hence, we obtain

$$(1.12) \quad \left| \int_{E_i} p \log p d\nu_i - \int_{E_i} p_i \log p_i d\nu_i \right| \rightarrow 0 \quad (i \rightarrow \infty),$$

under condition (A). When condition (B) is satisfied instead of (A), the above convergence (1.12) will be shown as follows: from Lemma 2.1 (i) of the author's paper [2], we have

$$(1.13) \quad \begin{aligned} \left| \int_{E_i} p \log p d\nu_i - \int_{E_i} p_i \log p_i d\nu_i \right| &\leq \int_{E_i} |p \log p - p_i \log p_i| d\nu_i \\ &\leq \int_{E_i} \left| 1 - \frac{p_i}{p} \right| (|p \log p| + p + p_i) d\nu_i \\ &\leq h_i(p, p_i) \left\{ \int_{E_i} |p \log p| g_i dm + \int_{E_i} p g_i dm + \int_{E_i} f_i dm \right\}. \end{aligned}$$

By the investigation analogous to that of the case of condition (A) above, the values of the members in the bracket of the last expression (1.13) are bounded for sufficiently large  $i$ . Hence, (1.12) follows from (1.13) and condition (B).

Thus our theorem is proved in the case when  $|I(f: g)| < \infty$ .

2°) Secondly, we shall prove the theorem in the case when

$$(1.14) \quad I(f: g) = \infty.$$

For this case it will be shown that

$$(1.15) \quad I(f_i: g_i) \rightarrow \infty \quad (i \rightarrow \infty).$$

For each positive integer  $N$ , we define the function such as

$$(1.16) \quad \text{and} \quad \begin{aligned} f^N(x) &= \begin{cases} f(x), & \text{on } E \cap \{x: p(x) \leq N\}, \\ Ng(x), & \text{on } E \cap \{x: p(x) > N\}, \\ 0, & \text{otherwise,} \end{cases} \\ f_i^N(x) &= \begin{cases} f_i(x), & \text{on } E_i \cap \{x: p_i(x) \leq N\}, \\ Ng_i(x), & \text{on } E_i \cap \{x: p_i(x) > N\}, \\ 0, & \text{otherwise} \end{cases} \quad (i=1, 2, \dots). \end{aligned}$$

Put  $p^N(x) = f^N(x)/g(x)$  on  $E$  and  $p_i^N(x) = f_i^N(x)/g_i(x)$  on  $E_i$  ( $i=1, 2, \dots$ ). Then, these functions will be definite except for the set of  $m$ -measure zero on each of the sets  $E$  and  $E_i$ 's, respectively.

For these functions we consider the mean information such as

$$(1.17) \quad \text{and} \quad \begin{aligned} I(f^N: g) &= \int_E p^N \log p^N d\nu, \\ I(f_i^N: g_i) &= \int_{E_i} p_i^N \log p_i^N d\nu_i \quad (i=1, 2, \dots). \end{aligned}$$

Since the integrands of the expressions of the right-hand sides of (1.17) are all bounded from above and measures  $\nu$  and  $\nu_i$ 's are all finite measures, it follows from (1.2) that

$$(1.18) \quad |I(f^N : g)| < \infty, \quad \text{and} \quad |I(f_i^N : g_i)| < \infty \quad (i=1, 2, \dots),$$

for any fixed  $N$ . Moreover, since  $p^N(x)$  and  $p_i^N(x)$ 's coincide with  $p(x)$  and  $p_i(x)$ 's respectively on the domains where they are less than or equal to  $N$  ( $\geq 1$ ), and  $p^N$  and  $p_i^N$ 's are all monotone nondecreasing functions of  $N$ , it holds that

$$(1.19) \quad \begin{aligned} &I(f^N : g) \uparrow I(f : g) \quad (N \rightarrow \infty), \\ &I(f_i^N : g_i) \uparrow I(f_i : g_i) \quad (N \rightarrow \infty) \quad \text{for any fixed } i \quad (i=1, 2, \dots). \end{aligned}$$

First, we shall show that, for any fixed  $N$ ,

$$(1.20) \quad I(f_i^N : g_i) \rightarrow I(f^N : g) \quad (i \rightarrow \infty).$$

Since condition (1.5) in the proof of case 1° is fulfilled for  $I(f^N : g)$  by virtue of (1.18), in order to show the convergence (1.20) it will be sufficient to confirm that functions  $f^N$ ,  $g$  and  $\{f_i^N, g_i^N\}$  ( $i=1, 2, \dots$ ) satisfy all assumptions (i)-(iii) and condition (A) or (B) of the present theorem. It will be evident that they satisfy all assumptions (i)-(iii). Since the definitions of functions  $p^N$  and  $p_i^N$ 's in (1.16) do not change the values of functions  $g$  and  $g_i$ 's, and it holds that

$$(1.21) \quad \begin{aligned} &d_i(p^N, p_i^N) \leq d_i(p, p_i), \quad \text{and} \\ &h_i(p^N, p_i^N) \leq h_i(p, p_i) \quad (i=1, 2, \dots; N=1, 2, \dots), \end{aligned}$$

it is seen that condition (A) or (B) is fulfilled for our present case if (A) or (B) is satisfied for  $f$ ,  $g$  and  $\{f_i, g_i\}$  ( $i=1, 2, \dots$ ), respectively. Hence, (1.20) holds true for any fixed  $N$ .

From (1.14) and (1.19) it follows that

$$(1.22) \quad I(f^N : g) \rightarrow \infty \quad (N \rightarrow \infty).$$

Therefore, for any  $M$  ( $>0$ ), there exists a positive integer  $N'$  such that

$$(1.23) \quad M+1 < I(f^{N'} : g).$$

It will be seen from (1.20) that, for this  $N'$  there exists a positive integer  $N''$  such that  $N'' \leq i$  implies

$$(1.24) \quad |I(f_i^{N'} : g_i) - I(f^{N'} : g)| < 1.$$

Hence, it follows from (1.19), (1.23) and (1.24) that

$$(1.25) \quad M < I(f_i : g_i) \quad \text{for } i \geq N'',$$

which implies (1.15).

This completes the proof of the theorem.

## 2. Corollaries

The result of Lemma 2.1 in Chapter 4 of S. Kullback [1] states that a necessary and sufficient condition for the convergence of the mean information  $I(f_i : f)$  to zero for  $i \rightarrow \infty$ , where  $f$  and  $f_i$ 's are generalized probability density functions, is given by

$$(2.1) \quad h(f, f_i) = \text{ess sup}_E \left| 1 - \frac{f_i}{f} \right| \rightarrow 0 \quad (i \rightarrow \infty),$$

where the set  $E$  is the carrier of  $f$  and  $f_i$ 's. In general, however, this condition is not a necessary condition, as will be seen by some simple examples. Corollary 2.1 below shows that our theorem 1.1 gives the same condition as (2.1) which is sufficient for the convergence  $I(f_i : f) \rightarrow 0$  for  $i \rightarrow \infty$ .

Theorem 2.2 in Chapter 4 of S. Kullback [1] is concerned with the convergence of  $I(f_i : g)$  to  $I(f : g)$  for  $i \rightarrow \infty$ , when the functions  $f$ ,  $g$  and  $f_i$ 's are generalized probability density functions with the same carrier. A necessary condition was given by the same one as (2.1) above under the assumption that  $I(f : g)$  is finite, but Corollary 2.2 below does not require the finiteness of  $I(f : g)$ .

The notation  $(f_i, g_i; f, g) \Rightarrow (f_i, f; f, f)$ , for example, in Corollary 2.1 below means that functions  $f_i$ ,  $f$ ,  $f$  and  $f$  are taken instead of  $f_i$ ,  $g_i$ ,  $f$  and  $g$  in Theorem 1.1.

COROLLARY 2.1.  $(f_i, g_i; f, g) \Rightarrow (f_i, f; f, f)$ .

Assume that

(i)  $f$  and  $\{f_i\}$  ( $i=1, 2, \dots$ ) are generalized probability density functions with respect to  $m$ , with  $D(f)=D(f_i)=E(m)$  ( $i=1, 2, \dots$ ). Then, if the condition

$$(2.2) \quad h(f, f_i) = \text{ess sup}_E \left| 1 - \frac{f_i}{f} \right| \rightarrow 0 \quad (i \rightarrow \infty)$$

is satisfied, then it holds that

$$(2.3) \quad I(f_i : f) \rightarrow 0 \quad (i \rightarrow \infty).$$

COROLLARY 2.2.  $(f_i, g_i; f, g) \Rightarrow (f, g; f, g)$ .

Assume that

(i)  $f, g$  and  $\{f_i\}$  ( $i=1, 2, \dots$ ) are generalized probability density functions with respect to  $m$ , with  $D(f)=D(g)=D(f_i)=E(m)$  ( $i=1, 2, \dots$ ).

Under this assumption, if the condition

$$(2.4) \quad h(f, f_i) = \operatorname{ess\,sup}_E \left| 1 - \frac{f_i}{f} \right| \rightarrow 0 \quad (i \rightarrow \infty)$$

is satisfied, then it holds that

$$(2.5) \quad I(f_i : g) \rightarrow I(f : g) \quad (i \rightarrow \infty).$$

COROLLARY 2.3.  $(f_i, g_i; f, g) \Rightarrow (f, f_i; f, f)$ .

Under the assumption of Corollary 2.1, if the condition (2.2) is satisfied, then it holds that

$$(2.6) \quad I(f : f_i) \rightarrow 0 \quad (i \rightarrow \infty).$$

COROLLARY 2.4.  $(f_i, g_i; f, g) \Rightarrow (f, g_i; f, g)$ .

Assume that

(i)  $f, g$  and  $\{g_i\}$  ( $i=1, 2, \dots$ ) are generalized probability density functions with respect to  $m$ , with  $D(f)=D(g)=D(g_i)=E(m)$  ( $i=1, 2, \dots$ ).

Then, the condition

$$(2.7) \quad h(g, g_i) = \operatorname{ess\,sup}_E \left| 1 - \frac{g_i}{g} \right| \rightarrow 0 \quad (i \rightarrow \infty)$$

implies that

$$(2.8) \quad I(f : g_i) \rightarrow I(f : g) \quad (i \rightarrow \infty).$$

Finally, we consider two types of truncation of the generalized probability density function; suppose  $\mu$  and  $\nu$  are two probability measures on the measurable space  $(R, S)$  which are absolutely continuous with respect to  $m$ , with densities  $f(x)$  and  $g(x)$  and with  $D(f)=D(g)=E(m)$ . Let  $\{E_i\}$  ( $i=1, 2, \dots$ ) be a sequence of sets in  $S$  such that  $E_i \subset E(m)$  ( $i=1, 2, \dots$ ) and  $\nu(E - E_i) \rightarrow 0$  as  $i \rightarrow \infty$ . Define

$$(2.9) \quad \text{and} \quad f^i(x) = \begin{cases} f(x) & \text{on } E_i \\ 0 & \text{otherwise} \end{cases}$$

$$g^i(x) = \begin{cases} g(x) & \text{on } E_i \\ 0 & \text{otherwise} \end{cases} \quad (i=1, 2, \dots),$$

and

$$(2.10) \quad f^{(i)}(x) = \begin{cases} f(x)/\mu(E_i) & \text{on } E_i \\ 0 & \text{otherwise} \end{cases}$$

and

$$g^{(i)}(x) = \begin{cases} g(x)/\nu(E_i) & \text{on } E_i \\ 0 & \text{otherwise} \end{cases} \quad (i=1, 2, \dots).$$

For these truncated probability density functions, Theorem 1.1 shows also that

$$(2.11) \quad I(f^i : g^i) \rightarrow I(f : g) \quad (i \rightarrow \infty),$$

and

$$(2.12) \quad I(f^{(i)} : g^{(i)}) \rightarrow I(f : g) \quad (i \rightarrow \infty).$$

DEPT. OF MATH., COLLEGE OF SCI. AND ENG.,  
NIHON UNIVERSITY

#### REFERENCES

- [1] S. Kullback, *Information Theory and Statistics*, John Wiley & Sons, 1959.
- [2] S. Ikeda, "Continuity and characterization of Shannon-Wiener information measure for continuous probability distributions," *Ann. Inst. Stat. Math.*, Vol. XI (1959).